1
2
3
4
5  **Precipitation Calibration Based on Frequency Matching Method**
6  **(FMM)**
7
8
9
10
11  Yuejian Zhu[1*] and Yan Luo[1,2]
12

13  1. *Environmental Modeling Center/NCEP/NWS/NOAA, College Park, MD*
14  2. *I. M. Systems Group, Inc. College Park, MD*

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35  To be submitted to *Weather and Forecasting*
36
37
38
39

40  *\* Corresponding author address:*
41
42  Yuejian Zhu,
43  Environmental Modeling Center/NCEP/NWS/NOAA
44  5830 University Research Court
45  College Park, MD 20740
46  E-mail: Yuejian.Zhu@noaa.gov

1                                          **Abstract**
2
3

4          A post-processing technique is employed to correct model bias for precipitation

5    fields in real time based on a comparison of the frequency distributions of observed and

6    forecast precipitation amounts. Essentially, a calibration is made by defining an

7    adjustment to the forecast value in such a way that the adjusted cumulative forecast

8    distribution over a moving time window dynamically matches the corresponding

9    observed distribution accumulated over a domain of interest, e.g., the entire contiguous

10   United States (CONUS), or different River Forecast Center (RFC) regions in our cases.

11   In particular, the Kalman Filter method is used to catch the flow-dependence and bias

12   information. Calibration is done on a point-wise basis for a specified domain. Using this

13   unique technique, the calibration of precipitation forecasts for the National Centers for

14   Environmental Prediction (NCEP) Global Ensemble Forecast System (GEFS) was

15   implemented into operations in May 2004. To further satisfy various users, especially for

16   hydro-meteorological and short-range weather forecast applications, a recent upgrade to

17   the May 2004's implementation has been made. It includes application of bias correction

18   for higher resolution forecasts with better analysis, and construction of a cumulative

19   frequency distribution based on each RFC region instead of the entire CONUS domain to

20   take realistic regional climate features into account. This study focuses on one degree

21   spatial and 6-hour temporal resolution out to a 384 hour (about 16-day) forecast to

22   provide detailed information on precipitation events, using the newly developed NCEP

23   Climatology-Calibrated Precipitation Analysis (CCPA) as the proxy for the truth. Mean

24   forecast errors and skill are evaluated with respect to CCPA over the CONUS and each

1    RFC for the period 2009-2010 for 6-hour accumulations at one-degree spatial resolution.

2    From this study it was found that this frequency matching algorithm substantially

3    improves NCEP GFS/GEFS model precipitation forecast biases over a wide range of

4    forecast amounts and produces more realistic precipitation patterns.   Moreover, this

5    approach improves the forecast prediction skill measured by most verification scores. In

6    addition, the skill of probabilistic quantitative precipitation forecast (PQPF) has been also

7    improved by applying this method to the individual GEFS ensemble members.

8

9

10

## 1. Introduction

There are many important applications that require a more accurate quantitative precipitation forecast (QPF) and probabilistic quantitative precipitation forecast (PQPF). One of these applications is the daily forecast. The QPF and ensemble based PQPF forecast products were implemented into NCEP operations in the late 1990's (Zhu et al. 1998, Zhu 2005). A better calibrated PQPF could benefit the short- and medium-range forecasts, and extend the forecast predictability (Eckel and Walters 1998, Zhu and Toth 1999). There are several studies on calibration of PQPF; some focus on the methodology (Krzysztofowicz and Sigrest 1998; Christopher et al. 2008) and others use reforecast information (Hamill et al. 2002, Fundel et al. 2009). An analog method has been developed by using large samples of reforecasts (Hamill and Whitaker 2006), and is experimentally run at ESRL and provides additional guidance for NCEP Weather Prediction Center (WPC) forecasters. Another important application of this method is for down-stream applications. Water management decisions are crucially dependent on forecast information regarding the possible future evolution of precipitation. On the other hand, hydrologic models need accurate precipitation forecasts from numerical weather prediction (NWP) as forcing inputs. Therefore, a realistic representation of the precipitation field in forecasts is very important. However, many studies have demonstrated systematic biases in the model precipitation products due to model deficiencies. It has long been recognized that model precipitation uncertainty affects the accuracy of hydrologic modeling (Demargne et al. 2013), because the performance of distributed, physically based hydrologic models depends greatly on the quality of the precipitation input data. For both these reasons, post-processing techniques have been

1  developed and applied to reduce these biases in the precipitation. Many studies have

2  demonstrated some success with precipitation forecasts through statistical post-

3  processing. For instance, Yuan et al. (2007) applied an artificial neural network as a

4  postprocessor to calibrate Probabilistic Quantitative Precipitation Forecasts (PQPF) from

5  the NCEP Regional Spectral Model (RSM) ensemble forecast system. Voison et al.

6  (2010) described two bias correction methods with spatial disaggregation (BCSD) and an

7  analog technique for downscaling and calibrating errors from ensemble precipitation

8  forecasts. In this study we developed a method for precipitation calibration in real time

9  called the "frequency matching method". Basically the methodology employed here is a

10  statistical adjustment based on cumulative frequency distributions of forecast and

11  observed precipitation amounts. Two steps are undertaken in calibration with frequency

12  matching. As first, it requires an observation dataset at the same spatial and temporal

13  resolution as the model forecast output and a reasonable number of days of prior forecasts

14  to construct their respective cumulative frequency distributions for forecasts and

15  observations. In addition, it introduces a time moving window for sampling appropriate

16  historical bias information that makes the cumulative frequency distribution of forecasts

17  match that of the observations. The second step in this method makes use of the

18  cumulative frequency distributions of observations and forecasts, in this way a frequency

19  match is performed between prior observations and forecasts. The resulting correction

20  factor is applied to adjust a target forecast value at each grid point and each grid point is

21  treated individually.

22      In this paper, we first briefly review the background of the 4 May 2004

23  implementation at NCEP. The 2004 implementation was first developed as a pioneer

1  version of precipitation calibration with frequency matching for application in

2  precipitation forecasts with 24 hour accumulations at 2.5 degree resolution (Zhu and Toth

3  2004). Just as with any other numerical weather prediction (NWP) model, Quantitative

4  Precipitation Forecasts (QPF) from the Global Forecast System (GFS) at NCEP suffer

5  from biases due to model deficiencies. Probabilistic Quantitative Precipitation Forecasts

6  (PQPF) based on the Global Ensemble Forecast System (GEFS) at NCEP are biased as

7  well due to imperfections in the model and ensemble formation. Typically, model

8  precipitation bias is dependent on the model version, lead time and location. In most

9  cases, small amounts of precipitation are over-forecasted while large amounts are under-

10 forecasted. By calibrating each member of the ensemble based on verification statistics

11 accumulated over the continental US (CONUS), the bias in QPF (first moment) is

12 practically eliminated, and the PQPF (second moment) is substantially improved. By

13 following the approach of the 2004 implementation with timely availability of higher

14 resolution model output and a better analysis, named the Climatology-Calibrated

15 Precipitation Analysis (CCPA, Hou et al. 2013), we pursue a similar application at one

16 degree resolution and every 6-hours out to 384 hours (about 16 days) globally.

