1 **Objective methods for thinning the frequency of reforecasts while meeting**

2 **post-processing and model validation needs**

3

4

5 Sergey Kravtsov,[a] Paul Roebber,[a] Thomas M. Hamill,[b] James Brown[c]

6 [a] *University of Wisconsin-Milwaukee (UWM), Milwaukee, WI*

7 [b] *NOAA Physical Sciences Laboratory, Boulder CO*

8 [c] *Office of Water Prediction, National Weather Service (OWP NWS)*

9

10

11

12 *Corresponding author*: Sergey Kravtsov, kravtsov@uwm.edu

13

14                                    ABSTRACT

15      This paper utilizes statistical and statistical-dynamical methodologies to select, from the full

16      observational record, a minimal subset of dates that would provide representative sampling of

17      local precipitation distributions across the contiguous US (CONUS). The CONUS region is

18      characterized by a great diversity of precipitation-producing systems, mechanisms and large-

19      scale meteorological patterns (LSMPs) which can provide favorable environment for local

20      precipitation extremes. This diversity is unlikely to be adequately captured in methodologies

21      which rely on grossly reducing the dimensionality of the  data — by representing it in terms of

22      a few patterns evolving in time — and thus requires data thinning techniques based on high-

23      dimensional dynamical or statistical data modeling.  We have built a novel high-dimensional

24      empirical model of temperature and precipitation capable of producing highly statistically

25      accurate surrogate realizations of the observed 1979–1999 (training-period) evolution of these

26      fields. This model also provides skillful hindcasts of precipitation over the 2000–2020

27      (validation) period.  We devised a subsampling strategy based on the relative entropy of the

28      empirical model's precipitation (ensemble) forecasts over CONUS and demonstrated that it

29      generates a set of dates that captures a majority of high-impact precipitation events while

30      substantially reducing a heavy-precipitation bias inherent in an alternative methodology based

31      on the direct identification of large precipitation events in the Global Ensemble Forecast

32      System (GEFS,  version 12) reforecasts. The impacts of data thinning on the accuracy of

33      precipitation statistical post-processing, as well as on the calibration and validation of the

34      Hydrologic Ensemble Forecast Service (HEFS) reforecasts are yet to be established.

35

SIGNIFICANCE STATEMENT

High-impact weather events are usually associated with extreme precipitation, which is notoriously difficult to predict even using highly resolved state-of-the-art numerical weather prediction models based on first physical principles. The same is true for statistical models that use past data to anticipate the future behavior likely to stem from an observed initial condition. Here we use both types of models to identify the timing of initial conditions, over the historical climate record, that are likely to produce extreme precipitation events. We show that the overall statistics of precipitation over contiguous US can be encapsulated in a greatly reduced set of initial conditions, which makes testing and validation of hydrological forecast models and the associated decision support much less computationally expensive.

## 1. Introduction

The statistical post-processing of weather forecasts has been shown to be extremely useful for ameliorating model biases and extracting usable forecast signal amidst the noise due to chaotic error growth and sampling due to limited ensemble size (Hamill and Whitaker 2006; Hamill et al. 2006, 2013, 2015; Scheuerer and Hamill 2015). Post-processed forecasts are typically more skillful and reliable, rendering them useful for automated decision support. Large sample sizes of reforecasts are particularly helpful in four particular situations: (a) the post-processing of rare events, (b) the post-processing of longer-lead events, where usable signal is small, noise is large, and forecasts are for time-averaged quantities. While the production of a long, complete time series of reforecasts is desirable for such situations, the computational expense of reforecasting scales linearly with the reforecast sample size. Objective methods that can indicate what subset of dates are the most important to generate reforecasts are greatly desired. Given the national forecast responsibilities of the National Weather Service (NWS), that subset of dates should ideally be large enough to provide the necessary training and validation data over the contiguous US (CONUS).

There are several challenges to be anticipated with designing a procedure for reforecast sub-sampling. One challenge of sub-selecting past dates is that they will be less useful for training if the dates are based on the existence of *observed* high-impact weather such as heavy precipitation. In such a case, the training data is biased toward the existence of high-impact

3

File generated with AMS Word template 1.0

66    events, and post-processed guidance will likely over-forecast them. Accordingly, we seek
67    methodologies for deciding on which dates to use that avoid the use of validating observations
68    but instead use only information such as the initial condition state or the existence of conditions
69    related to severe weather at a similar date noted in previous reforecasts.

70        Yet another challenge could be the under-sampling of more commonplace events. Were
71    such a reforecast sub-sampling procedure designed for a very limited geographic area, dry
72    weather or light/moderate precipitation could be drastically under-sampled, leading to poor-
73    quality guidance of more common weather events. However, suppose a methodology is
74    developed to identify past cases with high-impact weather separately for multiple regions
75    across the CONUS. We would anticipate that high-impact weather in one region would
76    coincide with more commonplace weather in other regions, thereby avoiding under-sampling
77    of more commonplace events when forming the overall sample. Thus, reforecasts conducted
78    from a union of the identified dates, we hypothesize, should be adequate for training of both
79    common and uncommon weather-forecast post-processing.

80        In subsampling, and thereby reducing, the number of historical dates on which
81    reforecasting is conducted, the "thinned" reforecasts must facilitate end-user applications, such
82    as hydrologic forecasting, watch/warning operations, and decision support. Here, it is
83    important to establish that the reforecast sample size can be reduced, materially (i.e., saving
84    meaningful computational resources), without an unacceptably negative impact on the quality
85    of the hydrologic forecasts and associated decision support, particularly for large and extreme
86    events. The NWS Office of Water Prediction (OWP) currently uses and plans to use
87    meteorological reforecasts for a wide variety of hydrologic modeling applications. For
88    example, the Hydrologic Ensemble Forecast Service (HEFS: Demargne et al. 2014) is used by
89    the thirteen River Forecast Centers (RFCs) of the NWS to produce reliable and skillful
90    hydrologic forecasts for, among other things, informing flood forecasting operations and
91    managing water resources. The HEFS ingests weather and climate forecasts from the various
92    meteorological models, including the Global Ensemble Forecast System (GEFS: Guan et al.
93    2021; Hamill et al. 2021; Zhou et al. 2021), and produces ensemble streamflow forecasts for
94    the short to the long range. The HEFS depends on a large sample of meteorological reforecasts
95    to: 1) downscale and bias-correct the precipitation and temperature forecasts used in the
96    hydrologic models; 2) validate the HEFS, particularly for large and extreme events; and 3)

File generated with AMS Word template 1.0

97 support myriad decision support applications and end-users, such as the New York City
98 Department for Environmental Protection (NYCDEP), who require hydrologic (re)forecasts to
99 help manage the NYC water supply.

100     Recent work by OWP suggests that the sensitivity of the HEFS to reforecast sample size
101 originates primarily from the need to validate the HEFS and provide guidance for large and
102 extreme events (refs?). This is not surprising, because the statistical modeling used in the HEFS
103 is relatively parsimonious, whereas decision makers are particularly interested in the accuracy
104 of the HEFS for large and extreme events. In order to demonstrate that a "thinned'
105 meteorological reforecast can adequately support validation and decision support with the
106 HEFS, it is important to conduct hydrologic reforecasting, both with and without data thinning,
107 and demonstrate that: (a) The HEFS can be calibrated using a thinned sample *without* an
108 unacceptable decline in forecast quality (e.g. without residual biases from under- or over-
109 sampling large and extreme events), as demonstrated through statistical validation, and; (b) any
110 increase in validation sampling uncertainty does not materially impact the ability of OWP to
111 guide strategic investments in the HEFS or to support decision makers in using historical
112 (validation) information, particularly for large and extreme events.

113     The methodologies described below should estimate probabilities of large and extreme
114 events (cases) across the US, but the underlying methodology may estimate probabilities for
115 subdomains of the US and then combine them. In this study, we will evaluate the importance
116 of a case based on the forecasts of precipitation exclusively. While hydrologic predictions can
117 be sensitive to other weather variables such as temperature and melting level, these are likely
118 to be second-order effects which will be ignored here to generate a benchmark solution.
119 Furthermore, this paper will only deal with the construction of an optimal thinned sample based
120 solely on the meteorological information; the actual hydrologic forecasting and validation will
121 be reported on in a future companion publication.

122     The rest of the paper is organized as follows. Section 2 provides scientific background that
123 illustrates our thought process in developing a novel statistical methodology to model
124 precipitation and introduces our proposed case selection techniques. These methodologies are
125 described in detail and their performance is evaluated in sections 3 and 4, respectively. Section
126 5 contains a summary of the paper, as well as some discussion and outlook. Some of the more

5

127 technical figures are placed in the Supplemental Information, which also includes a link to the
128 data sets used or generated in this study.

129

## 2. Background and proposed methodologies

*a. Statistical downscaling and prediction of precipitation*

132     Statistical prediction (and downscaling) methods for precipitation are based on the
133 (extensively studied) association between extreme precipitation and recurrent large-scale
134 meteorological patterns (LSMP), which provide favorable environment for smaller-scale
135 processes often underlying the extreme precipitation events (although not all such events are
136 tied to LSMP). Barlow et al. (2019) reviewed, among other things, the types of meteorological
137 synoptic systems and mechanisms for extreme precipitation LSMPs for the North America
138 region and found a great diversity of LSMPs depending on the geographical location and
139 season. LSMPs are distinct from teleconnection patterns in that the LSMPs are conditioned on
140 the occurrence of a specific event (here, extreme precipitation), whereas classical
141 teleconnections are not. The most intuitive way of defining the LSMP is through compositing,
142 although a variety of other methods are available, including regression-based and cluster-
143 analysis methods (Grotjahn et al. 2016). For example, Robertson et al. (2016) used *K*-means
144 cluster analysis (Robertson and Ghil 1999) of the reanalysis wind data over North America to
145 identify seven distinct large-scale circulation types and tie some of them to enhanced
146 probability of springtime flooding events in the Midwest of the US. We note here that while
147 identifying a small subset of large-scale recurrent patterns — independent of precipitation —
148 to classify weather states is an attractive methodology, it is apparently at odds with the extreme-
149 precipitation LSMPs' diversity mentioned above; hence, the practical utility of such
150 methodologies to downscale precipitation is likely to be quite limited.

151     Classical regression approaches such as canonical correlation analysis (CCA: Wilks 2011)
152 also have a limited applicability to short-term precipitation modeling due to non-Gaussian and
153 intermittent nature of precipitation; however, they may be suitable and have been utilized for
154 the prediction of *seasonal* rainfall both directly (Sinha et al. 2013) and as an auxiliary tool for
155 selecting external predictors in conjunction with alternative methodologies (Holsclaw et al.
156 2016). The most widely used class of the latter alternative methods for statistical modeling,

6

157 downscaling and prediction of precipitation involves, in one way or another, generalized linear

158 models (GLM: McCullagh and Nelder 1989) — an extension of classical linear regression

159 models to simulate the (conditional) expectation of a non-Guassian distributed variable (such

160 as precipitation) as a function of external predictors (exogenous variables) associated with non-

161 stationary forcing (seasonal, anthropogenic or otherwise related to the climate variability

162 external to the climate sub-system of interest) or, of most relevance to the present discussion,

163 with the occurrence of LSMPs. These models are typically constructed to estimate probability

164 of daily precipitation at a grid point (or weather station) level (for example, Furrer and Katz

165 2007), although some generalizations to multiple stations accounting for spatial correlations

166 between them are also available (Kenabatho et al. 2012). Manzanas et al. (2018) fitted separate

167 GLM models to downscale daily precipitation occurrence and, separately, daily precipitation

168 amount at each grid cell using upper-air predictors simulated by multi-model seasonal climate

169 hindcasts over the Philippines. They showed that this methodology can yield a significant

170 forecast skill improvement for seasonal precipitation prediction over that of raw forecasts in

171 cases where the dynamical model predicts large-scale exogenous variables better than it

172 predicts the precipitation itself.

173 An alternative approach to precipitation modeling over a spatially extended array of grid

174 points or stations — a Hidden Markov model approach — assumes the existence of a few

175 discrete "hidden" weather states that capture spatial dependencies of rainfall probabilities

176 within the region considered, with Markovian daily transitions between these states tied to

177 exogenous predictors via GLM regression; in the latter case these models are referred to as

178 non-homogeneous Markov models: NHMM (Robertson et al. 2004). Holsclaw et al. (2016)

179 developed a combined HMM-GLM approach, in which a weather state HMM model is

180 complemented by a GLM model that can modify individual (hidden) states at a station level in

181 response to external predictors (rather than the probabilities of transitions between fixed states,

182 as in NHMM). We speculate that this approach would also be challenging to adapt for faithful

183 modeling of extreme precipitation over the entire CONUS, where, once again, the heaviest tails

184 of local precipitation distributions are associated with a multitude of precipitation producing

185 systems (Barlow et al. 2019), rather than with a small number of weather states and/or

186 exogenous predictors.

187 *b. Present approaches*

7

188     To summarize the above discussion, neither classical linear regression-based methods nor
189    clustering or HMM methods are directly suitable for statistical modeling and prediction of
190    precipitation over the entirety of CONUS due to non-Gaussian and intermittent nature of
191    precipitation and a great diversity of precipitation-producing systems/mechanisms in this
192    region, respectively. GLM regression methods may work at a grid-point level but will still
193    require the choice of exogenous dynamical variables based on a subjective zoning of the area;
194    these methods are also incompatible with automated linear regularization and predictor-
195    selection techniques such as CCA or (closely related) partial least squares methods (PLS: Wold
196    et al. 1984).

