1
2
# Sensitivities of the NCEP Global Forecast System

4
5
6
7
8
9

Jih-Wang A.Wang[1,2], Prashant D. Sardeshmukh[1,2], Gilbert P. Compo[1,2], Jeffrey S.
Whitaker[2], Laura C. Slivinski[1,2], Chesley M. McColl[1,2], and Philip J. Pegion[1,2]

[1]University of Colorado, CIRES, Boulder CO
[2]NOAA Earth System Research Laboratory, Physical Sciences Division, Boulder CO

15
16
17
18
19
20

22 January 2019

22

23

24

*Corresponding Author*: Dr. Jih-Wang Aaron Wang, NOAA Earth System Research

Laboratory, Physical Sciences Division, 325 Broadway, Boulder, CO, 80305.

Email: Aaron.Wang@noaa.gov

28

29

30

31
32
33
34
35
36
37

38
39 **Abstract**
40
41 An important issue in developing a forecast system is its sensitivity to additional

42 observations for improving initial conditions, to the data assimilation (DA) method used,

43 and to improvements in the forecast model. These sensitivities are investigated here for the

44 Global Forecast System (GFS) of the National Centers for Environmental Prediction

45 (NCEP). Four parallel sets of 7-day ensemble forecasts were generated for 100 forecast

46 cases in mid-January to mid-March 2016. The sets differed in their 1) inclusion or

47 exclusion of additional observations collected over the eastern Pacific during the El Niño

48 Rapid Response (ENRR) field campaign, 2) use of a Hybrid 4D-EnVar versus a pure EnKF

49 DA method to prepare the initial conditions, and 3) inclusion or exclusion of stochastic

50 parameterizations in the forecast model. The Control forecast set used the ENRR

51 observations, hybrid DA, and stochastic parameterizations. Errors of the ensemble-mean

52 forecasts in this Control set were compared with those in the other sets, with emphasis on

53 the upper tropospheric geopotential heights and vorticity, mid-tropospheric vertical

54 velocity, column-integrated precipitable water, near-surface air temperature, and surface

55 precipitation. In general, the forecast errors were found to be only slightly sensitive to the

56 additional ENRR observations, more sensitive to the DA methods, and most sensitive to

57 the inclusion of stochastic parameterizations in the model, which reduced errors globally

58 in all the variables considered except geopotential heights in the tropical upper troposphere.

59 The reduction in precipitation errors, determined with respect to two independent

60 observational datasets, was particularly striking.

61

## 1. Introduction

The large improvement in weather prediction skill over the past several decades has been described as a "quiet revolution" resulting from many small steps rather than a few dramatic leaps (Bauer *et al.,* 2015). One has now apparently entered a stage of diminishing returns in skill improvement, with no clear guidance as to improving which aspects of current forecast systems will yield the greatest benefit. Broadly speaking, forecast systems have three basic elements: 1) the input observations, 2) the data assimilation (DA) method used to merge those observations with model-generated guess fields to create the forecast initial conditions, and 3) the forecast model itself. As forecast systems continue to evolve, their relative sensitivities to these three elements will evolve as well, and it will remain important to identify the element with the largest sensitivity to help set priorities in system development.

After decades of progress, both in-situ and remotely sensed observations available for forecast initialization have become plentiful, albeit with important gaps in the tropics and polar regions (see http://www.wmo.int/pages/prog/www/OSY/GOS.html). DA techniques have also improved, in both theory and implementation. In particular, two commonly used DA methods – Ensemble Kalman Filter (EnKF; Evensen, 2003) and Four-Dimensional Variational Data Assimilation (4DVar; Lewis and Derber, 1985; Courtier *et al.*, 1994) – and their various hybrids (e.g., 4D-EnVar; see Section 2.2) have matured in merging observations with model-generated first-guess fields to provide more accurate initial conditions for forecasts. The forecast models themselves have also improved, both in their

85 representation of dynamical and physical tendencies and their use of much higher

86 horizontal and vertical resolution (e.g., references in

87 http://www.emc.ncep.noaa.gov/GFS/ref.php). These developments, together with

88 expanding computing resources, now enable several operational weather forecasting

89 centers around the world to generate ensembles of high-quality 10-day global forecasts on

90 a 50 km or finer mesh every 12 hours.

91

92 Despite this, weather forecasts continue to be far from perfect. There is room for

93 improvement in each of the three basic forecast system elements. The question is in which

94 element to invest the most effort to gain the greatest benefit. A first step toward addressing

95 this is to identify the element to which the forecasts are most sensitive. We will adopt this

96 approach here for the Global Forecast System (GFS) used at the National Centers for

97 Environmental Prediction (NCEP). Specifically, we will focus on its forecast performance

98 and sensitivities in the mid–January to mid-March 2016 period during the mature phase of

99 the 2015-16 El Niño event. An intensive observational El Niño Rapid Response (ENRR)

100 field campaign was conducted by the National Oceanic and Atmospheric Administration

101 (NOAA) over the tropical and subtropical eastern Pacific during the period (Dole *et al.*,

102 2018), and the impact of the additional observations on GFS performance is of particular

103 interest.

104

105 Section 2 provides relevant details of the additional ENRR observations, followed by a

106 description of the numerical experiments performed to test the sensitivity of the GFS

107 forecasts. Briefly, four parallel sets of 7-day 80-member ensemble forecasts were generated

108  for 100 forecast cases in the period, differing in their 1) inclusion or exclusion of the

109  additional ENRR observations, 2) use of a Hybrid 4D-EnVar versus a pure EnKF DA

110  method to prepare the initial conditions, and 3) inclusion or exclusion of stochastic physical

111  parameterizations in the forecast model. The Control forecast set used the ENRR

112  observations, hybrid DA, and stochastic parameterizations. Section 3 compares the errors

113  of the ensemble-mean forecasts in this Control set with those in the other sets, with

114  emphasis on the errors of upper tropospheric geopotential heights and vorticity, mid-

115  tropospheric vertical velocity, column-integrated precipitable water, near-surface

116  temperature, and surface precipitation. A summary and concluding remarks follow in

117  Section 4, emphasizing that although only a limited set of GFS sensitivities were

118  investigated here, our methodology could also be fruitfully applied to investigate the

119  sensitivities of other forecast systems to their three basic elements.

120

121  **2. Additional observations and experimental design**

122

123  *2.1 ENRR Field Campaign*

124  As discussed by Dole *et al*. (2018), a strong El Niño event was projected to occur in the

125  northern winter and spring of 2015-16 based on observed tropical Pacific sea surface

126  temperature (SST) anomalies in the preceding summer. NOAA seized this opportunity to

127  undertake the ENRR field campaign to record the event while it was ongoing. The extra

128  observations collected included 1) dropsonde, radar, and microwave radiometer

129  observations from campaign flights (mostly within 180º-135ºW and between Honolulu and

130  the equator), 2) radiosonde and surface observations from campaign cruises (Honolulu to

5

San Diego), 3) radiosonde and surface observations from Kiritimati Island (1.9˚N, 157.4˚W), and 4) radar observations from the U.S. west coast. These ENRR observations, together with the far more numerous routine conventional and satellite observations over the globe, provide an excellent opportunity to examine the impact of such event-oriented field campaign observations on weather forecast skill. The upper-air radiosonde and dropsonde observations covered most of the ENRR campaign area; there were 22,510 humidity observations, 33,646 temperature observations, and 35,943 wind observations by radiosondes and dropsondes from January 20 to March 16, 2016. We focus here on the forecast impact of only the upper-air radiosonde and dropsonde observations from the campaign, referring to them as "the ENRR observations". Full details of the campaign can be found in Dole *et al*. (2018) and at https://www.esrl.noaa.gov/psd/enso/rapid_response/, as well as in Slivinski *et al.* (2018).

2.2 *Analyses – Initial Conditions and "Truth"*

For clean comparisons, we generated our own analyses to provide initial conditions for our 7-day forecasts. We used the same 64-level version of NCEP's GFS model (Environmental Modeling Center, 2003) operational in April 2016 but at a lower horizontal resolution (spectral truncation of 254, approximate grid spacing of 50 km) for all the analyses and forecasts. To generate the analyses using NCEP's Global DA system, we performed sequential 6-hourly forecast-analysis cycles comprising the following steps:

Step 1: Combine an 80-member ensemble of 0- to 6-hr forecasts with observations in that 6-hour window to generate an 80-member ensemble of preliminary analyses.

154    Step 2: Perform IAU (incremental analysis update; see below for more details) from

155        hr-0 to hr-6 to generate the "ultimate" analyses and continue running the 80-

156        member ensemble for the next 6-hr background (i.e, first guess) ensemble of

157        forecasts.

