

Scale sensitivities in model precipitation skill scores

Andrew.Loughe@noaa.gov

Stephen Weygandt
Stan Benjamin

Andy Loughe^{*}
Jennifer Mahoney

NOAA Forecast Systems Laboratory

^{*}CIRES, University of Colorado, Boulder, CO

The problem of verifying convective precipitation forecasts

- **Thunderstorms** produce precipitation patterns with **significant small-scale detail**
- High-resolution **numerical models** are increasingly able to produce **similar small-scale detail**

But....

- Detailed model fields often have **small phase errors** compared to observations
- Traditional **skill scores are often worse** for detailed models even though they produce more realistic forecasts

The present situation

- Realize there is **no single perfect verification score**
- Active research on many new verification approaches
 - **Spatial structures measures**
 - **Object oriented techniques**
 - **Scale dependent techniques**

However....

- Operational precipitation verification still frequently relies upon **ETS, bias**
- Models with different **grid resolution** and different **resolvable-scales** are still being compared

Goals of this Study

- Systematically *document* the scale-sensitivities known to exist for traditional skill scores by...
 - Comparing **equitable threat** and **bias scores** for models verified on different resolution grids
 - Examining spectra from various models and observations on different resolution grids

It is not our purpose to:

- Develop a “new” verification skill score
- Decide how much small-scale detail is acceptable in mesoscale models

Key Questions

- How is **Equitable Threat Score** affected by the amount of small-scale detail in the:
 - **forecast field?**
 - **verification field?**
- How does **model bias** affect this scale dependency?

Specifically....

- Are ETS values from models with **different grid spacing** and **different biases** directly comparable?
- Does a **smoother precipitation field** yield a higher ETS value when compared with a highly **detailed verification field?**

Threat Score and Bias

- Threat Score = Hits / (Hits + Misses + False Alarms)

Highlights events that actually occur, rather than those which do not

ETS is the threat score corrected for a chance forecast...

$$\text{ETS} = (\text{Hits} - \text{chance}) / (\text{Hits} + \text{Misses} + \text{False Alarms} - \text{chance})$$

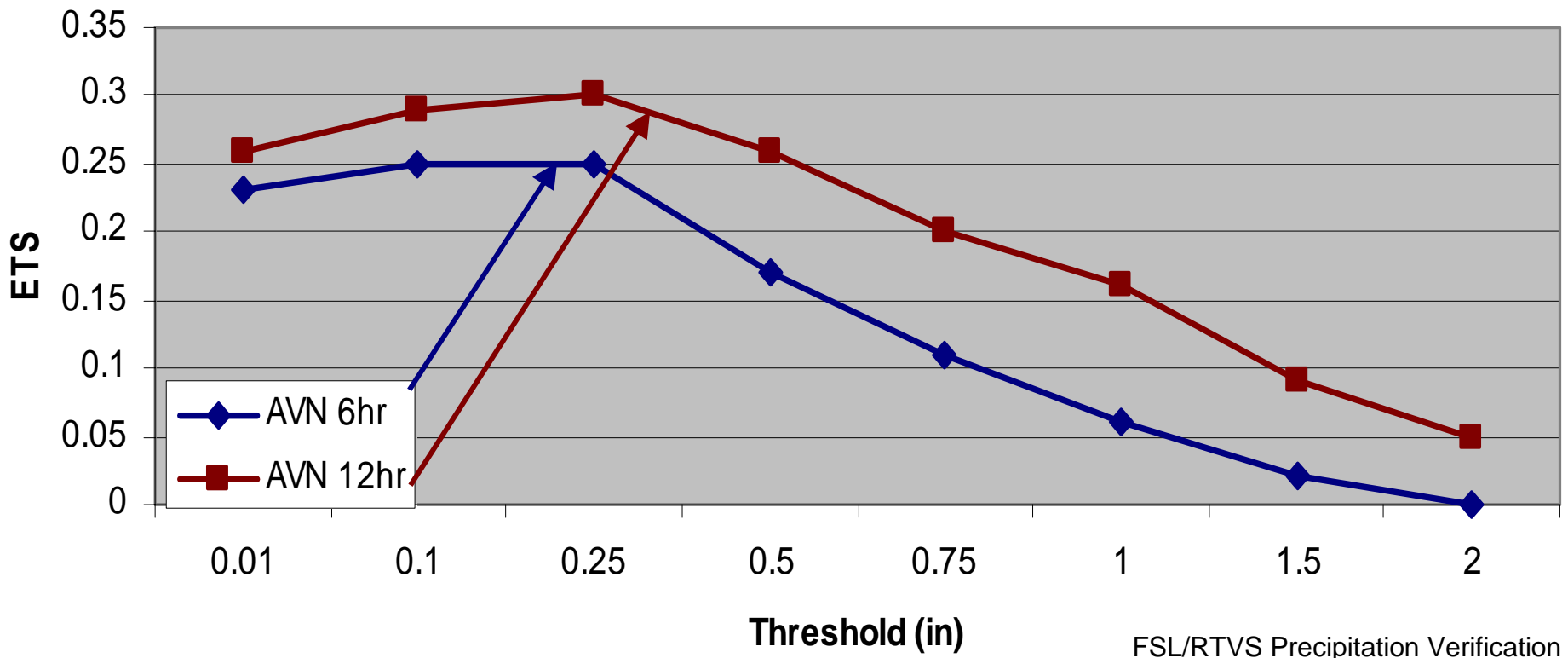
- Bias = Area Forecast / Area Observed

No dependence upon “hits!”

Smoothing of forecast fields over time

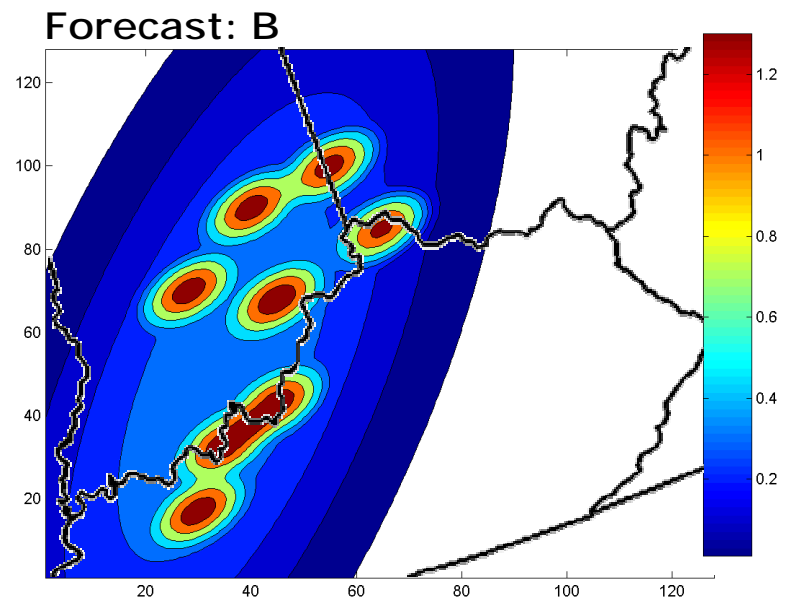
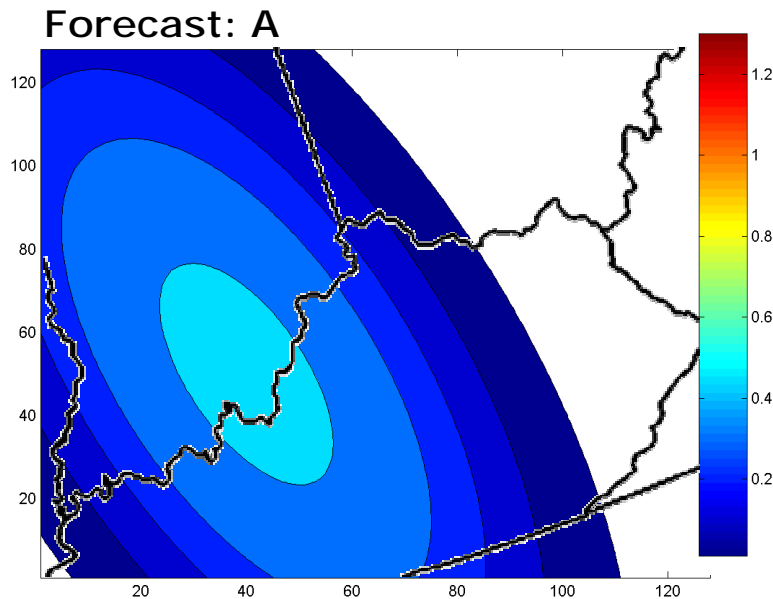
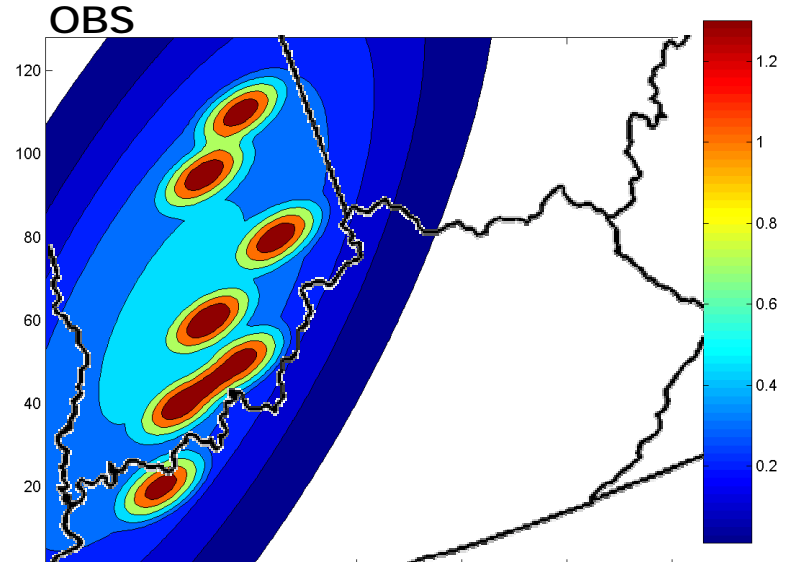
QPF verification statistics computed over a longer accumulation period are shown to be better than those computed over a shorter period

3 Years of AVN (GFS) Verification Statistics
6hr vs. 12hr Accumulation



Spatial smoothing of forecast fields also has been shown to result in higher skill scores...

