

1
2
3 **Bias Correction for Global Ensemble Forecast**
4

5
6
7 Bo Cui¹, Zoltan Toth², Yuejian Zhu³, Dingchen Hou³

8 ¹IMSG at Environmental Modeling Center, NCEP/NWS

9 ²Global Systems Division, ESRL/OAR

10 ³Environmental Modeling Center, NCEP/NWS
11
12
13
14

15 Accepted by Weather and Forecasting
16
17
18
19

20 Corresponding address:

21 Dr. Bo Cui

22 Environmental Modeling Center

23 NCEP/NWS/NOAA

24 5200 Auth Road,

25 Camp Springs, MD 20746

26 301-763-8000 ext 7594

27 E-Mail: Bo.Cui@noaa.gov
28
29
30

Abstract

1
2
3 The main task of this study is to introduce a statistical post-processing
4 algorithm to reduce the bias in the National Centers for Environmental
5 Prediction (NCEP) and Meteorological Service of Canada (MSC) ensemble
6 forecasts before they are merged to form a joint ensemble within the North
7 American Ensemble Forecast System (NAEFS). This statistical post-
8 processing method applies a Kalman filter type algorithm to accumulate the
9 decaying averaging bias and produces bias corrected ensembles for 35
10 variables. NCEP implemented this bias correction technique in 2006. The
11 NAEFS is a joint operational multi-model ensemble forecast system which
12 combines NCEP and MSC ensemble forecasts after bias correction.
13 According to operational statistical verification, both the NCEP and MSC bias
14 corrected ensemble forecast products are enhanced significantly.

15
16 In addition to the operational calibration technique, three other experiments
17 were designed to assess and mitigate ensemble biases on the model grid: a
18 decaying averaging bias calibration method with short samples, a climate
19 mean bias calibration method, and a bias calibration method using dependent
20 data. Preliminary results show that the decaying averaging method works well

1 for the first few days. After removing the decaying averaging bias, the
2 calibrated NCEP operational ensemble has improved probabilistic
3 performance for all measures until day 5. The reforecast ensembles from the
4 Earth System Research Lab Physical Sciences Division with and without the
5 climate mean bias correction were also examined. A comparison between the
6 operational and the bias-corrected reforecast ensembles shows that the climate
7 mean bias correction can add value, especially for week-2 probability
8 forecasts.

9
10
11
12
13
14
15
16
17
18
19
20

1. Introduction

Over the last decade, a global forecast model based global ensemble forecast system (such as the NCEP GEFS) has been found to be useful for medium-range probabilistic forecasting. Ensemble forecasting has been embraced as a practical way of estimating the uncertainty of weather forecasts and of making probabilistic forecasts (Toth and Kalnay 1993, 1997; Molteni et al. 1996; Houtekamer et al. 1996). However, ensemble forecasts still suffer from model and ensemble formation related shortcomings, i.e., imperfect model physics, initial conditions, and boundary conditions for regional ensembles. As Toth et al. (2003) indicated, ensemble forecasts contain systematic errors and these systematic errors remain and cause bias in the 1st and 2nd moments of the ensemble distribution. In order to make a skillful medium-range forecast, it is necessary to run post-processing algorithms to remove these systematic errors before the ensemble forecasts can be used.

A large variety of numerical weather prediction post-processing methods have been proposed and tested by many investigators (Gel 2007; Hacker and Rife 2007; Yussouf and Stensrud 2006, 2007; Cheng and Steenburgh 2007). These techniques are designed for deterministic forecasts and work well for near-

1 surface variables. There have been many attempts to post-process ensemble
2 forecasts to provide reliable probability forecasts. Recent work includes
3 ensemble model output statistics (Gneiting et al. 2005), gene-expression
4 programming (Bakhshaii and Stull 2009), and Bayesian Model Averaging
5 method (Raftery et al. 2005; Krzysztofowicz and Evans 2008). These
6 methods gained broad attention in the ensemble post-processing community.
7 However, these techniques are still in development. The need remains for
8 ensemble statistical post-processing. Further research in this direction is
9 desirable.

10

11 Post-processing of ensemble forecasts is a necessary and important step for
12 the daily operational runs at numerical weather prediction centers. Reliability,
13 accuracy and efficiency are the most important issues for daily operations. In
14 this paper, we first introduce a statistical post-processing algorithm to adjust
15 the 1st moment of ensemble forecasts. This statistical post-processing method
16 applies an adaptive (Kalman filter type, KF) algorithm to accumulate the
17 decaying averaging bias. In statistics the Kalman filter is a mathematical
18 method named after Rudolf E. Kalman (Kalman 1960). It is mainly used to
19 estimate system states that can only be observed indirectly or inaccurately by
20 the system itself and its process is carried out iteratively. The estimates

1 produced by this method tend to be closer to the true values than the original
2 measurements because the weighted average has a better estimated
3 uncertainty than either of the values that went into the weighted average
4 (http://en.wikipedia.org/wiki/Kalman_filter). The basic ideas of KF are
5 straightforward and the KF turns out to be useful for many applications in
6 science, engineering, and economics. We design a specific algorithm based on
7 the KF concept to estimate ensemble forecast errors and we call the bias
8 estimation and correction process the decaying averaging method.

9
10 This method was developed by the National Centers for Environmental
11 Prediction (NCEP) of the US National Weather Service (NWS) and was
12 implemented operationally in 2006 at NCEP to reduce the bias of the NCEP
13 and Canadian Meteorological Centre (CMC) ensemble forecasts. The
14 calibrated NCEP and CMC global ensembles are then merged to form a
15 joint ensemble within the North American Ensemble Forecast System
16 (NAEFS, Zhu et al. 2006). The NAEFS is an operationally joined multi-
17 model ensemble forecast system, which combines the NCEP and CMC
18 ensemble forecasts after bias correction (Zhu and Cui et al. 2007, 2008).

19

1 We also test three different calibration experiments that are designed to assess
2 and mitigate ensemble biases in the 1st (mean) moment of the ensemble on the
3 model grid with respect to analysis fields: a decaying averaging bias
4 calibration method with a short sample, a climate mean bias calibration
5 method, and a bias calibration method using dependent data. The decaying
6 averaging bias calibration method with a short sample is similar to the
7 technique implemented in the NCEP global ensemble forecast system
8 (GEFS), but uses training data for a fixed long period. The second calibration
9 method uses a climate mean bias to do the bias correction, which is supported
10 by a 25 year ensemble reforecast experiment. This reforecast experiment is
11 another ensemble run operationally at NCEP with a frozen analysis/modeling
12 system developed by scientists at ESRL/PSD (formerly CDC). The
13 ESRL/PSD reforecast is run operationally to produce a dataset of historical
14 weather forecasts generated with a fixed numerical model, the 1998 version
15 of NCEP's Global Forecast System (GFS, more information is available
16 online at <http://www.esrl.noaa.gov/psd/forecasts/reforecast/>). A reforecast for each
17 day since 1979 has been made with this GFS version, which is comprised of a
18 15-member ensemble forecast run out to 15 days (Hamill and Whitaker 2004,
19 2006). From the collected 25 years of reforecast data, climate mean forecast
20 errors are diagnosed and the reforecast data is calibrated by removing these

1 errors to increase the reforecast ensemble skill. The design and usages of the
2 second method aim at taking advantage of week-2 forecasts from the long
3 historical reforecast data.

4
5 The purpose of this study is to introduce and compare several statistical post-
6 processing methods to assess and mitigate ensemble biases in the 1st (mean)
7 moment of the ensemble. The decaying averaging method NCEP runs
8 operationally is described in Section 2. How and why a specific decaying
9 parameter is chosen will be discussed. The three different experiments, i.e.,
10 the decaying average method with a short sample, the ESRL/PSD reforecast
11 calibration method and the bias correction method using dependent data
12 (optimal calibrated ensemble) are described in Section 3. Some statistical
13 evaluation methods used in this paper are also reviewed in Section 3. Section
14 4 contains the evaluation of the several calibration methods. Results from the
15 calibrated NCEP/GEFS and CMC/GEFS will be compared with the raw
16 NCEP and CMC ensembles to evaluate the performance of the operational
17 bias correction method. The results from the three calibration experiments,
18 i.e., the decaying averaging bias correction of NCEP/GEFS, the climate mean
19 bias correction of ESRL/PSD reforecast and the optimal calibrated ensembles
20 are also compared. The relative merits of using the current best

1 analysis/modeling system with a small sample, versus the merits of an older
2 and frozen analysis/modeling system that has a long forecast sample for the
3 bias correction will be examined. In general, these two calibration methods
4 for the NCEP/GEFS and ESRL/PSD reforecast ensemble are incompatible
5 because of their different use of model systems. However, these comparisons
6 are not of the superiority of one method over another but to help us illustrate
7 the possibility improving the current operational decaying averaging method
8 to improve week-2 forecasts. Finally, the preliminary conclusion and future
9 plans are summarized in Section 5.

11 **2. The design of the NCEP bias correction method – decaying average**

12
13 NCEP implemented a statistical post-processing algorithm, i.e., the decaying
14 averaging bias correction method to calibrate global ensemble forecasts in
15 2006. The bias correction method is applied to the NCEP and CMC global
16 ensembles. The operational NCEP/GEFS system configuration is described in
17 Toth et al. (2004) and Wei et al. (2008). The 20-member ensembles are
18 produced at T126L28 horizontal and vertical resolution. The perturbations of
19 the initial conditions are from the Ensemble Transform with Rescaling (ETR)
20 technique. The operational CMC ensemble is described in Charron et al.