158

159    Step 3: Repeat Steps 1 through 2 for the next cycle.

160

161    In Step 1, we used either the Ensemble Kalman Filter method (EnKF; Evensen, 2003) or

162    the Hybrid Four-Dimensional Ensemble Variational method (Hybrid 4D-EnVar; Buehner

163    *et al.*, 2013; Kleist and Ide, 2015). The EnKF method is a Monte Carlo approximation of

164    the Kalman Filter. It uses a model ensemble of finite size to approximate the probability

165    distribution of predicted states, and updates the model-generated *a priori* state variables to

166    *a posteriori* variables by using the model ensemble covariance to estimate the Kalman gain

167    (Evensen, 2003). A reasonably large ensemble size is required for this purpose, and also to

168    avoid abrupt imbalances among the state variables being updated. The problem of abrupt

169    imbalances is partly overcome in Step 2 through an incremental analysis update (IAU;

170    Bloom *et al.*, 1996; Lei and Whitaker, 2016; Takacs *et al.*, 2018), which divides the

171    analysis increment from a preliminary analysis cycle into small portions and repeats the

172    background forecast by adding the portions as extra forcing to the forecast at every time

173    step. The final background forecast is the ultimate analysis, which closely resembles the

174    preliminary analysis at the end of the forecast-analysis cycle but does not have abrupt

175    imbalances, and is continued as the preliminary forecast for the next forecast-analysis

176    cycle. For the present study, each analysis that we used for model initialization and

177    verification purposes was the preliminary analysis (i.e., the output of EnKF or Hybrid DA

178    before application of the IAU forcing) in the current forecast-analysis cycle, but it had the

179    IAU forcing from the beginning of the experiment period (i.e., Jan 20, 2016; see Fig. 1 and

180    context) up to the previous forecast-analysis cycle. There are two options in the NOAA

181    EnKF code: the serial Ensemble Square Root Filter (EnSRF) and the Local Ensemble

182    Transform Kalman Filter (LETKF). The EnSRF used here is also implemented

183    operationally in the atmospheric GFS at NOAA. It is based on the serial EnSRF described

184    in Whitaker and Hamill (2002) and uses the parallel algorithm described in Anderson and

185    Collins (2007) for computational efficiency.

186

187    The Hybrid 4D-EnVar is a combination of EnKF and 4DVar (Four Dimensional

188    Variational method; Lewis and Derber, 1985; Courtier *et al.*, 1994) which aims (a) to

189    combine the time-varying ensemble covariances with static background error covariances

190    to estimate the total background error contribution to the cost function being minimized,

191    and (b) to eliminate the use of tangent-linear (TL) and adjoint (AD) models used in pure

192    4DVar (Wang *et al.*, 2008; Buehner *et al.*, 2013; Kleist and Ide, 2015).

193

194    In addition to the inclusion of a static background error covariance, the Hybrid 4D-EnVar

195    differs from the EnKF in the way 'covariance localization' is performed. Covariance

196    localization is a method for dealing with spurious covariances at large spatial lags that

197    result from using small ensemble sizes. In the Hybrid 4D-EnVar system, covariance

198    localization is performed in model space (Houtekamer and Mitchell, 2001) instead of

199    observation space (Gaspari and Cohn, 1999; see summary of both in Lei and Whitaker,

8

200    2015). This can significantly impact the assimilation of observations such as satellite

201    radiances, which involves using complicated forward observation operators to link the

202    model state to the radiances (Campbell *et al.*, 2009). In the global numerical weather

203    prediction (NWP) system of the National Weather Service (NWS), an 80-member EnKF

204    is run operationally to initialize the Global Ensemble Forecast System (GEFS) and to

205    provide ensemble covariances for the Hybrid 4D-EnVar data assimilation (Kleist and Ide,

206    2015) used by the Grid-point Statistical Interpolation (GSI) analysis system that generates

207    the high-resolution deterministic analysis for the high-resolution GFS forecasts. In our

208    analyses, we did not separately perform high-resolution deterministic analyses or forecasts;

209    instead, we substituted the ensemble mean as the deterministic solution so that the

210    interpolation from one resolution to another was avoided.

211

212    We performed the DA in Step 1 by using either the EnKF or Hybrid method, and either

213    including or excluding the ENRR observations, thus generating four separate sets of 80-

214    member ensemble analyses for the ENRR period. Given computing and storage constraints,

215    we worked mainly with the Hybrid-with-ENRR set (hereafter the Control analysis set), the

216    Hybrid-without-ENRR set (hereafter the Denial analysis set), and the EnKF-with-ENRR

217    observations (hereafter the EnKFonly analysis set). These three sets of analyses were then

218    used as initial conditions for three separate sets of 7-day 80-member ensemble forecasts.

219    For forecast verification, we could have used any one of these three analysis sets as "truth".

220    However, we chose the Control analysis set for this purpose as our "best" analysis product,

221    both because of its assimilation of all observations (including the ENRR observations) and

222    its improved quality resulting from the hybridization. Using the EnKFonly or Denial

223  analyses instead of the Control analyses for forecast verification did not affect any of our

224  findings for forecasts beyond 24 hours.

225

226  *2.3 Forecasts and Evaluations*

227  The three analysis sets were used to initialize three sets of 7-day forecasts every 12 hours

228  in the 57-day (20 January to 16 March) ENRR period. We will henceforth refer to these as

229  Control, Denial, and EnKFonly forecasts, respectively. Their performance was evaluated

230  by comparing them with the verifying Control analyses, and with independent

231  observational estimates in the case of precipitation. The impact of the ENRR observations

232  was gauged by comparing the skill of the Control and Denial forecasts, and the impact of

233  the DA method by comparing the skill of the Control and EnKFonly forecasts. Table 1 lists

234  these three sets of forecasts and their relevant characteristics.

235

236  All three forecast sets used stochastic parameterizations (SPs) to perturb the deterministic

237  physical tendencies in the model. The use of SPs in operational forecasts is usually

238  motivated by a need to increase the ensemble spread to make it more consistent with the

239  generally larger root-mean-square error (RMSE) of ensemble-mean forecasts. Such a

240  consistency is also implicitly assumed in the EnKF. The GFS SP module can employ three

241  different types of SPs, namely SPPT (Stochastically Perturbed Physical Tendencies;

242  Palmer *et al.*, 2009; Shutts *et al.*, 2011), SHUM (Stochastic HUMidity perturbations in the

243  boundary layer; Tompkins and Berner, 2008), and SKEB (Stochastic Kinetic Energy

244  Backscatter; Berner *et al.*, 2009), to increase the ensemble spread. The SPPT scheme has

245  the following general form for the tendency perturbation:

246
$$\dot{x}_p = (1 + r\mu)\dot{x}_c ,$$

247 where $\dot{x}_c$ and $\dot{x}_p$ are the physical tendencies of the state variable before and after applying

248 the stochastic perturbation, respectively; $r$ is a stochastic horizontal weight that is bounded

249 in the interval [-1,1] by using an inverse logit transform of a Gaussian distribution, and $\mu$

250 is a vertical weight that is 1 between the surface and 100hPa and is tapered to zero at 25hPa.

251 The horizontal weight $r$ can be represented in terms of spherical harmonics as

252
$$r = \sum_{mn} \hat{r}_{mn} Y_{mn} ,$$

253 where $\hat{r}_{mn}$ is the spherical harmonic coefficient of $r$ for total wavenumber $n$ and zonal

254 wavenumber $m$. This enables the tendency perturbation to be made scale-aware and

255 smoothed in space to the degree desired. Palmer *et al.* (2009) (see also Sardeshmukh, 2005)

256 represented $\hat{r}_{mn}$ as a combination of a first-order autoregressive AR(1) process and

257 spatially smoothed white noise as

258
$$\hat{r}_{mn}(t + \Delta t) = \phi \hat{r}_{mn}(t) + \sigma_n \eta_{mn}(t) ,$$

259 where $\Delta t$ is the model time step, $\phi = exp(-\Delta t / \tau)$ is the AR(1) coefficient, $\sigma_n$ is the

260 standard deviation (i.e., strength) of the tendency perturbation, and $\eta_{mn}(t)$ is a Gaussian

261 random number with zero mean and unit variance. $\sigma_n$ is a function of total wavenumber $n$

262 and spatial autocorrelation length scale $L$ such that the variance in grid space *Var(r)* is

263 uniform and the spatial pattern has a spatial autocorrelation corresponding to the equivalent

264 of a Gaussian function on the sphere (Palmer *et al.*, 2009; Sardeshmukh, 2005; Weaver and

265 Courtier, 2001). The SPPT scheme is applied to the tendencies of zonal wind, meridional

266 wind, specific humidity, and temperature induced by the GFS physics package, but not to

267 the tendencies induced by the clear-sky radiation scheme.

268

269   The SHUM perturbations are similar to the SPPT perturbations, except that they are applied

270   to the humidity itself and not the humidity tendency (although they may be interpreted as

271   perturbations to the humidity tendency integrated over a model time step), and only in the

272   lower troposphere. The formula is

273   $$q_p = (1 + r\mu)q_c \, ,$$

274   where $q_c$ and $q_p$ are the specific humidity before and after the stochastic perturbation

275   respectively. The vertical weight $\mu$ decays exponentially in pressure away from the surface.

276   The scheme additionally constrains the specific humidity to remain positive.

277

278   We used SPPT and SHUM perturbations (but not SKEB perturbations) in all three sets of

279   forecasts. We could have specified multiple values of the AR(1) e-folding time scale $\tau$,

280   spatial variance *Var(r)*, and spatial autocorrelation scale $L$ to avoid the early saturation of

281   ensemble spread at small scales. However, for simplicity we chose fixed values of $\tau$= 6

282   hours, *Var(r)* = 0.8 and $L$ = 500 km for the SPPT, and $\tau$= 6 hours, *Var(r)* = 0.005 and $L$ =

283   500 km for the SHUM perturbations.

284

285   Finally, in order to quantify the impact of the SPs, we generated a fourth set of 7-day

286   forecasts similar to the Control forecasts but without SPs (labeled noSP; see Table 1). As

287   with the other three forecast sets, the skill of the noSP forecasts was evaluated by

288   comparing with the verifying Control analyses, and the impact of the SPs was gauged by

289   comparing the skill of the Control and noSP forecasts.

