

Existing Verification Techniques

Traditional methods:

1. graphical summary
(scatter-plot, box-plot)
2. Continuous scores
(RMSE, correlation)
3. Categorical scores from
contingency tables (FBI,
ETS, POD)

Spatial methods:

1. Scale-decomposition
2. Neighbourhood-based
3. Field-morphing
4. Feature-based
5. Distance metrics for
binary images

There is no single technique which fully describes the complex
observation-forecast relationship

Key Q: what do we wish to know from our verification ?

Spatial verif inter-comparison: Gilleland et al (2009)

1. Scale-decomposition approaches

Briggs and Levine (1997), wavelet cont (MSE, corr);

Casati et al. (2004), Casati (2010), wavelet cat (HSS, FBI, scale structure)

Zepeda-Arce et al. (2000), Harris et al. (2001), Tustison et al. (2003), scale invariants parameters;

Casati and Wilson (2007), wavelet prob (BSS=BSSres-BSSrel, En2 bias, scale structure);

Jung and Leutbecher (2008), spherical harmonics, prob (EPS spread-error, BSS, RPSS);

Denis et al. (2002,2003), De Elia et al. (2002), discrete cosine transform, taylor diag;

Livina et al (2008), wavelet coefficient score. **De Sales and Xue (2010)**

1. Decompose forecast and observation fields into the sum of spatial components on different scales (wavelets, Fourier, DCT)
2. Perform verification on different scale components, separately (cont. scores; categ. approaches; probability verif. scores)

Account for the field coherent spatial structure:

- ➔ Assess scale structure
- ➔ Bias, error and skill on different scales
- ➔ Scale dependency of forecast predictability (no-skill to skill transition scale)

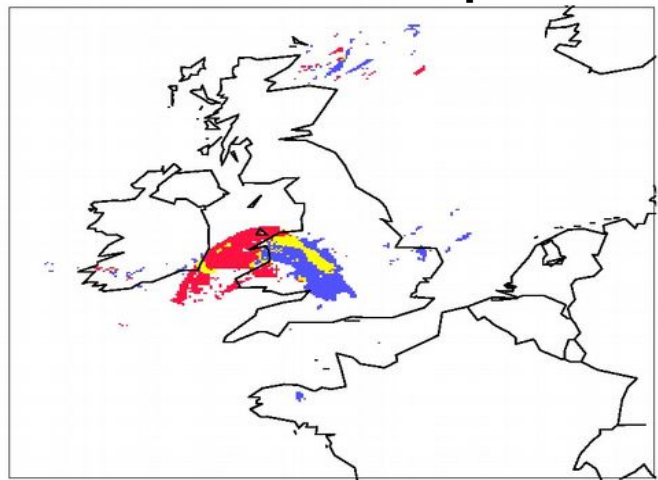
Tech Note: in the scale decomposition approaches the scale is obtained by a single-band pass filter. The scale is associated then to the feature size, provides feedback on physical processes associated to phenomena on different scales.

Intensity-scale verification technique Casati et al. (2004), Met Apps, vol. 11

The **intensity-scale** verification approach measures the skill as function of precipitation intensity and spatial scale of the error

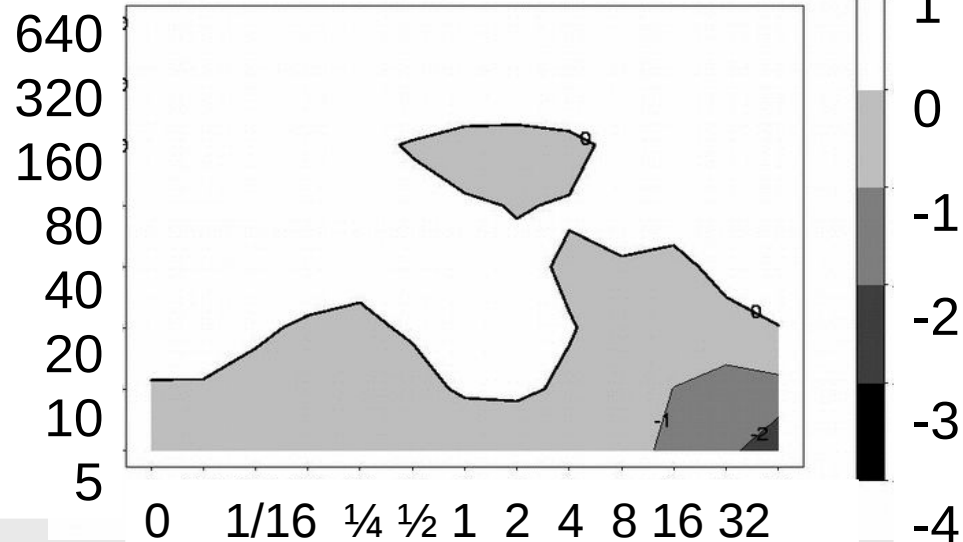
1. Intensity: **Threshold** => Categorical approach
2. Scale: **2D Wavelets** => decomposition of binary images
3. For each threshold and scale: skill score associated to the MSE of binary images = Heidke Skill Score