1 (2010). A single dynamical core, i.e., the Global Environmental Multiscale
2 (GEM) model is used to produce the ensembles. A multi-parameterization
3 approach and stochastic perturbations are used in order to sample model error
4 for the 20 members of the ensembles. Due to the different ensemble
5 configurations, the bias estimation and correction processes of NCEP/GEFS
6 and CMC/GEFS will be adjusted to meet their specific characteristics.

7

8 The operational environment requires that the ensemble post-processing
9 algorithms be relatively applicable and flexible for implementation. The
10 decaying averaging method applies an adaptive algorithm, and its application
11 includes two steps. The first step is to estimate the 1st moment bias with
12 respect to the analysis field, which is called the decaying averaging mean
13 error. The second step is to remove the error from the ensemble forecasts.
14 Both the bias assessment step and bias correction step are carried out
15 separately at each forecast lead time, on each individual grid point and for
16 each initial cycle.

17

18

19

20

1 a) *Bias estimation:*

2

3 The bias $b_{i,j}(t)$ for each lead-time t (a 6-hour interval up to 384 hours), and
4 each grid point (i,j) is defined as the difference between the analysis $a_{i,j}(t)$
5 and forecast $f_{i,j}(t)$ at the same valid time t_0 , on latest available analysis.

6

7
$$b_{i,j}(t) = f_{i,j}(t) - a_{i,j}(t) \quad (1)$$

8

9 b) *Decaying average:*

10

11 The average bias $B_{i,j}(t)$ will be updated by considering the prior period bias
12 $B_{i,j}(t-1)$ and current bias $b_{i,j}(t)$ by using the decaying average with weight
13 coefficient w .

14

15
$$B_{i,j}(t) = (1 - w) B_{i,j}(t-1) + w b_{i,j}(t) \quad (2)$$

16

17 This decaying average bias estimation method is a convenient way to consider
18 the most recent behavior of weather systems. Once initialized, the bias
19 estimate can be updated by considering just the current forecast error with
20 regard to the stored bias fields. The weight factor w controls how much

1 influence to give the most recent data. A w equal to 2% is used for the
2 NCEP/GEFS and CMC/GEFS bias accumulation which includes mainly the
3 past 50-60 days of information (Figure 1). Experiments with a choice of 2%
4 weight and other values (0.25%, 0.5%, 1%, 10%, respectively) have been
5 conducted. The details will be discussed in Section 4.

6
7 *c) Bias correction:*

8
9 The new bias corrected forecast $F_{i,j}(t)$ will be generated by applying the
10 decaying average bias $B_{i,j}(t)$ to current forecasts $f_{i,j}(t)$ at each lead-time and
11 each grid point.

$$13 \quad F_{i,j}(t) = f_{i,j}(t) - B_{i,j}(t) \quad (3)$$

14
15 Steps 1 to 3 allow users to accomplish the bias correction procedure for both
16 NCEP/GEFS and CMC/GEFS. Note that this procedure contains two options.
17 The first is that the NCEP/GEFS and CMC/GEFS can be grouped together
18 before post-processing, and then the bias correction is applied to the joint
19 ensemble. The second option is to apply the bias correction to the NCEP and
20 CMC ensembles separately and then the NCEP/GEFS and CMC/GEFS are

1 grouped together after post-processing. Although the first option is an easy
2 approach, it may not provide the best results since each participating
3 ensemble may have unique biases due to its own ensemble generation
4 configuration. Specific treatments are needed for each participating ensemble.
5 NCEP uses one model and perturbed initial conditions to create its ensemble
6 (Toth et al. 2004; Wei et al. 2008). The model related systematic errors grow
7 with lead time and it is assumed that the forecast errors obtained from the
8 ensemble mean can stand in for the systematic errors. The NCEP/GEFS
9 biases are estimated from the ensemble mean with respect to the NCEP
10 analysis, and the same bias estimation is applied to each ensemble member
11 during the calibration. On the other hand, the Canadian ensemble includes 20
12 perturbed forecasts and one control forecast. All are performed with the GEM
13 model but use different physics parameterizations, data assimilation cycles
14 and sets of perturbed observations (more information is available online at
15 http://www.weatheroffice.gc.ca/ensemble/index_e.html). Therefore, the individual
16 bias is estimated and used to correct each individual member independently.
17 For ease and efficiency, each participating ensemble calibrates its raw
18 forecast against its own analysis.

19

1 The bias correction procedure generates ensemble output on $1^\circ \times 1^\circ$
2 latitude/longitude grids for 35 selected variables. Table 1 lists all the variables
3 that are post-processed for both the NCEP/GEFS and CMC/GEFS. The
4 selection of these variables depends heavily on the assumption that the
5 forecast variable is well-represented by the Gaussian distribution. It will not
6 work very well for non-Gaussian-distributed variables, such as precipitation.
7 A new method is required for non-Gaussian-distributed variables.

8

9 **3. Experiments for comparing the usage of reforecast information**

10

11 In addition to the decaying averaging bias correction method used
12 operationally for the NCEP/GEFS and CMC/GEFS, three other post-
13 processing methods were designed in this study to assess and mitigate
14 ensemble biases in the 1st moment of the ensemble. The discrepancies among
15 these three methods are a way to estimate bias, which are through using
16 decaying averaging with a short sample, climate mean errors from the
17 ESRL/PSD reforecast and from dependent data, respectively. Discrepancies
18 contribute to the error estimation and the differences show the amount of
19 improvement.

20

1 *a) Bias estimation from decaying averaging with short training data*

2

3 The first method applies the Kalman filter type algorithm to get bias
4 estimations through the following procedure: a) A prior estimate starts the
5 procedure. At a given day T, calculate the time mean forecast errors between
6 days T - 46 and T - 17 to create an initial average. b) Updating. The average is
7 updated by setting it to the weighted average of the new forecast error at day
8 T - 16 with a weight of w, and to the previous average with a weight of 1-w
9 ($0 \leq w < 1$). c) Cycling. Repeat step b) every day from day T-15 to T-1.
10 Experiments with different decay weights w (1%, 2% and 10%, respectively)
11 were conducted and a detailed discussion of results can be found in Section 4.

12

13 This method is different from the operational decaying averaging technique as
14 it chooses a fixed length period of training data. The application of this
15 method is not applicable in operations since it requires a huge amount of disk
16 space to save 46 days of forecasts online. The design and testing of this
17 method was done before the final operational technique was selected. The
18 current operational decaying averaging technique is adapted from it and is
19 much simpler to implement. The bias estimate in operations is carried out
20 iteratively which takes the use of a long training dataset but doesn't require

1 this extra dataset to be saved on disk. The operational method is more flexible
2 in practice and avoids the issue of disk space required to save data for days T
3 - 46 through T – 1. However, results from the bias corrected ensemble using
4 46 days of training data are still useful. Comparisons between it and the other
5 calibration techniques will help to identify their strengths and weaknesses.

6
7 *b) Climate mean bias estimation from ESRL/PSD reforecast ensemble*

8
9 A second method to assess the bias is by using the climatological mean
10 forecast error, which is gotten from the ESRL/PSD 25-year reforecast
11 ensemble (from 1978 to 2003). Hamill et al. (2004) thought that it was not
12 effective to do a bias correction with only a short set of prior forecasts
13 because systematic errors may not be well established if only a few cases are
14 used in the tests, but that errors may be more obvious with the larger sample
15 afforded by a reforecast. With the Model Output Statistics (MOS) techniques
16 (Glahn and Lowry 1972; Carter et al. 1989; Vislocky and Fritsch 1995) and a
17 frozen forecast model, their results show that dramatic improvements in
18 medium- to extended-range probabilistic forecasts are possible by using
19 retrospective forecasts. Motivated by their success, especially for the week-2
20 probabilistic forecast, we introduced the climatological mean forecast error

1 into our bias correction and utilize it as the bias estimate. Following the 1st
2 moment bias correction procedure mentioned above, the climatological mean
3 forecast error is removed from the ESRL/PSD reforecast for each forecast
4 lead time and individual grid point. The reforecast ensembles with and
5 without the climatological bias correction are then examined and compared to
6 the NCEP operational ensembles calibrated from decaying averaging with
7 short training data. Please note that it is only for convenience that we classify
8 the two ensemble data sets as the operational and reforecast ensembles,
9 because the reforecast is also being run operationally at ESRL/PSD.

10

11 *c) Bias estimate using dependent data*

12

13 A third way to estimate the 1st moment bias of the ensemble is through the
14 calculation of a 31-day running mean forecast error centered on day T. The
15 implementation of this method is not feasible operationally but is used as an
16 optimal benchmark. The optimal scenarios, therefore, are compared to the raw
17 and calibrated ensembles to show how large the improvement in the ensemble
18 forecast could possibly be when using the 1st moment adjustment technique.

19

1 The three bias correction techniques discussed above are applied to the
2 NCEP/GEFS operational and ESRL/PSD reforecast ensembles, respectively.
3 The bias estimation and bias correction are carried out separately at each
4 forecast lead time and for each individual grid point. The bias correction is
5 applied to all ensemble member forecasts. The fields studied include the
6 NCEP/GEFS and ESRL/PSD reforecast ensemble 500hPa geopotential height
7 and 850hPa temperature for the period 1 March 2004 to 28 February 2005 for
8 the 00Z initial cycle. Other calibrated fields available from the operational
9 ensemble include the 2m temperature and 10m U and V components (not
10 shown in this paper). Each data source creates three different ensembles, the
11 raw, bias-corrected and optimal ensembles. For the operational NCEP/GEFS
12 ensemble, the three ensembles are named OPR_RAW, OPR_DAV2%
13 (removing the decaying average bias estimate) and OPR_OPT, respectively.
14 For the ESRL/PSD reforecast ensemble, the three ensembles are named
15 RFC_RAW, RFC_COR (removing the climatological mean bias estimate)
16 and RFC_OPT, respectively. All the ensemble forecasts and analyses are on
17 grid points with a spacing of 2.5 by 2.5 degrees globally. The NCEP
18 operational analysis is used for the bias estimation and verification
19 calculations.

