



AMERICAN METEOROLOGICAL SOCIETY

Weather and Forecasting

EARLY ONLINE RELEASE

This is a preliminary PDF of the author-produced manuscript that has been peer-reviewed and accepted for publication. Since it is being posted so soon after acceptance, it has not yet been copyedited, formatted, or processed by AMS Publications. This preliminary version of the manuscript may be downloaded, distributed, and cited, but please be aware that there will be visual differences and possibly some content differences between this version and the final published version.

The DOI for this manuscript is doi: [10.1175/WAF-D-13-00066.1](https://doi.org/10.1175/WAF-D-13-00066.1)

The final published version of this manuscript will replace the preliminary version at the above DOI once it is available.

If you would like to cite this EOR in a separate work, please use the following full citation:

Novak, D., C. Bailey, K. Brill, P. Burke, W. Hogsett, R. Rausch, and M. Schichtel, 2013: Precipitation and Temperature Forecast Performance at the Weather Prediction Center. *Wea. Forecasting*. doi:[10.1175/WAF-D-13-00066.1](https://doi.org/10.1175/WAF-D-13-00066.1), in press.



1
2
3 **Precipitation and Temperature Forecast Performance at the Weather Prediction Center**
4

5 David R. Novak, Christopher Bailey, Keith Brill, Patrick Burke, Wallace Hogsett, Robert
6 Rausch, and Michael Schichtel
7

8
9
10 @NOAA/NWS/NCEP, Weather Prediction Center,
11 College Park, Maryland
12
13
14
15
16
17
18

19 Submitted to *Weather and Forecasting*

20
21 June 10, 2013
22

23 Revised September 21, 2013
24 and
25 October 9, 2013
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

44 Corresponding author address: David R. Novak, NOAA/NWS Weather Prediction Center, 5830
45 University Research Court, Rm 4633, College Park, MD 20740.
46 E-mail: David.Novak@noaa.gov

47 ABSTRACT

48
49 The role of the human forecaster in improving upon the accuracy of numerical weather
50 prediction is explored using multi-year verification of human-generated short-range precipitation
51 forecasts and medium-range maximum temperature forecasts from the Weather Prediction
52 Center (WPC). Results show that human-generated forecasts improve over raw deterministic
53 model guidance. Over the past two decades, WPC human forecasters achieved a 20–40%
54 improvement over the NAM and GFS models for the 1-in (25.4-mm) 24 h⁻¹ threshold for the day
55 1 precipitation forecast, with a smaller, but statistically significant, 5–15% improvement over the
56 deterministic ECMWF model. Medium-range maximum temperature forecasts also exhibit
57 statistically significant improvement over GFS Model Output Statistics (MOS), and the
58 improvement has been increasing over the past five years. The quality added by humans for
59 forecasts of high-impact events varies by element and forecast projection, with generally large
60 improvements when the forecaster makes changes $\geq 8^{\circ}\text{F}$ (4.4°C) to MOS temperatures. Human
61 improvement over guidance for extreme rainfall events [3-in (76.2-mm) 24 h⁻¹] is largest in the
62 short-range forecast.

63 However, human-generated forecasts failed to outperform the most skillful downscaled,
64 bias-corrected ensemble guidance for precipitation and maximum temperature available near the
65 same time as the human-modified forecasts. Thus, as additional downscaled and bias-corrected
66 sensible weather element guidance becomes operationally available, and with the support of
67 near-real time verification, forecaster training, and tools to guide forecaster interventions, a key
68 test is whether forecasters can learn to make statistically significant improvements over the most
69 skillful of this guidance. Such a test can inform to what degree, and just how quickly the role of
70 the forecaster changes.

71 **1. Introduction**

72 As the skill of numerical weather prediction (NWP) and associated post-processed
73 guidance continues to improve, recent debate asks to what degree can human forecasters add
74 quality¹ to NWP (e.g., Mass 2003; Bosart 2003; Roebber et al. 2004; Reynolds 2003; Doswell
75 2004; Stuart et al. 2006; Homar et al. 2006; Novak et al. 2008; Ruth et al. 2009). The National
76 Centers for Environmental Prediction (NCEP) Weather Prediction Center (WPC²) has a broad
77 mission to serve as a center of excellence in quantitative precipitation forecasting (QPF),
78 medium range forecasting, winter weather forecasting, surface analysis, and the interpretation of
79 operational NWP. Historically, forecasters at the WPC have had access to a large portion of the
80 available model guidance suite, recently including multi-model ensemble information from
81 international partners. The WPC's unique national forecast mission coupled with its access to
82 state-of-the-art model guidance provides a rare opportunity to assess the quality added by
83 humans to ever-improving NWP guidance.

84 This work will examine multi-year historical and contemporary verification for short-
85 range deterministic precipitation forecasts and medium-range maximum temperature forecasts
86 generated at the WPC. Although humans can add substantial value to NWP through retaining
87 forecast continuity (run-to-run consistency), assuring element consistency (e.g., wind shifts with
88 fronts), and helping users make informed decisions (e.g., Roebber et al. 2010a), this work
89 focuses on the human role in improving forecast accuracy. In this respect, the current work
90 examines only one component of the forecaster's role, and is limited to analysis of just two
91 weather elements.

¹ Although "value-added" is often used colloquially, this work abides by the terms for forecast "goodness" defined in Table 1 of Murphy (1993), where "value" refers to the benefit realized by decision makers through the use of the forecasts and "quality" refers to the correspondence between forecasts and the matching observations.

² The Center's name was changed from the Hydrometeorological Prediction Center to the Weather Prediction Center on March 5, 2013.

92 The current work builds upon and extends previous analyses of WPC skill by Olson et al.
93 (1995), Reynolds (2003), and Sukovich et al. (2013), and points to future verification approaches
94 in the continuing history of NWP and the human forecaster. Section 2 presents analysis of QPF
95 while section 3 explores human improvement to medium range maximum temperature forecasts.
96 A discussion of the limitations of the work and implications of the verification for the future role
97 of the forecaster is presented in section 4.

98

99 **2. QPF**

100 *a. Production and verification method*

101 The WPC forecasters create deterministic QPFs at 6-h intervals through the day 3
102 forecast projection, and 48-h QPFs for days 4–5 and 6–7. The focus here is on the day 1–3
103 forecasts. The WPC deterministic QPF during the study period was defined as the most-likely,
104 areal-averaged amount mapped to a 32-km horizontal resolution grid. An example 24-h
105 accumulated QPF is shown in Fig. 1a. The forecast process for QPF involves forecaster
106 assessment of observations of moisture, lift, and instability and comparisons among deterministic
107 and ensemble forecasts of these parameters. Objectively post-processed model-based QPFs are
108 also available to WPC forecasters. Emphasis shifts from nowcasting based on observations in
109 the first 6–12 h of the forecast, to an increasing use of NWP as lead time increases. For example,
110 subjective blends of model guidance are used almost exclusively beyond 36-h. Forecasters
111 manually draw precipitation isohyets, which operational software converts to a grid. In areas of
112 complex topography, forecasters use monthly Parameter-elevation Regressions on Independent
113 Slopes Model (PRISM; Daly et al. 1994; Daly et al. 2008) output as a background.

114 The 24-h accumulated QPF was verified using a human quality-controlled (QC'ed)
115 analysis valid at 12 UTC. The analyst can choose a first-guess field from either the multisensor
116 Stage IV quantitative precipitation estimate mosaic analyses (Lin and Mitchell 2005) or Climate
117 Prediction Center (CPC) daily precipitation analysis (Higgins et al. 1996). The analyst QC's the
118 analysis based on gauge data and a review of radar data, and can adjust isohyets if necessary.
119 The QC'ed precipitation analysis is mapped onto a 32-km grid, matching the forecast grid.
120 Retrospective tests show that the relative skill difference between the WPC and NWP datasets
121 shown in this paper are not sensitive to the precipitation analysis used (e.g., the WPC QC'ed
122 analysis or Stage IV).

123 Conventional 2x2 contingency tables of dichotomous outcomes (e.g., Brill 2009) for
124 precipitation exceeding several thresholds are created by comparing each QPF to the
125 corresponding verifying analysis. The 2x2 contingency tables are used to calculate threat score
126 and frequency bias for the day 1 and day 3 forecast periods. The WPC forecast period naming
127 convention is shown in Fig. 2. Focus is placed on the threat score for the day 1 QPF valid at 12
128 UTC at the 1-in (25.4-mm) 24 h⁻¹ threshold. This threat score is reported to Congress as part of
129 the Government Performance and Results Act of 1993 (GPRA). Historically, the goal of the
130 WPC QPF was to improve upon the model guidance available during the interval of forecast
131 preparation (Reynolds 2003). Therefore the performance of the WPC QPF is compared against
132 model forecasts that are somewhat older (i.e., time lagged) than the WPC issuance time. The
133 latency of the WPC forecasts for the most frequently used model guidance are shown in Table 1.

134 The historical verification analysis was constrained to data available during the last ~50
135 years, which were largely deterministic forecasts. Bias-corrected forecasts were also generally
136 not available for verification purposes during the historical timeframe. Bias correction can

137 dramatically improve raw QPF guidance (e.g., Yussouf and Sensrud 2006; Brown and Seo
138 2010), and ensemble approaches can quantify predictability and reduce error. Thus, the
139 contemporary verification compares the WPC QPF to one created by an ensemble algorithm with
140 bias correction, issued near the time as the human-modified forecast. This product, the pseudo
141 bias corrected ensemble QPF (ENSBC), is based on the premise that the larger the uncertainty,
142 the smoother the forecast should be, whereas the smaller the uncertainty the more detailed the
143 forecast should be. During the study period the ENSBC was composed of a high-resolution
144 ensemble part comprising output from the deterministic NCEP North American Mesoscale
145 model (NAM; Janjic 2003), Global Forecast System (GFS; Caplan et al. 1997), and European
146 Center for Medium-Range Weather Forecasts (ECMWF; Magnusson and Kallen 2013), and a
147 full ensemble part composed of the high-resolution ensemble plus the Canadian GEM (Belair et
148 al. 2009), UK Met Office model (UKMET), and all members of the NCEP Short Range
149 Ensemble Forecast (SREF; Du et al. 2006). The product is objectively downscaled from 32 km
150 to 10 km (5 km over the west) using PRISM. A detailed description of the ENSBC algorithm is
151 provided in Appendix A.

152 Additionally, the historical analysis is limited to verification metrics with a long record to
153 facilitate historical context [threat score, frequency bias (referred to as bias hereafter), and mean
154 absolute error (MAE)]. Metrics such as the threat score have inherent limitations, including a
155 double penalty for false alarms (Baldwin et al. 2002) and bias sensitivity (Brill 2009; Brill and
156 Mesinger 2009). To address this issue, a bias-removed threat score is calculated using the
157 procedure based on probability matching (Ebert 2001) described by Clark et al. (2009). The
158 procedure uses probability matching to reassign the distribution of a forecast field with that of the
159 observed field, so that the modified forecast field has the same spatial patterns as the original
160 forecast, but has values adjusted so the distribution of their amplitudes exactly matches those of the

161 analysis. The end result is the removal of all bias. Because the NAM precipitation skill lags so
162 severely relative to WPC and other international guidance, and to simplify interpretation, the
163 bias-removed threat score calculation was not conducted for the NAM. This reduced skill is
164 likely due to the use of 6-h old boundary conditions from the GFS, an earlier data cutoff, as well
165 as a less advanced data assimilation system (G. DiMego and E. Rogers, personal
166 communication).

