

The Role in Verification in R20 Testing and Evaluation

Tara L. Jensen¹, and M. P. Mittermaier², T. Fowler¹, B. G. Brown¹, J. Halley Gotway¹, R. Bullock¹, J. K. Wolff¹, M. Harrold¹, E. Tollerud³, L. Nance¹, and Y. H. Kuo¹

1. NCAR/RAL, Developmental Testbed Center (DTC)

2. UK Met Office

3. CIRA, NOAA/ESRL, Developmental Testbed Center (DTC)

95th AMS – Fifth Conference on Research to Operations Transition
4-8 February 2015 Phoenix, AZ



Developmental Testbed Center



NCAR

Why verify?

- Purposes of verification (traditional definition)



- Administrative purpose
 - Monitoring performance
 - Choice of model or model configuration (has the model improved?)
- Scientific purpose
 - Identifying and correcting model flaws
 - Improved forecasts
- Economic purpose
 - Improved decision making
 - “Feeding” decision models or decision support systems



Identifying verification goals

What *questions* do we want to answer?

- Examples:
 - ✓ In what locations does the model have the best performance?
 - ✓ Are there regimes in which the forecasts are better or worse?
 - ✓ Is the probability forecast well calibrated (i.e., reliable)?
 - ✓ Do the forecasts correctly capture the natural variability of the weather?

Basic guide for developing verification studies

Identify multiple *verification attributes* that can provide answers to the questions of interest

Select *measures and graphics* that appropriately measure and represent the attributes of interest

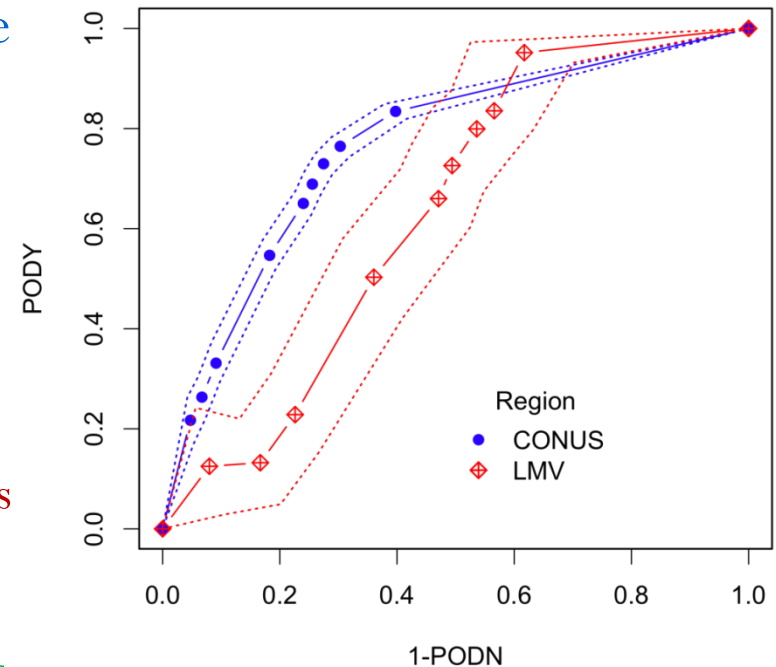
Identify a *standard of comparison* that provides a reference level of skill (e.g., persistence, climatology, old model)



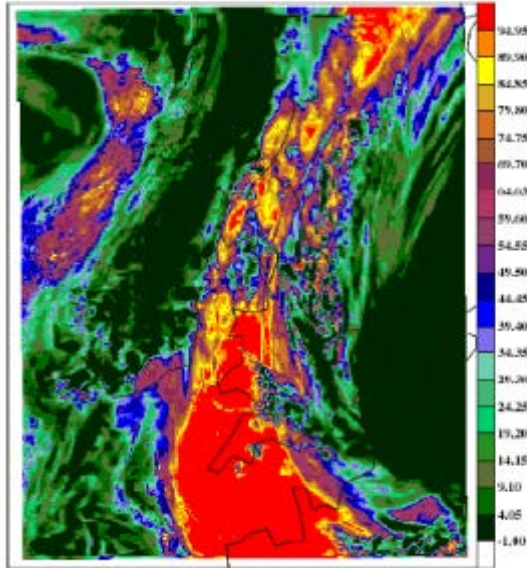
Confidence Intervals

Uncertainty in scores and measures should be estimated whenever possible!

