

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16

**Evaluation of CAPS Convection-Allowing FV3-LAM Ensembles during the
2022 HWT Spring Forecasting Experiment to Inform the Design of the
Rapid Refresh Forecast System (RRFS)**

Nathan Snook¹, Jun Park¹, Ming Xue^{1,2}, Keith A. Brewster^{1,2}, Marcus Johnson¹, Timothy
Supinie³, Xiao-Ming Hu¹, Jacob R. Carley⁴, Shun Liu⁴, and Ming Hu⁵

¹ *Center for Analysis and Prediction of Storms, Univ. of Oklahoma, Norman, OK*

² *School of Meteorology, Univ. of Oklahoma, Norman, OK*

³ *NOAA Storm Prediction Center, Norman, OK*

⁴ *NOAA Environmental Modeling Center, College Park, MD*

⁵ *NOAA Global Systems Laboratory, Boulder, CO*

Corresponding author:

Nathan Snook, nsnook@ou.edu

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

ABSTRACT

To inform the optimization of the future operational Rapid Refresh Forecast System (RRFS), the Center for Analysis and Prediction of Storms (CAPS) performed three sets of CONUS-domain 3-km FV3-LAM ensemble forecasts with various configurations during the 2022 NOAA Hazardous Weather Testbed Spring Forecasting Experiment in real time. The first set used different physics parameterizations, the second set included additional initial and lateral boundary condition perturbations, and the third introduced additional stochastic physics perturbations. This study evaluates precipitation, temperature, dewpoint, and wind forecasts, and compares them to those of the operational High-Resolution Ensemble Forecast (HREF) and Global Ensemble Forecast System (GEFS). Precipitation forecasts are verified against NCEP Stage-IV precipitation analyses, while temperature, dewpoint, and wind forecasts are verified against URMA surface analyses and radiosonde observations.

Overall, the ensemble configurations tested are generally suitable for predicting spring-season convective rainfall. The CAPS forecasts generally outperform GEFS in terms of ETS, frequency bias, and area under the ROC curve, approaching (but not exceeding) the performance of HREF in some metrics. Including stochastic physics perturbations resulted in forecasts objectively very similar to those without such perturbations, except for a small but consistent positive impact on ensemble spread of surface variables throughout the 84-hour forecast period. The CAPS forecasts have a near-neutral to slightly negative bias in total precipitation coverage. Forecasts using the NSSL microphysics scheme have more total rainfall than forecasts using the Thompson scheme and Stage-IV analyses, while forecasts using the NOAH-MP land surface model generally have lower total precipitation than other forecasts.

41 **1. Introduction**

42 As part of the transition to the Unified Forecast System (UFS) framework for operational
43 numerical weather prediction (NWP) systems, the National Oceanic and Atmospheric
44 Administration (NOAA) is developing a rapidly-updating convection-allowing ensemble
45 forecasting system—the Rapid Refresh Forecast System (RRFS; Banos et al 2022, Carley et
46 al. 2023)—which is planned to facilitate the retirement of a large portion of operational high-
47 resolution forecasting systems currently used by the National Weather Service (NWS). The
48 RRFS will use the limited area variant (Black et al. 2021) of the Finite-Volume Cubed-Sphere
49 (FV3) Model (FV3-LAM) (Lin 2004) and include an ensemble of forecasts optimized for short-
50 range prediction of high-impact weather. To assist in optimizing the ensemble design and
51 physics parameterization choices for the RRFS, the Center for the Analysis and Prediction of
52 Storms (CAPS) at the University of Oklahoma has been running experimental FV3-LAM
53 ensemble forecasts during the NOAA Hazardous Weather Testbed (HWT) Spring Forecasting
54 Experiment (SFE; Clark et al. 2020; Roberts et al. 2020) and the Hydrometeorology Testbed
55 (HMT) Flash Flood and Intense Rainfall Experiment (FFaIR) and Winter Weather Experiments
56 (NOAA 2023) since 2018—the ensemble membership and model configuration of these
57 forecast ensembles has varied from experiment to experiment. These experiments have
58 demonstrated the ability of convection-allowing FV3-LAM ensembles to generate skillful
59 forecasts, while also documenting systematic biases associated with different physics and land-
60 surface model combinations (Zhang et al. 2019; Snook et al. 2019; Supinie et al. 2022; Hu et
61 al. 2023).

62 During the 2022 HWT SFE, CAPS ran a 21-member FV3-LAM ensemble; member
63 configurations were selected to use various combinations of land surface models (LSMs), and
64 planetary boundary layer (PBL), surface layer, and microphysical schemes. The 2022 HWT
65 SFE configurations evolved from those used in 2020 and 2021 (Supinie et al. 2022; Hu et al.
66 2023). In the 2022 ensemble, initial conditions from experimental RRFS EnVar and EnKF
67 analyses (Carley et al. 2024) are used, instead of the GFS and Global Ensemble Forecast
68 System (GEFS) analyses used in prior years. Ensemble analyses from the experimental RRFS
69 EnKF run by the NOAA Global Systems Laboratory (GSL) are used to initialize perturbed-IC
70 ensemble members; unperturbed members are initialized using 3DEnVar analysis with
71 ensemble covariance provided by the EnKF system so that two systems are coupled. In this
72 sense the CAPS FV3-LAM forecasts of 2022 can be considered “hot start” forecast runs. The

73 experimental RRFS EnVar and EnKF analyses were run by GSL as hourly cycled DA systems
74 as prototypes of the planned operational RRFS.

75 Multiple suites of physics parameterization schemes are chosen for the CAPS FV3-LAM
76 forecasts, partly based on earlier performance evaluations (e.g., Supinie et al. 2022; Hu et al.
77 2023) and partly based on the likelihood of continued support and operational use by the NWS.
78 Stochastic physics perturbations are also included for five ensemble members. Model
79 configurations in this ensemble are designed to resemble current and proposed convection-
80 allowing model configurations (these configurations are discussed in greater detail below in
81 section 2a). The primary goal of this study is to examine the performance of these forecasts,
82 including deterministic forecast performance of selected members and ensemble forecast
83 performance of several sub-ensembles (i.e., subsets of the 21 members intended to investigate
84 different ensemble design considerations), with emphasis on precipitation forecasts.

85 Among the 21 forecast members, five use the same initial condition (IC) and lateral
86 boundary conditions (LBCs), with differences in physics parameterizations only, allowing for
87 the evaluation of physics performance. These same physics combinations are also used in other
88 sub-ensembles that introduce IC, LBC, and stochastic physics perturbations. The use of physics
89 diversity has been shown to help improve the spread and probabilistic forecast performance of
90 ensembles, especially at the mesoscale and convective scale (e.g., Stensrud et al. 2020; Clark
91 et al. 2008; Berner et al. 2011), and is used in the current NWS operational High-Resolution
92 Ensemble Forecast (HREF) system (Jirak et al. 2012; Roberts et al. 2019) which is also a multi-
93 model ensemble. Real-time HWT ensemble forecasts produced by CAPS over the years have
94 also used multiple physics configurations (e.g., Schwartz et al. 2008; Loken et al. 2019; Supinie
95 et al. 2022; Hu et al. 2023).

96 Inclusion of stochastic physics perturbations has also been shown to improve ensemble
97 characteristics; several perturbation approaches have been investigated in prior studies,
98 including stochastic kinetic energy backscatter (SKEB; Berner et al. 2009), stochastic
99 perturbations of physics tendencies (SPPT; Buizza et al. 1999, Palmer et al. 2009), stochastic
100 perturbed humidity (SHUM; Thompkins and Berner 2008), stochastic parameter perturbations
101 (SPP; Jankov et al., 2017, 2019), and cellular automata (Bengtsson et al. 2013). The current
102 NCEP GEFS uses SPPT and SKEB perturbations (Zhou et al. 2022). Jankov et al. (2017, 2019)
103 demonstrated that a single-physics ensemble with stochastic physics could perform comparably
104 to a multi-physics ensemble for regional systems. Kalina et al. (2021) found in experimental

105 configurations of the High-Resolution Rapid Refresh Ensemble (HRRRE) that SPP helped
106 increase ensemble spread in surface forecast variables. They also found that SPP increased the
107 reliability of near-term HRRRE precipitation forecasts but exacerbated a preexisting low-
108 frequency bias in the prediction of heavy rainfall in HRRRE.

109 Important questions regarding stochastic physics perturbations remain to be solved. One
110 such question is whether random perturbations introduced into physics parameters or physics
111 tendencies cause deterioration in individual members' forecast performance. Another question
112 is whether stochastic perturbations using a single physics suite can create sufficient spread
113 within an ensemble—if so, the cost of having to develop and maintain multiple physics
114 packages can be alleviated, allowing more resources to be devoted to the improvement of the
115 best-performing physics suite. Thus far, the findings concerning the optimal ensemble
116 configurations in terms of physics combinations and the use of stochastic perturbations are still
117 inconclusive, and the answers are likely dependent on the specific models used, their
118 implementation, and their applications. One subset of the CAPS ensemble examined in this
119 study uses a combination of SKEB, SPPT and SHUM stochastic perturbations, offering us an
120 opportunity to examine the effects of stochastic perturbations as well as to test the robustness
121 of their implementations in combination with multiple physics parameterizations. It should be
122 noted, however, that stochastic physics perturbations can be strongly impacted by the choice
123 of model dynamic core, and further evaluation would be needed in the context of any specific
124 future RRFS ensemble.

125 Though in this study we document and compare the performance of the CAPS ensembles
126 to existing operational forecasts as a baseline, the primary focus is on the relative performance
127 of the sub-ensembles and the effects of physics diversity, IC and LBC perturbations, and
128 stochastic physics perturbations. These experiments have helped to inform future operational
129 RRFS ensemble design; members of the RRFS development team directly consulted with
130 CAPS, and the results presented in this study were highly influential in the construction of the
131 initial version of RRFS.

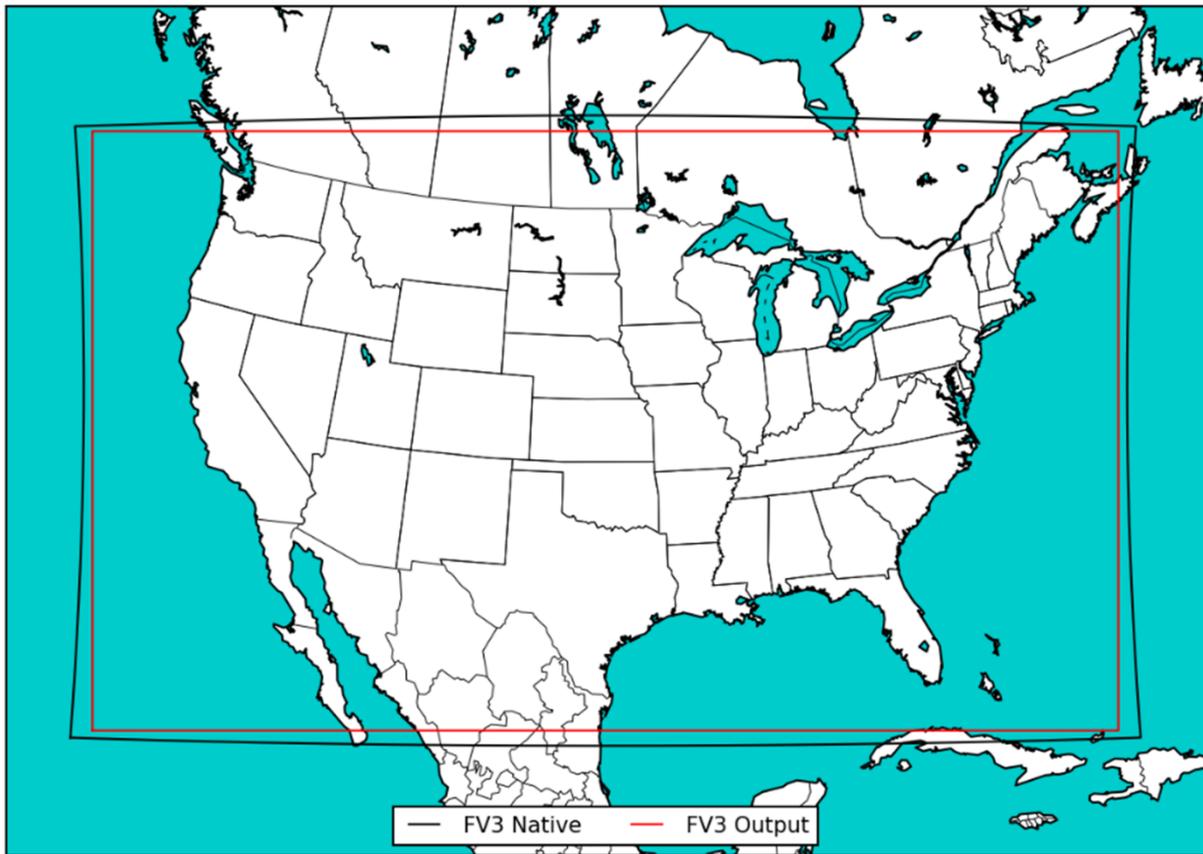
132 The remainder of this study is organized as follows: the specifics of the 2022 CAPS HWT
133 SFE ensemble forecasts, including the sub-ensemble design considerations and the verification
134 datasets and methodologies, are discussed in section 2. Objective and subjective verifications
135 of the ensembles using Stage-IV precipitation data, surface analyses, and sounding data are
136 presented in section 3. Conclusions and operational implications are discussed in section 4.

137

138 **2. FV3-LAM Configuration, Verification Data, and Methodology**

139 *a. Model configurations and ensemble design*

140 The forecasts produced by CAPS during 2022 HWT SFE (referred to hereafter as the
141 “CAPS ensemble”) use a model domain shown in Fig. 1, which covers the contiguous United
142 States (CONUS) with a horizontal grid spacing of ~3 km and 1821×1093 grid points together
143 with 65 vertical layers. Forecasts were generated on weekdays from 2 May 2022 through 3
144 Jun. 2022, and were initialized at 0000 UTC and run for 84 hours of forecast time, providing
145 1200-1200 UTC daily forecasts for days 1, 2, and 3 (12-36, 36-60, and 60-84 hours of forecast
146 time, respectively). Due to technical issues, the most common of which was unavailable or
147 incomplete RRFS EnKF initial condition data, complete forecast ensemble data were only
148 available for 13 of the 25 days during the 2022 HWT SFE period, as detailed below in section
149 2b. Post-processed model outputs were generated hourly using version 10.0.12 of the Unified
150 Post-Processing (UPP) package developed and maintained primarily by the NOAA EMC
151 (publicly available at <https://github.com/NOAA-EMC/UPP>); the UPP-processed outputs
152 include a standard suite of 2D forecast fields which are used for the verifications and
153 evaluations performed in this study. All computation was done on the Frontera supercomputer
154 operated by the Texas Advanced Computing Center (TACC).