197     Here we address these difficulties via a new methodology based on statistical modeling of
198    the so-called pseudo-precipitation field, which uses column integrated water vapor saturation
199    deficit as a negative complement to precipitation (Yuan et al. 2019). Pseudo-precipitation is
200    thus characterized by a more symmetric distribution than the actual precipitation, opening up
201    a possibility of utilizing standard linear regression methods for its modeling. Furthermore, in
202    contrast to classical precipitation field, pseudo-precipitation patterns provide, additionally,
203    information on both the synoptic-scale and anisotropic mesoscale environment (including
204    LSMPs) in which local precipitation occurs, making it ideally suited for linear inverse
205    modeling (LIM: Penland 1986; Penland and Sardeshmikh 1995) and related data-driven
206    modeling methodologies (Kravtsov et al. 2005, 2009, 2016, 2017). The LIMs exhibit sub-
207    seasonal forecasts skill comparable to that of state-of-the-art numerical weather prediction
208    (NWP) models (see, for example, Winkler et al. 2001) and, most importantly, are able to isolate
209    initial states associated with useful predictability of its own, as well as of NWP-model based
210    forecasts (Newman et al. 2003; Albers and Newman 2019). This property can be helpful for
211    identifying potentially predictable high-impact precipitation events — the main focus of the
212    present study. The proof-of-concept mesoscale-resolving regional inverse models of surface
213    temperature over CONUS have been developed and tested before (Kravtsov et al. 2017); these
214    models are complex enough (yet numerically efficient) to provide an overarching description
215    and forecast utilization of LSMPs associated with local weather extremes. We expect the same
216    statement to be true for the combined surface temperature/pseudo-precipitation modeling we
217    propose here.

File generated with AMS Word template 1.0

218    In addition to the above (main) purely statistical and numerically efficient methodology,

219    we will also develop and test a procedure for selecting an optimal thinned subsample of

220    representative dates by utilizing the GEFSv12 reforecasts of precipitation for the 2000–2019

221    period. This procedure would allow one to conduct a (greatly) reduced number of hydrologic

222    hindcasts to estimate the adequacy of the reduced sample for the post-processing, validation

223    and end-user needs. However, it is much more computationally demanding than the proposed

224    purely data-driven methodology insofar as it still requires, in the first place, the full-blown

225    meteorological reforecasts of the entire climate state to determine the thinned subsample, which

226    somewhat defies the purpose of data thinning. Full, every-day reforecasts were available for

227    the GEFS versions 10 (Hamill et al. 2013) and 12 (Guan et al. 2021), but such full records may

228    not be available in the future to be subsampled. Yet, the present dynamical/statistical *ad-hoc*

229    algorithm based on the GEFSv12 reforecasts can be considered a control against which to

230    evaluate our main statistical modeling methodology, and, in what follows, we describe this

231    algorithm first.

232

## 3. Data sets and methodological details

233

*a. Selecting reforecast case dates based on heavy precipitation in GEFSv12 reforecasts*

234

235    We argue here that a metric of an event's extremeness should be based on precipitation

236    magnitude as opposed to, say, the quantile of today's forecast relative to its climatological

237    distribution (for example, a 0.1-inch forecast in the desert may be an extreme event relative to

238    the local climatology but still of marginal significance to hydrologic applications). In the

239    present methodology, the importance of a case for potential selection was judged based on the

240    0–10-day total GEFSv12 ensemble-mean reforecast precipitation $P_{10}$, sampled daily over the

241    2000–2019 period for each of the 18 CONUS regions associated with distinct 2-digit

242    Hydrologic Unit Codes (HUC-2 units: https://nas.er.usgs.gov/hucs.aspx). Some case choices

243    were based on large ensemble-mean precipitation averaged over the entire HUC-2 unit, while

244    others were optimized on the top 20% of grid points inside that HUC-2 unit (at the 0.25º

245    resolution)  to emphasize smaller-scale impactful events.  A small number of cases were also

246    based  on  large  CONUS-wide  ensemble-mean  precipitation.  More  specifically,  the

247    subjectively chosen breakdown of cases was as follows:

248    (a) 30% of the total cases were optimized based on the maximum 10-day ensemble mean
249 precipitation in that HUC-2 unit. After choosing a case day on this criterion, an *ad-hoc* de-
250 weighting of the day before and the day after was applied so they are less likely to be
251 chosen. However, we find that the algorithm often chooses case days separated by at least 2
252 days (which can be easily adjusted if desired).

253    (b) 60% of the total cases are optimized based on the maximum 10-day ensemble-mean
254 precipitation at the 20 grid points within that HUC-2 that have the largest mean precipitation.

255    (c) The remaining 10% are chosen based on maximal CONUS-averaged ensemble-mean
256 precipitation.

257    In developing the above merged set of dates from across the subdomains, we chose the first
258 case date from each subdomain unless it was a repeat. Then we proceeded to the second ordered
259 case date in each subdomain, the third, and so forth, until we have reached $n$ total cases, where
260 $n$ is an adjustable pre-determined size of the thinned sample. The lists of presumed important
261 cases were developed separately for the warm (April–September) and cool season (October–
262 March), with $n = 520$.

263    The resulting procedure produces a list of dates with an irregular sampling in time, which
264 is to be expected if there exist long periods with no hydrologically significant activity
265 (assuming the GEFSv12 mean precipitation to be a reasonable proxy for such an activity)
266 which the algorithm aims to skip to provide more samples when there is strong forcing. The
267 clustering around the largest storms from multiple initial conditions/issued datetimes is
268 controlled, to an extent, by our de-weighting procedure, which involves a trade-off: on the one
269 hand, we don't want a lot of shared information between samples; on the other hand, we do
270 want to sample the largest events from several issued datetimes (and, hence, lead durations).
271 Other adjustable parameters include the total number of cases $n$ and the proportions of cases
272 associated with each of the case categories (a, b, c) above.

273    We will hereafter refer to the thinned sample produced by the above procedure as sample$_A$;
274 illustrative examples from this sample will be presented alongside with the results from our
275 alternative, purely data-driven methodology presented below.

276

*b. Selecting reforecast cases using EMR (Empirical Model Reduction) statistical model*

278    1) DATA SETS AND VARIABLES: INTRODUCING PSEUDO-PRECIPITATION

279    We analyzed data from the National Center for Environmental Prediction North American
280    Regional Reanalysis (NARR) (http://www.esrl.noaa.gov/psd/data/gridded/data.narr.html);
281    Messinger et al. (2006), using daily "observations" on a 349×277 grid with nominal horizontal
282    resolution of 32 km and 29 pressure levels, over the 1979–2020 period; about a third of these
283    data are from locations over land, leading to ~30000 data points in each of the ~365 (days per
284    year) ×42 years~15000 maps for a single-level field. The NARR data set has been widely used
285    in the climate downscaling community (see Zobel et al. 2018 and references therein). Bukovsky
286    and Karoly (2007) found that NARR provides faithful estimates of the observed precipitation
287    over CONUS, although some biases exist over Canada due to a relatively poor quality of the
288    assimilated data there.

289    We utilized NARR data sets for the (daily) accumulated precipitation $Pr$ and 2-m air
290    temperature $T_a$. We also used the air temperature $T$ and specific humidity $Q$ data at all available
291    pressure levels to compute the *air dryness D* related to the column-integrated water-vapor
292    saturation deficit (Yuan et al. 2019). In an air column of area $\delta A$, the mass of water vapor $\delta m$
293    to be added to achieve saturation throughout the column is

$$\delta m = -\delta A \int (\rho_v - \rho_{v,s})dz = \delta A \int (\rho_v - \rho_{v,s})\frac{dp}{\rho g} = \frac{\delta A}{g} \int (Q - Q_s)dp. \tag{1}$$

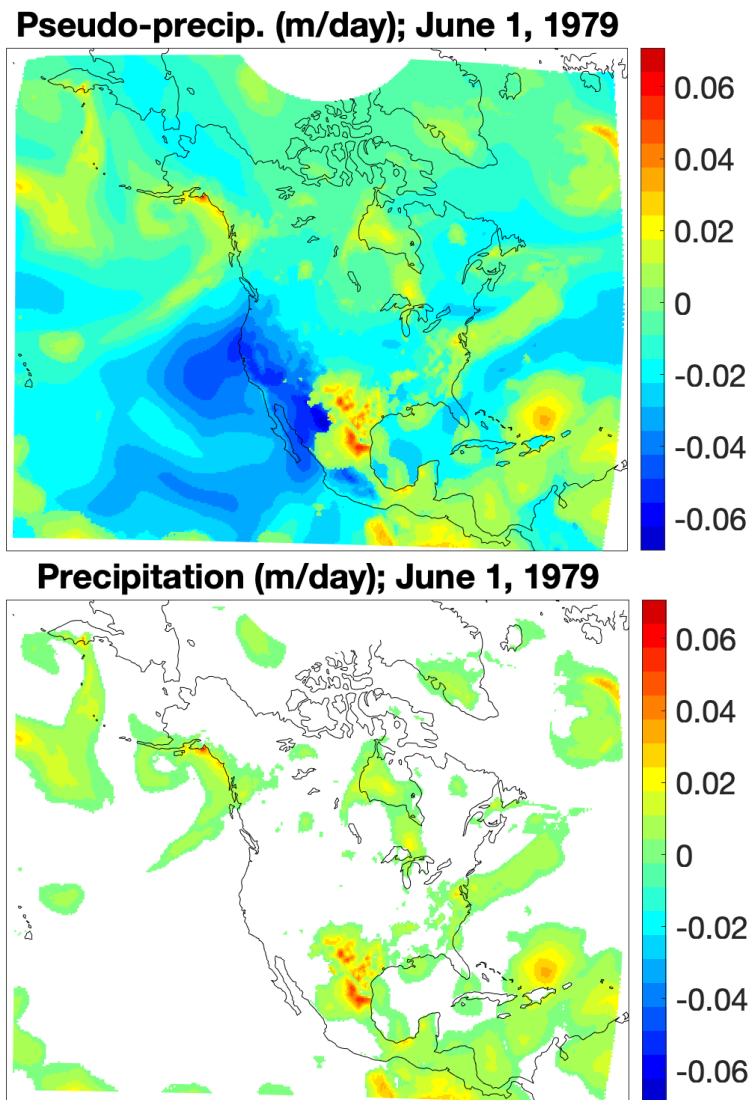294    Here $z$ is the geometric height, $p$ is the pressure, $\rho$ and $\rho_v$ are the dry-air and water-vapor
295    densities, respectively, the subscript $s$ denotes the quantities for saturated air and $g = 9.82$ m
296    s$^{-2}$ is the gravity acceleration. The specific humidity of saturated air $Q_s$ can be computed as
297    (Bolton 1980):

$$Q_s = \varepsilon \frac{e_s}{p}; \quad e_s = 6.112 \exp\left(\frac{17.67\,T}{T + 243.5}\right), \tag{2}$$

298    where $\varepsilon = 0.62198$ is the ratio of the molecular weights of water and dry air, $e_s$ is the
299    saturation water-vapor pressure, and air temperature $T$ is expressed in ºC. Air dryness $D$ is
300    defined as the equivalent water depth associated with the quantity $\delta m$ in (1):

File generated with AMS Word template 1.0

$$D = -\frac{\delta m}{\rho_w \delta A} = -\frac{1}{\rho_w g}\int (Q - Q_s)dp, \tag{3}$$

301    where $\rho_w = 1000$ kg m$^{-3}$ is water density. The air dryness in (3) can be thought of as a negative

302    complement to precipitation and used to construct the so-called pseudo-precipitation field $PP$,

303    which is, here, equal to the actual precipitation $Pr$ if $Pr > 0.001$ m day$^{-1}$ or to $Pr + D$

304    (essentially, the air dryness $D$) otherwise.



305

306

307    **Figure 1**: Pseudo-precipitation ($PP$) (top), as well as precipitation $Pr$ (bottom) on June 1, 1979, derived

308    from the NARR reanalysis (m). White areas in the bottom plot are either outside of the NARR domain

309    or, otherwise, have zero $Pr$.

File generated with AMS Word template 1.0

310

The *PP* field incorporates the information about both precipitation, which can exhibit small-scale intermittent structures, and multi-scale synoptic environment (see **Fig. 1**); it thus provides a promising, yet unexplored way to characterize and predict, statistically, wet and dry weather conditions. One of its attractive features is that the distribution of *PP*, unlike that of *Pr*, is a single-mode, two-tailed distribution, which makes *PP* more similar to other dynamical and thermodynamic variables describing atmospheric state. *This opens up a possibility for using standard methodologies developed previously for temperature and flow-field analysis and modeling (CCA, LIMs) to analyze and model pseudo-precipitation and, hence, its positive part associated with the actual precipitation.*

2) EMR MODELING OF PRECIPITATION

We here apply advanced methods for high-dimensional statistical data modeling to identify potentially predictable large/extreme precipitation events. This idea is rooted in the demonstrated ability of a sub-class of such inverse models — LIM models (section 2b) to "forecast the forecast skill" (Albers and Newman 2019).