	<u>A</u>	<u>B</u>
<u>MAE</u> :	0.157,	0.159
<u>RMSE</u> :	0.254,	0.309
<u>Bias</u> :	0.980,	0.980
<u>CSI</u> :	0.214,	0.161
<u>ETS</u> :	0.170,	0.102



From Mike Baldwin of NOAA/NWS/SPC OU/CIMMS

Double penalty

When forecast models resolve very small precipitation detail, they often suffer a **double penalty** when verified categorically on the observational grid.

In this example, the **10 km forecast** is penalized twice: once for not placing rain in the correct place (a miss), and once for placing rain in the wrong place (a false alarm).

The **20 km forecast** receives one hit and 3 false alarms, giving a higher ETS and bias.

$$ETS_{10km} = -0.03$$

$$ETS_{20km} = 0.20$$

FA			
	M		

10 km forecast

20 km forecast

FA	FA		
FA	H		

10 km observational grid

IHOP Real-time Modeling

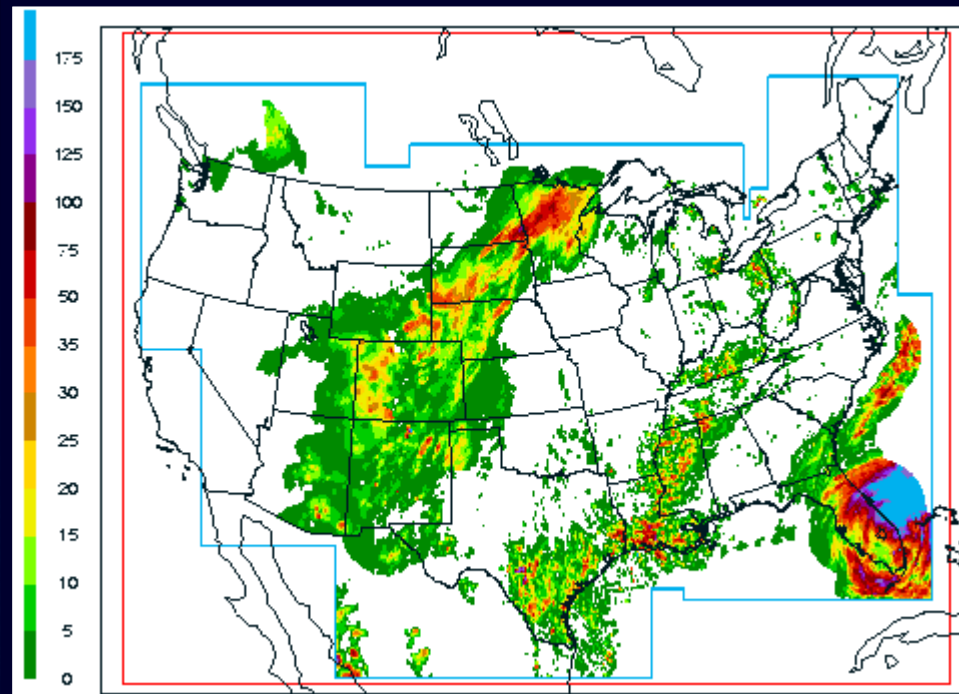
Experimental	RUC 10-km	(GD ensemble convection)
Operational	RUC 20-km	(GD ensemble convection)
Operational	Eta 12-km	(BMJ convection)
Experimental	LAPS MM5 12-km	(KF convection)

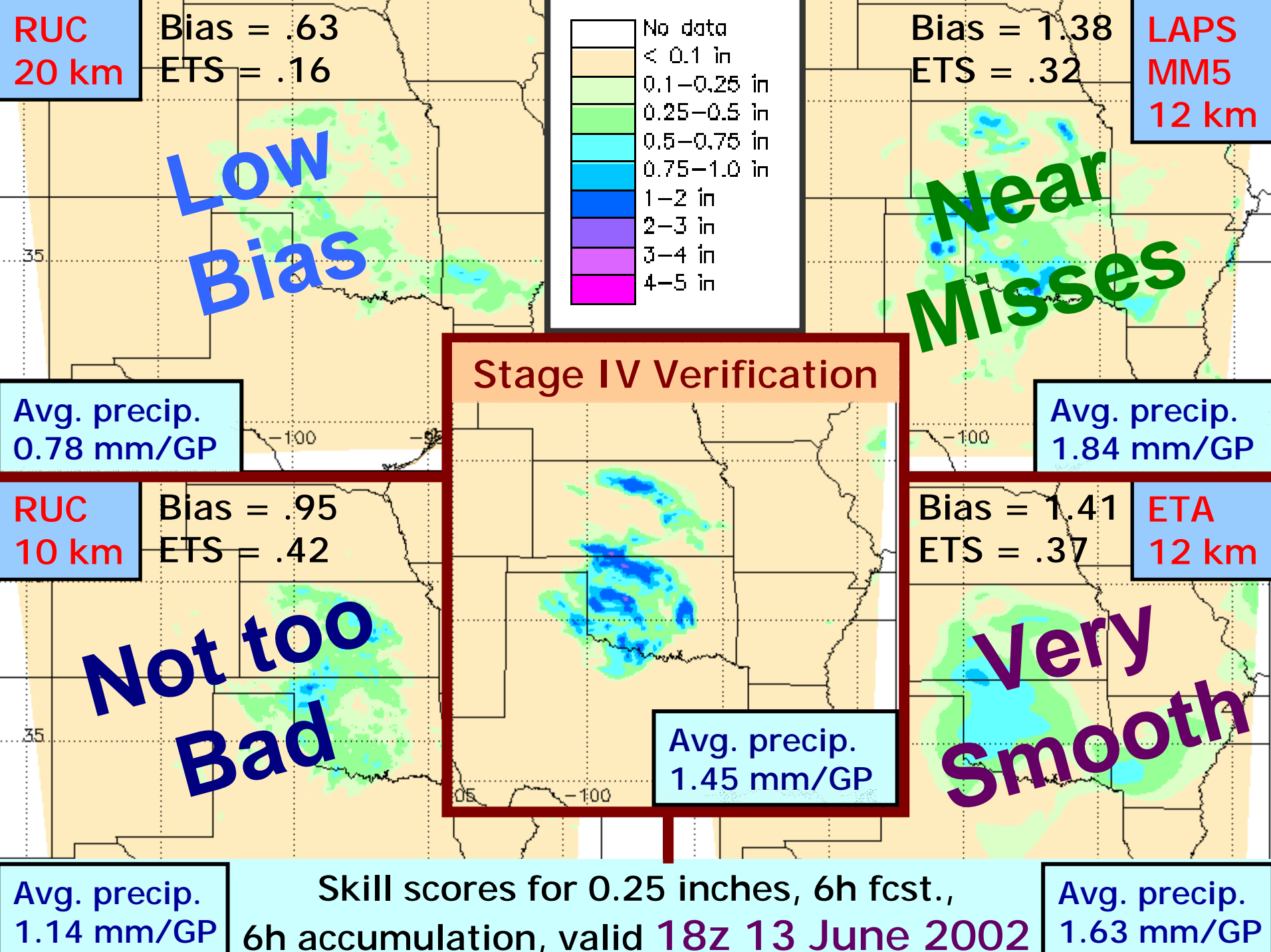
<u>Model</u>	<u>Native (km)</u>	<u>Coarsened (km)</u>		
RUC10	10	20	40	80
RUC20	20	20	40	80
ETA12	12	20	40	80
LMM12	12	20	40	80
Stage4 verif	4 (10)	20	40	80

This study is not a model bake off!

Observations: NCEP Stage IV Analysis

Mosaic of regional hourly and 6-hourly multi-sensor (radar+gauges) precipitation analysis at 4km.

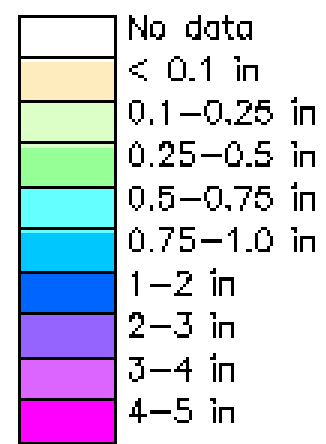




RUC
20 km

Bias = .63
ETS = .16

Low
Bias



Bias = 1.38
ETS = .32

LAPS
MM5
12 km

Near
Misses

Avg. precip.
0.78 mm/GP

Stage IV Verification

Avg. precip.
1.84 mm/GP

RUC
10 km

Bias = .95
ETS = .42

Not too
Bad

Bias = 1.41
ETS = .37

ETA
12 km

Very
Smooth

Avg. precip.
1.45 mm/GP

Avg. precip.
1.14 mm/GP

Skill scores for 0.25 inches, 6h fcst.,
6h accumulation, valid 18z 13 June 2002

Avg. precip.
1.63 mm/GP

Upscale forecasts and observations

EXPT. 1

Remap Stage IV and model data to common, coarser resolution grids

Compare scores from forecasts with different precipitation detail verified on their native grid