167 Finally, it is important to quantify the statistical significance of comparisons. To
168 accomplish this task the forecast verification system (fvs) software was used (described in
169 Appendix B). Assessment of statistical significance in fvs is accomplished using random
170 resampling following the method of Hamill (1999).

171 Thus, contemporary verification addressing these four issues (ensemble approaches,
172 bias-corrected guidance, bias sensitivity, and statistical significance) was conducted during the
173 latest years available (2011–12).

174

175 *b. Results*

176 Verification of the WPC QPF over the last 50 years (Fig. 3) is a testament to the
177 advancement of precipitation forecasts. Threat scores of the 1-in (25.4-mm) 24 h^{-1} threshold for
178 day 1 forecasts doubled during the period, while day 2 and 3 forecasts also continued to improve
179 (Fig. 3). Improvement has accelerated after 1995. This improvement is directly tied to the quality
180 of the NWP guidance. In fact, during the 1993–2012 period, the correlation of yearly values of
181 the day 1 threat score for the 1-in (25.4-mm) 24 h^{-1} threshold between WPC and the NAM and
182 WPC and the GFS was 0.91 and 0.88, respectively.

183 Although NWP serves as skillful guidance, verification over the past two decades shows
184 WPC human forecasters achieved a 20–40% improvement over the deterministic NAM and GFS
185 for the threat score of the 1-in (25.4-mm) 24 h⁻¹ threshold for the day 1 forecast (Fig. 4a). This
186 improvement was occurring during a period of advances in NWP skill. For example, the GFS 1-
187 in (25.4-mm) 24 h⁻¹ day 1 threat score in 1993 was 0.14, whereas in 2012 it was 0.25. Based on
188 the long term rate of model improvement, it would take ~13 years until the GFS attains a day 1
189 threat score equivalent to the current WPC threat score. This rate is nearly identical to the 14
190 years reported by Reynolds (2003) for the 2001 verification year.

191 The ECMWF precipitation forecast information became available to WPC forecasters in
192 the mid 2000s, and the first full year of formal verification was established in 2008. Verification
193 of the 1-in (25.4-mm) 24 h⁻¹ day 1 forecast over the 2008–2012 period shows that the WPC
194 forecast exhibits smaller 5–15% improvements over the very skillful deterministic ECMWF
195 model (Fig. 4a). However, WPC improvement over the ECMWF model has nearly doubled over
196 the past 5 years.

197 A complete picture of precipitation verification must include bias information. In recent
198 years the NAM, GFS, and ECWTF guidance have exhibited a low bias at the 1-in (25.4-mm) 24
199 h⁻¹ threshold, while the WPC has sustained a more favorable bias near 1.0 (Fig. 4a).
200 Contemporary verification using the bias-removed threat score shows that WPC has maintained a
201 statistically significant advantage over the ECMWF and GFS during 2011 and 2012 (Fig. 4b).
202 However, the ensemble-based post-processed QPF from the ENSBC was very competitive. In
203 fact, the ENSBC and WPC forecasts were statistically similar for the 1-in (25.4-mm) 24 h⁻¹
204 threshold during 2012 (Fig. 4b).

205 Mass (2003) and McCarthy et al. (2007) have asserted that the human is most effective
206 for the near-term forecast. However, the WPC percent improvement over the GFS at the 1-in
207 (25.4-mm) 24 h⁻¹ threshold for the day 3 forecast is similar to the percent improvement for this
208 threshold for the day 1 forecast (c.f. Figs. 4b and d). All guidance, including WPC, has a slight
209 low bias at the day 3 forecast (Figs. 4c,d). For both the day 1 and day 3 forecasts, the
210 competitive skill of the ECWMF forecast is evident, for which the human adds small, but
211 statistically significant positive skill. However, once again, the WPC is statistically similar to the
212 ENSBC at the 1-in (25.4-mm) 24 h⁻¹ threshold for the day 3 period during 2012 (Fig. 4d). Thus,
213 at least for precipitation at this threshold, the quality added by the forecaster does not appear
214 dependent on forecast projection.

215 Mass (2003), Bosart (2003), Stuart et al. (2006), and McCarthy et al. (2007) have
216 suggested that the human forecaster may be most adept at improving over NWP guidance for
217 high-impact events. The threat score for the 3-in (76.2 mm) 24 h⁻¹ threshold is arbitrarily used
218 here as a proxy for a high-impact event. The skill of both model and human forecasts at the 3-in
219 (76.2 mm) 24 h⁻¹ threshold is rather poor when compared to the 1-in threshold, illustrating the
220 challenge of forecasting extreme rainfall events (Fritsch and Carbone 2004; Sukovich et al.
221 2013). However, the day 1 WPC threat score exhibits a large improvement over select NWP
222 (Fig. 5a), with a slight dry bias. Contemporary verification accounting for bias shows WPC
223 significantly improved over the GFS in 2012 and ECMWF in both 2011 and 2012 at this
224 threshold. However, once again, WPC was similar in skill to the ENSBC product (Fig. 5b).

225 Skill comparisons for the 3-in (76.2 mm) 24 h⁻¹ threshold at the day 3 lead time reveals
226 generally less forecaster improvement, with similar model and WPC threat scores (Fig. 5c). In
227 fact, the GFS was superior to the WPC forecast in 2001 and 2003, and the ECMWF was superior

228 to the WPC forecast in 2009. All guidance, except the GFS, is severely under-biased. The
229 authors speculate that the GFS had frequent grid-point storms (e.g., Giorgi 1991) during the
230 verification period, which may have improved its bias, but degraded its threat score.
231 Contemporary verification shows the WPC bias-removed threat score is not statistically
232 significantly different than the corresponding threat scores from any of the competitive guidance
233 (Fig. 5d).

234 All of the above results suggest humans can make statistically significant improvement
235 over competitive deterministic model guidance for precipitation. The magnitude of quality-added
236 by the forecaster is generally not dependent on forecast projection for the 1-in (25.4-mm) 24 h^{-1}
237 threshold; however, human improvement for extreme rainfall events does appear dependent on
238 forecast projection, favoring larger human improvements over deterministic model guidance in
239 the short-range forecast. However, a downscaled, bias-corrected ensemble forecast available near
240 the same time as the human-modified forecast exhibits similar skill – even for extreme
241 precipitation events.

242

243 **3. Maximum temperature**

244 *a. Production and verification method*

245 WPC forecasters produce a 3–7 day forecast suite including gridded predictions of
246 sensible weather elements to support the National Digital Forecast Database (NDFD; Glahn and
247 Ruth 2003) (Fig. 1b), graphical depictions of the surface fronts and pressures (Fig. 1c), and
248 associated discussion of forecast factors and confidence. Two forecasters work in tandem to
249 complete this task and coordinate with users after assessment of NWP. Since 2004, forecasters
250 have used a graphical interface to apply weights to individual models and ensemble systems to

251 derive a most-likely sensible weather solution. The result of the forecaster's chosen blend can be
252 manually edited.

253 Before model data is weighted by the forecaster, the data is bias-corrected and
254 downscaled to a 5-km horizontal resolution. Bias correction of gridded model data is
255 accomplished using the NCEP decaying averaging bias-correction method of Cui et al. (2012),
256 applied as:

$$257 \quad B_{new} = (1 - w)B_{past} + wB_{current} , \quad (1)$$

258 where, $B_{current}$ is the latest calculated forecast error given by the difference between the forecast
259 and verifying analysis, B_{past} is the past accumulated bias, and B_{new} is the updated accumulated
260 bias. The NCEP 5-km resolution Real Time Mesoscale Analysis (RTMA; De Pondeva et al.
261 2011) was used as the verifying analysis. The weight factor, w , controls how much influence to
262 give the most recent bias behavior of weather systems. A w equal to 2% was used operationally.
263 Once initialized, the bias estimate can be updated by considering just the current forecast error
264 ($B_{current}$) and the stored average bias (B_{past}). The new bias-corrected forecast is generated by
265 subtracting B_{new} from the current forecasts at each lead time and each grid point.

266 Downscaling of coarse model data to a 5-km resolution grid is accomplished using a
267 decaying averaged downscaling increment (Cui et al. 2013). The downscaling increments are
268 created at each 6-h time step by differencing the coarse 1° resolution GFS analysis (GDAS) and
269 5-km resolution RTMA according to:

$$270 \quad D_{new} = (1 - w)D_{past} + wD_{current} , \quad (2)$$

271 where, $D_{current}$ is the latest calculated downscaling increment given by the difference between
272 GDAS and RTMA, D_{past} is the past accumulated downscale increment, and D_{new} is the updated
273 downscale increment. The weight factor, w , controls how much influence to give the most recent
274 difference. A w equal to 10% was used operationally. The 6-hour grids are then downscaled
275 using the mean downscaling increment for each 6-hour period. For maximum and minimum
276 temperature, at each grid point, the downscaled 6-hour grids are compared to each other to find
277 the highest (lowest) values for maximum (minimum) temperature over the 12–06 UTC period
278 (00–18 UTC period) to get a final maximum (minimum) temperature forecast grid. The verifying
279 maximum (minimum) temperature is taken as the highest (lowest) hourly value from the RTMA
280 at each grid point.

281 The resulting maximum and minimum temperatures are extracted from the 5-km grid to
282 448 points for the forecaster to edit where necessary. An objective analysis is performed on the
283 incremental changes made by the forecaster at the 448 points to create a difference grid. The
284 forecaster-edited difference grids are added to the forecaster-weighted output grids to get a final
285 adjusted 5-km forecast grid. Complete details of the methodology for all elements are
286 documented at:
287 http://www.wpc.ncep.noaa.gov/5km_grids/medr_5km_methodology_newparms.pdf

288 Both point and gridded verification are conducted. Points are verified by the respective
289 observed station information, while the RTMA is used to verify gridded fields. The fvs
290 (described in Appendix B) is used to calculate both point-based and gridded verification of
291 sensible weather elements, including determination of statistical significance.