*(i) General methodology*

The Empirical Model Reduction (EMR: Kravtsov et al. 2005, 2009, 2016, 2017) is a generalization of LIM data modeling methodology to incorporate memory effects in the postulated parametric form of this empirical model's evolution operator. The model construction usually takes place in a reduced phase space (for example, the space associated with $L$ leading Empirical Orthogonal Functions (EOFs) of the field(s) simulated, in which case the state of the system on a given day is described by the $L$-valued vector of PCs $\mathbf{x}$. The EMR emulator models the evolution of PCs using the following multi-level form (three levels are shown below):

$$d\mathbf{x} = \mathbf{x} \cdot \mathbf{A}^{(1)} + \mathbf{r}^{(1)},$$

$$d\mathbf{r}^{(1)} = \left[\mathbf{r}^{(1)}\ \mathbf{x}\right] \cdot \mathbf{A}^{(2)} + \mathbf{r}^{(2)}, \tag{4}$$

$$d\mathbf{r}^{(2)} = \left[\mathbf{r}^{(2)}\ \mathbf{r}^{(1)}\ \mathbf{x}\right] \cdot \mathbf{A}^{(3)} + \mathbf{r}^{(3)},$$

13

where the differentials on the left-hand side denote the daily increments of the corresponding variables. The first model level in isolation, with the residual $\mathbf{r}^{(1)}$ represented, at the simulation stage (see below), by the spatially correlated white noise, would make up a classical LIM model (for example, its 1-D analog would be the AR-1 red-noise model widely used to test for statistical significance of spectral peaks in a time series). Instead, in the EMR modeling, daily increments of the first-level residual $d\mathbf{r}^{(1)}$ are in turn modeled as a linear function of the extended predictor vector $\left[\mathbf{r}^{(1)}\ \mathbf{x}\right]$ to form the second level of the multi-level regression model (4). In the same way, the third level connects the daily increments of the second-level residual $d\mathbf{r}^{(2)}$ and the extended predictor vector $\left[\mathbf{r}^{(2)}\ \mathbf{r}^{(1)}\ \mathbf{x}\right]$ involving the variables from the previous two model levels.

The matrices of the model coefficients $\mathbf{A}$ and the level residuals are found by a regularized multiple linear regression (MLR) and depend on the seasonal cycle at the monthly resolution. While the residuals of the first and second level may involve serial correlations, the last level's residual $\mathbf{r}^{(3)}$ is typically white in time (otherwise, additional levels can be added). Note that while the model construction procedure is sequential from the first level down to the last level, the equations (4) — when rewritten as one equation containing the time-lagged variables — are formally equivalent to the autoregressive moving average model (ARMA: Box et al. 1994).

The model (4) can provide independent realizations of observations that are statistically very similar to the input data. At this stage of model simulation, the residual forcing at the third model level $\mathbf{r}^{(3)}$ is replaced by a random forcing, which can involve simultaneous or lagged spatial correlations between different PC "channels" and depend on the simulated state $\mathbf{x}$ (effectively *making the model nonlinear*). One can also use the EMR model for statistical forecasting of the out-of-sample data. Trivial linear transformation of the simulated PCs provides the data simulation or forecasts in the original physical space.

While the original LIM models, as well as the EMR methodology above, have been typically applied to fairly low-dimensional subsets of meteorological data, Kravtsov et al. (2015, 2017) demonstrated its applicability to larger or higher-resolution data sets such as regional surface temperature (Kravtsov et al. 2017) and precipitation. In the latter case, most relevant to the present project, the EMR modeling of combined $T_a$ and $PP$ fields resulting from an hourly, 16-km-resolution Japan regional reanalysis was successfully used by AIR

File generated with AMS Word template 1.0

364    Worldwide (Boston, MA) for flood-risk assessment over Japan (Boyko Dodov 2016, Director

365    of Flood Modeling, personal communication). We here build an analogous combined $T_a/PP$

366    daily EMR model over CONUS and utilize it to identify potentially predictable large and

367    extreme precipitation events to be included in the final thinned subsample.

368    *(ii) EMR application to NARR $T_a/PP$ data*

369       All model construction steps, including the identification of seasonal cycle and initial data

370    compression, were done using the NARR's 1979–1999 (training period) data. We built our

371    EMR model (4) in the phase space of 3000 common EOFs of the daily 2-m air temperature and

372    pseudo-precipitation (section 3b.1) anomalies with respect to the mean seasonal cycle

373    computed by the linear regression of raw daily data onto the first five harmonics of the annual

374    cycle. The maps of climatological standard deviation of these anomalies (over the 1979–1999

375    period) are shown in Supplemental **Fig. S1**. The EOF identification only used land grid points

376    (hence, the assessment of model performance should in principle also focus on the land region).

377    We first computed 1000 leading EOFs of $T_a$ and 3000 leading EOFs of $PP$ field, normalized

378    the corresponding individual PCs by the standard deviation of the leading PC of each field and

379    applied an additional EOF rotation to the data set of concatenated $T_a/PP$ individual normalized

380    PCs, finally retaining the leading 3000 common PCs so obtained. These PCs were again

381    normalized by the standard deviation of their own PC-1, while the corresponding dimensional

382    EOF patterns were found by regressing the individual fields onto these common PCs (note that

383    these patterns only represent the actual common EOFs over the land region and should be

384    interpreted as a teleconnection pattern over ocean). To initialize model forecasts performed

385    over the validation period (2000–2020), we projected the anomaly data there (again, with

386    respect to the 1979–1999 mean seasonal cycle) onto common $T_a/PP$ EOFs computed above.

387    For the back transformation, to produce the patterns in physical space from a map of individual-

388    day PC loadings (as obtained, for example, from our EMR model simulations), one is to simply

389    add all of the 3000 individual EOF patterns multiplied by the corresponding loadings, on top

390    of the mean seasonal cycle. The EOF truncation errors associated with the procedure above are

391    shown in the supplemental **Figs. S2** (training period) and **S3** (validation interval) and

392    demonstrate a fairly high accuracy (small errors) over CONUS for both $T_a$ and $PP$ data,

393    sufficient for a faithful representation of extreme hydroclimatic events in the region.

394　　　　The EMR model construction and simulation technical steps follow Kravtsov et al. (2017),

395　　except here we are only modeling the evolution of daily fields and thus disregard the sub-daily

396　　and monthly model tiers employed there. Note that all of the model operators in (4) are season-

397　　dependent at monthly resolution. For example, to estimate the model parameters for January,

398　　we consider the December–January–February (DJF) subset of daily data and use a regularized

399　　(PLS) version of multiple linear regression for each of the three model levels sequentially. At

400　　the simulation stage, the third-level residual $\mathbf{r}^{(3)}$ is simulated by pulling its randomized 5-day

401　　snippets from the library of actual residuals obtained during the model construction stage. This

402　　random forcing selection is also season-dependent, so that, for example, if the current time

403　　step is in January, the DJF subset of $\mathbf{r}^{(3)}$ library is used for that purpose. To avoid unnecessary

404　　discontinuities, the consecutive random forcing snippets were overlapped by two days and

405　　added with the weights $(\sqrt{3}/2, 1/2)$ and $(1/2, \sqrt{3}/2)$ before phasing out the previous snippet

406　　of $\mathbf{r}^{(3)}$ completely.

407　　　　We used the EMR model above in two ways: first to produce, from random initial

408　　conditions, 100 synthetic realizations of the 2-m air temperature and precipitation (positive

409　　pseudo-precipitation) 1979–1999 evolution and assess how well the model captures the

410　　observed statistical characteristics of these fields (section 4a). Second, we ran 0–10-day 100-

411　　member ensemble forecast of temperature and (pseudo) precipitation for each of the 2000–

412　　2020 initial conditions to assess the model's predictive skill (section 4b) and eventually utilized

413　　these forecasts to develop and test an innovative methodology for reforecast thinning (section

414　　4c). Since our interest here is in extreme precipitation events, we will focus below on the

415　　simulation of precipitation; the present EMR performance in modeling temperature will be

416　　considered elsewhere.

417　　3) CASE SELECTION USING EMR ENSEMBLE FORECASTS

418　　　　In principle, the EMR ensemble-mean hindcasts of the 0–10-day total precipitation $P_{10}$ can

419　　be processed in exactly the same way as the GEFSv12 reforecasts to produce an alternative

420　　representative subset of events of impact, as described in section 3a; the outcome of such a

421　　procedure, which results in the thinned sample we will refer to as sample$_B$, are briefly discussed

422　　at the very end of section 4c . However, a large size of the EMR hindcast ensemble (possible

423　　to achieve due to this model's numerical efficiency) makes it possible to develop an alternative

424　　methodology that involves relative entropy of the EMR hindcasts; this methodology will be

16

425    introduced below and described in detail in section 4c. We   will call the thinned sample

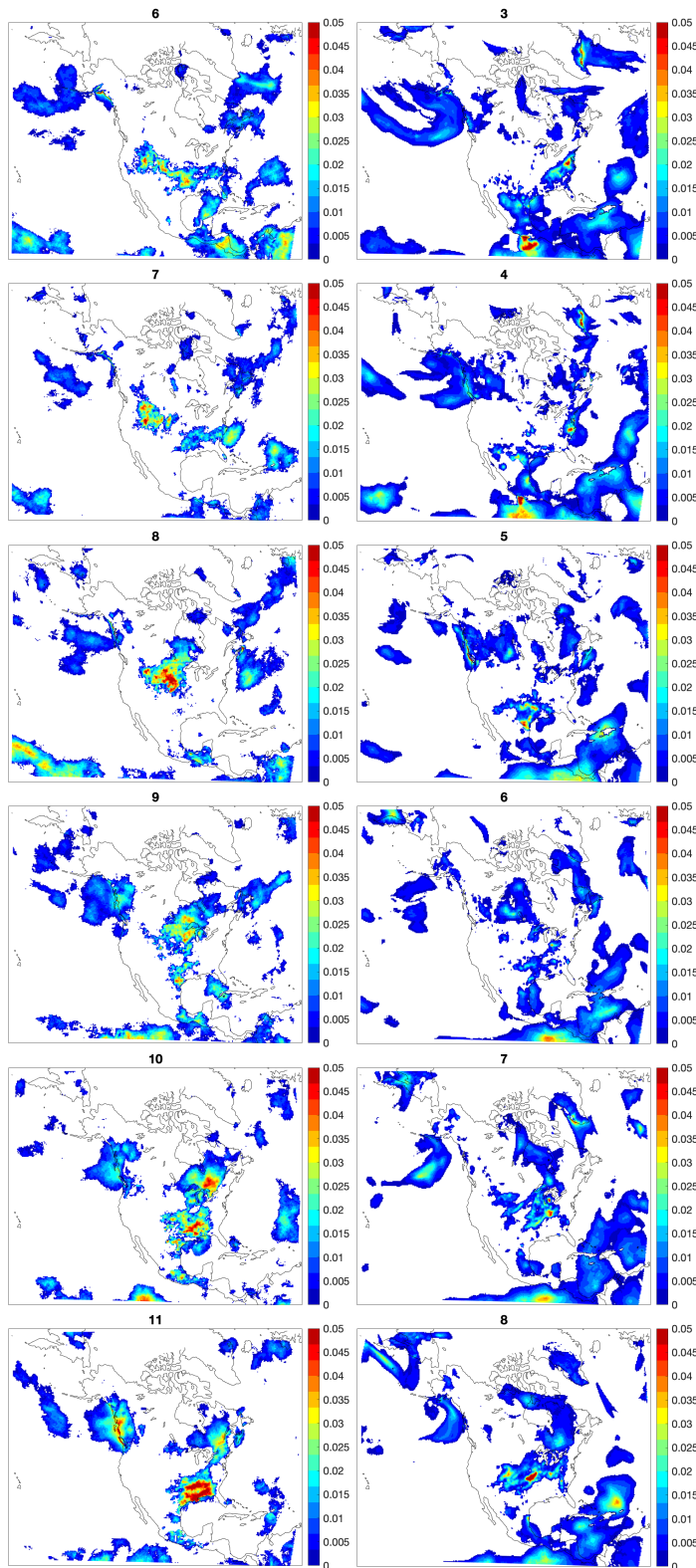426    produced by this EMR based method simply a "sample" or an "EMR-RE sample."



