

The Developmental Testbed Center Objective Evaluation Performed During the Hazardous Weather Testbed 2010 Spring Experiment.

Tara Jensen^{1*}, Steve Weiss², Jack Kain³, Michelle Harrold¹, Ming Xue⁴, Fanyou Kong⁴, Barb Brown¹, Patrick Marsh³, Adam Clark³, Kevin Thomas⁴, Mike Coniglio³, and Russ Schneider²



¹ NCAR/Research Applications Laboratory (RAL), Boulder, Colorado

² NOAA/Storm Prediction Center (SPC), Norman, Oklahoma

³ NOAA/National Severe Storms Laboratory (NSSL), Norman, Oklahoma

⁴ Center for Analysis and Prediction of Storms (CAPS), University of Oklahoma, Norman, Oklahoma

*Presenting Author E-mail: jensen@ucar.edu



NCAR



Abstract: The DTC objective evaluation during the 2010 HWT Spring Experiment complements the subjective evaluation that has traditionally taken place. With the addition of probabilistic verification capabilities in the DTC's Model Evaluation Tools (MET), both probabilistic products and deterministic forecasts will be evaluated this year. In addition to the severe convective weather component, the 2010 Spring Experiment objective evaluation plan includes evaluation of forecasts from WRF convection-allowing models for extreme precipitation events as well as aviation related thunderstorm indicators. This year's Spring Experiment ran from May 17 – June 18, 2010.

Approach

This year DTC evaluated (in near real-time):

- CAPS SSEF 4 km ensemble members as deterministic run (26 multi-model)
- CAPS 1km deterministic
- CAPS SSEF 4 km ensemble products (15 radar assimilation models used)
- NOAA/ESRL HRRR 3 km deterministic model
- NOAA/EMC NAM 12 km deterministic model (*baseline*)
- NOAA/EMC SREF 32-35 km ensemble products (21 member multi-model) (*baseline*)

•Observations: NSSL NMQ Q2 QPF and Reflectivity products

A full description of the each model contributed to HWT may be found on their website at: http://hwt.nssl.noaa.gov/Spring_2010/. Objective evaluation results are available at <http://verif.rap.ucar.edu/eval/hwt/2010/>.

Variables Evaluated :

FCST Field	Observation (NSSL Q2 fields)	Traditional (MET/Grid-Stat)	Spatial (MET/MODE)	Models
Prob of Exceed (0.5", 1", 2" over 3 and 6 hrs)	0.5", 1", 2" QPE over 3 and 6 hrs	Brier Score, Decomp of Brier score, Area under ROC	None	Ensemble products from CAPS and SREF
50 th Prob of Exceed (0.5", 1", 2" over 3 and 6 hrs)	0.5", 1", 2" QPE over 3 and 6 hrs	None	MMI, Intersection Area, Area Ratio, Centroid Distance, Angle Difference, % Objects and Area Matched, 50 th and 90 th percentile	Ensemble products from CAPS and SREF
QPF (0.25", 0.5", 1.0", 2" over 3 and 6 hrs)	0.25", 0.5", 1.0", 2" QPE over 3 and 6 hrs	GSS, CSI, FAR, PODY, FBAIS	See above	CAPS members, CAPS ens mean, SREF ens mean, HRRR, NAM
Sim. Comp. Refl (20,30,40,50 dBZ)	Composite refl (20,30,40,50 dBZ)	GSS, CSI, FAR, PODY, FBAIS	See above	CAPS members, CAPS ens products, HRRR, NAM
18 dBZ Echo Top (18, 25, 30, 35, 40, 45 kft)	18dBZ Echo Top (18, 25, 30, 35, 40, 45 kft)	GSS, CSI, FAR, PODY, FBAIS	See above	CAPS members, CAPS ens products, HRRR
Prob of 40dBZ echos	Composite reflectivity (40dBZ)	Brier Score, Decomp of Brier score, Area under ROC	None	Ensemble products from CAPS and SREF
50% Prob of 40dBZ echos	Composite reflectivity (40dBZ)	None	See above	Ensemble products from CAPS

Grid: All models were re-gridded to the 4 km Stage IV grid configuration.