Intense storm displaced



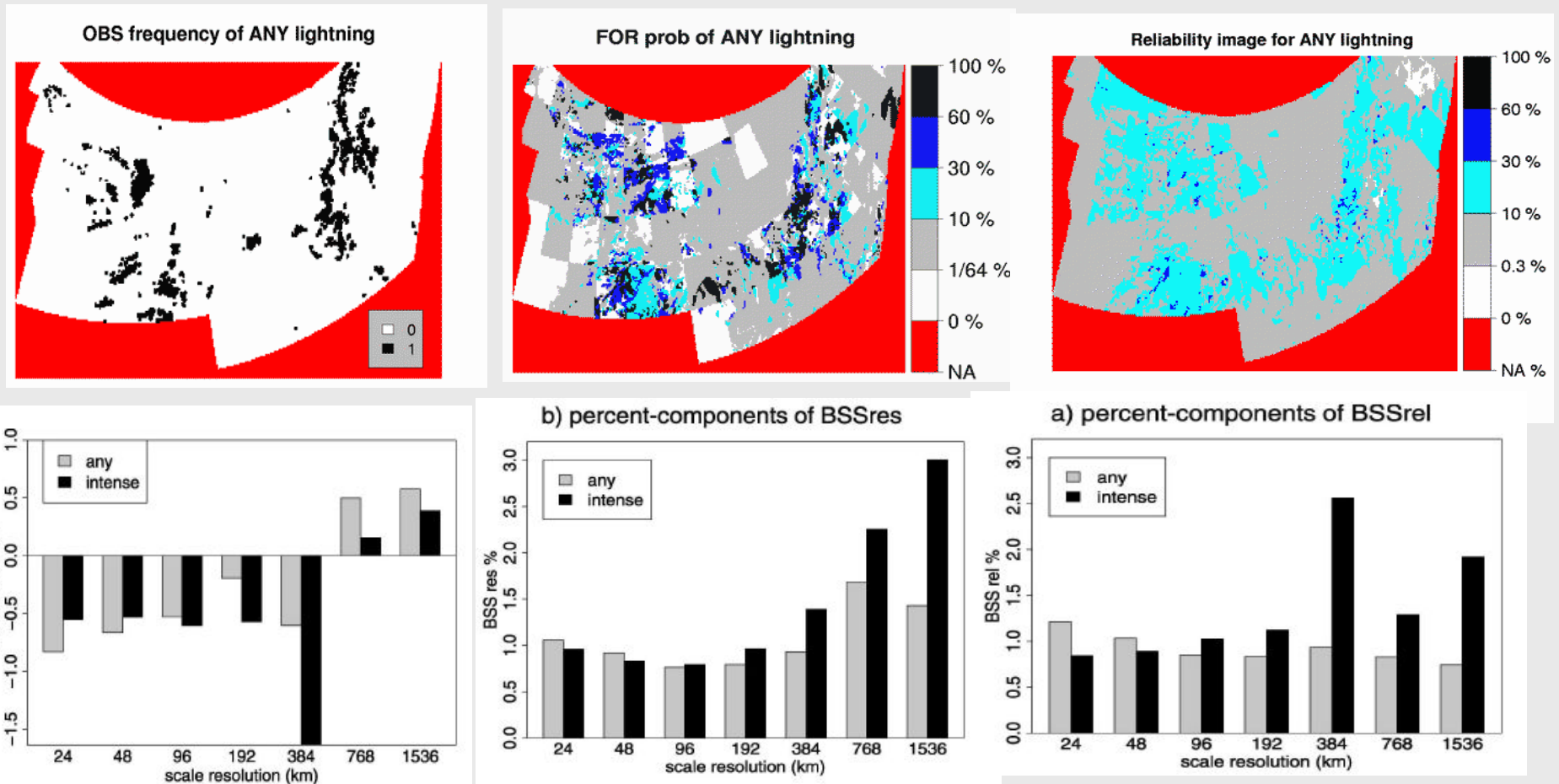
threshold = 1mm/h

scale (km)



threshold (mm/h)

