

1 **Money Doesn't Grow on Trees, But Forecasts Do: Forecasting Extreme**

2 **Precipitation with Random Forests**

3 Gregory R. Herman* and Russ S. Schumacher

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

5 *Corresponding author address: Gregory R. Herman, Department of Atmospheric Science, Col-
6 orado State University, 1371 Campus Delivery, Fort Collins, CO 80523.

7 E-mail: gherman@atmos.colostate.edu

ABSTRACT

8 Approximately eleven years of reforecasts from NOAA’s Second Genera-
9 tion Global Ensemble Forecast System Reforecast (GEFS/R) model are used
10 to train a contiguous United States (CONUS)-wide gridded probabilistic pre-
11 diction system for locally extreme precipitation, developed primarily using
12 the random forest (RF) algorithm. Locally extreme precipitation is quan-
13 tified for 24-hour precipitation accumulations in the framework of average
14 recurrence intervals (ARIs), with two severity levels: 1- and 10-year ARI ex-
15 ceedances. Forecasts are made from 0000 UTC forecast initializations for two
16 1200 UTC–1200 UTC periods: Days 2 and 3 comprising respectively forecast
17 hours 36–60 and 60–84. Separate models are trained for each of eight fore-
18 cast regions and for each forecast lead time. GEFS/R predictors vary in space
19 and time relative to the forecast point, and include not only the quantitative
20 precipitation forecast (QPF) output from the model, but also variables that
21 characterize the meteorological regime, including winds, moisture, and in-
22 stability. Numerous sensitivity experiments are performed to determine the
23 effects of the inclusion or exclusion of different aspects of forecast informa-
24 tion in the model predictors, the choice of statistical algorithm, and the effect
25 of performing dimensionality reduction via principal component analysis as a
26 pre-processing step. Overall, it is found that the machine learning (ML)-based
27 forecasts add significant skill over exceedance forecasts produced from both
28 the raw GEFS/R ensemble QPFs and from the European Centre for Medium-
29 Range Weather Forecasts’ (ECMWF) global ensemble across almost all re-
30 gions of CONUS. ML-based forecasts are found to be underconfident, while
31 raw ensemble forecasts are highly overconfident.

32 **1. Introduction**

33 Locally extreme precipitation can cause a variety of costly, disruptive, and endangering impacts,
34 including flooding, flash flooding, and landslides. In 2016 alone, these hazards combined caused
35 more than 120 fatalities and \$10 billion in damages over the United States (NWS 2017b). The pre-
36 diction of flash floods is a notoriously challenging forecast problem, requiring not only accurate
37 prediction of heavy rainfall magnitudes, but also of the spatiotemporal distribution of that rain-
38 fall; the hydrologic interactions between precipitation, terrain, and the land surface; and also of
39 antecedent precipitation and its effects on soil conditions. Forecasting precipitation processes re-
40 sponsible for most observed extreme rainfall over the contiguous United States (CONUS) is often
41 considered among the most challenging problems in contemporary numerical weather prediction
42 (NWP; e.g. Fritsch and Carbone 2004; Novak et al. 2014). Given that the rainfall forecast alone
43 presents such a considerable challenge, the additional hydrologic considerations in the flash flood
44 forecast problem present an even more daunting task. While recent advances in heavy rainfall and
45 flash flood forecasting have been made (e.g. Hapuarachchi et al. 2011; Novak et al. 2014; Barthold
46 et al. 2015), forecasts still struggle in many situations (e.g. Delrieu et al. 2005; Lackmann 2013;
47 Schumacher et al. 2013; Gochis et al. 2015; Nielsen and Schumacher 2016, among many others)
48 and substantial progress remains to be made.

49 Contemporary operational dynamical forecast models often struggle to accurately simulate the
50 physical processes responsible for its production. For example, models with parameterized con-
51 vection often have a variety of persistent errors and biases associated with their depiction of con-
52 vective systems, convective systems being responsible for the majority of flooding rains over much
53 of CONUS (e.g. Schumacher and Johnson 2006; Stevenson and Schumacher 2014; Herman and
54 Schumacher 2016a). These include a tendency to underpredict total rainfall from convective sys-

55 tems (e.g. Schumacher and Johnson 2008; Herman and Schumacher 2016a); produce systems
56 displaced too far to the north and west from where they are observed (e.g. Grams et al. 2006;
57 Wang et al. 2009; Clark et al. 2010); initiate convection too early (e.g. Davis et al. 2003; Wilson
58 and Roberts 2006; Clark et al. 2007); generate systems with too large an areal extent (e.g. Wilson
59 and Roberts 2006); and propagate them incorrectly, too slowly, or not at all (e.g. Davis et al. 2003;
60 Pinto et al. 2015). While convection-allowing models (CAMs) can better resolve the physical
61 processes responsible for heavy rainfall generation (e.g. Kain et al. 2006; Weisman et al. 2008;
62 Duda and Gallus 2013), they too can suffer from many of these biases (e.g. Kain et al. 2006; Lean
63 et al. 2008; Kain et al. 2008; Weisman et al. 2008; Herman and Schumacher 2016a). Furthermore,
64 although there is a plethora of CAM guidance out to the day-ahead time frame (out to 36 hours
65 to perhaps 48 hours after initialization), due to current computational constraints, there is almost
66 no operational CAM guidance running out to two days ahead, and nothing operational that runs
67 to three days ahead or beyond. Instead, global ensembles with parameterized convection serve
68 as the primary source of forecast information and uncertainty quantification at these lead times.
69 Even if CAMs were to integrate out to these time frames, it is not clear that the depictions would
70 provide more valuable forecast guidance due to the rapid upscale growth of smaller-scale errors
71 which accumulate with increasing forecast lead time (e.g. Zhang et al. 2003, 2007). Nevertheless,
72 there is considerable utility in skillful extreme precipitation forecasts at these longer lead times,
73 as preparative and mitigative actions may take time to execute and therefore be feasible multiple
74 days away from the event but not just hours before. For forecasting in the medium-range beyond
75 the day-ahead time scale, the aggregate consideration of these circumstances suggests the poten-
76 tial for considerable potential in applying statistical post-processing techniques to global ensemble
77 guidance; so doing may demonstrate the capability to alleviate or even eliminate many of these
78 dynamical model deficiencies. The increased skill at these lead times relative to even longer ones,

79 and operational production of extreme rainfall products at these lead times not generated even
80 farther out in time such as the Excessive Rainfall Outlooks produced by the Weather Prediction
81 Center (Barthold et al. 2015) further motivates a specific focus on the Day 2 and Day 3 lead times
82 for this application.

83 There is a long history of successful application of statistical post-processing to dynamical
84 model output (e.g. Klein et al. 1959; Glahn and Lowry 1972). Model Output Statistics (MOS;
85 e.g. Glahn and Lowry 1972), is a simple, effective multivariate linear regression technique relat-
86 ing a set of dynamical model predictors to sensible weather predictands such as minimum and
87 maximum temperature, wind speeds, and precipitation probability. This basic technique has long
88 demonstrated skill over both the underlying models and even human forecasters (e.g. Jacks et al.
89 1990; Vislocky and Fritsch 1997; Hamill et al. 2004; Baars and Mass 2005), but is inherently lim-
90 ited by the linear assumptions underlying the method. Statistical post-processing techniques have
91 also been successfully applied to QPFs, from early linear approaches (e.g. Bermowitz 1975; Anto-
92 lik 2000) to more contemporary techniques that can exploit more complex variable relationships,
93 including neural networks (e.g. Hall et al. 1999), reforecast analogs (e.g. Hamill and Whitaker
94 2006; Hamill et al. 2015), logistic regression (LR; e.g. Applequist et al. 2002), random forests (RF;
95 e.g. Gagne et al. 2014; Ahijevych et al. 2016; Gagne et al. 2017), and other parametric techniques
96 (e.g. Scheuerer and Hamill 2015). For other meteorological applications, other machine learning
97 algorithms, such as support vector machines (e.g. Zeng and Qiao 2011; Herman and Schumacher
98 2016b) and boosting (e.g. Herman and Schumacher 2016b; Hong et al. 2016) have also been suc-
99 cessfully applied. Related techniques have also been applied to forecasting related high-impact
100 phenomena, such as severe hail (Brimelow et al. 2006; Gagne et al. 2015) and tornadoes (Alvarez
101 2014). One of the most powerful aspects of machine learning algorithms—and RFs in particular—
102 is finding patterns and non-linear interactions in the supplied training data (e.g. Breiman 2001).

103 Depending on the extent and diversity of the data supplied in these experiments, trained RFs pose
104 the theoretical capability of diagnosing and automatically correcting for various kinds of model
105 biases, including context-dependent quantitative biases, such as QPF being systematically too high
106 or too low; spatial displacement biases in the placement of extreme precipitation features; and, to
107 some extent, temporal biases in the initiation or progression of extreme precipitation features.

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

117 **2. Data and Methods**

118 There are several successive steps applied in creating the final forecasts evaluated in this study.
119 A schematic overview of the forecast pipeline for the models trained in this study is depicted
120 in Figure 1. Many types of hydrometeorological information are first taken, then assembled in
121 a methodical manner, further pre-processed for subsequent analysis, analyzed using a statistical
122 machine learning algorithm, and finally, extreme precipitation forecast guidance is produced and
123 evaluated. This section details each of these steps in the model development and evaluation pro-
124 cess.

125 *a. Datasets*

126 Dynamical model data used for training the RF models in this study comes from NOAA's
127 Second-Generation Global Ensemble Forecast System Reforecast (GEFS/R; Hamill et al. 2013)
128 dataset. The GEFS/R is a global 11-member ensemble with parameterized convection and
129 T254L42 resolution—which corresponds to an effective horizontal grid spacing of ~ 55 km at 40°
130 latitude—initialized once daily at 0000 UTC back to December 1984. Perturbations are applied
131 only to the initial conditions, and are made using the ensemble transform with rescaling technique
132 (Wei et al. 2008). The ensemble system used to generate these reforecasts is nearly static through-
133 out its 30+ year period of coverage, though updates to the operational data assimilation system
134 over time have resulted in some changes in the bias characteristics of its forecasts over the period
135 of record (Hamill 2017). Some forecast fields are preserved on the native Gaussian grid ($\sim 0.5^\circ$
136 spacing), while others are available only on a $1^\circ \times 1^\circ$ grid. Temporally, forecast fields are archived
137 every three hours out to 72 hours past initialization, and are available every six hours beyond that.
138 This study employs an almost 11-year period of record to explore this forecast problem, using
139 daily initializations from January 2003 through August 2013.

140 In creating probabilistic extreme precipitation forecast guidance, the predictand must first be
141 concretely specified and a robust, consistent verification framework established. One of the many
142 challenges in heavy rainfall and flash flood forecasting is the considerable difficulty in verifying
143 events (e.g. Welles et al. 2007; Gourley et al. 2012; Barthold et al. 2015), as every approach has
144 its deficiencies and limitations. Flash flood reports, flash flood warnings, and quantitative precip-
145 itation estimate (QPE) exceedance of flash flood guidance (FFG)—a product issued routinely by
146 NWS River Forecast Centers (RFCs) which estimates the amount of precipitation over a prescribed
147 accumulation interval currently required at a given location to produce flash flooding—all attempt

148 to directly address the actual hydrologic impacts, but all suffer from a combination of intermittent
149 reporting (e.g. Pielke et al. 2002), lack of coverage (e.g. Marjerison et al. 2016), and substantial
150 regional differences in practices or methodology (e.g. Ntelekos et al. 2006; Schmidt et al. 2007;
151 Ashley and Ashley 2008; Villarini et al. 2010; Calianno et al. 2013). It is therefore attractive to
152 consider the problem from a simpler perspective by considering QPE exceedances of some tem-
153 porally static threshold, an approach which avoids these issues while knowingly neglecting other
154 important elements of the forecast problem. In particular, a fixed threshold (e.g. 50 mm hr^{-1}) can
155 be used as a proxy for flash flooding (e.g. Brooks and Stensrud 2000; Hitchens et al. 2013), as can
156 exceedances of thresholds defined relative to the local precipitation climatology (e.g. Schumacher
157 and Johnson 2006; Stevenson and Schumacher 2014; Herman and Schumacher 2016a), such as
158 average recurrence intervals (ARIs). An ARI defines a fixed frequency relative to the hydromete-
159 orological climatology of the region; in particular, it corresponds to the expected duration, given
160 the local climatology, between exceedances of a given threshold. For example, the 1-year ARI
161 for 24-hour precipitation accumulations describes the accumulation amount for which one would
162 expect the mean duration between exceedances of said amount to be one year. Past research has
163 shown that a fixed-frequency ARI-based framework has better correspondence with heavy pre-
164 cipitation impacts than the use of any fixed threshold across the hydrometeorologically diverse
165 regions of CONUS (e.g. Reed et al. 2007). From the perspective of forecast verification, defining
166 extreme precipitation with respect to a fixed threshold exceedance raises challenges when applied
167 uniformly across CONUS. Due to substantially varying frequency of event occurrence, skill dif-
168 ferences may be artificially reflected in an undesirable manner (e.g. Hamill and Juras 2006). The
169 ARI framework avoids this issue and provides reasonable correspondence with precipitation im-
170 pacts without considering the additional influential complications such as antecedent conditions,

171 local hydrology, and urban effects (e.g. Herman and Schumacher 2016a). Consequently, the ARI
172 framework is used to quantify extreme rainfall for these experiments.

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

186 In aggregating multiple QPE-to-ARI threshold grid point comparisons in a single predictand, the
187 forecasts issued correspond to neighborhood event probabilities, which is an increasingly popular
188 method of communicating probabilistic high-impact weather information in forecast operations
189 (e.g. Barthold et al. 2015; NWS 2017a). While the event frequency does increase relative to
190 the purported ARI, the fixed-frequency property and associatedly many of the aforementioned
191 desirable properties of the framework are approximately retained. For this study, focus is placed
192 exclusively on two 24-hour forecast periods: the 1200–1200 UTC period corresponding to forecast
193 hours 36–60 from the GEFS/R forecast fields and the subsequent 24-hour period encompassing
194 forecast hours 60–84, denoted respectively as Day Two and Day Three. At these times, there is

195 typically some knowledge to characterize the environmental conditions in which precipitation may
196 form, but it is beyond the current range of operational CAM guidance.

197 The National Centers for Environmental Prediction (NCEP) Stage IV Precipitation Analysis (Lin
198 and Mitchell 2005) is an operational QPE product created daily since December 2001. Stage IV
199 provides 24-hour analyses over the CONUS on a \sim 4.75 km grid. It uses both rain gauge observa-
200 tions and radar-derived rainfall estimates to generate an analysis, and is further quality controlled
201 via NWS River Forecast Centers (RFCs) to ensure stray radar artifacts and other spurious anoma-
202 lies do not appear in the final product. Despite some limitations (Herman and Schumacher 2016a;
203 Nelson et al. 2016), its analysis quality; resolution, allowing better ability to capture precipitation
204 extremes compared with other QPE products (e.g. Hou et al. 2014); and data record length make
205 it preferable to other precipitation analysis products, and is therefore be used as the precipitation
206 ‘truth’ for verification purposes in this study.

207 The ARI thresholds associated with the 1- and 10-year ARIs for 24-hour precipitation accumu-
208 lations are generated using the same methodology of Herman and Schumacher (2016a), where
209 CONUS-wide thresholds are produced by stitching thresholds from several sources. NOAA’s
210 Atlas 14 thresholds (Bonnin et al. 2004, 2006; Perica et al. 2011, 2013), an update from older
211 work and currently under development, are used wherever were available at the commencement of
212 this study. For five northwestern states—Washington, Oregon, Idaho, Montana, and Wyoming—
213 updated thresholds are not available, and derived NOAA Atlas 2 threshold estimates are used
214 instead (Miller et al. 1973). Additionally, in Texas and the Northeast—New York, Vermont, New
215 Hampshire, Maine, Massachusetts, Connecticut, and Rhode Island—Technical Paper 40 (TP-40;
216 Hershfield 1961) thresholds are used¹; everywhere else uses the Atlas 14 threshold estimates.

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

217 These 1- and 10-year thresholds are depicted in Figure 2 in panels (a) and (b), respectively. The 10-
218 year ARI thresholds in Figure 2b show a similar pattern to the 1-year ARI thresholds in the same
219 location shown in Figure 2a, but are substantially higher everywhere. More significantly, it is ap-
220 parent that at both severity levels, there are large regional disparities in the threshold magnitudes.
221 Over climatologically wet regions of CONUS, such as the Pacific coastal mountains and imme-
222 diately along the Gulf Coast, thresholds are as high as 100–150 mm and 250–300 mm for 1-year
223 and 10-year ARIs, respectively. Thresholds, particularly over central and eastern CONUS, tend
224 to decrease smoothly with increasing latitude and distance from major bodies of water. Sharper
225 variations are seen in areas of complex terrain over western CONUS. In the driest parts of the
226 arid southwest and intermountain west, thresholds can be as low as 10–15 mm and 25–30 mm for
227 the two ARI levels—a full order of magnitude difference from the largest thresholds at the same
228 intensity level. This highlights the stark contrast the use of the ARI framework has with the use of
229 a fixed threshold framework across all of CONUS.

230 Forecast models in this study are trained separately for eight distinct, yet cohesive and internally
231 fairly hydrometeorologically homogeneous regions of CONUS, using the delineation indicated
232 in Figure 3. Observed 1- and 10-year ARI exceedance events that occurred during the period of
233 record are depicted by region respectively in Figures 2 (c) and (d). There are important regional
234 differences in the seasonal cycles and climatology of ARI exceedances across CONUS. In the Pa-
235 cific Coast (PCST) region, the vast majority of exceedances at both the 1-year and 10-year severity
236 levels occur in the cool-season, and occur largely from atmospheric river events with large mois-
237 ture transport impinging on coastal topography (e.g. Rutz et al. 2014; Herman and Schumacher
238 2016a). This seasonality holds to a lesser extent in the neighboring Southwest (SW) region, with
239 some signal carrying over to the Rockies (ROCK) region as well. In the central and eastern re-
240 gions, the majority of events occur during the warm-season from more scattered convective-scale

241 processes, particularly in the months of May, June, and July (e.g. Schumacher and Johnson 2006;
242 Herman and Schumacher 2016a). Tropical cyclones can cause widespread and very significant
243 rainfall, and comprise a substantial portion of the extreme precipitation climatology, especially
244 in the Northeast (NE) and Southeast (SE) regions. Due to the spatial extent of their impacts and
245 immense rainfall totals they can produce, they form a much larger fraction of the climatology of
246 10-year ARI exceedances (Fig. 3d) than 1-year events (Fig. 3c). Additionally, the numbers are
247 lower than would be expected; by the explicit exceedance frequencies associated with the thresh-
248 olds, one would expect an average of one exceedance per point per year over the period of record
249 for the 1-year events (Fig. 3c) and 0.1 exceedances for 10-year events (Fig. 3d); in reality, event
250 counts are only about half of that. This is consistent with previous findings (e.g. Herman and
251 Schumacher 2016a), and likely in part attributable to limitations in the Stage IV product to capture
252 extremes (e.g. Nelson et al. 2016). There is also quite a bit of region-to-region variability in event
253 counts, particularly for 10-year exceedances, much of which is attributable to statistical variability
254 from having a short data record in relation to the event frequency.

255 *b. Predictor Assembly*

256 Input predictors, or features, to the random forests can be partitioned into two categories: model
257 predictors and background predictors. Model predictors, which constitute the vast majority of
258 the total number of inputs, come from atmospheric fields forecast in the GEFS/R which bear a
259 known physical relationship with extreme precipitation. A core set of $f=9$ fields used in this study
260 are: accumulated precipitation (APCP), convective available potential energy (CAPE), convective
261 inhibition (CIN), precipitable water (PWAT), surface temperature (T2M) and specific humidity
262 (Q2M), surface zonal (U10) and meridional winds (V10), and mean sea level pressure (MSLP).
263 Sensitivity experiments explore the use of additional upper-air atmospheric fields; a full list of

264 fields used in this study, their associated symbols used in this manuscript, and the grids on which
265 they are each archived are tabulated in Table 1. In addition, considering different atmospheric
266 fields in forecasting extreme precipitation, the spatiotemporal variations in these fields are consid-
267 ered as well. Spatially, predictors are structured in a forecast-point relative sense. In the control
268 model, GEFS/R forecast values up to $r=4$ grid boxes ($\sim 2^\circ$) latitudinally or longitudinally dis-
269 placed in any direction relative to the forecast point are considered. Temporally, portrayed fields
270 are considered at each archived time during the forecast interval, which corresponds to every three
271 hours during the Day Two period and every six hours during the Day Three period, for a total
272 of $t=9$ and $t=5$ forecast periods for the Day 2 and 3 periods, respectively. All told, this yields
273 $tf(2r + 1)^2$ model predictors, which yields respectively $M=6,561$ and $M=3,645$ model predictors
274 for the Day Two and Day Three control models. The other category of predictors, background
275 predictors (Table 2), are those which are solely associated with the forecast point, and have no
276 relation to the present meteorology. These include the location of the point, as well as the ARI
277 characteristics of the point and in the surrounding area.