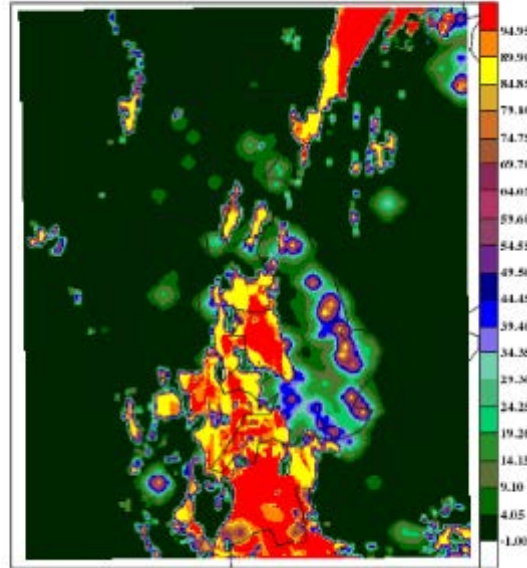
- Uncertainty arises from
 - Sampling variability
 - Observation error
 - Representativeness differences
- Erroneous conclusions can be drawn regarding improvements in forecasting systems and models without CIs
- Methods for *confidence intervals* and *hypothesis tests*
 - Parametric (i.e., depending on a statistical model)
 - Non-parametric (e.g., derived from re-sampling procedures, often called “bootstrapping”)



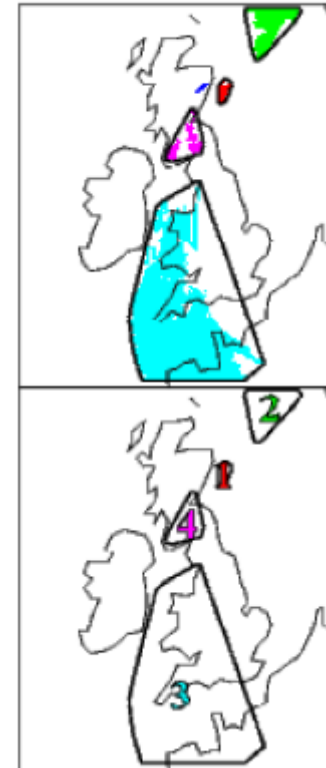
UK4 t+12h @ 1500 UTC



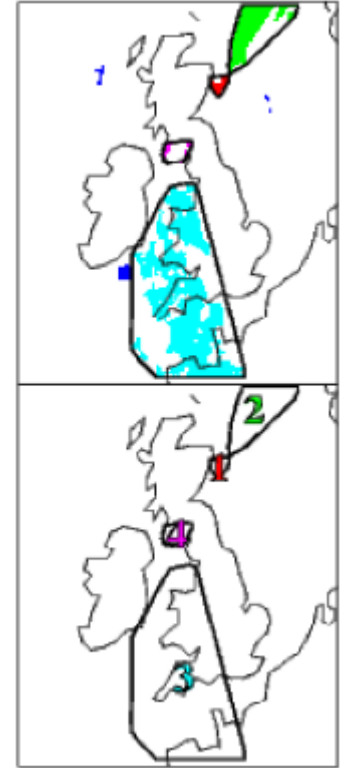
Analysis @ 1500 UTC



Forecast



Observation

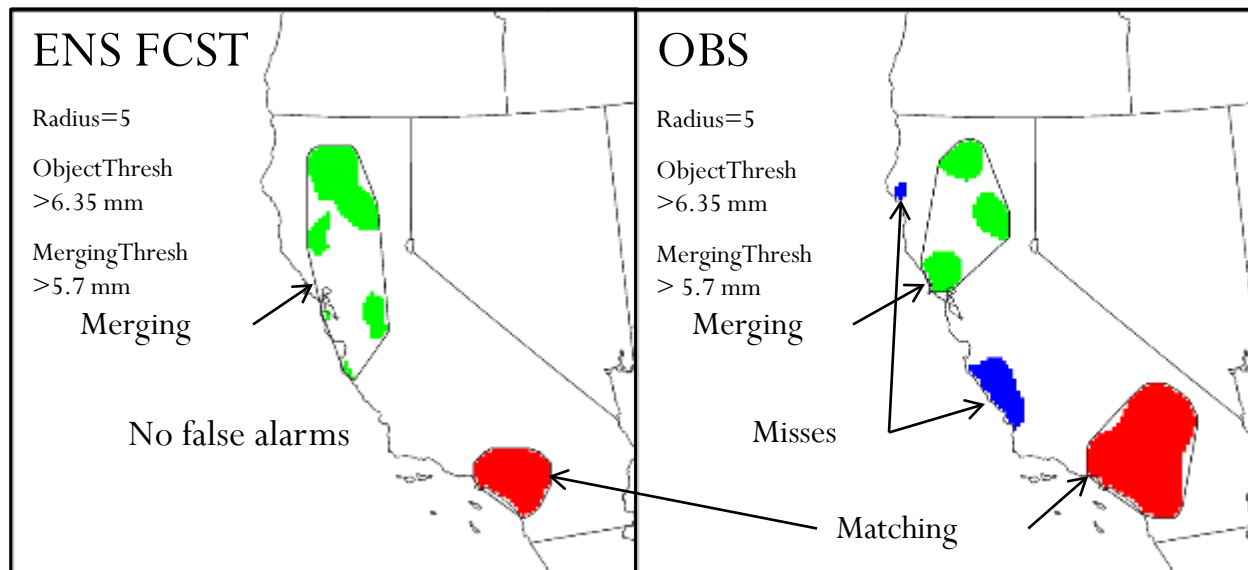
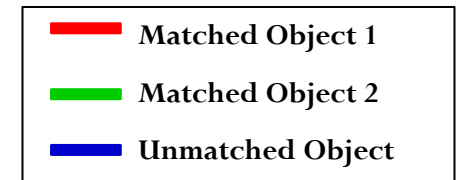