**Figure 2**: A JJA-season sequence of daily surface precipitation maps (m) from: (left) arbitrary [random] realization of EMR model; (right) NARR reanalysis. Day "1" in a panel caption would correspond to June 1, 1979. White areas in the bottom plot are either outside of the NARR domain or, otherwise, have zero *Pr*.
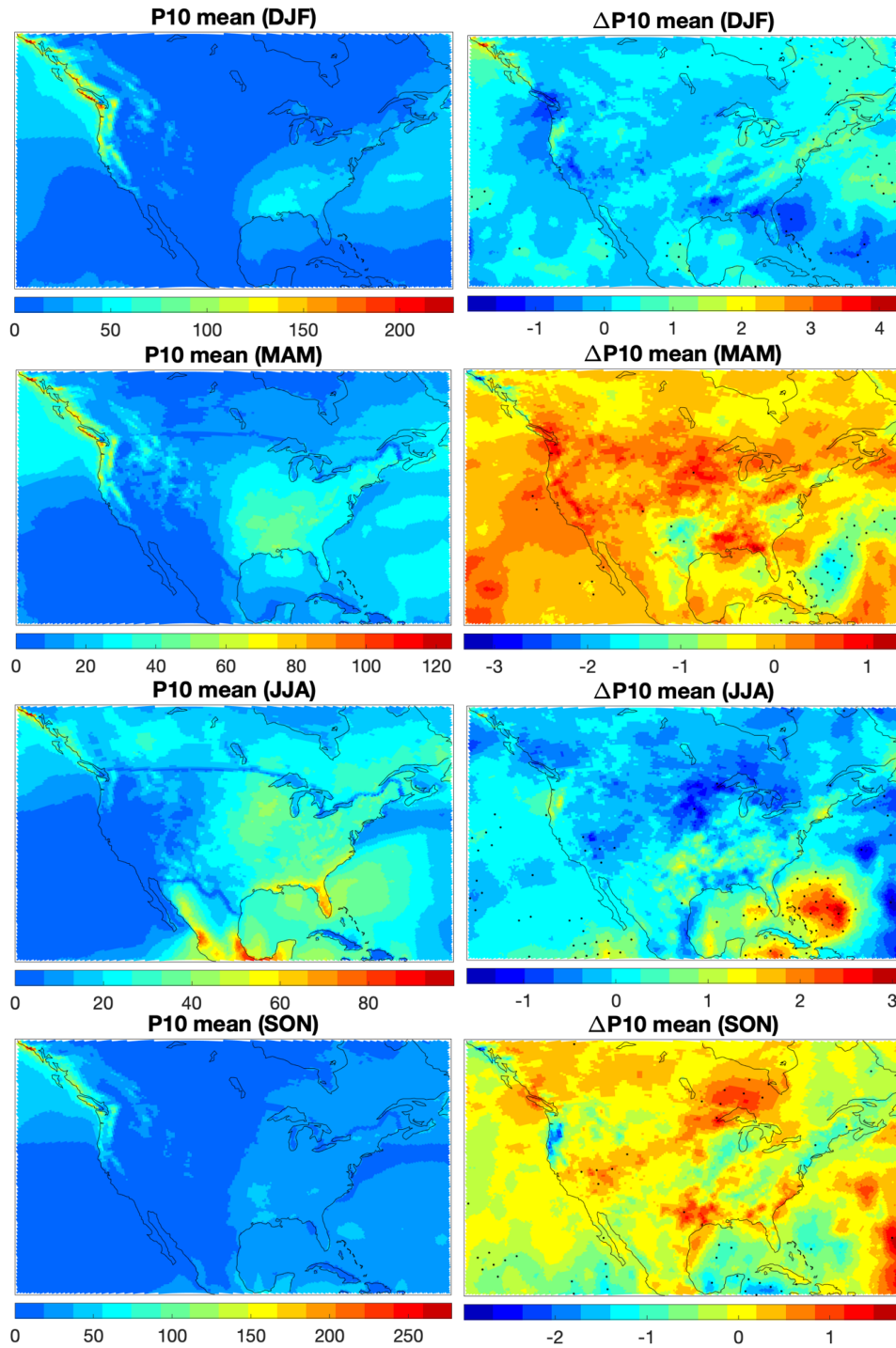
## 4. Results

*a. Using the EMR model as an emulator of daily precipitation evolution*

Preliminary inspection of the EMR-model daily precipitation simulations in physical space obtained by simply considering the positive pseudo-precipitation reveals model biases in the distribution of precipitation events (not shown). To eliminate these biases, we apply quantile mapping (for each of the DJF,  MAM, JJA, SON seasonal subsets) to each 1979–1999 model simulation of *pseudo-precipitation* to make the simulated local distributions of this quantity identical to those based on the original 1979–1999 NARR data. Specifically, the observed 1979–1999 and simulated 2000–2020 *PP* time series at a given grid point and for a given season (DJF, MAM,  JJA, SON) were sorted in the ascending order, upon which the sorted 2000–2020 simulated values were replaced by the sorted 1979–1999 observed values, then put back in the original order (cf. Hamill 2018). This procedure automatically ensures the identical local (i.e., a given grid point's) *precipitation* distributions between the model and NARR reanalysis as well. However, the spatiotemporal characteristics of sequences of daily precipitation maps are entirely due to dynamics embedded in the EMR model's propagator. Examples of such sequences for the warm and cold season are shown in **Fig. 2** and supplemental **Fig. S4**, respectively and give one a visual impression of how well the model matches the space–time structure of the observed stationary and propagating precipitation patterns; the external link to longer sequences is also available in the Supplemental Information.

We also compute,  for future use, daily time series of day 0–10 cumulative precipitation ($P_{10}$) and display its (seasonal) mean and 99th percentile in **Figs. 3** and **4**, respectively. Note that while the simulated local *daily* precipitation distributions are fixed due to quantile mapping, the simulated and observed distributions of $P_{10}$ can be different if the spatial scales or persistence/intermittency of the simulated precipitation differ from the observed characteristics. However, this does not seem to be the case here, with the simulated $P_{10}$ mean entirely consistent with observations (Fig. 3).  The simulated $P_{10}$'s  99th percentile (Fig.  4) is a slight overestimate compared to observations (including large areas over land), reflecting, perhaps, a slightly overly persistent local precipitation anomalies, but the overall match between the simulated and observed $P_{10}$ distributions is still very good.

**Figure 3**: 1979–1999 seasonal climatology of the 0–10-day total precipitation at the surface — $P_{10}$ (mm). Left: climatology based on an ensemble of 100 EMR model simulations; right: the difference between the simulated and NARR based $P_{10}$ climatology, with stippling indicating the regions over which this difference is of the same sign for more than 97 realizations (so, effectively, is statistically significant at the 5% level).

File generated with AMS Word template 1.0

463

**Figure 4**: The same as in Fig. 3, but for the 99$^{th}$ percentile of $P_{10}$.

465

466

467

468     *b. EMR model predictive skill*

469         To initialize the EMR model forecasts starting from a given day $n$ within the 2000–2020

470     validation interval, we assume that the observable state vectors $\mathbf{x}$ at days $n$, $n$–1 and $n$–2 are all

471     known. This, however, still requires us to solve for the values of the hidden-level variables $\mathbf{r}^{(1)}$

472     and $\mathbf{r}^{(2)}$ at the initial day $n$, which involves two pre-steps of the model (4) driven by a random

473     $\mathbf{r}^{(3)}$ forcing that ensure dynamical consistency [within the model (4)] of the hidden-state

474     variables with the observables $\mathbf{x}_n$, $\mathbf{x}_{n–1}$, $\mathbf{x}_{n–2}$. After these pre-steps, the model is integrated

475     forward in a normal way until the time $n$+10. This procedure is repeated for all of the available

476     initial conditions. Upon transformation back to physical space, the collection of *PP* forecasts

477     for a given lead time is, again, *quantile mapped* to the 1979–1999 local daily *PP* distributions;

478     finally, zeroing out the negative values of this quantile mapped *PP* forecast gives the final

479     forecast of the daily precipitation at this lead time, for each initial condition. Summing up the

480     precipitation forecasts for the days $n$ to $n$+10 makes up the final $P_{10}$ forecast for each initial

481     condition; we produced an ensemble of 100 such forecasts under different realizations of the

482     random forcing. Below we will focus on these $P_{10}$ forecasts when estimating the EMR model's

483     forecast skill.

484         We will also compare the EMR model forecasts with the benchmark damped persistence
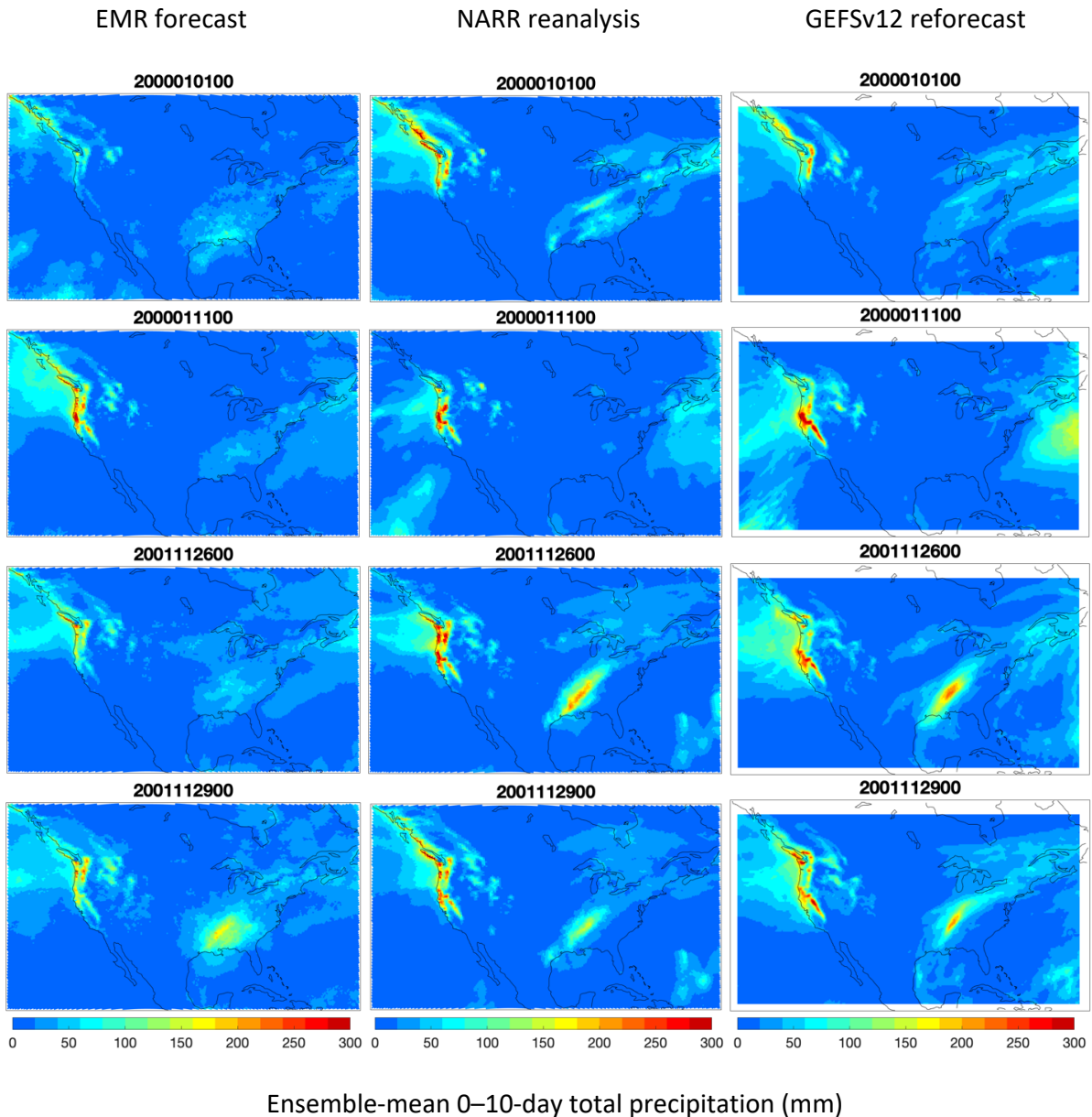
485     forecasts of daily precipitation:

486

487 $$p_{n+m} = r_m p_n + (1 - r_m)\bar{p}, \qquad (5)$$

488

489     where $r_m$ is the precipitation's lag-$m$ autocorrelation and $\bar{p}$ is the climatology, both computed

490     for each season's subset of the 1979–1999 NARR's daily precipitation data. The damped

491     persistence $P_{10}$ forecasts are obtained from (5) as the sum of $p_{n+m}$ for $m = \overline{0, 10}$.

492     1) DETERMINISTIC SKILL

493         We first discuss some traditional deterministic measures of skill by comparing the observed

494     $P_{10}$ values with their ensemble-mean EMR based prediction. **Figure 5** provides cool-season

495     examples of such a comparison for select cases of substantial observed $P_{10}$ episodes over

496     CONUS (see **Fig. S5** for analogous warm-season comparisons). Visual inspection confirms

497     reasonable EMR forecasts (left column) of the spatial scale, shape, location and magnitude of
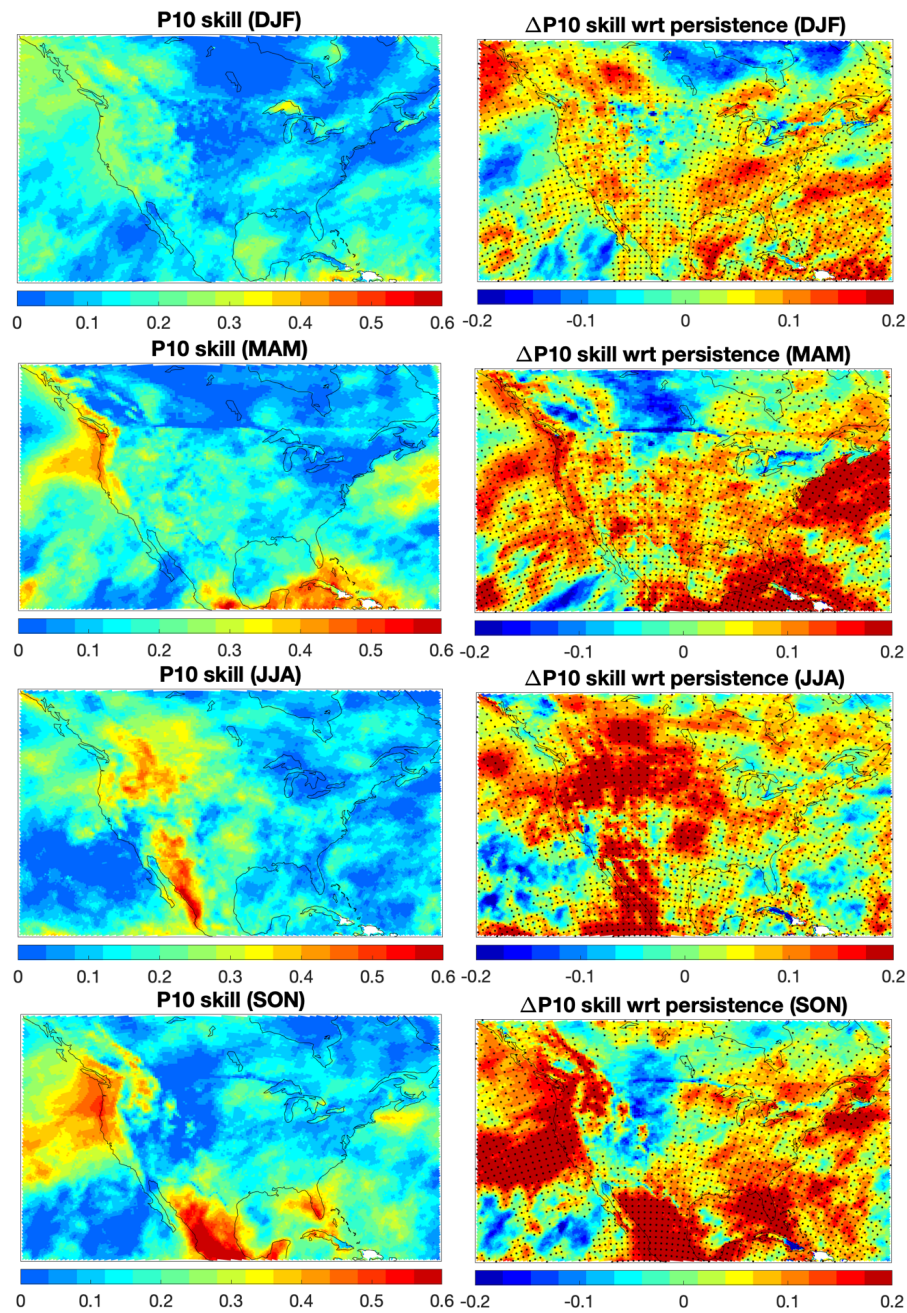
21

| EMR forecast | NARR reanalysis | GEFSv12 reforecast |
|---|---|---|



Ensemble-mean 0–10-day total precipitation (mm)