278 *c. Dimensionality Reduction*

279 There are a large number of model predictors, and they are also highly correlated—spatially,
280 temporally, and across variables. With millions of training examples and thousands or, in some
281 cases, tens of thousands of features, the forecast problem can become computationally intractable.
282 Further, having many highly correlated features can readily result in model overfitting—making
283 predictions based on noise affecting an individual native feature rather than the underlying signal—
284 a phenomenon commonly termed the “curse of dimensionality” (e.g. Friedman 1997). There are
285 numerous ways these concerns can be addressed; broadly speaking, the most common approaches
286 are either feature selection or feature extraction. In feature selection, a subset of initial predictors

287 are chosen that collectively bear the strongest predictive relationship with the predictand, whereas
288 in feature extraction, a smaller set of new predictors are derived from the original set. Both of
289 these procedures can be performed subjectively through manual means or objectively through au-
290 tomated means. In this case, all of the input predictors are believed to have a physical relationship
291 with extreme precipitation, and choosing only the most predictive fields (e.g. model QPF) and
292 discarding the rest risks removing valuable predictive information not contained in the retained
293 predictor set. The primary issue with the input predictors in this case is not that many may not
294 have any physical bearing on the predictand, but rather that each predictor represents a value at
295 a different point of a continuous field, or a different property at the same point, and are thus
296 necessarily highly correlated to one another. Furthermore, while one could conceivably extract
297 features using field averages or some other pre-determined method, this may not be optimal. For
298 example, it may be better to weight values closer to the forecast point more heavily, while still
299 retaining some information from the far-field predictors. Holistically assessed, this suggests that
300 an objective approach to feature extraction rates to be the most beneficial approach to dimension-
301 ality reduction. Though it has some limitations (e.g. Shlens 2014), principal components analysis
302 (PCA; Ross et al. 2008; Pedregosa et al. 2011) is a robust and frequently utilized approach for
303 dimensionality reduction by feature extraction. This creates a small set of uncorrelated predictors
304 that explain the signal in the forecast data and gives insight into the regional modes of atmospheric
305 variability as depicted in the GEFS/R model (explored in more depth in a companion paper), while
306 leaving the noise in lower-order principal components (PCs), acting in principle to both alleviate
307 overfitting and manage computational requirements.

308 *d. Machine Learning Algorithms and Sensitivity Experiments*

309 The primary statistical algorithm used in this study is random forests (RFs; Breiman 2001). RFs
310 are in essence an ensemble of *decision trees*, where traditionally each tree individually makes
311 a deterministic prediction about the outcome of the predictand; the relative frequencies of each
312 possible predictand outcome in the ensemble of trees are then used to make a probabilistic forecast.
313 Much further detail on tree and RF construction and mechanics can be found in Appendix A.
314 There are also several parameters which can be tuned to the particular forecast problem in order
315 to maximize model performance. Four-fold cross-validation is used for model development in
316 this study, whereby each model configuration examined is trained four times, once each on three-
317 quarters of the training data, and then evaluated on the final withheld quarter. Each quarter of
318 the training data is a temporally contiguous chunk to avoid issues of sample independence and
319 approximately mimic information that would be available in an operational context. All parameter
320 settings and sensitivity experiments are evaluated in this framework. The set of RF parameters
321 tuned is described in Appendix A, and the results presented in Appendix C.

322 In this study, there are a great deal of dynamical model data considered as input information
323 on which the RF can base a prediction. An important question in the design of statistical post-
324 processing algorithms in a resource-constrained environment is which of this information actually
325 improves forecasts, and the relationship between extent of additional dynamical model informa-
326 tion and statistical model skill. In order to obtain answers to these important questions, a suite of
327 sensitivity experiments are conducted as summarized in Table 3. Sensitivity to the inclusion of
328 horizontal variations in atmospheric fields is explored by varying the previously described predic-
329 tor radius parameter R from 0 to 4. The effect of including predictor information from additional
330 upper-air fields is explored by comparing the inclusion and exclusion of two sets of fields. The first

331 incorporates temperature, specific humidity, zonal and meridional winds at 850 and 500 hPa, and
332 850 hPa vertical velocity in the so-called Upper-Air Core predictor group, while further inclusion
333 of those same fields at 700 and 250 hPa—the so-called Upper-Air Extra predictor group—allows
334 a reasonably comprehensive view into the value of this forecast information. Investigation into
335 the utility of temporal resolution of simulated atmospheric fields in generating skillful forecasts
336 is also performed. Predictor density is three-hourly for Day 2 guidance and six-hourly for Day
337 3 guidance; models are additionally trained with predictors at twelve-hourly temporal density for
338 both lead times and six-hourly temporal density for the Day 2 forecast model to ascertain this as-
339 pect. Another important element—and one with implications for how operational centers allocate
340 their computational resources—is the extent of use of ensemble information and its implication
341 on RF model skill. Using forecast information from only the GEFS/R’s control member in model
342 training is compared with using the ensemble median from the full ensemble, and then further
343 with the use of the ensemble second-lowest and second-highest values for each atmospheric field
344 in conjunction with the median to evaluate the impact of this dimension of forecast information,
345 following the findings of Herman and Schumacher (2016b), which found relatively little sensitiv-
346 ity in performance with respect to how ensemble information is used, but using the near-minimum,
347 median, and near-maximum values outperformed using the mean and spread. Additionally, where
348 possible, models are trained with and without the aforementioned PCA pre-processing step, and an
349 assessment of the effect of this pre-processing step on model skill is made by comparing the two.
350 An additional sensitivity experiment explores the effect of region size on forecast skill, hypothesiz-
351 ing models trained for larger regions may exhibit higher skill due to more available training data.
352 This is performed by aggregating the ROCK and SW regions into a new WEST one, combining
353 the Southern Great Plains (SGP), Northern Great Plains (NGP), and Midwest (MDWST) regions
354 into a CENTRAL region, and collecting SE and NE regions into a single EAST region, while

355 leaving PCST, with its highly unique extreme precipitation climatology unperturbed. A final sen-
356 sitivity experiment investigates model performance as a function of model algorithm, specifically
357 by comparing with logistic regression (LR), a common and comparatively simpler alternative to
358 statistically deriving forecast probabilities. Further discussion of LR and other machine learning
359 alternatives to the RF algorithm is included in Appendix B.

360 *e. Model Evaluation*

361 Based on the parameter tuning and sensitivity experiment results, final model configurations are
362 selected. The final model is run over a completely withheld 4-year evaluation period spanning
363 September 2013–August 2017. The forecasts generated from the final model are compared with
364 those from the full ensemble of raw GEFS/R QPFs, as well as the full 50-member ECMWF global
365 ensemble, accessed from TIGGE (Molteni et al. 1996; Bougeault et al. 2010). The comparison
366 with the former provides an assessment of what improvement, if any, these models yield compared
367 with the raw guidance from which their forecasts are derived when evaluated in a real-time setting.
368 The latter, meanwhile, provides an assessment for how these forecasts compare with state-of-the-
369 science operational ensemble guidance available at these lead times. To make these comparisons,
370 the QPF from each ensemble member of the two ensembles is regridded onto the ~ 4.75 km Stage
371 IV HRAP grid on which the Atlas thresholds lie using a first-order conservative scheme (Ramshaw
372 1985). These regridded QPFs are then compared with the 1-year and 10-year ARI thresholds to
373 create deterministic exceedance forecasts with respect to the two thresholds for each ensemble
374 member. These binary grids are then upscaled to the GEFS/R grid using the same procedure as
375 the verification upscaling: any exceedance in the downscaled grid corresponds to an exceedance
376 at the nearest GEFS/R point in the upscaled grid. Since the predictand categories are necessarily
377 mutually exclusive, the 1-year ARI exceedance grids are modified so that any member forecasting

378 a 10-year ARI exceedance at a point is not forecasting a between 1-and-10-year exceedance at that
 379 same point and time period. The prevailing operational method of generating forecast probabili-
 380 ties from a dynamical ensemble—democratic voting, whereby the fraction of ensemble members
 381 forecasting the event is used as the forecast probability (e.g. Buizza et al. 1999; Eckel 2003)—is
 382 applied to each ensemble to generate the exceedance probabilities for the reference forecasts.

383 Skill, both in the final assessment of model performance as well as in all aforementioned sen-
 384 sitivity experiments, is quantified by means of the Rank Probability Skill Score (RPSS) with a
 385 climatological reference:

$$RPSS = 1.0 - \frac{\sum_{d=1}^D (\sum_{p=1}^P (\sum_{m=1}^K (\sum_{j=1}^m P_{jpd} - O_{jpd})^2))}{\sum_{d=1}^D (\sum_{p=1}^P (\sum_{m=1}^K (\sum_{j=1}^m P_{clim_j} - O_{jpd})^2))} \quad (1)$$

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

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

393 where there are $N = DP$ total forecasts, broken into C discrete probability bins with N_c forecasts
 394 being issued for each bin c. $\overline{O_{j^*}}$ denotes the climatological (based on the period of record) fre-
 395 quency of observing event category j^* and $\overline{O_{cj^*}}$ denotes the proportion of forecasts in probability
 396 bin c observing event category j^* , where j^* is the aggregation of event categories of at least j in
 397 the RPSS framework. $\overline{O_{j^*}} (1 - \overline{O_{j^*}})$, the so-called “uncertainty” term, also represents the BS of a

398 climatological forecast. Converting to a Brier skill score (BSS) framework by dividing out by this
 399 term:

$$BSS_{j^*} = 1.0 - \frac{BS_{j^*}}{BS_{clim_{j^*}}} = \frac{\frac{1}{N} \sum_{c=1}^C N_{cj^*} (\overline{O_{j^*}} - \overline{O_{cj^*}})^2}{\underbrace{\overline{O_{j^*}}(1 - \overline{O_{j^*}})}_{\text{"Resolution"}}} - \frac{\frac{1}{N} \sum_{c=1}^C N_{cj^*} (P_{cj^*} - \overline{O_{cj^*}})^2}{\underbrace{\overline{O_{j^*}}(1 - \overline{O_{j^*}})}_{\text{"Reliability"}}}} \quad (3)$$

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

401 3. Results: Sensitivity Experiments

402 Examining forecast skill as a function of time step between atmospheric field predictors (i.e. the
 403 CORE_LTIME models of Table 3; Fig. 4a), two striking findings concern 1) the large variations
 404 in forecast skill across regions and 2) the evidently low sensitivity of time step length and forecast
 405 skill within any given region. For the 3-hour time step, predictors are gathered from a total of 9
 406 forecast times; with the 6-hour step, 5 forecast times are used; and with the 12-hour time step, a
 407 total of 3 forecast times are used. The 12-hour time step therefore has one-third the total number
 408 of predictors as the model with the 3-hour time step, but still yields nearly identical forecast skill
 409 results. In most regions and forecast periods, there is a slight degradation in performance going
 410 from the 6- to 12-hour time step, but the difference is not generally statistically significant by a
 411 99% bootstrap skill score difference test (not shown). The one exception to this is in the PCST
 412 region, which has much higher skill overall than the other regions for both forecast periods, and
 413 exhibits somewhat higher sensitivity to the predictor time step than the other regions, particularly
 414 in going from 6-hours to 12-hours, with skill differences of approximately 0.01.

415 Similar to the temporal resolution findings, a general lack of sensitivity as a function of predictor
 416 spatial extent (Fig. 4b). This finding comes in stark contrast to that of Herman and Schumacher
 417 (2016b), which found great sensitivity of predictor spatial extent in forecasting airport flight rule
 418 conditions. Albeit weak, a slight improvement in skill for most forecast period, region combi-

419 nations can be noted with increasing predictor radius, often to the extent that the skill difference
420 between 0 and 4 grid box radii is statistically significant (not shown). Two regions in particular,
421 the NE and PCST, exhibit by far the most sensitivity to predictor spatial extent, with differences of
422 roughly 0.02 observed over the evaluated interval. Also of note is that a radius of 4 grid boxes—
423 the highest number evaluated—did not always yield the best performance results; most notably,
424 the Day 2 model for the NE region maximized skill at a radius of 2, with a slight deterioration
425 of forecast skill with increasing radius thereafter. In those regions where the GEFS/R cannot ex-
426 plicitly resolve the processes responsible for producing extreme precipitation, the RF is ultimately
427 making forecasts more on environmental factors; these do not vary drastically in time or space,
428 and thus a single number or small set of numbers at or immediately surrounding the forecast point
429 are sufficient to characterize the basic properties of the environment, and this is all that the RF is
430 really using for much of its predictions. However, in regions impacted more readily by larger scale
431 systems where the dynamical model can more directly simulate the precipitation processes such as
432 PCST and the NE, the spatial variations in atmospheric fields carry more signal rather than noise
433 and thus contribute more predictive value.

434 Like varying spatial and temporal density, there is relatively little sensitivity to the inclusion of
435 more atmospheric fields (Fig. 5a). Slight but consistent improvement is observed in adding the
436 core upper-air fields as predictors, but adding further levels beyond the core group was found to
437 not improve predictive skill, and actually resulted in a decrease in skill for the PCST, NE, ROCK,
438 and SE regions—those which are most affected by larger scale precipitation systems. Though still
439 rather small, somewhat more distinct sensitivity to type of ensemble information included (Fig.
440 5b) can be seen here across all regions, with improvements seen using only predictor information
441 from the GEFS/R the ensemble median versus the control member, and slight additional improve-
442 ment using the ensemble second-from-minimum and second-from-maximum in addition to the

443 ensemble median. The largest differences in magnitude are again for the PCST region, but in this
444 experiment, clear and statistically significant differences (not shown) are also seen for low skill,
445 convectively active regions such as MDWST.

446 Aggregating regions (Fig. 5c) results in a slight degradation in forecast skill. While predictor in-
447 formation is available to pinpoint the locations via latitude and longitude and then perform further
448 analysis from there, effectively performing the manual regional partitioning in the CTL_NPCA
449 model, these findings demonstrate that there is some—albeit limited—utility in manually parti-
450 tioning training data with distinct hydrometeorological relationships, rather than relying on the
451 machine learning algorithm to discern the distinction automatically. Comparing the impact of ap-
452 plying PCA pre-processing to the RF (Fig. 5d, leftmost two columns), performing PCA tends to
453 either improve performance, as is the case for the PCST, NE, SW, and MDWST regions, or make
454 little difference, as seen in the ROCK, NGP, SGP, and SE regions. The positive differences tend to
455 be larger in magnitude, both in relative and absolute senses, for Day 2 model versions compared
456 with Day 3. Forecasts produced through LR tend to be substantially worse than those generated by
457 RFs (Fig. 5d, center columns). However, the exact magnitude to which this is the case varies by
458 region; substantial differences in skill are seen between RF and LR forecasts for the SW, ROCK,
459 and SGP regions, while there is almost no skill difference between the Day 3 forecasts in the PCST
460 region. This may suggest the linear assumptions inherent to the LR algorithm perform better in
461 larger scale systems than in the more convectively active ones in which the responsible processes
462 are highly nonlinear, but this causality is not entirely clear. Finally, a weighted average of RF and
463 LR forecasts outperforms its component members for all regions and forecast periods. The extent
464 of overperformance is strongly tied to the skill difference between the RF and LR models; when
465 the skill difference is small, the value of the weighted average is comparatively large to when the
466 RF performs much better than LR (cf. Fig. 5d PCST and SW lines). Since these weighted av-

467 erages performed the best in cross-validation, a weighted average using each of the CTL_NPCA,
468 CTL_PCA, and CTL_LR models was chosen for the final model configuration.

469 **4. Results: Final Model Performance**

470 For both the final ML models and the forecasts from the raw QPFs of both the GEFS/R and
471 ECMWF (Fig. 6), a usually statistically significant² deterioration in forecast skill from Day 2 to
472 Day 3 is evident in each CONUS region over the four year test period. Forecast skill is significantly
473 higher in regions with extreme precipitation associated partially or primarily with synoptic scale
474 precipitation episodes, such as PCST, SW, and ROCK, rather than smaller scale convective systems
475 that characterize extreme precipitation in the NGP, SGP, and MDWST regions. At an extreme, the
476 NGP and SGP GEFS/R raw QPFs have no skill in predicting 1-year and 10-year ARI exceedances.
477 Especially for the ML models, the bigger Day 2 vs. Day 3 skill differences are also seen where the
478 skill is higher, again suggesting the direct forecasting of the precipitation as opposed to forecasts
479 more reflecting the forecast environment, either dynamically via parameterized convection in the
480 case of raw QPFs, or directly in the case of the ML model forecasts. Furthermore, the ML models
481 exhibit a larger skill deterioration between Days 2 and 3 than either of the raw ensemble forecast
482 sets.

²For all of the comparisons in this section, 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 has been included in the figures, much of the statistical significance difference results discussed in-text are omitted for the sake of concision.

483 Comparing the forecast systems, the ECMWF forecasts consistently and statistically signifi-
484 cantly outperform the GEFS/R forecasts at all lead times except in the SE region (Fig. 6). Encour-
485 agingly, the ML model forecasts are statistically significantly more skillful for all eight regions
486 and both lead times compared with the GEFS/R forecasts from which they are based. The post-
487 processing is thus clearly accomplishing its purpose of improving forecast skill. But it is also
488 apparent that the GEFS/R is not a state of the science model for extreme QPF prediction given its
489 lower skill compared with the ECMWF. The real test of the ML model then is how it compares
490 with current best operational guidance for these lead times, represented here with the ECMWF
491 ensemble. The comparison (Fig. 6) is generally quite favorable, with the Day 3 ML forecasts
492 outperforming the Day 2 ECMWF forecasts across all regions except ROCK and PCST. In the
493 non-western regions, the extent of overperformance is quite considerable when comparing equal
494 lead times, with skill score improvements of factors of two to three seen in many comparisons. In
495 the ROCK and PCST regions, the ML and ECMWF forecasts performed about equally at Day 2,
496 and ECMWF performed slightly better at Day 3. Overall though, the ML models demonstrated
497 ability to consistently outperform current operational model guidance, especially in convectively
498 active regions where there is no operational guidance that can dynamically resolve the physical
499 processes producing extreme precipitation.

500 Reliability diagrams of raw GEFS/R forecasts (Fig. 7) reveal highly overconfident probabilistic
501 exceedance forecasts for all regions, lead times, and severity levels as evidenced by the shallow
502 slope relative to the one-to-one line. The raw GEFS/R forecasts are relatively sharp, with more
503 than 0.01% of forecasts falling into each probability bin above 10%, and a vast majority of zero
504 probability forecasts (not shown). For all regions, there are cases where every ensemble member
505 has simultaneously predicted a 1-year exceedance (Fig. 7a,c), but the same is not true for 10-
506 year exceedance predictions in the northeastern regions: NE, NGP, and MDWST (Fig. 7b,d).

507 The ECMWF (Fig. 8) is also overconfident, but we see that it is also negatively biased for all
508 regions, lead times, and severity levels. Its degree of overconfidence is dampened compared with
509 the GEFS/R, and it is not as sharp, with fewer very occurrences of very high forecast probabilities
510 except in the westernmost regions of ROCK and PCST (Fig. 8a,c, inset panels). With 50 members
511 rather than 11, there is also substantially more resolution across the probability spectrum in the
512 ECMWF forecasts. By the very nature of how these forecasts are generated, quite a bit of sharpness
513 is inherent at the cost of reliability, since it is not possible for probabilities near the climatological
514 event frequency to be issued for either raw ensemble, particularly for the GEFS/R.

515 The reliability diagrams for the different components of the final model, the CTL_NPCAs,
516 CTL_PCA, and CTL_LR models are shown respectively in Figures 9, 10, and 11. The CTL_NPCAs
517 (Fig. 9) shows markedly different characteristics than either of the raw ensembles. In particu-
518 lar, all of the regions exhibit an underconfidence signal, with low probability events below about
519 2% for 1-year events (Fig. 9a,c) occurring with observed relative frequencies below the forecast
520 probabilities. The relative event frequencies are conversely appreciably higher than the forecast
521 probabilities would indicate for probabilities above 5%. The PCST probabilities are the most neg-
522 atively biased, while NE probabilities are the most positively biased, among the regions. Overall,
523 reliability is much better than for either raw ensemble, but this comes at the expense of sharpness.
524 Less than 1 in 10,000 forecasts are above about 20% for 1-year exceedance probabilities for any
525 CTL_NPCAs model (Fig. 9a,c, inset panels), or 5% for 10-year events (Fig. 9b,d, inset panels), and
526 maximum probabilities are in the 30–80% range for 1-year exceedance forecasts depending on the
527 lead time and region, compared with 100% for all lead times and regions in the raw ensembles.
528 For the 10-year events, maximum probabilities range from 10%–40%, again much lower than the
529 maximum probabilities from the raw ensembles, which were still at or near unity (e.g. Fig. 7b,
530 8b).

531 The CTL_PCA model (Fig. 10) exhibits very similar reliability characteristics to the CTL_NPCA
532 model (cf. Fig. 9, 10), including the underconfidence, reduced sharpness compared with the raw
533 GEFS/R, and different regional probability bias characteristics. The probability distribution for
534 1-year exceedance events is not markedly different between the Day 2 and Day 3 forecasts (cf.
535 Fig. 10a,c), but the relatively higher probabilities issued for 10-year exceedances in Day 2 do not
536 occur at the Day 3 lead times (cf. Fig. 10b,d). This is consistent with increasing confidence in
537 very extreme events with decreasing lead time; something seen very pronounced in the CTL_PCA
538 model and to a slightly lesser extent in the CTL_NPCA model, but to a much lesser extent in the
539 raw ensemble forecasts. The CTL_LR model (Fig. 11) exhibits some similarities and some dif-
540 ferences with the RF-based models. PCST forecasts are consistently the most negatively biased,
541 followed by ROCK and the SE, with NE region forecasts being the least negatively biased. How-
542 ever, unlike the RF-based forecasts, the LR model issues more high probabilities; at Day 2, for
543 example (Fig. 11a), forecasts in the highest probability bin were issued for most regions. At the
544 highest probabilities, the forecasts revert to being positively biased, as they are for events with
545 probabilities issued in the 0.01–1% range. At very low probabilities (not shown), LR-based fore-
546 casts are substantially more negatively biased than for RF-based forecasts, leading to considerable
547 overconfidence overall when considering that the vast majority of forecasts issued occur on this
548 low probability end of the spectrum. While LR (and regression in general) is effective at removing
549 bias in a global sense, since a single regression equation must necessarily apply globally to all
550 forecasts, it inherently cannot perform more localized, context-depend forms of bias correction,
551 leading to forecast probability-dependent model biases. Finally, the final ML model reliability
552 (Fig. 12) unsurprisingly reflects a blend of the component members, retaining some of the under-
553 confidence of the RF-based models while adding a bit of sharpness from the CTL_LR model in

554 regions where it verified skillfully enough in cross-validation (e.g. PCST, Fig. 5d) to garner much
555 weight.