Stage IV

Native

10-km

20-km

40-km

80-km

model native

← compare →

← compare →

← compare →

Model

Native

20-km

40-km

80-km

← compare →

Smooth forecasts only

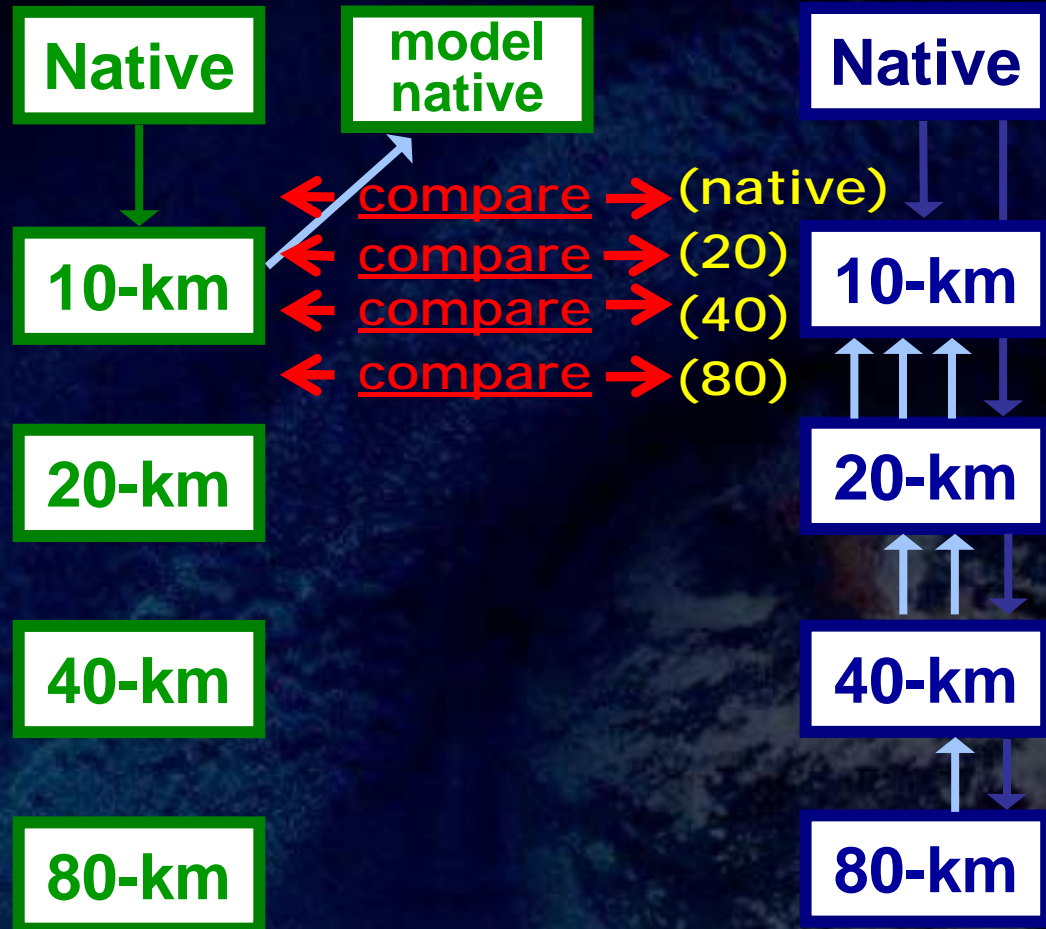
EXPT. 2

Stage IV

Model

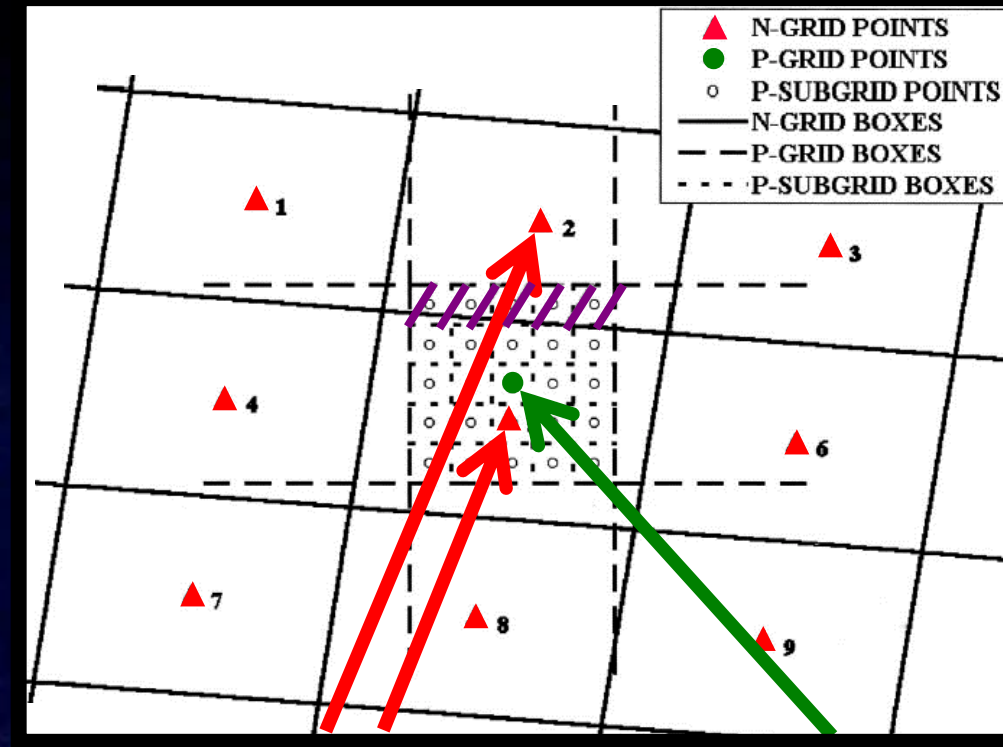
Remap native and coarsened forecast fields to 10-km Stage IV grid

Compare scores from forecasts with different precipitation detail verified against detailed Stage IV data



Grid Transformations

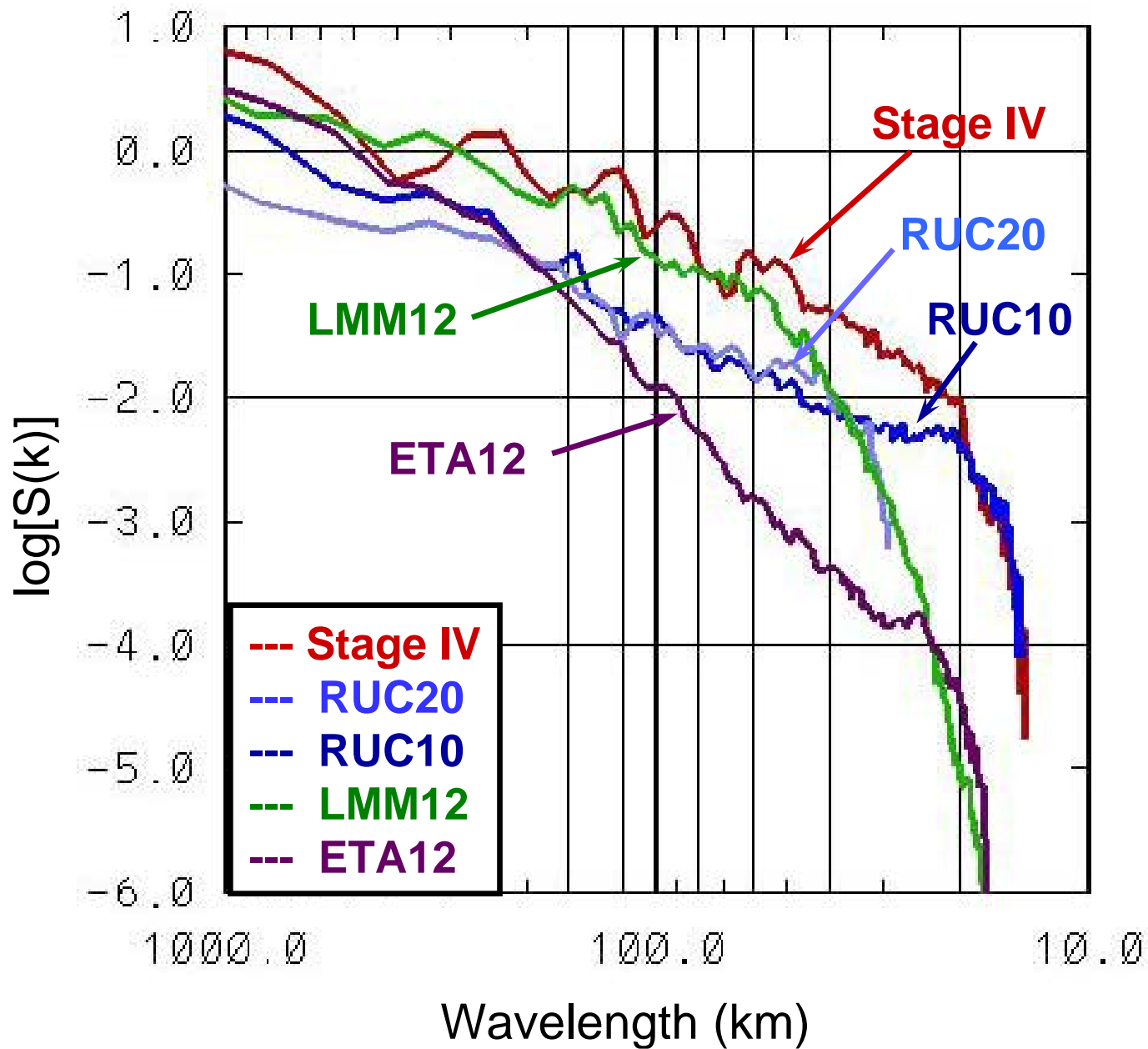
- NCEP “neighbor budget” (Baldwin 2000) used for all grid remappings
- Preserves total precip, minimizes edge smearing
- Less impact on skill scores than bilinear interp (Accadia et al., 2003)