Casati and Wilson (2007) MWR 135

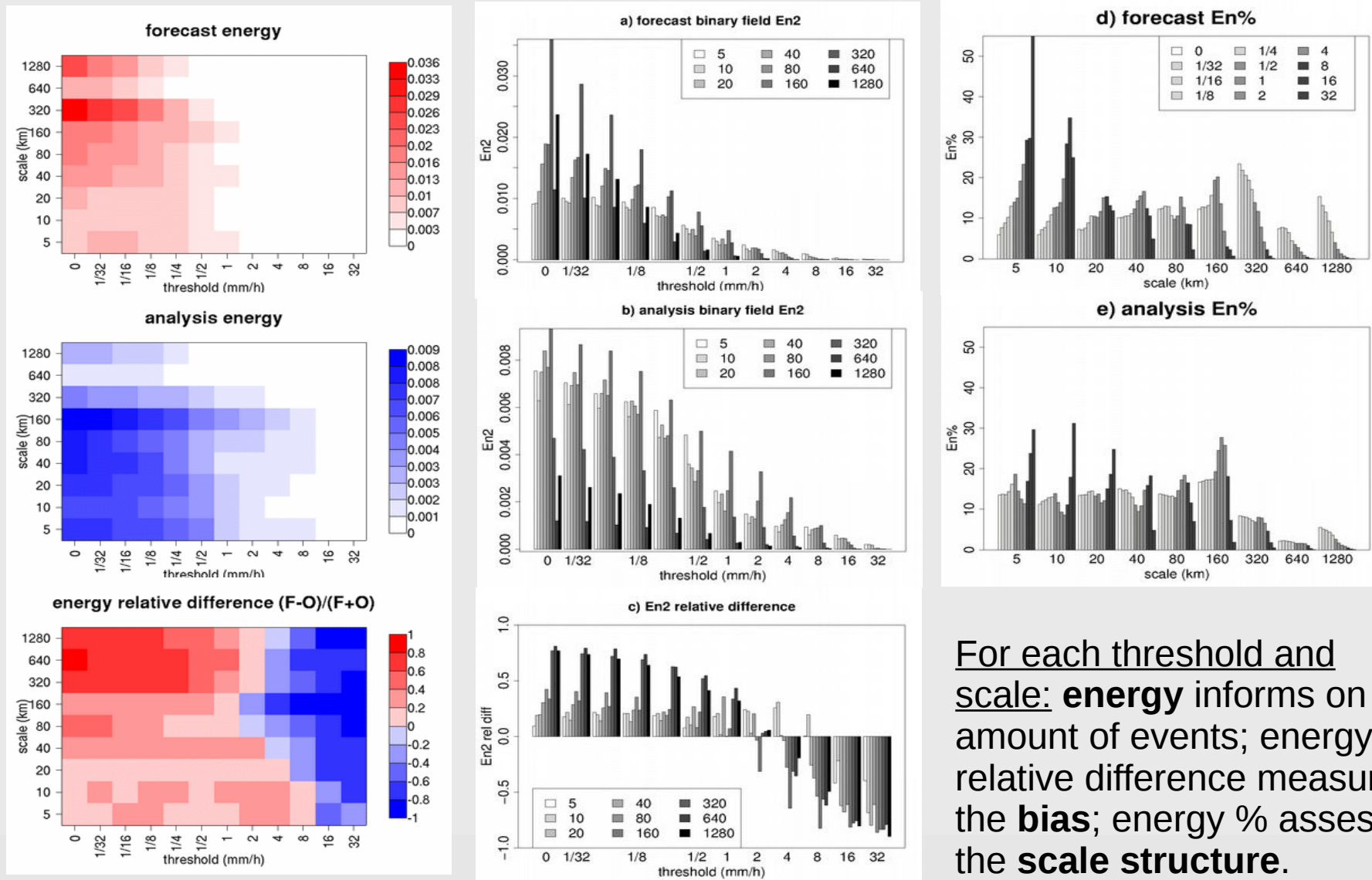


Skill on different scales: positive on large scales (> 700 km); negative on small scales (< 350 km); no-skill to skill transition scale ~ 500 km. **Bias on different scales:** overforecast of 400km feature (leads to poor skill). Assess scale structure.