17      To provide a better proxy of the truth for the precipitation field over CONUS at

18 high spatial and temporal resolutions, CCPA has been developed and evaluated at NCEP

19 by Hou et al. (2013). The dataset takes advantage of the higher climatological reliability

20 of the CPC dataset (Xie et al. 2010) and the higher temporal and spatial resolution of the

21 Stage IV dataset (Lin and Mitchell 2005). Thus, CCPA is reliable and quality controlled,

22 with a high spatial and temporal resolution. It is available as 6 hour accumulations from

23 2002 onwards. The CCPA data are first produced on the 4 km HRAP (Hydrologic

1    Rainfall Analysis Project) grid, the same as the NCEP Stage IV over CONUS, as a

2    primary product and then interpolated to 1, 0.5, 0.125 degree and NDGD (5km) grids by

3    a volume conservation scheme as by-products.  The 1 degree CCPA is applied in this

4    study as it exactly matches the model output grid.

5    We continue to investigate here the method that applies to the NCEP GFS/GEFS

6    precipitation model output with CCPA. Then we analyze aspects of the bias correction of

7    ensemble precipitation forecasts, including precipitation forecast skill and reliability. Our

8    objective is to produce bias-corrected precipitation ensemble forecasts through post

9    processing for near real time forecast applications.

10    The reminder of the paper is organized as follows. Section 2 describes the

11    frequency matching method for precipitation calibration. Section 3 reviews the

12    background of 2004 implementation. A few cases to demonstrate the success of this

13    method will be presented. Section 4 applies and evaluates the bias correction approach

14    for higher resolutions using CCPA, and in the last section we present our conclusions

15    with suggestions for future work that will further improve the calibration of precipitation.

16

17    **2. Methodology**
18
19    A systematic difference (or 'bias') between forecast and observed precipitation

20    amounts can be progressively removed using information provided by observations. In

21    this study, the bias information can be estimated through comparing forecast and

22    observed precipitation frequency distributions. The general frequency matching method

23    proceeds as follows. First, we conduct a bias assessment by constructing a cumulative

24    distribution function (CDF) for the preceding forecast and corresponding observed

1    precipitation amounts. Given a set of precipitation thresholds in ascending order, the CDF

2    is calculated as the count of numbers of grid points over a given domain where the

3    forecast or observed precipitation values exceed a threshold. The CDF is updated with the

4    Kalman filter method, which is similar to the bias correction method in the NAEFS

5    (Northern American Ensemble Forecast System) (Cui et al. 2011, 2013), expressed as:

6
7    $$\overline{CDF}_{i,j} = (1-W) * \overline{CDF}_{i,j-1} + W * CDF_{i,j} \qquad (1)$$
8
9    where $\overline{CDF}_{i,j}$ is the decaying averaged CDF at threshold i for day j, while $\overline{CDF}_{i,j-1}$ is the

10   prior decaying averaged CDF for day j-1. $CDF_{i,j}$ is the newly counted CDF at threshold i

11   for day j. $W$ is the decaying weight between 0 and 1, defined simply by an approximated

12   time moving window nd (nd cannot equal zero).

13   $$W = 1 / nd \qquad (2)$$

14   Here a time moving window (or decaying weight) is chosen to make a weighted average

15   of these CDFs over the domain depending on how far it is from the target forecast day,

16   which is illustrated in Figure 1. The higher the weight the faster the decaying speed

17   (which indicates there is a higher weight on the most recent data and less on the oldest

18   data) and vice versa. Our strategy is to specify prior forecast days (an approximated time

19   moving window, or decaying weight) for each grid point and each lead time as a pool for

20   sampling appropriate historical bias information from forecasts and observations. For

21   instance, a 50-day window (W=0.02) means training data is accumulated over the most

22   recent 50-day period with the most weight on the most recent data (See Figure 1 for

23   W=0.02).  Thus, the idea behind the adaptive method is to catch the dynamic flow-

24   dependence and statics of observations. The time moving window (or decaying weight)

25   can be tuned from short (or large) to long (or small) times (weights) to ensure the best

1    performance of the method. In our adaptation of the frequency matching method, there

2    are two ways to construct CDFs for forecasts and observations. We call the CDF based

3    on the whole CONUS domain the CONUS CDF, and the CDF based on each RFC region

4    (See Figure 2) is the RFC CDF. For each grid point within a specific domain (e.g.,

5    CONUS or any RFC) and for each forecast lead time, the observed and forecast CDFs are

6    derived using the same time moving window (or decaying weight). To be useful for

7    applications, this method needs to handle the initial CDFs, which is termed spin-up.

8           Second is the bias adjustment. In order to keep the spatial and temporal coherence

9    of a forecast as similar as possible to that of the observation, we match the cumulative

10    frequency distribution of the forecast to that of the observation using a frequency

11    matching algorithm. Here the updated CDFs from Equation (1) form cumulative

12    frequency distributions. As illustrated in Figure 3, according to this pair of distributions,

13    for a raw forecast value ("RAW") we find and assign an observed value that has the same

14    frequency within a given domain as the forecast value to the correspondent calibrated

15    forecast ("CAL"). Consequently, the bias information is estimated based on the paired

16    and updated CDFs for the forecast and corresponding observed values. For example, in

17    Figure 3 in the case of a CDF (forecast) greater than the CDF (observed), model

18    precipitation tends to be over-forecasted, so to match the frequency a correction factor of

19    less than one will be expected to reduce a forecast value.  In doing this matching process,

20    linear interpolation is applied twice in real calculations to derive a correction factor for

21    each grid point. Mathematically, given an array of thresholds $T_1$, $T_2$, …, $T_n$ in ascending

22    order as the abscissas, an array of observed CDFs $O_1, O_2$, …, $O_n$ and an array of forecast

23    CDFs $F_1, F_2$, …, $F_n$ as ordinates, an array of calibrated thresholds $T_1^*$, $T_2^*$, …, $T_n^*$ are

1    derived through the first linear interpolation. Consequently, the forecast CDF $F_i^*$ at $T_i^*$

2    ($i=1,\ldots,$ n)   is equal to observed CDF $O_i$   at $T_i$ ($i=1,\ldots,$ n) just as in what we call

3    frequency matching. That is:

4           $O_1(T_1) = F_1^*(T_1^*),$

5           $O_2(T_2) = F_2^*(T_2^*),$

6           ……

7           $O_n(T_n) = F_n^*(T_n^*).$

8    Next, a correction factor is calculated as the ratio of a calibrated threshold to its related

9    threshold, i.e., $Ri = T_i^* / T_i$, $i=1,\ldots,$ n. Once again, given an array of thresholds $T_1, T_2, \ldots,$

10   $T_n$ as the abscissas and the array of correction factors $R_1, R_2, \ldots, R_n$ as ordinates, for a

11   forecast value ("RAW") at any grid point a correction factor ("r"), the ratio of a

12   calibrated forecast value ("CAL") to its corresponding raw forecast value at a grid point,

13   is derived by linear interpolation. Then the correction factor ("r") is applied to the raw

14   forecast value ("RAW") to compute the final calibrated forecast value ("CAL= r *

15   RAW"). This correction is applied to each model grid point which implies that the

16   correction is a function of forecast value. No adjustment of a zero precipitation forecast

17   value is made in order to prevent an unrealistic negative precipitation value due to

18   interpolation.