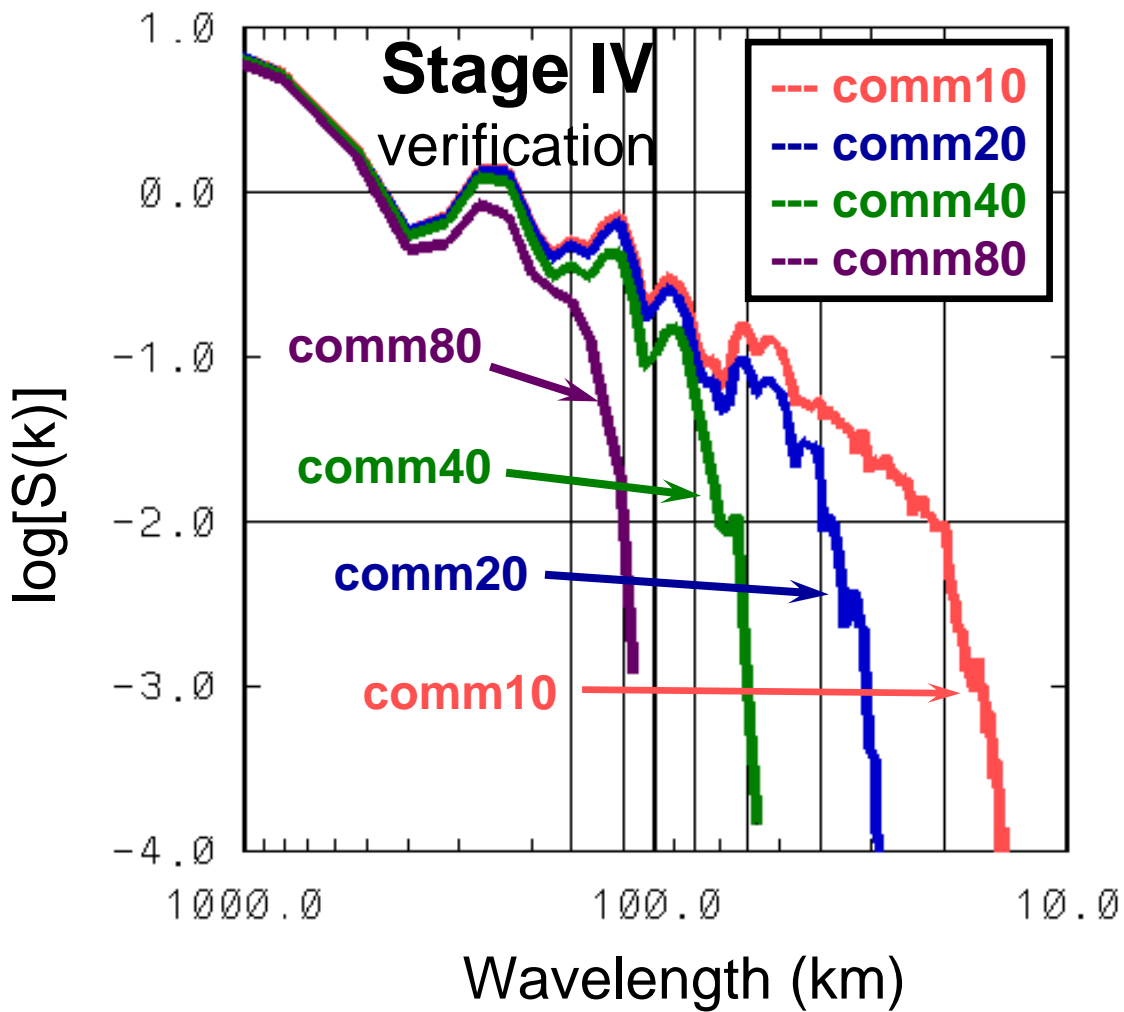


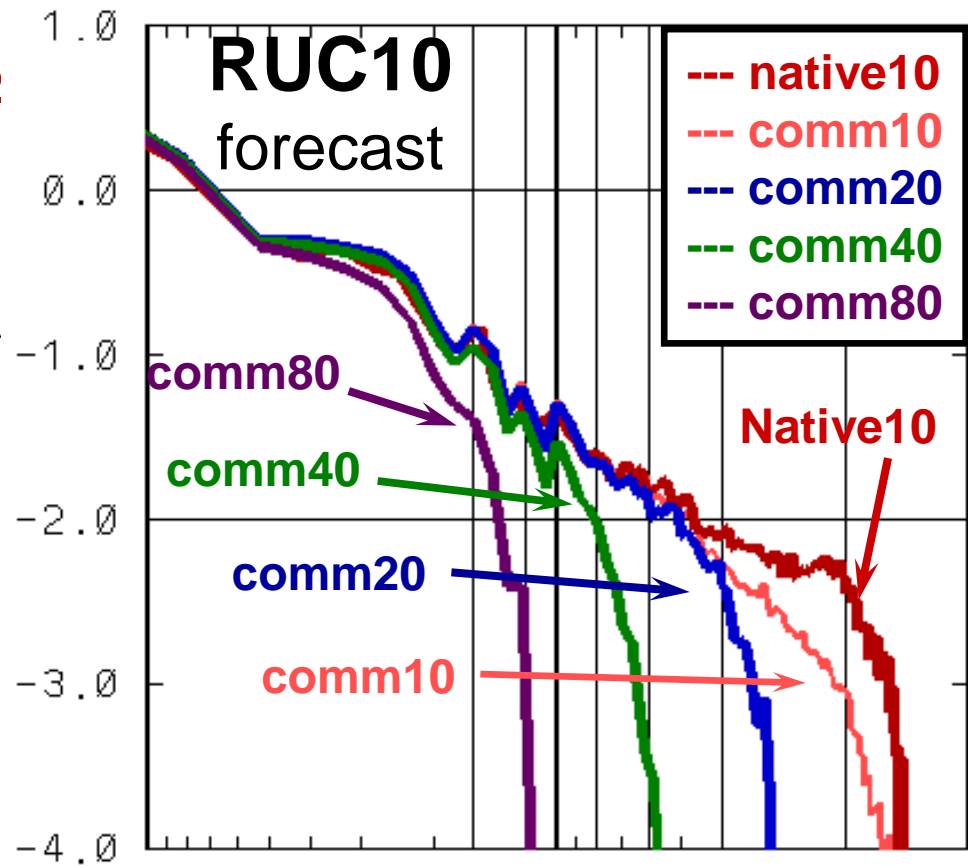
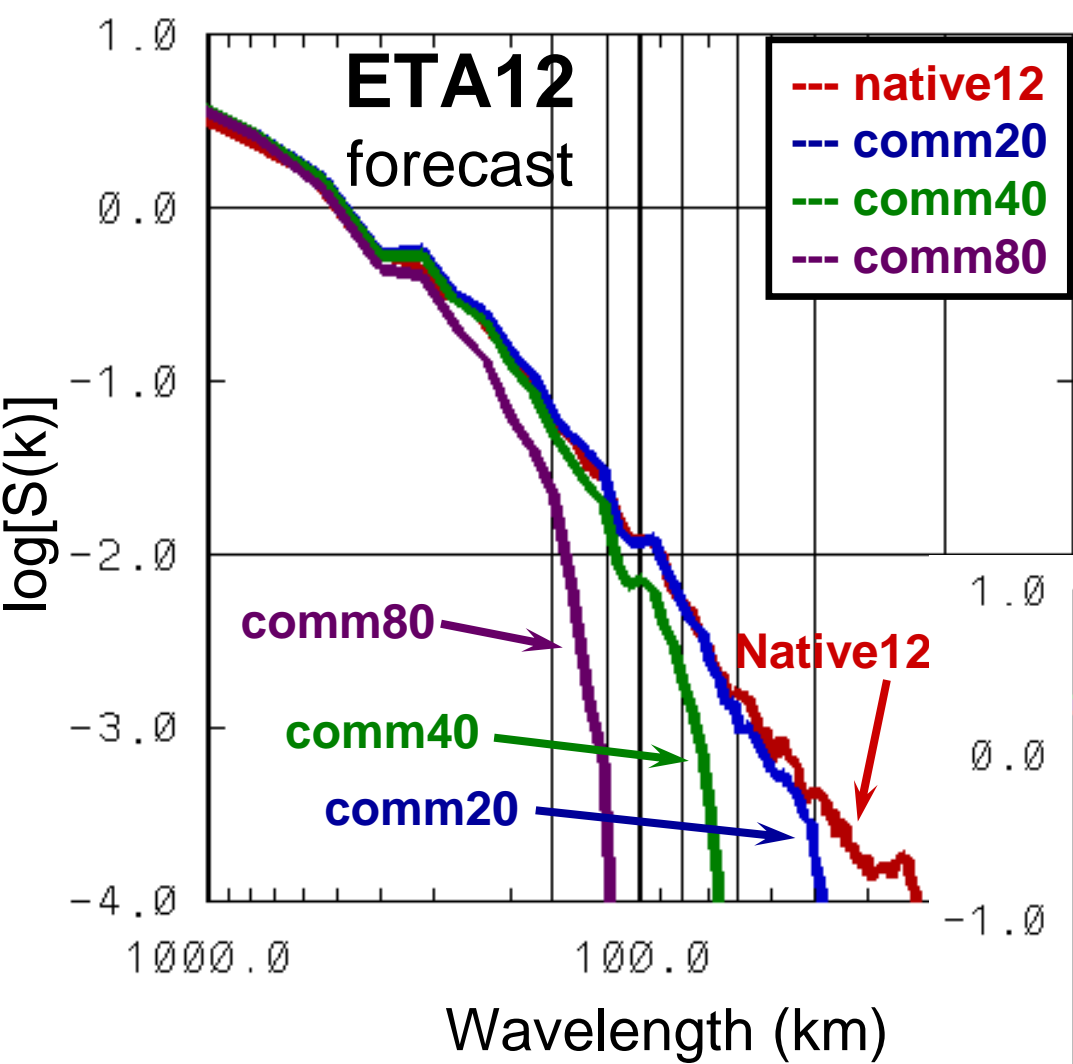
input
points