556 The relationship between the reliability analysis and skill via the Brier score decomposition
557 (Murphy 1973) quantitatively solidifies many of the general observations discerned by inspection
558 of the reliability diagrams. Though sharper than competing forecasts, the raw GEFS/R forecasts
559 consistently exhibit the worst resolution component contribution to forecast skill for all regions and
560 severity levels, both for Day 2 forecasts (Fig. 13a,c) and Day 3 forecasts (Fig. 14a,c) due to an
561 ability to actually distinguish events from non-events by resolving the responsible physical mech-
562 anisms. The final ML models exhibit better resolution term skill contributions than the ECMWF
563 ensemble forecasts, with the exception of the ROCK and NGP regions for 1-year events (Fig. 13a,
564 14a). Between the component models, resolution term skill tended to best for CTL_NPCA fore-
565 casts over the test period, particularly at the 10-year severity level (e.g. Fig. 13c) but the extent of
566 the difference tended to be relatively small and there were numerous instances where PCA-based
567 models exhibited more resolution. The weighted average consistently exhibited higher resolution
568 than any of the component members. With respect to the reliability contribution to skill (Fig.
569 13b,d Day 2; Fig. 14b,d Day 3), ECMWF forecasts were, perhaps surprisingly given the lack of
570 explicit calibration, the most reliable forecast set for all regions and lead times, while in many
571 cases the ML models had a more negative contribution to the total skill than the raw GEFS/R,
572 likely resulting from the underconfidence. The resolution term is at largest one and at least zero
573 in this decomposition, while the reliability term is at most zero. The magnitude of the resolution
574 terms is consistently several factors larger than the reliability term for all forecast sets, and the
575 differences in that term are generally of more absolute impact on the overall Brier skill scores.

576 Lastly, while by no means a comprehensive characterization of the system, a sample of real
577 cases over the test period are presented to illustrate some of the strengths and weaknesses of the

578 system. On the evening of 19 May 2015 and morning of 20 May 2015, a vigorous mesoscale
579 convective system developed over southern Oklahoma and northern Texas, producing very heavy
580 rainfall that contributed to historic flooding in the region during May 2015 (e.g. Wolter et al.
581 2016). Stage IV analysis (Fig. 15a) reveals that the 24-hour precipitation totals exceeded 1-year
582 ARI thresholds within much of an E/W band encompassing the region, with embedded areas of 10-
583 year exceedances along the state border region. While the ECMWF ensemble forecasts indicate
584 some possibility of extreme precipitation in that region during this time frame three days out (Fig.
585 15c), the probabilities are displaced too far to the south and west, and the probabilities of 10-year
586 exceedances are very low. There is some improvement in positioning with the Day 2 forecast (Fig.
587 15b), but remains too far west, with probabilities still quite low particularly at the 10-year ARI
588 level. Raw GEFS/R forecasts at Day 3 (Fig. 15e) indicate quite high risk for a 1-year exceedance
589 over a fairly narrow area, better positioned than the ECMWF ensemble at the same lead time but
590 still too far to the west. Outside of this area, the GEFS/R indicates almost no risk of an extreme
591 rainfall event, and also indicates no risk of a 10-year exceedance anywhere in the domain. The
592 Day 2 forecast (Fig. 15d) looks similar to the Day 3 outlook, except that the probabilities are
593 reduced somewhat in the target area, which also has incorrectly displaced further to the south and
594 west. The ML model depicts a much different picture. It exudes much less confidence, with lower
595 maximum probabilities compared with either raw ensemble, but non-zero exceedance probabilities
596 of both 1- and 10-year exceedances across much of the domain for both Days 3 (Fig. 15g) and Day
597 2 (Fig. 15f). Importantly, the model elevated probabilities compared with the raw guidance in the
598 place that extreme precipitation was actually observed (to the east of where it was forecast in the
599 GEFS/R). In fact, at Day 2 (Fig. 15f), the probability maximum is located right where the heaviest
600 precipitation actually occurred, displaced well to the north and east of where it was forecast in
601 the GEFS/R (Fig. 15d). Additionally, while still low, the 10-year event probabilities are correctly

602 much higher over the verifying area when compared with either raw ensemble, with maximum Day
603 2 probabilities of around 30% and 3% for 1-year and 10-year exceedances, respectively. Finally, in
604 contrast to the raw guidance, the ML model became increasingly confident in an event occurring
605 with decreasing lead time (cf. Fig. 15f,g).

606 A different mesoscale precipitation produced extreme precipitation over southwestern Wiscon-
607 sin, southeastern Minnesota, and northeastern Iowa during the evening hours and overnight hours
608 of 21 September and 22 September 2016, respectively. Much of the area experienced 1-year ARI
609 exceedances for the 24-hour period ending 1200 UTC 22 September 2016, and within the 1-year
610 exceedance area, there were many embedded cells that produced 10-year ARI exceedances (Fig.
611 16a). ECMWF forecasts at Day 3 indicated risk of extreme rainfall, even at the 10-year severity
612 level (Fig. 16c), but the location was poor, with exceedance probabilities high in eastern Min-
613 nesota and northern Wisconsin where extreme rainfall was not observed, and very low probabili-
614 ties in northeastern Iowa and southeastern Wisconsin where it was. Both the positioning and risk
615 of very extreme precipitation improved for the Day 2 forecast issuance (Fig. 16b), but probabilities
616 still remained too far to the north. The GEFS/R at Day 3 (Fig. 16e) indicated very little risk of
617 extreme precipitation in the area, with just one member correctly predicting a 1-year exceedance
618 in southeastern Minnesota. The risk of an event occurring within the domain increased for the Day
619 2 issuance, but the locations got worse, with maximum risk indicated in eastern Nebraska, western
620 Iowa, and northeastern Wisconsin, with the only 10-year prediction occurring in the latter location.
621 Somewhat like the raw GEFS/R, the ML model had only some indication of extreme precipitation
622 risk at Day 3 (Fig. 16g). However, it both has the higher probabilities (near 10% in both cases)
623 distributed over a much larger area, and indicates some risk of a 10-year event, with probabili-
624 ty maxima near 1.5%. Additionally, it has the maximum probability axis nearly collocated with
625 where heaviest precipitation occurred, well to the south of the ECMWF probabilities, albeit still

626 slightly too far to the north. The Day 2 forecast issuance (Fig. 16f) is largely similar. The two
627 main changes are a correctly increased risk in the area where the event actually verified, and an
628 incorrectly increased risk of heavy precipitation in eastern Nebraska where the raw GEFS/R had
629 heavy precipitation on Day 2 (Fig. 16d).

630 **5. Discussion and Conclusions**

631 An ML model based on RFs and LR is used to generate CONUS-wide probabilistic forecasts for
632 the exceedance of 1- and 10-year ARI thresholds for 24-hour precipitation accumulations during
633 the Day 2 and Day 3 periods. Approximately eleven years of GEFS/R forecasts, in particular the
634 ensemble median, are used to train these models, and forecasts are using simulated atmospheric
635 fields and varying in both space and time (Table 1), in addition to a variety of geographic and cli-
636 matological forecast predictors (Table 2). Separate models are trained for each of the two 24-hour
637 periods and for each of eight different regions of CONUS, as depicted in Figure 3. A variety of
638 sensitivity experiments are performed, as outlined in Table 3, to ascertain the utility of different
639 aspects of forecast information in predicting locally extreme precipitation. Finally, the final fore-
640 cast models were evaluated, and compared with forecasts based only on the ensemble of raw QPFs
641 from the GEFS/R and ECMWF. The ML models trained in this study demonstrably outperformed
642 the raw GEFS/R forecasts for all regions and forecast lead times (Fig. 6), often more than doubling
643 the forecast skill and adding substantially more than 24-hours lead time improvement in forecast
644 skill. With the exception of the PCST and ROCK regions, the same held for comparison of the
645 ML model forecasts with ECMWF ensemble forecasts as well. Both raw ensembles tended to
646 be negatively biased and overconfident in predicting extreme QPFs (Fig. 7,8), particularly at the
647 10-year ARI for central CONUS regions; this was reversed in the final ML model forecasts, which
648 were more reliable at higher probabilities, but generally underconfident (Fig. 12).

649 In general, unlike past studies (e.g. Herman and Schumacher 2016b), in most regions, the tem-
650 poral resolution and extent of spatially displaced predictors from the forecast point considered had
651 little to no impact on forecast skill (Fig. 4), in addition to the use of upper-level information and
652 additional ensemble information (Fig. 5). These results are suggestive of two findings. First, most
653 of the relevant information about predictors displaced spatiotemporally from the forecast point,
654 other atmospheric fields, or other ensemble member information, can be derived with at least mod-
655 erate accuracy using just the information from the ensemble median from a group of core set of
656 fields collocated and concurrent with the forecast; that is, these additional predictors contain only
657 limited independent forecast information, at least for this coarse dynamical model and this under-
658 dispersive ensemble configuration. It also suggests that, for the most part, the predictive ability is
659 coming primarily through a characterization of the overall environment, which can be reasonably
660 summarized with only a subset of predictors, rather than the simulated spatiotemporal variability
661 and full 3-D characterization of the atmospheric evolution in the underlying dynamical model.
662 This finding comes in contrast to similar studies of other forecast problems using the GEFS/R,
663 such as the Herman and Schumacher (2016b) study which investigated using the GEFS/R to cre-
664 ate ML-based probabilistic forecasts of cloud ceiling and visibility at different airports and found
665 considerable value in the inclusion of spatially displaced predictors. However, there is at least one
666 major exception; none of this really held for the PCST region; here, more complex models with
667 more predictors did notably improve forecast skill. This is perhaps in part because the physical
668 processes associated with extreme precipitation are much better resolved in the GEFS/R in this
669 region compared with the others, and so the added information adds usable forecast utility beyond
670 simply duplicatively characterizing the atmospheric environment for the forecast. The largest skill
671 difference of the sensitivity experiments came for most regions in changing algorithmic assump-

672 tions and processes (Fig. 5d); the simpler linear assumptions of LR tended to degrade forecast
673 skill compared with the more limited assumptions underlying the RF models.

674 The results of this study reveal that the application of more sophisticated statistical methods
675 and ML algorithms such as RFs can demonstrably improve forecasts of extreme precipitation and
676 potentially other rare, high-impact weather events in the medium range when compared with the
677 methods and techniques that are most prevalent in forecast operations today. One unique aspect
678 here is the scope of this model; while most past studies which employed these techniques for
679 numerical weather prediction have focused on a small domain, or just a sampling of points, the
680 models trained here demonstrate an ability to generate skillful, reliable forecasts year-round for
681 all of CONUS and a range of lead times. There are many forecast problems that remain to be
682 explored, but the results of this study and others strongly suggest that further development and
683 application of these data-intensive statistical techniques could substantially improve our forecasts
684 over the current state of the art, even compared with using more sophisticated dynamical models.
685 To that end, implementation of this methodology for operational use to assist Weather Prediction
686 Center forecasters with the development of their excessive rainfall outlooks is currently underway.
687 This forecast technique presents some advantages over purely dynamical approaches, as dynamical
688 models are inherently limited by two factors by which these statistical techniques are not. First, dy-
689 namical models require ever increasing computational resources for increasing model resolution;
690 constraints on computing power prevent sufficient resolution to directly resolve many small-scale
691 processes, many of which are observed in the highest impact weather phenomena. Second, dy-
692 namical models are limited by our physical understanding of the processes we are attempting to
693 simulate or forecast. Machine learning algorithms, in contrast, can detect predictive patterns in the
694 available information even in places where we do not know or understand the physical connection
695 between the information and the phenomenon which we wish to predict. While they are also lim-

696 ited in complexity by computational and data resources, the strict limits on resolvability are not
697 there: physical resolution can often be gained through post-processing of larger scale information.
698 There is thus ample reason to believe that further investigation of these techniques for NWP is a
699 worthwhile venture, and eventual implementation into forecast operations could help forecasters
700 with their tasks by skillfully synthesizing many different sources of forecast information to help
701 alleviate their often time-pressed schedules. This in turn can aid end-user preparedness and, in the
702 case of high-impact events, hopefully help to protect lives and property.

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

710 Some limitations of this work are worthy of note. Stage IV precipitation is used as truth for this
711 study; though there is not a clearly better verification source available, it does have its drawbacks.
712 It does have some spurious quality control issues, and often struggles in areas of complex terrain
713 due to radar beam blockage, interference, and limited gauge coverage (Herman and Schumacher
714 2016a; Nelson et al. 2016). Since the model is trained to forecast Stage IV QPE exceedances, this
715 can lead to some idiosyncrasies and other anomalies associated with the nuances and persistent
716 characteristics of the Stage IV product. One such anomaly is the persistent presence of very small
717 areas of exceedances in some regions of complex terrain during times of favorable convective
718 conditions. This can be removed by quality control procedures to an extent, but some artifacts do
719 remain. This happens most prominently in the terrain of western New Mexico, and a small region

720 has many more instances of ARI exceedances over both the training and test periods than any of
721 the corresponding areas. The ML-based models recognize this, and for the SW region consistently
722 issue much higher probabilities right around this region than in the immediate surroundings. In
723 one sense, this is correct—it is correctly predicting what it was trained to predict—but is still
724 undesirable behavior due to a disparity between “truth” in the study and the true extreme rainfall
725 risk. Solutions to this issue and related issues in other parts of the country must be explored in the
726 future so that the method realizes more operational utility. Additionally, while the choice of using
727 the ARI framework was intentional decision and provides numerous benefits, it is not an end-all
728 for predicting heavy precipitation impacts. While ARIs have *better* correspondence with impacts
729 than a fixed threshold, there are still regional discrepancies in which ARIs have optimal association
730 with impacts, and the framework employed here does not account for antecedent conditions, which
731 can be critical for assessing flash flood risk.

732 Additionally, the predictors for this study come from a very coarse and otherwise rather an-
733 tiquated global model. The GEFS/R was used for this study because, unlike almost any other
734 dynamical model, it has been nearly static for a very long period of record and has nearly station-
735 ary bias characteristics—an essential property for performing this kind of analysis. However, the
736 models trained herein are not working off of the ‘state of the art’ of flash flood predictors. The
737 longer range Day 2 and Day 3 lead times were chosen for this study in part because the discrep-
738 ancancy between GEFS/R forecast quality and ‘state of the art’ is smaller at these longer lead times
739 due to less convection-allowing guidance being available, and higher-resolution models degrading
740 in utility with increasing forecast lead time.

741 There are also some complications that must be considered for real-time implementation. As one
742 example, the regional models are trained completely independently of one another, with different
743 training data and different solutions. Consequently, they can occasionally give rather different

744 predictions on nearly identical inputs, resulting in undesirable probability discontinuities across
745 region boundaries. Appropriate methods for removing probability discontinuities in space must be
746 further explored.

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

762 This study also focused on a rather specific time interval and took all dynamical predictors
763 from a single, somewhat antiquated ensemble system. Future expansion both to the 12–36 hour
764 Day 1 period and beyond the Day 3 period will be explored, including predictors from more
765 contemporary CAM models and potentially including observations as well for the shorter lead
766 time forecasts. Operational models also tend to undergo periodic upgrades and thus do not remain
767 static like the ensemble system used here. The sensitivity of ML model performance to changes

768 in dynamical model bias characteristics that result from these upgrades is a question of considerable
769 operational relevance and an additional factor worthy of future investigation. It was also seen
770 that the ML models suffered to varying degrees from underconfidence and, in some instances,
771 negative bias. Methods of probability calibration of the ML model probabilities as a final post-
772 processing step should be explored in future work, and parameter choices reconsidered in light
773 of this additional calibration. Finally, this study only explored a subset of available machine
774 learning algorithms. Other choices, including adaptive learning algorithms, may be able to better
775 exploit predictor-predictand relationships, appropriately update to reflect changes in an underlying
776 dynamical model, and produce superior forecasts for the locally extreme precipitation and flash
777 flood forecast problem.

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

786 APPENDIX A

787 **Random Forests and Their Parameters**

788 As noted in the main text, RFs are simply an ensemble of decision trees. Decision trees consist
789 of a network of two types of nodes: decision nodes and leaf nodes. Decision nodes each have

790 exactly two children, which may be either decision nodes or leaf nodes, with a binary split based
791 on the numeric value of a single input predictor determining whether to traverse to the left or right
792 child. A leaf node has no children and instead, makes a categorical prediction of the outcome of
793 the input example based on the leaf's relationship to its ancestor nodes. For a given forecast, one
794 begins at a decision tree's root, traversing through its children based on the relative value of the
795 forecast's predictors to each decision node's threshold critical value for the predictor associated
796 with the node. This process is repeated until a leaf node is reached; its value corresponding to the
797 leaf becomes the tree's deterministic prediction.

798 Decision trees can be a powerful approach for a wide array of applications, but they also have
799 several significant drawbacks. In particular, they are very prone to overfitting (e.g. Brodley and
800 Utgoff 1995), fitting to the noise of the training data rather than just the underlying relationships.
801 They also don't convey any information about forecast uncertainty, as would be the case in a prob-
802 abilistic framework. RFs are used instead to alleviate these concerns by producing a probabilistic
803 forecast in a way that can significantly decrease error from overfitting the supplied training error
804 with only a slight increase to error from oversimplistic model assumptions, provided the trees are
805 sufficiently uncorrelated. The difficulty then revolves around generating a large set (forest) of
806 skillful decision trees that are not strongly correlated. The decision tree generating procedure de-
807 scribed above is deterministic: a given set of training data will always produce the same decision
808 tree. A forest of identical decision trees, of course, adds no value over using a single decision tree.
809 Two additional processes—tree bagging and feature bagging—are employed to produce unique
810 trees. Tree bagging produces unique trees through a straightforward bootstrapping procedure.
811 Specifically, a forest of size B is formed from the n training examples by creating B samples of
812 size n , with replacement, from the original training data, and running the decision tree algorithm
813 on each sample. Overfitting due to correlated trees can still occur under this approach, particularly

814 if a small subset of the original features are much more robust predictors of the verifying category
815 than the rest (Breiman 2001; Murphy 2012). To overcome this problem, feature bagging is also
816 employed, whereby only a random subset of the m original input predictors are considered at each
817 decision node; the size of the random subset is denoted here as Z ; $1 \leq Z \leq m$. This combination can
818 result in a set of B largely uncorrelated trees, each of which is individually fairly skillful.

819 With any machine learning algorithm, there are numerous considerations in the actual model
820 construction, which manifest themselves in tunable parameters. Compared with other machine
821 learning algorithms, such as gradient boosting or support vector machines, RFs are often praised
822 for their relative insensitivity to their parameters with respect to model performance, but it is
823 nevertheless important to explore the parameter space in order to realize the full utility of the
824 algorithm. The forest size B is perhaps the most obvious parameter. The general relationship
825 between model performance and B is well known and consistent across all prediction problems;
826 it starts quite low at very low B , initially increases rapidly with increasing B , and then slowly
827 asymptotes to some threshold performance limit as the relationships between input features have
828 been fully explored by the forest and the inclusion of new trees becomes redundant. Larger forest
829 sizes require more computational expense, so the goal is to select B such that it is small enough to
830 be computationally tractable but large enough to be near the performance limit. Another parameter
831 noted above is Z , the number of features to consider at each node split. If this number is too small,
832 model performance may suffer from only considering irrelevant or otherwise unproductive features
833 in the context of the node; if Z is too large, performance will also suffer because of underdispersive
834 trees producing an overfit forest solution. Another frequently explored parameter is the splitting
835 criterion evaluation function. Most commonly used are either the Gini impurity or the information
836 gain; past studies have shown that this choice is not important for many forecast problems. In this
837 study, information gain will be used, and can be expressed for a training set T , candidate splitting

838 feature x_a and candidate split value v_a as:

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

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

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

841 The chosen splitting feature and split value are selected among those considered which maximize
842 Equation A1 (e.g. Quinlan 1986; Murphy 2012). However, there are two other parameters that
843 have the most substantial influence on model performance. The first, denoted S , is the minimum
844 number of training examples required to split a node. Traditionally, RFs create a leaf only once a
845 node is ‘pure’, that is, all the remaining training examples associated with that node have the same
846 labels—event outcomes. In this way, each tree makes a categorical prediction of the predictand
847 outcome, and probabilities are generated only in counting the proportion of trees in the forest
848 making a particular forecast. However, this can make predictions from an individual tree very
849 susceptible to the outcome of a particular historical case, and in some cases result in substantial
850 overfitting. Instead, by increasing S , an RF can be allowed to make ‘impure’ leaves; at these nodes,
851 an individual tree makes a probabilistic prediction based on the proportion of remaining training
852 examples exhibiting each event class rather than continuing to split based on the remaining training
853 data. Making S too large, however, can result in underfitting—lumping data as indistinguishable
854 when there are in fact underlying discernible distinctions between remaining training examples
855 with different labels. The last parameter, denoted P , is not actually an RF algorithm parameter at
856 all. When PCA is performed, there is always a question about the number of components to retain.
857 Though there are some heuristics (e.g. North et al. 1982), there is no definitive method to know
858 *a priori* how many retained components P will produce the most skillful forecasts (Wilks 2011).