20

1 *d) Evaluation methods*

2

3 Several probabilistic (for ensemble distribution) and deterministic (for
4 ensemble mean) verification methods are used to evaluate the ensemble
5 forecast performance, such as Ranked Probability Skill Score (RPSS),
6 Relative Operating Characteristics (ROC) skill score, excessive outlier,
7 Pattern Anomaly Correlation coefficient (PAC), Root Mean Square error of
8 the ensemble mean (RMS) and relative economic value.

9

10 The RPSS score is one of the most important measures for evaluating the
11 performance of probabilistic forecasts (Toth et al. 2003). The higher the RPSS
12 score (the maximum is 1), the better the probabilistic system is by being both
13 reliable and exhibiting high resolution. The best probabilistic system would
14 be rewarded by a RPSS score of 1. The ROC skill score is a measure of the
15 forecast system resolution (Zhu et al. 1996, 2002). A ROC skill score of 1
16 corresponds to a perfect forecast while 0 indicates no skill above the sample
17 climatology. Continuous Ranked Probability Skill Score (CRPSS) measures
18 the reliability and resolution. For statistics over a long period, CRPSS is very
19 similar to RPSS (Zhu and Toth 2008). Therefore, either one of these two
20 measures are used, whichever is more convenient.

4. Results and discussion

Issues examined in this section include a) the choice of decaying averaging weight factors, b) the results of bias correction techniques applied to the NCEP/GEFS and CMC/GEFS, and c) comparisons of the three experimental calibration techniques.

a) Tests of different decaying averaging weights

The first issue when applying the decaying averaging method is the choice of decaying weight w . Different w 's had been chosen and tested. Figure 2 shows some of the decaying weight sensitivity test results. Among the six curves in Figure 2, the curve OPR_OPT is for the calibrated NCEP/GEFS result using dependent data and the curve OPR_RAW is for the raw NCEP ensemble forecast. The other four curves (OPR_RUN_DAV2%, OPR_RUN_DAV1%, OPR_RUN_DAV0.5%, OPR_RUN_DAV0.25%) are for RPSS of 500hPa geopotential height that are averaged from March 1, 2004 to February 28, 2005 for the Northern Hemisphere with decaying weights of 2%, 1%, 0.5% and 0.25%, respectively. All four calibrated ensembles show improvement compared with the raw ensemble OPR_RAW for all lead times. There is little

1 room for further improvement compared with OPR_OPT test for short lead
2 times until day 4. Though the four curves are close together, for short lead
3 times OPR_DAV2% is better than the other decaying weights. On the other
4 hand OPR_DAV0.25% is the best for week-2 forecasts. Other statistics are
5 calculated that show the 2% ensemble gets large improvements in ROC and
6 BSS scores over the Northern and Southern Hemispheres. The improvement
7 of these scores in summer is more significant than in spring. A higher
8 decaying weight (10%) is also investigated and compared with 2%. The
9 choice of a 10% weight works better for the tropics compared to 1% or 2%. In
10 general, the 2% weight works better for most regions and seasons (not
11 shown). For an optimal result, the decaying accumulated bias with a 2%
12 weight is used in operations and is updated every day for all 35 selected
13 variables.

14
15 For a better understanding and explanation of the above results, another issue
16 related to the post-processing algorithm design is quickly reviewed here - a
17 comparison between the decaying average and equal weight bias estimate
18 approaches. The equal weight approach also makes a 1st moment bias
19 calculation over some previous days but with equal weighting for each day.
20 We applied the equal weight bias estimate method to two seasons of the 2004

1 operational NCEP/GEFS ensemble and compared the results with the
2 decaying weight OPR_DAV2%. Results from the equal weight and decaying
3 weight are very similar (less than 2% weight for a longer forecast, not
4 shown). The reasons to choose the decaying weight not the equal weight
5 include: (a) the decay method has a higher weight for the latest information,
6 which is good for a flow-dependent system (short-term forecast), and (b) the
7 application of the decay weight method is operationally cost effective. There
8 is no need to save extra data on the central computer system, and bias
9 estimates can include more historical information through continuous updates
10 once the latest analysis is available. In general, the result from the decaying
11 average will be better than single average (equal weight) method.

12

13 *b) NCEP/GEFS and CMC/GEFS bias correction in operation*

14

15 Figure 3 shows the mean forecast error (solid) and the mean absolute forecast
16 error (dash) of 2-meter temperature for the Northern Hemisphere. Curves
17 E20s/E20m are for the raw ensembles and E20sb/E20mb are for the bias
18 corrected ensembles. All are from 20 member ensembles. NCEP/GEFS raw
19 ensemble has tendency of warm bias (curve E20s, Figure 3a) and CMC/GEFS
20 displays cold bias tendency (curve E20m, Figure 3b) at all lead time. After

1 bias correction, the mean errors of both NCEP and CMC ensembles are very
2 close to zero. Meanwhile the reduced mean absolute errors of bias corrected
3 ensembles improve that the change of mean error is not from the balance of
4 positive and negative values. This suggests the bias correction technique can
5 effectively alleviate ensemble forecast system from over predicting (too
6 warm) or under predicting (too cold) temperatures. The calibrate ensemble
7 forecasts become more close to the actual temperatures. Because of different
8 characteristics of NCEP and CMC ensemble systems, the extents of
9 improvement from bias correction are different for them.

10
11 Figure 4 shows CRPSS scores of the 2m temperature for 2008 summer,
12 verified over the Northern Hemisphere. Values added from bias correction are
13 noticeable for both NCEP/GEFS and CMC/GEFS. Figure 5 shows the RPSS
14 of 850hPa temperature for fall 2009. Both the NCEP/GEFS and CMC/GEFS
15 are improved due to calibration. Overall bias reductions globally (absolute
16 value) after calibration can reach up to 75% for days 0 to 3, 60% for days 2 to
17 8 and 45% for days 8 to 15 for 500hPa height field (not shown). More
18 evaluation results such as PAC, relative economic value (Zhu et al. 2002) are
19 also examined. They are available online at
20 <http://www.emc.ncep.noaa.gov/gmb/yzhu/html/opr/naefs.html> and are updated

1 seasonally. In general, there are skill gains for forecast days 1-6 due to bias
2 correction.

3
4 However, the calibration technique doesn't work well in all circumstances.
5 For some seasons and variables there is no skill improvement for the week-2
6 forecast. Figure 6 shows there is almost no improvement for the 1000hPa
7 height field RPSS score after day 7. Previous studies indicated there remains
8 room for improvement in the week-2 forecast, as seen when comparing the
9 calibrated and the optimal ensembles displayed in Figure 2. How to improve
10 the current calibration technique? Do we need a hind-cast for the calibration
11 of the week-2 forecast? These questions will be discussed in the next section.

12
13 *c) Calibration techniques comparison*

14
15 Figure 7 shows the annual mean RPSS scores of the 500hPa geopotential
16 heights verified over the Northern Hemisphere. Of the three operational
17 NCEP/GEFS ensembles, the one with optimal bias correction (OPR_OPT)
18 gets the highest RPSS scores among the six curves. The decaying average
19 bias correction algorithm also works well. The RPSS of the OPR_DAV2% is
20 improved versus the OPR_RAW for all lead times, but especially in the short

1 range, which can be judged from the small distance between the two curves
2 OPR_DAV2% and OPR_OPT.

3

4 For the three reforecast ensembles, it is not surprising to note that the optimal
5 bias corrected ensemble (RFC_OPT) shows the best performance when
6 compared with the raw and climatological bias corrected ensembles
7 (RFC_RAW and RFC_COR). A comparison between the RFC_RAW and
8 RFC_COR shows that the RFC_COR shows a noticeable RPSS improvement
9 versus the RFC_RAW, especially for the week-2 forecasts. Using the
10 climatological mean bias estimate, it is possible to make probabilistic week-2
11 forecasts more skillful than the raw reforecast. Though the reforecast
12 ensemble uses an older version model and relatively poorer quality initial data
13 than the operational ensemble (Figure 7), the RFC_COR has an even better
14 performance than the OPR_RAW and OPR_DAV2% after day 10, indicating
15 the effectiveness of a large data sample in improving week-2 forecasts.

16

17 Figure 8 shows the annual mean RMS error of the ensemble mean forecasts
18 for 500hPa height, verified over the Northern Hemisphere. The six curves
19 coming from the six ensembles for forecast days 1-6 are divided into two
20 clusters that belong to the operational and reforecast ensemble forecast

1 groups, respectively. Of the three operational ensembles, the OPR_OPT has
2 the lowest RMS error among the six ensembles. The OPR_DAV2% has
3 reduced RMS errors for the first week compared with the OPR_RAW but its
4 RMS becomes larger for week-2 forecasts. However, the two similar curves
5 of the OPR_OPT and OPR_DAV2% for the first week suggest that there is
6 only a limited opportunity for future improvement in bias correction for the
7 first few days. The big distance between the OPR_OPT and OPR_DAV2%
8 curves for week-2 indicates that the OPR_DAV2% calibration technique has
9 the potential to improve extended forecasts.