19           This calibration technique with frequency matching should work with any model

20   output as long as observations are available and are processed to be at model grid points.

21   However, our experience with this technique indicates some important considerations

22   must be addressed. That is, precaution must be taken about the selections of thresholds

23   and number of decay days, particularly when the CDF is calculated for each RFC rather

1   than the CONUS because there will be a much smaller sample size as implied from Table

2   1. For example, an insufficient amount of non-zero sample data is very likely to cause

3   more than two equal values of zero as CDFs for adjacent highest thresholds, though this

4   situation is not allowed in this method as it may lead to a failure in the interpolation.  To

5   deal with this problem, selecting a reasonable range of thresholds is necessary to produce

6   non-equal CDF values. Another solution is choosing a proper number of decaying days.

7   If the number of decaying days is too small it will be problematic since there will not be

8   sufficient sample data, especially when dry-climate regions experience a long duration

9   drought. Therefore, there is an inevitable trade-off as to the number of decaying days

10  when tuning for optimal calibration performance. It is believed that potential difficulties

11  in CDF construction in dry regions are related to the small number of days with

12  precipitation, imposing a practical challenge to this method. When the above statistical

13  deficiencies and operational limitations are avoided, the method should be

14  computationally realistic and feasible for real-time implementation.

15

16  **3.  Background review**
17
18      In the 2004 implementation (Zhu and Toth 2004) the calibration system was

19  designed to apply a bias correction globally to all 00 UTC forecasts, including high and

20  low resolution control forecasts and all ensemble member forecasts for 24 hour amounts

21  at 2.5 degree resolution. The operational NCEP GFS/GEFS forecast system runs four

22  times per day (00, 06, 12 and 18 UTC) and produces 1 degree global ensemble

23  precipitation forecast products for 6-hour accumulations.   It contains twenty-two

24  ensemble members - a high resolution GFS run, low resolution GEFS control run and ten

1    pairs of perturbed runs using the ET method (Wei et al. 2006, 2008). Once generated,

2    precipitation forecasts from the 00 UTC cycle only are processed into 24 hour

3    accumulations and aggregated to 2.5 degree resolution prior to bias correction. Technical

4    information about NCEP's latest GEFS ensemble forecast system is available online (Zhu

5    et al. 2012)

6        Bias assessment is approached separately for the GFS high resolution and

7    ensemble control (low resolution) forecasts at each lead time to save computational time.

8    Data are sampled from prior forecasts and observations with a 30-day average of the

9    whole CONUS domain as the cold start sampling. Later, the corresponding decaying

10   weight used is 1/30. The observations with 24 hour accumulations come from the US

11   RFC rain gauge network with about 10,000 observation station reports after re-gridding

12   to the common 2.5 degree model grid. A set of thresholds of 0.2, 2.0, 5.0, 10.0, 15.0,

13   25.0, 35.0, 50.0, 75.0 mm/day were carefully selected for the 24 hour accumulation

14   amount to ensure no failure of interpolation in the calibration procedure, as detailed in

15   Section 3. The bias assessment based on CONUS CDF may be applied to the global

16   domain, when assuming that the bias information over CONUS is much the same as over

17   other parts of the globe, which may not be an optimum application. This application can

18   be improved when global precipitation observations become available in real time.  The

19   calibration system runs once daily at the 00 UTC cycle and typically the daily runs are

20   completed within a minute in real time on a supercomputer.

21       The evaluation period for this implementation was chosen to be 1 Dec. 2000 − 28

22   Feb. 2001. Comparisons of the calibrated forecast against the raw forecast in terms of

23   some scores were made and are shown in Figure 4.  Figure 4(a) presents equitable threat

1    scores (ETS) and bias scores at the 2.0mm threshold for each forecast lead time. Figure

2    4(b) provides the 36-60 hour reliability diagram at the 2.5mm threshold, validated for all

3    grid points in the CONUS.  The calibrated forecast shows a remarkably improved bias

4    score over CONUS at all thresholds. Not only is the bias reduced, the post-processing

5    through frequency matching helped increase the probabilistic forecast skill, such as with

6    the Brier score (not shown). There was a much reduced PQPF (mean) bias in the

7    calibrated forecasts, indicating a dramatic improvement in reliability relative to the raw

8    forecasts. The reliability curve approaches the diagonal line, which indicates that the

9    biases in PQPF were removed to some degree. The Brier score was also improved

10   dramatically at all lead times. In general, these calibrated forecasts were much more

11   skillful than the raw forecasts at all lead times.

12
13   **4. Applications and evaluations**
14
15        In this section, we expand on earlier work to upgrade the calibration system and

16   make it capable of bias correction at higher temporal and spatial resolutions. More

17   specifically, we use the current application at 1 degree resolution with 6 hour

18   accumulations. This section describes how the higher resolution precipitation forecasts

19   are calibrated so that their cumulative frequency distribution matches that of the

20   observations.

21        The operational NCEP GFS/GEFS 6-hourly precipitation forecast (up to a 384

22   hour lead time) has a spatial resolution of one degree latitude and longitude and runs up

23   to real time. There is a high resolution GFS run, a low resolution GEFS control run and

24   20 ensemble members for each forecast. Unlike in the 2004 implementation, here all 6-

25   hourly one-degree forecasts for the four cycles are directly bias corrected with respect to

1  the gridded precipitation analysis CCPA at the same resolution as the forecasts. To be

2  more realistic and better capture regional climate regimes, 12 RFC CDFs are derived for

3  each lead time to construct cumulative frequency distributions. For each category of the 9

4  thresholds (0.2, 1, 2, 3.2, 5, 7, 10, 15, 25 mm/6hr), a CDF is calculated as the number of

5  grid points over each RFC where the forecasts or observed precipitation amounts are

6  greater than the threshold. Again, to reduce the computational burden, we only derive one

7  set of CDFs from the high resolution GFS run and another set of CDFs from the low

8  resolution GEFS control run.  Then the latter set of CDFs is applied to the 20 ensemble

9  members since all of them are low resolution forecasts from the same forecast model,

10  resulting in 2 rather than 22 sets of CDFs per lead time per threshold per RFC region. In

11  each bias correction run, there are a total 1536 forecast-observation CDF pairs for 64

12  forecast lead times for the low high resolution runs and 768 pairs for 30 forecast lead

13  times for the high resolution runs, summed for a total of 9 thresholds and 12 RFC

14  regions. Because there is no high quality global precipitation analysis available, the CDF

15  of CONUS (is sum of 12 CDFs for all the RFCs) is using for a grid point outside of

16  CONUS, therefore  the quality of precipitation calibration for these areas is limited.. Bias

17  information is sampled with an approximate 50-80 day moving time window, and thus

18  the decaying weight selected is 0.02. The bias correction is applied four times per day to

19  each 6-hourly forecast at each grid point globally and to each forecast lead time

20  independently.