target
point

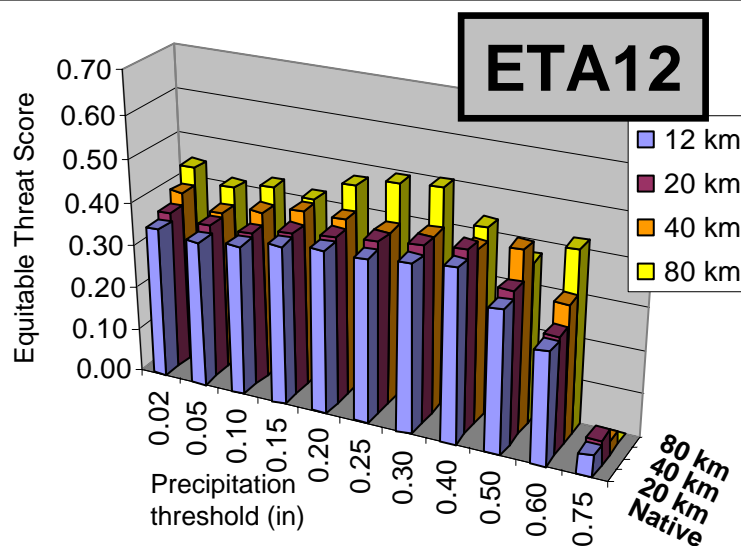
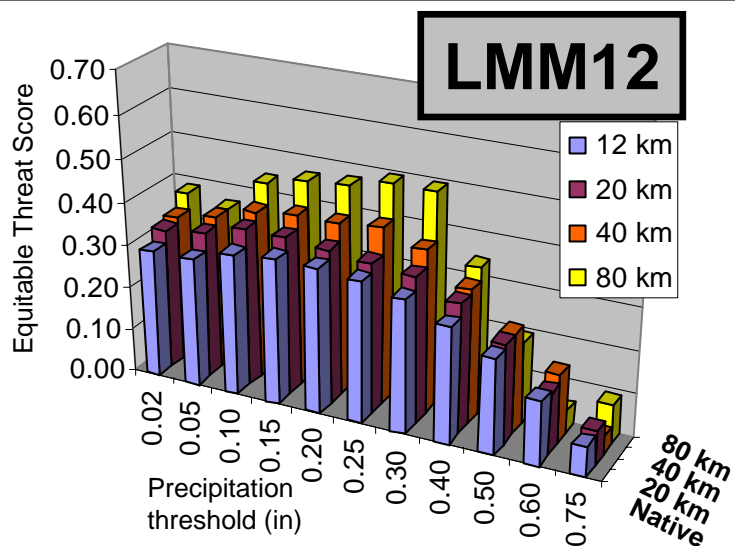
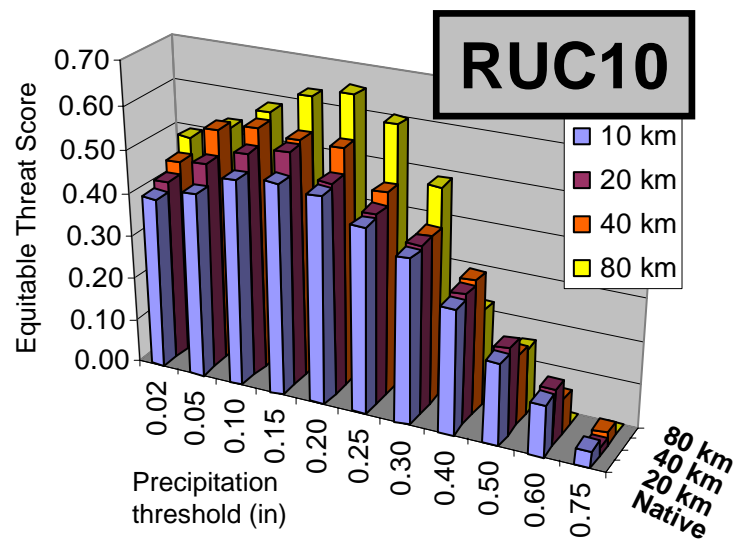
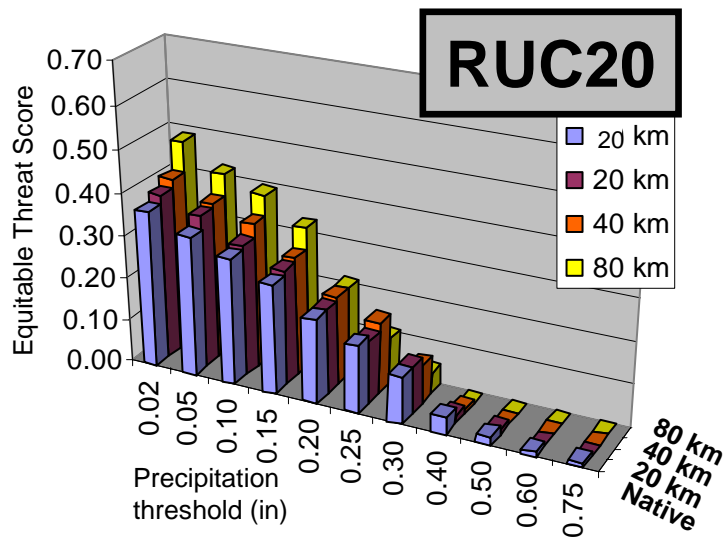
- Sub-divide each target grid-box into 25 sub-boxes (5x5)
- Nearest neighbor from input grid to each sub-box point
- Target values = simple average of 25 sub-box values



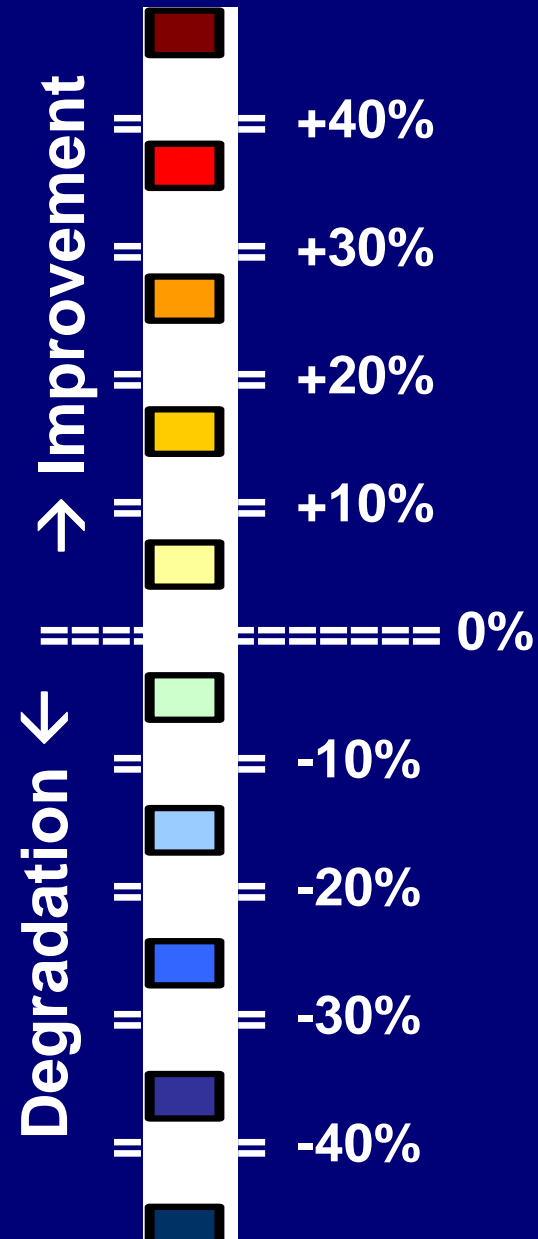




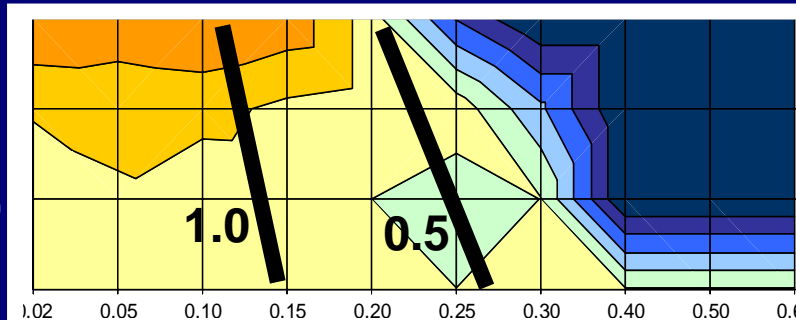
Expt. 1 Results: Upscale model and verification



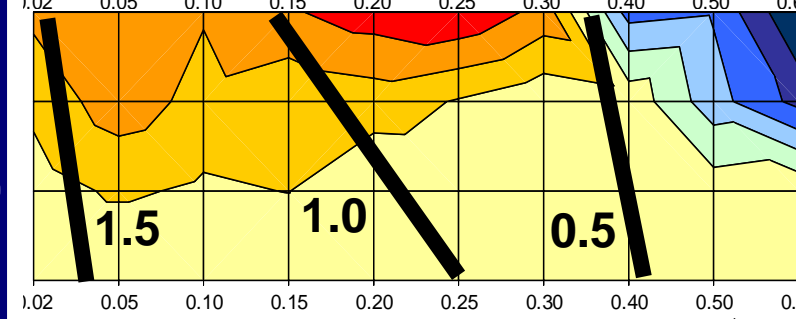
Expt. 1: ETS % change relative to native grid