10

11 Of the three reforecast ensembles, the RFC_COR has smaller RMS values
12 than the RFC_RAW for all lead times, even for week-2. A comparison
13 between the operational and reforecast ensembles shows that the operational
14 ensemble mean (OPR_RAW) has a much lower RMS error than the
15 ESRL/PSD hindcast (RFC_RAW). The RFC_RAW short-range error is
16 around 50% larger than the OPR_RAW. Though the reforecast runs start from
17 relatively poorer quality initial data than the operational ensemble, the
18 RFC_COR works for short range forecasts and its curve with reduced RMS
19 error comes close to the OPR_RAW curve after day 10. Both Figures 7 and 8
20 show that the decaying averaging with a 2% weight and 45 days of training

1 data works very well in the short range. All measures are improved until day
2 5.

3
4 Figures 9 and 10 are the RPSS and ROC skill scores of 850hPa temperature
5 for the summer of 2004, verified over the Northern Hemisphere. Results are
6 similar to those shown for the 500hPa height. The OPR_DAV2% has better
7 performance than the OPR_RAW and the RFC_COR also shows noticeable
8 improvement versus the RFC_RAW. Notice that there is poor performance in
9 the short range for the RFC_COR versus the RFC_RAW, starting from the
10 initial time. We think it was caused by the bias between the reanalysis, which
11 was used during the climatological bias estimation, and the operational
12 analysis which was used for the ensemble evaluation. The climatological bias
13 estimate is calculated from the difference between the daily forecast and
14 verification climatologies as a function of forecast lead time. The
15 climatologies are computed from 31 day running means using data from 1979
16 to 2003. Therefore, there is an existing bias between the operational analysis
17 and the verification climatologies. Such a bias can be mitigated or removed
18 through bias correction.

19

1 Other bias corrected variables include 2m temperature and 10m U and V
2 components (not shown). Based on the results from different variables, some
3 tentative conclusions have been obtained.

4 The decaying averaging with a 2% weight and 45 days of operational training
5 data works very well over the short range (almost as well as the “optimal”),
6 which makes its application possible for frequent updates of the DA/NWP
7 modeling system. On the other hand the climatological mean bias correction
8 can add value, especially for week-2 probability forecasts. Since the
9 operational analysis/modeling system that supports the NCEP/GEFS
10 undergoes frequent (once or twice a year) changes, it would be a very large
11 computing problem if the reforecast method requiring the same model used
12 for operational forecasting was also used for the reforecasting. No such long-
13 term archive based on the most recent analysis/modeling system is available
14 for the reforecast ensemble. The generation of a large hind-cast ensemble is
15 expensive but may be helpful. The use of up-to-date data assimilation/NWP
16 techniques is imperative at all ranges.

17
18
19
20

5. Summary and future plans

A statistical post-processing algorithm, i.e., the decaying average method has been applied to the NCEP/GEFS and CMC/GEFS to generate calibrated ensemble forecasts. The implementation of this technique is expected to improve NCEP and CMC global ensemble forecasts in order to provide more accurate NAEFS products. Due to the different ensemble configurations, calibration strategies applied to the NCEP and CMC ensembles have been adjusted. The NCEP/GEFS is created by using one model with perturbed initial conditions. We assume the biases from one model have a kind of similarity, and ensemble mean biases are thought to be able to represent these systematic errors. Therefore, the NCEP/GEFS uses the ensemble mean bias to calibrate each member. For the CMC/GEFS, the bias for each individual ensemble member is calculated and used for that member. Both the NCEP/GEFS and CMC/GEFS benefit from the application of bias correction. Several studies have shown the NAEFS, when compared to the CMC and NCEP ensemble system, shows significant improvements both in terms of reliability and resolution (Zhu and Toth 2008; Candille 2009). Even with the attractive properties of the decaying average method, its limitations and performance for some variables in week-2 forecasts in some seasons

1 represents a major drawback. There is room for future improvement from (a)
2 adjusting weights to allow a longer training time, and (b) take advantage of
3 reforecasts/hindcasts.

4

5 To further improve the current operational bias correction technique, three
6 other experiments were designed and assessed using annual retrospective
7 experiments from March 1, 2004 to February 28, 2005. Results show that the
8 decaying average bias estimation method with a short sample works well for
9 the first few days. The calibrated NCEP/GEFS ensemble, after removing time
10 mean forecast errors for the most recent period, has an improved probabilistic
11 performance for all measures until day 5. The reforecast ensembles from the
12 ESRL/PSD with and without a climate mean bias correction are also
13 examined. A comparison between the NCEP/GEFS and ESRL/PSD bias-
14 corrected ensemble forecasts shows that a climate mean bias correction can
15 add value, especially for week-2 probability forecasts. This conclusion is very
16 similar to the studies by Hamill et al. in multiple papers (Hamill et al. 2004,
17 2006).

18

19 The major drawback of climate mean bias correction is the need for a long
20 training dataset, and since a reforecast works best with a frozen model the

1 database must be completely rebuilt whenever the model is updated, which
2 requires large computing resources. In this way, routine improvements to the
3 model are incorporated in the reforecast-based products as soon as they are
4 implemented. However, due to the good performance of the climate mean
5 bias correction, the current reforecast ensemble uses an old low-resolution
6 version of the model system and it is worth the effort to generate the
7 reforecast dataset and apply it to the ensemble post-processing. NCEP has
8 plans with ESRL/PSD to jointly implement a real time hindcast experiment in
9 the 2011-2012 time frame and utilize additional resources to generate a set of
10 historical ensemble reforecasts (20 years). The operational forecast model
11 will be applied to the reforecast configurations. Our post-processing study
12 will benefit from this new high-resolution reforecast dataset. Using the
13 reforecast dataset, we will be able to test our post-processing methodology
14 and compare it with the calibration method developed by ESRL/PSD. New
15 bias correction methods developed under THE Observing-system Research
16 and Predictability EXperiment (THORPEX) project will also be considered
17 for use in the NAEFS statistical post-processing system.

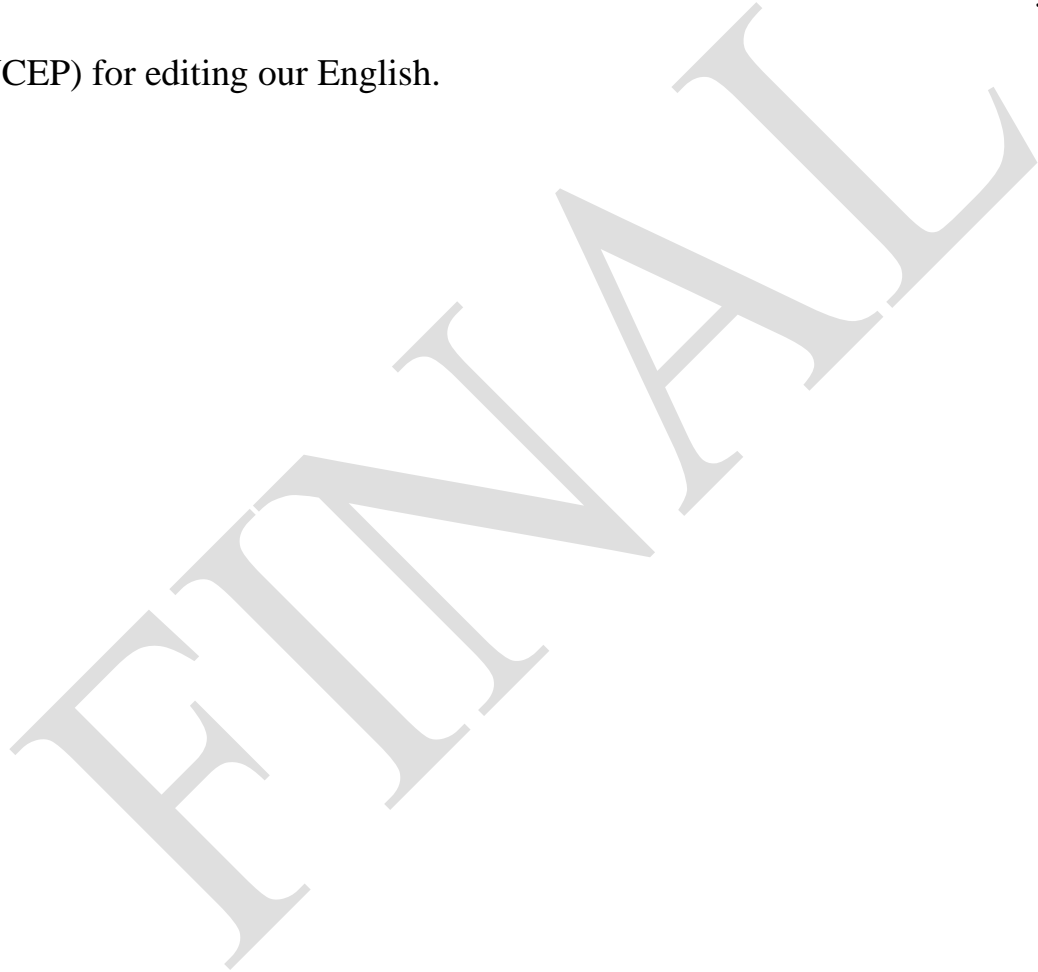
18

19

20

1 **Acknowledgments:** We gratefully acknowledge the support of Jeff Whitaker
2 (NOAA Earth System Research Lab, Physical Sciences Division), Tom
3 Hamill (NOAA Earth System Research Lab, Physical Sciences Division),
4 Richard Verret (CMC), and Richard Wobus (NCEP). This study would have
5 been much more difficult without their assistance. We also thank Mary Hart
6 (NCEP) for editing our English.

7
8
9
10
11
12
13
14
15
16
17
18
19
20



1 **References:**

2

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

5 Candille, G., 2009: The Multiensemble Approach: The NAEFS Example.
6 *Mon. Wea. Rev.*, 137, 1655-1665.

7 Carter, G. M., J. P. Dallavalle, and H. R. Glahn, 1989: Statistical forecasts
8 based on the National Meteorological Center's numerical weather
9 prediction system. *Wea. Forecasting*, 4, 401-412.