155

156 Fig. 1. Map showing the extent of the native grid (black box) and output grid (red box)
 157 used for the 2022 CAPS FV3-LAM ensemble. The output grid is interpolated from the native
 158 FV3-LAM grid which uses the Extended Schmidt Gnomonic grid (a cubed-sphere projection)
 159 to a Lambert conformal map projection.

160

161 The CAPS ensemble consists of 21 FV3-LAM forecast members (summarized in Table 1);
 162 diversity among ensemble members is achieved via a combination of variation in physics
 163 configurations, IC and LBC perturbations, and the use of stochastic physics perturbations. It
 164 should be noted that, while physics, IC/LBC, and stochastic perturbations are all intended to
 165 increase ensemble spread, the goal is not merely to maximize spread; over-dispersion among
 166 ensemble members is just as undesirable as under-dispersion. In practice, however, many
 167 convection-allowing ensembles tend to be under-dispersive, hence the focus in this study on
 168 methods for increasing ensemble dispersion. All members use a version of FV3-LAM checked
 169 out from <https://github.com/NOAA-GSL/ufs-weather-model> on 30 March 2022 (tagged as
 170 “BaselineC-20220331”). The Rapid Radiative Transfer Model for GCMs (RRTMG; Mlawer
 171 et al. 1997) radiation parameterization is used in all members. For members using the NOAA
 172 or NOAA-MP LSM, soil temperature and moisture were initialized via interpolation from
 173 RRFS initial conditions containing RUC LSM variables. For members using the NSSL

174 microphysical scheme, microphysical variables not present in IC data are initialized via
 175 estimation using a gamma distribution. No cumulus parameterization scheme is used.

176

Experiment	Microphysics	PBL	Surface	LSM	IC/BC Source
<i>Baseline member using GFS initial conditions</i>					
M0B0L0_PG	Thompson	MYNN	MYNN	NOAH	GFS/GFS
<i>Multi-Physics configurations: same initial and boundary conditions (P members)</i>					
M0B0L0_P	Thompson	MYNN	MYNN	NOAH	RRFS CNTL/GFS
M1B0L0_P	NSSL	MYNN	MYNN	NOAH	RRFS CNTL/GFS
M0B0L1_P	Thompson	MYNN	GFS	NOAH-MP	RRFS CNTL/GFS
M1B2L2_P	NSSL	TKE-EDMF	GFS	RUC	RRFS CNTL/GFS
M0B2L1_P	Thompson	TKE-EDMF	GFS	NOAH-MP	RRFS CNTL/GFS
<i>Physics + initial condition perturbation ensemble (PI sub-ensemble)</i>					
M0B0L0_PI	Thompson	MYNN	MYNN	NOAH	RRFS01/GEFS m1
M0B1L0_PI	Thompson	Shin-Hong	GFS	NOAH	RRFS02/GEFS m2
M0B2L1_PI	Thompson	TKE-EDMF	GFS	NOAH-MP	RRFS03/GEFS m3
M0B0L1_PI	Thompson	MYNN	GFS	NOAH-MP	RRFS04/GEFS m4
M0B2L2_PI	Thompson	TKE-EDMF	GFS	RUC	RRFS05/GEFS m5
M1B0L0_PI	NSSL	MYNN	MYNN	NOAH	RRFS06/GEFS m6
M1B1L0_PI	NSSL	Shin-Hong	GFS	NOAH	RRFS07/GEFS m7
M1B2L1_PI	NSSL	TKE-EDMF	GFS	NOAH-MP	RRFS08/GEFS m8
M1B0L1_PI	NSSL	MYNN	GFS	NOAH-MP	RRFS09/GEFS m9
M1B2L2_PI	NSSL	TKE-EDMF	GFS	RUC	RRFS10/GEFS m10
<i>Physics + initial condition + stochastic perturbation ensemble (PSI sub-ensemble)</i>					
M0B0L0_PSI	Thompson	MYNN	MYNN	NOAH	RRFS01/GEFS m1
M1B0L0_PSI	NSSL	MYNN	MYNN	NOAH	RRFS06/GEFS m6
M0B0L1_PSI	Thompson	MYNN	GFS	NOAH-MP	RRFS04/GEFS m4
M1B2L2_PSI	NSSL	TKE-EDMF	GFS	RUC	RRFS10/GEFS m10
M0B2L1_PSI	Thompson	TKE-EDMF	GFS	NOAH-MP	RRFS03/GEFS m3

177 Table 1. Member configurations for the 2022 CAPS FV3-LAM ensemble. All members
 178 were run using FV3-LAM, and use the RRTMG radiation parameterization scheme. Member
 179 names contain information regarding the configuration of each member: the number after
 180 “M” indicates the microphysical scheme, the number after “B” indicates the PBL scheme,
 181 and the number after “L” indicates the land surface model. Suffixes indicate which sub-
 182 ensemble a member belongs to: “P” for members which are part of the core configurations,
 183 “PI” for members in the initial condition perturbation ensemble which uses initial and
 184 boundary conditions derived from GEFS RRFS, and “PSI” for members in the sub-ensemble
 185 using stochastic physics perturbations in addition to the perturbations of the “PI” members.
 186 The baseline configuration member, which uses GFS initial conditions, has a “PG” suffix.

187

188 The CAPS ensemble includes one member using a baseline configuration, and three groups
 189 of members investigating different aspects of ensemble design, including five members with
 190 core physics configurations and two sub-ensembles investigating different ensemble

191 perturbation strategies. The baseline configuration, labelled “M0B0L0_PG” in Table 1, uses
192 the Thompson microphysical scheme, MYNN PBL and surface layer schemes, and the NOAH
193 LSM, and uses ICs and LBCs from operational GFS. A group of five members were run using
194 core configurations focusing on physics diversity (hereafter referred to as the “physics” or “P”
195 members, the latter designation referring to the suffixes in the naming convention used in Table
196 1). Each of these members uses identical ICs and LBCs, but employs a unique combination of
197 microphysics, PBL schemes, and LSMs (Table 1). The physics combinations examined by the
198 P members are designed to resemble current and proposed convection-allowing model
199 configurations, including configurations similar to those of the GFS, the High-Resolution
200 Rapid Refresh (HRRR), the NCEP North American Model (NAM), and the Warn-on-Forecast
201 System (WoFS). LBCs used by the P members match those of the baseline configuration, and
202 ICs are obtained from the experimental RRFS prototype ensemble variational (EnVar) system.

203 Within the core configurations of the P members, two microphysics schemes are
204 examined—the National Severe Storms Laboratory (NSSL) scheme (Mansell et al. 2010;
205 Mansell and Ziegler 2013), and the Thompson scheme (Thompson et al. 2004, 2008;
206 Thompson and Eidhammer 2014). Three PBL schemes are examined, including the MYNN
207 scheme (Nakanishi and Niino 2009), the scale-aware Shin-Hong scheme (Shin and Hong
208 2015), and the TKE-EDMF scheme (Han and Bretherton 2019). The LSMs investigated
209 include the NOAH (Ek et al. 2003), NOAH-MP (Niu et al. 2011), and RUC (Smirnova et al.
210 2016) LSMs. For further discussion of the implementation of these schemes within the FV3-
211 LAM version used in this study, we refer the reader to Supinie et al. (2022).

212 One sub-ensemble adds IC and LBC diversity (hereafter referred to as the “IC/LBC
213 perturbation” or the “PI” sub-ensemble). In the PI sub-ensemble, initial conditions for each
214 member are provided by the first 10 members of an experimental RRFS EnKF data assimilation
215 (DA) system produced by GSL (Liu et al. 2023), and LBCs are obtained from GEFS member
216 forecasts. This sub-ensemble also includes physics diversity, with five of the ten members
217 using physics combinations matching those of the five physics sub-ensemble members.

218 The other sub-ensemble includes five members with stochastic physics perturbations in
219 addition to physics and IC/LBC diversity (hereafter referred to as the “stochastic” or the “PSI”
220 sub-ensemble). Stochastic perturbations are generated using a combination of SPPT, SKEB
221 and SHUM, using a unique random seed for each member. SPPT perturbs the physics tendency
222 terms in the equations of the U and V wind components, temperature, and water vapor mixing

223 ratio using a standard deviation of 0.7 in the red noise process used to generate stochastic
224 perturbations (Buizza et al. 1999). SKEB perturbations are applied to the U - and V -component
225 wind fields with a standard deviation of 0.5 m s^{-1} . SHUM is applied within the boundary layer
226 to perturbations of the 3-dimensional specific humidity field using a normalized perturbation
227 coefficient of 0.006. For all stochastic methods, perturbations have a characteristic spatial
228 length scale of 150 km and a characteristic time scale of six hours, and are updated hourly.
229 These parameters were inherited from defaults in the UFS Short-range Weather Application
230 package. We believe that these settings are suitable for applications involving convective
231 rainfall, and note that they are similar, at least in terms of time and length scales, to those used
232 in prior studies (e.g., Jankov et al. 2017). The five PSI sub-ensemble members have the same
233 physics combinations as the P sub-ensemble members. The physics combinations of the five
234 PSI sub-ensemble members also match those of five corresponding PI ensemble members,
235 allowing for a clean comparison of these PI and PSI members to evaluate the impact of
236 stochastic physics perturbations.

237 *b. Data and methods for forecast evaluation*

238 In this study, we verify and evaluate 0-84 hour forecasts of the CAPS ensemble and its sub-
239 ensembles, with a particular focus on forecasts of precipitation. The operational High
240 Resolution Ensemble Forecast (HREF) and GEFS forecasts are used as references for
241 comparison. The HREF is an “ensemble of opportunity” convection-allowing ensemble
242 forecasting system made up of 5 forecasting systems running at similar 3-4 km grid spacings;
243 through time-lagging HREF achieves 10 ensemble members (Jirak et al. 2012; Roberts et al.
244 2020), and is the current operational standard at convection-allowing resolution over the
245 CONUS. GEFS is included because HREF forecasts only run for 48 hours (allowing for only
246 36-h forecasts after the 12-h time-lag members are accounted for); a substantial dry bias in
247 precipitation at higher intensity thresholds is expected of the GEFS given its much coarser grid
248 spacing of 0.25 degrees. Evaluation of GEFS and HREF forecasts is not a goal of this study;
249 they are included purely as baselines against which to evaluate the CAPS forecast sub-
250 ensembles.

251 The evaluations presented in this study were generated using the `grid_stat` and
252 `ensemble_stat` modules of version 10.0.0 of the Model Evaluation Tools (MET) software
253 package (Brown et al. 2021, DTC 2023). The observation datasets used as “truth” for
254 precipitation forecast evaluation are the NCEP Stage-IV 1-h, 6-h, and 24-h precipitation

255 analyses (Nelson et al. 2016), which combine WSR-88D radar-estimated rainfall with rain
256 gauge observations (including manual error correction). Because Stage-IV data are only valid
257 over the United States, the verification is constrained to the land areas of the CONUS.

258 Surface temperature and wind forecasts are verified against the NOAA UnRestricted
259 Mesoscale Analysis (URMA). URMA is very similar to the Real Time Mesoscale Analysis
260 (RTMA; de Pondecia et al. 2011; Morris et al. 2020), but is generated 6 hours later to allow
261 inclusion of late-arriving observations not used in RTMA (de Pondecia et al. 2015). While
262 URMA does involve interpolation of observations to a grid we note that it does so via a
263 sophisticated, anisotropic 2DVAR analysis scheme and is designed to have high fidelity to
264 surface observations, making it suitable for validation and verification of near-surface forecast
265 fields (Supinie et al. 2022), and has previously been used operationally by NWS as a “truth
266 analysis” (de Pondecia et al. 2015). Verifications of vertical profiles for temperature, dewpoint,
267 and wind forecasts are performed against radiosonde observations (RAOB) at ~70 NWS
268 operational sites within the CONUS. For verification, all forecast data are interpolated to the
269 grid of the analyses—either the Stage-IV grid (which has horizontal grid spacing of
270 approximately 4.76 km) or the 2.5-km URMA grid—using MET’s mass-conserving
271 “BUDGET” interpolation method.

272 We examine the member-by-member performance of the CAPS ensemble using equitable
273 threat score (ETS), frequency biases, and performance diagrams for precipitation forecasts at
274 varying accumulated precipitation intensity thresholds, calculated using a 2×2 contingency
275 table computed from model precipitation forecasts and observed Stage-IV precipitation
276 accumulations. Ensemble consensus forecasts of precipitation are also evaluated, including the
277 probability-matched mean (PM mean; Ebert 2001), and the patchwise localized PM mean
278 (LPM mean; Snook et al. 2019, 2020) of the CAPS ensemble and its sub-ensembles. The simple
279 ensemble mean is not included because of well-known biases arising from the smoothing effect
280 of ensemble averaging. Accumulation periods of 6-h and 24-h are considered. We use both
281 fixed and percentile thresholds for evaluating precipitation forecasts. An accumulation
282 threshold of 1 mm is used as a proxy for geographic extent of precipitation, and a threshold of
283 25 mm is used to focus on prediction of moderate to heavy precipitation. For percentile
284 thresholds, the 99th percentile of accumulated precipitation is used to focus on heavy
285 precipitation, while the 90th percentile (for 24-h accumulations) and the 95th percentile (for 6-
286 h accumulations) are used to include regions of lighter precipitation. The use of percentile

287 thresholds allows us to eliminate the impact of frequency bias (Carafo et al. 2021). The
288 ensemble mean RMSEs and ensemble spread of predicted temperature, dewpoint, and zonal
289 wind are calculated against RAOB profile data to examine the performance of various
290 ensembles. Similar comparisons are done for 2-m temperature, 2-m dewpoint, and 10-m zonal
291 winds against URMA analyses.