Courtesy of Marion Mittermaier
The Met Office

Spatial Methods

Method for Object-based Diagnostic Evaluation (MODE)

How it works



Comparing objects can tell you things about your forecast like . . .

This:

30% Too Big
(area ratio=1.3)

Shifted west 1 km
(centroid distance = 1km)

Rotated 15°
(angle diff = 15%)

Peak Rain 1/2" too much
(diff in 90th percentile of intensities = 0.5)

Instead of this:

POD = 0.35

FAR = 0.7235

CSI = 0.1587

Example of Advanced Techniques in R20



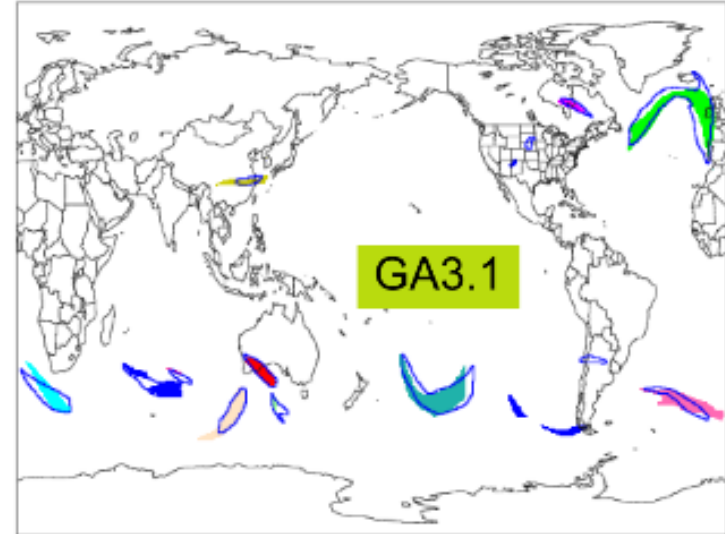
Temporal evolution

- Older N320 trial
250 hPa winds > 60 m/s
at forecast lead time of t+96h
from the 12Z initialisation
compared to EC analyses
- Differences in the size of
forecast and analysed
objects is not overshadowed
by growth of synoptic
forecast error.

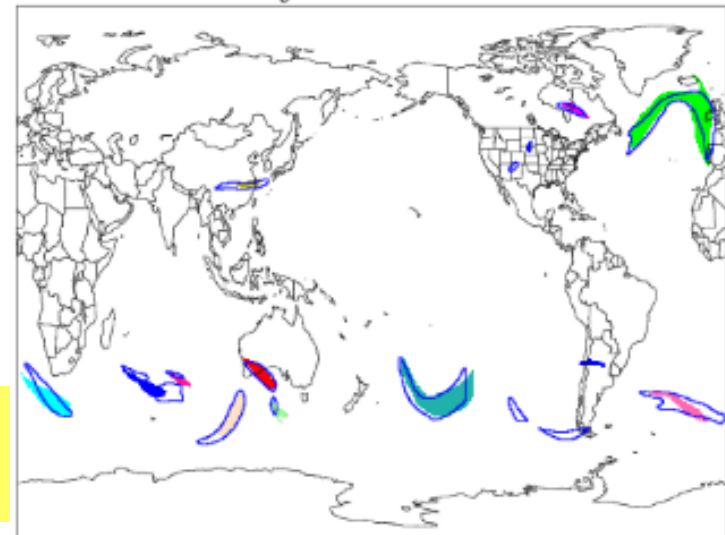
*Slide courtesy of Marion Mittermaier
The Met Office – August 2014*