10 Charron, M., G. Pellerin, L. Spacek, P. L. Houtekamer, N. Gagnon, H. L.
11 Mitchell, and L. Michelin, 2010: Toward Random Sampling of Model
12 Error in the Canadian Ensemble Prediction System. *Mon. Wea. Rev.*,
13 138, 1877-1901.

14 Cheng, W. Y. Y., and W. J. Steenburgh, 2007: Strengths and weaknesses of
15 MOS, running-mean bias removal, and Kalman filter techniques for
16 improving model forecasts over the western United States. *Wea.*
17 *Forecasting*, 22, 1304-1318.

18 Gel, Y. R., 2007: Comparative analysis of the local observation-based (LOB)
19 method and the nonparametric regression-based method for gridded

1 bias correction in mesoscale weather forecasting. *Wea. Forecasting*, 22,
2 1243-1256.

3 Glahn, H. R., and D. A. Lowry, 1972: The use of model output statistics
4 (MOS) in objective weather forecasting. *J. Appl. Meteor*, 11, 1203-
5 1211.

6 Gneiting, T., A. E. Raftery, A. H. Westveld, and T. Goldman, 2005:
7 Calibrated Probabilistic Forecasting Using Ensemble Model Output
8 Statistics and Minimum CRPS Estimation. *Mon. Wea. Rev.*, 133, 1098-
9 1118.

10 Hacker, J., and D. L. Rife, 2007: A practical approach to sequential
11 estimation of systematic error on near-surface mesoscale grids. *Wea.*
12 *Forecasting*, 22, 1257-1273.

13 Hamill, T. M., J. S. Whitaker, and X. Wei, 2004: Ensemble re-forecasting:
14 improving medium-range forecast skill using retrospective forecasts.
15 *Mon. Wea. Rev.*, 132, 1434-1447.

16 Hamill, T. M., J. S. Whitaker, and S. L. Mullen, 2006: Reforecasts: An
17 important dataset for improving weather predictions. *Bull. Amer.*
18 *Meteor. Soc.*, 87, 33-46.

19 Houtekamer, P. L., L. Lefaivre, and J. Derome, 1996: The RPN ensemble
20 prediction system. *Proceedings, ECMWF Seminar on Predictability*.

- 1 Vol. II. ECMWF, 121–146. [Available from ECMWF, Shinfield Park,
2 Reading, Berkshire RG2 9AX, United Kingdom.]
- 3 Kalman, R. E., 1960: A new approach to linear filtering and prediction
4 problems. *Trans. ASME-J. Basic Eng.*, 82, 35-45.
- 5 Krzysztofowicz, R., and W. B. Evans, 2008: Probabilistic forecasts from the
6 National Digital Forecast Database, *Wea. Forecasting*, **23**, 270–289.
- 7 Charron, M., G. Pellerin, L. Spacek, P. L. Houtekamer, N. Gagnon, H. L.
8 Mitchell, and L. Michelin, 2010: Toward Random Sampling of Model
9 Error in the Canadian Ensemble Prediction System. *Mon. Wea. Rev.*,
10 138, 1877-1901.
- 11 Molteni, F., R. Buizza, T. N. Palmer, and T. Petroliagis, 1996: The ECMWF
12 ensemble prediction system: Methodology and validation. *Quart. J.*
13 *Roy. Meteor. Soc.*, 122, 73–119.
- 14 Raftery, A. E., T. Gneiting, F. Balabdaoui, and M. Polakowski, 2005: Using
15 Bayesian model averaging to calibrate forecast ensembles. *Mon. Wea.*
16 *Rev.*, **133**, 1155–1174.
- 17 Toth, Z., and E. Kalnay, 1993: Ensemble forecasting at NMC: The generation
18 of perturbations. *Bull. Amer. Meteor. Soc.*, 74, 2317–2330.
- 19 Toth, Z., and E. Kalnay, 1997: Ensemble forecasting at NCEP and the
20 breeding method. *Mon. Wea. Rev.*, 125, 3297–3319.

- 1 Toth, Z., O. Talagrand, G. Candille, and Y. Zhu, 2003: Probability and
2 ensemble forecasts. Book of: Forecast Verification: A practitioner's
3 guide in atmospheric science. *Ed.: I. T. Jolliffe and D. B. Stephenson.*
4 *Wiley*, 137-163.
- 5 Toth, Z., Y. Zhu, and R. Wobus, 2004: March 2004 Upgrades of the NCEP
6 global ensemble forecast system. [Available online at
7 http://www.emc.ncep.noaa.gov/gmb/ens/ens_imp_news.html]
- 8 Vislocky, R. L., and J. M. Fritsch, 1995: Improved model output statistics
9 forecasts through model consensus. *Bull. Amer. Meteor. Soc*, 76, 1157-
10 1164.
- 11 Yussouf, N., and D. J. Stensrud, 2006: Prediction of near-surface variables at
12 independent locations from a bias-corrected ensemble forecasting
13 system. *Mon. Wea. Rev.*, 134, 3415-3424.
- 14 Wei, M., Z. Toth, R. Wobus, and Y. Zhu, 2008: Initial perturbations based
15 on the ensemble transform (ET) technique in the NCEP global
16 operational forecast system. *Tellus*, 60A, 62–79.
- 17 Zhu, Y., G. Iyengar, Z. Toth, S. Tracton, and T. Marchok, 1996: "Objective
18 Evaluation of the NCEP Global Ensemble Forecasting System"
19 *Preprints of the 15th AMS Conference on Weather Analysis and*

1 *Forecasting*, 19-23 August 1996, Norfolk, Virginia. Amer. Meteor.
2 Soc.

3 Zhu, Y., Z. Toth, R. Wobus, D. Richardson, and K. Mylne, 2002: The
4 economic value of ensemble-based weather forecasts. *Bull. Amer.*
5 *Meteor. Soc.*, 83, 73–83.

6 Zhu, Y., Z. Toth, R. Wobus, M. Wei, and B. Cui, 2006: May 2006 upgrade of
7 the GEFS and first implementation of NAEFS systems. [Available
8 online at http://www.emc.ncep.noaa.gov/gmb/ens/ens_imp_news.html]

9 Zhu, Y., B. Cui, and Z. Toth, 2007: December 2007 upgrade of the NCEP
10 global ensemble forecast system (NAEFS). [Available online at
11 http://wwwt.emc.ncep.noaa.gov/gmb/yzhu/imp/i200711/IMP_PLAN_final_v08_brief.pdf]
12 [ef.pdf](http://wwwt.emc.ncep.noaa.gov/gmb/yzhu/imp/i200711/IMP_PLAN_final_v08_brief.pdf)]

13 Zhu, Y., and B. Cui, 2008: GFS bias correction. [Available online at
14 http://www.emc.ncep.noaa.gov/gmb/yzhu/imp/i200711/1-GFS_bc.pdf]

15 Zhu, Y., and Z. Toth, 2008: Ensemble Based Probabilistic Forecast
16 Verification. Preprints, *19th Conference on Probability and Statistics*.
17 New Orleans, Louisiana, Amer. Meteor. Soc. 2.2

18
19
20

1 **Figure Captions:**

2 Figure 1. Historical information (days) used for different decaying average
3 weights (0.01, 0.02 and 0.05). Accumulated area (under the curve) is equal to
4 1.0.

5 Figure 2. Ranked Probability Skill Score (RPSS) of 500hPa geopotential
6 height averaged from March 1, 2004 to February 28, 2005 for the Northern
7 Hemisphere with decaying weights 0.25%, 0.5%, 1% and 2%.

8
9 Figure 3. Mean error (solid) and mean absolute error (dash) of 2-meter
10 temperature averaged from June 1, to August 31, 2008 for the Northern
11 Hemisphere for (a) NCEP/GEFS and (b) CMC/GEFS. E20s/E20m is for the
12 raw ensemble forecast. E20sb/E20mb is for the bias corrected ensemble
13 forecast.

14
15 Figure 4. Continuous Ranked Probability Skill Scores (CRPSS) of 2-meter
16 temperature averaged from June 1, to August 31, 2008 for the Northern
17 Hemisphere for (a) NCEP/GEFS, (b) CMC/GEFS. E20s/E20m is for the raw
18 ensemble forecast. E20sb/E20mb is for the bias corrected ensemble forecast.

19

1 Figure 5. Ranked Probability Skill Score (RPSS) of 850hPa temperature
2 averaged from September 1, to November 30, 2009 for the Northern
3 Hemisphere for (a) NCEP/GEFS, and (b) CMC/GEFS. E20s/E20m is for the
4 raw ensemble forecast. E20sb/E20mb is for the bias corrected ensemble
5 forecast.

6
7 Figure 6. NCEP/GEFS Continuous Ranked Probability Skill Scores (CRPSS)
8 of 1000hPa height averaged from December 1, 2007 to February 29, 2008 for
9 the Northern Hemisphere. E20s is for the raw ensemble forecast. E20sb is for
10 the bias corrected ensemble forecast.

11
12 Figure 7. Ranked probability skill score (RPSS) of Northern Hemisphere
13 500hPa geopotential height from March 1, 2004 to February 28, 2005
14 comparing the NCEP operational forecast (OPR) and ESRL re-forecast
15 (RFC). OPR_RAW is the NCEP operational raw ensemble forecast,
16 OPR_OPT is the NCEP operational forecast using optimal bias correction,
17 OPR_DAV2% is the NCEP operational forecast using a 2% weight for bias
18 correction, RFC_RAW is the raw reforecast, RFC_OPT is the reforecast after
19 optimal bias correction, RFC_COR is the reforecast after removing the
20 climatological mean bias.