292 Following Supinie et al. (2022), we choose for the sake of simplicity not to use a
293 neighborhood-based contingency table for many of our evaluations. When a neighborhood is
294 used, we apply the neighborhood maximum ensemble probability (NMEP; Schwartz and
295 Sobash 2017). Employing increasing neighborhood radii generally increased forecast skill, but
296 resulted in similar relative performance among ensembles and sub-ensembles, as we discuss
297 below in section 3b. The forecast evaluations discussed in this study use forecasts initialized
298 at 0000 UTC on each weekday from 13 May 2022 to 3 Jun. 2022 excluding 18 May, 30 May,
299 and 2 June (for which complete data were not available), totaling 13 days of data. While this
300 represents a small sample size and somewhat limits the robustness of the results of this study,
301 technical limitations precluded the generation of additional forecast runs: forecasts were not
302 run prior to 13 May due to unavailability of experimental RRFS initial conditions from GSL,
303 while a combination of missing RRFS initial conditions and other technical issues resulted in
304 incomplete forecast output on later missing days.

305

306 **3. Forecast evaluation results**

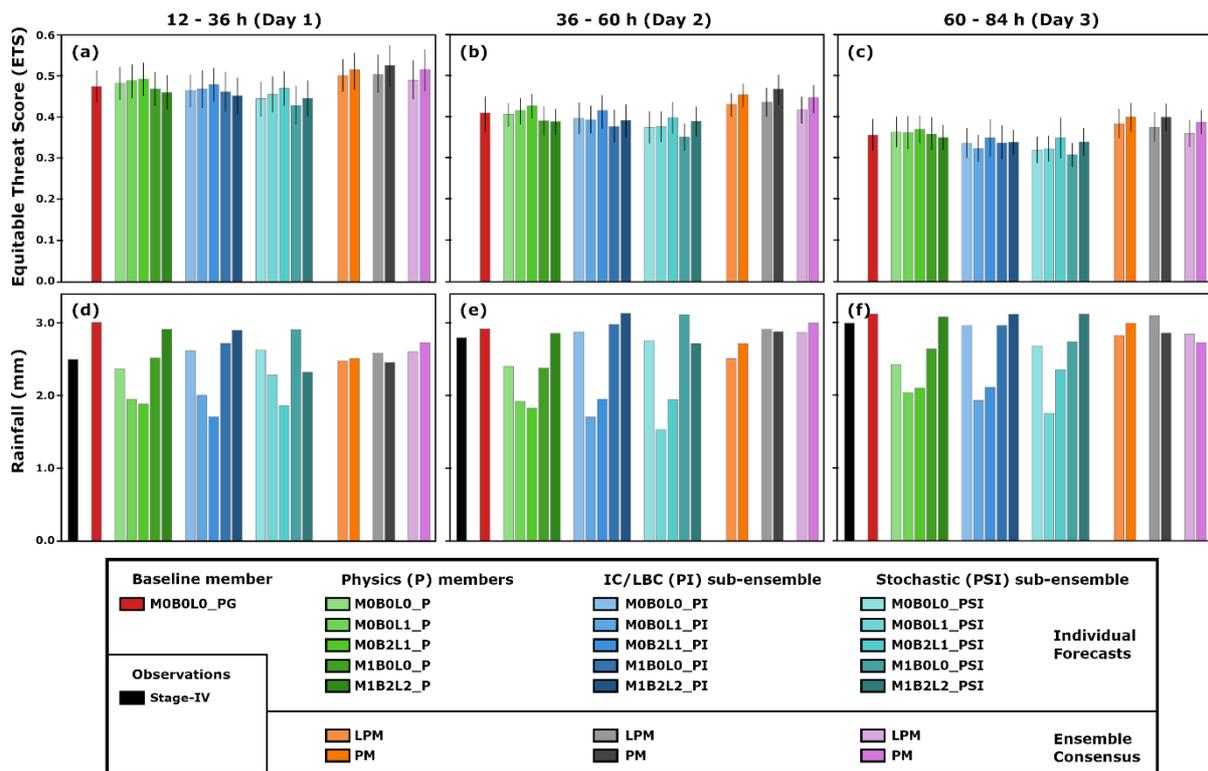
307 *a. Deterministic and probabilistic evaluations of precipitation forecasts*

308 There is little difference among individual members in ETS of precipitation exceeding the
309 90th percentile threshold (Fig. 2); all members exhibit ETS declining from between 0.4 and 0.5
310 on Day 1 (Fig. 2a) to between 0.3 and 0.4 on Day 3 (Fig. 2c). Within each sub-ensemble, the
311 members do not show significant differences, but the M0B2L1 member of each sub-ensemble
312 consistently exhibits the highest ETS for Days 1, 2, and 3. As expected, the PM and LPM
313 outperform individual members in terms of ETS.

314 Individual ensemble members do, however, exhibit substantial differences in their 90th
315 percentile values (Fig. 2d-f). The 90th percentile of observed Stage-IV precipitation ranges
316 from around 2.5 mm for Day 1 (Fig. 2d) to near 3.0 mm for Day 3 (Fig. 2f)—this variation in
317 observed 90th percentile is due to day-to-day variation in total precipitation accumulation. In

12

318 contrast, the 90th percentile of individual member forecasts ranges from around 1.7 to 3.1 mm—
 319 for most members, the 90th percentile of precipitation is slightly to substantially below the
 320 observed 90th percentile. In all sub-ensembles and at all lead-times, the members using the
 321 NOAH-MP LSM (“L1” in the ensemble member naming scheme; see Table 1) consistently
 322 exhibit the lowest 90th percentile precipitation value, while members using the NSSL
 323 microphysics scheme (“M1” members; see Table 1) tend to exhibit the highest. When sub-
 324 ensemble consensus measures are considered, however, the 90th percentile of sub-ensemble
 325 PM and LPM means are generally in good agreement with the observed 90th percentile value
 326 (Fig. 2d-f). Comparing the PI and PSI sub-ensembles, the inclusion of stochastic perturbations
 327 has little impact on ETS or 90th percentile value; no significant differences are noted.
 328 Furthermore, the consistent performance across sub-ensembles in terms of relative ETS and
 329 90th percentile value among physics configurations suggests that the choice of physics
 330 configuration has a stronger impact on the precipitation forecasts at the 90th percentile threshold
 331 than either IC/LBC perturbations or stochastic perturbations.

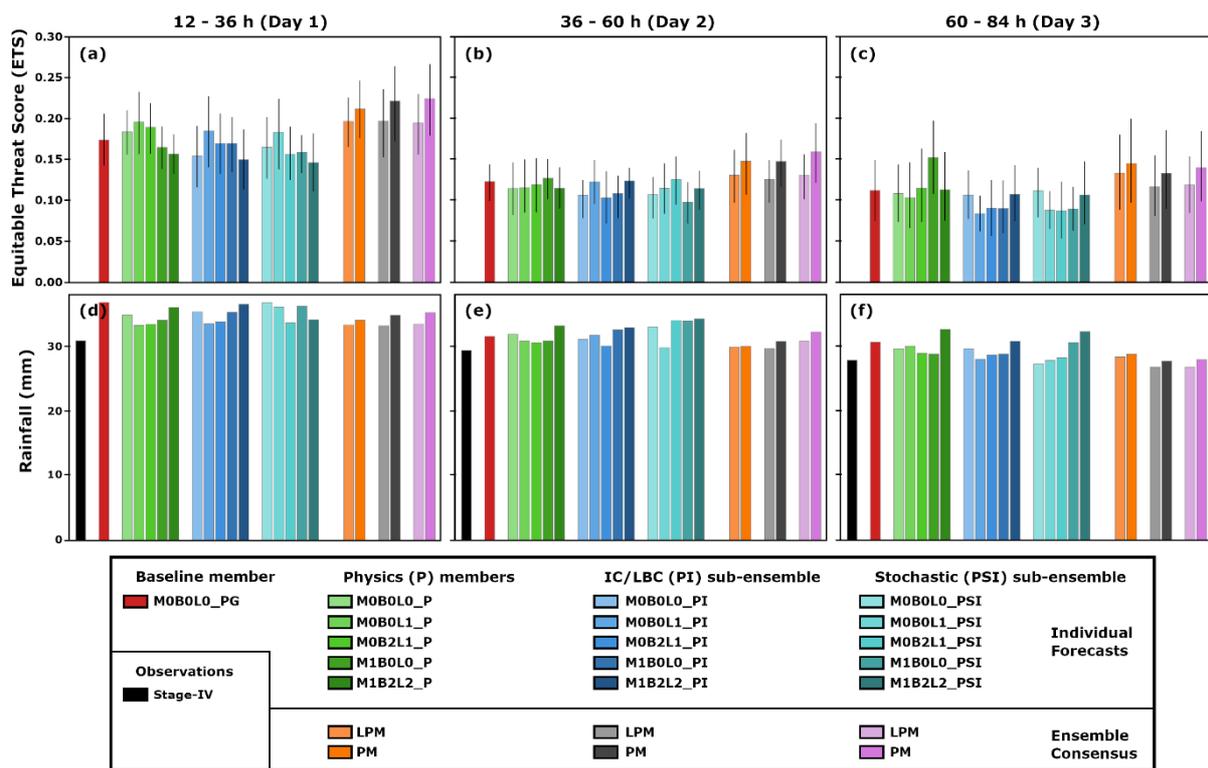


332

333 Fig. 2. (a-c) 90th percentile equitable threat score (ETS) and (d-f) 90th percentile rainfall
 334 amount for 24-h accumulated precipitation forecasts for (a, d) Day 1 (12-36 h), (b, e) Day 2
 335 (36-60 h), and (c, f) Day 3 (60-84 h) forecasts, aggregated over all available cases from the
 336 2022 HWT SFE period for the CAPS ensemble. Error bars for each member in panels (a-c)
 337 indicate the 5.0-95.0 percentile range using bootstrap resampling (10,000 resamples). The
 338 black bar in panels (d-f) indicates the 90th percentile of Stage-IV observed rainfall.

339

340 When we consider the 99th percentile for 24-h accumulated precipitation (Fig. 3), overall
 341 ETS values are somewhat lower than those at the 90th percentile, ranging from approximately
 342 0.15-0.20 on Day 1 (Fig. 3a) to around 0.10 for most members on Day 3 (Fig. 3c). As at the
 343 90th percentile threshold, differences among members are not significant, and the performance
 344 of individual members within the PI and PSI sub-ensembles is similar. Unlike at the 90th
 345 percentile, however, the 99th percentile of 24-h accumulated precipitation is higher than that of
 346 the observed Stage-IV precipitation for nearly all members on all three days (Fig. 3d-f). Taken
 347 together with the 90th percentile statistics presented in Fig. 2d-f, we find that most members
 348 have precipitation distributions which are lighter than observed for light to moderate rainfall,
 349 but which have heavier than observed accumulations in the most intense precipitation regions.
 350 These results align with the findings of prior studies, such as Zhou et al. (2019), which have
 351 indicated that FV3-based models tend to over-predict heavy precipitation.



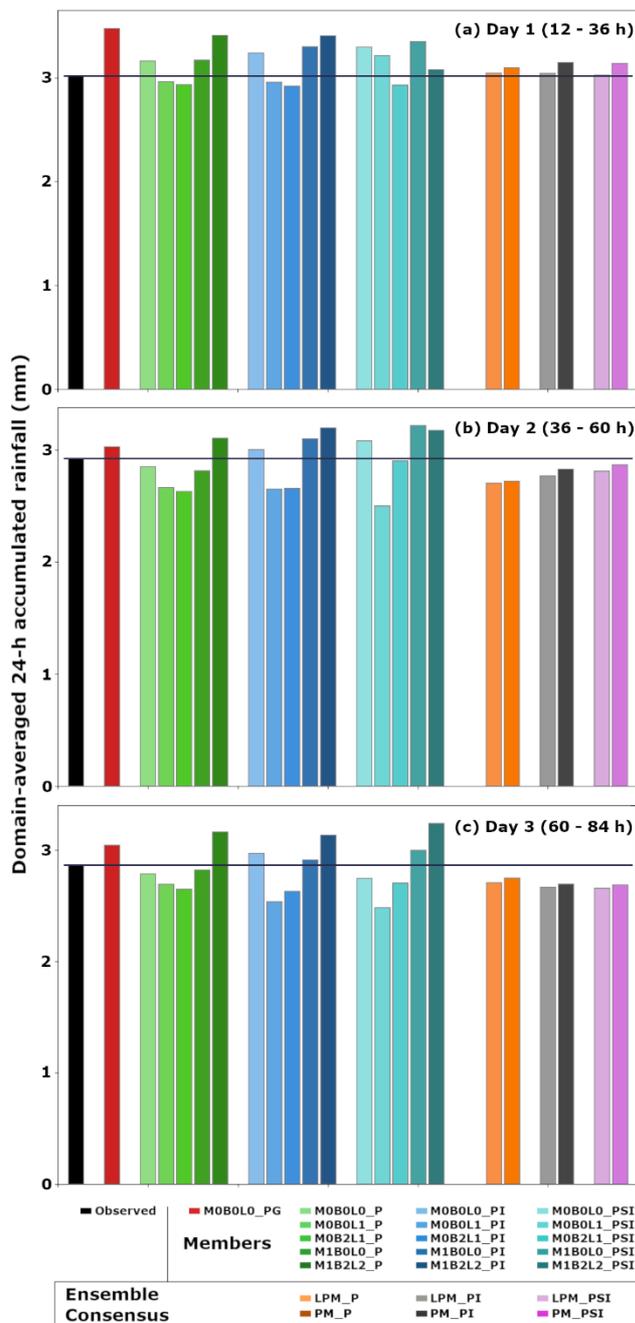
352

353 Fig. 3. As Fig. 2, but for the 99th percentile instead of the 90th.

354

355 The conclusions drawn from Figs. 2, 3 regarding precipitation distributions can be
 356 contextualized via an analysis of average total 24-h accumulated rainfall per Stage-IV grid cell
 357 (Fig. 4). In terms of total 24-h accumulated rainfall, the “M1” members, which use the NSSL