**Figure 5**: Examples of (cool season) $P_{10}$ forecasts using EMR model (left) and GEFSv12 system (right), along with the actual $P_{10}$ maps based on NARR reanalysis (middle). Units are mm. The forecast initialization time (the same across each row) is shown in panel captions in the YYYYMMDDHH format.
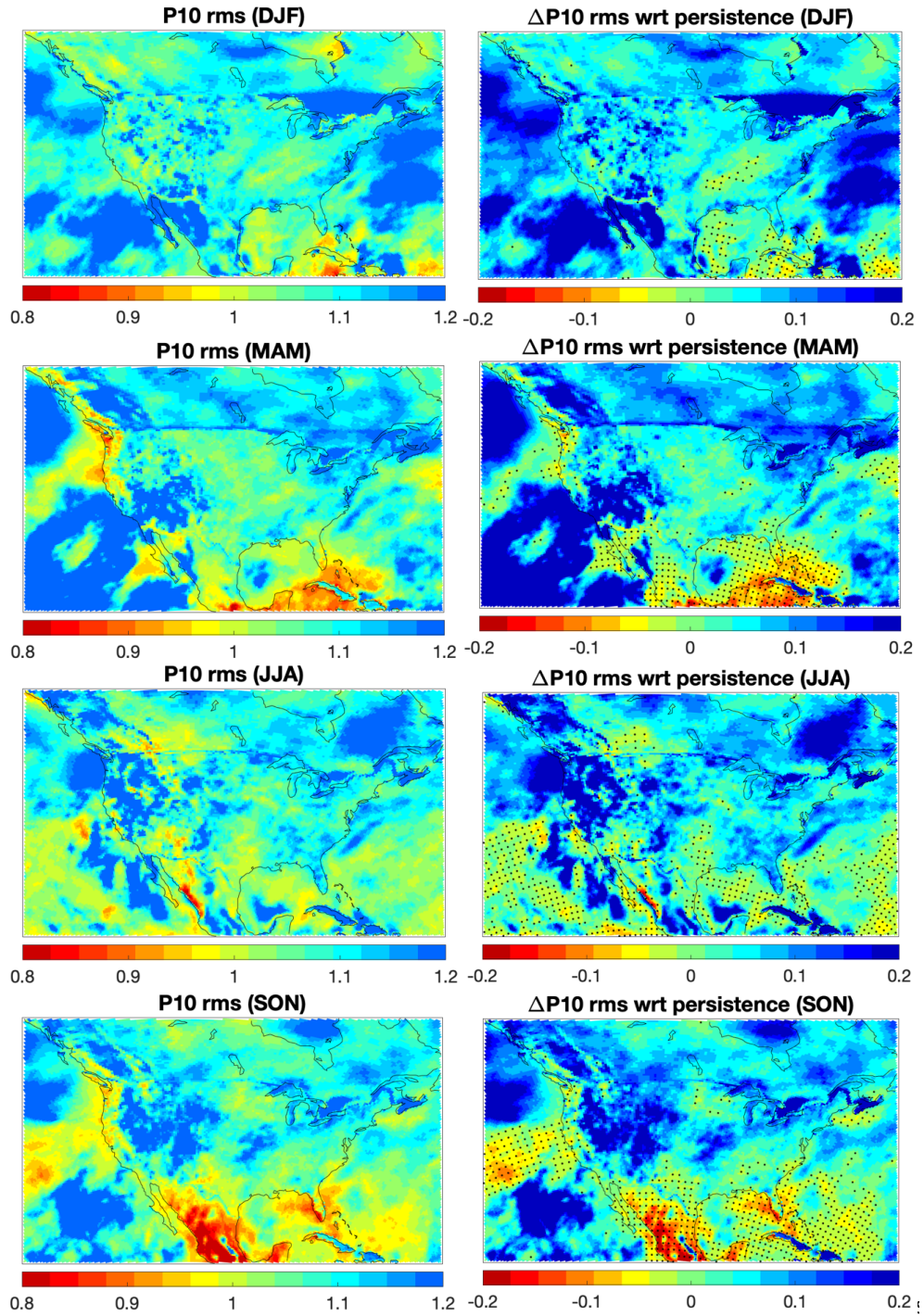
the observed large $P_{10}$ events (middle column), qualitatively similar to analogous GEFSv12 forecasts (right column). The overall correlations between the observed and forecasted $P_{10}$ time series (for each season) [**Fig. 6**, left], while positive, are fairly low, at the 0.2–0.3 level in most areas, with the exception of a few season-dependent regions reaching potentially useful levels

22

File generated with AMS Word template 1.0

508  of 0.5–0.6. However, these correlations are consistently higher than those for the damped-

509  persistence forecasts (Fig. 6, right).



510

**Figure 6**: The EMR model precipitation forecast skill. Left: Correlation between (1-day lead-time)
EMR forecast (ensemble-mean of 100 members) and daily $P_{10}$ time series from NARR reanalysis, for
each season. Right: The difference between forecast skill of the EMR model and (daily) damped
persistence forecast of $P_{10}$ (see text for details). Stippling indicates the areas of positive differences,
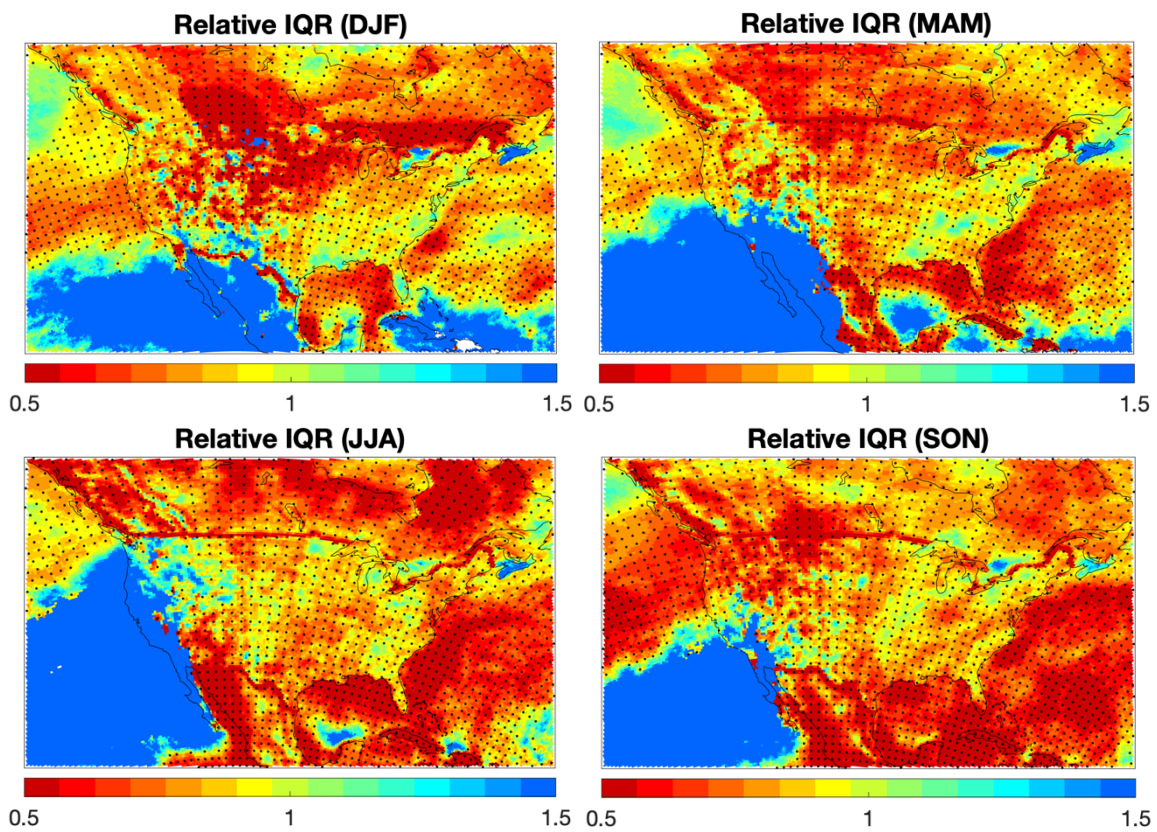where the EMR forecast beats the damped persistence forecast).

23

**Figure 7**: EMR model's $P_{10}$ forecast (2000–2020) root-mean-square (rms) error relative to (1979–1999) climatological standard deviations, for each season. Note the inverted color scale; otherwise, the same layout and conventions as in Fig. 6.
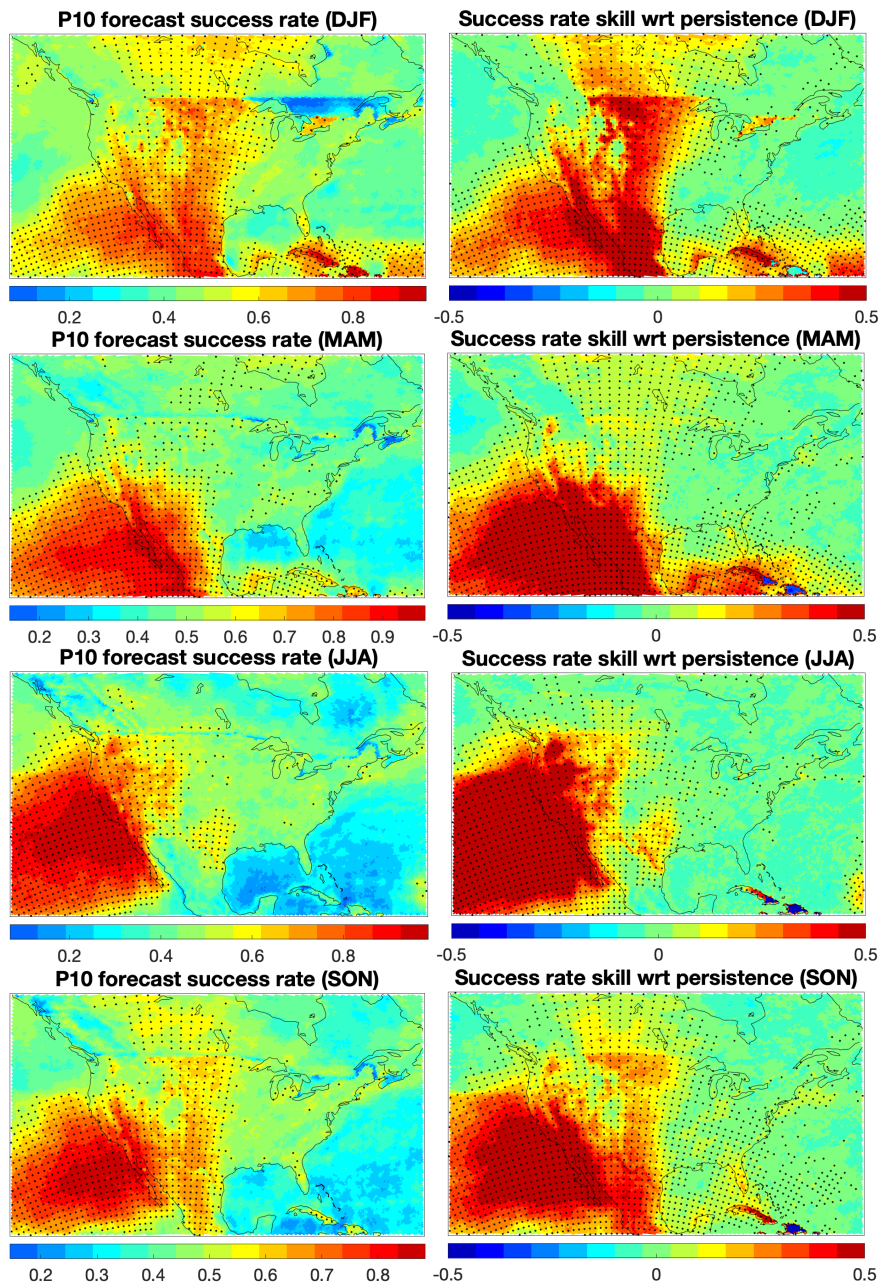
The root-mean-square (rms) distance between the observed and forecasted $P_{10}$ time series (**Fig. 7**) is generally close to the $P_{10}$'s climatological standard deviation, with EMR model forecasts beating damped persistence forecasts in some of the southern areas but performing similar to damped persistence forecasts elsewhere.