292

293 *b. Results*

294 Historical verification of maximum temperature at 93 points across the nation shows the
295 marked improvements in medium range temperature forecast skill over time. Today's 7-day
296 maximum temperature forecast is as accurate as a 3-day forecast in the late 1980s (Fig. 6).
297 Comparison of the 00 UTC GFS MOS forecast to the 20 UTC "final" daily issuance of the WPC
298 forecast shows the human forecaster improves upon GFS MOS (Fig. 6). Before 1998 WPC
299 forecasters were verified relative to a version of MOS termed "Kleins" (Klein and Glahn 1974).
300 Starting in 1998, WPC forecasters were verified relative to modern MOS (Glahn et al. 2009), and
301 MOS was used as the starting point for their forecasts. Differences between the Kleins and MOS
302 are apparent, with WPC forecasters improving more against Kleins (Fig. 6). The long term (30-
303 year) trend shows the human is improving less over the NWP. However, within the last seven
304 years, the WPC forecasts are improving over MOS on the order of 5% (Fig. 6). This
305 improvement may be related to a change in forecast methodology in 2004, whereby forecasters
306 use a graphical interface to apply weights to individual models and ensemble systems to derive a
307 most likely sensible weather solution. Further, ECMWF guidance became available reliably to
308 forecasters by 2008.

309 It is necessary to account for the 13-h latency between the WPC final forecast issuance
310 (19 UTC) and 00 UTC GFS MOS (Table 2). WPC issues a preliminary forecast that substantially
311 reduces this latency. Comparison of the preliminary WPC forecast issuance to the 00 and 12
312 UTC MOS is examined. This analysis also uses the full expanded set of 448 points over the
313 contiguous United States (CONUS). The results are summarized as an aggregate of monthly
314 scores averaged during the 2007–2012 period (60 months) for maximum temperature. WPC
315 accomplishes a 7–9% improvement over 00 UTC MOS with an 8-h latency, and a 4–5%
316 improvement over 12 UTC MOS with a human forecast issued 4 h *prior* to MOS (Fig. 7). Both

317 results are statistically significant at the 90% level for all days. Using a linear trend over the past
318 decade, it would take ~5 additional years for the 12 UTC GFS MOS to improve to the accuracy
319 of earlier-issued human maximum temperature forecasts.

320 One hypothesis for the improvement over MOS is that the human forecaster is adept at
321 recognizing when MOS is in large error, and thus makes large changes from MOS. Figure 8
322 shows that for frequent small changes the human forecaster makes small improvements over 12
323 UTC MOS (~5%). However, for infrequent large deviations from 12 UTC MOS [i.e., >8°F
324 (4.4°C)], forecasters usually make changes in the correct direction, exhibiting average percent
325 improvements near 15%.

326 Gridded verification allows examination of how the human gridded forecasts compare to
327 downscaled, bias-corrected international model guidance and gridded MOS (GMOS; Glahn et al.
328 2009). The WPC final forecasts are statistically significantly better than all raw downscaled
329 international model guidance and GMOS (Fig. 9a). However, bias-correction substantially
330 improves the maximum temperature model guidance – so much so that the bias-corrected
331 ECMWF ensemble mean is statistically significantly superior to the WPC gridded forecast for
332 days 5–7 (Fig. 9b).

333 Given that surface pressure patterns influence temperature and precipitation patterns,
334 further verification of the WPC mean sea-level pressure (PMSL) forecasts for days 3–7 was
335 conducted for 2012. Verification of anomaly correlation of the deterministic ECMWF and GFS,
336 and their respective ensemble system means is shown in Fig. 10. WPC has a higher anomaly
337 correlation score than all guidance at all time ranges; however, WPC is only statistically
338 significantly superior to all these gridded datasets at the day 6 forecast projection. The
339 deterministic ECMWF, which is available near the time of the final WPC forecast issuance,

340 exhibits similar skill to WPC at days 3 and 4. The 00 UTC ECMWF ensemble mean at day 7 is
341 also similar to WPC skill.

342

343 **4. Discussion and summary**

344 Analysis of multi-year verification of short-range precipitation forecasts and medium-
345 range maximum temperature forecasts from the Weather Prediction Center (WPC) are compared
346 to automated NWP guidance. Results show that human-generated forecasts improve over raw
347 deterministic model guidance when verified using both traditional methods as well as
348 contemporary methods. However, perhaps the more compelling result is that on the basis of a
349 statistical analysis of two recent years, human-generated forecasts failed to outperform the most
350 skillful downscaled, bias-corrected ensemble guidance for precipitation and maximum
351 temperature available near the same time as the human-modified forecasts.

352 Specifically, historical verification results show that the human-generated WPC QPFs
353 improve upon deterministic raw model guidance, and that the percent improvement has been
354 relatively constant over the past two decades (e.g. Fig. 4a). Medium range maximum temperature
355 forecasts also exhibit improvement over MOS. The improvement has been increasing during the
356 2005–12 period. The quality added by humans for forecasts of high-impact events varies by
357 element and forecast projection, with generally large improvements when the forecaster makes
358 changes $\geq 8^{\circ}\text{F}$ (4.4°C) to MOS temperatures in the medium range forecast. Human improvement
359 for extreme rainfall events [3-in (76.2-mm) 24 h^{-1}] is dependent on forecast projection, favoring
360 larger human improvements in the short-range forecast. Contemporary verification confirms that
361 the human forecaster makes small, but statistically significant improvement over competitive
362 deterministic model guidance for precipitation and maximum temperature.

363 However, human-generated forecasts failed to outperform the most skillful downscaled,
364 bias-corrected ensemble guidance for precipitation and maximum temperature available near the
365 same time as the human-modified forecasts. Such downscaled, bias-corrected ensemble guidance
366 represents the most skillful operational benchmark. Thus, it is premature to claim superiority by
367 the human forecaster until such forecasts are statistically significantly better than the most
368 skillful guidance. In fact, these results raise the question of whether human-generated forecast
369 superiority has ended.

370 Indeed, as computer resources advance, models will explicitly simulate more processes,
371 and more and better observations will be used by improved data assimilation systems. These
372 advances will lead to improved NWP guidance. Additionally, more sophisticated post-processing
373 of raw model guidance, including bias-correction and downscaling, will improve automated
374 forecasts of sensible weather elements. Roebber et al. (2004) cite the human ability to interpret
375 and evaluate information as an inherent advantage over algorithmic automated processes.
376 However, artificial intelligence algorithms continue to strive to simulate such human decisions –
377 for example, developing methods to automate selective consensus of ensemble members (e.g.,
378 Etherton 2007), or applying artificial neural network and evolutionary programming approaches
379 that “learn” through time (e.g., Bakhshaii and Stull 2009; Roebber 2010b). Given this future
380 environment, it is difficult to envision the human forecaster adding quality in terms of forecast
381 accuracy.

382 On the other hand, there is a distinction between long-term statistical verification (the
383 primary the focus of this paper) and critical deviations from skillful guidance in local regions and
384 cases. Contemporary post-processing approaches are best at correcting repeatable, systematic
385 errors, but struggle when the forecast sample size is small for unusual weather scenarios. The

386 forecaster's decision to deviate from skillful automated guidance in these unusual weather
387 scenarios often comes with substantial societal consequences, such as whether a snowstorm will
388 affect a city (Bosart 2003), or whether a killing freeze will occur. Thus, it is especially critical
389 that the forecaster make the very best decision in these scenarios. Figure 8 shows that when
390 forecasters make large changes from MOS, the deviations are generally in the correct direction,
391 providing evidence of skill in recognizing opportunities to deviate from MOS temperatures.
392 Obviously, more evidence of this skill for other variables, benchmarked against more skillful
393 datasets, and filtered to examine only the most critical weather scenarios is needed to more
394 conclusively demonstrate the forecaster's skill at these deviations.

395 Bosart (2003) contends that as more and more automation occurs, forecasters' skill at
396 recognizing critical opportunities to deviate from guidance may atrophy. Thus, a key component
397 of assuring the forecaster continues to add quality to NWP is keeping the forecaster engaged in
398 the forecast process. Indeed, the WPC forecasters appear to have learned how to improve over
399 the ECWMF precipitation forecasts over the past 5 years (Figs. 4a,c), perhaps learning when to
400 deviate from the skillful guidance. From the authors' experience a key to this improvement is
401 greater emphasis on using the most skillful datasets as the forecaster's starting point, and
402 encouraging changes only when confidence is high. Further, improvement can be gained with
403 greater availability of near-real time verification, using the most skillful guidance as the
404 benchmark. Finally, investment in training forecasters in the strengths and weaknesses of the
405 most skillful guidance, and providing tools to guide forecaster modifications may lead to further
406 forecaster improvements. An example of such a tool is ensemble sensitivity analysis, which can
407 indicate the source of upstream uncertainties for a given forecast parameter. As demonstrated by
408 Zheng et al. 2013, in theory, this tool allows forecasters to identify and monitor the sensitive

409 areas using available observations (satellite, aircraft or other types) in real time to assess the
410 likelihood of future scenarios.

411 Emphasis on the most skillful downscaled, bias-corrected guidance with supporting near-
412 real-time verification, forecaster training, and tools to guide forecaster interventions has only
413 recently been established at WPC, but has already resulted in forecasters making high-order
414 forecast decisions. These high-order decisions include the removal of outlier forecast guidance
415 that degrades the consensus forecast (e.g. a spurious tropical cyclone), adjusting for regime-
416 dependent biases that are not corrected (or that are introduced) in the post-processing, and
417 perhaps most importantly, deciding when to substantially deviate from the skillful guidance.
418 Thus, as additional downscaled and bias-corrected sensible weather element guidance becomes
419 operationally available, and with the support of near-real time verification, forecaster training,
420 and tools to guide forecaster interventions, a key test is whether forecasters can learn to make
421 statistically significant improvements over the most skillful of this guidance. Such a test can
422 inform to what degree, and just how quickly the role of the forecaster changes.

423 Given that only one component of the forecaster's role (accuracy) was considered and
424 only deterministic short-range QPF and medium range maximum temperature forecasts were
425 assessed, the above results must not be over generalized. Downscaling and bias-correcting of a
426 full suite of sensible weather elements is not an operational reality yet, as challenges remain with
427 elements such as wind, sky cover, ceiling, and visibility to name a few. Additionally, the
428 contemporary verification was limited to two years. Further the financial cost/benefit of human
429 involvement in the forecast process was not considered in the above analysis. Finally, a critical
430 question facing the forecasting community is if and how a forecaster may add quality to
431 ensemble guidance of many variables (e.g., Roebber et al. 2004, Novak et al. 2008). Thus, a

432 more complete investigation of the human’s role in improving upon NWP using other metrics,
433 elements, time ranges, and formats (probabilistic) is encouraged, and may lead to new paradigms
434 for human involvement in the forecast process.

435

436 *Acknowledgements.* This work benefited from insightful discussions with Lance Bosart, Brian
437 Colle, Brad Colman, Ed Danaher, Larry Dunn, Jim Hoke, Cliff Mass, Paul Roebber, David
438 Schultz, Neil Stuart, and Jim Steenburgh, as well as stimulating email exchanges on the
439 University of Albany “map” listserv. Mark Klein assisted with gathering necessary data. Jim
440 Hoke and Mike Eckert assisted with a previous version. Two anonymous reviewers provided
441 constructive comments leading to improvements in the presentation of this work. The views
442 expressed are those of the authors and do not necessarily represent a NOAA/NWS position.