358 microphysical scheme, stand out as having more total rainfall than the “M0” members, which
 359 use the Thompson microphysical scheme—particularly for M1B2L2. This pattern is consistent
 360 across sub-ensembles (i.e., it is present for “M1” members in the “P”, “PI”, and “PSI” sub-
 361 ensembles), and results in the “M1” members consistently producing more total rainfall than
 362 the Stage-IV analysis (Fig. 4). The “L1” members, which use the NOAH-MP land-surface
 363 model, are consistently among the lowest in terms of total rainfall.
 364



365

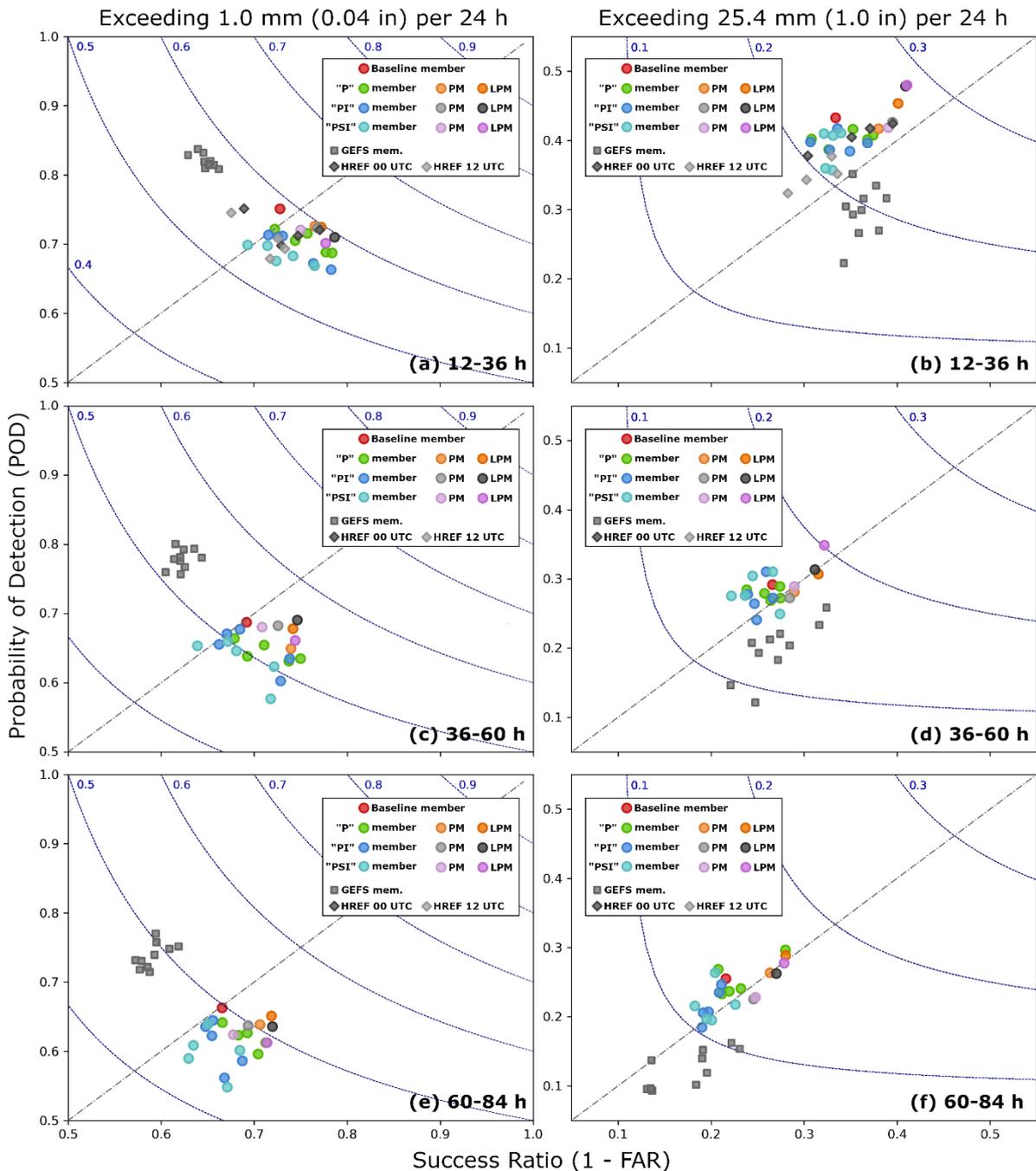
366 Fig. 4. Domain-averaged total 24-h accumulated precipitation for (a) 12-36 h, (b) 36-60
367 h, and (c) 60-84 h of forecast time for selected CAPS FV3 members and ensemble/sub-
368 ensemble consensus forecasts (PM and LPM mean), compared to that of the Stage-IV
369 analysis (“observed”). Color-coding of the plotted bars matches that used in Figs. 2 and 3.
370 The horizontal black line in each panel is plotted at the observed value (to facilitate visual
371 comparison of forecasts with observations).

372

373 Performance of 24-h rainfall accumulation forecasts is also examined in Fig. 5 via
374 performance diagrams (Roebber 2009) using thresholds of 1.0 mm and 25.4 mm. Forecast
375 performance improves toward the upper right of the diagram; forecasts along the one-to-one
376 diagonal are unbiased. The performance diagrams shown in Fig. 5 use fixed rather than
377 percentile-based precipitation thresholds, allowing forecast bias assessment. At the 1.0 mm
378 threshold (Fig. 5a, c, e), GEFS members exhibit a high bias, while the CAPS FV3 and HREF
379 members and CAPS sub-ensemble PM and LPMs exhibit near-neutral or slightly low biases.

380 At the 25.4 mm (1.0 inch) threshold (Fig. 5b, d, f), forecast accuracy is generally quite
381 similar between the PI and PSI ensembles, with neither clearly outperforming the other. At
382 this higher threshold, HREF members perform similarly to or slightly better than the PI and
383 PSI members; both the CAPS sub-ensembles and the HREF exhibit a slight over-prediction
384 bias on day 1 (Fig. 5b), and close-to-unbiased performance on days 2 and 3 (Fig. 5d, f). GEFS
385 performs noticeably worse in both CSI and bias compared to the CAPS sub-ensembles and
386 HREF, which is unsurprising given its much lower resolution. The overall performance of the
387 CAPS sub-ensembles is similar to that noted by Johnson et al. (2023) for CAPS forecasts
388 produced during the 2018 HWT SFE.

389



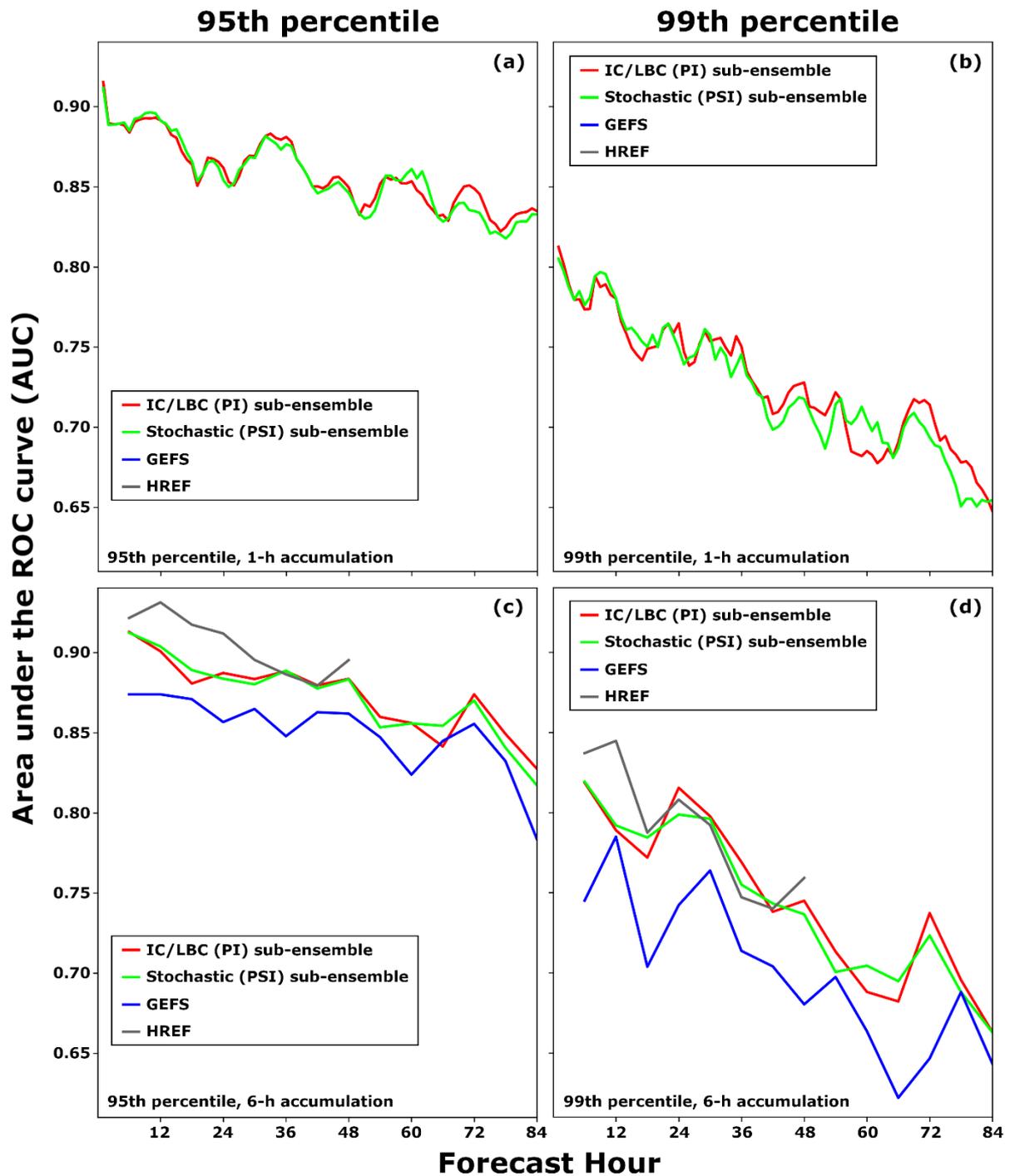
390

391 Fig. 5. Performance diagrams (zoomed-in) for forecasts of 24-h accumulated
 392 precipitation exceeding (a, c, e) 0.04 inches or (b, d, f) 1.00 inches at forecast lead times of
 393 (a, b) 12-36 h, (c, d) 36-60 h, and (e, f) 60-84 h using all available days from the 2022 HWT
 394 SFE period. The dark blue labelled contours indicate lines of constant critical success index
 395 (CSI), while the diagonal dotted line indicates a constant bias of 1.0. Color-coding of
 396 ensemble members is indicated by the legends in panel (a). Due to the shorter run-length of
 397 HREF, HREF data are shown only in panels (a) and (b). To facilitate readability, the left and
 398 right columns are zoomed in different portions of the performance diagram.

399

400 The relative performances of the five-member PI and PSI sub-ensembles are compared to
401 one another and to those of the five members of HREF without time-lagging and the five GEFS
402 members which provided boundary conditions for the PSI sub-ensemble (see Table 1) in terms
403 of area under the receiver operating characteristic (ROC) curve (AUC) (Mason 1982) for
404 accumulated precipitation in Fig. 6. Shown are verification for 1-h (Fig. 6a, b) and 6-h (Fig.
405 6c, d) accumulation at 95th and 99th percentile thresholds. Due to limitations in the temporal
406 resolution and format of available data, 1-h verification is shown only for the PI and PSI sub-
407 ensembles. AUC measures resolution (i.e., the ability to distinguish between events and non-
408 events), and is constructed by calculating the area beneath the curve when plotting probability
409 of detection (POD) versus probability of false detection (POFD) (e.g., Clark et al. 2021).

410 At both percentile thresholds, the CAPS PI and PSI sub-ensembles exhibit qualitatively and
411 quantitatively similar performance in terms of AUC, with only very minor differences between
412 the two (Fig. 6). Overall patterns of forecast skill are similar for 1-h and 6-h accumulations,
413 though the 1-h accumulations exhibit additional evidence of diurnal variation in forecast skill
414 in the PI and PSI sub-ensembles, with higher forecast skill between 0600 and 1200 UTC and
415 lower forecast skill between 1800 and 0000 UTC; this pattern is particularly notable at the 95th
416 percentile threshold (Fig. 6a). When the CAPS PI and PSI sub-ensembles are compared to the
417 operational HREF and GEFS ensembles, HREF exhibits superior performance in terms of AUC
418 up to 12 hours of forecast time, and then performs similarly to or better than the CAPS sub-
419 ensembles through the remainder of its forecast period (ending at 48 h of forecast time). The
420 CAPS sub-ensembles do outperform GEFS, particularly at the 99th percentile threshold (Fig.
421 6d). The superior performance of the CAPS sub-ensembles and HREF compared to GEFS is
422 not surprising, given the large resolution difference (~3 km vs. ~25 km). The inferior
423 performance of coarser-resolution global models in predicting precipitation (especially heavy
424 precipitation) compared to convection-allowing models has been well-documented in previous
425 studies (e.g., Zhu et al. 2018).