290

291     To summarize, the Control, Denial, EnKFonly and noSP forecasts were each 7-day 80-

292     member ensemble forecasts, started twice a day at 00Z and 12Z in the 57-day ENRR

293     period. There were thus 114 forecast cases in each set. The forecast output frequency was

294     3 hours (i.e. 3, 6, 9, ..., 168 hours). To ensure the same number of forecast verifications for

295     all forecast lead times, we only evaluated forecasts valid between January 27 and March

296     16. As illustrated in Fig. 1, this verification period spans 50 days and contains 100

297     verification cases (with each case corresponding to one initialization time) for each forecast

298     lead time. Overall, for each forecast lead time we thus had 4 sets × 80 forecasts × 100 cases

299     = 32,000 forecasts of all model variables at all grid points. We shall show below that these

300     large sample sizes enable us to quantify the impacts of the ENRR observations, DA

301     methods, and SPs on the forecast skill with statistical confidence.

302

303     **3. Forecast Evaluation and Comparisons**

304

305     *3.1 Forecast Errors*

306

307     We define the forecast error as the RMSE of the $M$=80 member ensemble-mean forecast

308     with respect to the 80-member ensemble-mean Control analysis, determined over all

309     $N$=100 forecast cases as

310
$$RMSE(t) = \left\{ \frac{1}{N} \sum_{n=1}^{N} V'^{2}_{n,t} \right\}^{1/2},$$

311     where

$$312 \qquad V'_{n,t} = V_{f,n,t} - V_{a,n} = \frac{1}{M}\sum_{m=1}^{M} V_{f,n,t}^{m} - \frac{1}{M}\sum_{m=1}^{M} V_{a,n}^{m}$$

313   Here subscript $t$ refers to forecast lead time, $f$ and $a$ to the forecast or verifying analysis of

314   variable $V$, $n$ to the forecast case number, and $m$ to the ensemble member number. This

315   expression was used to calculate *RMSE(t)* for selected variables at each grid point. An

316   analogous expression, with the area-weighted gridpoint values of $V'^{2}_{n,t}$ averaged

317   additionally over the globe as well as over some specific regions, was used to calculate

318   global and regional values of *RMSE(t)*. We focus here on the forecast errors of geopotential

319   height at 200 hPa ($Z_{200hPa}$), relative vorticity at 200 hPa ($\xi_{200hPa}$), vertical velocity at 500

320   hPa ($\omega_{500hPa}$), column-integrated precipitable water (PWAT), and 2-meter air temperature

321   ($T_{2m}$). The RMSEs for a few additional variables were also examined but are not shown

322   here due to their similar behavior.

323

324   For precipitation, we compared forecasts of 12-hour accumulated precipitation values

325   (AP12HR) with two independent observational datasets: the NASA (National Aeronautics

326   and Space Administration) GPM (Global Precipitation Measurement) dataset (Huffman *et*

327   *al.,* 2014) and the PERSIANN (Precipitation Estimation from Remotely Sensed

328   Information using Artificial Neural Networks) dataset (Sorooshian *et al.,* 2014; Ashouri *et*

329   *al.,* 2015). For brevity, we only show the comparison with the NASA GPM dataset, since

330   the comparison with the PERSIANN dataset yielded similar results.

331

332   Fig. 2 shows the area-weighted global RMSEs of the Control, Denial, EnKFonly, and noSP

333   forecasts of $Z_{200hPa}$, $\xi_{200hPa}$, $\omega_{500hPa}$, PWAT, and $T_{2m}$ at 12-hourly intervals up to 7 days (hr-

334      168), as well as the RMSEs of AP12HR between 20ºS and 20ºN and between 60ºS and

335      60ºN. The initial (hr-0) error of the Denial forecasts reflects the difference between the

336      Control and Denial analyses (not shown). The Control forecasts have slightly smaller errors

337      than the Denial forecasts until hr-24 but show no discernible impact thereafter, at least in

338      this global metric, of including the ENRR observations in the initial conditions.

339

340      In contrast, the global RMSEs of the EnKFonly forecasts are larger than those of the

341      Control and Denial forecasts throughout the forecast period. Indeed, the EnKFonly

342      forecasts are worse than the Control forecasts beyond Day 1 even when both are verified

343      against the EnKFonly analyses (not shown) instead of the Control analyses as in Fig. 2. We

344      should stress that this result does not imply that an EnKF method is inferior to a Hybrid

345      method in general. One can think of several ways in which our particular implementation

346      of the EnKF algorithm could have been improved, such as by adjusting the vertical

347      covariance localization of the satellite radiance observations, by improving the balance

348      constraints on analysis increments, and by increasing the ensemble size of the ensemble

349      Kalman Filter. Nevertheless, Fig. 2 clearly demonstrates the greater sensitivity of the

350      forecast errors to initial conditions prepared using different DA methods than to the

351      inclusion or exclusion of the ENRR observations in those initial conditions.

352

353      The global RMSEs of the Control forecasts are smaller than those of noSP forecasts for

354      $\omega_{500hPa}$, $\xi_{200hPa}$, and PWAT throughout the 7-day forecast range, demonstrating the

355      beneficial impact of including SPs in the model. Similar reductions in ensemble-mean

356      forecast errors have been reported in other forecast systems (e.g., Leutbecher *et al.,* 2017).

357     The global RMSEs of the noSP forecasts are larger than those of the EnKFonly forecasts

358     after Day 3 for $\omega_{500hPa}$, Day 6 for $\xi_{200hPa}$, and Day 5 for PWAT. In other words, beyond Day

359     3 these forecasts errors are more sensitive to including or not including SPs in the forecast

360     model than they are to the use of the Hybrid versus EnKF DA method to prepare the

361     forecast initial conditions. The $\omega_{500hPa}$ errors saturate by about Day 6 (Fig. 2c), but

362     interestingly the PWAT errors do not saturate even by Day 15 (not shown). The

363     precipitation errors (Fig. 2f) saturate at an intermediate lead time of about Day 7. Although

364     $\omega_{500hPa}$ and PWAT are both important for determining precipitation strength, the near-

365     simultaneity of $\omega_{500hPa}$ and precipitation error saturation suggests that $\omega_{500hPa}$ has a stronger

366     control than PWAT on determining precipitation variations on the time scales of synoptic

367     weather (see also Sardeshmukh *et al.*, 2015).

368

369     The error growth curves of $T_{2m}$ (Fig. 2e) and precipitation (Fig. 2f) in the Control, Denial,

370     EnKFonly, and noSP forecasts have a similar general character to that of the other

371     variables, with little or no sensitivity to the ENRR observations, considerably higher

372     sensitivity to the choice of the Hybrid versus EnKF DA method, and greatest sensitivity to

373     the use of SPs in the model. For all variables in Fig. 2 except $Z_{200hPa}$, the Control forecasts

374     are the best and the noSP forecasts are the worst by Day 7. The impact of the SPs is

375     evidently cumulative over time, resulting by Day 7 in a reduction of the precipitation

376     forecast error in the Control forecasts by ~4.3% in the 20ºS-20ºN latitude domain and by

377     ~3% in the 60ºS-60ºN latitude domain.

378

379     Note that the errors of the 12-hour accumulated precipitation amounts in all four forecast

380     sets, measured with respect to the observational GPM values, are already quite large ($> 6.5$

381     mm) at hr-12. The GPM precipitation is a blend of radar-reflection and radiance based

382     precipitation estimates from multiple satellites, and is calibrated against in-situ ground

383     observations. For a cleaner comparison with the precipitation forecasts, we integrated the

384     30-minute $0.1°$ resolution GPM values to 12-hr $0.5°$ resolution values. Given that

385     precipitation is a positive semi-definite quantity, its substantial error even at short forecast

386     ranges suggests that there are precipitation events of which locations and large magnitude

387     ($> 100$mm accumulations in 12 hours) are not captured by our forecasts.

388

389     The general conclusions drawn from the global forecast error growth curves in Fig. 2 are

390     also valid for limited regions. To illustrate this, Fig. 3 shows the RMSEs of $\omega_{500hPa}$ in the

391     Northern Hemisphere ($20°$N-$90°$N), Southern Hemisphere ($20°$S-$90°$S), Tropics ($20°$S-

392     $20°$N), and the contiguous United States (CONUS; $125°$W-$66°$W, $24°$N-$50°$N). The errors

393     saturate in the Northern Hemisphere, Southern Hemisphere, and Tropics by Day 7, and

394     nearly saturate in the CONUS region by the end of Day 7. Geographically, the errors are

395     largest in the extratropical storm track regions and in areas of tropical deep convection

396     (Fig. 4a). They are particularly large over the CONUS region, not surprisingly because the

397     region overlaps strongly with the northern hemispheric storm track at those longitudes, but

398     also possibly because of erroneous model representations of the influence of the Rocky

399     Mountains on synoptic weather systems.

400

401    A beneficial impact of the ENRR observations on the regional $\omega_{500hPa}$ forecasts is not

402    discernible in Fig. 3 beyond Day 1, which reflects an average of small differences of mixed

403    signs between the Control and Denial forecasts. For instance, small positive and negative

404    impacts on Day 7, likely not statistically significant, are scattered around the globe (Fig.

405    4b) with no coherent geographical structure. On the other hand, using the Hybrid versus

406    the EnKF initial conditions leads to smaller Day-7 errors in many though not all regions

407    (Fig. 4c). However, including SPs in the model unambiguously reduces the $\omega_{500hPa}$ error

408    almost everywhere on the globe (Fig. 4d). The improvement is particularly clear in the

409    Northern Hemisphere storm track and tropical convective regions.

410

411    Given the strong link between $\omega_{500hPa}$ and precipitation on synoptic time scales, the results

412    for the precipitation errors in the Control forecasts and how they differ from the errors in

413    the other three forecast sets (Fig. 5) are highly consistent with the results for the $\omega_{500hPa}$

414    errors in Fig 4. Similar to the $\omega_{500hPa}$ errors, the precipitation errors are least sensitive to

415    including or excluding the ENRR observations, more sensitive to the choice of the Hybrid

416    versus EnKF DA method used to initialize the forecasts, and most sensitive to using or not

417    using the SPs in the forecast model.