Domains: The 00 UTC models were evaluated over three regions: the entire domain, the static VORTEX-2 domain provided by CAPS at the 12 UTC initialization time, and a regional, movable area-of-interest domain selected by HWT Spring Experiment participants each day. Figure 1 depicts examples of these domains.



Figure 1. Examples of the three evaluation domains used in the DTC objective evaluation of HWT Spring Experiment models. Left: Full domain represents two-thirds CONUS; Middle: VORTEX-2 domain; and Lower Right: regional area-of-interest domain that moves daily; referred to as daily domain.

Details:
DOMAIN: VORTEX-2
Field: 20 dBZ reflectivity (REFC)
CAPS SSEF Ensemble PM Mean
CAPS SSEF 1 km Model
CAPS SSEF ARW-CN (control w/o radar assimilation)
CAPS SSEF ARW-CO (control w/o radar assimilation)
HRRR
NAM
PM Mean = Probability Matching used to compute mean REFC (reference: Ebert (2000), MWR)
Statistics aggregated over entire 5 week experiment

Preliminary Results: Severe Storms

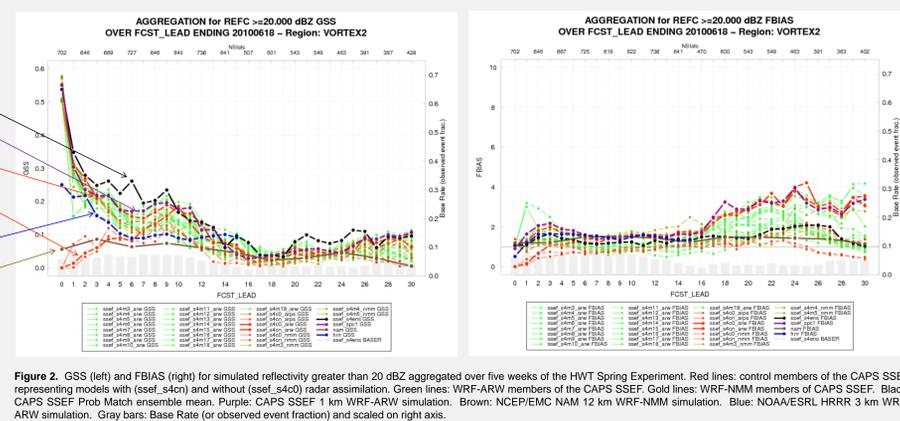
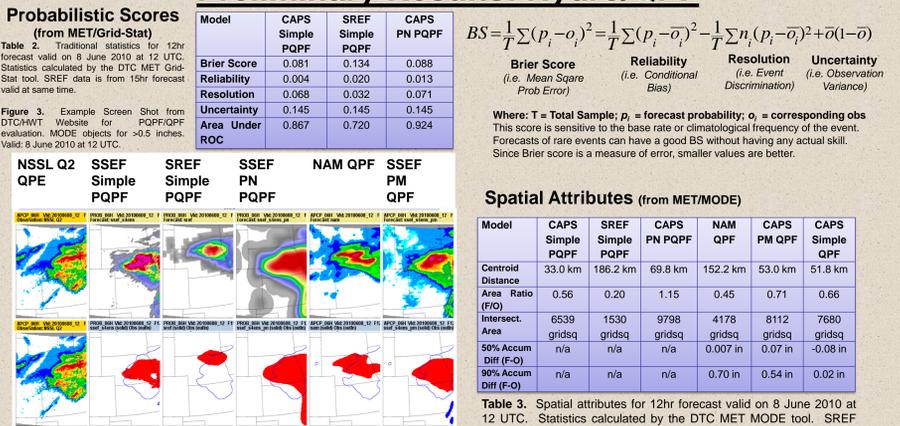