443

444

445

446

447

448

449

450

451

452

453

454

455

APPENDIX A

456

Description of Pseudo Bias Corrected Ensemble QPF

457

The pseudo bias corrected ensemble QPF (ENSBC) is a series of 6-h accumulations

458

posted at 6-h intervals. Each 6-h QPF is computed in three phases:

459

1. Calculate the weighted ensemble mean (WEM).

460

2. Perform the pseudo bias correction (PBC).

461

3. Apply downscaling based on data obtained from the PRISM precipitation climatology.

462

The first phase assumes that the larger the uncertainty, the smoother the forecast should be,

463

whereas the smaller the uncertainty the more detailed the forecast should be. Two ensemble

464

means are computed. The high resolution ensemble mean is the mean of an ensemble made up of

465

relatively high-resolution deterministic single model runs (NAM, GFS, ECMWF). The full

466

ensemble mean is the mean of a high-resolution ensemble consisting of the same deterministic

467

runs along with the GEM and UKMET, and a standard ensemble system (e.g., NCEP Short-

468

Range Ensemble Forecast or NCEP Global Ensemble Forecast System). The maximum QPF

469

from the high-resolution ensemble is added as an additional member. The members of the high

470

resolution ensemble are equally weighted in the warm season, but not in the cold season

471

(October through April), and the weights are adjusted periodically with reference to verification.

472

The members of the full ensemble are equally weighted. The spread of the full ensemble is

473

obtained to compute a normalized spread, $\hat{\sigma}$, which is the full ensemble spread divided by the

474

full ensemble mean, with a small amount added to prevent division by zero. A weight value, w ,

475

is computed at each grid point:

476

$$w = \frac{\hat{\sigma}}{\hat{\sigma}_{\max}}, \quad (\text{A1})$$

477 where $\hat{\sigma}_{\max}$ is the domain maximum of the normalized spread. Then the WEM is computed at
478 each grid point:

$$479 \quad WEM = w\mu + (1-w)\mu_{hr}, \quad (A2)$$

480 where μ is the full ensemble mean and μ_{hr} is the high resolution ensemble mean. Thus,
481 where the forecast uncertainty as measured by the normalized spread is relatively large the WEM
482 is weighted toward the full ensemble mean; whereas, at points with lower normalized spread and
483 less uncertainty, the WEM is weighted toward the high resolution ensemble mean.

484 In the next phase, the WEM is passed to the PBC, which has nine tuning parameters, is
485 perpetually evolving, and undergoes fairly regular (about every six weeks or so) adjustments
486 based on verification. Here, the PBC is described in general terms.

487 For WEM 6-h precipitation amounts less than about 6—9 mm the PBC algorithm uses
488 the 10th percentile QPF from the full ensemble to reduce frequency bias (areal coverage). A
489 weighting function, ω , is applied to modify the WEM according to

$$490 \quad WEM = \omega WEM + (1-\omega) QPF_{10}, \quad (A3)$$

491 where QPF_{10} is the 10th percentile QPF from the multi-model ensemble. The weighting function
492 linearly increases to one as WEM values increase from 0 up to 6—9 mm, with higher limits for
493 longer forecast projections.

494 For WEM precipitation amounts greater than ~10 mm, the WEM is compared to the high
495 resolution ensemble mean, which is assumed to have better bias characteristics than the WEM
496 based on the findings of Ebert (2001). The algorithm iterates over an arbitrary list of increasing
497 precipitation thresholds, computes the bias of the volume of QPF exceeding the threshold for the
498 WEM relative to the high resolution ensemble mean over the entire domain, and then applies a

499 correction factor to bring this volumetric bias to unity for QPF exceeding the threshold. The
500 correction factor is constrained to range between .5 and 2.0. As the threshold value increases,
501 the high resolution ensemble mean is nudged toward the 90th percentile amount from the full
502 ensemble. This is intended to augment bias for higher thresholds, at which ensemble means tend
503 to be under-biased. The successive bias corrections alter the amount of precipitation but not its
504 placement.

505 The final phase is a downscaling based on PRISM and accomplished using correction
506 factors that vary monthly. Although more sophisticated downscaling techniques exist (Voisin et
507 al. 2010), they are too complex and computationally demanding for the development and
508 computing resources available to the WPC. This simple terrain correction is based on 5-km
509 PRISM data over the western third of the CONUS and 10-km resolution data elsewhere. The
510 method has some similarity to the terrain correction scaling used in Mountain Mapper (Henkel
511 and Peterson 1996). The PRISM data are first remapped to the 32-km WPC QPF grid,
512 preserving area averages. These values are then placed back on the high-resolution PRISM grid
513 via bilinear interpolation. Then the ratios of the original PRISM data to the back-interpolated
514 data are computed. Finally, the ratios are moved to the 32-km resolution by assigning the
515 nearest-neighboring value from the high-resolution grid. A monthly varying lower bound
516 ranging from .3 in the cold season to .9 in the warm season is imposed on the ratios. The
517 downscaling coefficients are replaced with values smoothed using a 9-point smoother at points
518 where the values are less than 1. Multiplication of the pseudo bias corrected QPF by the
519 downscaling factor completes the ENSBC processing.

520 As various model data become available, the ENSBC is executed ten times per day to
521 provide guidance for WPC forecast operations. However, a special configuration of ENSBC

522 execution is performed to create a competitive, realistic benchmark for the WPC QPF suite of
523 Day 1--3 forecasts. This configuration releases products in the same order as the WPC manual
524 forecasts for two “final” cycles per day: 00 and 12 UTC. The execution schedule permits
525 creation of products using the same models available to WPC forecasters, but without the human
526 time handicap; therefore, the automated product suite is about an hour earlier than the WPC
527 official delivery deadline for Day 1, almost two hours earlier for Day 2, and nearly four hours
528 earlier for the Day 3 forecasts. It should be noted that WPC forecasters often send products well
529 in advance of the deadlines, especially for Day 3. All comparisons to ENSBC in the main text
530 are against this benchmark.

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

APPENDIX B

547

A Description of the WPC-EMC Forecast Verification System (fvs)

548 The fvs performs three functions:

549

1. Retrieves and combines data records read from one or more Verification Statistics
550 DataBase (VSDB) text files under the control of user defined search conditions.

550

551

2. Computes performance metrics from the combined data.

552

3. Displays the performance metrics and optional statistical significance box-whisker
553 elements graphically or in a text formatted output.

553

554

The VSDB records in the text files are created by comparisons of forecast objects to observed

555

objects. This comparison is typically, but not necessarily, a forecast grid to analysis grid, a

556

forecast grid to observation points, or point forecasts to point observations. The software

557

systems used to generate VSDB record files are quite varied and not part of fvs. A single VSDB

558

record usually contains summary statistics for comparisons at multiple analysis or observation

559

points over an area or spatial volume. The summary statistics are either means or fractions. For

560

example, for verification of standardized anomalies, the following means along with the data

561

count are written in the VSDB record: means of forecast and observed anomalies, means of

562

squares of forecast and observed anomalies, and the mean of the product of forecast and

563

observed anomalies. With the data count, these means can be converted to partial sums that are

564

combined in step 1 outlined above. Another example applies to verification of dichotomous

565

events such as QPF exceeding a specific threshold for which a 2 x 2 contingency table is

566

required. In this case, each VSDB record contains fractions of forecasts exceeding the threshold,

567

observations exceeding the threshold, and both exceeding the threshold (hits). Again,

568 multiplication by the data count turns these fractions into values that can be added in combining
569 the data according to user specified search conditions.

570 In addition to the data values, each VSDB record contains information identifying the
571 forecast source, forecast hour, valid time, verification area or volume, verifying analysis,
572 parameter, and the statistic type. The statistic type is important because it determines what set of
573 performance metrics can be computed once the VSDB records have been retrieved and
574 combined. The user-defined search conditions are important because they inform the fvs as to
575 the independent variable associated with the display of the performance metrics. Any of the
576 identifier fields or combinations of them may be specified as the independent variable, so, the fvs
577 will search for and combine VSDB records as a function of different values (string or numeric)
578 for selected identifier information. The fvs will also perform consistency checks (event
579 equalization) under user direction to assure equal comparisons of multiple forecast sources. If
580 consistency checking is in force, the fvs saves the uncombined data from the search of VSDB
581 records in a binary file. The uncombined data are used in random resampling following the
582 method of Hamill (1999) if the user requests displays of box-whisker objects to depict statistical
583 significance of differences of any performance metric for paired comparisons of different
584 forecast sources.

585 Once step 1 is finished, the resulting data may be used to compute a variety of
586 performance metrics, depending on the statistic type. The fvs performs step 2 and step 3
587 seamlessly, first computing the requested metric, then generating the display. If box-whisker
588 objects are requested, the resampling is done separately at each point along the abscissa of the
589 graphical depiction during the display process. Numerous user-specified parameters are

590 provided to allow the user to control labels, text fonts, bar, line or marker characteristics, and
591 colors for the objects appearing in the graphical display.

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634

REFERENCES

Bakhshaii, A., and R. Stull, 2009: Deterministic ensemble forecasts using gene-expression programming. *Wea. Forecasting*, **24**, 1431–1451

Baldwin, M. E., S. Lakshmivarahan, and J. S. Kain, 2002: Development of an “events-oriented” approach to forecast verification. Preprints, *19th Conf. on Weather Analysis and Forecasting and 15th Conf. on Numerical Weather Prediction*, San Antonio, TX, Amer. Meteor. Soc., 210-213.

Bélair, S., M. Roch, A.-M. Leduc, P. A. Vaillancourt, S. Laroche, and J. Mailhot, 2009: Medium-range quantitative precipitation forecasts from Canada’s new 33-km deterministic global operational system. *Wea. Forecasting*, **24**, 690–708.

Brill, K. F., 2009: A general analytic method for assessing sensitivity to bias of performance measures for dichotomous forecasts. *Wea. Forecasting*, **24**, 307–318.

———, and F. Mesinger, 2009: Applying a General Analytic Method for Assessing Bias Sensitivity to Bias-Adjusted Threat and Equitable Threat Scores. *Wea. Forecasting*, **24**, 1748–1754.

Bosart, L. F., 2003: Whither the weather analysis and forecasting process? *Wea. Forecasting*, **18**, 520–529.

Brown, J. D. and D-J Seo, 2010: A nonparametric post-processor for bias correcting ensemble forecasts of hydrometeorological and hydrologic variables. *J. of Hydrometeorology*, **11(3)**, 642–665.

Caplan, P., J. Derber, W. Gemmill, S.-Y. Hong, H.-L. Pan, and D. Parrish, 1997: Changes to the 1995 NCEP operational medium-range forecast model. *Wea. Forecasting*, **12**, 581–594.

635 Clark, A. J., W. A. Gallus Jr., M. Xue, and F. Kong, 2009: A comparison of precipitation
636 forecast skill between small convection-allowing and large convection-parameterizing
637 ensembles. *Wea. Forecasting*, **24**, 1121–1140.