1 Figure 8. The same as Figure 7, but for Root Mean Square (RMS) errors from
2 the ensemble mean.

3

4 Figure 9. The same as Figure 6, but for Ranked Probabilistic Skill Score
5 (RPSS) of 850hPa temperature for June, July and August 2004.

6

7 Figure 10. The same as Figure 6, but for Relative Operational Characteristics
8 (ROC) score of 850hPa temperature for June, July and August 2004.

9

10

11

12

13

14

15

16

17

18

19

20

21

22

1 Table 1. List of post-processed variables for the NCEP/GEFS and CMC/GEFS ensembles.

2

Ensemble	CMC GEFS (20 member) & NCEP GEFS (20 member, control and GFS)
GRID	1° × 1°
DOMAIN	Global
FORMAT	WMO Grib Format
HOURS	6 hourly out of 384 hours
GZ, TT, U, V	200, 250, 500, 700, 850, 925, 1000hPa
MSLP, Sfc Pres	Mean Sea Level Pressure, Surface Pressure
TT, Tmax, Tmin, U, V	2m Temperature, 2m Maximum and Minimum Temperature, 10m U and V

3

4

5

6

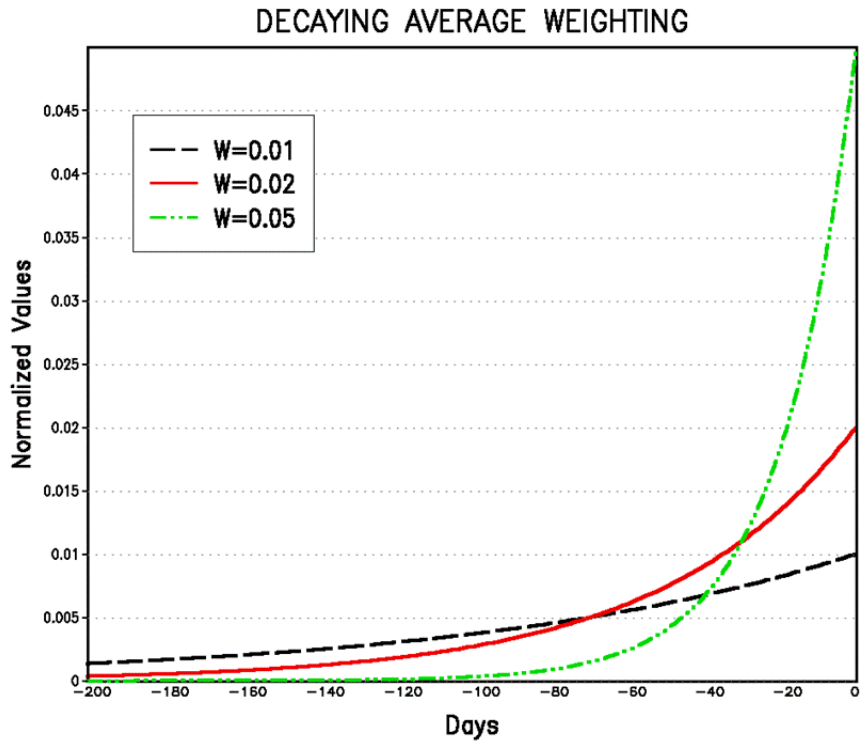
7

8

9

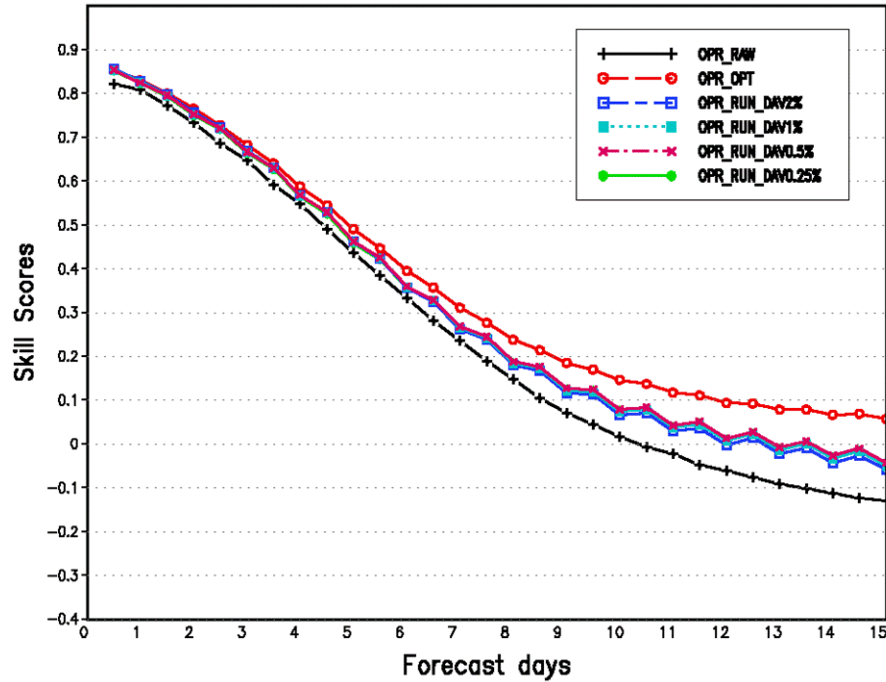
10

FINAL

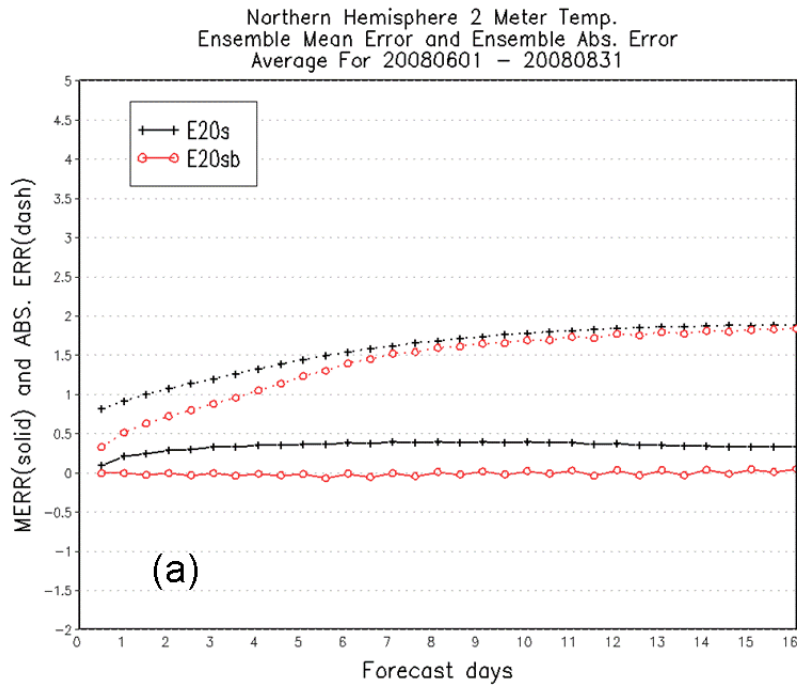


1
2
3
4
5

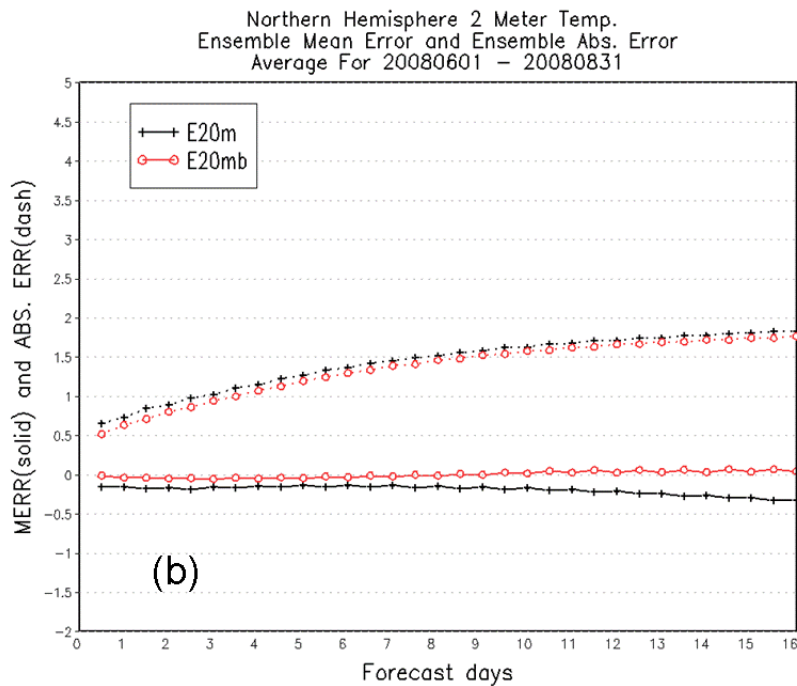
Figure 1. Historical information (days) used for different decaying average weights (0.01, 0.02 and 0.05). Accumulated area (under the curve) is equal to 1.0.