Figure 2. GSS (left) and FBAIS (right) for simulated reflectivity greater than 20 dBZ aggregated over five weeks of the HWT Spring Experiment. Red lines: control members of the CAPS SSEF representing models with (sref_s4cn) and without (sref_s4co) radar assimilation. Green lines: WRF-ARW members of the CAPS SSEF. Gold lines: WRF-NMM members of CAPS SSEF. Black: CAPS SSEF Prob Match ensemble mean. Purple: CAPS SSEF 1 km WRF-ARW simulation. Brown: NCEP/EMC NAM 12 km WRF-NMM simulation. Blue: NOAA/ESRL HRRR 3 km WRF-ARW simulation. Gray bars: Base Rate (or observed event fraction) and scaled on right axis.

Discussion:
 Based on Gilbert Skill Score, (an assessment of how well the forecast "yes" events correspond to observed "yes" events while taking into account possibility of a random chance "yes")
 CAPS SSEF PM Mean and 1 km Deterministic runs appear to provide advantage over baselines HRRR and NAM.
 Radar assimilation appears to provide advantage over no assimilation during 0-6hr fcst
 However some of that advantage may be due to a positive FBAIS, which is
 forecast events observed events
 Increased Forecast Events improves potential for a Hit

Details:
DOMAIN: DAILY
Fields:
 Quant. Precip Est. (QPE)
 Quant. Precip. Fcst. (QPF)
 Probability QPF (PQPF)
 MODE Objects for >0.5 inch
Simple PQPF = Ensemble Relative Frequency at Gridpnt
PN PQPF = Neighborhood method applied before computing Ens. Rel. Freq.
PM QPF = Probability Matching used to compute QPF intensity
Statistics for Single Valid Time 8 Jun 2010 – 12 UTC

Preliminary Results: Hydro/QPF



Discussion:
 Based on the Brier Score SREF Simple and PN PQPF appear to have more skill than SREF Simple PQPF.
 Based on Reliability and Resolution The SREF PN PQPF shows slightly more ability to resolve events but has slightly less reliability.
 Based on Area Under the Curve With a higher Area under the Receiver Operating Characteristic Curve SREF PN QPN may be deemed the slightly more skillful for this application.
 Based on Centroid Distance, Area Ratio, Intersection Area SREF Ensemble Products appear to provide better QPF location guidance than SREF and NAM for this case.
 Based on 50th Percentile and 90th Percentile Intensity Differences SREF Simple Ensemble Mean provided the best guidance for this case.
 Aggregations of these individual case products will become available on the website near the end of the summer.

Details:
DOMAIN: DAILY
Fields:
 Observed Reflectivity (REFC)
 Forecasted Simulated REFC
 Observed 18 dBZ Radar Echo Top (RETOP) height
 Forecasted RETOP height
 MODE Objects for >25,000 ft
Microphysics Schemes:
 Thompson
 WRF Single Moment (WSM)
 WRF Double Moment (WDM)
 Morrison
 All Models in example have radar assimilation in Initial Conditions
PM RETOP = Probability Matching used to compute QPF intensity
Statistics aggregated over entire 5 week experiment

Preliminary Results: Aviation Weather

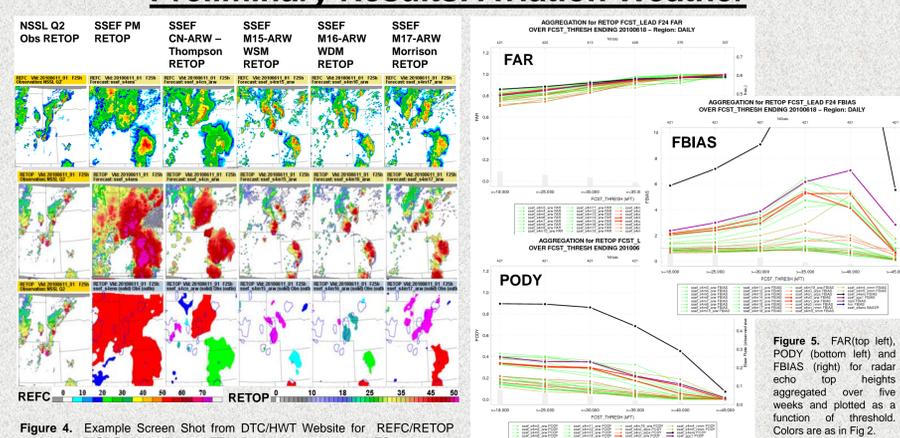


Figure 4. Example Screen Shot from DTC/HWT Website for REFC/RETOP evaluation. MODE objects for >25KFT Valid: 11 June 2010 at 01 UTC.