21      The operational forecasts initialized daily at 00 UTC from 1 March 2009 through

22  28 February 2010 will be assessed. These forecasts produced with the same modeling

23  suite were used to produce the calibrated forecasts. Both sets of forecasts will be

1    examined out to 384 hours with precipitation accumulation output available every 6

2    hours. Although the method we developed can apply to global forecasts, in this study our

3    evaluation domain is the CONUS, which allows evaluations of this method using the one

4    degree CCPA dataset. The evaluation focuses on the biases and skill levels of the

5    calibrated ensemble precipitation forecasts with respect to raw forecasts. We analyze

6    several examples and present some verification statistics. The verification statistics will

7    be stratified by either lead time or threshold.

8        Figure 5 shows one application of this calibration for the high resolution GFS

9    forecast for selected forecast lead times (78-hr, 84-hr, 90-hr and 96-hr). The comparison

10   is of 6 hourly accumulated precipitation (mm) initialized at 00UTC 24 January 2010 for

11   the raw GFS forecast (left), calibrated forecast (middle) and observation (CCPA, right).

12   Apparently the GFS over-forecasted for the CONUS in general, and the calibrated

13   forecast reduced the forecast amount accordingly. Figure 6 shows the ensemble PQPF

14   (same time period) for the 0.254mm/6 hours threshold where raw ensemble PQPF is on

15   the left, calibrated PQPF (CPQPF) is in the middle and the observation is on the right.

16   The forecast area of the PQPF is reduced; the quantity (value) of PQPF is smaller in the

17   calibrated PQPF, which matches better with the observations.

18       To demonstrate the benefits from this calibration, several different scores have

19   been presented for the seasonal and yearly averages. The bias scores and ETS for

20   CONUS for the period of 1 December 2009 - 28 February 2010 are shown in Figure 7

21   and Figure 8. Figure 7 (a) is for the 0-6 hour forecast bias of the different thresholds, and

22   Figure 7 (b) shows forecast lead times out to 180 hours for greater than 0.2mm/6 hours.

23   The numbers above the thresholds in Figure 7 (a) indicate the sample size of the one by

1   one degree forecast box we have verified. Overall, the bias is reduced and ETS is

2   increased in the calibrated forecasts for both the GFS and GEFS control for all lead times,

3   and the improvement of ETS tends especially to be more effective for shorter lead times.

4   Similar improvements in bias scores and ETS are also observed for the RFC regions,

5   such as the MBRFC and NERFC shown in Figure 9 and Figure 10, respectively, although

6   they exhibit slightly larger diurnal variability.

7       The RMSE (root mean square error) and ABSE (absolute error) of CONUS for

8   the period of 1 March 2009 – 28 February 2010 (one year) for every 6-hr accumulated

9   precipitation forecast are shown in Figures 11 and 12. Figure 11 is for the GFS forecast

10  and Figure 12 is for the GEFS control forecast. Based on this year of statistics, RMSE is

11  reduced significantly for the GFS, but not for the (lower resolution) GEFS control. This

12  difference might be related to the model resolutions and model versions (the operational

13  GFS model version is slightly different from GEFS for this period due to different

14  implementation times). In particular, the higher resolution model produces larger errors

15  compared to the lower resolution model due to resolution and forecast sharpness (Figure

16  11 and Figure 12). It may be better to separately verify the forecast intensity and pattern

17  (or position). The results could be different if different verification methods are applied,

18  such as MODE (the Method for Object-Based Diagnostic Evaluation; Davis et al. 2006a,

19  2006b). However, the RMSEs are very similar after calibration for both the higher and

20  lower resolution model forecasts. Meanwhile, for this one year of statistics, ABSE is

21  reduced for both the GFS and GEFS control at all lead times.

22      For the ensemble forecast, RMSE and ABSE of the ensemble mean, ensemble

23  spread and CRPS (Continuous Ranked Probability Score; Zhu and Toth 2008) have been

1 calculated for the period of 1 March 2009 – 28 February 2010 and displayed in Figure 13.

2 This is one year verification against CCPA for every 24-hr accumulated precipitation

3 forecast. The results indicate that 1) the RMSE is marginally reduced (similar to the

4 ensemble control in Figure 11) and the ABSE for ensemble mean is significantly

5 reduced; 2) CRPS is improved; and 3) ensemble spread is increased for longer lead time

6 forecasts. The improved spread and CRPS could be explained as a by-product of the

7 frequency matching method. The algorithm not only matches the precipitation frequency

8 (reducing the bias), but also adjusts the amount of precipitation forecasted by each

9 ensemble member (adjusting the distribution). A comparison of the Brier scores between

10 raw and calibrated forecasts is also shown in Figure 14. The Brier score is negatively

11 oriented, which means the smaller the score value the better the results. As expected, the

12 score is reduced after bias correction (dotted curves) for all lead times.

13 **5.  Conclusions and future plans**

14      The frequency matching method is developed and applied to the NCEP QPF and

15 PQPF forecasts for the first precipitation calibration since 2004. The latest version will be

16 implemented in 2013 for finer temporal (every 6 hours out to 16 days) and spatial (1*1

17 degree) resolutions. The prior CDFs of the forecast and observation can be easily

18 generated from the GFS/GEFS precipitation forecast and CCPA through applying the

19 Kalman filter method (or decaying average). The performance of this method has been

20 investigated with respect to one year of operational GEFS precipitation products. Results

21 show that model bias has been effectively reduced and some skill scores have been

22 improved in the calibrated forecasts. The good performance of the bias correction is

23 obviously due to the fact that it can dynamically catch systematic model biases in most

1    cases.  Another attractive advantage of this method is that it saves a significant amount of

2    both computer and human resources. Unlike other statistical post-processing methods, it

3    is not heavily reliant on a huge amount of data for model bias training, so it takes up

4    much less disk space on the computer systems and is able to update the model bias when

5    a model is upgraded as well.

6         One important issue is the validity of the frequency matching method. The

7    method used in this study is based on a certain knowledge of model bias information

8    drawn from past verification statistics. Remember that as mentioned in Section 2, this

9    method is not perfect as it is unable to make adjustments to areas that have no

10   precipitation; therefore, this kind of dry bias can never be removed, though this is also the

11   case with other traditional precipitation bias correction methods. Generally, model

12   forecasts include two kinds of errors, intensity error and pattern errors. This method

13   appears to have a positive impact on intensity error dominated cases. However, it has a

14   neutral or negative impact on pattern error dominated cases (Figure 12), causing a poorer

15   sampling of bias information. In this case bias is reduced at the expense of an increase in

16   random error. Further investigation is needed to fully understand the performance of this

17   method and to determine where and when it has a significantly positive impact and the

18   usefulness of the calibrated products.

19        In this study the decaying average weight is constantly selected as 0.02 (except

20   for the 2000-2001 application) for all lead times. Actually the decaying average weights

21   really depend on experiments which could range from 0.01 to 0.5. In general, the weight

22   is varied for different forecast lead times; a larger weight is good for short lead times

23   which can catch up quick moving systems and a smaller weight is more favorable for

1  long lead time forecasts (not shown). Therefore, choosing an optimum weight for each

2  lead time could be a constructive way to improve the calibration system in the future.

3  Meanwhile, the weight is varied for geographical locations and seasons. There are two

4  improvements we are expecting to validate through future study. One is an optimum

5  weight, which will need large samples for experiments. The weights should be a function

6  of lead time, location and season. The second is a down-scaling process to produce a

7  much finer resolution forecast (5km and 2.5km resolutions).