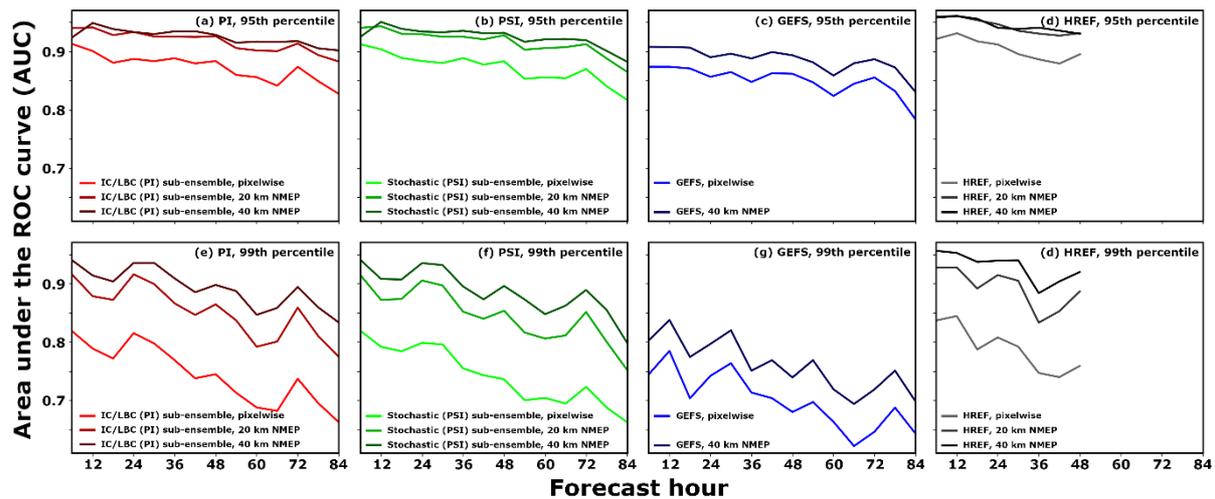


426

427 Fig. 6. Area under the receiver operating characteristic curve (ROC AUC) for forecasts
 428 of (a, b) 1-h accumulated precipitation and (c, d) 6-h accumulated precipitation exceeding (a,
 429 c) the 95th percentile or (b, d) the 99th percentile at various forecast lead times. Shown are the
 430 CAPS PI and PSI 5-member sub-ensembles compared to the five members of GEFS which
 431 provided boundary conditions for the CAPS FV3 PSI sub-ensemble and the 0000 UTC
 432 HREF. Values shown here are calculated using all days during the 2022 HWT SFE period
 433 with valid CAPS ensemble data. Due to the shorter run-length of HREF, HREF is shown
 434 only for forecast times of 6 to 36 hours.

435

436 The impact of using a neighborhood during probabilistic forecast generation is examined
 437 in Fig. 7, which compares the AUC of forecasts generated using a pixelwise (no neighborhood)
 438 ensemble approach to those generated using ensemble NMEP with neighborhood radii of 20
 439 or 40 km. Only 0 and 40 km radii are used for GEFS, as a neighborhood of 20 km is smaller
 440 than the grid spacing of GEFS. Across neighborhood radii, HREF consistently exhibits the
 441 highest AUC, followed by the CAPS PI and PSI sub-ensembles which exhibit similar AUC
 442 performance to one another, slightly worse than HREF. GEFS consistently exhibits the lowest
 443 AUC of the ensembles examined. The impact of using a neighborhood via NMEP is uniformly
 444 positive in terms of AUC, but the increase in AUC is greater for the CAPS sub-ensembles and
 445 the HREF than it is for GEFS (Fig. 7c, g), likely because the relatively coarse GEFS grid results
 446 in fewer localized heavy rainfall areas for which forecasts would be greatly improved by using
 447 a 40 km neighborhood. The CAPS sub-ensembles and the HREF, which both have native
 448 horizontal grid spacing of around 3 km, exhibit similar magnitudes of increase in AUC going
 449 from a 0 km to 20 km to 40 km neighborhood.



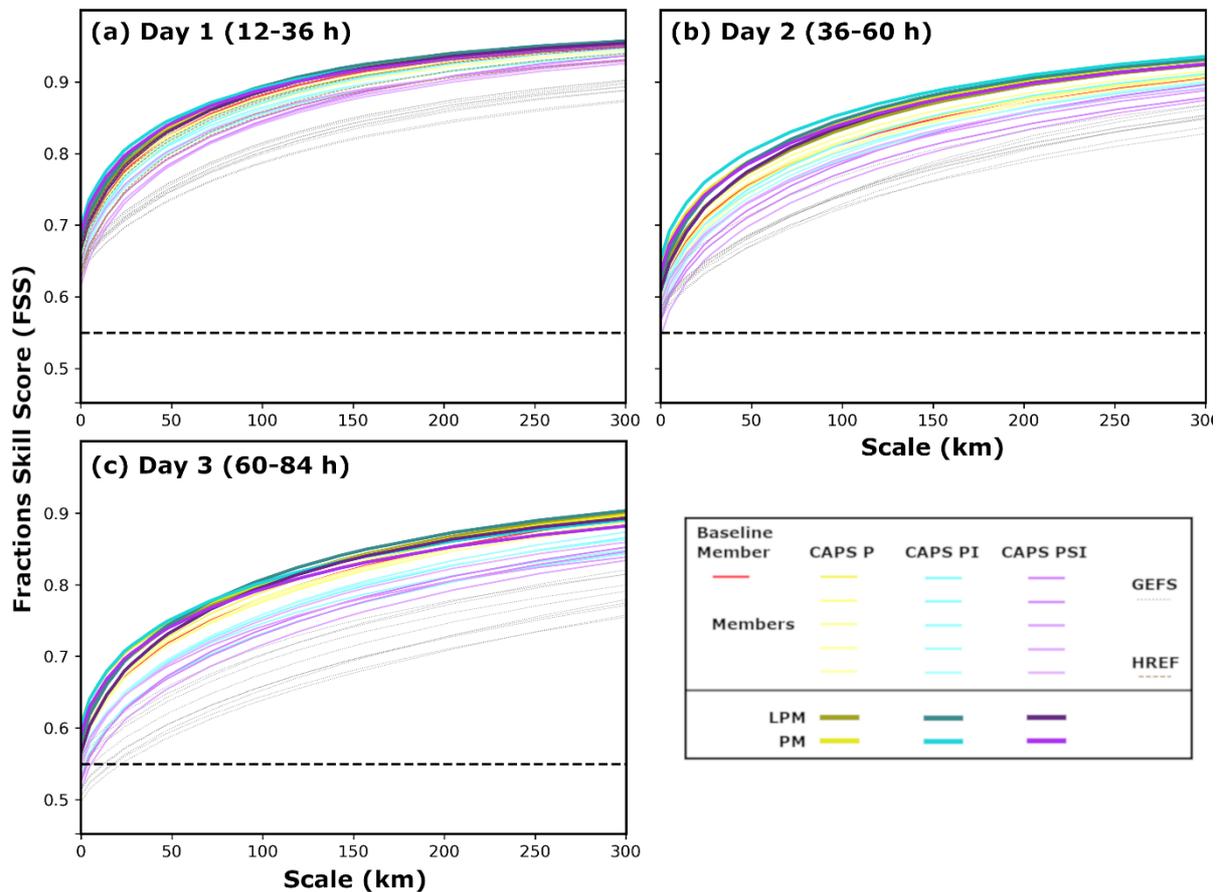
450

451 Fig. 7. Area under the receiver operating characteristic curve (ROC AUC) for forecasts
 452 of 6-h precipitation accumulation exceeding (a-d) the 95th percentile and (e-h) the 99th
 453 percentile at forecast lead times ranging from 6 hours to 84 hours for (a, e) the 5-member
 454 CAPS PI and (b, f) PSI sub-ensembles, as well as (c, g) the five members of GEFS which
 455 provided boundary conditions for the CAPS FV3 PSI sub-ensemble and (d, h) the 0000 UTC
 456 HREF. Curves are plotted using NMEP with neighborhood radii of 0, 20, and 40 km for PI,
 457 PSI, and HREF, and with radii of 0 and 40 km for GEFS. As in Fig. 6, values are calculated
 458 using all days during the 2022 HWT SFE period with valid CAPS ensemble data.

459

460 Fractions Skill Score (FSS; Roberts and Lean 2008) is another neighborhood-focused
 461 forecast skill metric which evaluates the degree of similarity between a forecast and

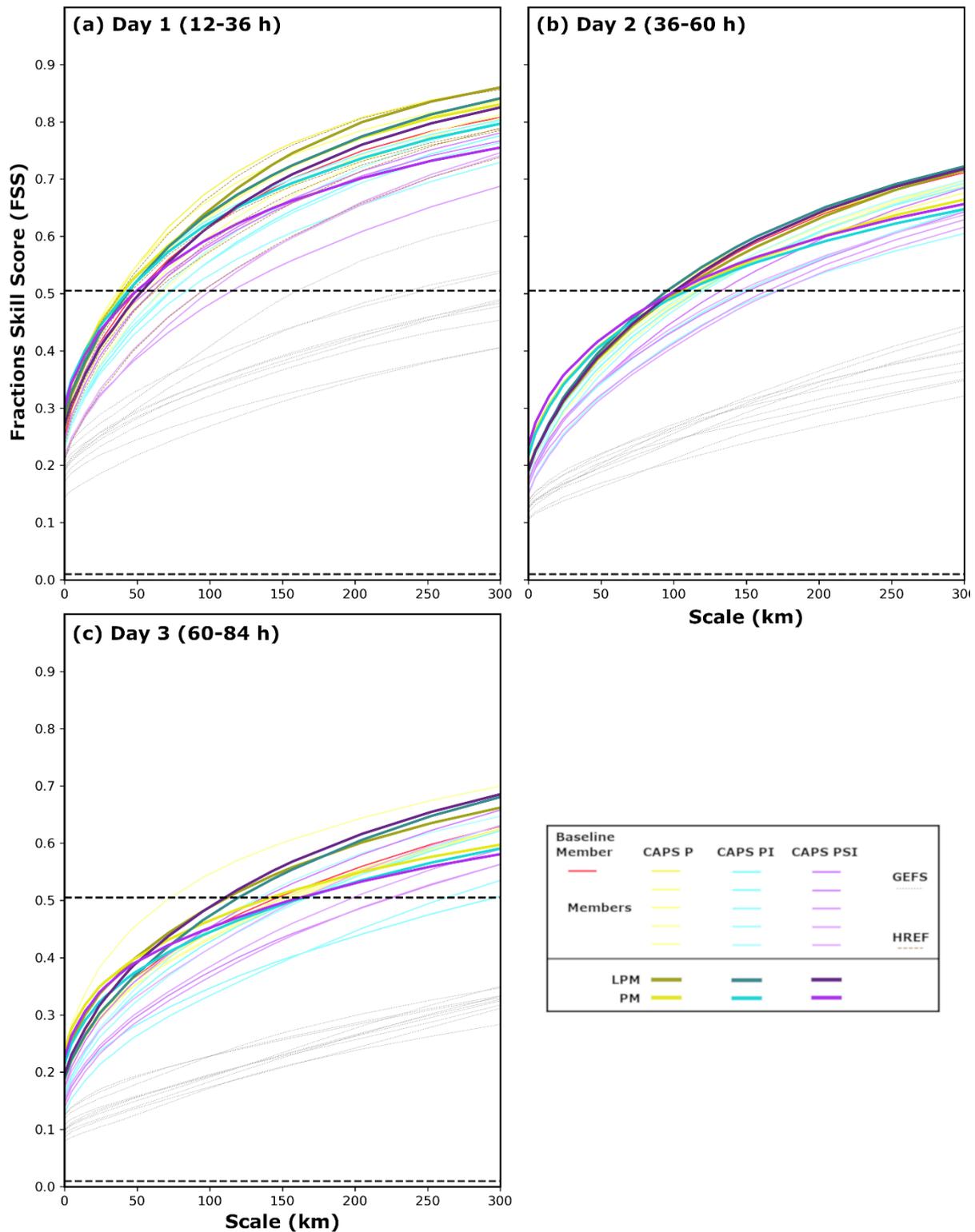
462 observations at varying spatial scales. FSS varies between 0.0 and 1.0, with higher values
 463 indicating higher forecast skill. A forecast can be considered useful at a given scale if the FSS
 464 at that scale exceeds that which would be obtained at the grid scale from a forecast with
 465 probability at each point equal to the base rate of the event being evaluated (Roberts and Lean
 466 2008). At the 90th percentile (Fig. 8), all CAPS forecast members exhibit skillful forecasts of
 467 24-h accumulated precipitation in terms of FSS, with modest variation among members. In
 468 terms of ensemble consensus measures, the PM exhibits slightly higher skill than the LPM in
 469 terms of FSS at the 90th percentile, with the PI sub-ensemble slightly outperforming the PSI
 470 sub-ensemble. Individual HREF members perform similarly to the best-performing CAPS
 471 members on Day 1 (Fig. 8a), while GEFS members consistently exhibit lower FSS than either
 472 CAPS or HREF members, although even GEFS is skillful at all scales on Day 1 and 2, and at
 473 all but the smallest scales on Day 3.



475 Fig. 8. Fractions Skill Score (FSS) for forecasts of 24-h precipitation accumulation
 476 exceeding the 90th percentile for (a) Day 1, (b) Day 2, and (c) Day 3 forecasts for individual
 477 CAPS ensemble members, the PM and LPM of the CAPS P, PI and PSI sub-ensembles, and
 478 corresponding members of the HREF and GEFS ensembles for comparison. The horizontal
 479 dashed line indicates the minimum useful value of FSS.

480

481 At the 99th percentile (Fig. 9), there is more substantial variation in FSS across the CAPS
482 members, although neither PI nor PSI members consistently outperform one another, while
483 HREF members perform comparably to better-performing members of the CAPS sub-
484 ensembles (Fig. 9a). GEFS performs notably worse than either the CAPS or HREF members.
485 By Day 2, no GEFS member exhibits a skillful FSS at any scale up to 300 km (Fig. 9b, c).