418

419    Fig. 6 shows the errors of near-surface air temperature ($T_{2m}$) in the Control forecasts and

420    how they differ from the errors in the other three forecast sets. Note that the prescribed SST

421    boundary conditions are updated daily in the analyses but not in the 7-day forecasts. Still,

422    because the SSTs vary little and the $T_{2m}$ values over the ocean are tightly linked to them,

423    the $T_{2m}$ RMSE over the oceans remains relatively small over the 7-day forecast range. Also,

424 because the prescribed SSTs are identical in all the four forecast sets, the differences of the

425 $T_{2m}$ errors over the oceans among the forecast sets are small as well. The Control forecast

426 errors are larger over land and largest in high latitudes (Fig. 6a). The differences between

427 the RMSEs of the Control and Denial forecasts are also large over high-latitude land, but

428 with mixed signs (Fig. 6b). The impact of the choice of the Hybrid over the EnKF DA

429 method is stronger than the impact of the ENRR observations (cf. Figs. 6c and 6b).

430 Including the SPs again has the largest impact (Fig. 6d), with an unambiguous reduction

431 of the $T_{2m}$ error almost everywhere, but especially over land areas.

432

433 Using SPs is clearly beneficial for the $\omega_{500hPa}$, precipitation, and $T_{2m}$ forecasts over most of

434 the globe. For upper tropospheric geopotential heights ($Z_{200hPa}$), however, the benefit is not

435 so clear-cut. The impact is negligible in the extratropics and negative in the tropics, as

436 shown in Fig. 7 for the same four regions as in Fig. 3. The Control and Denial forecast

437 errors are again very similar, except in the CONUS region where the Control errors are

438 slightly smaller than the Denial errors on Days 3-5 (Fig. 7d). Perhaps this is to be expected,

439 given that the CONUS region is downstream of the region of the ENRR observations. We

440 also show below in Section 3.2 that even though the positive impact of the ENRR

441 observations is weak, there is a recognizable enhancement of El Niño-related features over

442 North America in $Z_{200hPa}$ due to the ENRR observations.

443

444 It is evident that the $Z_{200hPa}$ RMSE sensitivity to the DA methods is different in the Northern

445 Hemisphere, Southern Hemisphere and Tropics (cf. Figs. 7a, 7b, 7c). Using the Hybrid

446 versus the EnKF method has a large positive impact on the $Z_{200hPa}$ forecasts in the Southern

447    Hemisphere, a weaker positive impact in the Northern Hemisphere, but a negative impact

448    in the Tropics starting from about Day 2. Interestingly, using the Control (Hybrid DA)

449    versus the EnKFonly analyses as initial conditions also increases the positive tropical bias

450    of the Day-7 $Z_{200hPa}$ Control forecasts (cf. Figs. 9a, 9c). The EnKFonly analyses have lower

451    $Z_{200hPa}$ than the Control analyses in the tropics, resulting from several methodological

452    differences in the EnKF algorithm, including (a) covariance localization of satellite

453    radiances (see Lei *et al.* (2019) for a recent study); (b) lack of additional balance constraints

454    on analysis increments; (c) no static background error covariances; and (d) use of

455    maximum likelihood versus minimum variance estimation as in 4D-EnVar. While both

456    Control and EnKFonly forecasts develop positive tropical biases over 7 days, the

457    EnKFonly forecasts are closer to the truth and have smaller RMSEs. The forecast model

458    drift toward higher $Z_{200hPa}$ in the tropics is worthy of further investigation. With regard to

459    the impact of SPs on the $Z_{200hPa}$ forecasts, their positive impact does not become clear in

460    the global RMSE metric until the end of Day 7 (Fig. 2a), because of cancellations between

461    the positive impacts in the extratropics and negative impacts in the tropics seen in Fig. 8d.

462

463    Fig. 8 shows the Day-7 errors of the Control $Z_{200hPa}$ forecasts and how they differ from the

464    errors in the other three forecast sets. The impact of the ENRR observations is relatively

465    small in the tropics and mixed in the extratropics (Fig. 8b). Using the Hybrid versus EnKF

466    initialization yields a similarly mixed impact in the extratropics, and a small but clear

467    degradation in the tropics (Fig. 8c). Using the SPs in the forecast model yields a more

468    consistent beneficial impact in the extratropics, but also a much stronger degradation of the

469    $Z_{200hPa}$ forecasts in the tropics (Fig. 8d). Interestingly, this degradation occurs not just over

470 the tropical convective areas but also over clear-sky areas in the descending branch of the

471 Pacific Walker cell, in which one would expect scant local SPPT tendencies of radiative

472 heating.

473

474 *3.2 Forecast biases*

475

476 Thus far, we have considered GFS forecast sensitivities to the ENRR observations, data

477 assimilation method, and stochastic parameterizations in terms of RMSE measures of

478 ensemble-mean forecasts. It is also relevant to consider how these three factors affect the

479 mean forecast drift, i.e., the systematic bias at each forecast lead time of the ensemble-

480 mean forecasts averaged over all 100 forecast cases. Fig. 9a shows the biases of the Day-7

481 $Z_{200hPa}$ Control forecasts. Note that unlike the RMSEs, which are positive at all locations,

482 the biases can be positive or negative. Some prominent features in Fig. 9a, such as the

483 positive biases over North America, East Asia, Europe, and the tropics, and the negative

484 biases over the northwest Pacific, northeast Pacific, and northeastern U.S., appear early in

485 the forecasts and are evident throughout the 7-day forecasts (not shown).

486

487 The other panels of Fig. 9 show the systematic differences of the ensemble-mean $Z_{200hPa}$

488 Control forecasts from the ensemble-mean forecasts in the other three forecast sets. They

489 may also be interpreted as the impacts of the ENRR observations (Fig. 9b), Hybrid vs.

490 EnKF initial conditions (Fig. 9c), and stochastic parameterizations (Fig.9d) on the Control

491 forecast biases. The impact of the ENRR observations is apparently to intensify El Niño-

492 related features in the Day-7 $Z_{200hPa}$ forecasts: a low along the Canadian West Coast and

21

493    U.S. Pacific Northwest, a high to the west of the Great Lakes, and another high off the

494    Northeast U.S. coast. Although this impact is not statistically significant (see Fig. 11), it is

495    not inconsistent with the response to an anomalous equatorial heat source located east of

496    the dateline (Ting and Sardeshmukh, 1993) during El Niño events. The impact is likely due

497    to a slight but systematic strengthening of the tropical upper tropospheric convective

498    outflow in the Control analyses using the ENRR wind observations (Slivinski *et al.,* 2018)

499    and consequently the Rossby wave source associated with the El Niño-related tropical

500    heating (Sardeshmukh and Hoskins, 1988).

501

502    The impacts of the DA method and SPs on the ensemble-mean $Z_{200hPa}$ Control forecast

503    biases in Fig. 9c are much larger than those of the ENRR observations. Both increase the

504    ensemble-mean $Z_{200hPa}$ in the tropics and subtropics, and contribute to the positive bias of

505    the Control $Z_{200hPa}$ forecasts over these large regions covering more than 50% of the globe.

506    The negative impact of the SPs is especially strong and remarkable, considering that the

507    Control forecast biases are determined with respect to analyses which include SPs in the

508    DA model. This degradation is evident as early as Day 1 in the tropics, spreading thereafter

509    to higher latitudes (not shown). A preliminary diagnosis suggests that it originates largely

510    from a nonlinear response of convection to the SHUM perturbations, which are themselves

511    unbiased (i.e., have zero mean). The impact of using the Hybrid versus EnKF initial

512    conditions is more mixed in this regard, with alternating positive and negative impacts

513    along the Northern Hemisphere extratropical jet stream waveguide.

514
515    Fig. 10 shows similar bias results for $\omega_{500hPa}$ in an identical format to Fig. 9. To focus on

516    larger-scale features, we smoothed the fields using the spatial filter described in

517     Sardeshmukh and Hoskins (1984), retaining scales corresponding to total spherical

518     wavenumbers 15 and lower. Even so, the fields remain noisy, but with a clear suggestion

519     of a wave-train of alternating positive and negative Control forecast biases along the

520     extratropical jet stream waveguide. This wave-train is also evident in the other panels of

521     Fig. 10 showing the bias impacts of the ENRR observations, using the different DA

522     methods, and SPs. Inspection of maps similar to those in Fig 10, but for earlier forecast

523     lead times (not shown) reveal this wave-train to be a remarkably robust eastward

524     propagating feature of the Control forecast biases and bias impacts. Note that the bias

525     impacts of the ENRR observations and DA method stem only from differences in the

526     forecast initial conditions, whereas the bias impacts of the SPs result from changes to the

527     forecast model. The impact of the ENRR observations occurs initially as westward

528     propagating tropical waves that provide perturbations in sensitive regions for exciting the

529     mid-latitude wave-train. The impact of the DA method is stronger than that of the ENRR

530     observations, because the systematic differences between the Hybrid and EnKF DA (see

531     Section 2.2 for the DA method description) are larger than those between the Control and

532     Denial analyses. The impact of the SPs is different in being much stronger in the tropics,

533     and with a slower emergence of the midlatitude wave-train. This slower emergence is not

534     unexpected, since the SPs provide new perturbations throughout the forecast and prevent

535     the occurrence of coherent optimal conditions for exciting the wave-train.