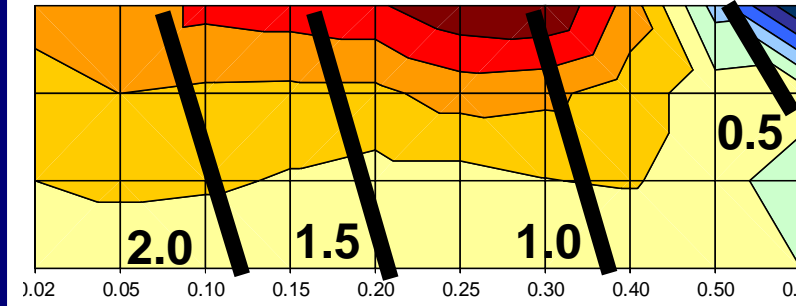
(a)
RUC20



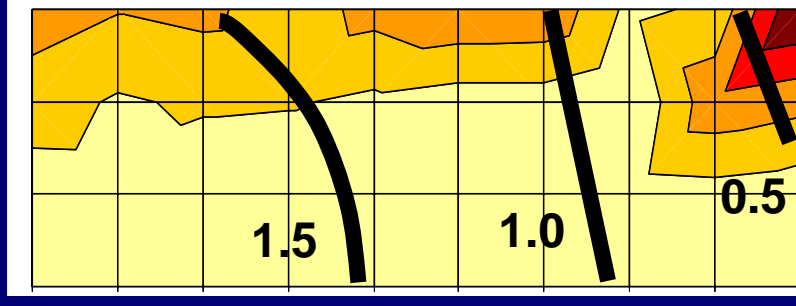
(b)
RUC10



(c)
LMM12



(d)
ETA12

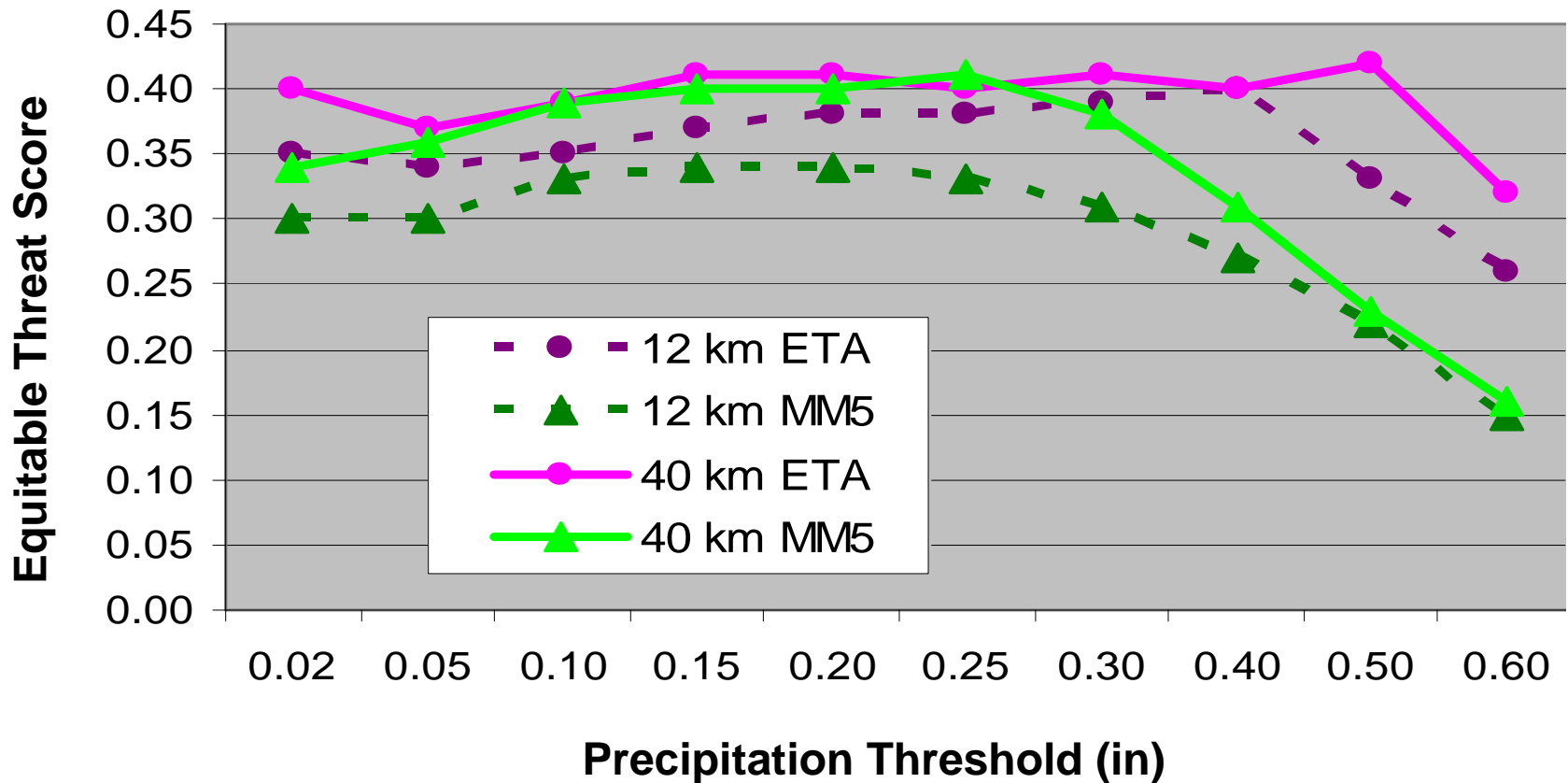


80 km
40 km
20 km
Native
80 km
40 km
20 km
Native
80 km
40 km
20 km
Native
80 km
40 km
20 km
Native

Verification Resolution (km)

0.02 .05 .10 .15 .20 .25 .30 .40 .50 .60
Precipitation Threshold (in)

Expt. 1 Results: Upscale model and verification



LMM12 (near misses) improves more with upscaling than ETA12 (very smooth)

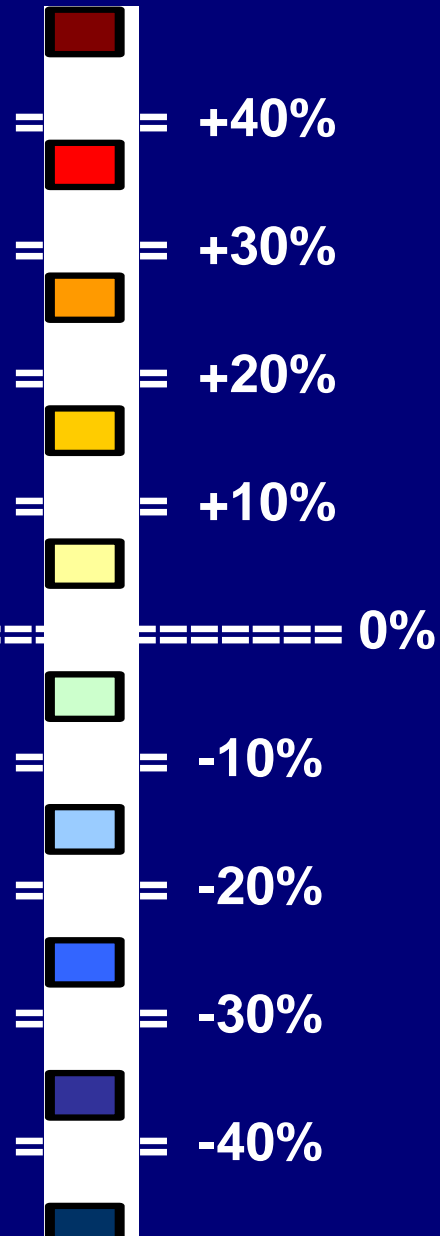
Summary of Expt. 1 Results

- ETS **improves for all models** and **most thresholds** as forecast and verification fields are upscaled
- For detailed forecasts, a **precipitation threshold cutoff** exists above which **forecast degradation** occurs with upscaling
- The **cutoff threshold** shifts to lower amounts with further upscaling, and is **correlated with bias ~ 0.5**
- For **smooth forecasts**, **less ETS improvement** with upscaling occurs and no cutoff threshold exists

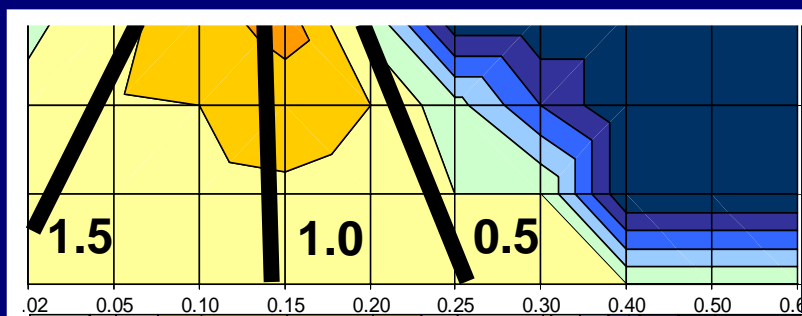
How do these results change, when only the forecast is smoothed?