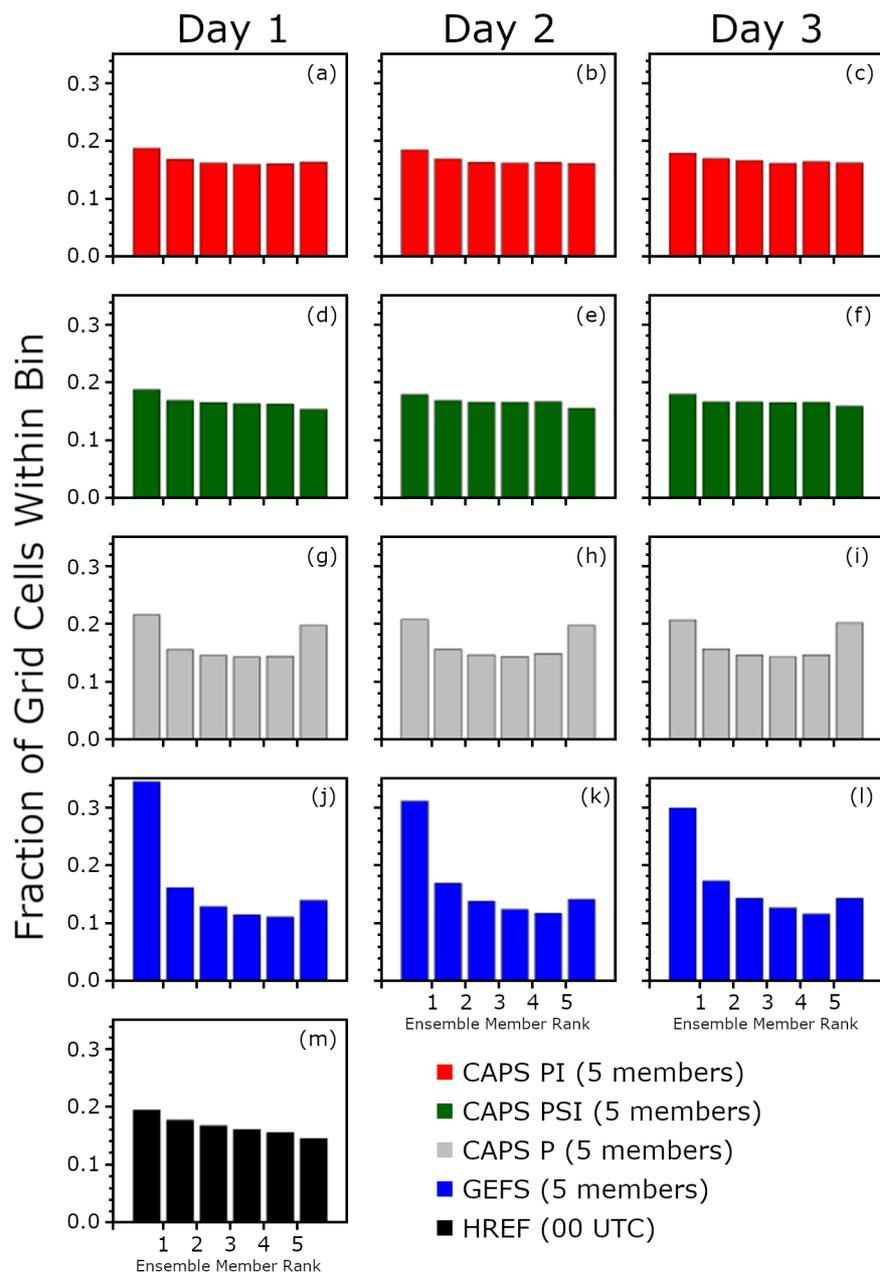
486

487 Fig. 9. As Fig. 8, but for the 99th percentile of 24-h accumulated rainfall instead of the
 488 90th percentile. The lower horizontal dashed line in each panel indicates the base rate, while
 489 the upper horizontal dashed line indicates the minimum useful value of FSS.

490

491 In Fig. 10, rank histograms are presented for forecasts of 24-h accumulated precipitation.
492 Rank histograms measure ensemble reliability. Ideally, the fraction of points within each rank
493 should be similar, leading to a flat rank histogram (Hamill 2001). The PI and PSI sub-
494 ensembles (Fig. 10a-f) exhibit nearly flat rank histograms for 24-h accumulated precipitation
495 for days 1-3, with a slight downward slope towards higher ranks, indicating slight over-
496 prediction of precipitation. When the CAPS P members are treated as a sub-ensemble, having
497 no initial condition perturbations and relying entirely upon physics diversity to generate
498 ensemble spread, under-dispersion is evident in the form of a U-shaped rank histogram (Fig.
499 10g-i), underscoring the importance of the initial condition perturbations applied in the PI and
500 PSI sub-ensembles. Given the similarity between the PI and PSI sub-ensembles, it is likely
501 that the initial condition perturbations impart greater impact than the stochastic perturbations.
502 In contrast to the CAPS sub-ensembles, the GEFS 5-member sub-ensemble (Fig. 10j-l) shows
503 a substantial over-forecasting bias on all days, with the strongest bias on day 1 (Fig. 10j). The
504 HREF 0000 UTC members (Fig. 10m) exhibit better rank histogram behavior than GEFS, but
505 with slightly greater over-forecasting bias than the CAPS sub-ensembles (evidenced by the
506 decreasing occurrence of grid cells at increasing ranks in HREF, compared to the flatter rank
507 histogram of the CAPS PI and PSI sub-ensembles).

508



509

510 Fig. 10. Rank histograms of 24-h accumulated precipitation for (a-c) the CAPS PI sub-
 511 ensemble, (d-f) the CAPS PSI sub-ensemble, (g-i) the CAPS P sub-ensemble, (j-l) the five
 512 members of the GEFS which provided boundary conditions for the CAPS PSI 5-member sub-
 513 ensemble, and (m) the 00 UTC members of the HREF. Data are shown for forecast lead
 514 times of 12-36 h (left column), 36-60 h (center column), and 60-84 h (right column).

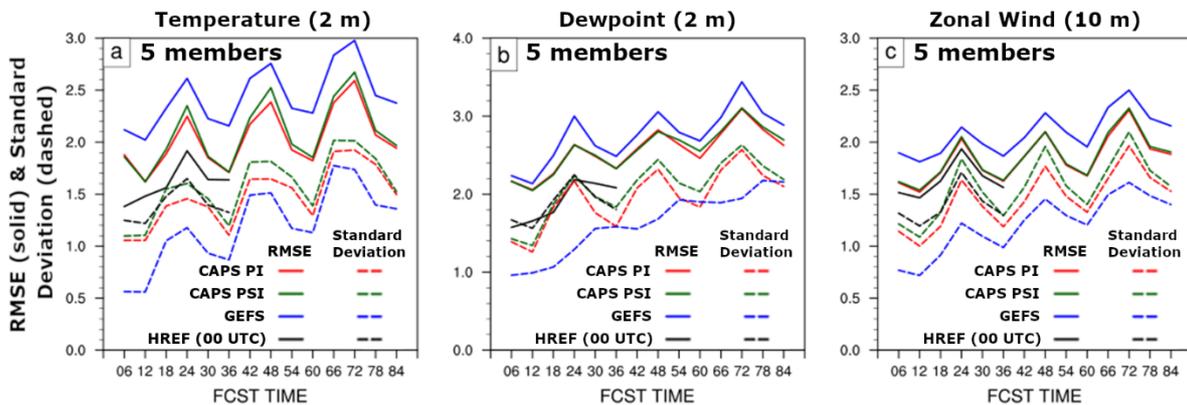
515

516 *b. Evaluations of wind, temperature, and dewpoint forecasts*

517 For 2 meter temperature (T), 2 meter dewpoint (T_d), and 10 meter zonal wind (U), the 5-
 518 member CAPS PI and PSI sub-ensembles consistently exhibit lower RMSEs and greater
 519 ensemble spread than GEFS, while HREF has the lowest RMSEs and highest spread

25

520 throughout that the period for which HREF forecasts are available (Fig. 11). RMSE and spread
 521 curves exhibit a strong diurnal variation for T , T_d , and U , with maxima typically occurring
 522 around 0000 UTC (during the afternoon and evening hours over the CONUS when convective
 523 activity is typically near its peak). RMSEs of T_d and U are very similar between the PI and PSI
 524 sub-ensembles, while RMSE of T is slightly higher in the PSI sub-ensemble during the diurnal
 525 maximum. While statistics are not shown for the meridional wind component, performance is
 526 qualitatively similar to that of U (not shown). For all three variables ensemble spread is
 527 consistently slightly higher for the PSI sub-ensemble than for the PI sub-ensemble, indicating
 528 that stochastic perturbations do help improve the ensemble spread, with only small increases
 529 in RMSE at a few times. A small improvement in spread to RMSE ratio in the PSI sub-
 530 ensemble compared to the PI sub-ensemble is present throughout the forecast period out to the
 531 maximum forecast lead-time of 84 h.



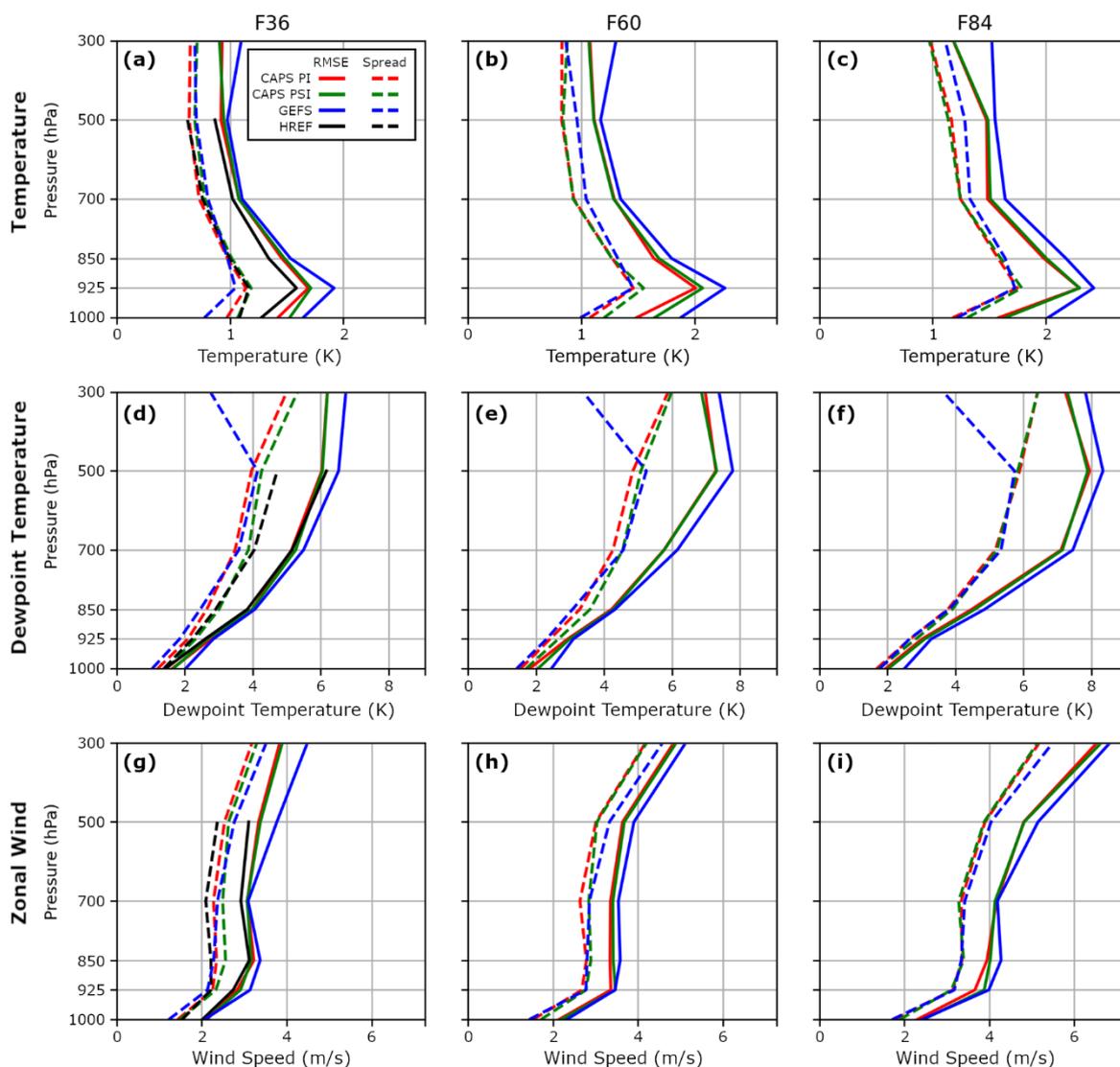
532

533 Fig. 11. RMSE and ensemble spread for forecasts of (a) 2-meter temperature (K), (b) 2-
 534 meter dewpoint (K), and (c) 10-meter zonal (u) wind component (m s^{-1}). Shown are the 5-
 535 member CAPS PI and PSI 5-member sub-ensembles compared to the five members of GEFS
 536 which provided boundary conditions for the CAPS FV3 PSI sub-ensemble and the 0000 UTC
 537 HREF.

538

539 Vertical profiles of RMSE and spread calculated against sounding data are plotted for T ,
 540 T_d , and U at 36, 60, and 84 hours of forecast time in Fig. 12. While some differences between
 541 the performance of the CAPS PI and PSI sub-ensembles were noted near the surface (Fig. 11),
 542 in Fig. 12 the CAPS PI and PSI sub-ensembles generally exhibit similar performance, both
 543 performing slightly better than the GEFS. The most notable difference between the PI and PSI
 544 sub-ensembles is a slightly improved spread to RMSE ratio in the PSI sub-ensemble, primarily
 545 below 500 hPa. At 36 hours, HREF exhibits slightly smaller RMSEs for T and U (Figs. 12a,g),
 546 slightly smaller spread for U (Fig. 12g) and larger spread for T_d (Fig. 12d) than other ensembles.

547 Overall, HREF has better error-spread consistency, especially for T_d . Spread in temperature is
 548 slightly higher in GEFS than in the HREF or CAPS sub-ensembles at 60 and 84 hours, albeit
 549 with larger RMSE; this may reflect better tuning in the stochastic physics used by GEFS than
 550 in the CAPS PSI sub-ensemble.



551
 552 Fig. 12. Vertical RMSE and spread profiles of (a-c) temperature, (d-f) dewpoint, and (g-
 553 i) zonal wind component at forecast lead times of 36 h (left column), 60 h (center column),
 554 and 84 h (right column) for the 5-member CAPS PI and PSI sub-ensembles, as well as for the
 555 0000 UTC HREF and five members of GEFS. HREF data are available only up to 500 hPa.

556

557 4. Summary and discussion

558 During the 2022 NOAA HWT SFE, CAPS ran an ensemble of 21 FV3-LAM forecasts in
 559 real-time using a variety of configurations with different combinations of microphysical

560 schemes, LSMs, and PBL parameterizations. These forecasts used 3-km grid spacing over the
561 contiguous United States and were subdivided into three sub-ensembles: a “P” sub-ensemble
562 with physics diversity only, a “PI” sub-ensemble that included perturbations to initial and
563 lateral boundary conditions in addition to physics diversity, and a “PSI” sub-ensemble that
564 included stochastic physics perturbations in addition to IC/LBC and physics diversity.
565 Precipitation, temperature, dewpoint, and wind forecasts from this ensemble and its sub-
566 ensembles are evaluated using the operational HREF and GEFS as references for comparison.