Model Evaluations Tool (MET) used to
identify objects and synthesize attributes

Forecast Objects with Observation Outlines



Observation Objects with Forecast Outlines



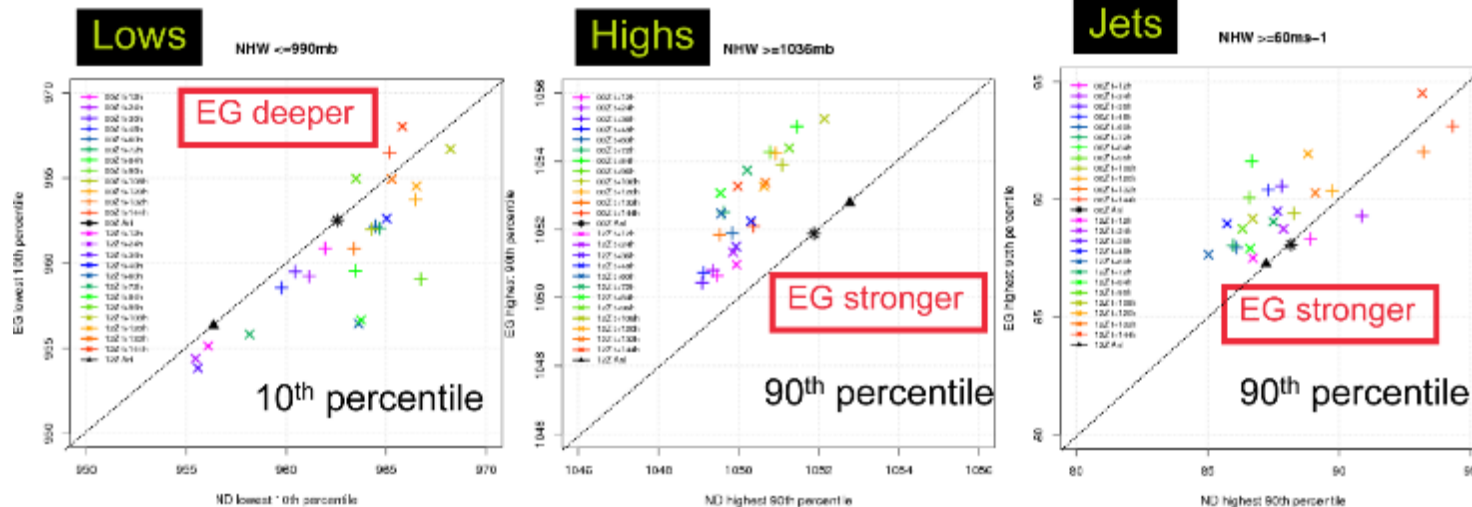
Example of Advanced Techniques in R20



Met Office

Object intensities

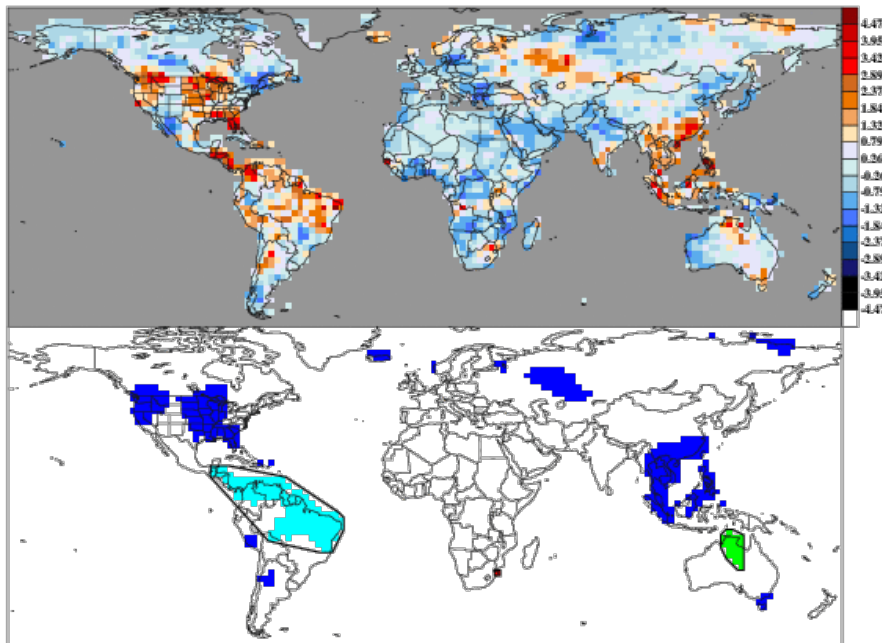
EC analyses
N768 EG v N512 ND



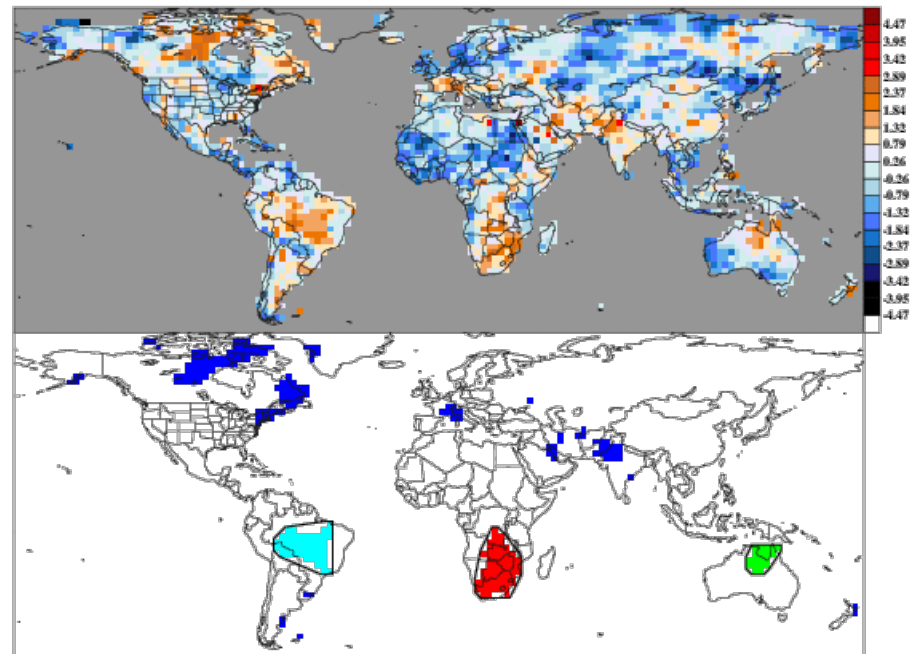
- Do not look at absolute min/max values in objects. Use the **10th or 90th percentile** as a more reliable estimate of how the intensity distribution has shifted/changed.
- Lows are deeper, highs and jets are stronger → **sharper gradients** and a **more active** energetic model.
- Differences in the 00Z and 12Z analyses.

Use of MODE on Climate...

MODE: pdsisc at SFC vs pdsisc at SFC
Forecast



MODE: pdsisc at SFC vs pdsisc at SFC
Observation



Scorecarding and NWP Indexes

AFWA Configuration Testing | DTC

Upper Air Dew Point Temperature		Annual				Summer				Winter			
		f12	f24	f36	f48	f12	f24	f36	f48	f12	f24	f36	f48
BCRMSE	850	RRTMG	RRTMG	RRTMG	RRTMG	--	--	RRTMG	RRTMG	RRTMG	--	RRTMG	--
	700	RRTMG	--	--	--	--	--	--	--	RRTMG	RRTMG	--	--
	500	--	--	--	--	--	--	--	RRTMG	--	--	--	--