Expt. 2 (Smooth model only): ETS % change

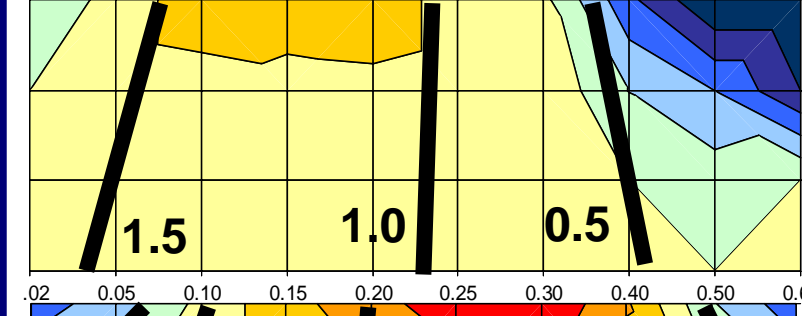
↑ Improvement
↓ Degradation ←



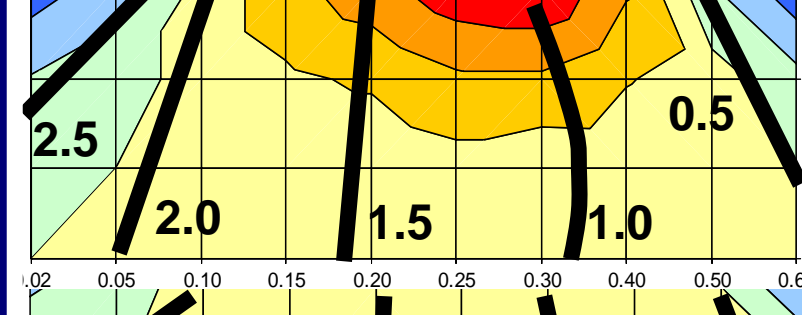
(a)
RUC20



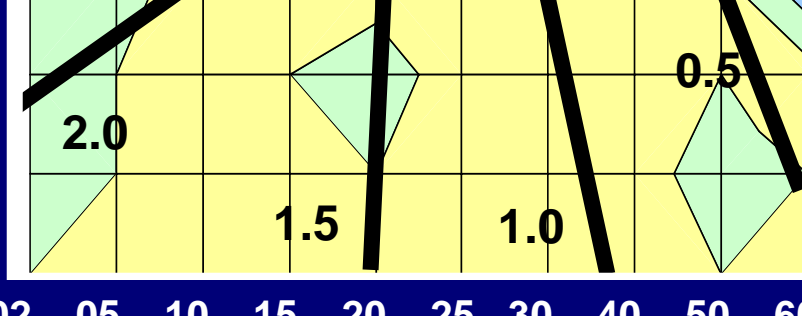
(b)
RUC10



(c)
LMM12



(d)
ETA12

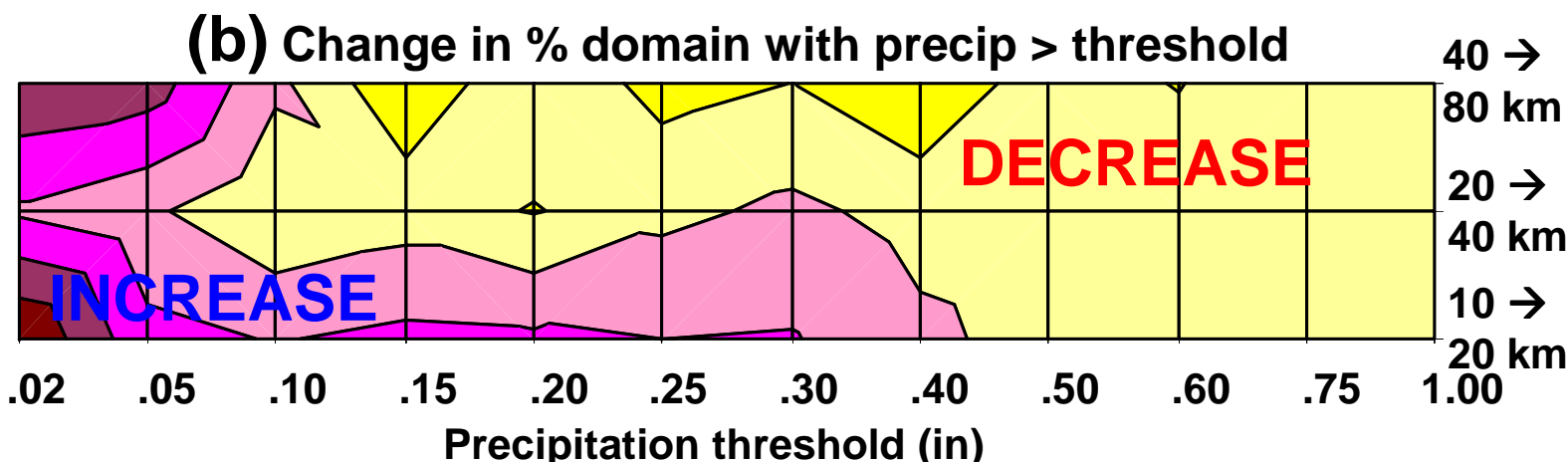
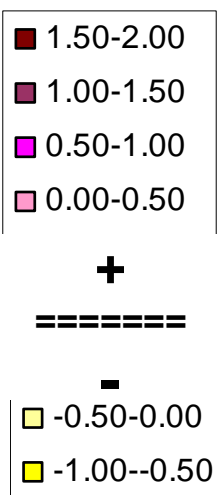
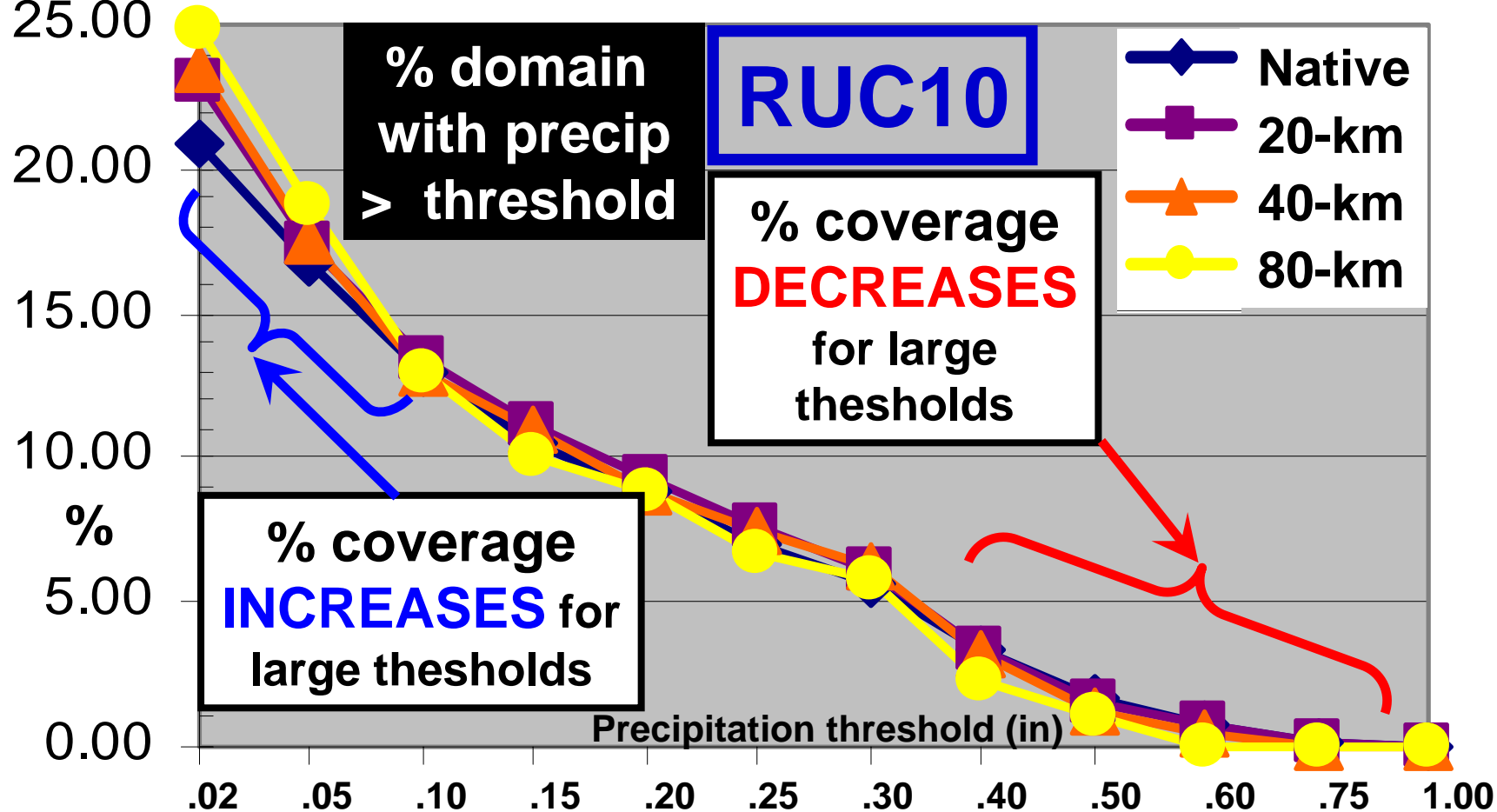