638 Cui, B., Z. Toth, Y. Zhu, and D. Hou, 2012: Bias correction for global ensemble forecast. *Wea.*
639 *Forecasting*, **27**, 396–410.

640 ———, Zhu, Y. and Toth, Z., 2013: Development of statistical post-processor for NAEFS.
641 Submitted to *Weather and Forecasting*.

642 Daly, C., R. P. Neilson, and D. L. Phillips, 1994: A statistical-topographic model for mapping
643 climatological precipitation over mountainous terrain. *J. Appl. Meteor.*, **33**, 140–158.

644 ———, M. Halbleib, J. I. Smith, W. P. Gibson, M. K. Doggett, G. H. Taylor, J. Curtis, and P. A.
645 Pasteris, 2008: Physiographically sensitive mapping of climatological temperature and
646 precipitation across the conterminous United States. *Int. J. Climatol.*, **28**, 2031–2064.

647 De Pondeca, M. S. F. V., and Coauthors, 2011: The real-time mesoscale analysis at NOAA’s
648 National Centers for Environmental Prediction: Current status and development. *Wea.*
649 *Forecasting*, **26**, 593–612.

650 Doswell, C. A. III, 2004: Weather forecasting by humans—heuristics and decision making. *Wea.*
651 *Forecasting*, **19**, 1115–1126.

652 Du, J., J., J. McQueen, G. DiMego, Z. Toth, D. Jovic, B. Zhou, and H. Chuang, 2006: New
653 dimension of NCEP Short-Range Ensemble Forecasting (SREF) System: Inclusion of WRF
654 Members. Preprints, *WMO Expert Team Meeting on Ensemble Prediction System*, Exeter,
655 United Kingdom, World Meteorological Organization. [Available online at
656 <http://wwwt.emc.ncep.noaa.gov/mmb/SREF/reference.html>].

657 Ebert, E. E., 2001: Ability of a poor man's ensemble to predict the probability and distribution of
658 precipitation. *Mon. Wea. Rev.*, **129**, 2461–2480.

659 Etherton, B. J., 2007: Preemptive forecasts using an ensemble kalman filter. *Mon. Wea. Rev.*,
660 **135**, 3484–3495.

661 Fritsch, J. M., and R. E. Carbone, 2004: Improving quantitative precipitation forecasts in the
662 warm season: A USWRP research and development strategy. *Bull. Amer. Meteor. Soc.*, **85**,
663 955–965.

664 Giorgi, F., 1991: Sensitivity of simulated summertime precipitation over the western United
665 States to different physics parameterizations. *Mon. Wea. Rev.*, **119**, 2870–2888.

666 Glahn, H. R., and D. P. Ruth, 2003: The new digital forecast database of the National Weather
667 Service. *Bull. Amer. Meteor. Soc.*, **84**, 195–201.

668 ———, K. Gilbert, R. Cosgrove, D. P. Ruth, and K. Sheets, 2009: The gridding of MOS. *Wea.*
669 *Forecasting*, **24**, 520–529.

670 Hamill, T. M., 1999: Hypothesis tests for evaluating numerical precipitation forecasts. *Wea.*
671 *Forecasting*, **14**, 155–167.

672 Henkel, A., and C. Peterson, 1996: Can deterministic quantitative precipitation forecasts in
673 mountainous regions be specified in a rapid, climatologically consistent manner with
674 Mountain Mapper functioning as the tool for mechanical specification, quality control, and
675 verification? *Extended Abstracts, Fifth National Heavy Precipitation Workshop*, State
676 College, PA, NWS/NOAA, 31 pp. [Available from Office of Climate, Water, and Weather
677 Services, W/OS, 1325 East-West Hwy., Silver Spring, MD 20910.]

678 Higgins, R. W., J. E. Janowiak, and Y. -P. Yao, 1996: A gridded hourly precipitation data base
679 for the United States (1963-1993). *NCEP/Climate Prediction Center Atlas 1*, U.S.
680 Department of Commerce, NOAA/NWS, 47 pp.

681 Homar, V., D. J. Stensrud, J. J. Levit, and D. R. Bright, 2006: Value of human-generated
682 perturbations in short-range ensemble forecasts of severe weather. *Wea. Forecasting*, **21**,
683 347–363.

684 Janjic, 2003: A nonhydrostatic model based on a new approach. *Meteor. Atmos. Phys.*, **82**, 271–
685 285.

686 Klein, W. H., and H. R. Glahn, 1974: Forecasting local weather by means of model output
687 statistics. *Bull. Amer. Meteor. Soc.*, **55**, 1217–1227.

688 Lin, Y. and K. Mitchell, 2005: The NCEP Stage II/IV hourly precipitation analyses:
689 development and applications. Preprints. *19th Conf. on Hydrology, San Diego, CA, Amer.*
690 *Meteor. Soc.*, Vol. 1, 2–5.

691 Magnusson, L., and E. Kallen, 2013: Factors influencing skill improvements in the ECMWF
692 forecast system. *Mon. Wea. Rev.*, **141**, 3142–3153.

693 Mass, C. F., 2003: IFPS and the future of the national weather service. *Wea. Forecasting*, **18**,
694 75–79.

695 McCarthy, P. J., D. Ball, and W. Purcell, 2007: Project phoenix: Optimizing the machine–person
696 mix in high-impact weather forecasting. Preprints, *22nd Conference on Weather Analysis and*
697 *Forecasting*, Amer. Meteor. Soc. Park City, UT.

698 Murphy, A. H., 1993: What is a good forecast? An essay on the nature of goodness in weather
699 forecasting. *Wea. Forecasting*, **8**, 281–293.

700 Novak, D. R., D. R. Bright, and M. J. Brennan, 2008: Operational forecaster uncertainty needs
701 and future roles. *Wea. Forecasting*, **23**, 1069–1084.

702 Olson, D. A., N. W. Junker, and B. Korty, 1995: Evaluation of 33 years of quantitative
703 precipitation forecasting at the NMC. *Wea. Forecasting*, **10**, 498–511.

704 Reynolds, D., 2003: Value-added quantitative precipitation forecasts: how valuable is the
705 forecaster? *Bull. Amer. Meteor. Soc.*, **84**, 876–878.

706 Roebber, P. J., D. M. Schultz, B. A. Colle, and D. J. Stensrud, 2004: Toward improved
707 prediction: High-resolution and ensemble modeling systems in operations. *Wea. Forecasting*,
708 **19**, 936–949.

709 ———, M. Westendorf, and G. R. Meadows, 2010a: Innovative weather: a new strategy for
710 student, university, and community relationships. *Bull. Amer. Meteor. Soc.*, **91**, 877–888.

711 ———, 2010b: Seeking consensus: A new approach. *Mon. Wea. Rev.*, **138**, 4402–4415.

712 Ruth, D. P., B. Glahn, V. Dagostaro, and K. Gilbert, 2009: The performance of MOS in the
713 digital age. *Wea. Forecasting*, **24**, 504–519.

714 Stuart, N.A., and Coauthors, 2006: The future of humans in an increasingly automated forecast
715 process. *Bull. Amer. Meteor. Soc.*, **87**, 1497–1501.

716 ———, D. M. Schultz, G. Klein, 2007: Maintaining the role of humans in the forecast process:
717 Analyzing the psyche of expert forecasters. *Bull. Amer. Meteor. Soc.*, **88**, 1893–1898.

718 Sukovich, E., F. M. Ralph, D. R. Novak, F. E. Barthold, and D. Reynolds, 2013: Extreme
719 quantitative precipitation forecast performance benchmarks over 11 years. *Wea.*
720 *Forecasting*, in press.

721 Voisin, N., J. C. Schaake, and D. P. Lettenmaier, 2010: Calibration and downscaling methods
722 for quantitative ensemble precipitation forecasts. *Wea. Forecasting*, **25**, 1603–1627.

723 Yussouf, N., and D. J. Stensrud, 2006: Prediction of near-surface variables at independent
724 locations from a bias-corrected ensemble forecasting system. *Mon. Wea. Rev.*, **134**, 3415–
725 3424.

726 Zheng, M., E. K. M. Chang, and B. A. Colle, 2013: Ensemble sensitivity tools for assessing
727 extratropical cyclone intensity and track predictability. *Wea. Forecasting*, in press.

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758
759
760
761
762
763
764
765
766

767
768

769

770

771

772

773

774

775

776

777

778

TABLES

Table 1. Timing of the availability of day 1 QPF guidance from the WPC, GFS, NAM, and ECMWF. The elapsed time between when guidance is available and when the WPC forecast is available (WPC latency) is shown in the right column.

Guidance Source	Time Available	WPC Latency
Overnight WPC	10 UTC	
00 UTC GFS	05 UTC	5 h
00 UTC NAM	03 UTC	7 h
00 UTC ECMWF	07 UTC	3 h
Overnight ENSBC	09 UTC	1 h

779 Table 2. Timing of the availability of medium range forecast guidance from the WPC and GFS.
 780 The elapsed time between when guidance is available and when the WPC forecast is available
 781 (WPC latency) is shown in the right column.

Guidance Source	Time Available	WPC Latency Final (prelim)
WPC Final (prelim)	19 UTC (14 UTC)	
00 UTC GFS MOS	06 UTC	13 h (8 h)
00 UTC ECMWF	08 UTC	11 h (6 h)
00 UTC ECMWF Ensemble	10 UTC	9 h (4 h)
12 UTC GFS MOS	18 UTC	1 h (-4 h)

782
 783
 784
 785
 786
 787
 788
 789
 790
 791
 792
 793
 794

795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817

FIGURE CAPTIONS

Fig. 1. Examples of WPC forecasts of (a) QPF, (b) medium range maximum temperature, and (c) medium range pressures and fronts. Examples are from different days.

Fig. 2. Timeline showing the WPC forecast period naming convention for the overnight issuance, including the forecast projection (h), time (UTC), and day 1, day 2, and day 3 designations.

Fig. 3. Time series of annual WPC threat scores for the 1-in (25.4-mm) 24-h⁻¹ threshold for the day 1 (red), day 2 (green), and day 3 (blue) forecasts from 1960–2012. Percent areal coverage of the 1-in (25.4-mm) 24-h⁻¹ threshold over the contiguous United States over the year is shown by the thin black line. Linear threat score trends are shown in respective colors. The linear trends are divided into two periods to account for increasing improvement after 1995. (Data updated yearly at: <http://www.WPC.ncep.noaa.gov/images/WPCvrf/WPC10yr.gif>)

Fig. 4. (a,c) WPC QPF percent improvement (bars) over the NAM (green), GFS (blue), and ECMWF (purple) for the (a) day 1 and (c) day 3 24-h accumulated precipitation threat score for the 1-in (25.4-mm) 24 h⁻¹ threshold during the 2001–2012 period. The frequency bias of each data set is shown as diamonds. (b,d) As in (a, c) except calculated using bias-removed threat score and including the ENSBC product. Statistically significant differences from WPC at the 90% level are marked by the black asterisk.