Bias	850	AFWA	AFWA	AFWA	AFWA	AFWA	AFWA	AFWA	RRTMG	--	AFWA	AFWA	AFWA
	700	AFWA	AFWA	AFWA	AFWA	RRTMG	RRTMG	RRTMG	RRTMG	AFWA	AFWA	AFWA	AFWA
	500	AFWA	AFWA	AFWA	AFWA	AFWA	AFWA	AFWA	AFWA	AFWA	AFWA	AFWA	AFWA

Model Evaluations Tool (MET) used to calculate the statistics and scripting used to formulate scorecard

Statistical Significance (light shading)

- Differences pass the test

Practical Significance (dark shading)

- Which SS differences are greater than the observation uncertainty

NWP Index

Example – AFWA GO Index – a weighted average of the RMSE values for wind speed, dewpoint temperature, temperature, height, and pressure at several levels in the atmosphere.

Table 10-1. Variables, levels, and weights used to compute the GO Index.

Variable	Level	Weights by lead time			
		12 h	24 h	36 h	48 h
Wind speed	250 hPa	4	3	2	1
	400 hPa	4	3	2	1
	850 hPa	4	3	2	1
	Surface	8	6	4	2
Dewpoint temperature	400 hPa	8	6	4	2
	700 hPa	8	6	4	2
	850 hPa	8	6	4	2
	Surface	8	6	4	2
Temperature	400 hPa	4	3	2	1
	Surface	8	6	4	2
Height	400 hPa	4	3	2	1
Pressure	Mean sea level	8	6	4	2

Model Evaluations Tool (MET) has this capability to compute GO Index but is flexible enough to compute other user-defined NWP Indexes