8

16 **References:**

17
18 Cui, B., Z. Toth, Y. Zhu and D. Hou, 2012: Bias Correction for Global Ensemble

19     Forecast. *Weather and Forecasting* Vol. 27 396-410

20 Cui, B., Y. Zhu, Z. Toth and D. Hou, 2013: Statistical Post Process for NAEFS. Summited

21     to *Weather and Forecasting.*

22 Davis, C.A., B.G. Brown, and R.G. Bullock, 2006a: Object-based verification of

23     precipitation forecasts, Part I: Methodology and application to mesoscale rain

24     areas. *Monthly Weather Review*, 134, 1772-1784.

1  Davis, C.A., B.G. Brown, and R.G. Bullock, 2006b: Object-based verification of

2     precipitation forecasts, Part II: Application to convective rain systems. *Monthly*

3     *Weather Review*, 134, 1785-1795.

4  Demargne, J., L. Wu, S. Regonda, J. Brown, H. Lee M. He, D-J Seo, R. Hartman, M.

5     Fresch, and Y. Zhu, 2013: The Science of NOAA's Operational Hydrologic

6     Ensemble Forecast Service. Submitted to *Bulletin of American Meteorological*

7     *Society.*

8  Eckel, F. A., and M. K. Walters, 1998: Calibrated Probabilistic Quantitative Precipitation

9     Forecasts Based on the MRF Ensemble, *Weather and Forecasting,* 13, 1132-1147

10 Fundel, F., A. Walser, M. A. Liniger, C. Frei, and C. Appenzeller, 2009: Calibrated

11    Precipitation Forecasts for a Limited-Area Ensemble Forecast System Using

12    Reforecasts, *Monthly Weather Review,* 138, 176-189

13 Hamill, T. M., and J. S. Whitaker, 2006: Probabilistic quantitative precipitation forecasts

14    based on reforecast analogs: theory and application, *Monthly Weather Review*, 134,

15    3209-3229

16 Hamill, T. M., R. Hagedorn, and J. S. Whitaker, 2008: Probabilistic Forecast Calibration

17    Using ECMWF and GFS Ensemble Reforecasts. Part II: Precipitation, *Monthly*

18    *Weather Review,* 136, 2620-2632

19 Hou, D., M. Charles, Y. Luo, Z. Toth, Y. Zhu, R. Krzysztofowicz, Y. Lin, P. Xie, D. J.

20    Seo, M. Pena and B. Cui, 2013: Climatology-Calibrated Precipitation Analysis at

21    Fine Scales: Statistical Adjustment of STAGE IV towards CPC Gauge-Based

22    Analysis, *Journal of Hydrometeorology* (in press).

1    Krzysztofowicz, R. and A. A. Sigrest, 1999: Calibration of Probabilistic Quantitative

2        Precipitation Forecasts, *Weather and Forecasting,* 14, 427-442

3    Lin, Y., and K. E. Mitchell, 2005: The NCEP Stage II/IV Hourly Precipitation Analyses:

4        Development and Applications. Preprints, *19th Conf. on Hydrology,* San Diego,

5        CA, *Amer. Meteor. Soc.,* 1.2. [Available online at

6        http://ams.confex.com/ams/pdfpapers/83847.pdf]

7    Christopher, C. P., A. T. Ferro, I. T. Jolliffe, and D. B. Stephenson, 2008: Calibration of

8        Probabilistic Forecasts of Binary Events, *Monthly Weather Review,* 137 1142-1149

9    Voisin N., J. C. Schaake, and D. P. Lettenmaier, 2010: Calibration and Downscaling

10       Methods for Quantitative Ensemble Precipitation Forecasts. *Weather and*

11       *Forecasting,* 25, 1603-1627.

12    Yuan, H., X. Gao, S. L. Mullen, S. Sorooshian, J. Du, and H. H. Juang, 2007: Calibration

13       of Probabilistic Quantitative Precipitation Forecasts with an Artificial Neural

14       Network. *Weather and  Forecasting,* 22, 1287-1303.

15    Xie P, M. Chen M and W. Shi, 2013: CPC unified gauge analysis of global daily

16       precipitation. (to be submitted )

17    Wei, M., Z. Toth, R. Wobus, and Y. Zhu, C. H. Bishop, X. Wang, 2006: Ensemble

18       Transform Kalman Filter-based ensemble perturbations in an operational global

19       prediction system at NCEP. *Tellus* 58A, 28-44

20    Wei, M., Z. Toth, R. Wobus, and Y. Zhu, 2008: Initial Perturbations Based on the

21       Ensemble Transform (ET) Technique in the NCEP Global Operational Forecast

22       System. Tellus 59A, 62-79

1    Zhu, Y., Z. Toth, E. Kalnay, and S. Tracton, 1998: Probabilistic Quantitative Precipitation

2           Forecasts based on the NCEP Global Ensemble. Preprints, 12th Conf. on

3           Numerical Weather Prediction, Phoenix, AZ, *Amer. Meteor. Soc.,* 286–289.

4    Zhu, Y. and Z. Toth, 1999: Calibration of Probabilistic Quantitative Precipitation

5           Forecast, Preprints of the 17th AMS Conference on Weather Analysis and

6           Forecasting 13-17 September 1999, Denver, CO, *Amer. Meteor. Soc.*

7    Zhu, Y. and Z. Toth, 2004: May 2004 Implementation of a QPF Bias-Correction

8           Algorithm, [ Available online at:

9           http://www.emc.ncep.noaa.gov/gmb/ens/ens_imp_news.html ]

10   Zhu, Y., 2005: Ensemble forecast: A New Approach to Uncertainty and Predictability,

11          *Advance in Atmospheric Sciences,* Vol. 22, No. 6, 781-788

12   Zhu, Y. and Z. Toth, 2008: Ensemble based Probabilistic Forecast Verification, Preprints,

13          19th Conference on Predictability and Statistics, 20-24 January 2008, New

14          Orleans, Louisiana, *Amer. Meteor. Soc.*

15   Zhu, Y., D. Hou, M. Wei, R. Wobus, J. Ma, B. Cui and S. Moorthi, 2012: [ Available

16          online at: http://www.emc.ncep.noaa.gov/gmb/yzhu/html/imp/201109_imp.html ]

17

18

19

20

21

22

23

1

2

3

4    Table 1. Total grid point counts of CONUS and each RFC at 1 degree spacing.

| index | RFC | count |
|-------|-------|-------|
| *1* | CNRFC | 67 |
| *2* | CBRFC | 83 |
| *3* | MBRFC | 152 |
| *4* | ABRFC | 56 |
| *5* | WGRFC | 75 |
| *6* | NCRFC | 105 |
| *7* | LMRFC | 51 |
| *8* | OHRFC | 47 |
| *9* | NERFC | 31 |
| *10* | MARFC | 21 |
| *11* | SERFC | 61 |
| *12* | NWRFC | 95 |
| *Total* | CONUS | 844 |

5

6

7

DECAYING AVERAGE WEIGHTING

1

2  Figure 1. Decaying averaged weight as a function of preceding days (weighting function
3  for decaying average of preceding days). The dashed curve denotes a weight beginning
4  with a maximum value of 0.01 at day 0. The solid curve denotes a weight beginning with
5  a maximum value of 0.02 and the dotted curve denotes a weight beginning with a
6  maximum value of 0.03 at day 0. All curves gradually approach  zero depending how far
7  away the preceding days are from day 0. The larger the weight at day 0 the faster the
8  decaying speed, which indicates greater weight on the most recent data and less on the
9  oldest data.
10

11

12

13

14

15
16

1
2
3 Figure 2. The domains of the twelve River Forecast Centers (RFC). Note that the CCPA
4 analysis covers the twelve RFCs over the Contiguous United States (CONUS).
5

6

7

8

9

10

11

12

13

14

15

**Precipitation Distribution**



1

Figure 3. Schematic of the frequency matching algorithm demonstrated as precipitation
distributions normalized by observation frequency varying with threshold. The dashed
line is for observed precipitation and the solid line is for forecast precipitation. See text
for details.