536

537     The bias results in Figs. 9 and 10 have a dynamically meaningful interpretation in at least

538     the extratropics. The extratropical wave-train is highly reminiscent of the most unstable (or

539     least damped) perturbation eigenmode of the extratropical circulation investigated by Hall

540    and Sardeshmukh (1998). On the other hand, since almost any perturbation can set off such

541    an unstable eigenmode with arbitrary amplitude and phase, its appearance in our bias

542    impact statistics makes it harder to distinguish among our estimated bias sensitivities to the

543    ENRR observations, DA methods, and SPs and to establish their statistical significance.

544

545    Indeed, it turns out that the bias impacts in Figs. 9b, 9d, 10b, and 10d are generally not

546    statistically significant in the extratropics. This is shown in Fig.11 for $Z_{200hPa}$ and $\omega_{500hPa}$ in

547    terms of the Student's $t$ scores of the estimated bias differences. The details of these

548    significance calculations are provided in Appendix A. The impact of the ENRR

549    observations on the Day-7 forecast biases is insignificant almost everywhere on the globe.

550    While the bias impacts of the hybrid DA are significant in some scattered areas in the

551    extratropics, the bias impacts of the SPs are generally insignificant outside the tropics.

552    However, they are both highly significant in the tropics.

553

554    **4. Summary and concluding remarks**

555

556    In our forecast sensitivity experiments, the impact of the ENRR observations on the

557    RMSEs of the ensemble-mean forecasts was relatively large at short forecast lead times

558    (about 1 day) whereas the impact of using the Hybrid versus EnKF DA method lasted

559    throughout the forecast period (7 days). This was evident for all the six variables examined

560    ($Z_{200hPa}$, $\xi_{200hPa}$, $\omega_{500hPa}$, PWAT, $T_{2m}$, and AP12HR). The impact of the SPs was to reduce

561    the RMSEs of the ensemble-mean forecasts of all these variables, except $Z_{200hPa}$ in the

562    tropics. Furthermore, this generally positive impact of the SPs grew with forecast lead time.

24

563 The mechanisms through which SPs reduce the errors of ensemble-mean forecasts are

564 worthy of a more detailed investigation, which will be reported elsewhere.

565

566 To varying degrees, the ENRR observations, DA method, and SPs also impacted the

567 forecast biases. The impact of the ENRR observations was the weakest and not statistically

568 significant over most of the globe. The impacts of the DA method were statistically

569 significant in the tropics and in some scattered areas in the extratropics, while the impacts

570 of the SPs were highly significant and generally concentrated in the tropics. The impact of

571 the SPs was stronger than that of the DA method.

572

573 In summary, our goal in this study was to assess the relative sensitivities of global GFS

574 forecasts during late winter/early spring 2016 to the additional ENRR observations

575 collected during the period, to the DA method used to provide the forecast initial

576 conditions, and to the use of SPs in the forecast model. Of these, the sensitivity to the

577 additional ENRR observations, in terms of both biases and RMSEs of the ensemble-mean

578 forecasts, was found to be the weakest, and that to the SPs the strongest, in the 100 forecast

579 cases investigated. The generally positive impact of the SPs on the ensemble-mean

580 forecasts, and also their strongly negative impact on the tropical $Z_{200hPa}$ forecasts, are

581 noteworthy and require further investigation.

582

583 Modern forecast systems are sensitive to many system elements, and our investigation was

584 certainly not meant to be exhaustive in this regard. Rather, our goal was to provide a sense

585 of the relative sensitivities to the three principal types of development activities that are of

586    current interest at major forecasting centers: collecting and using more observations,

587    developing better data assimilation methods, and improving the forecast models.

588

589    As far as we are aware, our study is the first to perform sensitivity tests of sufficient size

590    simultaneously on all the three basic elements of an ensemble forecast system to produce

591    statistically meaningful results for intercomparisons. Even so, the generalizability of our

592    results is limited. For example, our result that the additional ENRR observations did not

593    significantly improve the GFS forecast skill does not necessarily imply that additional

594    observations will have little impact on forecast skill in general. It is well known that short-

595    range forecasts of high-impact weather events benefit from additional in-situ observations

596    (e.g., NOAA Sensing Hazards with Operational Unmanned Technology project). Clearly,

597    the impact of additional observations depends on their relative augmentation of pre-

598    existing observational networks as well as on the types and scales of target weather events.

599

600    Our investigation of forecast sensitivities to DA methods was likewise not exhaustive, as

601    we only compared one implementation of the Hybrid 4D-EnVar to one implementation of

602    the EnKF. We might have obtained different results by using, for example, a different

603    relative weighting of the static and time-varying background error covariances in the cost

604    function of the Hybrid filter (see Section 2.2), or by further optimizing the EnKF

605    parameters. Adopting another distinct DA method might also have yielded different results

606    in this regard.

607

608    Perhaps the strongest robust conclusion of our study is that utilizing even simple types of

609    stochastic parameterizations (SPs) in the forecast model can have stronger and generally

610    beneficial impacts on forecast skill than tinkering with other elements of current forecast

611    systems. However, even this conclusion comes with a caveat that we did not exhaustively

612    investigate forecast sensitivities to other types of stochastic parameterizations.

613    Nonetheless, the main positive result from including stochastic parameterizations seems

614    clear.

615

616    We end with a cautionary note that state-of-the-art forecast systems are now sufficiently

617    advanced and finely tuned that establishing the impacts of forecast system changes on

618    forecast skill with statistical confidence requires careful numerical experimentation with

619    large forecast ensemble sizes. The fact that even with 8,000 (= 100 forecast cases × 80

620    ensemble members for each case) 7-day forecasts in each of our four forecast sets (Control,

621    Denial, EnKFonly, noSP), the apparently large impacts on the extratropical biases in Figs.

622    9 and 10 turned out to be not statistically significant in the Northern Hemisphere upper

623    tropospheric waveguide provides a sobering reminder in this regard.

624
625
626

627 **Appendix A**

628

629 To test the statistical significance of the forecast differences in Figs. 9 and 10, we used the

630 Student's $t$ test (see Fig. 11 for their $t$ values), assuming that the variables are normally

631 distributed. Specifically, at each gridpoint we computed the t-statistic

632

633
$$t = \frac{\overline{x_1} - \overline{x_2}}{\left(\frac{\sigma_1^2}{n_1^*} + \frac{\sigma_2^2}{n_2^*}\right)^{1/2}},$$

634

635 where $\overline{x_1}$ and $\overline{x_2}$ are the means of 8,000 (= 100 forecast cases × 80 ensemble

636 members/forecast case) valid forecast values from two different forecast sets, $\sigma_1^2$ and $\sigma_2^2$

637 are the variances of the 8,000 values in the two forecast sets, and $n_1^*$ and $n_2^*$ are the

638 estimated degrees of freedom (DOF) or effective sample sizes.

639

640 The DOF are smaller than 8,000, because the $I$=80 ensemble values for each forecast case

641 are not truly independent, and the $J$=100 forecast cases also have some serial dependence

642 since they are initialized only 12 hours apart. We estimated the DOF as follows. Let $z_{ij}$ be

643 the forecast from the $i$-th ensemble member and $j$-th forecast case. One can group $z_{ij}$ by

644 ensemble member or case number so that

645
$$\{z_{ij}\} = \{x_i\} = \{y_j\},$$

646 where $x_i$ is the case series of the $i$-th ensemble member, and $y_j$ is the ensemble member

647 series of the $j$-th case. One can think of $x$ and $y$ as the row and column vectors, respectively,

648 of the matrix $z$. Then one can write

649
$$Var\left(\sum_{i=1}^{I} x_i\right) = \sum_{i=1}^{I} Var(x_i) + \sum_{i \neq k} Cov(x_i, x_k).$$

650 This variance has two contributions: 1) the sum of the variances of the individual ensemble

651 members, and 2) the sum of covariances between any two distinct ensemble members. This

652 may also be expressed as

653
$$Var\left(\sum_{i=1}^{I} x_i\right) = Var(IM_x) = I^2 Var(M_x),$$

654 where $M_x = \frac{1}{I}\sum_{i=1}^{I} x_i$ is the case series of the ensemble means. By combining the two

655 equations above, and assuming that all the $z_{ij}$ are independent and identically distributed

656 (i.i.d.), the variance of the ensemble-mean forecasts, from the Law of Large Numbers

657 (LLN), may be written as

658
$$Var(M_x) = \frac{\sum_{i=1}^{I} Var(x_i) + \sum_{i \neq k} Cov(x_i, x_k)}{I^2} = \frac{Var(z_{ij})}{I}.$$

659 However, the $z_{ij}$ are not independent, because of the non-zero covariance between any two

660 distinct ensemble members ( $\sum_{i \neq k} Cov(x_i, x_k) \neq 0$ ). If positive, this covariance makes the

661 ratio

662
$$r_x = \frac{[\sum_{i=1}^{I} Var(x_i)]/I^2}{Var(z_{ij})/I} = \frac{[\sum_{i=1}^{I} Var(x_i)]/I}{Var(z_{ij})}$$

663 less than 1. The DOF in the ensemble member dimension (i.e. the effective ensemble size)

664 is then not $I$ but $I \times r_x$ since

665
$$Var(M_x) = \frac{Var(z_{ij})}{I \times r_x}$$

666 agrees with the LLN. Similarly, the ratio

29

$$r_y = \frac{[\sum_{j=1}^{J} Var(y_j)]/J}{Var(z_{ij})},$$

667

668     provides an estimate of the dependency among the different forecast cases. The overall

669     DOF is then $(I \times r_x) \times (J \times r_y) = 8,000 \times r_x \times r_y$.

670

671     Fig. A1 shows maps of $Var(z_{ij})$, $\sum_{i=1}^{I} Var(x_i)/I$, and $\sum_{j=1}^{J} Var(y_j)/J$ for the spatially

672     smoothed Day-7 $\omega_{500hPa}$ Control forecasts. If all the forecasts were independent, the three

673     maps would be identical. The results show that $r_x$ is a nearly uniform 0.8 everywhere over

674     the globe, while $r_y$ is generally between 0.3 and 0.9. The overall DOF $\omega_{500hPa}$ in the Control

675     forecasts is thus generally between 2,500 and 5,000 for our samples of size 8,000.