567 For 6-h accumulated precipitation, the CAPS FV3-LAM forecasts exhibited 90th percentile
568 values ranging from around 60% to 120% of the observed 90th percentile in Stage-IV rainfall,
569 with more relative impact from model physics configuration than from the inclusion of initial
570 and lateral boundary condition perturbations or stochastic perturbations. Members using the
571 NOAA-MP LSM exhibited the lowest 90th percentile rainfall values, while members using the
572 NSSL microphysics scheme exhibited among the highest 90th percentile values. Similar
573 patterns in the relative precipitation intensity of NOAA-MP and NSSL members were also
574 noted in terms of domain-averaged total rainfall. At the 99th percentile threshold, almost all
575 members had 99th percentile values slightly above that of Stage-IV observations, with relatively
576 little variation among members. In general, many members exhibited underprediction at low
577 rainfall rates and overprediction at relatively high rainfall rates. Difference in ETS for
578 individual forecast members, as well as among members in the PI and PSI sub-ensembles, were
579 generally small at both percentile thresholds. Sub-ensemble consensus, measured using the
580 PM and LPM, showed good agreement with observed Stage-IV rainfall both in terms of
581 percentile thresholds and domain-averaged total rainfall accumulation.

582 Overall, the greatest positive impact of including stochastic perturbations (which were
583 present in the CAPS PSI sub-ensemble but not in the PI sub-ensemble) was the improved
584 spread-error ratio noted in forecasts of 2-m temperature and dewpoint and 10-m zonal wind.
585 Some minor improvement in spread-error ratio is also noted for temperature, dewpoint, and
586 zonal wind verified against soundings, extending up to around 500 hPa on forecast days 1 and
587 2, with dewpoint exhibiting the greatest improvement. When the PI and PSI sub-ensembles
588 were compared using other metrics, including AUC for 6-h precipitation forecasts,
589 performance diagrams, rank histograms, and FSS for 24-h accumulated rainfall, and domain-
590 averaged total rainfall, there was very little noticeable difference in the performance of the PI
591 and PSI sub-ensembles. These sub-ensemble evaluations do, however, highlight other

592 important aspects of CAPS forecast performance, such as the improved performance of the
593 LPM mean over the PM mean in terms of FSS at higher precipitation thresholds, suggesting
594 that the patchwise LPM algorithm being used is successful in its goal of retaining localized
595 structures in the precipitation field. We note that some prior studies have found greater benefit
596 when applying stochastic perturbations. For example, Jankov et al. (2019) found that including
597 stochastic perturbations in a multi-physics CAM ensemble did improve the RMSE of 10-m
598 wind forecasts, and Zhang (2019) reported a positive impact of stochastic perturbations for
599 CAM forecasts of precipitation. In light of these prior studies, it is possible that, with additional
600 tuning of the stochastic perturbation methods used, greater benefit might be obtainable from
601 inclusion of stochastic perturbations.

602 Compared to the HREF and GEFS forecasts used as baseline comparisons for evaluation,
603 the CAPS FV3-LAM forecasts typically performed slightly worse than (or at best no better
604 than) HREF over the 0-48 hour forecast timeframe when HREF forecasts were available,
605 though the CAPS FV3-LAM forecasts did consistently outperform GEFS throughout the 84
606 hour forecast period. The better performance of the CAPS FV3-LAM forecasts compared to
607 GEFS is unsurprising given the much higher spatial resolution of the FV3-LAM forecasts.

608 As was noted in the introduction, one of the motivations for this study is to inform the
609 design of future convection-allowing NWP ensembles, in particular the planned RRFS
610 ensemble using FV3-LAM. Overall, the configurations of FV3-LAM examined in this study
611 appear to be generally suitable for predicting spring season convective rainfall, outperforming
612 GEFS across a variety of forecast quality metrics, and approaching the performance of HREF
613 in some metrics. HREF has been tuned and optimized over many years of operational use, and
614 is widely used for prediction of severe weather and high-impact precipitation events, so the
615 ability of an experimental FV3-based ensemble with limited tuning to approach the
616 performance of HREF is encouraging. We recommend that future studies continue
617 investigation into the bias behavior of the NOAA-MP LSM; reducing the biases inherent in
618 this scheme will help improve the performance of future operational ensembles. The use of
619 stochastic perturbations enhanced ensemble spread, but further optimization of their
620 configurations is needed for them to contribute more significantly towards improving the
621 ensemble spread without increasing biases of the ensembles. Future investigations of
622 stochastic physics perturbations could also examine the impacts of varying stochastic
623 perturbation settings, or applying stochastic perturbation methods (such as SKEB, SPPT, and

624 SHUM) individually or in different combinations. Ideally, if stochastic perturbations alone are
625 able to outperform a multi-physics ensemble, efforts can be focused on developing and
626 maintaining a single well-performing physics suite. Finally, we note that while this study
627 generally focused on features aggregated over 6 or 24 hours, future studies will examine FV3-
628 LAM ensemble forecast performance in terms of more localized, shorter-duration features and
629 severe weather hazards.

630

631 *Acknowledgments*

632 This work was primarily supported by the NOAA UFS R2O program under grant
633 NA21OAR4320204, with additional support from NOAA Testbed grants NA19OAR4590141
634 and NA22OAR4590522. The CAPS FV3-LAM forecasts were run in real-time on Frontera at
635 the Texas Advanced Computing Center (TACC) at the University of Texas-Austin, a National
636 Science Foundation supported high performance computing center. The authors would like to
637 thank the NOAA staff involved in organizing and leading the HWT SFE, including Adam Clark
638 and Israel Jirak. The authors would also like to thank Philip Pegion for assistance with the
639 configurations of stochastic perturbations. The authors also thank Dr. Curtis Alexander for
640 facilitating the use of GSL-produced experimental RRFS analysis initial conditions.

641

642 *Data Availability Statement.*

643 The FV3-LAM and HREF model output evaluated in this study, has been archived by the
644 Center for Analysis and Prediction of Storms (CAPS) and can be downloaded via the Open
645 Science Framework at <https://osf.io/u38yp/>. GFS forecasts are available from UCAR's
646 Research Data Archive (RDA) repository (<https://rda.ucar.edu/datasets/ds084.1/>), while GEFS
647 forecasts can be downloaded from NOAA's registry of open data on AWS at
648 <https://registry.opendata.aws/noaa-gefs/>. The Unrestricted Mesoscale Analysis (URMA) data
649 provided by NOAA EMC which used for forecast evaluation can be accessed via a NOAA
650 open data registry at <https://registry.opendata.aws/noaa-rtma/>. Stage-IV precipitation analyses
651 can be downloaded from UCAR's Earth Observing Laboratory archive
652 (<https://data.eol.ucar.edu/dataset/21.093>).

653

654

REFERENCES

655 Banos, I. H., W. D. Mayfield, G. Q. Ge, L. F. Sapucci, J. R. Carley, and L. Nance, 2022:
656 Assessment of the data assimilation framework for the Rapid Refresh Forecast System
657 v0.1 and impacts on forecasts of a convective storm case study. *Geoscientific Model*
658 *Development*, **15**, 6891-6917, <https://doi.org/10.5194/gmd-15-6891-2022>.

659 Bengtsson, L., M. Steinheimer, P. Bechtold, and J.F. Geleyn, 2013: A stochastic
660 parametrization for deep convection using cellular automata. *Quart. J. Roy. Meteor. Soc.*,
661 **139(675)**, 1533-1543, doi:10.1002/qj.2108

662 Berner, J., G. Shutts, M. Leutbecher, and T. Palmer, 2009: A spectral stochastic kinetic
663 energy backscatter scheme and its impact on flow- dependent predictability in the
664 ECMWF ensemble prediction system. *J. Atmos. Sci.*, **66**, 603–626,
665 doi:10.1175/2008JAS2677.1.

666 Berner, J., S. Y. Ha, J. P. Hacker, A. Fournier, and C. Snyder, 2011: Model Uncertainty in a
667 Mesoscale Ensemble Prediction System: Stochastic versus Multiphysics Representations.
668 *Mon. Wea. Rev.*, 139, 1972-1995, 10.1175/2010mwr3595.1.

669 Black, T. L., Abeles, J. A., Blake, B. T., Jovic, D., Rogers, E., Zhang, X., et al., 2021: A
670 limited area modeling capability for the Finite-Volume Cubed-Sphere (FV3) dynamical
671 core and comparison with a global two-way nest. *Journal of Advances in Modeling Earth*
672 *Systems*, **13**, e2021MS002483, <https://doi.org/10.1029/2021MS002483>

673 Brown, B., and Coauthors, 2021: The Model Evaluation Tools (MET): More than a Decade
674 of Community-Supported Forecast Verification. *Bull. Amer. Meteor. Soc.*, **102**, E782–
675 E807, <https://doi.org/10.1175/BAMS-D-19-0093.1>.

676 Buizza, R., M. Miller, and T. N. Palmer, 1999: Stochastic representation of model
677 uncertainties in the ECMWF ensemble prediction system. *Quart. J. Roy. Meteor. Soc.*,
678 **125**, 2887–2908, doi:10.1002/qj.49712556006.

679 Cafaro, C., and Coauthors, 2021: Do Convection-Permitting Ensembles Lead to More
680 Skillful Short-Range Probabilistic Rainfall Forecasts over Tropical East Africa?. *Wea.*
681 *Forecasting*, **36**, 697–716, <https://doi.org/10.1175/WAF-D-20-0172.1>.

682 Carley, J. R., M. E. Pyle, C. R. Alexander, and S. Weygandt, 2023: On the Development of
683 NOAA’s Rapid Refresh Forecast System. *Research Activities in Earth System Modelling.*
684 *Working Group on Numerical Experimentation*. Report No. 53. WCRP Report

685 No.6/2023, 5.07-05.08.
686 https://wgne.net/bluebook/uploads/2023/docs/05_Carley_Jacob_RapidRefreshModel.pdf
687 Carley, J. R., M. E. Pyle, C. R. Alexander, and S. Weygandt, 2024: Development and Status
688 of NOAA's Rapid Refresh Forecast System. *14th Conference on Transition of Research*
689 *to Operations*. Baltimore, MD, 23 Jan. 2024, Amer. Meteor. Soc., 3B.1,
690 <https://ams.confex.com/ams/104ANNUAL/meetingapp.cgi/Paper/430385>
691 Clark, A. J., W. A. Gallus, and T.-C. Chen, 2008: Contributions of Mixed Physics versus
692 Perturbed Initial/Lateral Boundary Conditions to Ensemble-Based Precipitation Forecast
693 Skill. *Mon. Wea. Rev.*, 136, 2140-2156, 10.1175/2007mwr2029.1.
694 Clark, A. J., and Coauthors, 2011: Probabilistic Precipitation Forecast Skill as a Function of
695 Ensemble Size and Spatial Scale in a Convection-Allowing Ensemble. *Mon. Wea.*
696 *Rev.*, **139**, 1410–1418, <https://doi.org/10.1175/2010MWR3624.1>.
697 Clark, A. J., and Coauthors, 2020: A real-time, simulated forecasting experiment for
698 advancing the prediction of hazardous convective weather. *Bull. Am. Meteorol. Soc.*, **101**,
699 E2022–E2024, <https://doi.org/10.1175/BAMS-D-19-0298.1>.
700 De Pondeva, M. S. F. V., and Coauthors, 2011: The real-time mesoscale analysis at NOAA's
701 National Centers for Environmental Prediction: Current status and development. *Wea.*
702 *Forecasting*, **26**, 593–612, <https://doi.org/10.1175/WAF-D-10-05037.1>.
703 DTC 2023: Developmental Testbed Center: Model Evaluation Tools (MET). Accessed 7
704 Nov. 2023. <https://dtcenter.org/community-code/model-evaluation-tools-met>.
705 Ebert, E. E., 2001: Ability of a poor man's ensemble to predict the probability and
706 distribution of precipitation. *Mon. Wea. Rev.*, **129**, 2461-2480.
707 Ek, M. B., K. E. Mitchell, Y. Lin, E. Rogers, P. Grunmann, V. Koren, G. Gayno, and J. D.
708 Tarpley, 2003: Implementation of Noah land surface model advances in the National
709 Centers for Environmental Prediction operational mesoscale Eta model. *J. Geophys. Res.*
710 *Atmos.*, **108**, 1–16, <https://doi.org/10.1029/2002jd003296>.
711 Hamill, T. M., 2001: Interpretation of rank histograms for verifying ensemble forecasts. *Mon.*
712 *Wea. Rev.*, 129, 550, [https://doi.org/10.1175/1520-](https://doi.org/10.1175/1520-0493(2001)129<0550:IORHFV>2.0.CO;2)
713 [0493\(2001\)129<0550:IORHFV>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<0550:IORHFV>2.0.CO;2).