Quantify **contribution of different scales to reliability and resolution** components in the Brier Skill Score ($BS = rel - res + unc$, $BSS = BSSres - BSSrel$).

Intensity-scale verification technique

Casati (2010) Wea & For, vol. 25



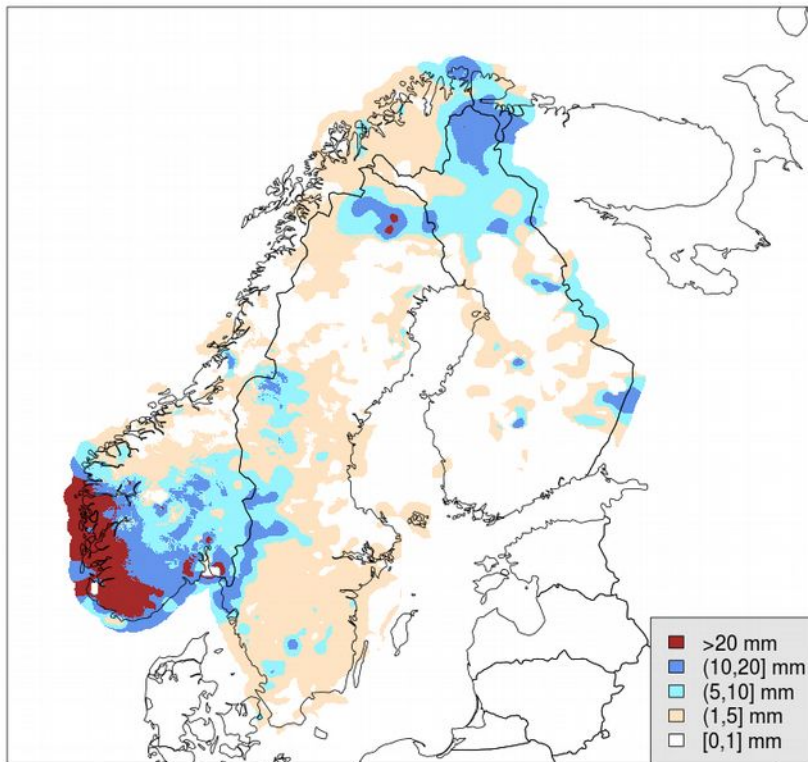
For each threshold and scale: **energy** informs on the amount of events; energy relative difference measures the **bias**; energy % assess the **scale structure**.

UERRA: Scale-Separation verification technique

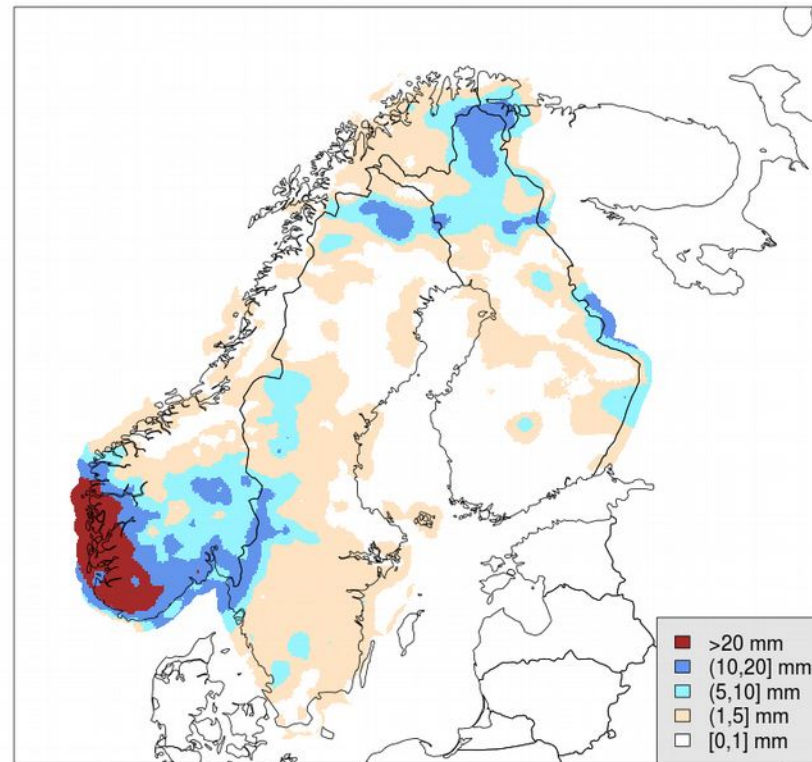
- **RRA**: EURO4M MESAN data – European high-res surface reanalysis
- **Ref**: NGCD
- **Var**: RR, daily prec
- Time period considered: JJA 2008