Overall, the deterministic measures of skill suggest, at best, a modest performance of the EMR model in forecasting $P_{10}$. This, however, may be in part due to unsuitability of these measures to describe the forecast quality of a discontinuous and highly intermittent — in space and time — state variable such as precipitation. In particular, considering the ensemble-mean forecast only completely disregards much of the useful information associated with the entire ensemble of forecasts.



**Figure 8**: The (average 2000–2020) EMR model's $P_{10}$ forecast interquartile range (IQR) — based on an ensemble of 100 forecasts — relative to the (1979–1999) climatological IQR of $P_{10}$, for each season. The ratios below unity (stippling) indicate an enhanced forecast utility relative to that of climatology forecast. Note the inverted color scale.

**Figure 9**: Left: The EMR forecast success rate defined as the fraction of $P_{10}$ forecasts (over all initial conditions, in each season separately) for which the actual $P_{10}$ value from the NARR reanalysis is within the IQR of (100-member) ensemble forecasts; stippling shows the areas with success rate exceeding the value of 0.5 (associated with the climatology forecast). Right: the difference between the EMR success rate and the success rate associated with the damped persistence forecast combined with the IQR of the EMR model (see text for details); stippling denotes the areas of positive differences (EMR model beats damped persistence forecast).

2) PROBABILISTIC CHARACTERISTICS OF SKILL

547       A perhaps more suitable measure of skill for precipitation should involve probabilistic

548   characteristics associated with ensemble forecasts of this quantity. An example of such a

549   measure is shown in **Fig. 8**, which plots the climatological ratio of the interquartile range (IQR)

550   of the EMR model forecasts to the climatological IQR of $P_{10}$. This quantity is related to the so-

551   called potential predictability (see Kleeman 2002 and references therein), with the values less

552   than 1 (this value corresponds to climatology forecast) and increasingly closer to zero

553   indicating a progressively more reliable forecast. Based on this measure, the EMR model

554   provides potentially useful forecasts throughout the region of interest, including CONUS.

555       While providing a measure of forecast utility, the potential predictability does not directly

556   compare the forecast with the actual observed precipitation value for the time of forecast. To

557   do so, we here introduce an additional forecast skill measure — the forecast success rate — by

558   counting the frequency of forecasts for which the observed $P_{10}$ value is within the IQR range

559   of the EMR forecast ensemble. The EMR model forecast success rate has large areas with

560   values exceeding 0.5 (the observed value of $P_{10}$ is within the IQR of EMR forecasts 50% of

561   the time or more) and sometimes nearing the value of 1  (**Fig.  9**, left). We also combined the

562   damped persistence forecasts of $P_{10}$ with the mean and IQR range of the corresponding EMR

563   forecast to compute the success rate associated with the damped persistence forecast: in

564   particular, the "range" associated with a damped persistence forecast $f_p$ was set to be $f_p -$

565   $\Delta_m, f_p + \Delta_p$, where $\Delta_m$ and $\Delta_p$ are the offsets  between the EMR model's ensemble mean  and

566   its 25th and 75th percentiles,  respectively. We verified that the damped persistence forecast

567   success rate defined in this way is substantially lower than the EMR model's success rate  (Fig.
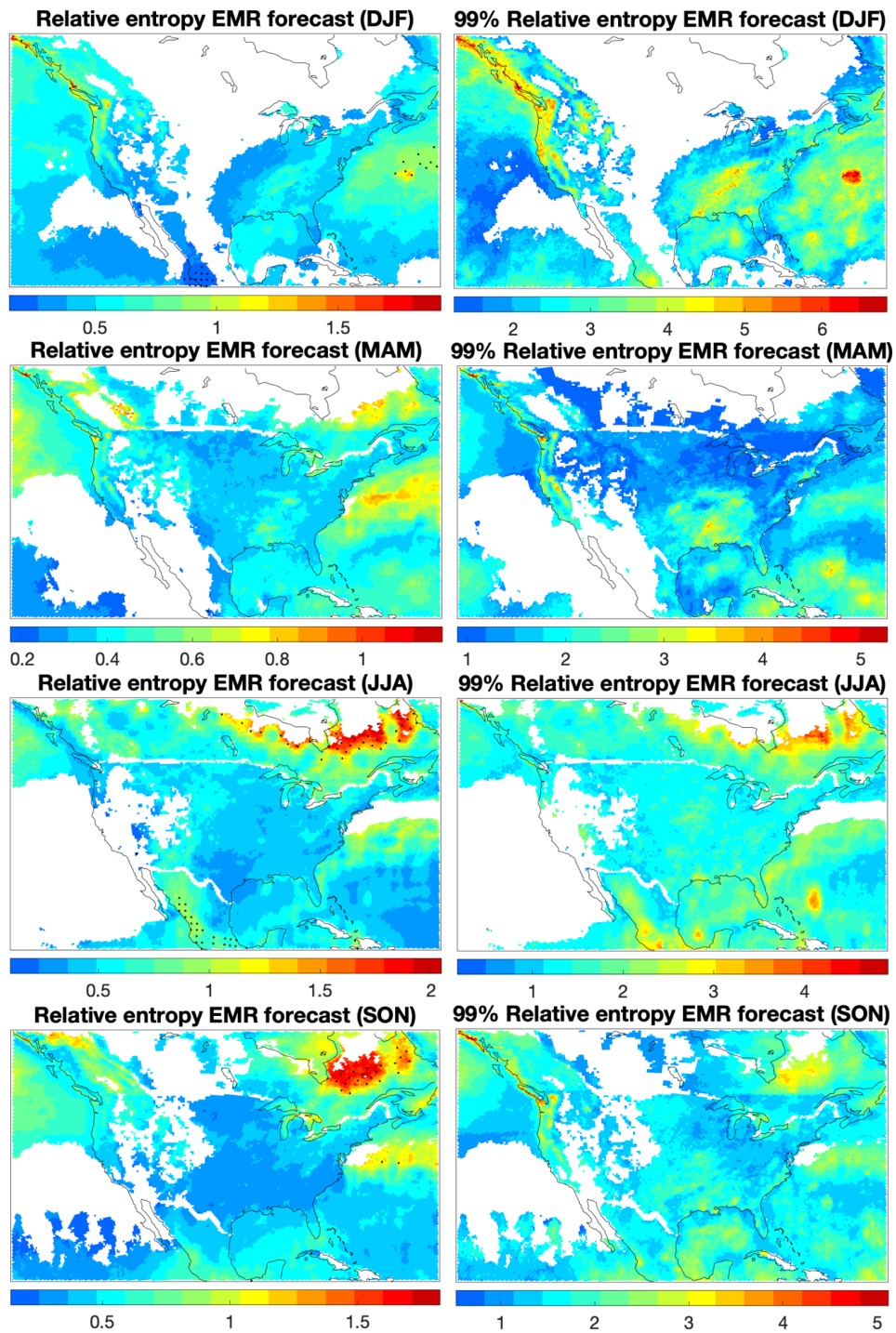
568   9,  right).

569       Hence,  the EMR model produces reliable (low-dispersion) forecasts that tend to track the

570   observed precipitation (signal), much more so than the damped persistence forecasts. Kleeman

571   (2002) argued that a forecast's relative entropy

572

573
$$R = \sum_i p_i \ln \frac{p_i}{q_i}, \qquad (6)$$

574   where $p_i$ is climatological distribution and $q_i$ is that for the prediction, can be very useful in

575   characterizing prediction utility as it naturally captures both the signal and dispersion compo-

File generated with AMS Word template 1.0

**Figure 10**: Relative entropy of EMR forecasts. Left: the expectation (climatology), with stippling
showing the areas where this expectation exceeds that associated with the damped persistence forecast
(see text for details); right: the 99[th] percentile. Note that the relative entropy here was only computed
and shown over the grid points at which the 99[th] percentile of $P_{10}$ exceeded 50 mm (cf. Fig. 7, left);
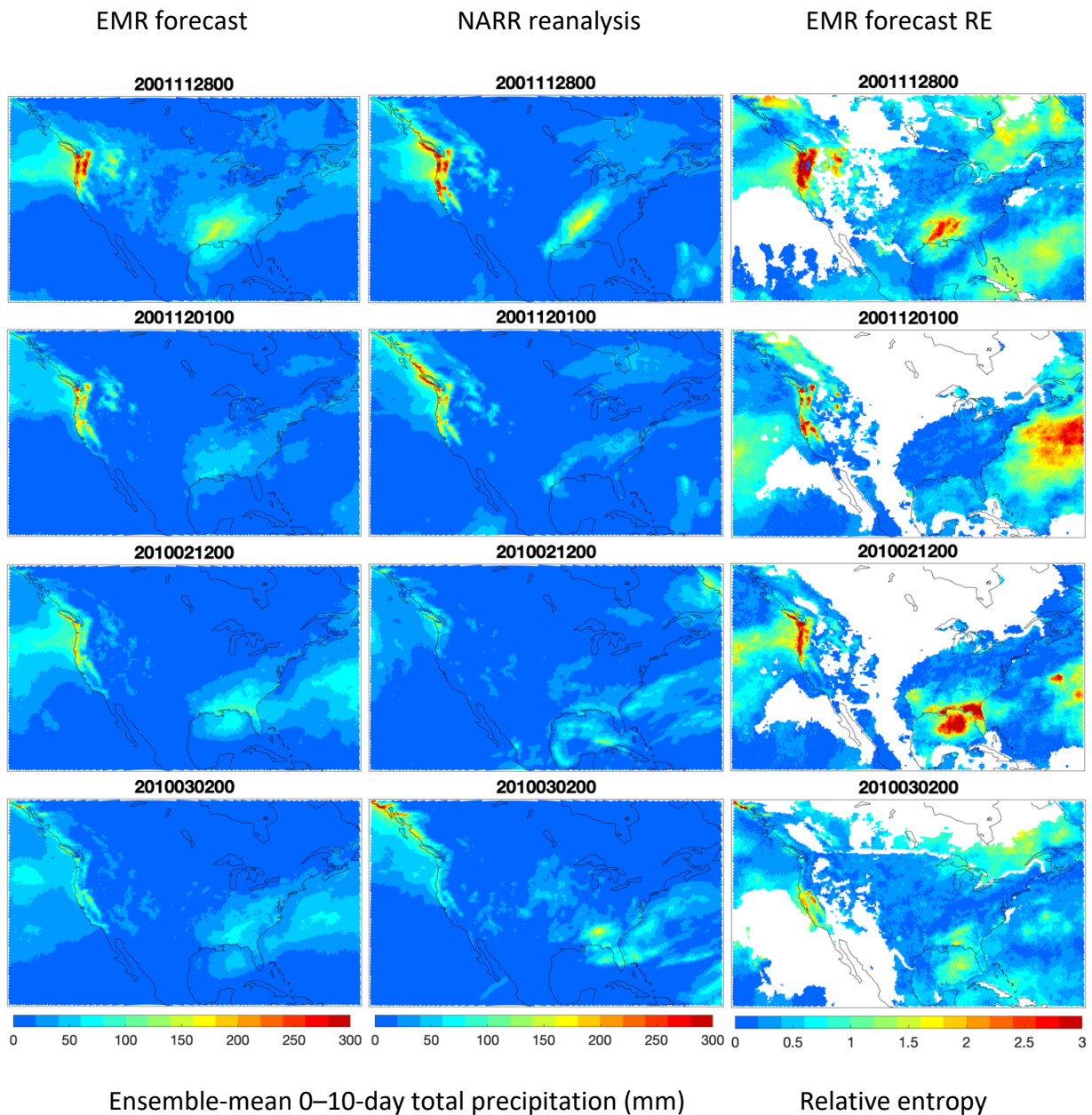the areas in which this is not the case are colored white.

28

583 nents of skill. The relative entropy measures how different the forecast distribution is from a
584 climatological distribution. However, the *expectation* of $R$ (characterizing climatological
585 difference between forecasts and observations) would tend to be lower for the forecast schemes
586 that are more skillful than others. For example, the climatological relative entropy associated
587 with the damped persistence forecasts is expected to be higher than that for the EMR forecasts.
588 This is indeed the case (**Fig. 10**, left) [note that the relative entropy here was only computed
589 and shown over the grid points at which the $99^{th}$ percentile of $P_{10}$ exceeded 50 mm (cf. Fig. 7,
590 left)]. Yet, over time, the relative entropy associated with individual $P_{10}$ forecasts can greatly
591 exceed its climatological value (Fig. 10, right). In section 4c below, we will develop a
592 subsampling strategy in which the forecasts with large values of the quantity $R$ are tagged to
593 define and sample potential large and extreme precipitation events.