Fig. 5. (a,c) Comparison of the threat score (bars) and frequency bias (diamonds) for the 3-in (76.2-mm) 24 h⁻¹ threshold for (a) day 1 and (c) day 3 forecasts during the 2001–2012 period.

818 (b,d) As in (a,c) except using bias-removed threat score (bars) and including the ENSBC
819 product. Statistically significant differences in threat score from WPC at the 90% level are
820 marked by the black asterisk.

821
822 Fig. 6. Time series comparison of the WPC (solid) and 00 UTC GFS MOS (dashed) maximum
823 temperature forecast Mean Absolute Error (MAE) ($^{\circ}\text{F}$) at 98 major stations. Data are missing
824 between 1996 and 1997.

825
826 Fig. 7. (a) Comparison of 2007–2012 time-averaged maximum temperature Mean Absolute Error
827 for WPC and 00 UTC GFS MOS and WPC and 12 UTC GFS MOS for the day 3, 5, and 7
828 forecast projections. (b) WPC percent improvement over the 00 and 12 UTC GFS MOS.

829
830 Fig. 8. (top) WPC final forecast percent improvement over the 12 UTC GFS MOS at stations that
831 were adjusted from MOS during 2012. Percent improvement (left axis) for changes from $\geq 1-10^{\circ}$
832 F are displayed for day 4 to 7 forecasts. (bottom) Corresponding percentage of points adjusted
833 out of a maximum of 448 points (right axis).

834
835 Fig. 9. Comparison of 5-km gridded maximum temperature Mean Absolute Error from WPC
836 and (a) raw and (b) downscaled and bias-corrected 00 UTC ECMWF, ECMWF ensemble, GFS,
837 and GEFS over the CONUS during 2012. The RTMA is used as the verifying analysis. Due to
838 missing data, a homogeneous sample of 321 days is used in (a) and 313 days in (b). Statistically
839 significant differences from WPC at the 90% level are shown as asterisks.

840

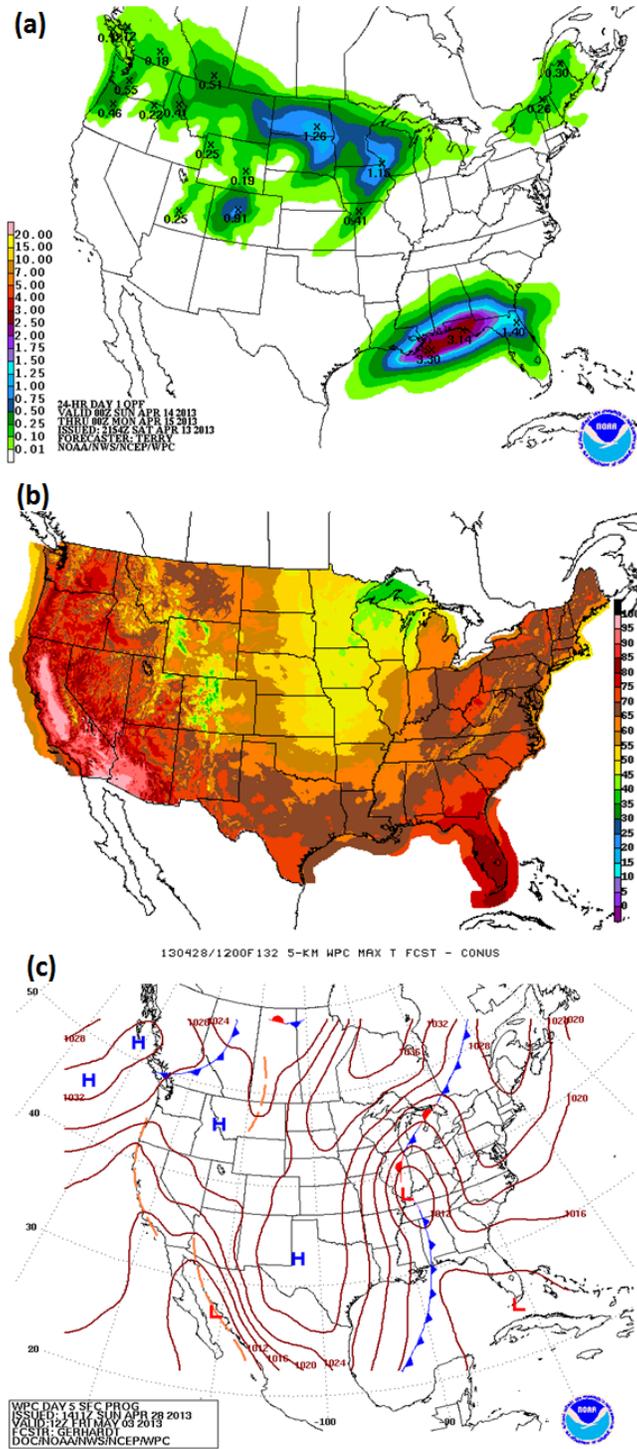
841 Fig. 10. Comparison of PMSL forecast anomaly correlation for the WPC final forecast and
842 various international model guidance. The 90% confidence interval relative to the WPC forecast
843 is shown as a box.

844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883

884

FIGURES

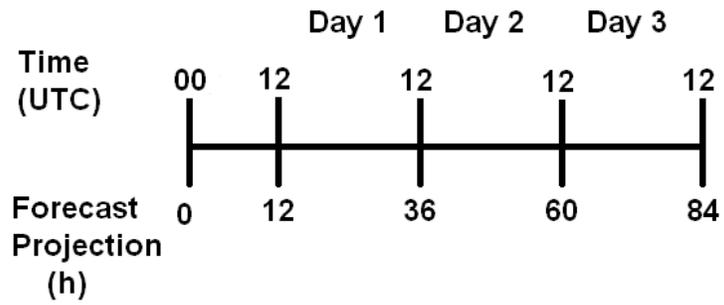
885
886
887



888
889

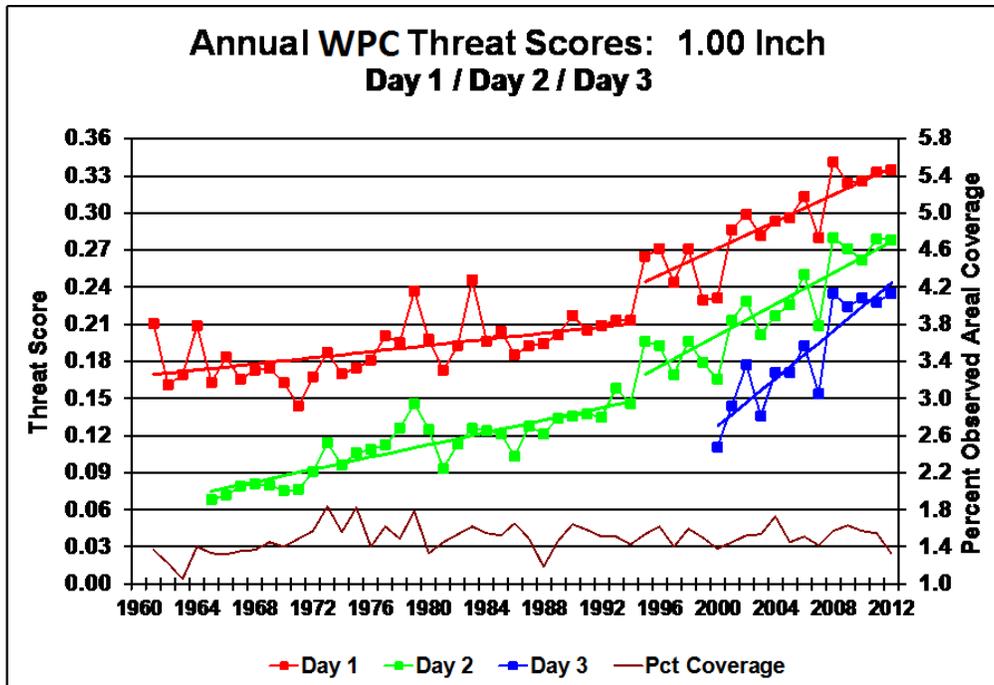
890 Fig. 1. Examples of WPC forecasts of (a) QPF, (b) medium range maximum temperature, and (c)
891 medium range pressures and fronts. Examples are from different days.

892
893
894



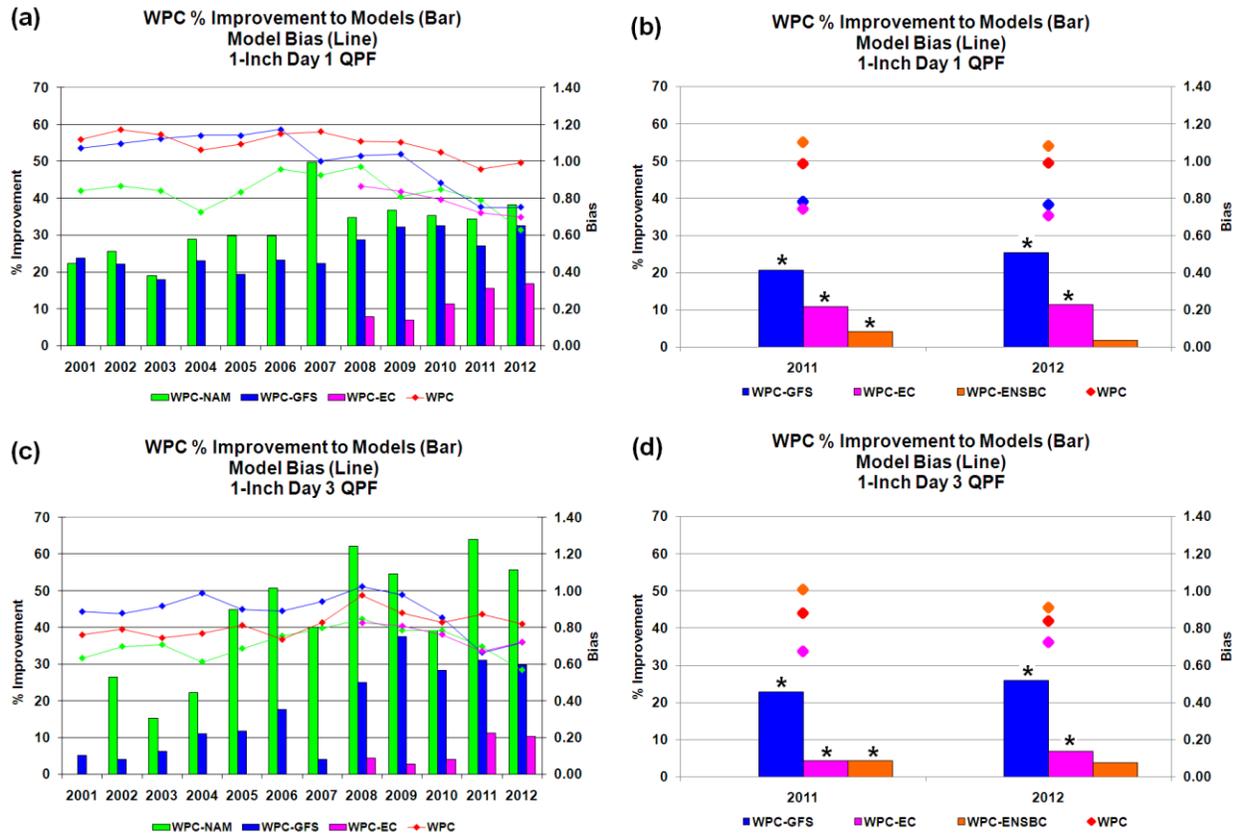
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914

Fig. 2. Timeline showing the WPC forecast period naming convention for the overnight issuance, including the forecast projection (h), time (UTC), and day 1, day 2, and day 3 designations.