859 If P is too small, valuable forecast data is discarded and predictive performance consequently
860 suffers. However, if it is too large, the retained PCs eventually become essentially just noise,
861 and the RF, by fitting to these predictor values in the training data, will yield an overfit model
862 that does not generalize to unseen data. Experiments that will not be discussed herein revealed
863 that using information gain to determine splits and letting $B=1000$ produced skill near that of an
864 infinitely large forest, and skill was insensitive to modifications of these settings, including modest
865 increases in the forest size beyond this point. However, the Z-S-P parameter space are explored
866 for the models trained and those results are presented in Appendix C.

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

876 APPENDIX B

877 **Logistic Regression and Other Algorithms**

878 Other machine learning algorithms do not extrapolate well to the high dimensionality of the fore-
879 cast problem explored here. While time to train a model is not of primary concern for operational
880 forecasting, since it is performed only once (or periodically) offline, there are nevertheless some
881 practical considerations; models that take months or longer to train would be unlikely to be re-

882 alistic choices, for example. The “online” forecasting component, that is, the time required to
883 take a new forecast, input it into a trained model and receive a forecast, is of operational concern,
884 but all of the forecast techniques considered here can produce forecasts in a matter of minutes,
885 and the small differences are not considered to be of practical concern. Using the random for-
886 est classification heuristic of considering the square root of the total number of features at each
887 node split (Geurts et al. 2006), the computational complexity of training a random forest of size
888 B from N training examples with F features ($N > F$) may be expressed as $O(B\sqrt{FN}\log(N))$, and
889 may be readily parallelized across trees or within trees. Some algorithms are quadratic or, in the
890 case of support vector machines (Cortes and Vapnik 1995), even cubic in the number of training
891 examples, and do not parallelize as readily. Others, such as logistic regression (LR), are linear
892 in the number of training examples, but require matrix multiplication yielding a computational
893 complexity $O(NF^2)$. PCA pre-processing acts to make learning algorithms more computationally
894 manageable; consequently, even otherwise potentially intractable approaches, such as LR with the
895 $\sim 10,000$ native features used here, become feasible after applying the PCA pre-processing step.
896 It furthermore serves the goal of alleviating the so-called “curse of dimensionality” and reducing
897 overfitting.

898 One sensitivity experiment compares model performance as a function of the model algorithm
899 by comparing skill of forecasts produced by RFs with those produced with LR. LR is in many
900 senses a simpler model than an RF, since the structural form of the relationship between the pre-
901 dictors and the predictand is predefined before training. RFs, in contrast, make few assumptions
902 about the relationships between the predictors and the predictand, allowing more diverse diagnoses
903 of underlying relationships. However, this lack of assumptions can result in some degree of over-
904 fitting which fundamentally cannot be mitigated through parameter tuning. As an application of
905 the generalized linear model, LR assumes a linear predictors-predictand relationship via the logit

906 function and is thus in one aspect considerably more limited than an RF. In LR, a single regression
 907 equation, or K equations for a multiclass problem with K categories, is computed to represent
 908 the probability of the outcome being category k given the set of input predictors \mathbf{x} . In particular,
 909 verifying probabilities are computed using the softmax function:

$$P(y = k|\mathbf{x}) = \frac{e^{\mathbf{x}^T \mathbf{w}_k}}{\sum_{j=1}^K e^{\mathbf{x}^T \mathbf{w}_j}} \quad (\text{B1})$$

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

$$J(\mathbf{Y}, \mathbf{P}(\mathbf{w})) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \frac{1}{CN} \sum_{i=1}^N \sum_{k=1}^K \mathbf{Y}_{i,k} \log(\mathbf{P}_{i,k}) \quad (\text{B2})$$

922 where C represents the extent of regularization, with smaller values indicating that large weights
 923 are penalized more than with larger values of C . Alternative approaches to regularization exist
 924 (e.g. Murphy 2012), and are explored to some degree in sensitivity experiments of Appendix C.

APPENDIX C

Results: Parameter Tuning

926

927 RF model parameters were tuned for each region and lead time separately through the 4-fold
928 cross-validation procedure employed throughout the study. Overall, the optimal parameters were
929 found not to vary with the two different lead times, but did vary for two of the parameters as
930 a function of forecast region, at least to an extent; the full results appear in Table C1. For the
931 S parameter—the number of predictors considered for each node split, the default heuristic of
932 the square root of the total number of features was found to maximize RPSS for all regions and
933 lead times. In all instances where both were tested, unbalanced sampling from the event classes
934 in proportion to their true observed frequencies outperformed balanced equal sampling from each
935 event class, in contrast to Ahijevych et al. (2016) and others; the finding appeared to be attributable
936 to biased probabilities produced from the balanced sampling technique. For the Z parameter, the
937 minimum number of remaining training examples in an impure parameter subspace required to
938 perform a further node split, was generally found to be around 120. Lesser values maximized skill
939 in the western regions, with values of 30 maximizing skill in the SW and ROCK regions, and Z=4
940 producing the best skill over PCST. A couple of the larger regions of the east, SE and MDWST,
941 maximized RPSS with a value of 240, although the sensitivity between Z=120 and Z=240 was
942 small for all regions. For P in the CTL_PCA models, skill was generally maximized with P=30,
943 that is, retaining the 30 PCs which explain the most variance of the entire GEFS/R predictor set.
944 For most regions, there was very limited sensitivity in the P=30–40 interval—although there was
945 larger sensitivity outside this interval—and P=40 was found to produce slightly better skill in the
946 NGP region. The PCST region was again the main exception, where P=60 was found to maximize
947 cross-validation RPSS.

948 LR model parameters were tuned using an identical framework to ascertain the type of regular-
949 ization, either based on a L1 norm which penalizes non-zero weights, or L2 norm—described in
950 Appendix B—which penalizes large magnitude weights. L2 regularization was consistently found
951 to produce superior results, perhaps because the number of retained PCs was already taken from
952 the P parameter in the RF experiments, acting to nullify many potential non-zero weights of higher
953 numbered PCs. Unlike the RF experiments, there were occasionally some large differences in the
954 obtained optimal regularization parameter value C between lead times within the same region.
955 Generally, models performed better with more regularized solutions, but there were some notable
956 exceptions, with the Day 2 NGP model and Day 3 NE model obtaining optimal C parameter values
957 on the other end of the spectrum.

958 **References**

- 959 Ahijevych, D., J. O. Pinto, J. K. Williams, and M. Steiner, 2016: Probabilistic forecasts of
960 mesoscale convective system initiation using the random forest data mining technique. *Wea.*
961 *Forecasting*, **31** (2), 581–599.
- 962 Alvarez, F. M., 2014: Statistical calibration of extended-range probabilistic tornado forecasts with
963 a reforecast dataset. Ph.D. thesis, SAINT LOUIS UNIVERSITY.
- 964 Antolik, M. S., 2000: An overview of the National Weather Service’s centralized statistical quan-
965 titative precipitation forecasts. *J. Hydrol.*, **239** (1), 306–337.
- 966 Applequist, S., G. E. Gahrs, R. L. Pfeffer, and X.-F. Niu, 2002: Comparison of methodologies for
967 probabilistic quantitative precipitation forecasting*. *Wea. Forecasting*, **17** (4), 783–799.
- 968 Ashley, S. T., and W. S. Ashley, 2008: Flood fatalities in the United States. *J. Appl. Meteor.*
969 *Climatol.*, **47** (3), 805–818.

970 Baars, J. A., and C. F. Mass, 2005: Performance of National Weather Service forecasts compared
971 to operational, consensus, and weighted model output statistics. *Wea. Forecasting*, **20** (6), 1034–
972 1047.

973 Barthold, F. E., T. E. Workoff, B. A. Cosgrove, J. J. Gourley, D. R. Novak, and K. M. Mahoney,
974 2015: improving flash flood forecasts: The HMT-WPC flash flood and intense rainfall experi-
975 ment. *Bull. Amer. Meteor. Soc.*, **96** (11), 1859–1866.

976 Bermowitz, R. J., 1975: An application of model output statistics to forecasting quantitative pre-
977 cipitation. *Mon. Wea. Rev.*, **103** (2), 149–153.

978 Bonnin, G. M., D. Martin, B. Lin, T. Parzybok, M. Yekta, and D. Riley, 2006: Precipitation-
979 frequency atlas of the United States. *NOAA Atlas 14*, **2**.

980 Bonnin, G. M., D. Todd, B. Lin, T. Parzybok, M. Yekta, and D. Riley, 2004: Precipitation-
981 frequency atlas of the United States. *NOAA Atlas 14*, **1**.

982 Bougeault, P., and Coauthors, 2010: The thorpex interactive grand global ensemble. *Bull. Amer.*
983 *Meteor. Soc.*, **91** (8), 1059–1072.

984 Breiman, L., 2001: Random forests. *Mach. Learn.*, **45** (1), 5–32, doi:10.1023/A:1010933404324.

985 Brimelow, J. C., G. W. Reuter, R. Goodson, and T. W. Krauss, 2006: Spatial forecasts of maximum
986 hail size using prognostic model soundings and HAILCAST. *Wea. Forecasting*, **21**, 206–219,
987 doi:10.1175/WAF915.1.

988 Brodley, C. E., and P. E. Utgoff, 1995: Multivariate decision trees. *Machine learning*, **19** (1),
989 45–77.

990 Brooks, H. E., and D. J. Stensrud, 2000: Climatology of heavy rain events in the United States
991 from hourly precipitation observations. *Mon. Wea. Rev.*, **128** (4), 1194–1201.

- 992 Buizza, R., A. Hollingsworth, F. Lalaurette, and A. Ghelli, 1999: Probabilistic predictions of
993 precipitation using the ecmwf ensemble prediction system. *Wea. Forecasting*, **14** (2), 168–189.
- 994 Calianno, M., I. Ruin, and J. J. Gourley, 2013: Supplementing flash flood reports with impact
995 classifications. *J. Hydrol.*, **477**, 1–16.
- 996 Clark, A. J., W. A. Gallus Jr, and T.-C. Chen, 2007: Comparison of the diurnal precipitation
997 cycle in convection-resolving and non-convection-resolving mesoscale models. *Mon. Wea. Rev.*,
998 **135** (10), 3456–3473.
- 999 Clark, A. J., W. A. Gallus Jr, and M. L. Weisman, 2010: Neighborhood-based verification of
1000 precipitation forecasts from convection-allowing NCAR WRF model simulations and the oper-
1001 ational NAM. *Wea. Forecasting*, **25** (5), 1495–1509.
- 1002 Cortes, C., and V. Vapnik, 1995: Support-vector networks. *Mach. Learn.*, **20** (3), 273–297, doi:
1003 10.1007/BF00994018.
- 1004 Davis, C. A., K. W. Manning, R. E. Carbone, S. B. Trier, and J. D. Tuttle, 2003: Coherence
1005 of warm-season continental rainfall in numerical weather prediction models. *Mon. Wea. Rev.*,
1006 **131** (11), 2667–2679.
- 1007 Delrieu, G., and Coauthors, 2005: The catastrophic flash-flood event of 8–9 September 2002 in
1008 the Gard Region, France: a first case study for the cévennes–vivaraïis mediterranean hydrome-
1009 teorological observatory. *J. Hydrometeor.*, **6** (1), 34–52.
- 1010 Duda, J. D., and W. A. Gallus, 2013: The impact of large-scale forcing on skill of simulated
1011 convective initiation and upscale evolution with convection-allowing grid spacings in the WRF*.
1012 *Wea. Forecasting*, **28** (4), 994–1018.

- 1013 Eckel, F. A., 2003: Effective mesoscale, short-range ensemble forecasting. Ph.D. thesis, University
1014 of Washington.
- 1015 Friedman, J. H., 1997: On bias, variance, 0/1loss, and the curse-of-dimensionality. *Data mining
1016 and knowledge discovery*, **1** (1), 55–77.
- 1017 Fritsch, J. M., and R. Carbone, 2004: Improving quantitative precipitation forecasts in the warm
1018 season: A USWRP research and development strategy. *Bull. Amer. Meteor. Soc.*, **85** (7), 955–
1019 965, doi:10.1175/BAMS-85-7-955.
- 1020 Gagne, D. J., A. McGovern, S. E. Haupt, R. A. Sobash, J. K. Williams, and M. Xue, 2017: Storm-
1021 based probabilistic hail forecasting with machine learning applied to convection-allowing en-
1022 sembles. *Wea. Forecasting*, **32**, 1819–1840.
- 1023 Gagne, D. J., A. McGovern, and M. Xue, 2014: Machine learning enhancement of storm-scale en-
1024 semble probabilistic quantitative precipitation forecasts. *Wea. Forecasting*, **29** (4), 1024–1043.
- 1025 Gagne, D. J., II, A. McGovern, J. Brotzge, M. Coniglio, J. Correia Jr, and M. Xue, 2015: Day-
1026 ahead hail prediction integrating machine learning with storm-scale numerical weather models.
1027 *AAAI*, 3954–3960.
- 1028 Geurts, P., D. Ernst, and L. Wehenkel, 2006: Extremely randomized trees. *Machine learning*,
1029 **63** (1), 3–42.
- 1030 Glahn, H. R., and D. A. Lowry, 1972: The use of model output statistics (MOS) in objec-
1031 tive weather forecasting. *J. Appl. Meteor.*, **11** (8), 1203–1211, doi:10.1175/1520-0450(1972)
1032 011(1203:TUOMOS)2.0.CO;2.
- 1033 Gochis, D., and Coauthors, 2015: The great Colorado flood of September 2013. *Bull. Amer. Me-
1034 teor. Soc.*, **96** (9), 1461–1487.

- 1035 Gourley, J. J., J. M. Erlingis, Y. Hong, and E. B. Wells, 2012: Evaluation of tools used for moni-
1036 toring and forecasting flash floods in the United States. *Wea. Forecasting*, **27** (1), 158–173.
- 1037 Graham, R. A., and R. H. Grumm, 2010: Utilizing normalized anomalies to assess synoptic-scale
1038 weather events in the western united states. *Wea. Forecasting*, **25** (2), 428–445.
- 1039 Grams, J. S., W. A. Gallus Jr, S. E. Koch, L. S. Wharton, A. Loughe, and E. E. Ebert, 2006: The
1040 use of a modified Ebert-McBride technique to evaluate mesoscale model QPF as a function of
1041 convective system morphology during IHOP 2002. *Wea. Forecasting*, **21** (3), 288–306.
- 1042 Hall, T., H. E. Brooks, and C. A. Doswell III, 1999: Precipitation forecasting using a neural
1043 network. *Wea. Forecasting*, **14** (3), 338–345.
- 1044 Hamill, T. M., 2017: Changes in the systematic errors of global reforecasts due to an evolving data
1045 assimilation system. *Mon. Wea. Rev.*, (2017).
- 1046 Hamill, T. M., G. T. Bates, J. S. Whitaker, D. R. Murray, M. Fiorino, T. J. Galarneau Jr, Y. Zhu,
1047 and W. Lapenta, 2013: NOAA’s second-generation global medium-range ensemble reforecast
1048 dataset. *Bull. Amer. Meteor. Soc.*, **94** (10), 1553–1565, doi:10.1175/BAMS-D-12-00014.1.
- 1049 Hamill, T. M., and J. Juras, 2006: Measuring forecast skill: Is it real skill or is it the varying
1050 climatology? *Quart. J. Roy. Meteor. Soc.*, **132**, 2905–2924.
- 1051 Hamill, T. M., M. Scheuerer, and G. T. Bates, 2015: Analog probabilistic precipitation fore-
1052 casts using GEFS reforecasts and climatology-calibrated precipitation analyses. *Mon. Wea. Rev.*,
1053 **143** (8), 3300–3309.
- 1054 Hamill, T. M., and J. S. Whitaker, 2006: Probabilistic quantitative precipitation forecasts based on
1055 reforecast analogs: Theory and application. *Mon. Wea. Rev.*, **134** (11), 3209–3229.

- 1056 Hamill, T. M., J. S. Whitaker, and X. Wei, 2004: Ensemble reforecasting: Improving medium-
1057 range forecast skill using retrospective forecasts. *Mon. Wea. Rev.*, **132** (6), 1434–1447.
- 1058 Hapuarachchi, H., Q. Wang, and T. Pagano, 2011: A review of advances in flash flood forecasting.
1059 *Hydrological Processes*, **25** (18), 2771–2784.
- 1060 Herman, G. R., and R. S. Schumacher, 2016a: Extreme precipitation in models: An evaluation.
1061 *Wea. Forecasting*, **31**, 1853–1879, doi:10.1175/WAF-D-16-0093.1.
- 1062 Herman, G. R., and R. S. Schumacher, 2016b: Using reforecasts to improve forecasting of fog and
1063 visibility for aviation. *Wea. Forecasting*, **31** (2), 467–482.
- 1064 Hershfield, D. M., 1961: Technical paper no. 40: Rainfall frequency atlas of the United States.
1065 *Weather Bureau, Department of Commerce, Washington, DC.*
- 1066 Hitchens, N. M., H. E. Brooks, and R. S. Schumacher, 2013: Spatial and temporal characteristics
1067 of heavy hourly rainfall in the United States. *Mon. Wea. Rev.*, **141** (12), 4564–4575.
- 1068 Hong, T., P. Pinson, S. Fan, H. Zareipour, A. Troccoli, and R. J. Hyndman, 2016: Probabilis-
1069 tic energy forecasting: Global energy forecasting competition 2014 and beyond. *International*
1070 *Journal of Forecasting*, **32** (3), 896–913.
- 1071 Hou, D., and Coauthors, 2014: Climatology-calibrated precipitation analysis at fine scales: Sta-
1072 tistical adjustment of stage iv toward cpc gauge-based analysis. *J. Hydrometeor.*, **15** (6), 2542–
1073 2557.
- 1074 Jacks, E., J. B. Bower, V. J. Dagostaro, J. P. Dallavalle, M. C. Erickson, and J. C. Su, 1990: New
1075 NGM-based MOS guidance for maximum/minimum temperature, probability of precipitation,
1076 cloud amount, and surface wind. *Wea. Forecasting*, **5** (1), 128–138.

- 1077 Junker, N. W., M. J. Brennan, F. Pereira, M. J. Bodner, and R. H. Grumm, 2009: Assessing the
1078 potential for rare precipitation events with standardized anomalies and ensemble guidance at the
1079 hydrometeorological prediction center. *Bull. Amer. Meteor. Soc.*, **90** (4), 445–453.
- 1080 Kain, J. S., S. J. Weiss, J. J. Levit, M. E. Baldwin, and D. R. Bright, 2006: Examination of
1081 convection-allowing configurations of the WRF model for the prediction of severe convective
1082 weather: The SPC/NSSL Spring Program 2004. *Wea. Forecasting*, **21**, 167–181, doi:10.1175/
1083 WAF906.1.
- 1084 Kain, J. S., and Coauthors, 2008: Some practical considerations regarding horizontal resolution
1085 in the first generation of operational convection-allowing NWP. *Wea. Forecasting*, **23** (5), 931–
1086 952.
- 1087 Klein, W. H., B. M. Lewis, and I. Enger, 1959: Objective prediction of five-day mean temperatures
1088 during winter. *J. Meteor.*, **16** (6), 672–682.
- 1089 Lackmann, G. M., 2013: The south-central US flood of May 2010: Present and future. **26** (13),
1090 4688–4709.
- 1091 Lean, H. W., P. A. Clark, M. Dixon, N. M. Roberts, A. Fitch, R. Forbes, and C. Halliwell, 2008:
1092 Characteristics of high-resolution versions of the Met Office Unified Model for forecasting con-
1093 vection over the United Kingdom. *Mon. Wea. Rev.*, **136** (9), 3408–3424.
- 1094 Lin, Y., and K. E. Mitchell, 2005: 1.2 the NCEP Stage II/IV hourly precipitation analyses: Devel-
1095 opment and applications. *19th Conf. Hydrology, American Meteorological Society, San Diego,*
1096 *CA, USA, Citeseer.*
- 1097 Marjerison, R. D., M. T. Walter, P. J. Sullivan, and S. J. Colucci, 2016: Does population affect the
1098 location of flash flood reports? *J. Appl. Meteor. Climatol.*, **55** (9), 1953–1963.

1099 Miller, J., R. Frederick, and R. Tracey, 1973: NOAA Atlas 2. *Precipitation-frequency atlas of the*
1100 *western United States*, **3**.

1101 Molteni, F., R. Buizza, T. N. Palmer, and T. Petroliaigis, 1996: The ecmwf ensemble prediction
1102 system: Methodology and validation. *Quart. J. Roy. Meteor. Soc.*, **122 (529)**, 73–119.

1103 Murphy, A. H., 1973: A new vector partition of the probability score. *J. Appl. Meteor.*, **12 (4)**,
1104 595–600.

1105 Murphy, K. P., 2012: *Machine learning: a probabilistic perspective*. MIT press.

1106 Nelson, B. R., O. P. Prat, D.-J. Seo, and E. Habib, 2016: Assessment and implications of NCEP
1107 Stage IV quantitative precipitation estimates for product intercomparisons. *Wea. Forecasting*,
1108 **31 (2)**, 371–394.