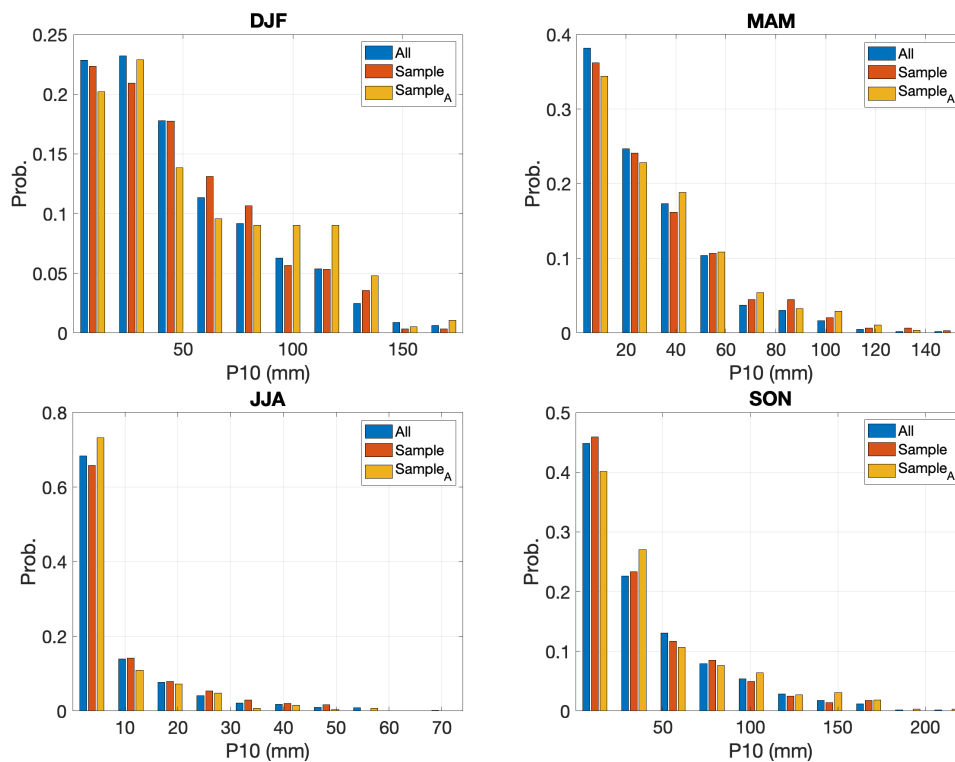594 *c. EMR based probabilistic algorithm for thinning reforecast sample size*

595 Note that the cases displayed in Figs. 5 and S5 were selected using the *ad hoc* algorithm
596 based on heavy precipitation in GEFSv12 reforecasts (section 3a; Sample$_A$) (the multi-page
597 image files with analogous maps for other selected cases are available through a webpage
598 referenced in the Supplementary Information). As mentioned before, the same algorithm was
599 applied to the EMR model's ensemble-mean $P_{10}$ forecasts (which are also available through
600 the supplementary website); see a brief discussion at the end of this section. We here also
601 developed and applied an alternative strategy, which selects the dates based on the large value
602 of the EMR forecasts' relative entropy. In particular, we computed, for each day, the average
603 among the top 10% relative-entropy grid point values over CONUS (which were also pre-
604 selected to have the seasonal 1979–1999 $P_{10}$'s $99^{th}$ percentile exceeding 5 cm, thus excluding
605 the white areas in Fig. 10); each day in the record was then ranked based on its relative entropy
606 score. Upon selecting 40% of the highest-score dates from the first and 40% of the highest-
607 score dates from the second half of the original 2000–2020 sample (thereby eliminating
608 possible effects of any long-term relative entropy trends), we edited out the member with a
609 higher $R$ from all the pairs of consecutive high-relative-entropy days identified above, and then
610 from the pairs separated by two days. This procedure results in the identification of 1095 cases
611 separated by at least two days out of the total 7671 days comprising the 2000–2020 period,
612 which we argue to be an optimal subset including the majority of the high-impact events and
613 yet also representative of the climatological $P_{10}$ distribution. If more frequent sampling is

29

File generated with AMS Word template 1.0

614  required, the  additional dates for reforecasts can be added at random from the remainder of
615  the record.



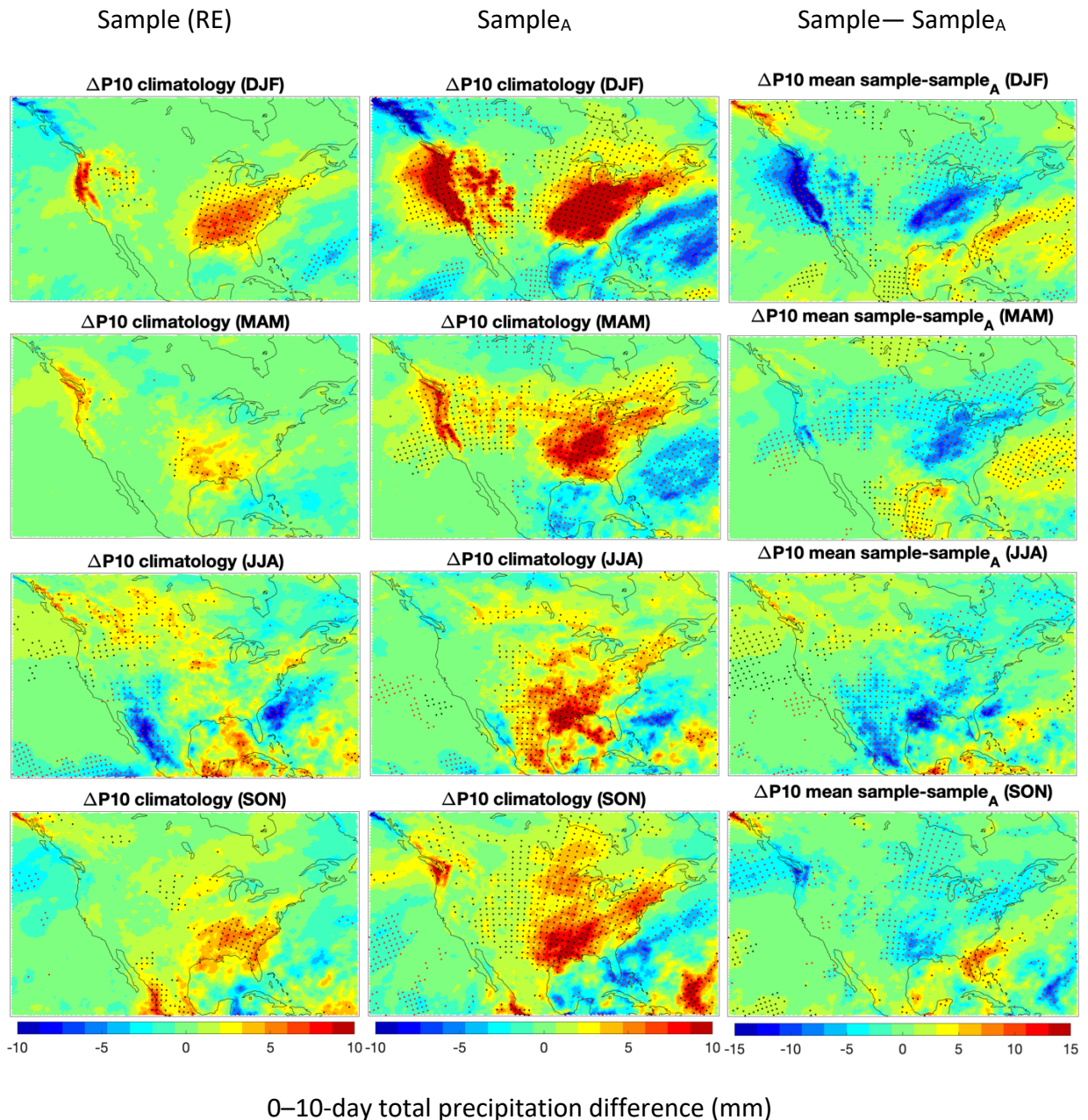Ensemble-mean 0–10-day total precipitation (mm)     Relative entropy

616

**Figure 11**: Examples of (cool season) $P_{10}$ forecasts using EMR model (left), along with the actual $P_{10}$
maps based on NARR reanalysis (middle). Units are mm. The right column shows the corresponding
map of the relative entropy. The forecast initialization time (the same across each row) is shown in
panel captions in the YYYYMMDDHH format. Note that, similar to Fig. 10, the relative entropy in
the right-column plots was only computed and shown over the grid points at which the 99[th] percentile
of $P_{10}$ exceeded 50 mm (cf. Fig. 7, left); the areas in which this is not the case are colored white.

623    The size of the latter sample is also consistent with that of the sample$_A$, which has 520 cases

624    per each of the semi-annual cool and warm seasons (over 2000–2019 period, with the following

625    breakdowns: DJF – 188, MAM – 276, JJA – 247, SON – 329 cases). The corresponding

626    breakdowns for the present sample are: DJF – 282, MAM – 290, JJA – 240, SON – 283 cases,

627    featuring a more uniform seasonal distribution of cases, with more  DJF cases and fewer SON

628    cases compared to the GEFSv12 based subsample. The two samples turn out to be largely

629    independent,  with only 198 (~20%) matching dates over the 2000–2019 period. A few

630    examples of the $P_{10}$ observed and predicted maps based on the present sample are shown in

631    **Figs. 11** and **S6** (and others are available through the Supplementary website). The third

632    column of these figures shows the distribution of the EMR forecasts' relative entropy on a

633    given day, which tends to track the areas of large and extreme precipitation (recall that the

634    relative-entropy-based selection criterion was only applied over CONUS, rather than over a

635    larger region of the NARR reanalysis).



636

637    **Figure 12**: The $P_{10}$'s probability density function (PDF) estimates at 47.4 N, 122.4 W (Seattle, WA)

638    based, for each season, on the entire daily $P_{10}$ data (blue), and two thinned subsamples ~1/7 the size

639    of the whole available data: a subsample based on relative entropy of EMR forecasts (EMR-RE)

640    (sample, red) and the one (sample$_A$, yellow) based on ensemble-mean GEFSv12 $P_{10}$ forecasts

641    associated with significant precipitation events over CONUS (see text for details).

31

Sample (RE)          Sample$_A$          Sample— Sample$_A$

0–10-day total precipitation difference (mm)

**Figure 13**: The differences between the estimates of $P_{10}$'s climatological mean based on: EMR-based thinned sample and the entire seasonal $P_{10}$ data (left column); GEFSv12-based thinned sample$_A$ and the entire seasonal $P_{10}$ data (middle column); EMR-based thinned sample and GEFSv12-based thinned sample$_A$ (right column). Stippling shows areas where the differences are statistically significant at the 5% level according to the two-sided bootstrap test involving surrogate random subsamples of the same size as either the EMR-based thinned sample or GEFSv12-based thinned sample$_A$.

32

File generated with AMS Word template 1.0

651    To assess relative performance of the two methods, we computed distributions of $P_{10}$

652    associated with each sample and compared them with the climatological distribution of $P_{10}$.

653    An example of these distributions in **Fig. 12** demonstrates that the EMR based sample provides

654    a better match to the NARR based $P_{10}$ climatological distribution than the GEFSv12 based

655    subsample, which tends to be excessively heavy tailed (DJF panel gives a particularly clear

656    example of this for the location chosen).  The positive bias of the GEFSv12 based sample

657    (perhaps natural, given the  selection criterion built on the direct occurrence of the large or

658    extreme precipitation) is also  evident in the maps of the climatological mean (**Fig. 13**) and (to

659    a somewhat lesser extent) in the maps of the 99$^{th}$ percentile (**Fig. S7**), of the distributions based

660    on full and subsampled data.  **Overall, the present sample has a distribution of 0–10-day**

661    **total precipitation that is closer to the distribution based on the full data compared to**

662    **that of GEFSv12 based sample$_A$, while capturing the majority of high-impact**

663    **precipitation events.** It should be noted,  however, that the ultimate test of the success of the

664    subsampling will be the accuracy of postprocessed precipitation guidance based on the sample

665    at hand, and not the fidelity against the NARR data.  For example, heavier precipitation periods

666    preferentially sampled by the GEFSv12 algorithm by design may be particularly important for

667    establishing the statistical relationships in situations with heavy precipitation that are of

668    greatest interest.

669    Finally, we note here that the thinned sample$_B$ obtained using the same algorithm as for the

670    GEFSv12 data, but applied to the EMR precipitation forecasts, produced results inferior of

671    those associated with either the EMR-RE sample or the GEFSv12-based sample$_A$ in terms of

672    the similarity of climatological precipitation distributions based on the thinned and full

673    available data samples (**Figs. S8** and **S9**). This may be due to the fact that the EMR forecasts

674    of $P_{10}$ have a smaller deterministic skill than analogous high-end GEFSv12 reforecasts.

675

## 5. Summary and discussion

677    In this study, we developed a novel methodology for multi-scale statistical modeling of

678    precipitation by utilizing the Empirical Model Reduction (EMR) technique (Kravtsov and co-

679    authors 2005–2017) applied to the NARR reanalysis. The key element of the new algorithm is

680    the usage of the pseudo-precipitation $PP$ — whose positive values are associated with the

681    actual precipitation, while negative values represent the column integrated water vapor

33

File generated with AMS Word template 1.0

682 saturation deficit — as a part of the climate state vector to be simulated  by the EMR model.

683 The *PP* field thus carries information about both the mesoscale precipitation features and

684 synoptic-scale environmental background (large-scale meteorological patterns: LSMP)

685 potentially conducive to high-impact precipitation events. This EMR model was found to

686 provide a seamless spatiotemporal statistical description of the precipitation-producing weather

687 systems across a wide range of spatial scales over the entirety of CONUS and to possess a

688 significant predictive skill,  especially in a probabilistic sense.