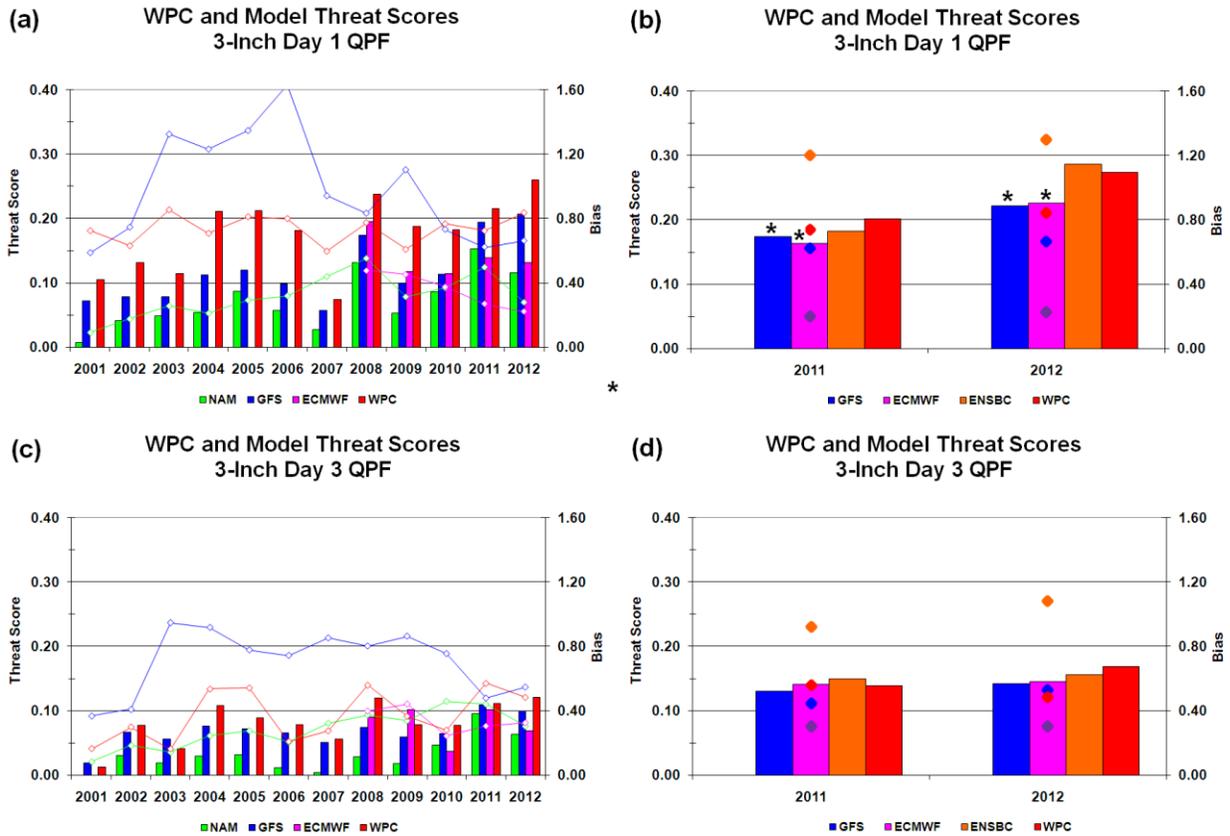


915
 916 Fig. 3. Time series of annual WPC threat scores for the 1-in (25.4-mm) 24-h⁻¹ threshold for the
 917 day 1 (red), day 2 (green), and day 3 (blue) forecasts from 1960–2012. Percent areal coverage of
 918 the 1-in (25.4-mm) 24-h⁻¹ threshold over the contiguous United States over the year is shown by
 919 the thin black line. Linear threat score trends are shown in respective colors. The linear trends are
 920 divided into two periods to account for increasing improvement after 1995. (Data updated yearly
 921 at: <http://www.WPC.ncep.noaa.gov/images/WPCvrf/WPC10yr.gif>)

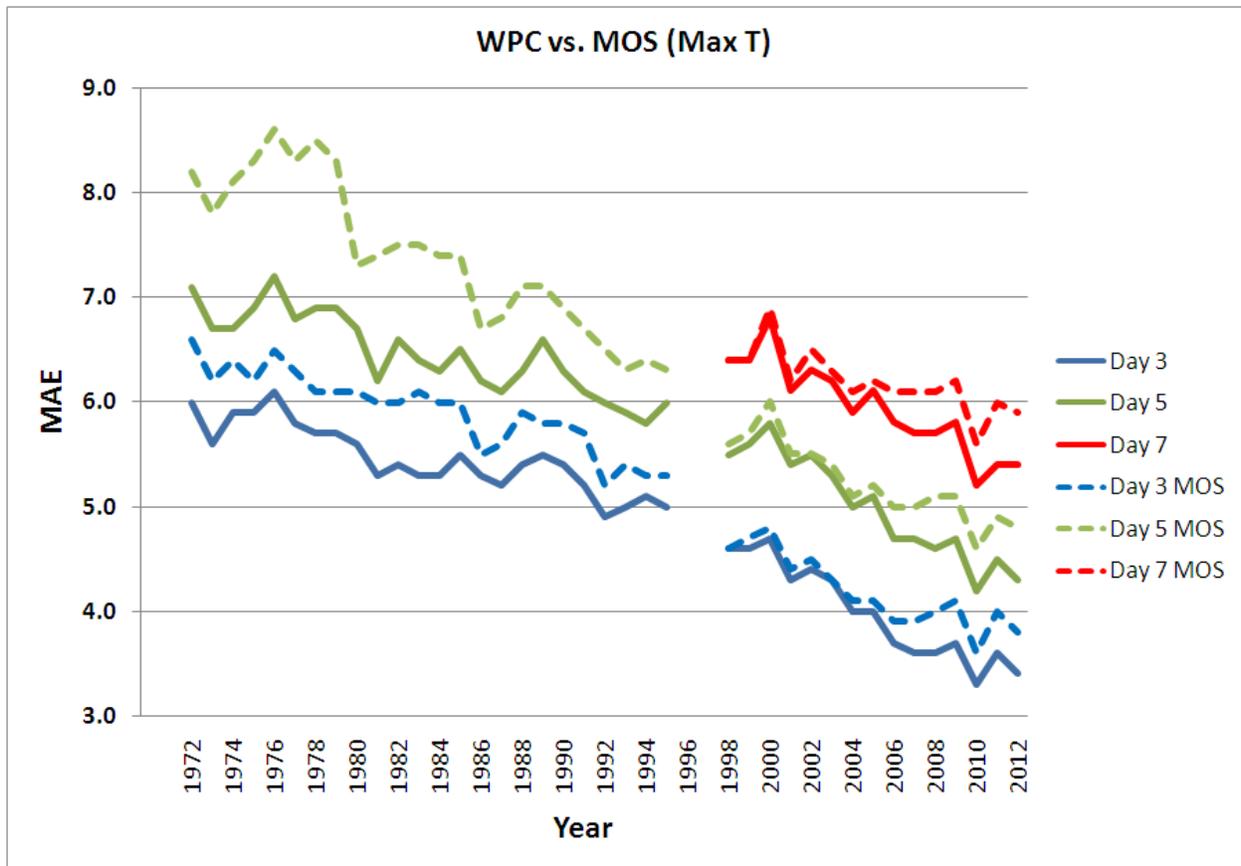
922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940



941
 942 Fig. 4. (a,c) WPC QPF percent improvement (bars) over the NAM (green), GFS (blue), and
 943 ECMWF (purple) for the (a) day 1 and (c) day 3 24-h accumulated precipitation threat score for
 944 the 1-in (25.4-mm) 24 h⁻¹ threshold during the 2001–2012 period. The frequency bias of each
 945 data set is shown as diamonds. (b,d) As in (a, c) except calculated using bias-removed threat
 946 score and including the ENSBC product. Statistically significant differences from WPC at the
 947 90% level are marked by the black asterisk.
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962