1109 Nielsen, E. R., G. R. Herman, R. C. Tournay, J. M. Peters, and R. S. Schumacher, 2015: Double
1110 impact: When both tornadoes and flash floods threaten the same place at the same time. *Wea.*
1111 *Forecasting*, **30**, 1673–1693, doi:10.1175/WAF-D-15-0084.1.

1112 Nielsen, E. R., and R. S. Schumacher, 2016: Using convection-allowing ensembles to understand
1113 the predictability of an extreme rainfall event. *Mon. Wea. Rev.*, **144**, 3651–3676.

1114 North, G. R., T. L. Bell, R. F. Cahalan, and F. J. Moeng, 1982: Sampling errors in the estimation
1115 of empirical orthogonal functions. *Mon. Wea. Rev.*, **110 (7)**, 699–706.

1116 Novak, D. R., C. Bailey, K. F. Brill, P. Burke, W. A. Hogsett, R. Rausch, and M. Schichtel,
1117 2014: Precipitation and temperature forecast performance at the Weather Prediction Center.
1118 *Wea. Forecasting*, **29 (3)**, 489–504.

1119 Ntelekos, A. A., K. P. Georgakakos, and W. F. Krajewski, 2006: On the uncertainties of flash flood
1120 guidance: Toward probabilistic forecasting of flash floods. *J. Hydrometeor.*, **7 (5)**, 896–915.

1121 NWS, 2017a: Service change notice 17-100. National Centers for Environmental Predic-
1122 tion, Weather Prediction Center, [Available online at [http://www.nws.noaa.gov/os/notification/
1123 scn17-100wpc_excessive_rainfall.htm](http://www.nws.noaa.gov/os/notification/scn17-100wpc_excessive_rainfall.htm)].

1124 NWS, 2017b: Summary of natural hazard statistics in the United States. National Weather Service,
1125 Office of Climate, Weather, and Water Services, [Available online at [http://www.nws.noaa.gov/
1126 om/hazstats.shtml](http://www.nws.noaa.gov/om/hazstats.shtml)].

1127 Pedregosa, F., and Coauthors, 2011: Scikit-learn: Machine learning in Python. *J. Mach. Learn.*
1128 *Res*, **12**, 2825–2830.

1129 Perica, S., and Coauthors, 2011: Precipitation-frequency atlas of the United States. *NOAA Atlas*
1130 *14*, **6**.

1131 Perica, S., and Coauthors, 2013: Precipitation-frequency atlas of the United States. *NOAA Atlas*
1132 *14*, **9**.

1133 Pielke, R. A., M. W. Downton, and J. B. Miller, 2002: *Flood damage in the United States, 1926-*
1134 *2000: a reanalysis of National Weather Service estimates*. University Corporation for Atmo-
1135 spheric Research Boulder, CO.

1136 Pinto, J. O., J. A. Grim, and M. Steiner, 2015: Assessment of the High-Resolution Rapid Refresh
1137 models ability to predict mesoscale convective systems using object-based evaluation. *Wea.*
1138 *Forecasting*, **30** (4), 892–913.

1139 Quinlan, J. R., 1986: Induction of decision trees. *Machine learning*, **1** (1), 81–106.

1140 Ramshaw, J. D., 1985: Conservative rezoning algorithm for generalized two-dimensional meshes.
1141 *Journal of Computational Physics*, **59** (2), 193–199.

- 1142 Reed, S., J. Schaake, and Z. Zhang, 2007: A distributed hydrologic model and threshold
1143 frequency-based method for flash flood forecasting at ungauged locations. *Journal of Hydrol-*
1144 *ogy*, **337** (3), 402–420.
- 1145 Ross, D. A., J. Lim, R.-S. Lin, and M.-H. Yang, 2008: Incremental learning for robust visual
1146 tracking. *International Journal of Computer Vision*, **77** (1-3), 125–141.
- 1147 Rutz, J. J., W. J. Steenburgh, and F. M. Ralph, 2014: Climatological characteristics of atmospheric
1148 rivers and their inland penetration over the western United States. *Mon. Wea. Rev.*, **142** (2),
1149 905–921.
- 1150 Scheuerer, M., and T. M. Hamill, 2015: Statistical postprocessing of ensemble precipitation fore-
1151 casts by fitting censored, shifted gamma distributions*. *Mon. Wea. Rev.*, **143** (11), 4578–4596.
- 1152 Schmidt, J. A., A. Anderson, and J. Paul, 2007: Spatially-variable, physically-derived flash flood
1153 guidance. *AMS 21st Conference on Hydrology, San Antonio, TX B*, Vol. 6.
- 1154 Schumacher, R. S., A. J. Clark, M. Xue, and F. Kong, 2013: Factors influencing the develop-
1155 ment and maintenance of nocturnal heavy-rain-producing convective systems in a storm-scale
1156 ensemble. *Mon. Wea. Rev.*, **141** (8), 2778–2801.
- 1157 Schumacher, R. S., and R. H. Johnson, 2006: Characteristics of US extreme rain events during
1158 1999–2003. *Wea. Forecasting*, **21** (1), 69–85.
- 1159 Schumacher, R. S., and R. H. Johnson, 2008: Mesoscale processes contributing to extreme rainfall
1160 in a midlatitude warm-season flash flood. *Mon. Wea. Rev.*, **136** (10), 3964–3986.
- 1161 Shlens, J., 2014: A tutorial on principal component analysis. *arXiv preprint arXiv:1404.1100*.

- 1162 Stevenson, S. N., and R. S. Schumacher, 2014: A 10-year survey of extreme rainfall events in the
1163 central and eastern United States using gridded multisensor precipitation analyses. *Mon. Wea.*
1164 *Rev.*, **142** (9), 3147–3162.
- 1165 Villarini, G., W. F. Krajewski, A. A. Ntelekos, K. P. Georgakakos, and J. A. Smith, 2010: Towards
1166 probabilistic forecasting of flash floods: The combined effects of uncertainty in radar-rainfall
1167 and flash flood guidance. *J. Hydrol.*, **394** (1), 275–284.
- 1168 Vislocky, R. L., and J. M. Fritsch, 1997: Performance of an advanced MOS system in the 1996-97
1169 National Collegiate Weather Forecasting Contest. *Bull. Amer. Meteor. Soc.*, **78** (12), 2851.
- 1170 Wang, S.-Y., T.-C. Chen, and S. E. Taylor, 2009: Evaluations of NAM forecasts on midtropo-
1171 spheric perturbation-induced convective storms over the US northern plains. *Wea. Forecasting*,
1172 **24** (5), 1309–1333.
- 1173 Wei, M., Z. Toth, R. Wobus, and Y. Zhu, 2008: Initial perturbations based on the ensemble trans-
1174 form (ET) technique in the NCEP global operational forecast system. *Tellus A*, **60** (1), 62–79.
- 1175 Weisman, M. L., C. Davis, W. Wang, K. W. Manning, and J. B. Klemp, 2008: Experiences with
1176 0-36-h explicit convective forecasts with the WRF-ARW model. *Wea. Forecasting*, **23** (3), 407–
1177 437.
- 1178 Welles, E., S. Sorooshian, G. Carter, and B. Olsen, 2007: Hydrologic verification: A call for action
1179 and collaboration. *Bull. Amer. Meteor. Soc.*, **88** (4), 503–511.
- 1180 Wilks, D. S., 2011: *Statistical Methods in the Atmospheric Sciences*. Academic press, 676 pp.
- 1181 Wilson, J. W., and R. D. Roberts, 2006: Summary of convective storm initiation and evolution
1182 during IHOP: Observational and modeling perspective. *Mon. Wea. Rev.*, **134** (1), 23–47.

- 1183 Wolter, K., J. K. Eischeid, L. Cheng, and M. Hoerling, 2016: What history tells us about 2015 us
1184 daily rainfall extremes. *Bull. Amer. Meteor. Soc.*, **97** (12), S9–S13.
- 1185 Zeng, J., and W. Qiao, 2011: Support vector machine-based short-term wind power forecasting.
1186 *Power Systems Conference and Exposition (PSCE), 2011 IEEE/PES*, IEEE, 1–8.
- 1187 Zhang, F., N. Bei, R. Rotunno, C. Snyder, and C. C. Epifanio, 2007: Mesoscale predictability
1188 of moist baroclinic waves: Convection-permitting experiments and multistage error growth dy-
1189 namics. *J. Atmos. Sci.*, **64** (10), 3579–3594.
- 1190 Zhang, F., C. Snyder, and R. Rotunno, 2003: Effects of moist convection on mesoscale predictabil-
1191 ity. *J. Atmos. Sci.*, **60**, 1173–1185, doi:10.1175/1520-0469(2003)060<1173:EOMCOM>2.0.CO;
1192 2.

1193 **LIST OF TABLES**

1194 **Table 1.** Summary of dynamical model fields examined in this study, including the ab-
 1195 breviated symbol to which each variable is referred throughout the paper, a
 1196 description of each variable, the predictor group with which the field is associ-
 1197 ated in the manuscript text, and the highest resolution for which the field can
 1198 be obtained from the GEFS/R. 57

1199 **Table 2.** List of background predictors used in this study, and their associated symbols
 1200 and descriptions. 58

1201 **Table 3.** Summary of the models trained in this study, and the corresponding names
 1202 designated to the models. ‘X’ indicates the process is performed or the infor-
 1203 mation is used; a lack of one indicates the opposite. MEDIAN corresponds
 1204 to the ensemble median, CTRL corresponds to the ensemble control member’s
 1205 fields, and CNFDB uses the median in addition to the second-from-lowest and
 1206 second-from-highest member values for each field. Horizontal radius is listed
 1207 in grid boxes from forecast point; timestep denotes the number of hours be-
 1208 tween GEFS/R forecast field predictors. Slashes indicate the first number ap-
 1209 plies to the Day 2 version of the model, while the latter number applies to the
 1210 Day 3 version. Letters enclosed by parentheses indicate sub-versions of mod-
 1211 els, with one parameter changed to the value adjacent to the letter. Asterisks
 1212 indicate a model applies only to Day 2, and not Day 3. Otherwise, models ap-
 1213 ply to all eight forecast regions and have both Day 2 and Day 3 versions. Those
 1214 models with bolded names are incorporated into the weighted blend of the final
 1215 model configuration. 59

1216 **Table C1.** Optimal RF parameters obtained in cross-validation for the Z-S-P parameter
 1217 space. SQRT indicates the square root of the total number of predictors; sym-
 1218 bols are otherwise as described in the manuscript text. Evaluated values were
 1219 1, 2, 4, 8, 16, 30, 60, 120, 240, 480 for Z, and 20, 25, 30, 40, 50, 60, 70, 80, 90,
 1220 100 for P. 60

1221 **Table C2.** Optimal LR parameters obtained in cross-validation for the C parameter and
 1222 regularization type for all lead times and regions. Evaluated for C were 0.0001,
 1223 0.0008, 0.0060, 0.0464, 0.359, 2.78, 21.54, 167.8, 1291, and 10000. 61

Symbol	Description	Predictor Group	Grid
APCP	Precipitation accumulation in past (3) 6 hours	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 two meters above ground	Core	Native Gaussian
T2M	Air temperature two meters above ground	Core	Native Gaussian
U10	Zonal-component of 10-meter wind	Core	Native Gaussian
V10	Meridional-component of 10-meter wind	Core	Native Gaussian
Q300	Specific humidity at 300 hPa	Upper-Air Extra	1°x1°
Q500	Specific humidity at 500 hPa	Upper-Air Core	1°x1°
Q700	Specific humidity at 700 hPa	Upper-Air Extra	1°x1°
Q850	Specific humidity at 850 hPa	Upper-Air Core1	1°x1°
T250	Temperature at 250 hPa	Upper-Air Extra	1°x1°
T500	Temperature at 500 hPa	Upper-Air Core	1°x1°
T700	Temperature at 700 hPa	Upper-Air Extra	1°x1°
T850	Temperature at 850 hPa	Upper-Air Core	1°x1°
U250	Zonal-component of 250 hPa wind	Upper-Air Extra	1°x1°
U500	Zonal-component of 500 hPa wind	Upper-Air Core	1°x1°
U700	Zonal-component of 700 hPa wind	Upper-Air Extra	1°x1°
U850	Zonal-component of 850 hPa wind	Upper-Air Core	1°x1°
V250	Meridional-component of 250 hPa wind	Upper-Air Extra	1°x1°
V500	Meridional-component of 500 hPa wind	Upper-Air Core	1°x1°
V700	Meridional-component of 700 hPa wind	Upper-Air Extra	1°x1°
V850	Meridional-component of 850 hPa wind	Upper-Air Core	1°x1°
W850	Vertical velocity (omega) at 850 hPa	Upper-Air Core	1°x1°

1224 TABLE 1. Summary of dynamical model fields examined in this study, including the abbreviated symbol to
1225 which each variable is referred throughout the paper, a description of each variable, the predictor group with
1226 which the field is associated in the manuscript text, and the highest resolution for which the field can be obtained
1227 from the GEFS/R.

Symbol	Description
RPT1_LOCAL_MEDIAN	Median of 1-year RPTs whose closest GEFS/R grid point is the forecast point.
RPT1_LOCAL_MIN	Minimum of 1-year RPTs whose closest GEFS/R grid point is the forecast point.
RPT1_LOCAL_MAX	Maximum of 1-year RPTs whose closest GEFS/R grid point is the forecast point.
RPT10_LOCAL_MEDIAN	Median of 10-year RPTs whose closest GEFS/R grid point is the forecast point.
RPT10_LOCAL_MIN	Minimum of 10-year RPTs whose closest GEFS/R grid point is the forecast point.
RPT10_LOCAL_MAX	Maximum of 10-year RPTs whose closest GEFS/R grid point is the forecast point.
RPT1_REGIONAL_MEDIAN	Median of 1-year RPTs that lie within the domain from which model predictors are drawn.
RPT1_REGIONAL_MIN	Minimum of 1-year RPTs that lie within the domain from which model predictors are drawn.
RPT1_REGIONAL_MAX	Maximum of 1-year RPTs that lie within the domain from which model predictors are drawn.
RPT10_REGIONAL_MEDIAN	Median of 10-year RPTs that lie within the domain from which model predictors are drawn.
RPT10_REGIONAL_MIN	Minimum of 10-year RPTs that lie within the domain from which model predictors are drawn.
RPT10_REGIONAL_MAX	Maximum of 10-year RPTs that lie within the domain from which model predictors are drawn.
LAT	Latitude of forecast point.
LON	Longitude of forecast point.

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

Model Name	CTL_NPC	CTL_PCA	UAC_PCA	UAF_PCA	CORE_CNFD	CORE_CTRL	CORE_LSPACE	CORE_LTIME	CTL_LR
Algorithm	RF	RF	RF	RF	RF	RF	RF	RF	LR
PCA Pre-Processed		X	X	X	X				X
Uses Core Fields	X	X	X	X	X	X	X	X	X
Uses UAC Fields			X	X					
Uses UAE Fields				X					
Ensemble Information	MEDIAN	MEDIAN	MEDIAN	MEDIAN	CNFD	CTRL	MEDIAN	MEDIAN	MEDIAN
Horizontal Radius	4	4	4	4	4	4	0 (a), 1 (b), 2 (c), 3 (d)	4	4
Timestep	3/6	3/6	3/6	3/6	3/6	3/6	3/6	12 (a), 6 (b*)	3/6

1228 **TABLE 3.** Summary of the models trained in this study, and the corresponding names designated to the models.
1229 ‘X’ indicates the process is performed or the information is used; a lack of one indicates the opposite. MEDIAN
1230 corresponds to the ensemble median, CTRL corresponds to the ensemble control member’s fields, and CNFD
1231 uses the median in addition to the second-from-lowest and second-from-highest member values for each field.
1232 Horizontal radius is listed in grid boxes from forecast point; timestep denotes the number of hours between
1233 GEFS/R forecast field predictors. Slashes indicate the first number applies to the Day 2 version of the model,
1234 while the latter number applies to the Day 3 version. Letters enclosed by parentheses indicate sub-versions
1235 of models, with one parameter changed to the value adjacent to the letter. Asterisks indicate a model applies
1236 only to Day 2, and not Day 3. Otherwise, models apply to all eight forecast regions and have both Day 2 and
1237 Day 3 versions. Those models with bolded names are incorporated into the weighted blend of the final model
1238 configuration.

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

1239 Table C1. Optimal RF parameters obtained in cross-validation for the Z-S-P parameter space. SQRT
1240 indicates the square root of the total number of predictors; symbols are otherwise as described in the
1241 manuscript text. Evaluated values were 1, 2, 4, 8, 16, 30, 60, 120, 240, 480 for Z, and 20, 25, 30, 40, 50, 60, 70,
1242 80, 90, 100 for P.

Region	Regularization	C Parameter, Day 2	C Parameter, Day 3
ROCK	L2	0.0001	0.0001
NGP	L2	10000	0.0008
MDWST	L2	0.0001	0.0464
NE	L2	2.78	10000
PCST	L2	0.0001	0.0001
SW	L2	0.359	0.0001
SGP	L2	0.0008	0.0464
SE	L2	21.54	0.0008

1243 Table C2. Optimal LR parameters obtained in cross-validation for the C parameter and regularization type
1244 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,
1245 1291, and 10000.

LIST OF FIGURES

1247 **Fig. 1.** Schematic representation of the forecast process for this study. GEFS/R forecasts are taken,
1248 assembled across fields, space, and time to form a training matrix, and past observations
1249 are used to associate a label with each forecast initialization, forecast day, forecast point
1250 triplet. The training matrix optionally undergoes pre-processing through principal compo-
1251 nent analysis, and then is input to one or more machine learning algorithms. From here,
1252 probabilistic ARI exceedance forecasts may be readily generated. 64

1253 **Fig. 2.** Return period thresholds at the (a) 1-year and (b) 10-year ARI levels over CONUS for a 24-
1254 hour accumulation interval. Climatology of observed exceedances of the (c) 1-year, 24-hour
1255 ARI thresholds and (d) 10-year, 24-hour ARI thresholds between January 2003 and August
1256 2013 based on Stage IV Precipitation Analysis. Pie charts indicate the monthly distribution
1257 of event occurrence within each study region as shown in Figure 3. Numbers above the pie
1258 charts indicate the mean number of exceedances per point per year within the region (a priori
1259 1 and 0.1 for 1-year and 10-year ARIs, respectively). 65

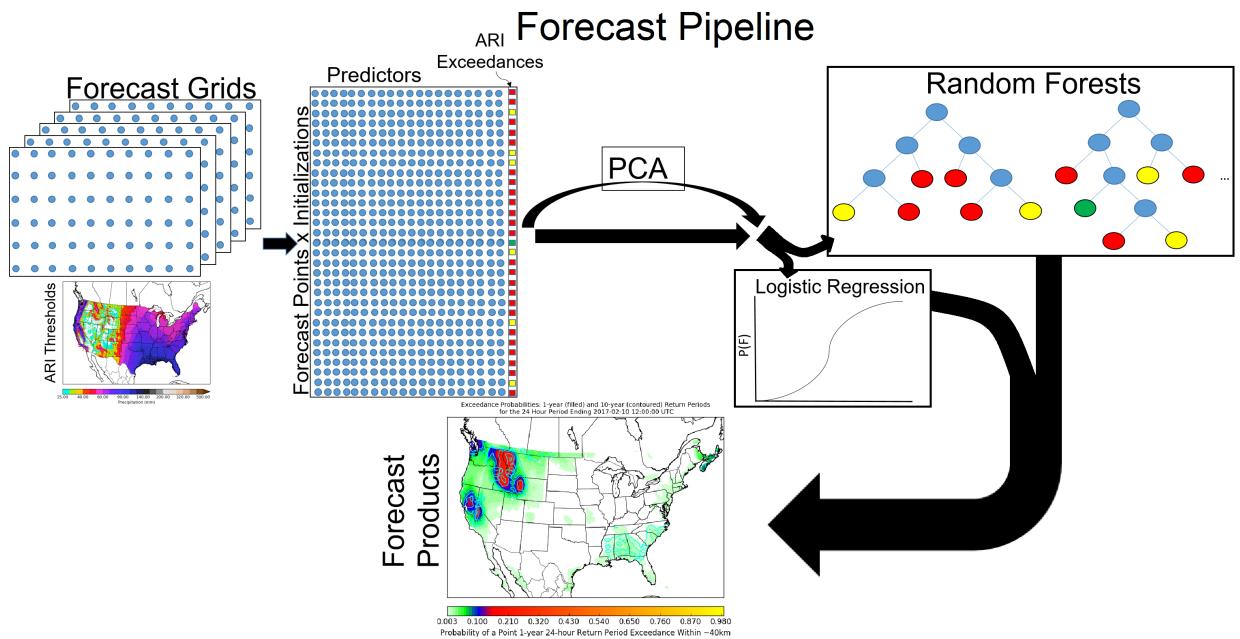
1260 **Fig. 3.** Map depicting the regional partitioning of CONUS used in this study, and the labels ascribed
1261 to each region. 66

1262 **Fig. 4.** Sensitivity experiment RPSS results for (a) the CORE_LTIME models, as a function of
1263 the timestep between incorporation of new atmospheric field forecast values, and (b) the
1264 CORE_LSPACE models, as a function of the radius of predictor information incorporated,
1265 both including both Day 2 and Day 3 versions of the model and for each region studied.
1266 Lines correspond to a particular day, region pair as indicated in the respective panel legends.
1267 Error bars in both panels correspond to 90% confidence bounds obtained by bootstrapping. 67