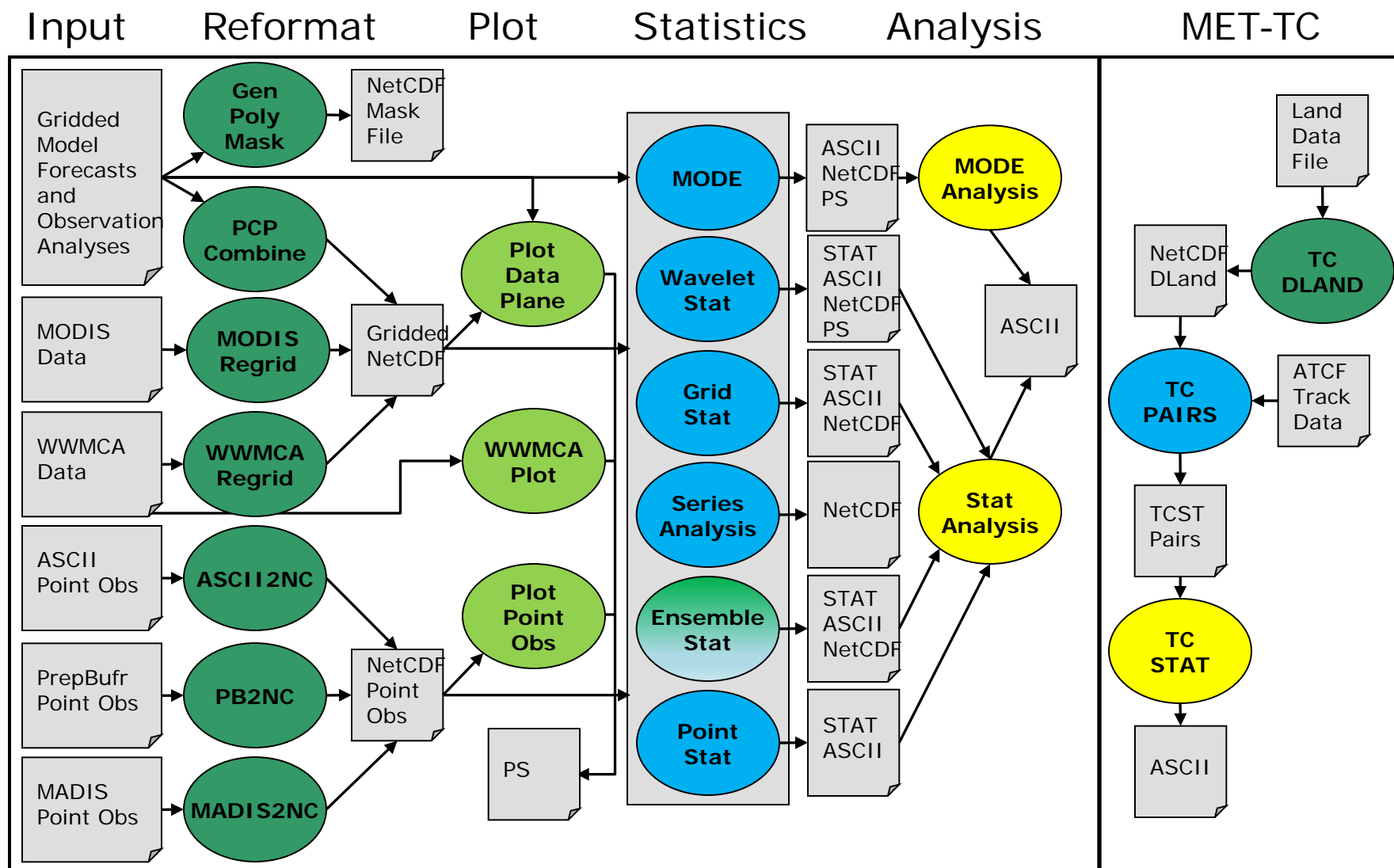


Provides a standard set of tools to facilitate
reliable, consistent verification across institutions

MET Package

- MET is community code supported by DTC that is free to download (registration required)
 - Approximately 2550 registered users
 - 124 countries
 - Universities, Government, Private Companies, Non-Profits
- Download MET release and compile locally.
 - Register and download: www.dtcenter.org/met/users
- Language:
 - Primarily in C++ with calls to some Fortran libraries
- Supported Platforms and Compilers:
 - Linux with GNU compilers
 - Linux with Portland Group (PGI) compilers
 - Linux with Intel compilers
- In-person tutorials given yearly –
NEXT TUTORIAL: FEB 2-3 in Boulder, CO
It's not too late to register! Contact me for how to do so.

MET v5.0 Tools (Release: August 2014)



Details on Categorical and Continuous Statistics

Continuous	Categorical / Multi-Categorical
<p>Forecast Mean</p> <p>Forecast Standard Deviation</p> <p>Observation Mean</p> <p>Observation Standard Deviation</p> <p>Pearson Correlation Coefficient (aka Correlation)</p> <p>Spearman's Rank Correlation</p> <p>Kendall's Tau statistic</p> <p>Number of ranks used in Kendall's tau</p> <p>Number of tied forecasts in Kendall's tau</p> <p>Number of tied observations in Kendall's tau</p> <p>Mean error</p> <p>Standard Deviation of error</p> <p>10th, 25th, 50th, 75th, 90th Percentile of Error</p> <p>Inner Quartile Range</p> <p>Multiplicative Bias (aka Bias)</p> <p>Mean Absolute Error</p> <p>Mean Square Error</p> <p>Bias-corrected Mean Square Error</p> <p>Root Mean Square Error</p> <p>Mean Absolute Deviation</p> <p>24 Statistics</p>	<p>Total number of matched pairs</p> <p>Contingency Table Counts</p> <p>Forecast rate</p> <p>Hit rate</p> <p>Observation rate</p> <p>Base rate</p> <p>Forecast mean</p> <p>Accuracy</p> <p>Frequency Bias</p> <p>Probability of Detection – Yes</p> <p>Probability of Detection – No</p> <p>Probability of False Detection (aka False Alarm Rate)</p> <p>False Alarm Ratio</p> <p>Critical Success Score (aka Threat Score)</p> <p>Gilbert Skill Score (aka Equitable Threat Score)</p> <p>Bias-Adjusted Gilbert Skill Score</p> <p>Odds Ratio</p> <p>Log-Odds Ratio</p> <p>Odds-Ratio Skill Score</p> <p>Hanssen-Kuipers Discriminant</p> <p>Heidke Skill Score</p> <p>Extreme Dependency Score</p> <p>Symmetric Extreme Dependency Score</p> <p>Extreme Dependency Index</p> <p>Symmetric Extreme Dependency Index</p> <p>25 Statistics</p>