- 714 Hu, X.-M., Park, J., Supinie, T., Snook, N. A., Xue, M., Brewster, K. A., Brotzge, J., &
715 Carley, J. R. (2023). Diagnosing Near-Surface Model Errors with Candidate Physics
716 Parameterization Schemes for the Multiphysics Rapid Refresh Forecast System (RRFS)
717 Ensemble during Winter over the Northeastern United States and Southern Great Plains,
718 *Monthly Weather Review*, **151**(1), 39-61, <https://doi.org/10.1175/MWR-D-22-0085.1>
- 719 Han, J., and C. S. Bretherton, 2019: TKE-based moist eddy-diffusivity mass-flux (EDMF)
720 parameterization for vertical turbulent mixing. *Weather and Forecasting*, **34**, 869–886,
721 <https://doi.org/10.1175/WAF-D-18-0146.1>.
- 722 Jankov, I., and et al., 2017: A performance comparison between multiphysics and stochastic
723 approaches within a North American RAP ensemble. *Monthly Weather Review*, **145**,
724 1161–1179, doi:10.1175/MWR-D-16-0160.1
- 725 Jankov, I., J. Beck, J. Wolff, M. Harrold, J. B. Olson, T. Smirnova, C. Alexander, and J.
726 Berner, 2019: Stochastically perturbed parameterizations in an HRRR-based ensemble.
727 *Monthly Weather Review*, **147**, 153–173, doi:10.1175/MWR-D-18-0092.1
- 728 Jirak, I. L., S. J. Weiss, and C. J. Melick, 2012: The SPC storm-scale ensemble of
729 opportunity: Overview and results from the 2012 Hazardous Weather Testbed Spring
730 Forecasting Experiment. 26th Conf. on Severe Local Storms, Nashville, TN, Amer.
731 Meteor. Soc., P9.137,
732 <https://ams.confex.com/ams/26SLS/webprogram/Paper211729.html>.
- 733 Johnson, M., M. Xue, and Y. Jung, 2023: Biases and skill of four two-moment bulk
734 microphysics schemes in convection-allowing forecasts for the 2018 Hazardous Weather
735 Testbed Spring Forecasting Experiment period. *Wea. Forecasting.*, **38**, 1621–
736 1642, <https://doi.org/10.1175/WAF-D-22-0171.1>.
- 737 Kalina, E. A., and et al. 2021: A Progress Report on the Development of the High-Resolution
738 Rapid Refresh Ensemble. *Weather and Forecasting*, **36**(3), 791-804. doi:10.1175/WAF-
739 D-20-0098.1
- 740 Lin, S. J., 2004: A “vertically Lagrangian” finite-volume dynamical core for global models.
741 *Mon. Wea. Rev.*, 132, 2293–2307, [https://doi.org/10.1175/1520-0493\(2004\)132<2293:AVLFDC>2.0.CO;2](https://doi.org/10.1175/1520-0493(2004)132<2293:AVLFDC>2.0.CO;2).
- 743 Liu, Shun, M. Hu, T. Lei, X. Zhang, S. Yokota, B. Blake, E. Rogers, C. R. Martin, T. T.
744 Ladwig, D. Dowell, C. Zhou, J. R. Carley, C. R. Alexander, and D. T. Kleist, 2023: Data

745 Assimilation Testing and Development Toward Implementation of the Rapid Refresh
746 Forecast System. *27th Conference on Integrated Observing and Assimilation Systems for*
747 *the Atmosphere, Oceans, and Land Surface (IOAS-AOLS)*, Denver, CO, Amer. Meteor.
748 Soc., 86, <https://ams.confex.com/ams/103ANNUAL/meetingapp.cgi/Paper/414100>.

749 Loken, E. D., J. C. Adam, X. Ming, and K. Fanyou, 2019: Spread and skill in mixed- and
750 single-physics convection-allowing ensembles. *Wea. Forecasting*, **34**, 305-330,
751 <https://doi-org.ezproxy.lib.ou.edu/10.1175/WAF-D-18-0078.1>.

752 Mansell, E. R., and C. L. Ziegler, 2013: Aerosol effects on simulated storm electrification
753 and precipitation in a two-moment bulk microphysics model. *J. Atmos. Sci.*, **70**, 2032–
754 2050, <https://doi.org/10.1175/JAS-D-12-0264.1>.

755 Mansell, E. R., C. L. Ziegler, and E. C. Bruning, 2010: Simulated electrification of a small
756 thunderstorm with two-moment bulk microphysics. *J. Atmos. Sci.*, **67**, 171–194,
757 <https://doi.org/10.1175/2009JAS2965.1>.

758 Mason, I. B., 1982: A model for assessment of weather forecasts. *Aust. Meteor. Mag.*, **30**,
759 291-303, [https://doi.org/10.1175/1520-0493\(1979\)107<0207:ORPFTY>2.0.CO;2](https://doi.org/10.1175/1520-0493(1979)107<0207:ORPFTY>2.0.CO;2).

760 Mlawer, E. J., S. J. Taubman, P. D. Brown, M. J. Iacono, and S. A. Clough, 1997: Radiative
761 transfer for inhomogeneous atmospheres: RRTM, a validated correlated-k model for the
762 longwave. *J. Geophys. Res. Atmos.*, **102**, 16663–16682,
763 <https://doi.org/10.1029/97JD00237>.

764 Morris, M. T., J. R. Carley, E. Colon, A. Gibbs, M. S. F. V. De Pondeva, and S. Levine,
765 2020: A Quality Assessment of the Real-Time Mesoscale Analysis (RTMA) for Aviation.
766 *Weather and Forecasting*, **35**, 977-996, <https://doi.org/10.1175/Waf-D-19-0201.1>.

767 Nakanishi, M., and H. Niino, 2009: Development of an Improved Turbulence Closure Model
768 for the Atmospheric Boundary Layer. *J. Meteorol. Soc. Japan*, **87**, 895–912,
769 <https://doi.org/10.2151/jmsj.87.895>.

770 Nelson, B. R., O. P. Prat, D. J. Seo, and E. Habib, 2016: Assessment and implications of
771 NCEP stage IV quantitative precipitation estimates for product intercomparisons.
772 *Weather and Forecasting*, **31**, 371–394, <https://doi.org/10.1175/WAF-D-14-00112.1>.

773 Niu, G. Y., and Coauthors, 2011: The community Noah land surface model with
774 multiparameterization options (Noah-MP): 1. Model description and evaluation with

775 local-scale measurements. *J. Geophys. Res. Atmos.*, **116**, 1–19,
776 <https://doi.org/10.1029/2010JD015139>.

777 NOAA, 2023: WPC Hydrometeorology Testbed. Accessed 6 Oct. 2023,
778 <https://www.wpc.ncep.noaa.gov/hmt/>

779 Palmer, T. N., R. Buizza, F. Doblas-Reyes, T. Jung, M. Leutbecher, G. J. Shutts, M.
780 Steinheimer, and A. Weisheimer, 2009: Stochastic parametrization and model uncertainty.
781 *ECMWF Tech. Memo.*, **598**, 42 pp, doi:10.21957/ps8gbwbdv

782 Ponca, M., S. Levine, J. Carley, Y. Lin, Y. Zhu, and J. Purser, 2015: Ongoing
783 improvements to the NCEP Real Time Mesoscale Analysis (RTMA) and UnRestricted
784 Mesoscale Analysis (URMA) and NCEP/EMC. NOAA, 2 pp., [https://www.wcrp-](https://www.wcrp-climate.org/WGNE/BlueBook/2015/individual-articles/01_Ponca_Manuel_et_al_RTMA.pdf)
785 [climate.org/WGNE/BlueBook/2015/individual-](https://www.wcrp-climate.org/WGNE/BlueBook/2015/individual-articles/01_Ponca_Manuel_et_al_RTMA.pdf)
786 [articles/01_Ponca_Manuel_et_al_RTMA.pdf](https://www.wcrp-climate.org/WGNE/BlueBook/2015/individual-articles/01_Ponca_Manuel_et_al_RTMA.pdf).

787 Roberts, B., B. T. Gallo, I. L. Jirak, A. J. Clark, D. C. Dowell, X. Wang, and Y. Wang, 2020:
788 What Does a Convection-Allowing Ensemble of Opportunity Buy Us in Forecasting
789 Thunderstorms? *Weather and Forecasting*, **35**, 2293–2316, [https://doi.org/10.1175/waf-d-](https://doi.org/10.1175/waf-d-20-0069.1)
790 [20-0069.1](https://doi.org/10.1175/waf-d-20-0069.1).

791 Roberts, B., I. L. Jirak, A. J. Clark, S. J. Weiss, and J. S. Kain, 2019: Post-Processing and
792 Visualization Techniques for Convection-Allowing Ensembles. *Bull. Amer. Meteor. Soc.*,
793 **100**, 1245–1258, <https://doi.org/10.1175/BAMS-D-18-0041.1>

794 Roberts, B., B. T. Gallo, I.
795 L. Jirak, A. J. Clark, D. C. Dowell, X. Wang, and Y. Wang, 2020: What does a
796 convection-allowing ensemble of opportunity buy us in forecasting thunderstorms?
Weather and Forecasting, **35**, 2293–2316, <https://doi.org/10.1175/WAF-D-20-0069.1>.

797 Roberts, N.M. and H.W. Lean, 2008: Scale-selective verification of rainfall accumulations
798 from high-resolution forecasts of convective events. *Monthly Weather Review*, **136**, 78-
799 97.

800 Roebber, P. J., 2009: Visualizing Multiple Measures of Forecast Quality. *Weather and*
801 *Forecasting*, **24**, 601–608, <https://doi.org/10.1175/2008waf2222159.1>.

802 Shin, H. H., and S.-Y. Hong, 2015: Representation of the subgrid-scale turbulent transport in
803 convective boundary layers at gray-zone resolutions. *Monthly Weather Review*, **143**, 250–
804 271, <https://doi.org/10.1175/MWR-D-14-00116.1>.

805 Smirnova, T. G., J. M. Brown, S. G. Benjamin, and J. S. Kenyon, 2016: Modifications to the
806 Rapid Update Cycle land surface model (RUC LSM) available in the weather research
807 and forecasting (WRF) model. *Monthly Weather Review*, **144**, 1851–1865,
808 <https://doi.org/10.1175/MWR-D-15-0198.1>.

809 Snook, N., F. Kong, A. Clark, B. Roberts, K. A. Brewster, and M. Xue, 2020: Comparison
810 and Verification of Pointwise and Patchwise Localized Probability Matched Mean
811 Algorithms for Convective Ensemble Forecasting. *Geophysical Research Letters*, **47**,
812 <https://doi.org/2020GL087839>.

813 Snook, N., F. Kong, K. Brewster, M. Xue, B. Albright, S. Perfater, K. Thomas, and T.
814 Supinie, 2019: Evaluation of Convection-Permitting Precipitation Forecast Products using
815 WRF, NMMB, and FV3 for the 2016-2017 NOAA Hydrometeorology Testbed Flash
816 Flood and Intense Rainfall Experiments. *Weather and Forecasting*, **34**, 781–804,
817 <https://doi.org/10.1175/WAF-D-18-0155.1>.

818 Stensrud, D. J., J. W. Bao, and T. T. Warner, 2000: Using initial condition and model physics
819 perturbations in short-range ensemble simulations of mesoscale convective systems. *Mon.*
820 *Wea. Rev.*, 128, 2077-2107, Doi 10.1175/1520-0493(2000)128<2077:Uicamp>2.0.Co;2.

821 Supinie, T. A., Park, J., Snook, N., Hu, X.-M., Brewster, K. A., Xue, M., & Carley, J. R.,
822 2022: Cool-Season Evaluation of FV3-LAM-Based CONUS-Scale Forecasts with
823 Physics Configurations of Experimental RRFs Ensembles. *Monthly Weather*
824 *Review*, **150(9)**, 2379-2398, <https://doi.org/10.1175/MWR-D-21-0331.1>

825 Schwartz, C. S., J. S. Kain, S. J. Weiss, M. Xue, D. R. Bright, F. Kong, K. W. Thomas, J. J.
826 Levit, M. C. Coniglio, and M. S. Wandishin, 2010: Toward Improved Convection-
827 Allowing Ensembles: Model Physics Sensitivities and Optimizing Probabilistic Guidance
828 with Small Ensemble Membership. *Wea. Forecasting*, 25, 263-280,
829 [10.1175/2009waf2222267.1](https://doi.org/10.1175/2009waf2222267.1).

830 Tompkins, A. M., and J. Berner, 2008: A stochastic convective approach to account for
831 model uncertainty due to unresolved humidity variability. *J. Geophys. Res.*, **113**, D18101,
832 [doi:10.1029/2007JD009284](https://doi.org/10.1029/2007JD009284)

833 Thompson, G., and T. Eidhammer, 2014: A Study of Aerosol Impacts on Clouds and
834 Precipitation Development in a Large Winter Cyclone. *J. Atmos. Sci.*, **71**, 3636–
835 3658, <https://doi.org/10.1175/JAS-D-13-0305.1>.

836 Thompson, G., R. M. Rasmussen, and K. Manning, 2004: Explicit Forecasts of Winter
837 Precipitation Using an Improved Bulk Microphysics Scheme. Part I: Description and
838 Sensitivity Analysis. *Monthly Weather Review*, **132**, 519–542,
839 [https://doi.org/10.1175/1520-0493\(2004\)132<0519:EFOWPU>2.0.CO;2](https://doi.org/10.1175/1520-0493(2004)132<0519:EFOWPU>2.0.CO;2).

840 Thompson, G., P. R. Field, R. M. Rasmussen, and W. D. Hall, 2008: Explicit Forecasts of
841 Winter Precipitation Using an Improved Bulk Microphysics Scheme. Part II:
842 Implementation of a New Snow Parameterization. *Monthly Weather Review*, **136**, 5095–
843 5115, <https://doi.org/10.1175/2008MWR2387.1>.

844 Zhang, C., M. Xue, T. A. Supinie, F. Kong, N. Snook, K. W. Thomas, K. Brewster, Y. Jung,
845 L. M. Harris, and S.-J. Lin, 2019: How Well Does the FV3 Model Predict Precipitation at
846 a Convection-Allowing Resolution? Results from CAPS Forecasts for the 2018 NOAA
847 Hazardous Weather Testbed with Different Physics Combinations. *Geophysical Research*
848 *Letters*, **46**, 3523-3531.

849 Zhang, X., 2019: Multiscale Characteristics of Different-Source Perturbations and Their
850 Interactions for Convection-Permitting Ensemble Forecasting during SCMREX. *Mon.*
851 *Wea. Rev.*, **147**, 291–310, <https://doi.org/10.1175/MWR-D-18-0218.1>.

852 Zhou, L., S. Lin, J. Chen, L. M. Harris, X. Chen, and S. L. Rees, 2019: Toward Convective-
853 Scale Prediction within the Next Generation Global Prediction System. *Bull. Amer.*
854 *Meteor. Soc.*, **100**, 1225–1243, <https://doi.org/10.1175/BAMS-D-17-0246.1>.

855 Zhou, X., Y. Zhu, D. Hou, B. Fu, W. Li, H. Guan, E. Sinsky, W. Kolczynski, X. Xue, Y. Luo,
856 J. Peng, B. Yang, V. Tallapragada, and P. Pegion, 2022: The Development of the NCEP
857 Global Ensemble Forecast System Version 12. *Wea. Forecasting*, **37**, 1069-1084,
858 [10.1175/WAF-D-21-0112.1](https://doi.org/10.1175/WAF-D-21-0112.1).

859 Zhu, K., M. Xue, B. Zhou, K. Zhao, Z. Sun, P. Fu, Y. Zheng, X. Zhang, and Q. Meng, 2018:
860 Evaluation of real-time precipitation forecasts during 2013-2014 summer seasons over
861 China at a convection-permitting resolution: spatial distribution, propagation, diurnal
862 cycles and skill scores. *J. Geophys. Res.*, **123**, 1037-1064.

863