676

677     The variance of the ensemble members is clearly representative of the total variance over

678     the whole globe, except that the magnitude is smaller because the ensemble members are

679     still not completely independent by Day 7 (Fig. A1 middle). On the other hand, the case

680     variance is not as representative, and the variance ratios are especially noisy in tropical

681     areas (Fig. A1 bottom).

682

683 **Appendix B**

684

685 The RMSEs in this study were defined as the square root of case-mean and area-mean

686 squared errors of ensemble-mean forecasts with respect to *truth* (see Section 2.2 and 3.1).

687 Because parametric forms of the probability distributions of RMSEs or RMSE differences

688 (hereafter ΔRMSEs) are generally unknown, we used a Bootstrap method to estimate the

689 sampling distributions of ΔRMSEs to assess the significance of ΔRMSEs obtained between

690 any two forecast sets. To this end we combined the 100 forecast cases in each set into a

691 pool of 200 cases. By randomly drawing with replacement from the pool, two new separate

692 100-case samples were made, and their ΔRMSE was calculated. Repeating this process

693 1000 times yielded 1000 values of ΔRMSE for estimating the sampling ΔRMSE

694 distribution. The statistical significance of the actual ΔRMSE was then judged by whether

695 it ranked above the 97.5 percentile or below the 2.5 percentile of this constructed ΔRMSE

696 distribution for a two-sided statistical test. This process was repeated for each 12-hourly

697 forecast lead time up to 168 hours (7 days).

698

699 Figs. B1-B3 show global and regional ΔRMSEs between the Control and the other three

700 (Denial, EnKFonly, and noSP) forecasts, corresponding to Figs. 2, 3, and 7 respectively,

701 as well as the 97.5% and 2.5% percentiles of the ΔRMSEs of their respective sampling

702 distributions. Fig. B1 shows that the Control global RMSEs are significantly smaller than

703 the Denial only for $\xi_{200hPa}$ and $\omega_{500hPa}$ in the first 24 hours of the forecasts, confirming that

704 the ENRR observations only benefit short-term forecasts at smaller spatial scales. The

705 general pattern in Figs. B1-B3 shows that Hybrid initialization (Control forecasts)

31

706     significantly lowers the RMSEs in the first few days, compared to EnKF initialization

707     (EnKFonly forecasts). Also, using SPs (Control forecasts) significantly lowers the RMSEs

708     in the later part of the 7-day forecast evolution, compared to not using SPs (noSP forecasts).

709     The exceptions are AP12HR $\Delta$RMSEs between $60^{\circ}$S and $60^{\circ}$N (Fig. B1f), which do not

710     ever exceed the confidence interval, and $Z_{200hPa}$ $\Delta$RMSE$_{Control-noSP}$ (Fig. B3d), which shows

711     larger errors when using SPs especially in the tropics.

712

713

714

724 **References**

725

726 Anderson, J. L., and N. Collins (2007), Scalable implementations of ensemble filter 247

727     algorithms for data assimilation, Journal of Atmospheric and Oceanic Technology,

728     248 24 (8), 1452–1463, doi:10.1175/JTECH2049.1.

729 Ashouri H., K. Hsu, S. Sorooshian, D. K. Braithwaite, K. R. Knapp, L. D. Cecil, B. R.

730     Nelson, and O. P. Prat, 2015: PERSIANN-CDR: Daily Precipitation Climate Data

731     Record from Multi-Satellite Observations for Hydrological and Climate Studies. Bull.

732     Amer. Meteor. Soc., doi: 10.1175/BAMS-D-13-00068.1.

733 Bauer, P., A. Thorpe, and G. Brunet, 2015: The quiet revolution of numerical weather

734     prediction. *Nature*, **525**, 47–55, doi:10.1038/nature14956

735 Berner, J., G. Shutts, M. Leutbecher, and T. Palmer, 2009: A spectral stochastic kinetic

736     energy backscatter scheme and its impact on flow- dependent predictability in the

737     ECMWF ensemble prediction system. J. Atmos. Sci., 66, 603–626,

738     doi:10.1175/2008JAS2677.1.

739 Bloom, S. C., L. L. Takacs, A. M. da Silva, and D. Ledvina, 1996: Data Assimilation

740     Using Incremental Analysis Updates. *Mon. Weather Rev.*, 124, 1256-1271.

741 Buehner, M., J. Morneau, and C. Charette, 2013: Four-dimensional ensemble-variational

742     data assimilation for global deterministic weather prediction. *Nonlinear Processes*

743     *Geophys.*, **20**, 669–682.

744 Campbell, W.F., C.H. Bishop, and D. Hodyss, 2010: Vertical Covariance Localization for

745     Satellite Radiances in Ensemble Kalman Filters. *Mon. Wea. Rev.,* **138**, 282–290,

746     https://doi.org/10.1175/2009MWR3017.1.

747    Courtier, P., J.-N. Thépaut, and A. Hollingsworth, 1994: A strategy for operational

748       implementation of 4D-Var, using an incremental approach. *Quarterly Journal of the*

749       *Royal Meteorological Society*, Vol. **120**, Issue 519, Part B, 1367-1387.

750    Dole, R.M., and Co-Authors, 2018: Advancing Science and Services during the 2015-16

751       El Niño: The NOAA El Niño Rapid Response Field Campaign. *Bulletin of the*

752       *American Meteorological Society*, Vol. **99**(5), 975-1002, DOI:10.1175/BAMS-D-16-

753       0219.1.

754    Efron, B., 1982: *The Jackknife, the Bootstrap, and Other Resampling Plans*. Society for

755       Industrial and Applied Mathematics, Philadelphia, PA, 92 pp.

756    Efron, B. and R. J. Tibshirani, 1993: *An Introduction to the Bootstrap.* Chapman & Hall,

757       New York, 456 pp.

758    Environmental Modeling Center, 2003: The GFS Atmospheric Model. NCEP Office Note

759       442, Global Climate and Weather Modeling Branch, EMC, Camp Springs, Maryland.

760    Evensen, G., 2003: The Ensemble Kalman Filter: theoretical formulation and practical

761       implementation. *Ocean Dynamics*, 53:343-367, doi:10.1007/s10236-003-0036-9.

762    Gaspari, G., and S. E. Cohn, 1999: Construction of correlation functions in two and three

763       dimensions. *Quart. J. Roy. Meteor. Soc.*, **125**, 723-757, doi:10.1002/qj.49712555417.

764    Hall, N.M.J., and P. D. Sardeshmukh, 1998: Is the Time-Mean Northern Hemisphere

765       Flow Baroclinically Unstable? *J. Atmos. Sci.*, Vol. **55**, 41-56.

766    Houtekamer, P. L., and H. L. Mitchell, 1998: Data assimilation using an ensemble

767       Kalman filter technique. *Mon. Wea. Rev.*, **126**, 796-811, doi:10.1175/1520-

768       0493(1998)126<0796:DAUAEK>2.0.CO;2.

769    Huffman, D., Bolvin, D. Braithwaite, K. Hsu, R. Joyce, P. Xie, 2014: Integrated Multi-

770    satellite Retrievals for GPM (IMERG), version 4.4. NASA's Precipitation Processing

771    Center, accessed 31 March, 2015, ftp://arthurhou.pps.eosdis.nasa.gov/gpmdata/

772 Kleist, D.T. and K. Ide, 2015: An OSSE-Based Evaluation of Hybrid Variational–

773    Ensemble Data Assimilation for the NCEP GFS. Part II: 4DEnVar and Hybrid

774    Variants. *Mon. Wea. Rev.,* **143**, 452–470, https://doi.org/10.1175/MWR-D-13-

775    00350.1

776 Lei, L., and J. S. Whitkaer, 2015: Model space localization is not always better than

777    observation space localization for assimilation of satellite radiances. *Mon. Wea. Rev.*,

778    **143**, 3948-3955, doi: 10.1175/MWR-D-14-00413.1.

779 Lei, L., and J. S. Whitaker, 2016: A Four-Dimensional Incremental Analysis Update for

780    the Ensemble Kalman Filter. *Mon. Wea. Rev.*, **144**, 2605-2621, DOI: 10.1175/MWR-

781    D-15-0246.1.

782 Lei, L., J. S. Whitaker and C. Bishop, 2018: Improving assimilation of radiance

783    observations by implementing model space localization in an ensemble Kalman filter.

784    *Journal of Advances in Modeling Earth Systems*, **10**, 12, 3221-3232. DOI:

785    10.1029/2018MS001468.