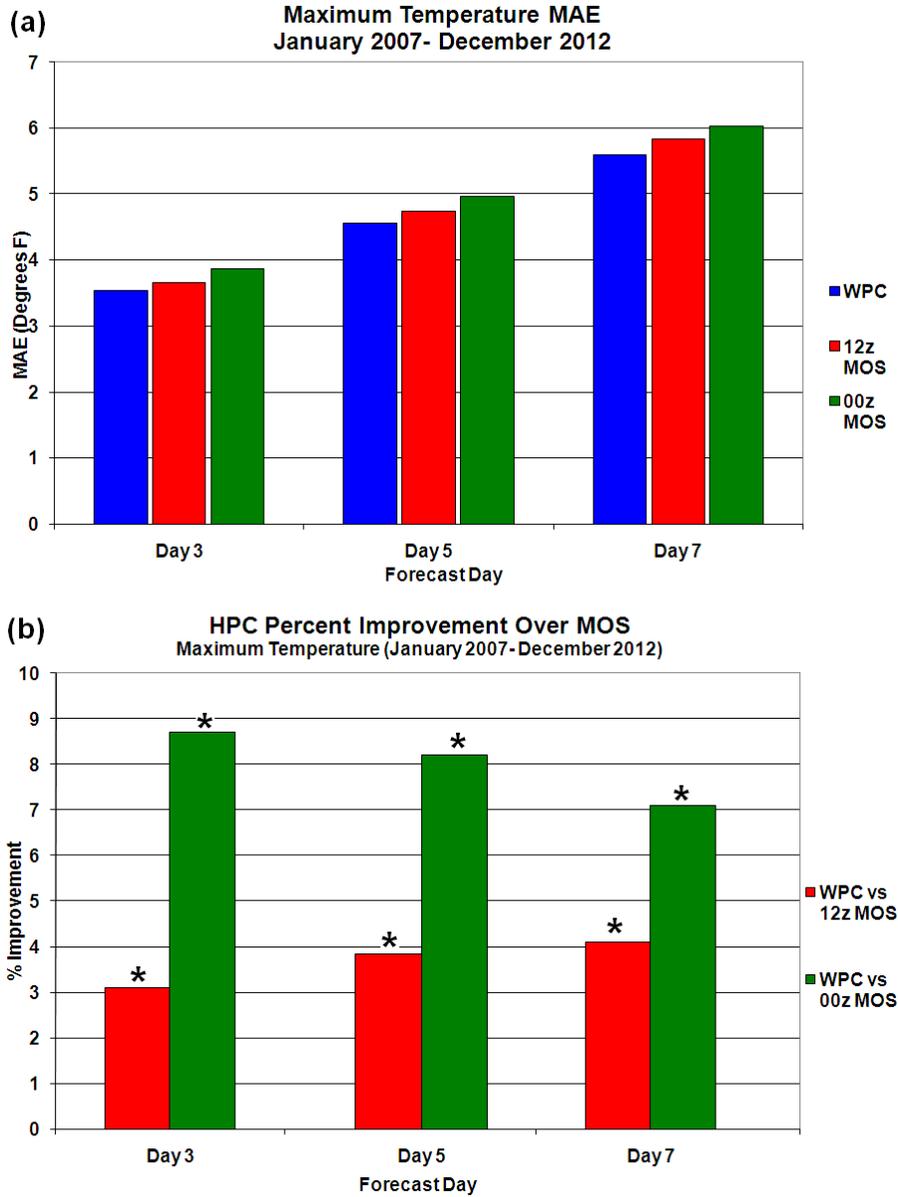


963
 964 Fig. 5. (a,c) Comparison of the threat score (bars) and frequency bias (diamonds) for the 3-in
 965 (76.2-mm) 24 h^{-1} threshold for (a) day 1 and (c) day 3 forecasts during the 2001–2012 period.
 966 (b,d) As in (a,c) except using bias-removed threat score (bars) and including the ENSBC
 967 product. Statistically significant differences in threat score from WPC at the 90% level are
 968 marked by the black asterisk.
 969



970
 971 Fig. 6. Time series comparison of the WPC (solid) and 00 UTC GFS MOS (dashed) maximum
 972 temperature forecast Mean Absolute Error (MAE) (°F) at 93 major stations. Data are missing
 973 between 1996 and 1997.
 974

975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989



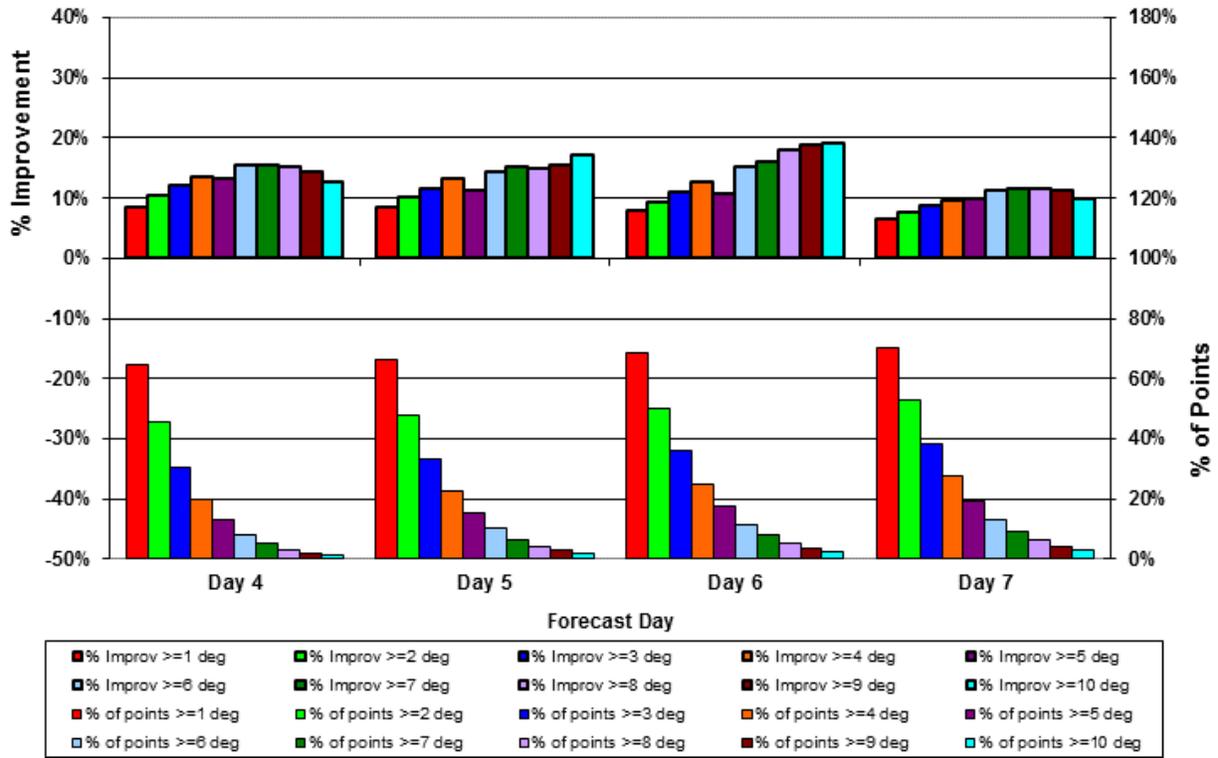
990
991

992 Fig. 7. (a) Comparison of 2007–2012 time-averaged maximum temperature Mean Absolute Error
993 for WPC and 00 UTC GFS MOS and WPC and 12 UTC GFS MOS for the day 3, 5, and 7
994 forecast projections at 448 points. (b) WPC percent improvement over the 00 and 12 UTC GFS
995 MOS. Statistically significant differences in percent improvement from WPC at the 90% level
996 are marked by the black asterisk.

997
998
999

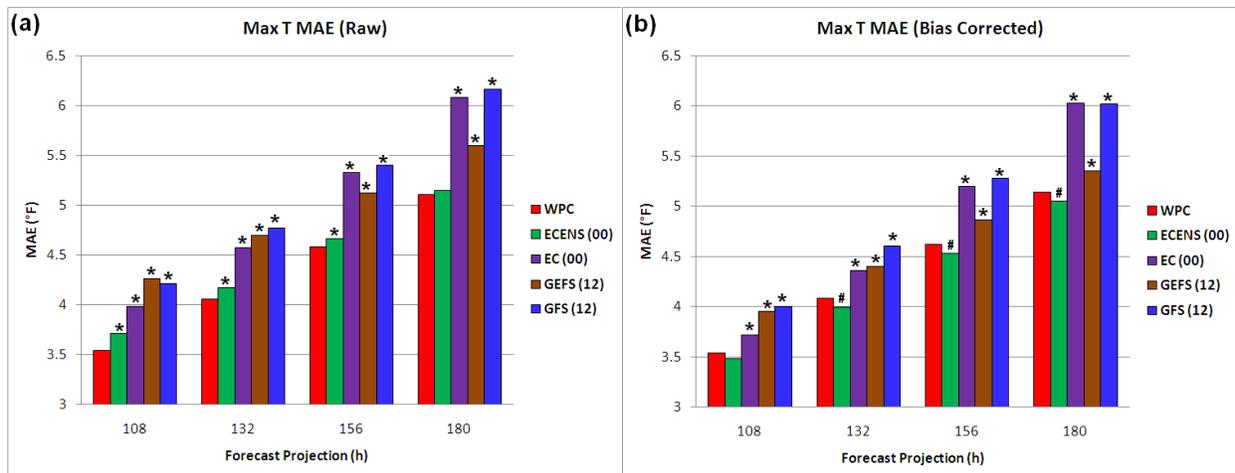
1000
1001
1002
1003

WPC Percentage Improvement Over 12z MOS
Max T MAE (Adjusted Stations Only) - January 2012 - December 2012



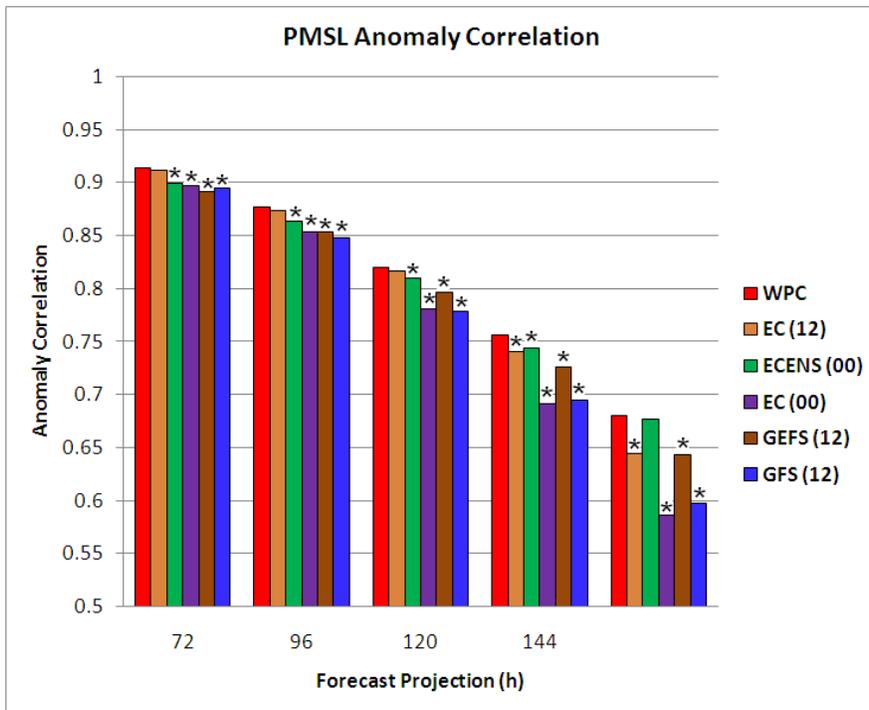
1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018

Fig. 8. (top) WPC final forecast percent improvement over the 12 UTC GFS MOS at stations that were adjusted from GFS MOS during 2012. Percent improvement (left axis) for changes from $\geq 1-10^\circ$ F are displayed for day 4 to 7 forecasts. (bottom) Corresponding percentage of points adjusted out of a maximum of 448 points (right axis).



1019 Fig. 9. Comparison of 5-km gridded maximum temperature Mean Absolute Error from WPC
 1020 (red bar) and (a) raw and (b) downscaled and bias-corrected 00 UTC ECMWF, ECMWF
 1021 ensemble, and 12 UTC GFS, and GEFS over the CONUS during 2012. The RTMA is used as the
 1022 verifying analysis. Due to missing data, a homogeneous sample of 321 days is used in (a) and
 1023 313 days in (b). Statistically significant positive (negative) differences from WPC at the 90%
 1024 level are shown as asterisks (number sign).
 1025

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035



1036
 1037
 1038
 1039

Fig. 10. Comparison of PMSL forecast anomaly correlation for the WPC final forecast and various international model guidance for 2012. Statistically significant differences from WPC at the 90% level are shown as asterisks.