6

7

8

9

4 (a)



North America
00Z01DEC2000 − 00Z28FEB2001
24 hrs avg (Threshold  >= 2.0 mm/24 hrs)

4(b)



Reliability Diagram ( fhr  36− 60 )

2

Figure 4. Examples from the 2004 implementation. Results were selected for the period
of 1 December 2000 - 28 Feburary 2001. (a) Averaged Equitable Threat Score (ETS) and
bias scores of the raw GFS (mrf) and GEFS control (ctl) forecasts and their calibrated
forecasts at a threshold of 2.0mm/day. (b) Reliability of the 2.5mm/day GEFS raw (ens,
red) and calibrated (ens_br, blue) forecasts at 36-60 hour lead time. The inset histogram
denotes the frequency of forecast usage of each probability bin.
9

Figure 5. Comparisons of 6 hourly accumulated precipitation (mm) initialized at 0000 UTC 24 January 2010 from the raw GFS forecasts (left column) and calibrated forecasts (middle column) against the CCPA product (right column) that are valid at corresponding time periods.

Ens Prob of Precip Amount Exceeding 0.01 inch (0.254 mm/6hrs)
Ini: 2010012400



Figure 6. GEFS probabilities (left), calibrated probabilities (middle) of the 6-hr
precipitation amount exceeding 0.01 inch initialized at 00 UTC 24 January 2010, and
CCPA precipitation estimates (right) for 6-h precipitation that are valid at the
corresponding time periods.

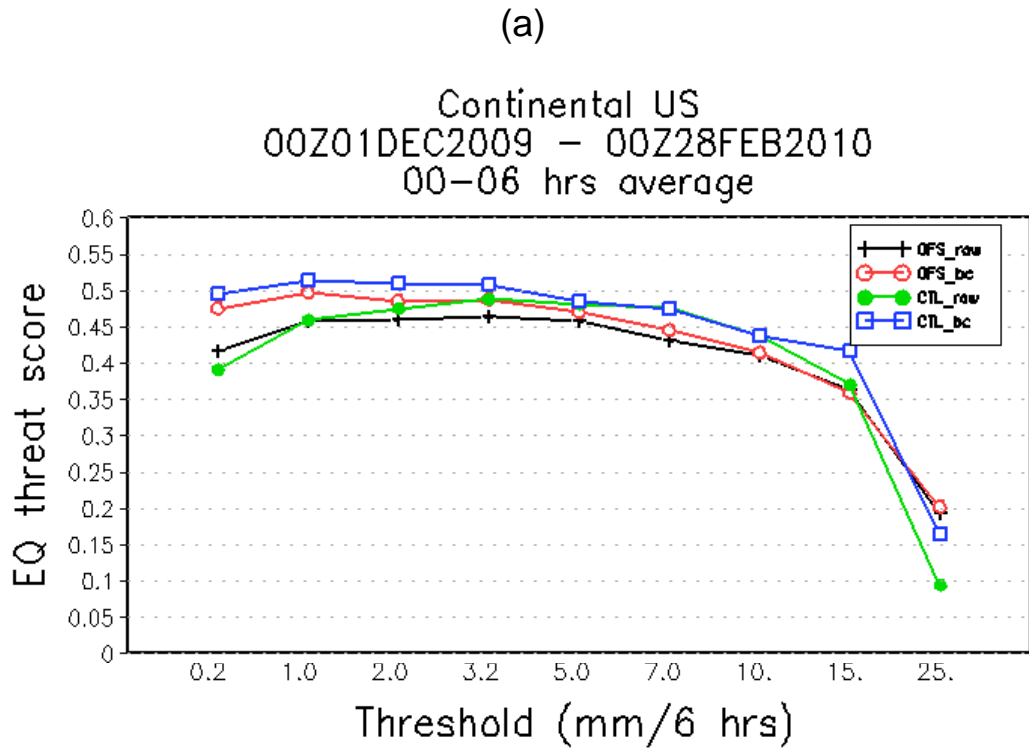1                                      (a)
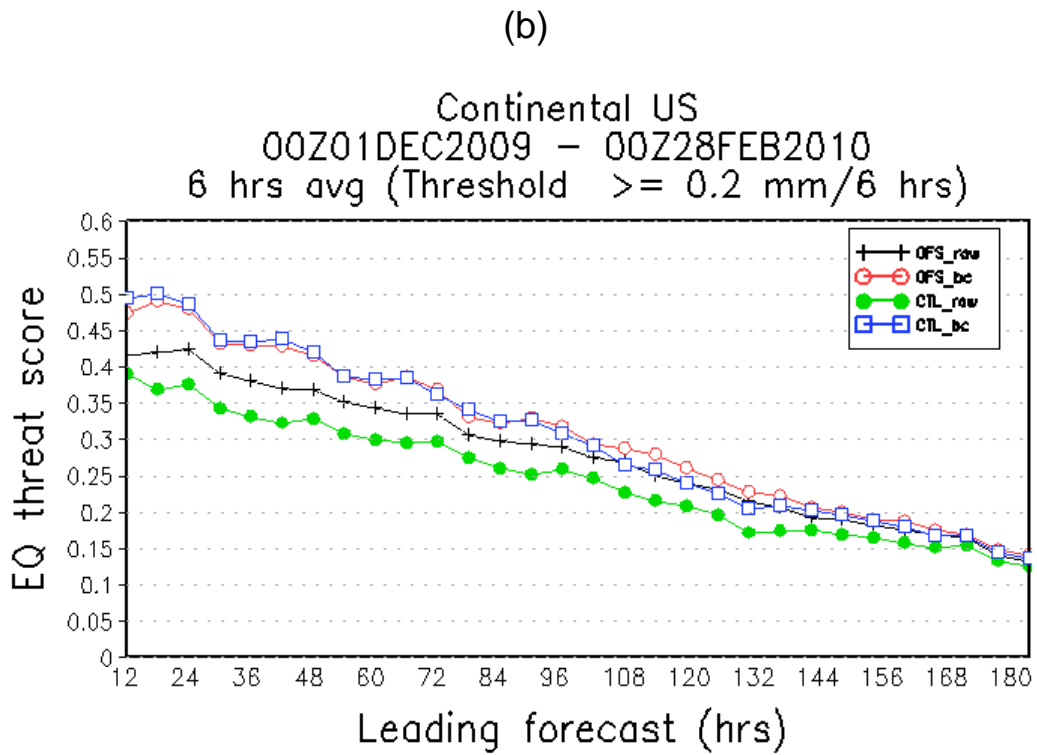