1
 2 Figure 2. Ranked Probability Skill Score (RPSS) of 500hPa geopotential height averaged
 3 from March 1, 2004 to February 28, 2005 for the Northern Hemisphere with decaying
 4 weights 0.25%, 0.5%, 1% and 2%.
 5
 6
 7
 8

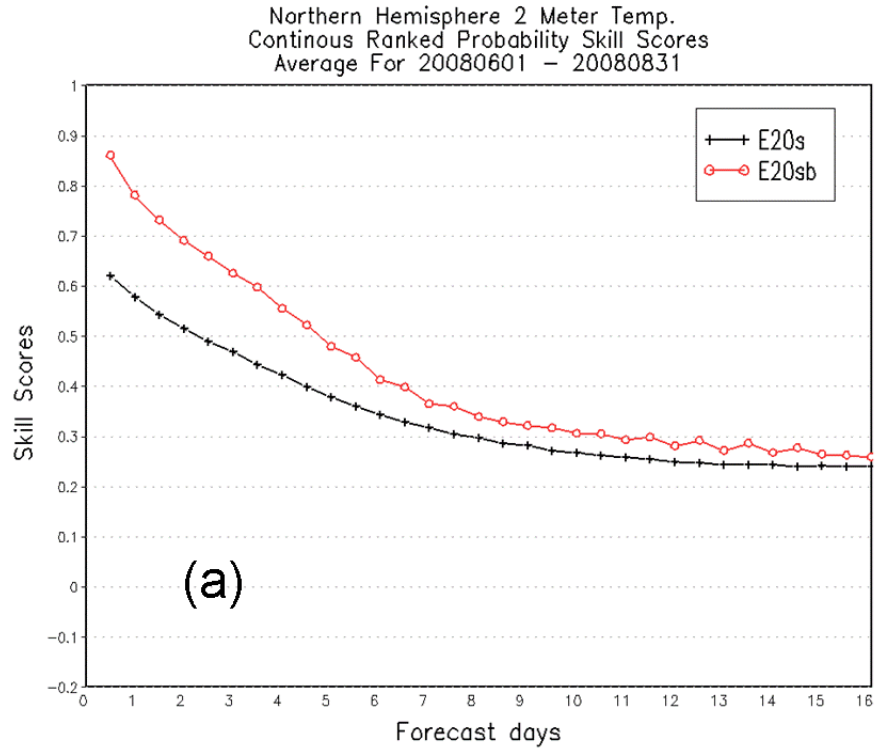


1

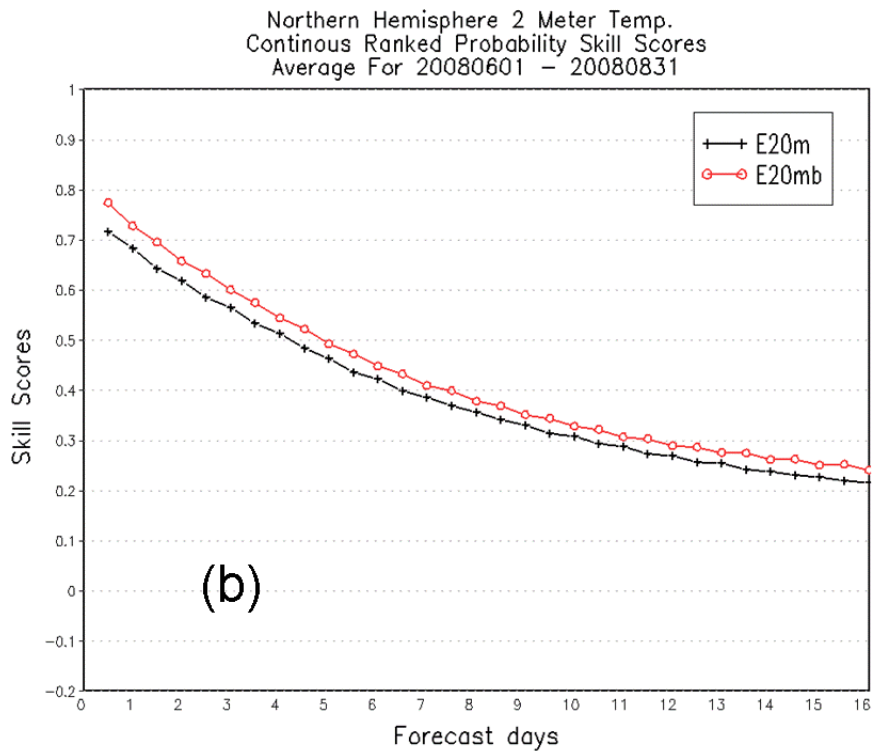


2

3 Figure 3. Mean error (solid) and mean absolute error (dash) of 2-meter temperature
 4 averaged from June 1, to August 31, 2008 for the Northern Hemisphere for (a)
 5 NCEP/GEFS and (b) CMC/GEFS. E20s/E20m is for the raw ensemble forecast.
 6 E20sb/E20mb is for the bias corrected ensemble forecast.



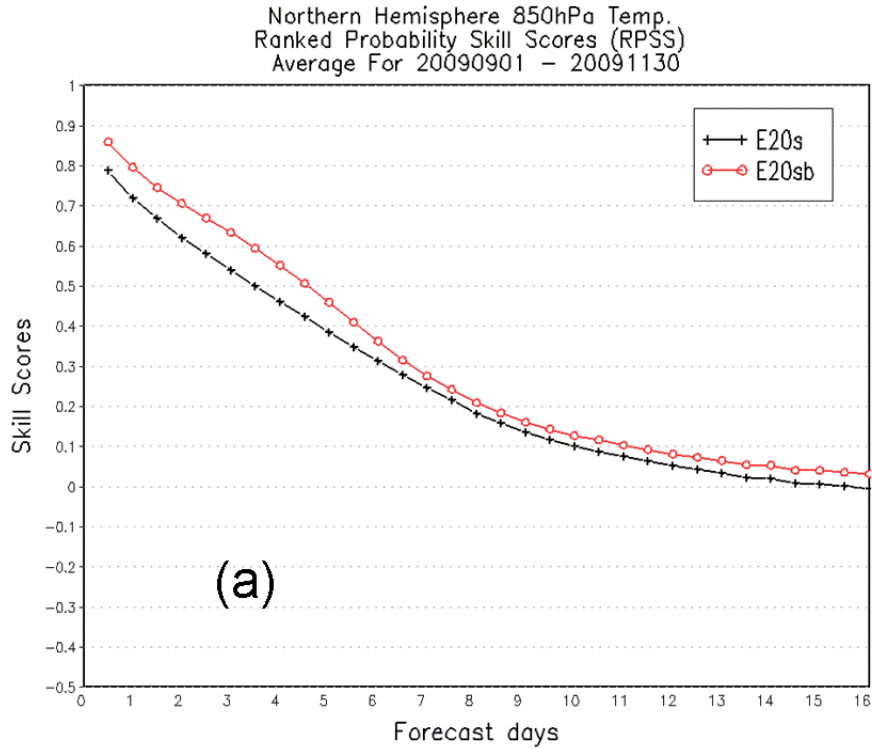
1



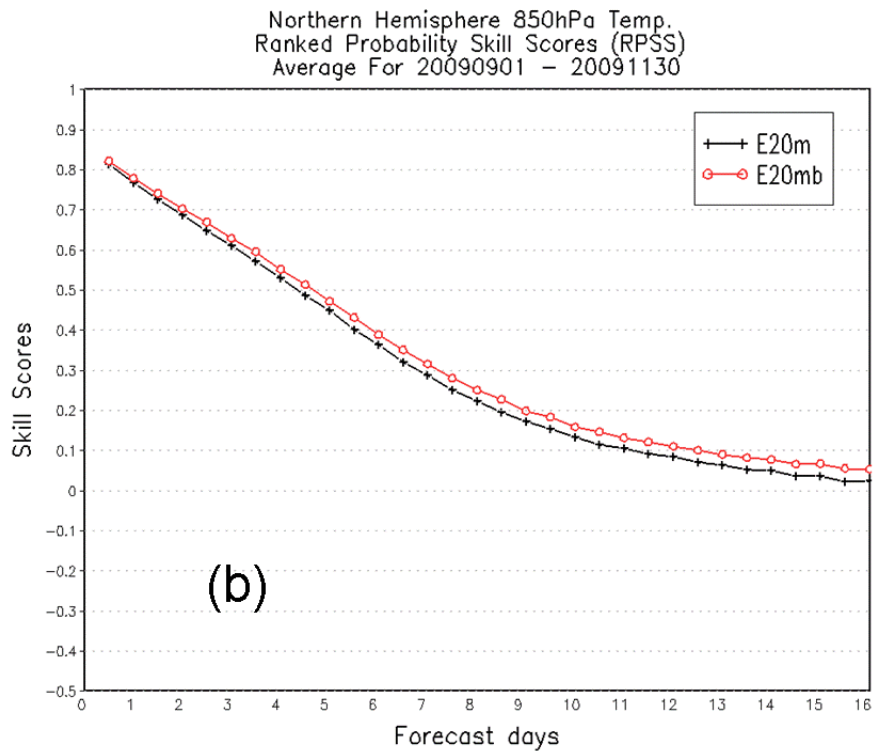
2

3

4 Figure 4. Continuous Ranked Probability Skill Scores (CRPSS) of 2-meter temperature
 5 averaged from June 1, to August 31, 2008 for the Northern Hemisphere for (a)
 6 NCEP/GEFS, (b) CMC/GEFS. E20s/E20m is for the raw ensemble forecast.
 7 E20sb/E20mb is for the bias corrected ensemble forecast.