786 Leutbecher, M., and Co-authors, 2017: Stochastic representations of model uncertainties

787    at ECMWF: state of art and future vision. *Quarterly J. Royal Met. Soc*., **143**, 707,

788    2315-2339. DOI: 10.1002/qj.3094

789 Lewis, J. M. and J. C. Derber, 1985: The use of adjoint equations to solve a variational

790    adjustment problem with advective constraints. *Tellus*, **37A**, 4, 309-322.

791    Palmer, T. N., R. Buizza, F. Doblas-Reyes, T. Jung, M. Leutbecher, G. J. Shutts, M.

792        Steinheimer, and A.Weisheimer, 2009: Stochastic parametrization and model

793        uncertainty. ECMWF Tech. Memo. 598, 42 pp.

794    Sardeshmukh, P. D., G. P. Compo, M. C. Penland, 2015: Need for Caution in Interpreting

795        Extreme Weather Statistics. *J. Climate*, Vol. **28**, 23, 9166-9187.

796    Sardeshmukh, P. D., 2005: Issues in stochastic parameterization. *Proc. Workshop on*

797        *Representation of Sub-Grid Processes Using Stochastic-Dynamic Models*, Shinfield

798        Park, Reading, ECMWF, 5–12.

799    Sardeshmukh, P. D. and Hoskins, B. J., 1984: Spatial Smoothing on the Sphere. Mon.

800        Wea. Rev., Vol. **112**, 2524-2529.

801    Sardeshmukh, P.D., and B. J. Hoskins, 1988: The Generation of Global Rotational Flow

802        by Steady Idealized Tropical Divergence. *J. Atmos. Sci.*, Vol. **45**, No. 7, 1228-1251.

803    Shutts, G., M. Leutbecher, A. Weisheimer, T. Stockdale, L. Isaksen, and M. Bonavita,

804        2011: Representing model uncertainty: stochastic parameterizations at ECMWF,

805        *ECMWF Newsletter*, **129**, 19-24.

806    Sorooshian, S., K. Hsu,D. Braithwaite,H. Ashouri, and NOAA CDR Program (2014):

807        NOAA Climate Data Record (CDR) of Precipitation Estimation from Remotely

808        Sensed Information using Artificial Neural Networks (PERSIANN-CDR), Version 1

809        Revision 1. NOAA National Centers for Environmental Information.

810        doi:10.7289/V51V5BWQ [access date: April 27th, 2017]

811    Slivinski, L.C., G.P. Compo, J.S. Whitaker, P.D. Sardeshmukh, J.H.A. Wang, K.

812        Friedman, and C. McColl, 2018: What is the impact of additional tropical observations

813        on a modern data assimilation system? *Mon. Wea. Rev.*, in review.

814    Takacs, L. L., M. J. Suárez, and R. Todling, 2018: The Stability of Incremental Analysis

815    Update. *Mon. Wea. Rev.,* **146**, 3259–3275, https://doi.org/10.1175/MWR-D-18-

816    0117.1.

817    Ting, M., and P.D. Sardeshmukh, 1993: Factors Determining the Extratropical Response

818    to Equatorial Diabatic Heating Anomalies. *J. Atmos. Sci.*, Vol. **50**., No. 6, 907-918.

819    Tompkins, A. M., and J. Berner, 2008: A stochastic convective approach to account for

820    model uncertainty due to unresolved humidity variability. *J. Geophys. Res.*, Vol. **113**,

821    D18101, doi:10.1029/2007JD009284.

822    Weaver, A., and P. Courtier, 2001: Correlation modelling on the sphere using a

823    generalized diffusion equation. *Quart. J. Roy. Meteorol. Soc.*, **127**, 1815–1846.

824    Whitaker, J. S., and T. M. Hamill (2002), Ensemble data assimilation without 289

825    perturbed observations, *Mon. Weather Rev.*, **130** (7), 1913–1924, doi:10.1175/1520-

826    0493(2002)130<1913:EDAWPO>2.0.CO;2.

827

828    **Figure Captions**

829

830    **Figure 1**. Schematic depiction of the 7-day forecasts generated and verification period

831    used. Each arrow represents one forecast case, and only the portion in the verification

832    period is evaluated for this study. Note that there are 80 members in the ensemble

833    forecast for each forecast case.

834

835    **Figure 2**. Global RMSEs of the Control (solid gray), Denial (dashed blue), EnKFonly

836    (dotted green) and noSP forecasts (dash-dot red), determined with respect to the Control

837    analyses for (a) 200hPa heights ($Z_{200hPa}$), (b) 200hPa vorticity ($\xi_{200hPa}$), (c) 500hPa vertical

838    p-velocity ($\omega_{500hPa}$), (d) precipitable water (PWAT), and (e) 2-meter air temperature ($T_{2m}$).

839    (f) The RMSE of 12-hr accumulated precipitation (AP12HR) averaged in the 20ºS to 20ºN

840    domain (thin upper curves) and the 60ºS to 60ºN domain (thick lower curves), determined

841    with respect to NASA GPM observational dataset. Note the ordinate for the precipitation

842    RMSE starts at 6 mm.

843

844    **Figure 3**. Domain $\omega_{500hPa}$ RMSEs of the Control, Denial, EnKFonly and noSP forecasts

845    with respect to the Control analyses in the (a) Northern Hemisphere (20ºN-90ºN), (b)

846    Southern Hemisphere (20ºS-90ºS), (c) Tropics (20ºS-20ºN), and (d) Contiguous United

847    States (CONUS; 125ºW-66ºW, 24ºN-50ºN).

848

849    **Figure 4**. (a) The $\omega_{500hPa}$ RMSEs of the Day-7 Control forecasts; (b) The differences of the

850    $\omega_{500hPa}$ RMSEs between the Day-7 Control and Denial forecasts; (c) Similar to (b), but

851    between the Control and EnKFonly forecasts; (d) Similar to (b), but between the Control

852    and noSP forecasts.

853

854

855    **Figure 5**. (a) The AP12HR RMSEs of the Control forecasts with respect to independent

856    NASA GPM product at the end of Day 7; (b) The AP12HR RMSE differences between the

857    Control and Denial forecasts at the end of Day 7; (c) Similar to (b), but between the Control

858    and EnKFonly forecasts; (d) Similar to (b), but between the Control and noSP forecasts.

859    The valid geographic domain is between 60ºS and 60ºN. If there exist only missing values

860    in a grid box ($0.5°\times0.5°$) at any moment during the verification period, that box is painted

861    gray in (b)-(d).

862

863    **Figure 6**. As in Fig. 4, except for $T_{2m}$.

864

865    **Figure 7**. As in Fig. 3, but for $Z_{200hPa}$.

866

867    **Figure 8**. (a) The $Z_{200hPa}$ RMSEs of the Control forecasts at the end of Day 7; (b) The

868    $Z_{200hPa}$ RMSE differences between the Control and Denial forecasts at the end of Day 7;

869    (c) Similar to (b), but between the Control and EnKFonly forecasts; (d) Similar to (b), but

870    between the Control and noSP forecasts.

871

872    **Figure 9**. (a) Bias of case-mean ensemble-mean Day-7 $Z_{200hPa}$ Control forecasts with

873    respect to the Control analyses; (b) Difference of case-mean ensemble-mean Control and

874    Denial forecasts; (c) Difference of case-mean ensemble-mean Control and EnKFonly

875    forecasts; (d) Difference of case-mean ensemble-mean Control and noSP forecasts. Note

876    that the contour interval in panel (a) is 4.5 times that in the other panels.

877

878    **Figure 10**. As in Fig. 9 except for $\omega_{500hPa}$. Note that the contour interval in panel (a) is five

879    times that in the other panels. The additional thick black curves in the extratropical

880    Northern Hemisphere enclose the region of 200hPa mean zonal winds stronger than 30m/s

881    in the Control analysis, which is a good proxy of the extratropical baroclinic waveguide.

882

883 **Figure 11**. Left panels: The Student's t scores for the Day-7 $Z_{200hPa}$ bias differences

884 between (top) the Control and Denial forecasts, (middle) the Control and EnKFonly

885 forecasts, and (bottom) the Control and noSP forecasts. A value of ±1.645 is 10%

886 significant in two-tailed test, ±1.96 is 5% significant, and ±2.576 is 1% significant. Right

887 Panels: Similar to left panels but for $\omega_{500hPa}$ fields. The thick black 30m/s contour of the

888 200hPa zonal winds in the Northern Hemisphere shows the approximate location of the

889 upper tropospheric jet stream waveguide, as in Fig. 10.

890

891 **Figure A1.** (top) The total variance of the spatially smoothed Day-7 $\omega_{500hPa}$ Control

892 forecasts; (middle) the sum of the variances within the individual ensemble members

893 across the cases, divided by group size 100; (bottom) the sum of the variances within the

894 individual cases across the ensemble members, divided by group size 80 (color shaded),

895 and the ratio of the values of the sum of the variances to the total variance (contours). The

896 contour interval in the bottom panel is 0.1, and the 1 contour is thickened. The variance

897 ratio in the middle panel is ~0.79 almost uniformly over the globe and hence no contour is

898 plotted. Note that if all the forecasts were independent, the values in the middle and bottom

899 panels would be equal to those in the top panel.