2
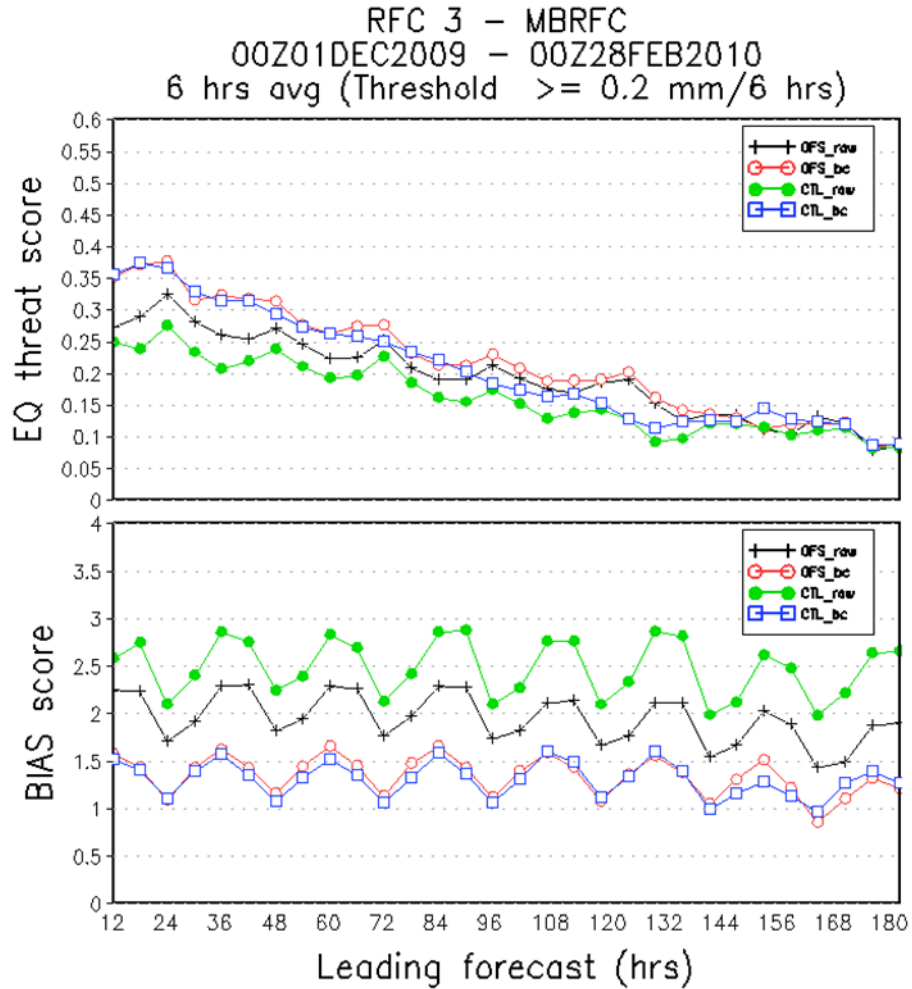3                                      (b)
4



5
6
7   Figure 7. Bias scores of raw forecasts (GFS-black; GEFS/CTL-green) and calibrated
8   forecasts (GFS-red; GEFS/CTL-blue) with increasing lead times for 6-h precipitation
9   averaged between 1 December 2009 and 28 February 2010 (a) as a function of threshold

1  and (b) at a 0.2mm threshold. The numbers in the plot above x-axis are total number of
2  cases verified.
3
4                                        (a)



Continental US
00Z01DEC2009 – 00Z28FEB2010
00–06 hrs average

5
6
7                                        (b)
8



Continental US
00Z01DEC2009 – 00Z28FEB2010
6 hrs avg (Threshold  >= 0.2 mm/6 hrs)

9

1   Figure 8. Same as Figure 7, but for Equitable Threat Scores (ETS).
2
3
4
5



RFC 3 — MBRFC
00Z01DEC2009 — 00Z28FEB2010
6 hrs avg (Threshold >= 0.2 mm/6 hrs)

6
7
8   Figure 9. Comparison of raw forecasts (GFS-black; GEFS/CTL-green) and calibrated
9   forecasts (GFS-red; GEFS/CTL-blue) with increasing lead times for 6-h precipitation
10  averaged between 1 December 2009 and 28 February 2010 and analyzed for the MBRFC
11  region for (a) Equitable Threat Score and (b) Bias score at a 0.2mm threshold.
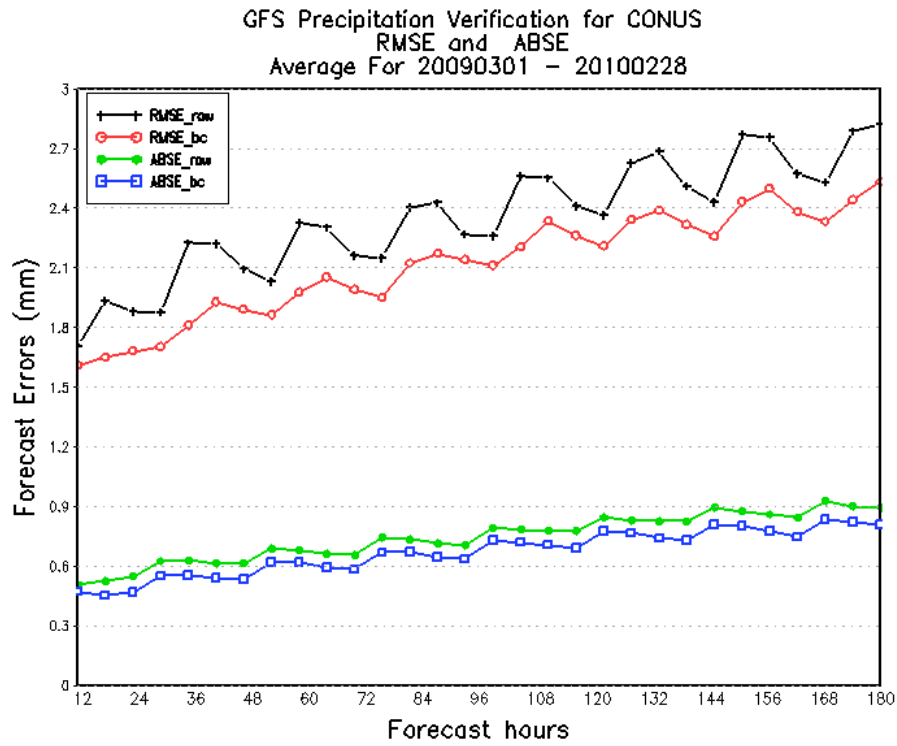12
13
14
15
16
17

RFC 9 — NERFC
00Z01DEC2009 — 00Z28FEB2010
6 hrs avg (Threshold  >= 0.2 mm/6 hrs)