689     We defined the events-of-impact in terms of the relative entropy (Kleeman 2002) of the

690 EMR based ensemble hindcasts of the 0–10-day total surface precipitation $P_{10}$ over the 2000–

691 2020 period and identified an optimal (arguably minimal) subset of dates proved to provide

692 local precipitation distributions consistent with those based on the full data set. By contrast, an

693 alternative statistical methodology for selecting such dates based directly on the magnitude of

694 $P_{10}$ in high-end ensemble-mean reforecasts of precipitation produced subsamples with a more

695 substantial heavy-precipitation bias. Thinning the frequency of reforecasts — the task that

696 motivated this research in the first place — is extremely important in a variety of hydrological

697 modeling applications to be described in a future companion paper.

698     Note that our selecting reforecast cases for their presumed importance in one metric (here,

699 0–10-day precipitation) may bias the sampling properties for different kinds of important

700 extreme events, which might include hurricanes, mixed precipitation events, severe weather,

701 extreme surface temperatures or winds, among others. For example, heavy precipitation events

702 are forecast better using the quantile approach with respect to precipitable water than the

703 absolute magnitude of the precipitable water (refs?). Such biases, however, would be a

704 limitation of any method that seeks to limit the reforecast sample size.

705     Another possible limitation of the EMR methodology developed here is that the EMR

706 model is trained on the earlier data, while the ongoing climate change may skew the more

707 recent historical record in various ways, introducing a bias into EMR forecasts associated with

708 the latter record.  For our present application, we believe that such biases associated with the

709 $P_{10}$ statistics are relatively small, as evidenced by a fairly uniform in time distribution of dates

710 in our thinned samples (so that, for example, the number of important cases identified in the

711 first and second halves of the 2000–2020 record is similar).

712     Our new EMR methodology for statistical modeling of precipitation is fundamentally

713 different from more traditional techniques (which typically work with individual precipitation

714 records at a local level and/or postulate *ad hoc* connections with a limited number of large-

34

File generated with AMS Word template 1.0

715  scale predictors:  see section 2a)  in that it automatically accounts for spatiotemporal multi-
716  scale structure of precipitation dynamics, thereby providing a unified framework to model
717  diverse precipitation environments. At the same time, it is still extremely numerically  efficient
718  and thus easily permits large-ensemble simulations/forecasts which are essential for monitoring
719  and fully utilizing probabilistic characteristics of precipitation, in contrast to full-blown
720  dynamical models necessarily limited in the number of ensemble members due to prohibitive
721  computational expenses.

722      This paper showcases just one application of the new EMR precipitation model to the
723  problem of thinning the frequency of reforecasts. Follow-up work will look into how the
724  various sampling strategies affect precipitation forecast calibration and hydrologic forecast
725  accuracy. We also plan to further test the EMR model's potential in a wider range of related
726  problems around the statistical/dynamical analysis of precipitation and its predictability.

727

728

733

734  *Data Availability Statement.*

735  The NARR reanalysis data is available at https://psl.noaa.gov/data/gridded/data.narr.html.
736  GEFSv12 data may be accessed at https://noaa-gefs-
737  retrospective.s3.amazonaws.com/index.html. This manuscript also has a supplementary
738  website with data and figures generated during this study, as described in detail in the
739  Supplemental Information.  All MATLAB/Python scripts associated with this project are
740  available from the authors by request.

741

742

File generated with AMS Word template 1.0

REFERENCES

Albers, J. R., and M. Newman, 2019: A priori identification of skillful extratropical subseasonal forecasts. *Geophys. Res. Lett.*, **46**, https://doi.org/10.1029/2019GL085270.

Barlow, M., W. J. Gutowski Jr., J. Gyakum, et al., 2019: North American extreme precipitation events and related large-scale meteorological patterns: a review of statistical methods, dynamics, modeling, and trends. *Climate Dyn.*, **53**, 6835–6875, https://doi.org/10.1007/s00382-019-04958-z

Bolton, D., 1980: The computation of equivalent potential temperature. *Mon. Wea. Rev.*, **108**, 1046–1053.

Box, G. E. P., G. M. Jenkins, and G. C. Reinsel, 1994: *Time Series Analysis, Forecasting and Control*. 3d Ed., Prentice Hall, 592 pp.

Bukovsky, M. S., and D. J. Karoly, 2007: A brief evaluation of precipitation from the North American Regional Reanalysis. *J. Hydrometeorol.*, **8**, 837–846, doi:10.1175/JHM595.1

Demargne, J., L. Wu, S. K. Regonda, J. D. Brown, H. Lee, M. He, D. Seo, R. Hartman, H. D. Herr, M. Fresch, J. Schaake, and Y. Zhu, 2014: The Science of NOAA's Operational Hydrologic Ensemble Forecast Service. *Bull. Amer. Meteor. Soc.*, **95**(1), 79–98. Retrieved Aug 31, 2021, from https://journals.ametsoc.org/view/journals/bams/95/1/bams-d-12-00081.1.xml

Furrer, E., and R. Katz, 2007: Generalized linear modeling approach to stochastic weather generators. *Climate Res.*, **34**, 129–144, doi:10.3354/cr034129.

Grotjahn, R., R. Black, R. Leung, M. W. Wehner, M. Barlow, M. Bosilovich, A. Gershunov, W. J. Gutowski Jr, J. R. Gyakum, R. W. Katz, Y.-Y. Lee, Y.-K. Lim, Prabhat, 2016: North American extreme temperature events and related large scale meteorological patterns: A review of statistical methods, dynamics, modeling and trends. *Climate Dyn.*, **46**, 1151–1184. https://doi.org/10.1007/s00382-015-2638-6

Guan, H., and others, 2021: The GEFSv12 reforecast dataset for supporting sub-seasonal and hydrometeorological applications.  In preparation.

Hamill, T. M., 2018:  Practical Aspects of Statistical Postprocessing.  Chapter 7 in the book Statistical Postprocessing of Ensemble Forecasts (Elsevier Press).

772    Hamill, T. M., and J. S. Whitaker, 2006: Probabilistic quantitative precipitation forecasts based
773        on reforecast analogs: theory and application. *Mon. Wea. Rev.*, **134**, 3209–3229.

774    Hamill, T. M., J. S. Whitaker, and S. L. Mullen, 2006: Reforecasts, an important dataset for
775        improving weather predictions. *Bull. Amer. Meteor. Soc.*, **87**, 33–46.

776    Hamill, T. M., G. T. Bates, J. S. Whitaker, D. R. Murray, M. Fiorino, T. J. Galarneau, Jr., Y.
777        Zhu, and W. Lapenta, 2013: NOAA's second-generation global medium-range ensemble
778        reforecast data set. *Bull Amer. Meteor. Soc.*, **94**, 1553–1565.

779    Hamill, T. M., M. Scheuerer, and G. T. Bates, 2015: Analog probabilistic precipitation
780        forecasts using GEFS Reforecasts and Climatology-Calibrated Precipitation Analyses.
781        *Mon. Wea. Rev.*, **143**, 3300–3309.

782    Hamill, T. M., and others, 2021: The reanalysis for the Global Ensemble Forecast System,
783        version 12. Mon. Wea. Rev., accepted/minor.

784    Holsclaw, T. , A. M. Greene, and A. W. Robertson, 2016: A Bayesian Hidden Markov Model
785        of daily precipitation over South and East Asia. *J. Hydrometeorol.*, **17**, 3–25, doi:
786        10.1175/JHM-D-14-0142.1.

787    Kenabatho, P. K., N. R. McIntyre, R. E. Chandler, and H. S. Wheater, 2012: Stochastic
788        simulation of rainfall in the semi-arid Limpopo basin, Botswana. *Int. J. Climatol.*, **32**,
789        1113– 1127, doi:10.1002/joc.2323.

790    Kleeman, R., 2002: Measuring dynamical prediction utility using relative entropy. *J. Atmos.*
791        *Sci.*, **59**, 2057–2072. https://journals.ametsoc.org/view/journals/atsc/59/13/1520-
792        0469_2002_059_2057_mdpuur_2.0.co_2.xml

793    Kravtsov, S., D. Kondrashov, and M. Ghil, 2005: Multi-level regression modeling of nonlinear
794        processes: Derivation and applications to climatic variability. *J. Climate*, **18**, 4404–4424.

795    Kravtsov, S., M. Ghil, and D. Kondrashov, 2010: *Empirical Model Reduction and the Modeling*
796        *Hierarchy in Climate Dynamics and the Geosciences. Stochastic Physics and Climate*
797        *Modeling*, T. Palmer and P. Williams, Eds., Cambridge University Press, pp. 35–72.

798    Kravtsov, S., N. Tilinina, Y. Zyulyaeva, and S. Gulev, 2016: Empirical modeling and stochastic
799        simulation of sea-level pressure variability. *J. Appl. Meteor. Climat.,* **55**, 1197–1219, doi:
800        http://dx.doi.org/10.1175/JAMC-D-15-0186.1.

37

801  Kravtsov, S., P. Roebber, and V. Brazauskas, 2017: A virtual climate library of surface
802      temperature over North America for 1979–2015. *Scientific Data*, **4**, 170,155EP,
803      doi:10.1038/sdata.2017.155.

804  Lindsay, R. K., M. A. Kohler, J. L. H. Paulhus, 1975: *Hydrology for Engineers*, 2nd ed.,
805      McGraw-Hill.

806  Manzanas, R., A. Lucero, A. Weisheimer, and J. M. Gutierrez, 2018: Can bias correction and
807      statistical downscaling methods improve the skill of seasonal precipitation forecasts?
808      *Climate Dyn.*, **50**, 1161–1176, doi: 10.1007/s00382-017-3668-z.

809  McCullagh, P., and J. Nelder, 1989: *Generalized Linear Models*. Chapman and Hall, 532 pp.

810  Mesinger, F., et al., 2006: North American Regional Reanalysis. *Bull. Amer. Meteor. Soc.* **87**,
811      343–360.

812  Newman, M., P. D. Sardeshmukh, C. R. Winkler,  and J. S. Whitaker, 2003: A study of
813      subseasonal    predictability.    *Mon.    Wea.    Rev.*,    **131**,    1715–1732,
814      https://doi.org/10.1175//2558.1.

815  Penland, C., 1996: A stochastic model of Indo-Pacific sea surface temperature anomalies.
816      *Physica D*, **98**, 534–558, https://doi.org/ 10.1016/0167-2789(96)00124-8.

817  Penland, C., and P. D. Sardeshmukh, 1995: The optimal growth of tropical sea surface
818      temperature anomalies. *J. Climate*, **8**, 1999–2024, https://doi.org/10.1175/1520-
819      0442(1995)008 <1999:TOGOTS>2.0.CO;2

820  Robertson, A. W., and M. Ghil, 1999: Large-scale weather regimes and local climate over the
821      western United States, *J. Climate*, **12**, 1796–1813.

822  Robertson, A. W., Y. Kushnir, U. Lall, and J. Nakamura, 2016: Weather and climatic drivers
823      of extreme flooding events over the Midwest of the United States. In: *Extreme Events:*
824      *Observations, Modeling, and Economics, Geophysical Monograph*, vol. 214, 1[st] Ed.  M.
825      Chavez, M. Ghil, J. Urrutia-Fucugauchi, Eds., American Geophysical Union, John Wiley
826      & Sons, Inc.

827  Robertson, A. W., S. Kirshner, and P. Smyth, 2004: Downscaling of daily rainfall occurrence
828      over northeast Brazil using a hidden Markov model. *J. Climate*, **17**, 4407–4424,
829      doi:10.1175/ JCLI-3216.1.

830  Scheuerer, M., and T. M. Hamill, 2015: Statistical post-processing of ensemble precipitation
831      forecasts by fitting censored, shifted Gamma distributions. *Mon. Wea. Rev.*, **143**, 4578–
832      4596.

833  Sinha, P., U. C. Mohanty, S. C. Kar, S. K. Dash, A. W. Robertson, and M. K. Tippett, 2013:
834      Seasonal prediction of the Indian summer monsoon rainfall using canonical correlation
835      analysis of the NCMRWF global model products. *Int. J. Climatol.,* **33**, 1601–1614, doi:
836      10.1002/joc.3536.

837  Von Mises, R., 1964: *Mathematical Theory of Probability and Statistics*. Academic Press, New
838      York.

839  Wilks, D. S., 2011: *Statistical Methods in Atmospheric Sciences*, 3rd ed., Academic Press, 704
840      pp.

841  Winkler, C. R., M. Newman, and P. D. Sardeshmukh, 2001: A linear model of wintertime low-
842      frequency variability. Part I: Formulation and forecast skill. *J. Climate*, **14**, 4474–4494,
843      https://doi.org/10.1175/1520-0442(2001)014<4474:ALMOWL>2.0.CO;2.

844  Wold, S., et al., 1984: *Chemometrics, Mathematics and Statistics in Chemistry*. Reidel
845      Publishing Company, Dordrecht, Holland.

846  Yuan, H., P. Schultz, E. I. Tollerud, D. Hou, Y. Zhu, M. Pena, M. Charles, and Z. Toth, 2019:
847      Pseudo-precipitation: A continuous precipitation variable. *NOAA Technical Memorandum*
848      *OAR GSD-62*, https://doi.org/10.25923/3h37-gp49.

849  Zhou, X., and others 2021: The Introduction of the NCEP Global Ensemble Forecast System
850      Version 12, Mon. Wea. Rev., submitted

851  Zobel, Z., J. Wang, D. J. Wuebbles, and V. R. Kotamarthi, 2018: Evaluations of high-resolution
852      dynamically downscaled ensembles over the contiguous United States. *Climate Dyn.*, **50**,
853      863–884, doi: 10.1007/s00382-017-3645-6.