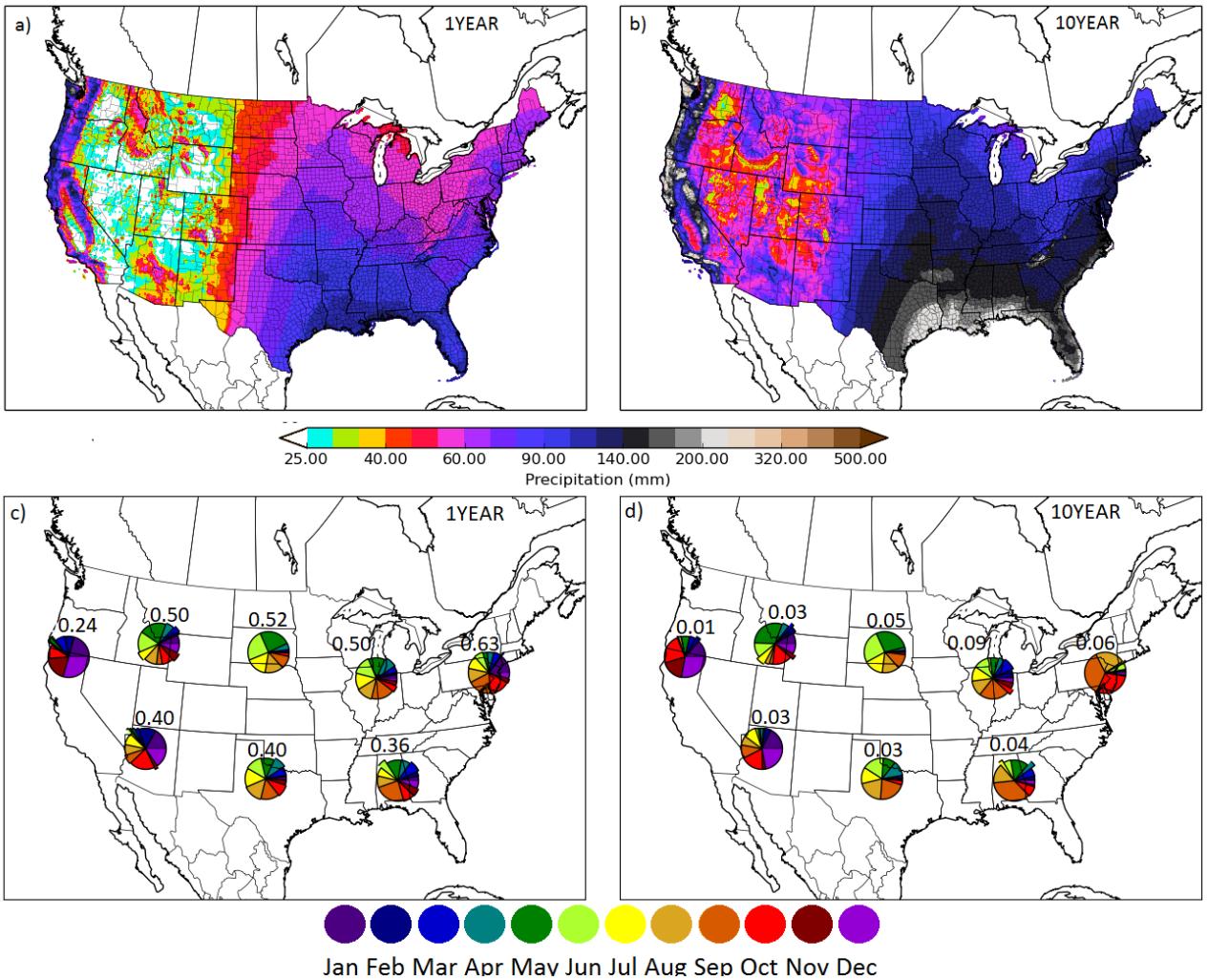
1268 **Fig. 5.** Sensitivity experiment RPSS results. Panel (a) as a function of the atmospheric fields in-
1269 cluded as input to the RF algorithm, for Day 3 forecast and broken out by region. From left
1270 to right, the columns correspond to results using: 1) just the ‘Core’ atmospheric field group,
1271 2) both the ‘Core’ and ‘Upper Air Core’ groups, 3) the ‘Core’, ‘Upper Air Core’, and ‘Up-
1272 per Air Extra’ groups. For more information on which fields are included in each predictor
1273 group, consult Table 1. Panel (b) as a function of the type of GEFS/R information used as
1274 input predictors to the RF algorithm, for Day 3 forecasts and broken out by region. From left
1275 to right, the columns correspond to results using: 1) just the forecast fields from the GEFS/R
1276 control member, 2) the ensemble median forecast values from the full ensemble, 3) the en-
1277 semble median, 2nd-from-minimum, and 2nd-from-maximum forecast values from the full
1278 ensemble. Panel (c) as a function of region aggregation, with the left column using the eight
1279 regions depicted in Figure 3, and the right column using training data which aggregates data
1280 from seven of the eight original regions into three regions, as described in the text. Panel
1281 (d) as a function of model algorithm for different forecast days and regions as indicated
1282 in the figure legend. From left to right, columns correspond to results of the CTL_NPCA
1283 model, CTL_PCA model, CTL_LR model, and a weighted combination of CTL_PCA and
1284 CTL_LR models as described in the paper text. For all panels, error bars correspond to 90%
1285 confidence bounds obtained by bootstrapping. 68

1286 **Fig. 6.** Final RPSS results obtained over the four year test period spanning September 2013–August
1287 2017, broken out by region. Red bars correspond to the results of the final forecast models
1288 trained in this study, while gray bars depict results from the raw GEFS/R QPF probabilities
1289 derived from the full ensemble. Dark bars illustrate Day 2 performance results, while lighter
1290 colors show results for Day 3. Error bars correspond to 90% confidence bounds obtained by
1291 bootstrapping. 69

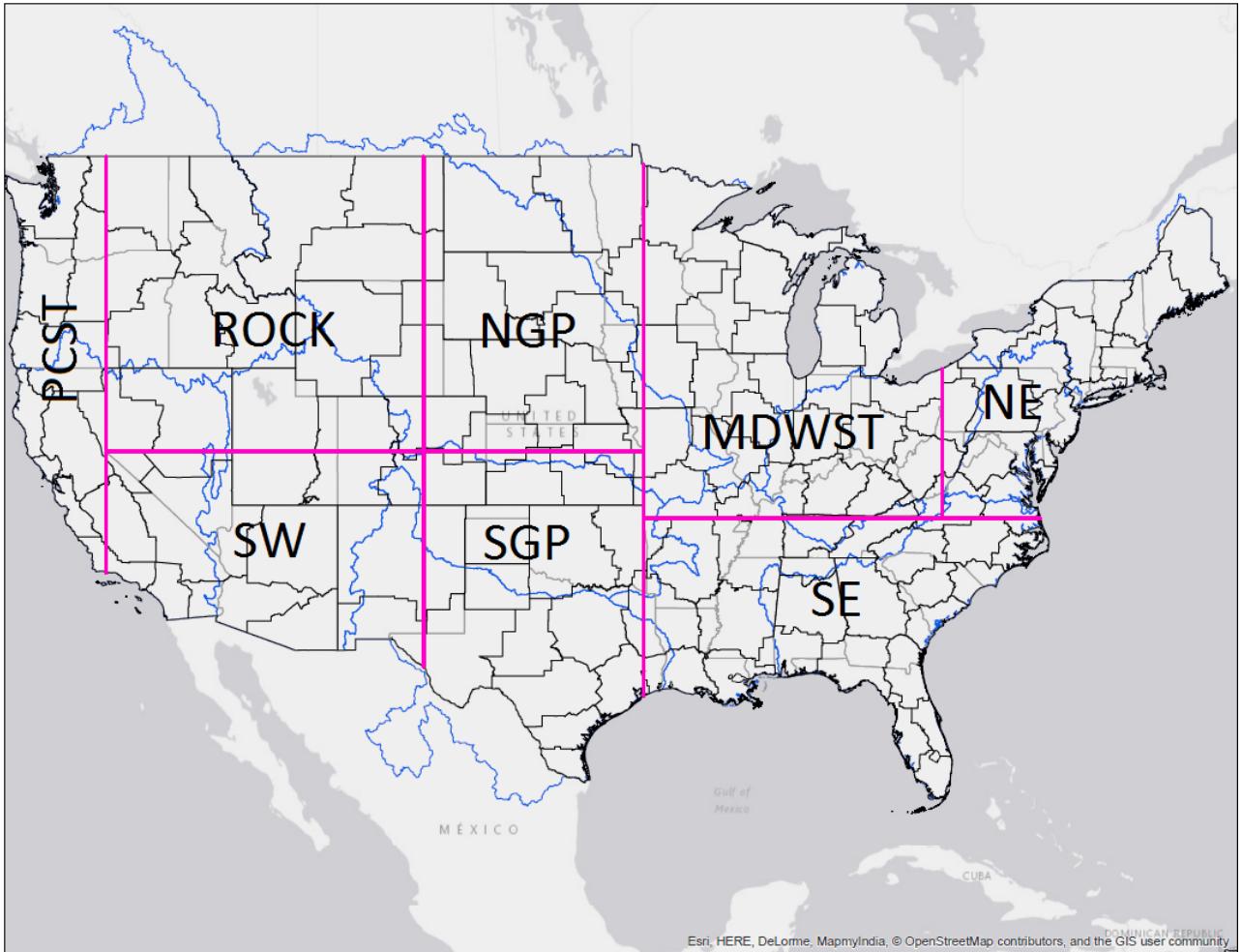
1292	Fig. 7.	Reliability diagrams for forecasts generated from raw QPFs of the full GEFS/R ensemble. Colored opaque lines with circular points indicate observed relative frequency as a function of forecast probability; the solid black line is the one-to-one line, indicating perfect reliability. Colors 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. Panel (a) shows Day 2 forecast results for 1-year ARI exceedance forecasts, (b) to Day 2 10-year ARI exceedance forecasts, (c) to Day 3 1-year exceedance forecasts, and panel (d) to Day 3 10-year ARI exceedance forecasts. 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. 70	70
1305	Fig. 8.	Same as Figure 7, but for ECMWF ensemble forecasts. 71	71
1306	Fig. 9.	Same as Figure 7, but for forecasts generated from the CTL_NPCA model. 72	72
1307	Fig. 10.	Same as Figure 7, but for forecasts from the CTL_PCA model. 73	73
1308	Fig. 11.	Same as Figure 7, but for forecasts from the CTL_LR model. 74	74
1309	Fig. 12.	Same as Figure 7, but for the final forecast model. 75	75
1310	Fig. 13.	Modified Murphy (1973) decomposition results, following equation 3 in text. Panel (a) depicts the equation 3 “resolution” term for all models and regions for Day 2 forecasts at the 1-year severity level, panel (b) depicts the “reliability” term results for the same forecasts and severity level. Panels (c) and (d) are analogous to panels (a) and (b), but for 10-year 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). 76	76
1316	Fig. 14.	Same as Figure 13, but for Day 3 forecasts. 77	77
1317	Fig. 15.	Case study depicting forecasts from the final ML model and both reference ensembles for the 24-hour period ending 1200 UTC 20 May 2015. (a) 24-hour Stage IV QPE ending at 1200 UTC 20 May 2015 in filled contours. The unfilled cyan contours enclose areas with 1-year 24-hour ARI exceedances, and the unfilled yellow contours enclose areas of 10-year 24-hour ARI exceedances. (b) ECMWF ensemble neighborhood ARI exceedance probabilities in the filled (1-year) and unfilled (24-year) contours for the 36–60 hour forecast initialized 0000 UTC 18 May 2015 and (c) for the 60–84 hour forecast initialized 0000 UTC 17 May 2015. Panels (d) and (e) depict analogous fields as panels (b) and (c), respectively, except for forecasts from the raw GEFS/R QPFs. Panels (f) and (g) similarly show respectively 36–60 and 60–84 hour forecasts, except for from the final version of the ML model trained in this study. 78	78
1328	Fig. 16.	Same as Figure 15, but for the 24-hour period ending 1200 UTC 22 September 2016. 79	79



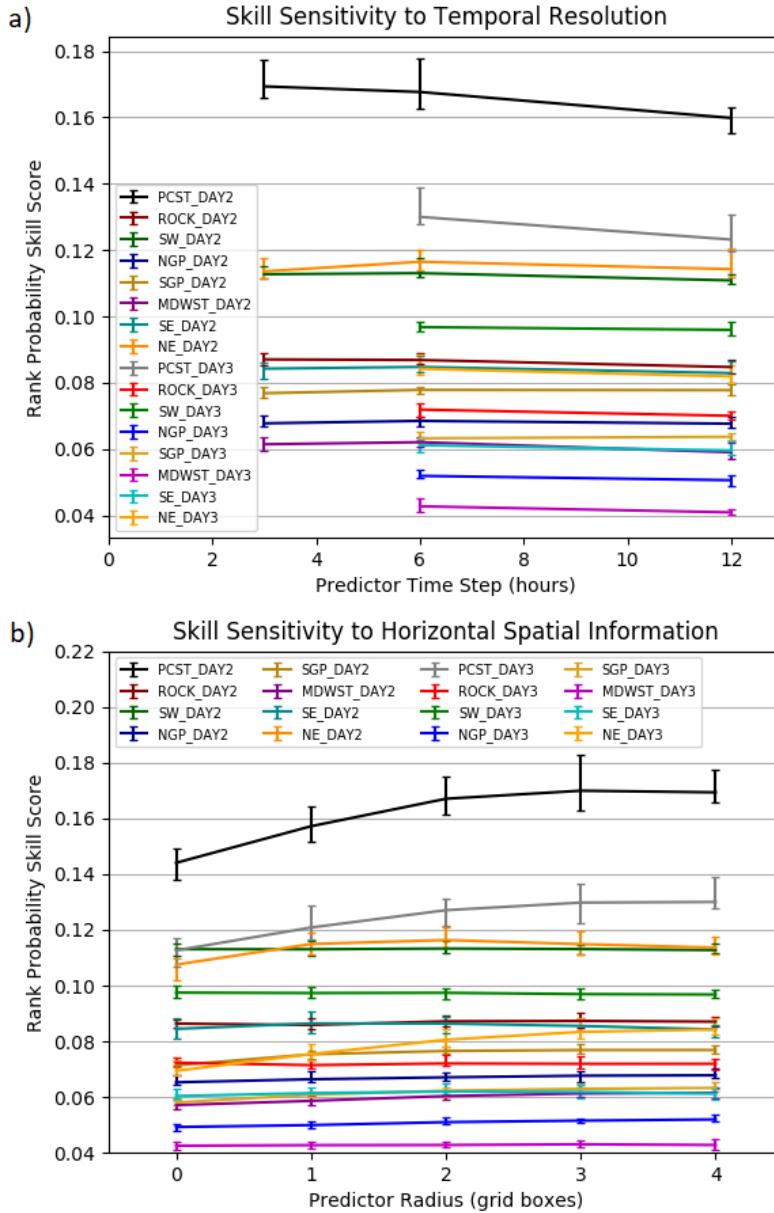
1329 FIG. 1. Schematic representation of the forecast process for this study. GEFS/R forecasts are taken, assembled
 1330 across fields, space, and time to form a training matrix, and past observations are used to associate a label with
 1331 each forecast initialization, forecast day, forecast point triplet. The training matrix optionally undergoes pre-
 1332 preprocessing through principal component analysis, and then is input to one or more machine learning algorithms.
 1333 From here, probabilistic ARI exceedance forecasts may be readily generated.



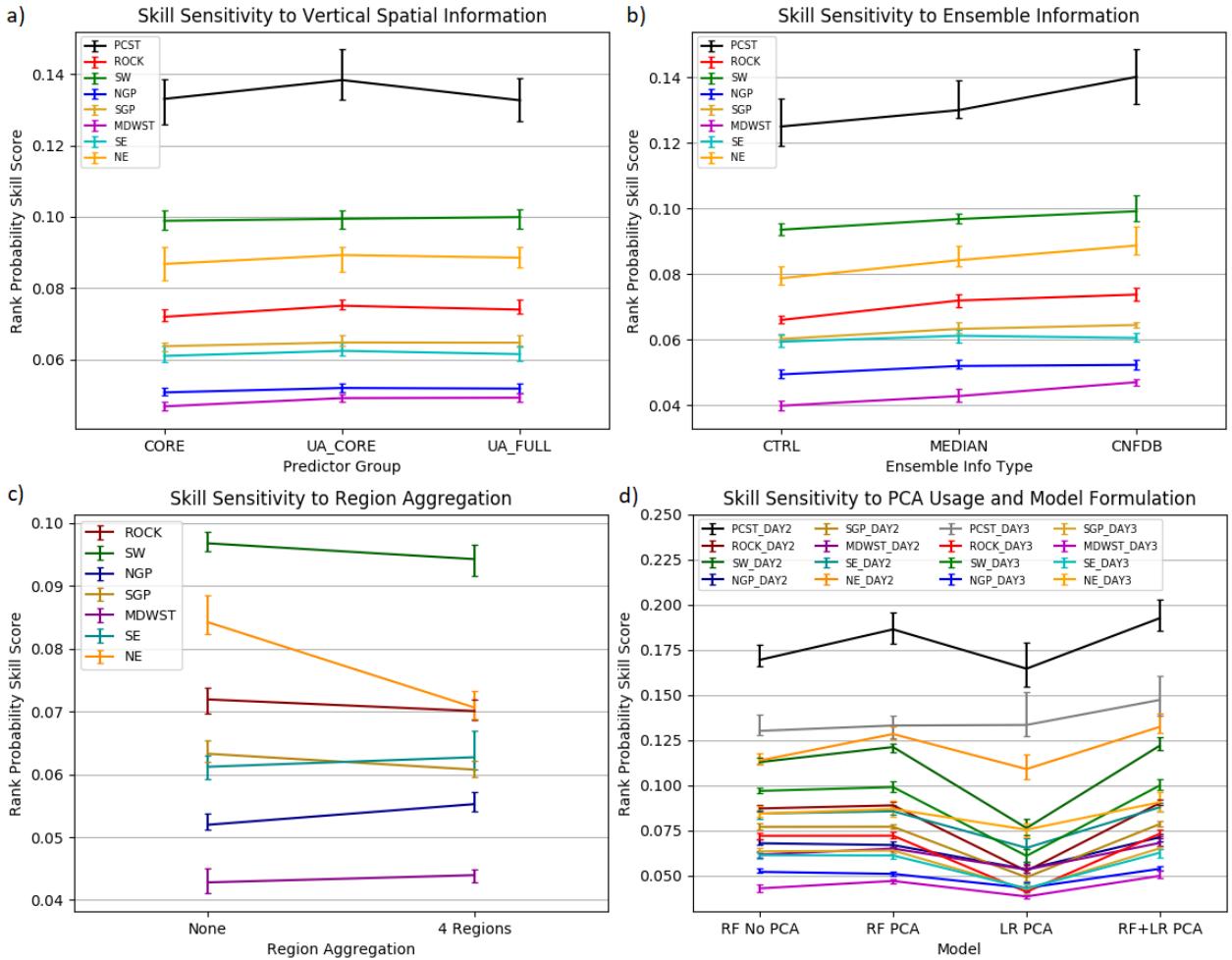
1334 FIG. 2. Return period thresholds at the (a) 1-year and (b) 10-year ARI levels over CONUS for a 24-hour
 1335 accumulation interval. Climatology of observed exceedances of the (c) 1-year, 24-hour ARI thresholds and
 1336 (d) 10-year, 24-hour ARI thresholds between January 2003 and August 2013 based on Stage IV Precipitation
 1337 Analysis. Pie charts indicate the monthly distribution of event occurrence within each study region as shown in
 1338 Figure 3. Numbers above the pie charts indicate the mean number of exceedances per point per year within the
 1339 region (a priori 1 and 0.1 for 1-year and 10-year ARIs, respectively).



1340 FIG. 3. Map depicting the regional partitioning of CONUS used in this study, and the labels ascribed to each
 1341 region.