Discussion:
 Based on Visual Inspection (of Fields and MODE Objects to left) SREF Probability Matched Mean appears to over-estimate RETOP spatial coverage
 Thompson scheme members appears to be significantly over-predicting the stratiform cloud shield (green in REFC over predicted leading to much larger area for RETOP also).
 Thompson members appear to dominate Probability Matching (example not shown here – see attached plot)
 Other Members seem to have better Frequency Bias but are over developing southern convective region. This may be a timing error.
 Based on FAR (top middle) SREF Probability Matched Mean appears to have high FAR
 Based on FBAIS (bottom middle) SREF Probability Matched Mean appears to have significantly Probability of Detection but this is most likely due to high FBAIS
 Based on FBAIS (right) SREF Probability Matched Mean appears to over-estimate RETOP spatial coverage

Summary

The DTC/HWT Objective Evaluation goals include are to 1) to provide objective evaluations of the experimental forecasts; 2) to supplement and compare to subjective assessments of performance; and 3) to expose the forecasters and researchers to both new and traditional approaches for evaluating forecasts. In total, 30 models and 3 ensemble product methods were evaluated in near real-time for the 5 week HWT 2010 Spring Experiment. Three forecast challenges were addressed: Severe Weather, QPF, and Aviation Weather. The DTC evaluated 1-2 variables for each area.

Some preliminary objective results:

- Models using radar assimilation methods exhibited improved skill during 0-6 hr lead times
- CAPS Convective Allowing Models with radar assimilation generally exhibited improved skill over operational baselines (such as NAM, HRRR, and SREF) but this improvements may be a by-product of increased frequency bias and hence may result in increase false alarms.

Preliminary observations about Ensemble Products:

- Simple arithmetic mean field only available for QPF field – it showed improved skill over other methods on some days but not all. Further investigation is recommended.
- Probability Matched mean product exhibits higher skill for reflectivity (REFC) and quantitative precipitation forecast (QPF) fields but appears to produce a sizable frequency bias (FBAIS) for radar echo top (RETOP), and hence increased false alarm ratio (FAR) that may make the field unusable
- Probability Neighborhood product was evaluated for QPF only. While the product exhibits a smoother Probability QPF field, which in general appears to increase traditional skill, the tendency for a sizable frequency bias may decrease its utility. The true test is which product was found to be the most useful. Further investigation of subjective logs is needed.

Acknowledgements

The Developmental Testbed Center is funded by the National Oceanic and Atmospheric Administration, Air Force Weather Agency and the National Center for Atmospheric Research. The CAPS research was supported by an allocation of advanced computing resources provided by the National Science Foundation. The computations were performed on Athena (a Cray XT4) at the National Institute for Computational Science (NICS; <http://www.nics.tennessee.edu>).

Dedicated work by many individuals led to the success of SE2010. Paul Oldenburg, John Halley Gotway, Randy Bullock, and Nancy Rehak developed the DTC evaluation system. David Ahijevych, Jamie Wolff, and Isidora Jankov, also from the DTC, led discussions related to forecast verification during SE2010 daily activities. At the SPC, HWT operations were made possible by technical support from Israel Jirak, Chris Melick, and Andy Dean. At the NSSL, Ryan Sobash provided valuable scientific and technical support. David Novak, Faye Barthold (from NOAA/Hydrometeorological Prediction Center), and Jason Levit (from NOAA/Aviation Weather Center) provided much needed guidance on meaningful evaluation criteria.