**New in
MET v5.0**

RUN:

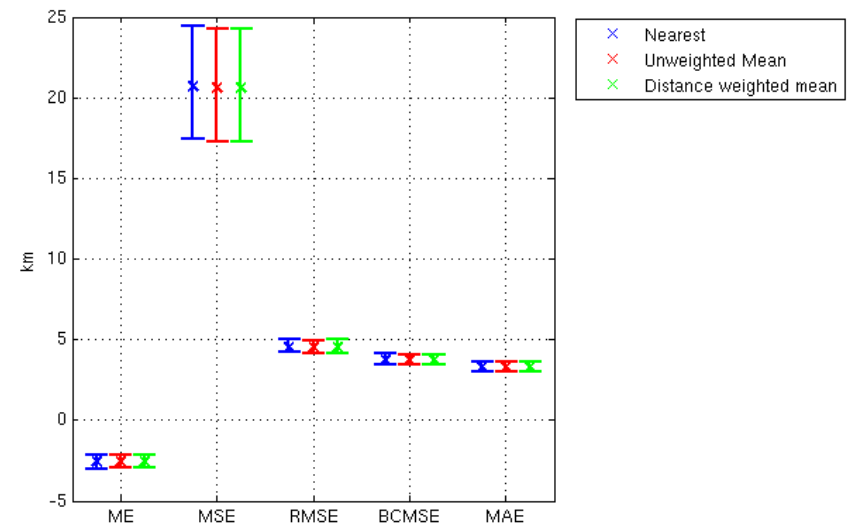
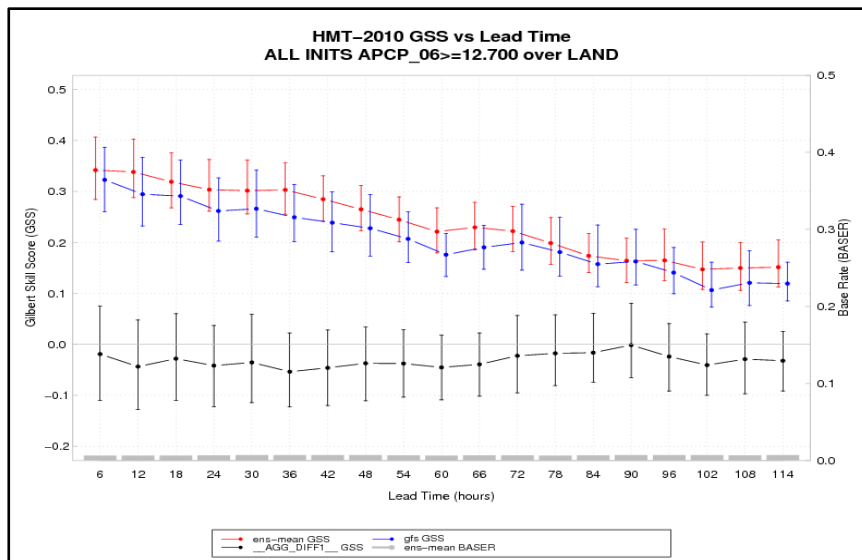
POINT-STAT
OR
GRID-STAT

Neighborhood and Ensemble/Probability Statistics

Neighborhood		Ensemble/Probability	
RUN: GRID-STAT	<p>Neighborhood Contingency Table Statistics (see previous slide)</p> <p>Fractions Brier Score Fraction Skill Score Asymptotic Fractions Skill Score Uniform Fractions Skill Score Forecast Event Frequency Observed Event Frequency</p>	<p>Ensemble Mean and Std Dev fields Ensemble Mean \pm 1 Std Dev fields Ensemble Min and Max fields Ensemble Range field Ensemble Valid Data Count field Ensemble Relative Frequency (probability)</p> <p>Ranked Histograms (if Obs Field Provided) PIT Histogram Ensemble Spread-Skill (if Obs Field Provided)</p> <p>Neighborhood Contingency Table Statistics (see previous slide) Brier Score Reliability Resolution Uncertainty Area Under ROC Calibration Refinement Likelihood Base Rate Probailiby Integral Transform (PIT) Reliability points ROC Curve Points</p>	RUN: ENSEMBLE STAT
	<p>Wavelet Decomposition</p> <p>Mean squared error for each scale Intensity skill score Forecast Energy Squared Observed Energy Squared Base Rate (not scale dependent) Frequency Bias</p>		
RUN: WAVELET- STAT			RUN: POINT-STAT OR GRID-STAT