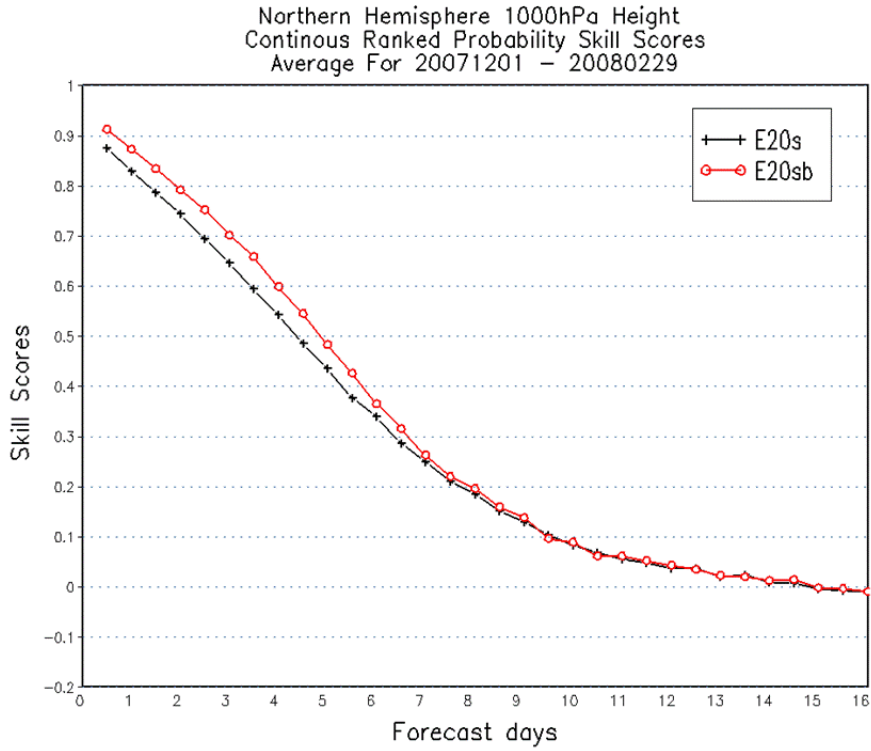


1



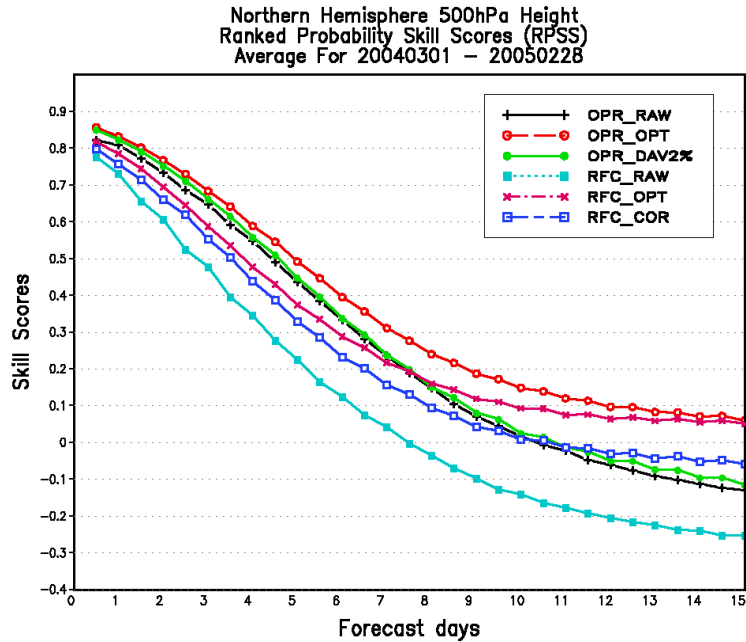
2

3 Figure 5. Ranked Probability Skill Score (RPSS) of 850hPa temperature averaged from
4 September 1, to November 30, 2009 for the Northern Hemisphere for (a) NCEP/GEFS,
5 and (b) CMC/GEFS. E20s/E20m is for the raw ensemble forecast. E20sb/E20mb is for the
6 bias corrected ensemble forecast.



1
2
3
4
5
6
7
8
9
10

Figure 6. NCEP/GEFS Continuous Ranked Probability Skill Scores (CRPSS) of 1000hPa height averaged from December 1, 2007 to February 29, 2008 for the Northern Hemisphere. E20s is for the raw ensemble forecast. E20sb is for the bias corrected ensemble forecast.



1
2
3
4
5
6
7
8
9
10
11
12

Figure 7. Ranked probability skill score (RPSS) of Northern Hemisphere 500hPa geopotential height from March 1, 2004 to February 28, 2005 comparing the NCEP operational forecast (OPR) and ESRL re-forecast (RFC). OPR_RAW is the NCEP operational raw ensemble forecast, OPR_OPT is the NCEP operational forecast by using optimal bias correction, OPR_DAV2% is the NCEP operational forecast using a 2% weight for bias correction, RFC_RAW is the raw reforecast, RFC_OPT is the reforecast after optimal bias correction, RFC_COR is the reforecast after removing the climatological mean bias.

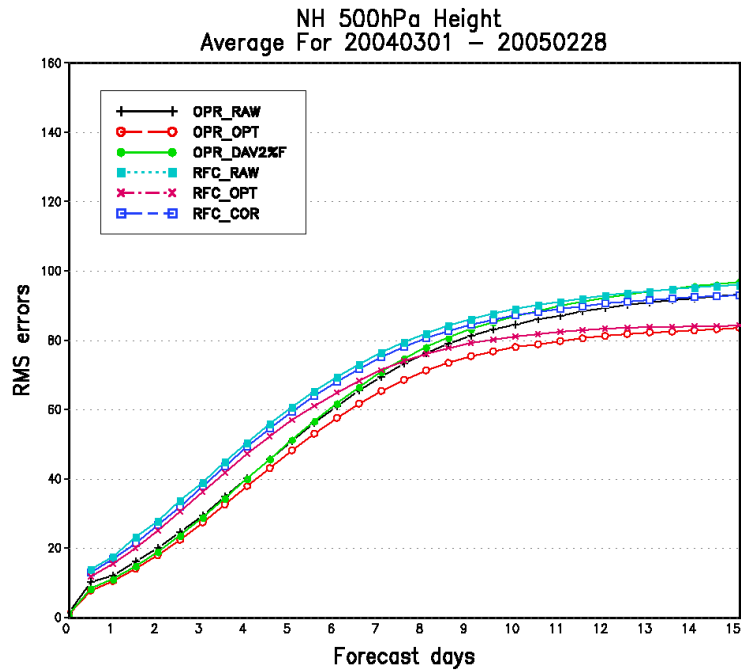


Figure 8. The same as Figure 7, but for Root Mean Square (RMS) errors from the ensemble mean.

1
2
3
4
5
6
7
8
9
10
11
12

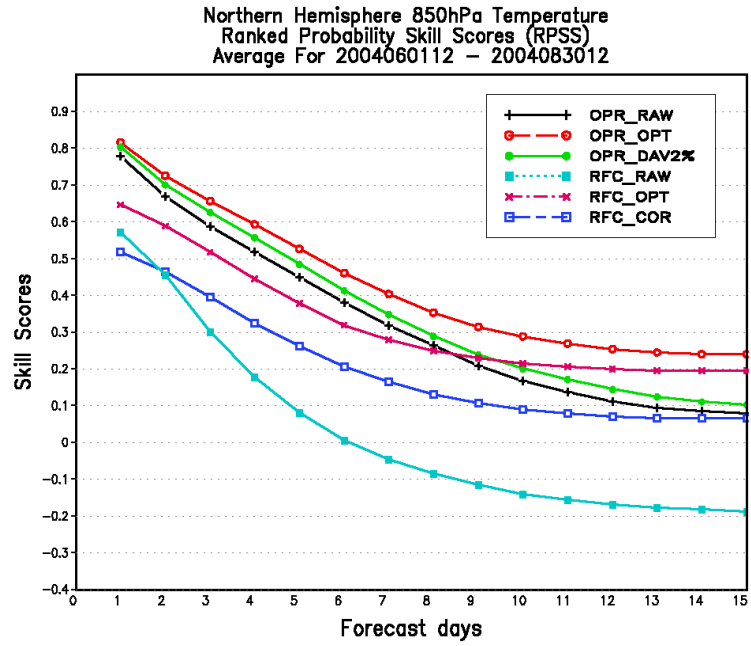
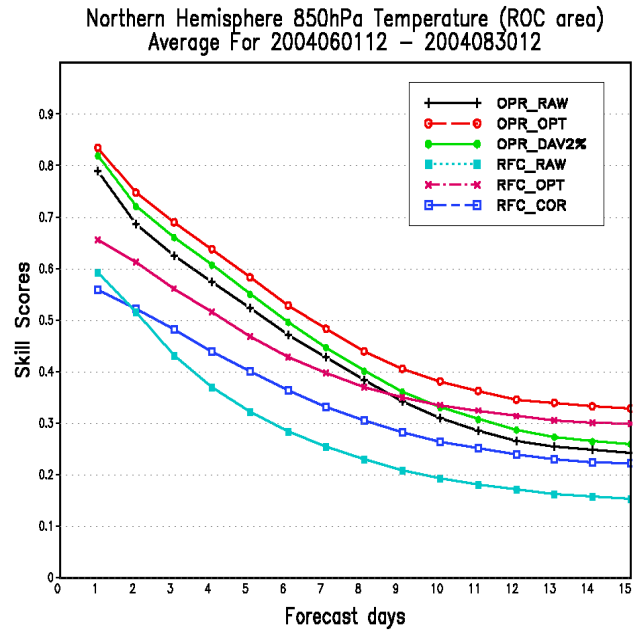


Figure 9. The same as Figure 7, but for Ranked Probabilistic Skill Score (RPSS) of 850hPa temperature for June, July and August 2004.

1
2
3
4
5
6
7
8
9
10
11
12



1
2
3
4
5
6
7

Figure 10. The same as Figure 7, but for Relative Operational Characteristics (ROC) score of 850hPa temperature for June, July and August 2004.