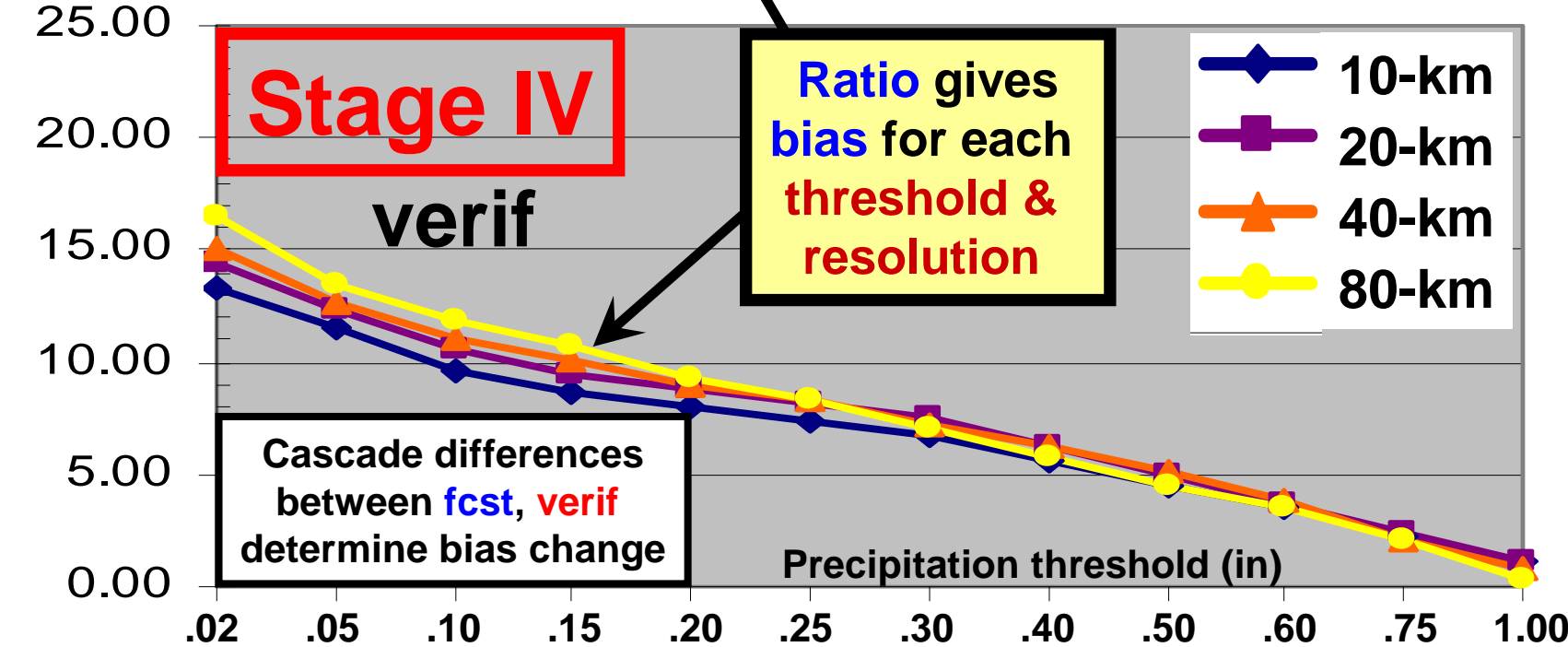
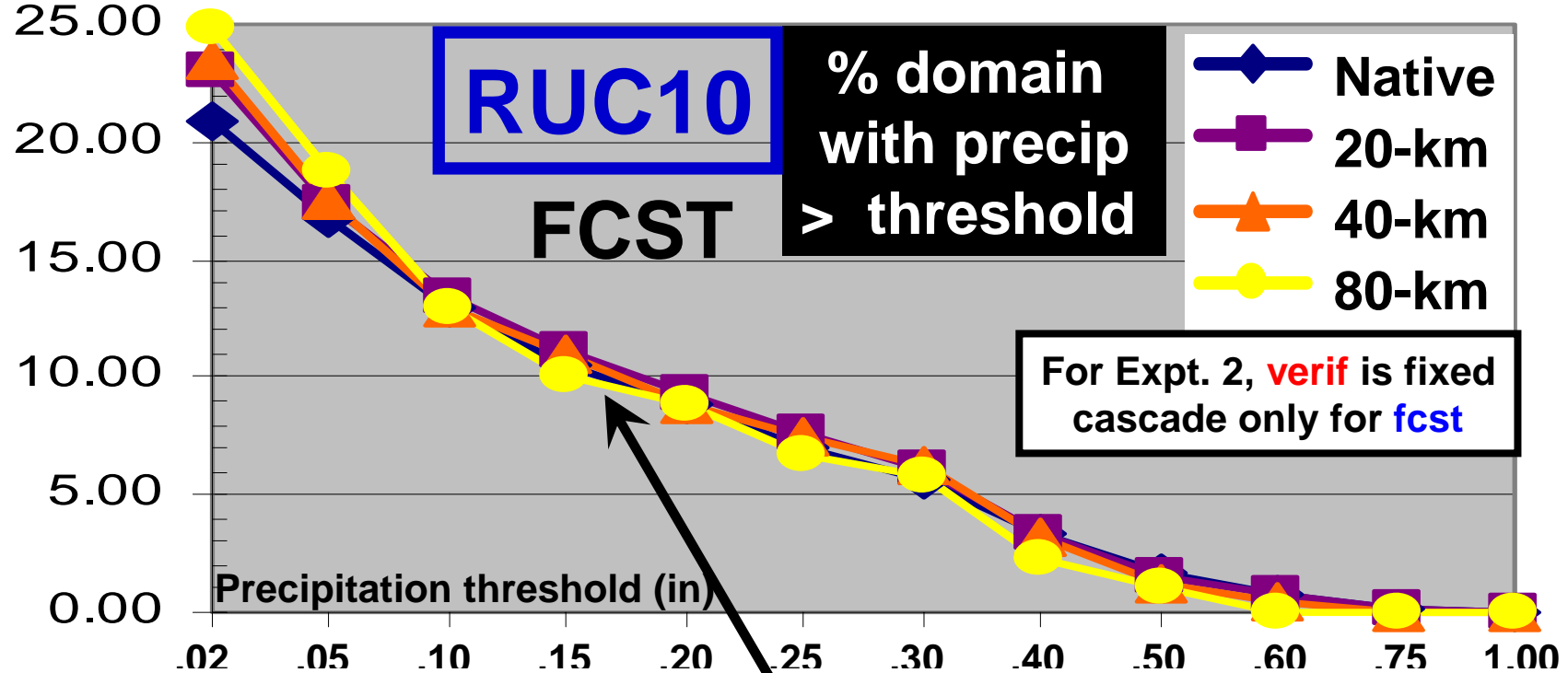
Verification Resolution (km)

Summary of Expt. 2 Results

- Even when verified against a fixed detailed field, **smoothing the forecast improves the ETS score**
- Bias **decreases for the highest thresholds** and **increases for the lowest thresholds**
- Upper cutoff threshold (bias ~ 0.5) remains, **ETS falls for low thresholds** as bias exceeds 2.0 for smoothed fields
- For smooth forecasts, very little change in ETS (no changes for either forecast or observations)

For ETS, smoother is better (either forecast or observations), with current model skill*





What controls ETS and bias changes?

- As forecast and observations are smoothed, local maxima are reduced, and larger precipitation amounts spread to nearby points
- Result is an overall **cascade** of precipitation from higher thresholds to lower thresholds
- **ETS**: Small-scale near misses suddenly become **hits!**
- **Bias**: Increase in coverage for low thresholds, decrease in coverage for high thresholds

Precipitation cascade is largely controlled by:

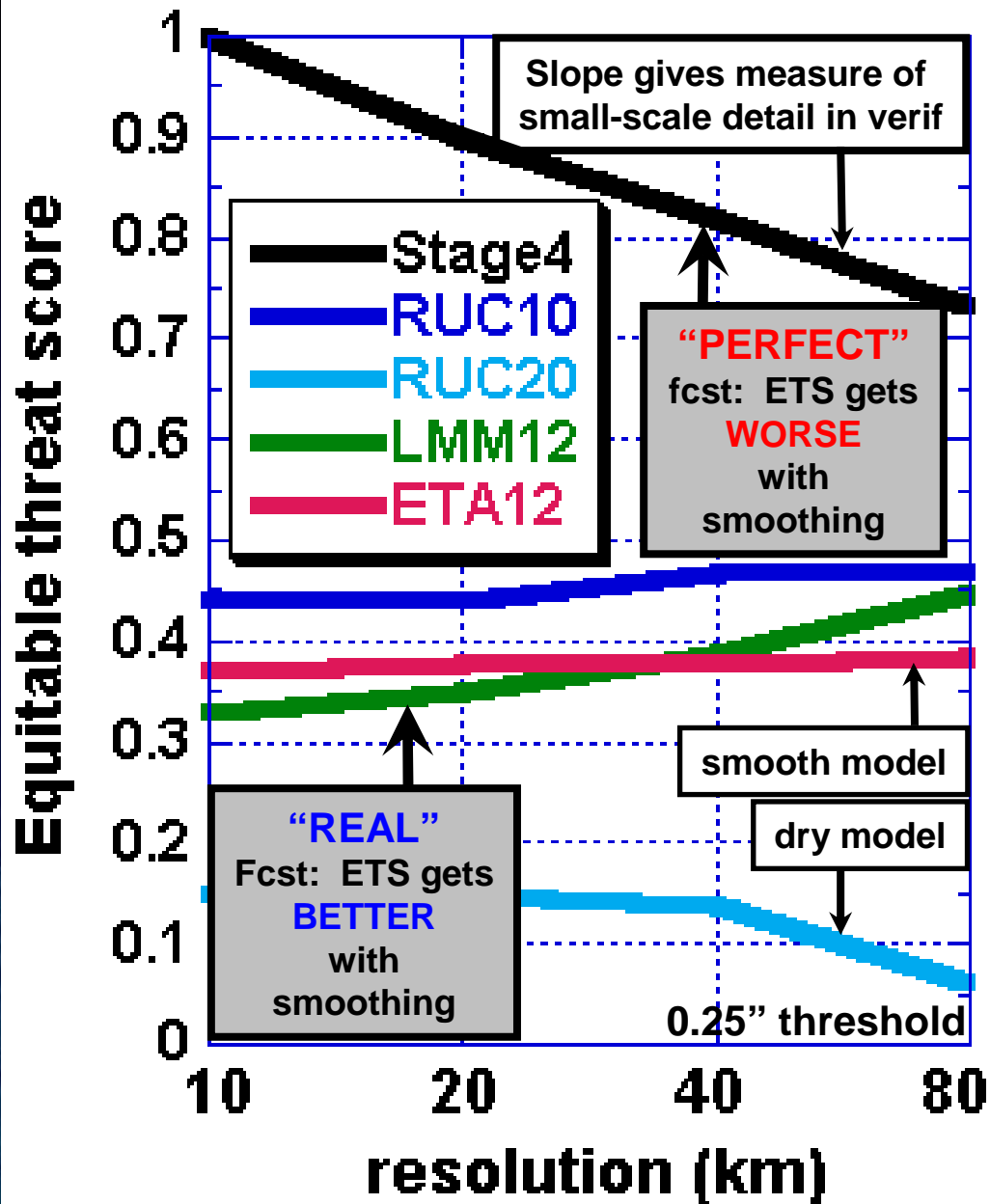
- **Small scale detail (spectra)**
- **Total precipitation volume**

* What if model skill was better ?

- ETS rewards **gridpoint** matches
- Details must be in the correct location
- **Models are not that good yet!**

ETS for coarsened “perfect” forecast gives upper-bound on ETS for a given amount of detail

ETS for upscaled forecasts verified on common 10-km grid



Conclusions

- Forecasts on **different native grids** are **not** directly comparable (coarser grid has the advantage)
- Forecasts with **different degrees of small-scale detail**, even if on the same grid, are **not** comparable (smoother field has the advantage)

ETS comparisons should only be made for precipitation fields with **similar spectra and bias, compared on matched grid resolutions (**using the same verification field**)**

Better verification measures?

- Spatial structure measures
- Object Oriented measures
- Scale dependent techniques

There is no:

- **one-size fits all verification score**
- **optimal amount of model detail**

Highly detailed forecasts often better duplicate observed spatial and temporal structures, contain more information for use in the model post-processing



That's all folks!