<http://exporter.nsc.liu.se/620eed0cb2c74c859f7d6db81742e114/>

NGCD, 24h PRec, 200806200000



SMHI-MESAN, 24h PRec, 200806191800

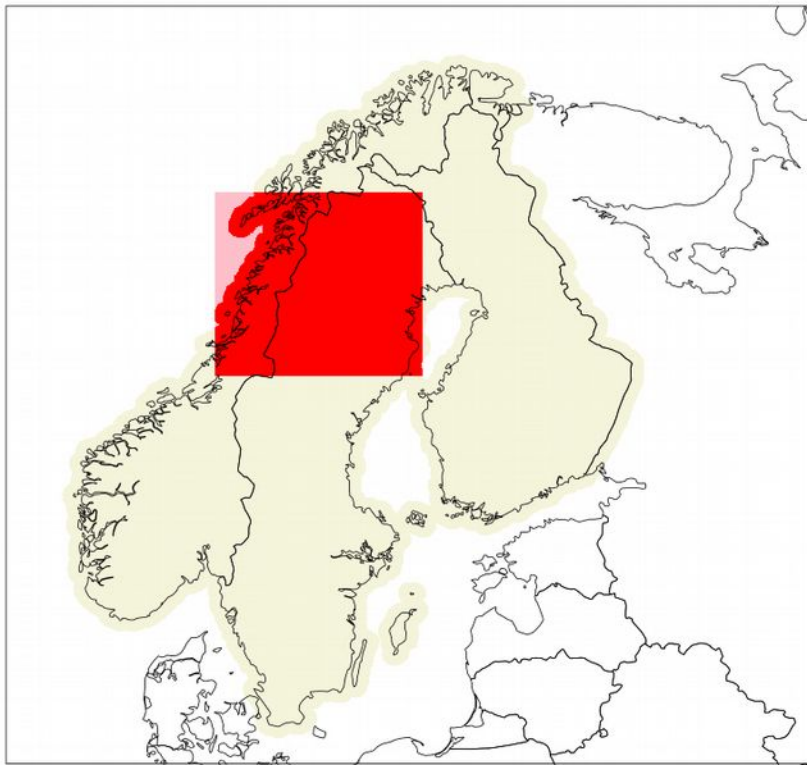


!! NGCD and SMHI-MESAN are both based on ECA&D data !!
VERIFICATION of NOT INDEPENDENT datasets

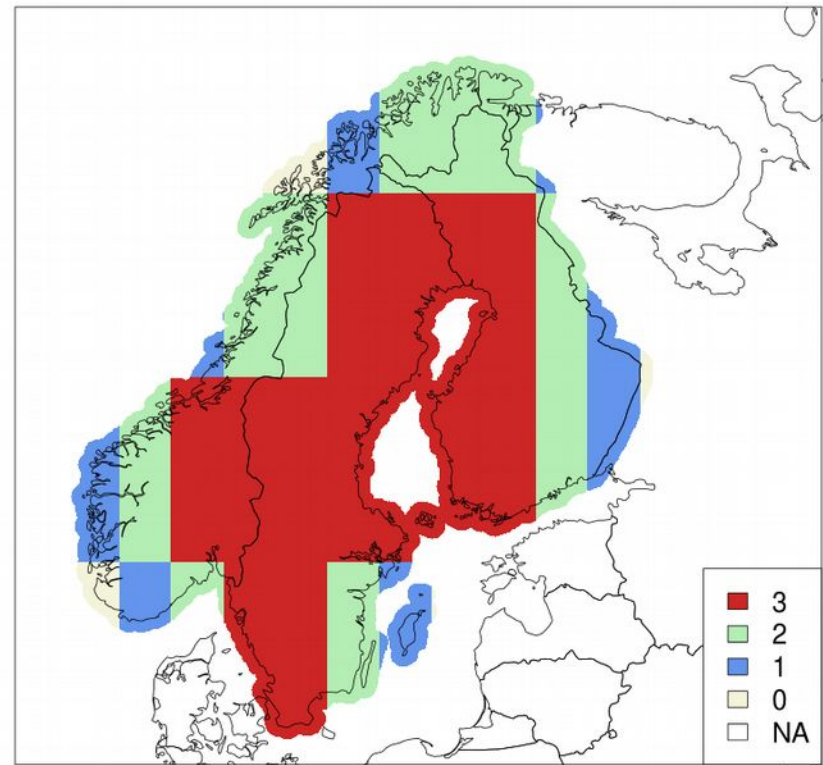
Change of Support problem: reprojection/regridding...