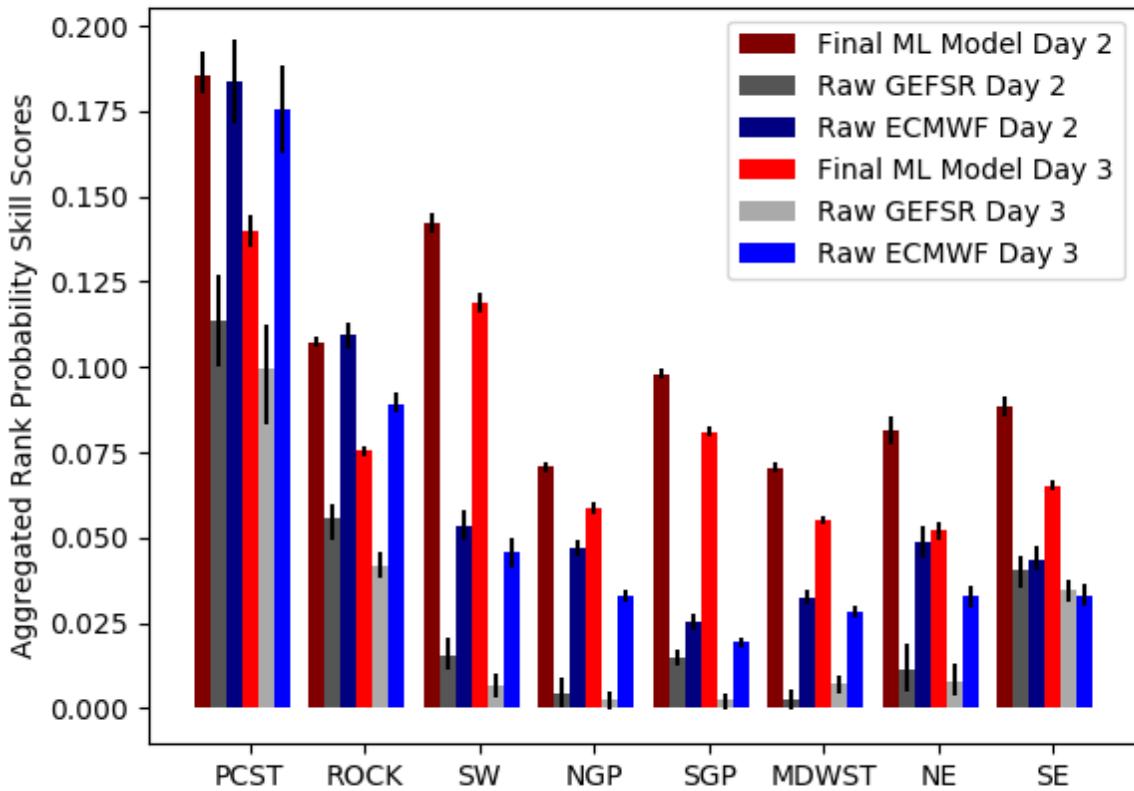


1342 FIG. 4. Sensitivity experiment RPSS results for (a) the CORE.LTIME models, as a function of the timestep
 1343 between incorporation of new atmospheric field forecast values, and (b) the CORE.LSPACE models, as a func-
 1344 tion of the radius of predictor information incorporated, both including both Day 2 and Day 3 versions of the
 1345 model and for each region studied. Lines correspond to a particular day, region pair as indicated in the respective
 1346 panel legends. Error bars in both panels correspond to 90% confidence bounds obtained by bootstrapping.

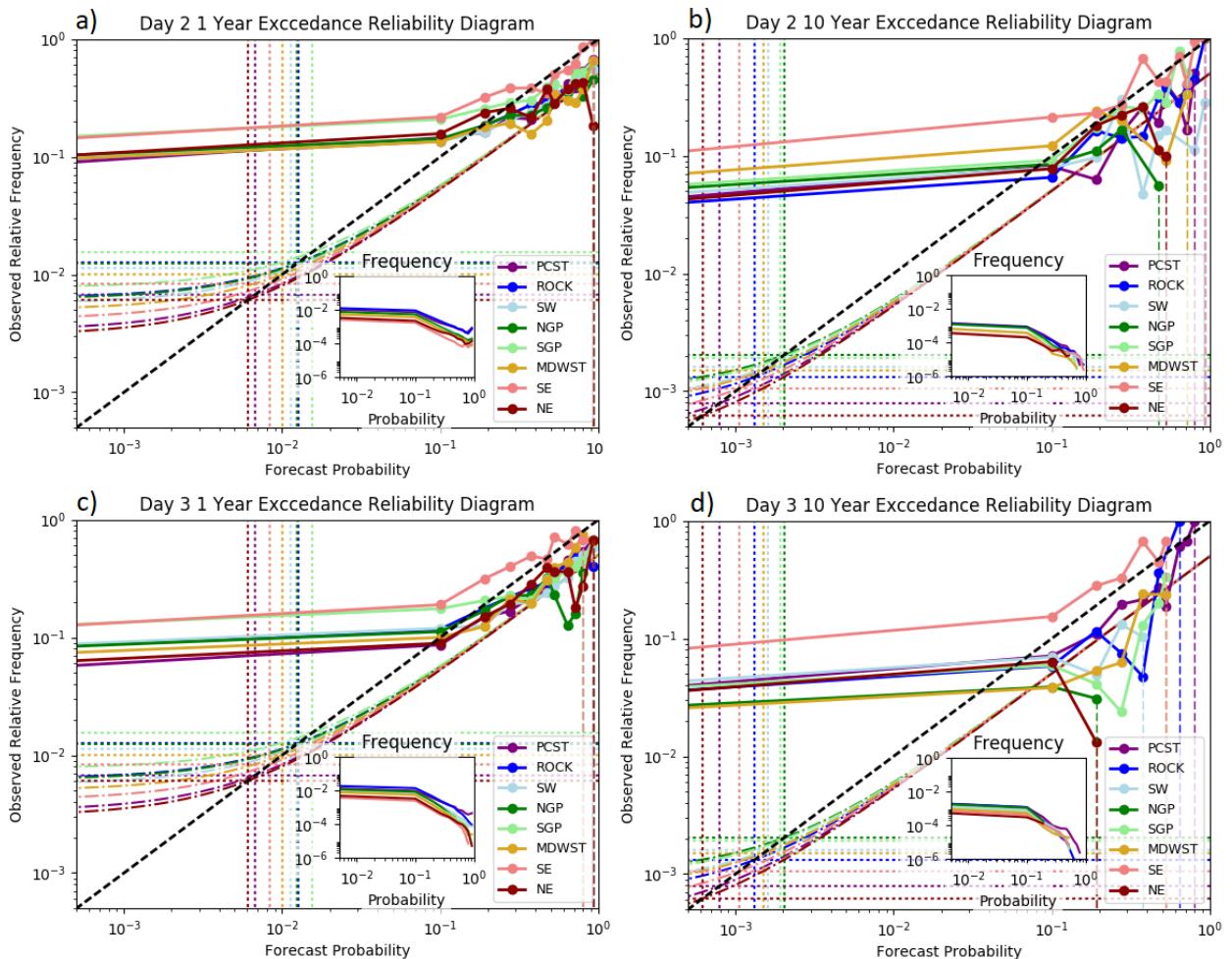


1347 FIG. 5. Sensitivity experiment RPSS results. Panel (a) as a function of the atmospheric fields included as input
 1348 to the RF algorithm, for Day 3 forecast and broken out by region. From left to right, the columns correspond to
 1349 results using: 1) just the ‘Core’ atmospheric field group, 2) both the ‘Core’ and ‘Upper Air Core’ groups, 3) the
 1350 ‘Core’, ‘Upper Air Core’, and ‘Upper Air Extra’ groups. For more information on which fields are included in
 1351 each predictor group, consult Table 1. Panel (b) as a function of the type of GEFS/R information used as input
 1352 predictors to the RF algorithm, for Day 3 forecasts and broken out by region. From left to right, the columns
 1353 correspond to results using: 1) just the forecast fields from the GEFS/R control member, 2) the ensemble median
 1354 forecast values from the full ensemble, 3) the ensemble median, 2nd-from-minimum, and 2nd-from-maximum
 1355 forecast values from the full ensemble. Panel (c) as a function of region aggregation, with the left column using
 1356 the eight regions depicted in Figure 3, and the right column using training data which aggregates data from
 1357 seven of the eight original regions into three regions, as described in the text. Panel (d) as a function of model
 1358 algorithm for different forecast days and regions as indicated in the figure legend. From left to right, columns
 1359 correspond to results of the CTL_NPC model, CTL_PCA model, CTL_LR model, and a weighted combination
 1360 of CTL_PCA and CTL_LR models as described in the paper text. For all panels, error bars correspond to 90%
 1361 confidence bounds obtained by bootstrapping.

Final Model Skill Score Report



1362 FIG. 6. Final RPSS results obtained over the four year test period spanning September 2013–August 2017,
 1363 broken out by region. Red bars correspond to the results of the final forecast models trained in this study, while
 1364 gray bars depict results from the raw GEFS/R QPF probabilities derived from the full ensemble. Dark bars
 1365 illustrate Day 2 performance results, while lighter colors show results for Day 3. Error bars correspond to 90%
 1366 confidence bounds obtained by bootstrapping.



1367 FIG. 7. Reliability diagrams for forecasts generated from raw QPFs of the full GEFS/R ensemble. Colored
 1368 opaque lines with circular points indicate observed relative frequency as a function of forecast probability; the
 1369 solid black line is the one-to-one line, indicating perfect reliability. Colors correspond to the performance of the
 1370 forecasts over different regions, as indicated in the legend in the lower-right of each panel. Inset panels indicate
 1371 the total proportion of forecasts falling in each forecast probability bin, using the logarithmic scale on the left
 1372 hand side of each panel; lines are again colored by region in accordance with the legend. Panel (a) shows Day 2
 1373 forecast results for 1-year ARI exceedance forecasts, (b) to Day 2 10-year ARI exceedance forecasts, (c) to Day 3
 1374 1-year exceedance forecasts, and panel (d) to Day 3 10-year ARI exceedance forecasts. All axes are logarithmic
 1375 as labeled. Colored dotted lines indicate the climatological event probability for each region for the ARI level of
 1376 the corresponding panel, while the dash-dotted lines indicate no skill lines for the color-corresponding region.

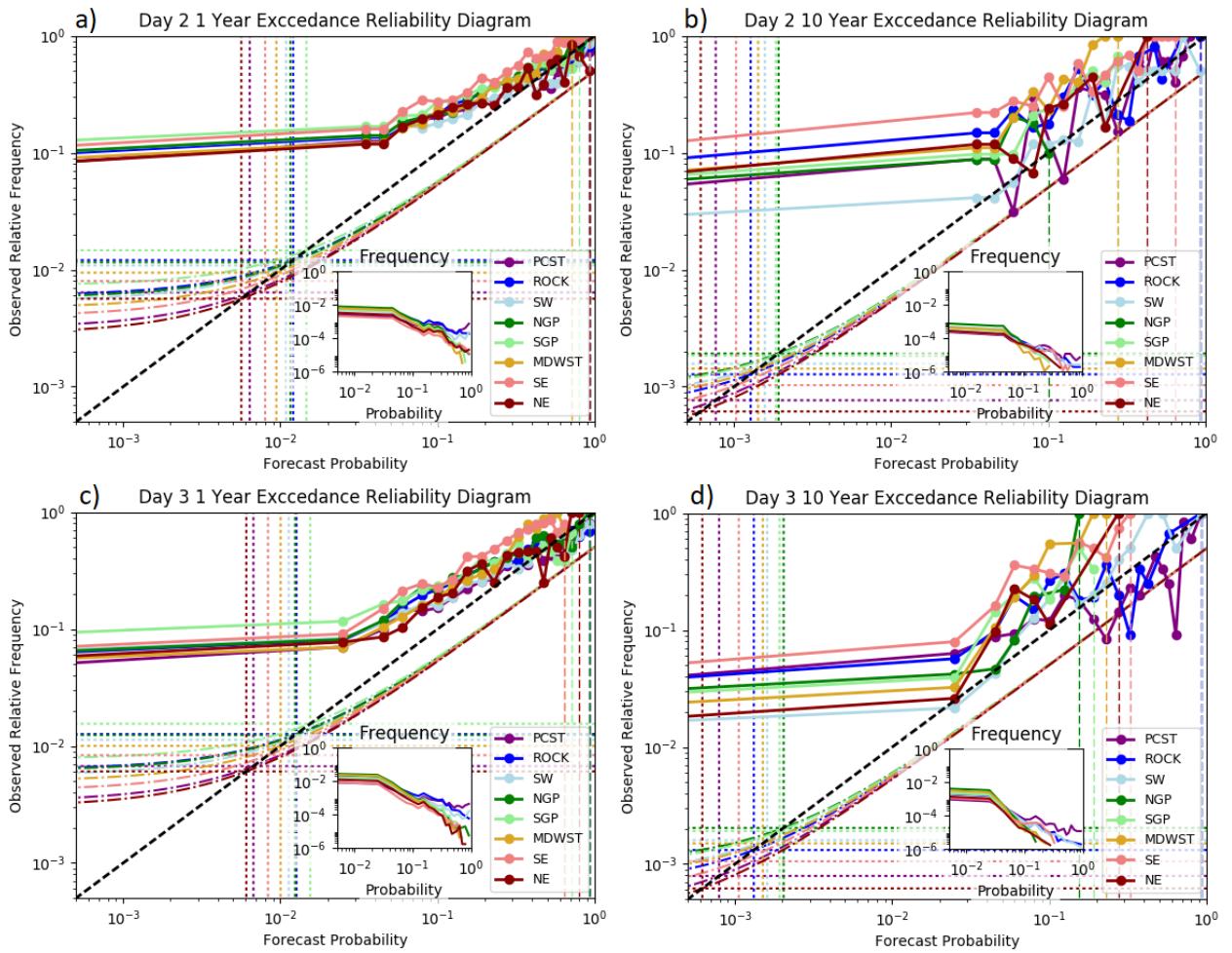


FIG. 8. Same as Figure 7, but for ECMWF ensemble forecasts.

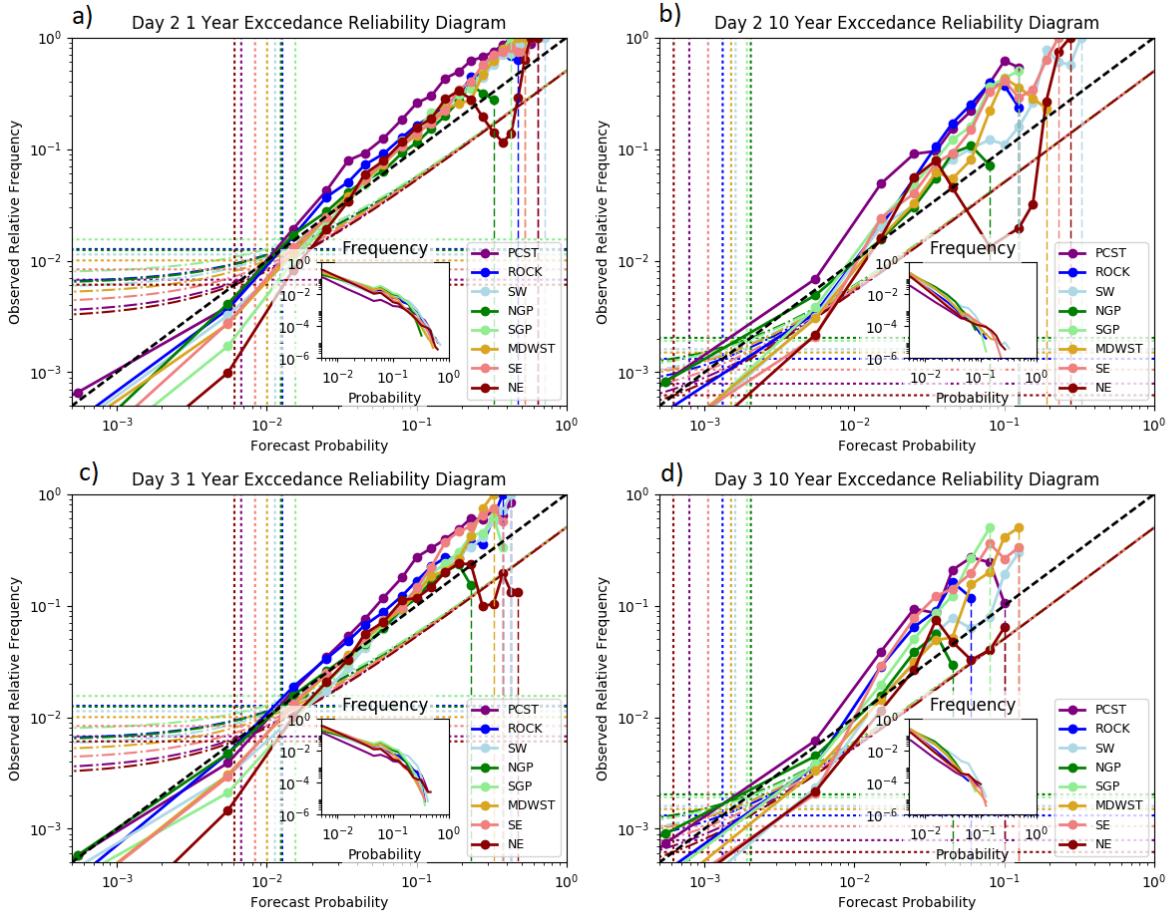


FIG. 9. Same as Figure 7, but for forecasts generated from the CTL_NPCA model.

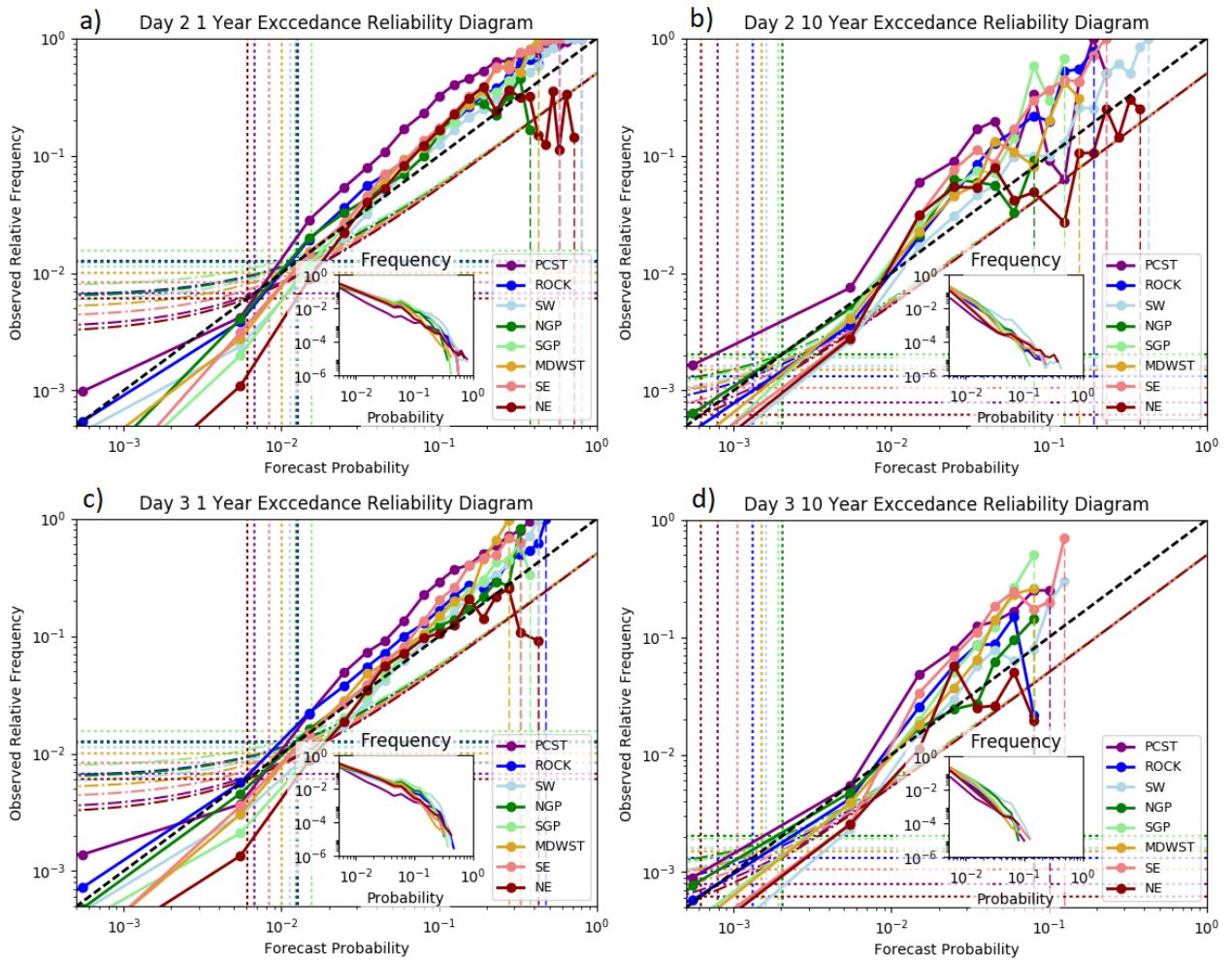


FIG. 10. Same as Figure 7, but for forecasts from the CTL_PCA model.