900

901 **Figure B1.** Global RMSE differences between the Control and Denial forecasts (solid

902 blue), between the Control and EnKFonly forecasts (solid green), and between the Control

903 and noSP forecasts (solid red) for (a) 200hPa geopotential heights ($Z_{200hPa}$), (b) 200hPa

904 vorticity ($\xi_{200hPa}$), (c) 500hPa vertical p-velocity ($\omega_{500hPa}$), (d) precipitable water (PWAT),

905 and (e) 2-meter air temperature ($T_{2m}$). (f) Similar to panel (a)-(d), except for 12-hr

906     accumulated precipitation (AP12HR) RMSE differences averaged in the 20°S to 20°N (thin

907     curves) and the 60°S to 60°N (thick curves) latitude domains. The dotted lines represent

908     the 2.5% (below ΔRMSE=0) and 97.5% (above ΔRMSE=0) of the constructed

909     distributions for Control-Denial (blue), Control-EnKFonly (green), and Control-noSP

910     (red), derived from the Bootstrap method.

911

912     **Figure B2.** Similar to Fig. B1, except for $\omega_{500hPa}$ in (a) Northern Hemisphere, (b) Southern

913     Hemisphere, (c) Tropics, and (d) Contiguous United States. See Fig. 3 and context for

914     domain definitions.

915

916     **Figure B3.** Similar to Fig. B2, except for $Z_{200hPa}$.

917
918

919

920
921   Table 1: List of forecast ensembles generated
922
923

| Label | Initial Condition | Data Assimilation Method | Forecast Model |
|---|---|---|---|
| Control | Includes ENRR obs | Hybrid | Includes Stochastic Physics |
| Denial | Excludes ENRR obs | Hybrid | Includes Stochastic Physics |
| EnKFonly | Includes ENRR obs | EnKF | Includes Stochastic Physics |
| noSP | Includes ENRR obs | Hybrid | No Stochastic Physics |

924
925
926
927

928
929
930
931
932
933
934
935
936
937

**Figure 1**. Schematic depiction of the 7-day forecasts generated and verification period used. Each arrow represents one forecast case, and only the portion in the verification period is evaluated for this study. Note that there are 80 members in the ensemble forecast for each forecast case.

951
952
953 **Figure 2**. Global RMSEs of the Control (solid gray), Denial (dashed blue), EnKFonly
954 (dotted green) and noSP forecasts (dash-dot red), determined with respect to the Control
955 analyses for global (a) 200hPa heights ($Z_{200hPa}$),(b) 200hPa vorticity ($\xi_{200hPa}$), (c) 500hPa
956 vertical p-velocity ($\omega_{500hPa}$), (d) precipitable water (PWAT), and (e) 2-meter air
957 temperature ($T_{2m}$). (f) The RMSE of 12-hr accumulated precipitation totals in the 20$^o$S to

45

958    20ºN domain (thin upper curves) and the 60ºS to 60ºN domain (thick lower curves),
959    determined with respect to NASA GPM observational dataset. Note the ordinate for the
960    precipitation RMSE starts at 6 mm.

961
962
963
964



965
966    **Figure 3**. Domain $\omega_{500hPa}$ RMSEs of the Control, Denial, EnKFonly and noSP forecasts
967    with respect to the Control analyses in the (a) Northern Hemisphere (20ºN-90ºN), (b)
968    Southern Hemisphere (20ºS-90ºS), (c) Tropics (20ºS-20ºN) and (d) Contiguous United
969    States (CONUS; 125ºW-66ºW, 24ºN-50ºN).

970
971
972
973

974

975



976

977

978 **Figure 4**. (a) The $\omega_{500hPa}$ RMSEs of the Day-7 Control forecasts; (b) The differences of the
979 $\omega_{500hPa}$ RMSEs between the Day-7 Control and Denial forecasts; (c) Similar to (b), but
980 between the Control and EnKFonly forecasts; (d) Similar to (b), but between the Control
981 and noSP forecasts.

982

983

984

985

986

991 **Figure 5**. (a) The AP12HR RMSEs of the Control forecasts with respect to independent
992 NASA GPM product at the end of Day 7; (b) The AP12HR RMSE differences between the
993 Control and Denial forecasts at the end of Day 7; (c) Similar to (b), but between the Control
994 and EnKFonly forecasts; (d) Similar to (b), but between the Control and noSP forecasts.
995 The valid geographic domain is between 60ºS and 60ºN. If there exist only missing values
996 in a grid box (0.5º×0.5º) at any moment during the verification period, that box is painted
997 gray in (b)-(d).

1002



1003
1004
1005 **Figure 6**. As in Fig. 4, except for $T_{2m}$.
1006
1007
1008

1009
1010
1011     **Figure 7**. As in Fig. 3, but for $Z_{200hPa}$.
1012
1013
1014
1015
1016
1017

1018



1019
1020
**Figure 8**. (a) The $Z_{200hPa}$ RMSEs of the Control forecasts at the end of Day 7; (b) The $Z_{200hPa}$ RMSE differences between the Control and Denial forecasts at the end of Day 7; (c) Similar to (b), but between the Control and EnKFonly forecasts; (d) Similar to (b), but between the Control and noSP forecasts.

1025
1026
1027
1028

1029



Mean $Z_{200mb}$, fhr=168

(a) Control Forecast Error

(c) Forecast Diff, Control − EnKFonly

−63 −45 −27 −9 9 27 45 63 m

−14 −10 −6 −2 2 6 10 14 m

(b) Forecast Diff, Control − Denial

(d) Forecast Diff, Control − noSP

−14 −10 −6 −2 2 6 10 14 m

−14 −10 −6 −2 2 6 10 14 m

1030
1031
1032
1033  **Figure 9**. (a) Bias of case-mean ensemble-mean Day-7 $Z_{200hPa}$ Control forecasts with
1034  respect to the Control analyses; (b) Difference of case-mean ensemble-mean Control and
1035  Denial forecasts; (c) Difference of case-mean ensemble-mean Control and EnKFonly
1036  forecasts; (d) Difference of case-mean ensemble-mean Control and noSP forecasts. Note
1037  that the contour interval in panel (a) is 4.5 times that in the other panels.
1038
1039
1040
1041

1042



Smoothed Mean $\omega_{500mb}$, fhr=168

(a) Control Forecast Error

(c) Forecast Diff, Control − EnKFonly

(b) Forecast Diff, Control − Denial

(d) Forecast Diff, Control − noSP

1043
1044
1045 **Figure 10**. As in Fig. 9, except for $\omega_{500hPa}$. Note that the contour interval in panel (a) is five
1046 times that in the other panels. The additional thick black curves in the extratropical
1047 Northern Hemisphere enclose the region of 200hPa mean zonal winds stronger than 30m/s
1048 in the Control analysis, which is a good proxy of the extratropical baroclinic waveguide.
1049
1050
1051
1052
1053

**Figure 11**. Left panels: The Student's t scores for the Day-7 $Z_{200hPa}$ bias differences between (top) the Control and Denial forecasts, (middle) the Control and EnKFonly forecasts, and (bottom) the Control and noSP forecasts. A value of $\pm1.645$ is 10% significant in two-tailed test, $\pm1.96$ is 5% significant, and $\pm2.576$ is 1% significant. Right Panels: Similar to left panels but for $\omega_{500hPa}$ fields. The thick black 30m/s contour of the 200hPa zonal winds in the Northern Hemisphere shows the approximate location of the upper tropospheric jet stream waveguide, as in Fig. 10.

1065
1066

Smoothed $\omega_{500mb}$, Control, fhr=168

$\sigma^2_{total}$



$\Sigma\sigma^2_x/100$



$\Sigma\sigma^2_y/80$



0.2    0.4    0.8    1.6    3.2    6.4    12.8    $10^{-3}(Pa/s)^2$

1067
1068
1069  **Figure A1**. (top) The total variance of the spatially smoothed Day7$\omega_{500hPa}$ Control forecasts;
1070  (middle) the sum of the variances within the individual ensemble members across the cases,
1071  divided by group size 100; (bottom) the sum of the variances within the individual cases
1072  across the ensemble members, divided by group size 80 (color shaded), and the ratio of the

1073  values of the sum of the variances to the total variance (contours). The contour interval in
1074  the bottom panel is 0.1, and the 1 contour is thickened. The variance ratio in the middle
1075  panel is ~0.79 almost uniformly over the globe and hence no contour is plotted. Note that
1076  if all the forecasts were independent, the values in the middle and bottom panels would be
1077  equal to those in the top panel.
1078

1079

1080



1081
**Figure B1**. Global RMSE differences between the Control and Denial forecasts (solid
blue), between the Control and EnKFonly forecasts (solid green), and between the Control
and noSP forecasts (solid red) for (a) 200hPa geopotential heights ($Z_{200hPa}$), (b) 200hPa

1085 vorticity ($\xi_{200hPa}$), (c) 500hPa vertical p-velocity ($\omega_{500hPa}$), (d) precipitable water (PWAT),
1086 and (e) 2-meter air temperature ($T_{2m}$). (f) Similar to panel (a)-(d), except for 12-hr
1087 accumulated precipitation (AP12HR) RMSE differences in the 20ºS to 20ºN (thin curves)
1088 and the 60ºS to 60ºN (thick curves) latitude domains. The dotted lines represent the 2.5%
1089 (below ΔRMSE=0) and 97.5% (above ΔRMSE=0) of the constructed distributions for
1090 Control-Denial (blue), Control-EnKFonly (green), and Control-noSP (red), derived from
1091 the Bootstrap method.
1092

1093
1094
1095



1096
1097 **Figure B2.** Similar to Fig. B1, except for $\omega_{500hPa}$ in (a) Northern Hemisphere, (b) Southern
1098 Hemisphere, (c) Tropics, and (d) Contiguous United States. See Fig. 3 and context for
1099 domain definitions.
1100

**Figure B3.** Similar to Fig. B2, except for $Z_{200hPa}$.

1101
1102
1103
1104
1105