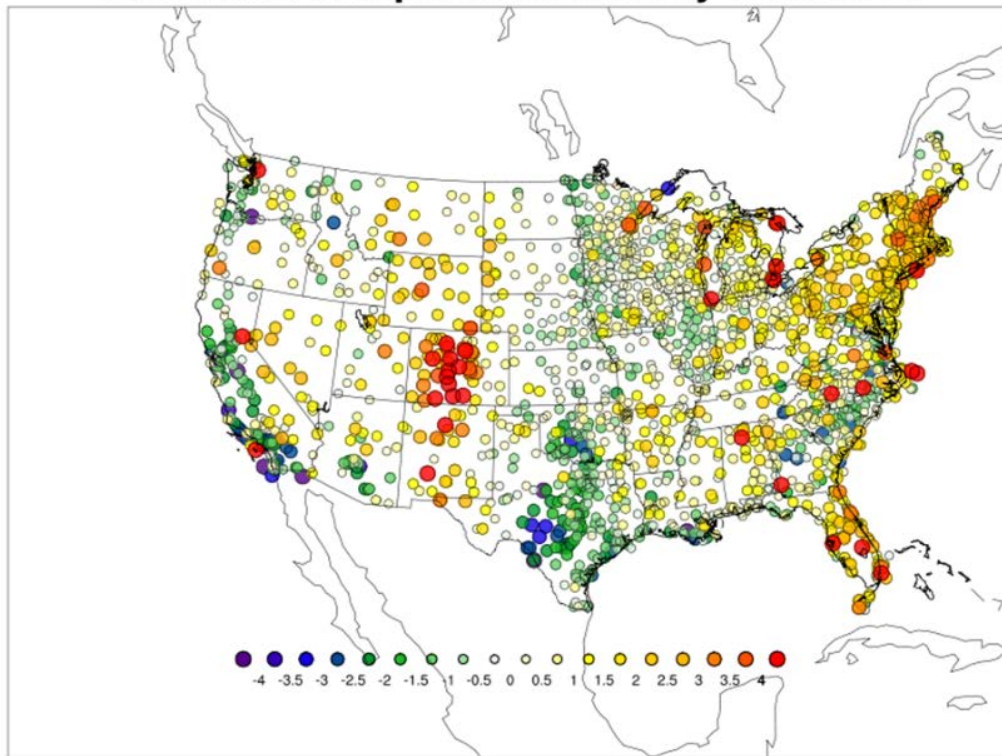
Reason to use MET: Easy use of Confidence Intervals

- Normal Approximation CI
 - Calculated for all statistics for which this is appropriate
- Bootstrapped CI
 - Can be turned on in config file
 - Number of repetitions are user defined
- Interpolation for Point Data – Nearest Neighbor, Unweighted Mean, Distance Weighted Mean



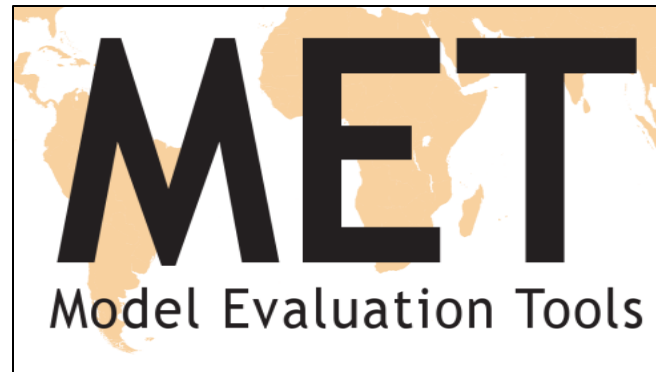
Series Analysis tools for geographic representation of scores

Dew Point Temperature Bias by Station ID



Config=AFWAOC_WRFv3.5 Season=WINTER Init=00UTC Fcst Hr=42h

- Accumulates statistics separately for each grid location over a series
 - Time
 - Height
 - Other series
- Accumulate over
 - Stations
 - Grids



Support for MET has been provided by
AFWA, NOAA and NCAR
through the Developmental Testbed Center (DTC)



NCAR



Developmental Testbed Center

Thank You and Further Information

JNT: <http://www.ral.ucar.edu/jnt>

DTC: <http://www.dtcenter.org>

MET: <http://www.dtcenter.org/met/users>

Email: jensen@ucar.edu
met-help@ucar.edu

Support for the Developmental Testbed Center (DTC),

is provided by

NOAA, AFWA
NCAR and NSF



NCAR