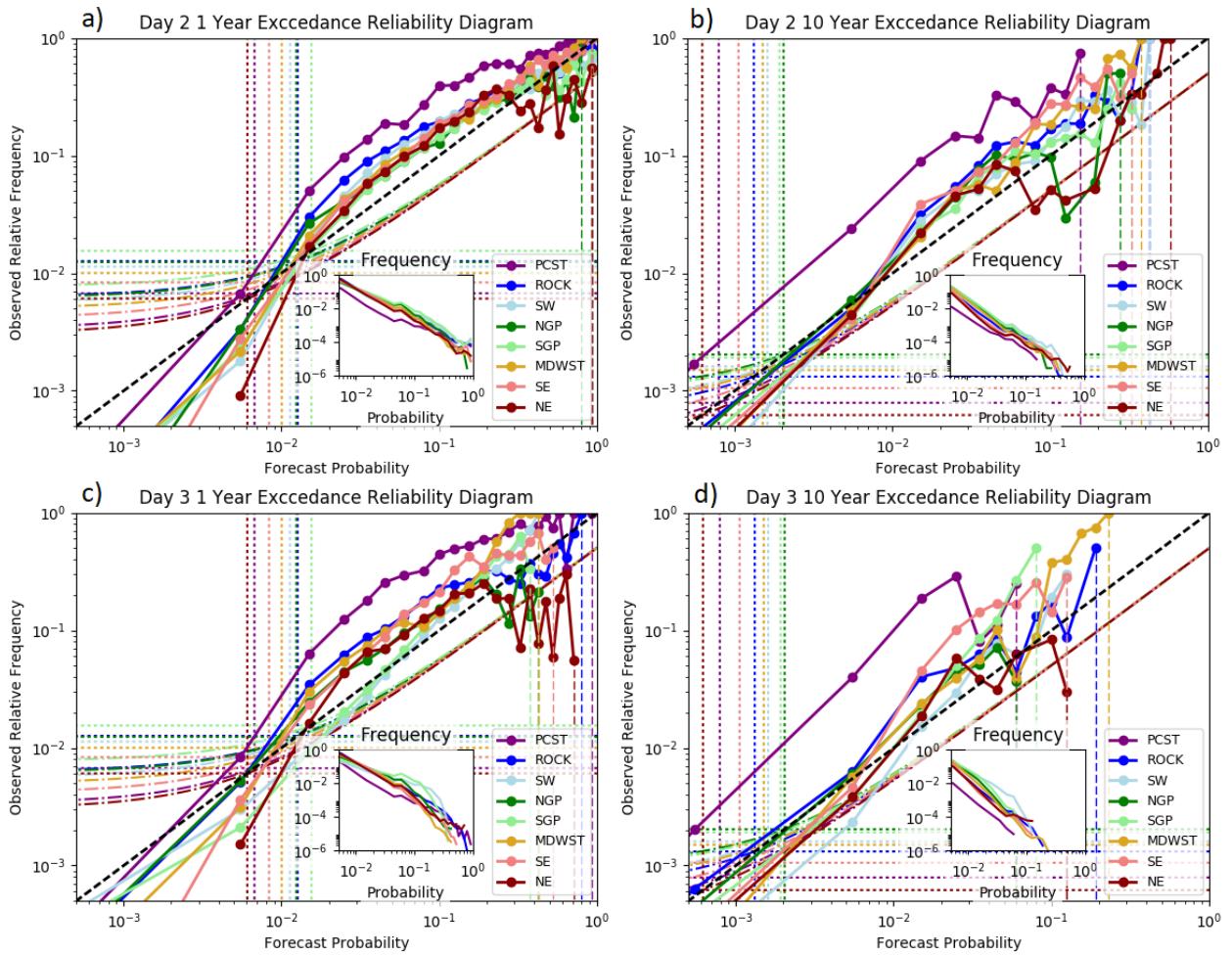


FIG. 11. Same as Figure 7, but for forecasts from the CTL_LR model.

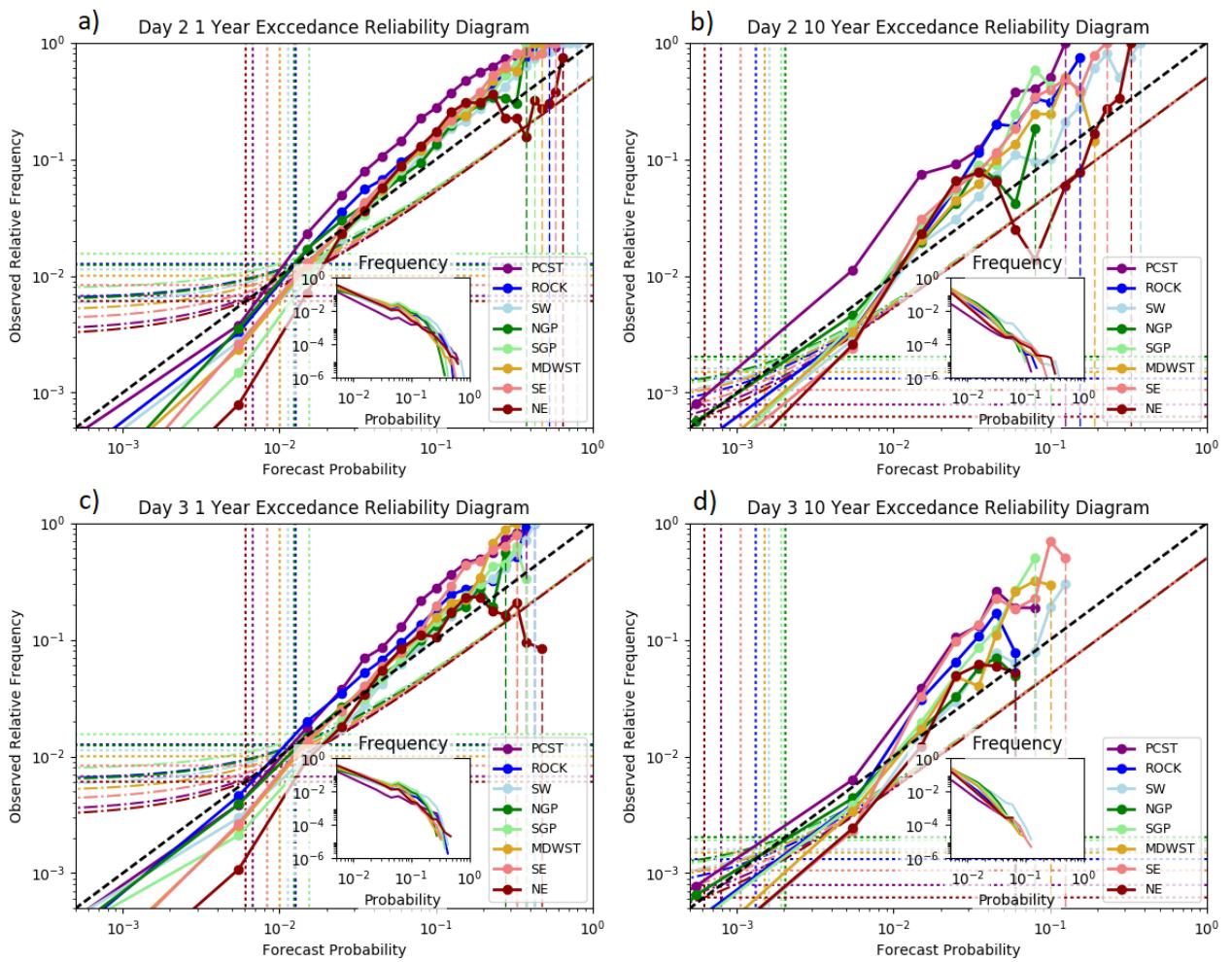
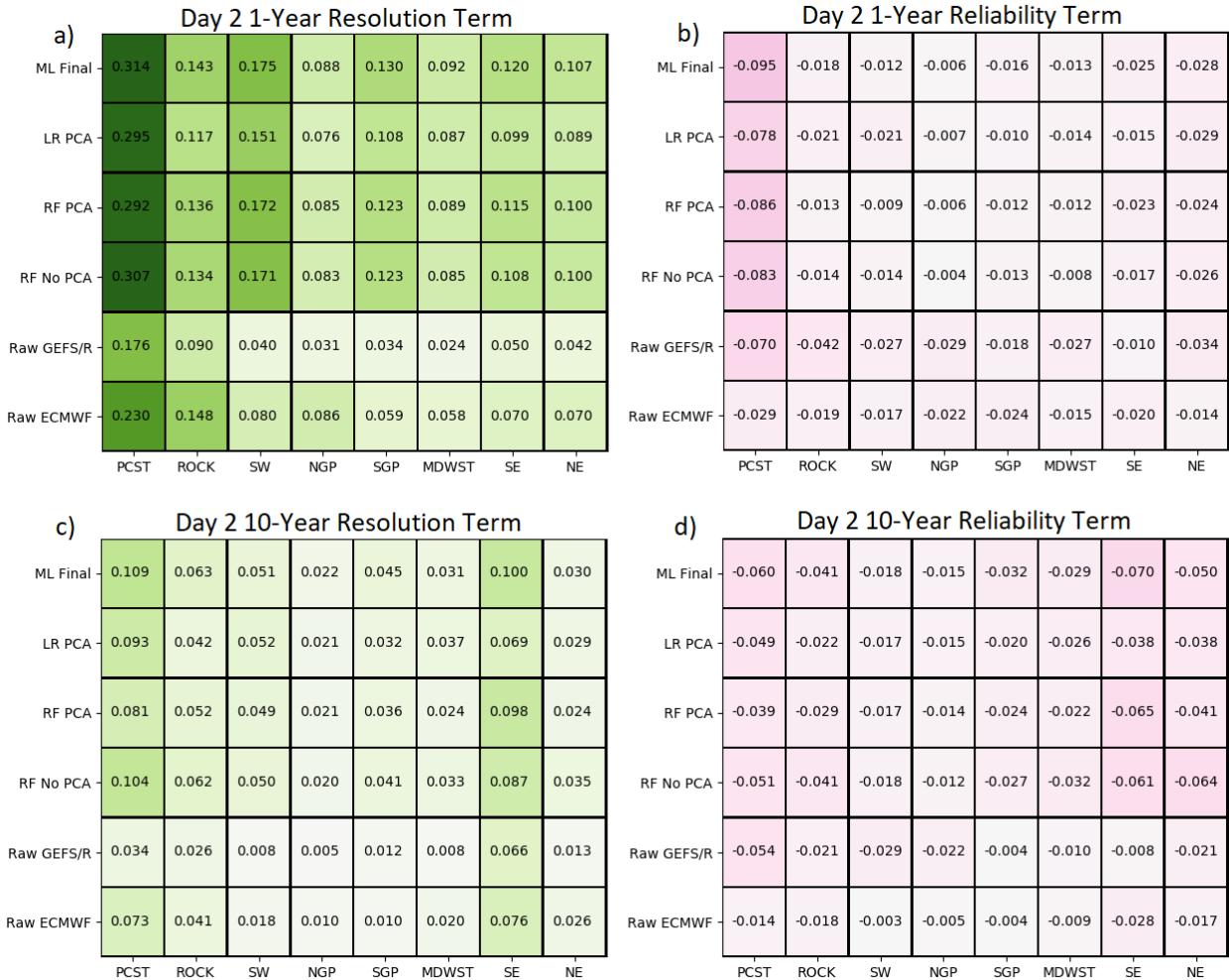


FIG. 12. Same as Figure 7, but for the final forecast model.



1377 FIG. 13. Modified Murphy (1973) decomposition results, following equation 3 in text. Panel (a) depicts
 1378 the equation 3 “resolution” term for all models and regions for Day 2 forecasts at the 1-year severity level,
 1379 panel (b) depicts the “reliability” term results for the same forecasts and severity level. Panels (c) and (d) are
 1380 analogous to panels (a) and (b), but for 10-year ARI exceedance forecasts. Numeric values indicate the value of
 1381 the corresponding term of the table, as indicated by the model label (row) and region (column).

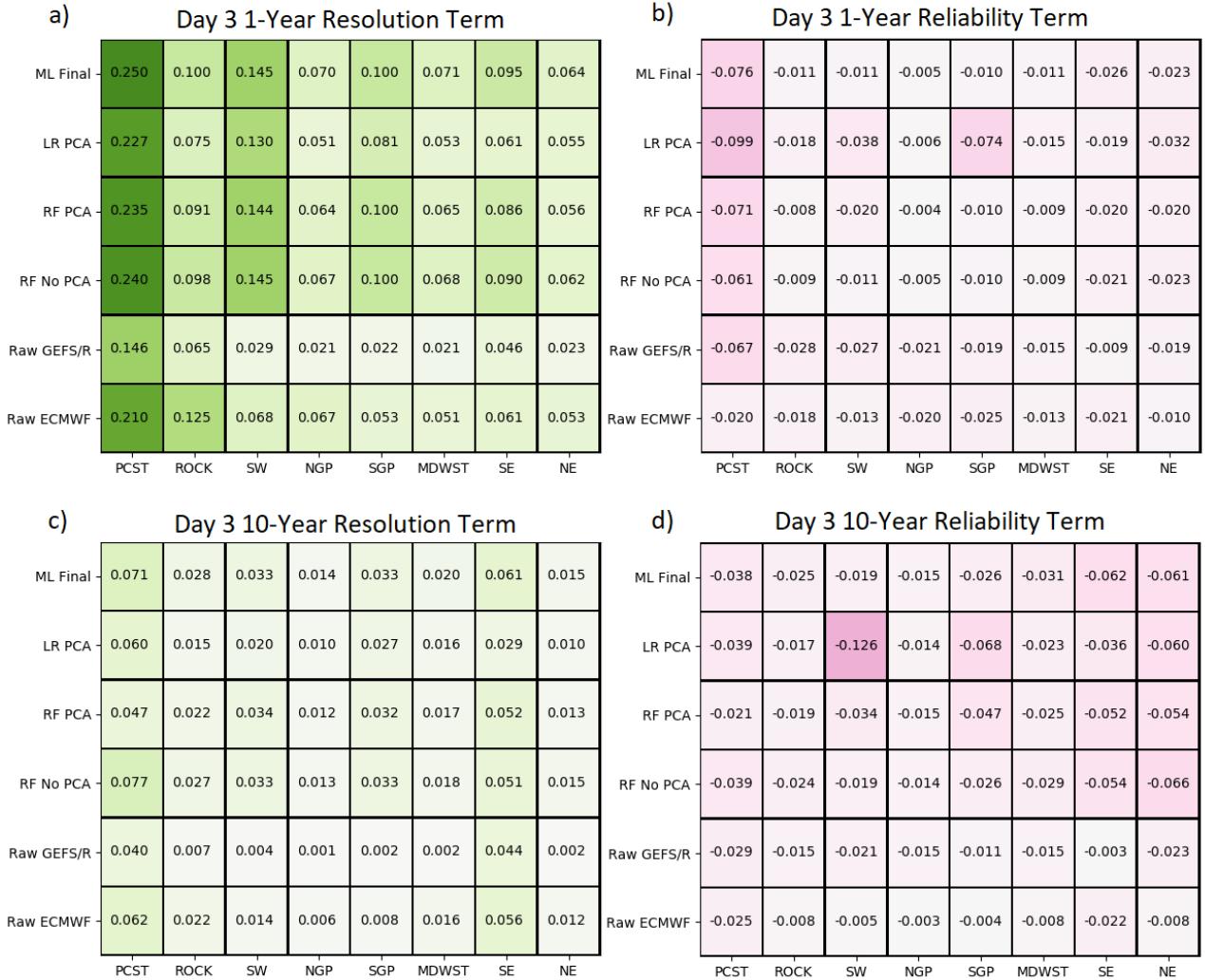
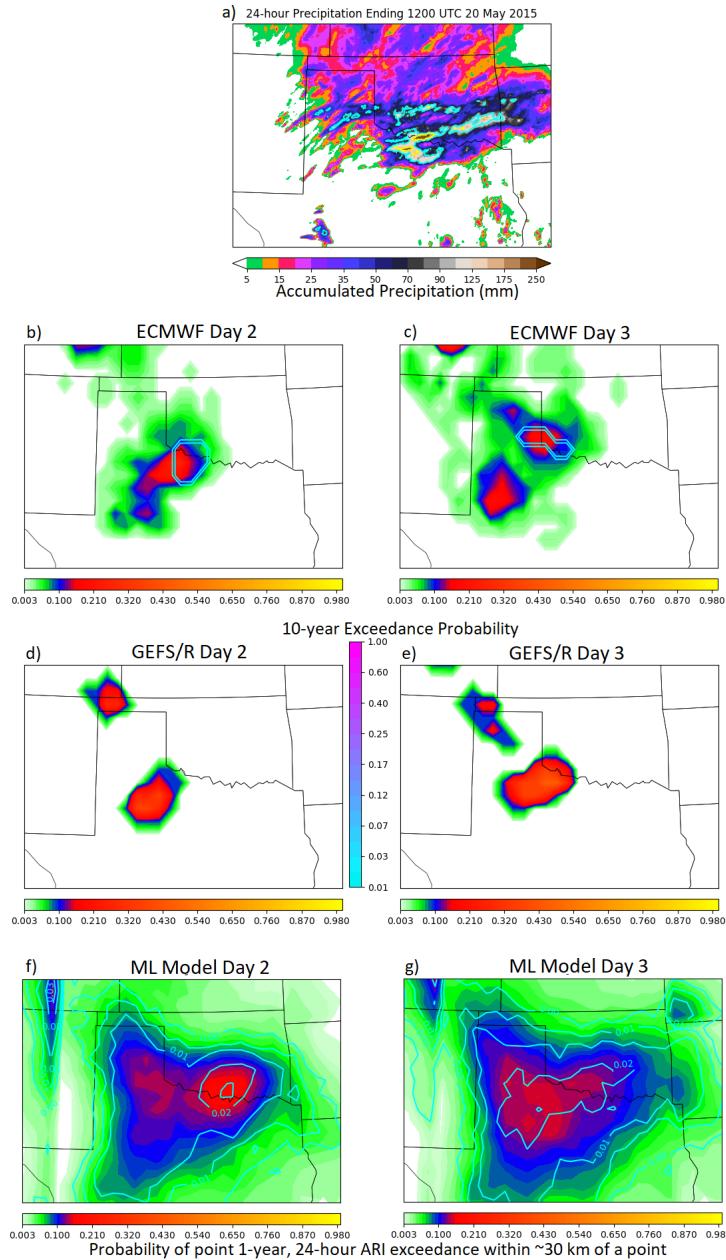


FIG. 14. Same as Figure 13, but for Day 3 forecasts.



1382 FIG. 15. Case study depicting forecasts from the final ML model and both reference ensembles for the 24-
 1383 hour period ending 1200 UTC 20 May 2015. (a) 24-hour Stage IV QPE ending at 1200 UTC 20 May 2015 in
 1384 filled contours. The unfilled cyan contours enclose areas with 1-year 24-hour ARI exceedances, and the unfilled
 1385 yellow contours enclose areas of 10-year 24-hour ARI exceedances. (b) ECMWF ensemble neighborhood ARI
 1386 exceedance probabilities in the filled (1-year) and unfilled (24-year) contours for the 36–60 hour forecast initial-
 1387 ized 0000 UTC 18 May 2015 and (c) for the 60–84 hour forecast initialized 0000 UTC 17 May 2015. Panels
 1388 (d) and (e) depict analogous fields as panels (b) and (c), respectively, except for forecasts from the raw GEFS/R
 1389 QPFs. Panels (f) and (g) similarly show respectively 36–60 and 60–84 hour forecasts, except for from the final
 1390 version of the ML model trained in this study. 78

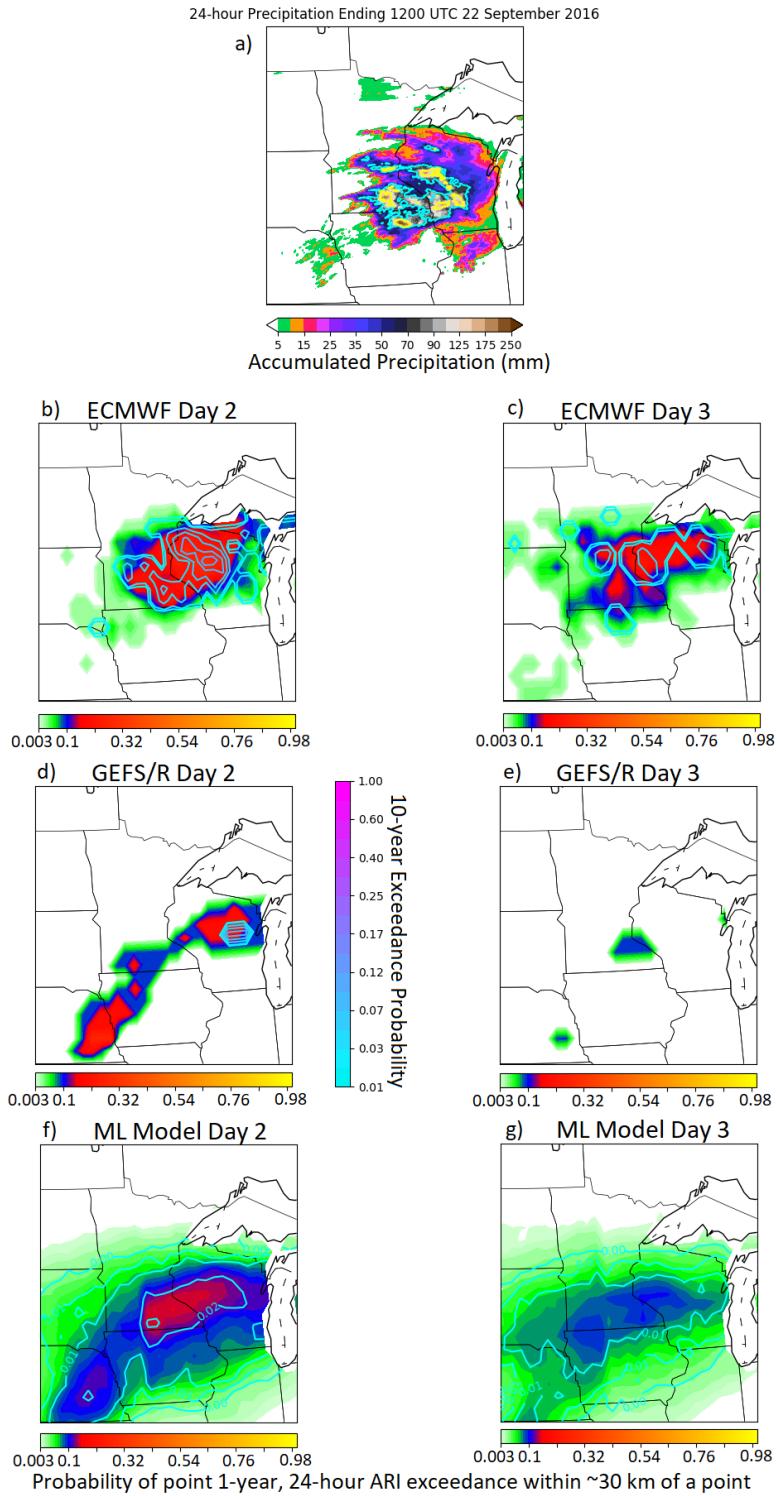


FIG. 16. Same as Figure 15, but for the 24-hour period ending 1200 UTC 22 September 2016.