1
2      Figure 10. Same as Figure 9, but for the NERFC region.
3

1



2
3
4
5  Figure 11. RMSE with increasing lead times for 6-h precipitation from the GFS high
6  resolution raw forecasts (black) and calibrated forecasts (red); and ABSE with increasing
7  lead times for 6-h precipitation from the GFS raw forecasts (green) and calibrated
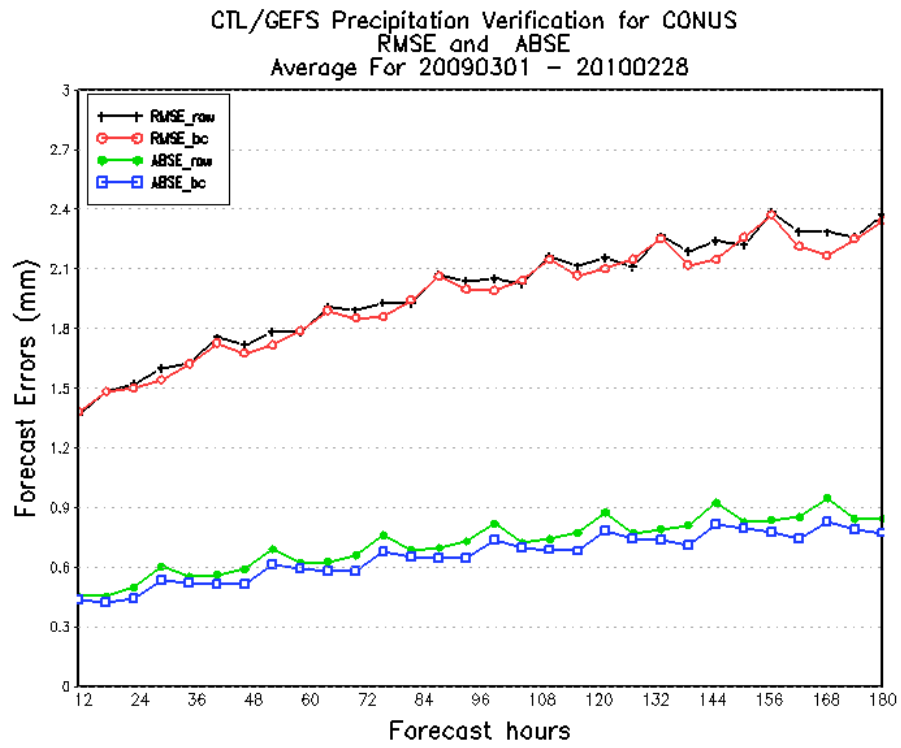8  forecasts (blue).
9

Figure 12. Same as Figure 11, but for GEFS/CTL forecasts.

Ensemble Precipitation Verification for CONUS
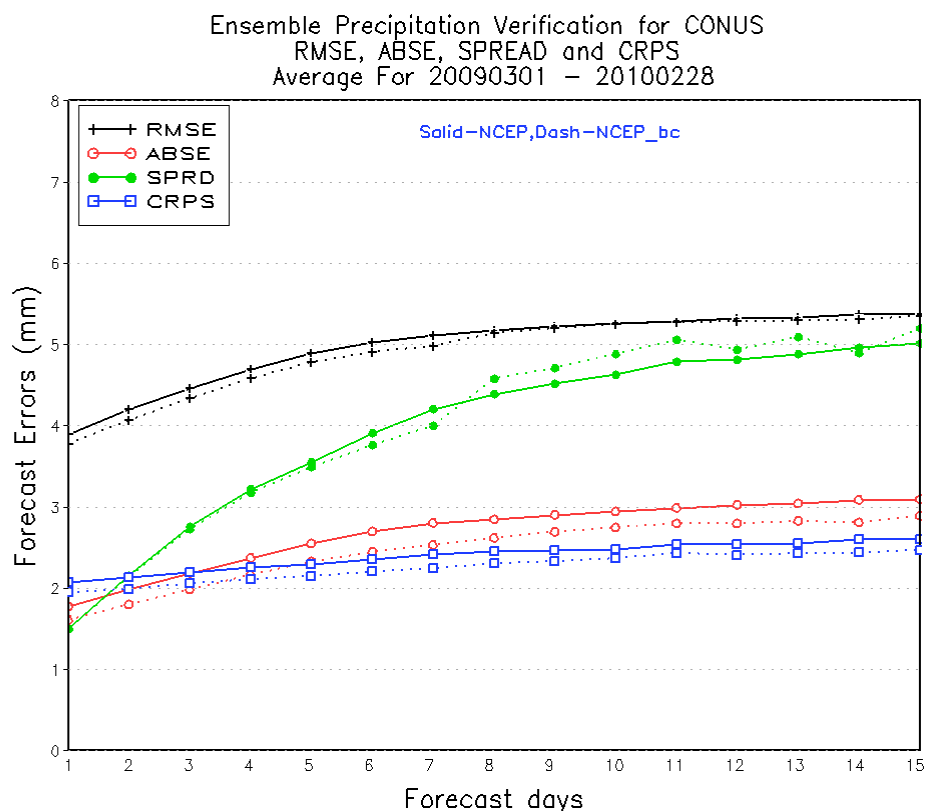RMSE, ABSE, SPREAD and CRPS
Average For 20090301 – 20100228

Figure 13. The RMSE (black), ABSE(red), SPRD (green) and CRPS (blue) with increasing lead times for 24-h precipitation from the GEFS ensemble mean (RMSE and ABSE) and ensemble members (spread and CRPS) raw forecasts (solid) and calibrated forecasts (dotted).
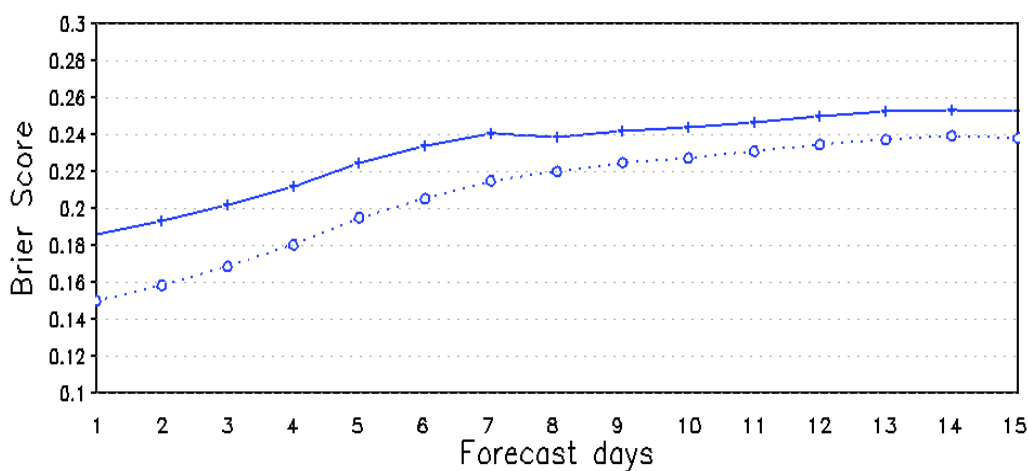


Figure 14. The Brier score at a 0.2mm threshold with increasing lead times for 24-h precipitation from the GEFS ensemble mean raw forecast (solid) and calibrated forecasts (dotted).