- SMHI-MESAN has been reprojected on the NGCD grid (nearest neighbor interpolation)
- Given that: NGCD covers land only & discrete wavelet transform need $2^n \times 2^n$ domains (i.e. dyadic domain)
---> verification over several overlapping dyadic sub-domains 512 Km x 512 Km. *Tiling approach*: smoothes out the effect due to discreteness of the wavelet transform support.

Dyadic Sub-Domains



composite of Dyadic Sub-Domains



case studies: RR for 1 day in summer 2008

Energies are evaluated for each scale component

$$\mathbf{x} = \sum_{l=1}^L W_l^m(\mathbf{x}) + W_L^f(\mathbf{x})$$

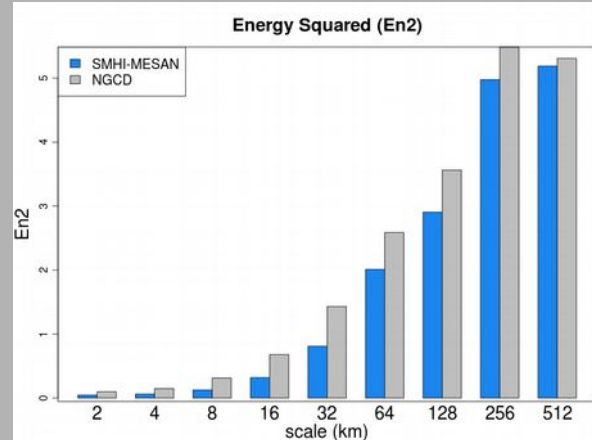
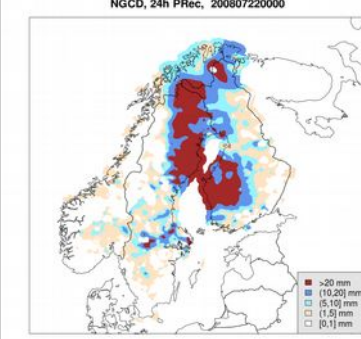
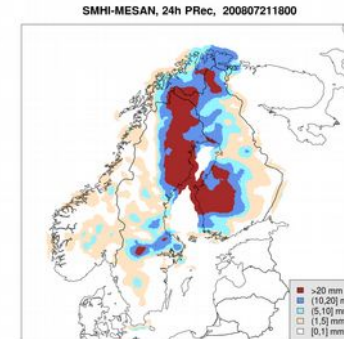
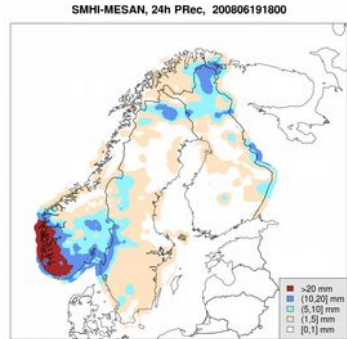
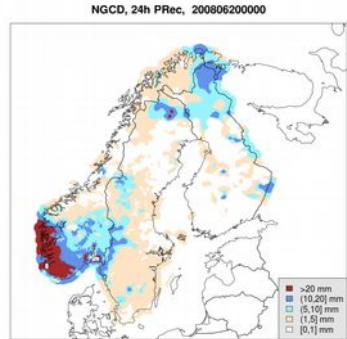
$$\text{En}^2(\mathbf{x}(t)) = \overline{\mathbf{x}(t)^2}$$

$$\text{En}^2(\mathbf{x}(t)) = \sum_{l=1}^L \text{En}^2[W_l^m(\mathbf{x}(t))] + \text{En}^2[W_L^f(\mathbf{x}(t))]$$

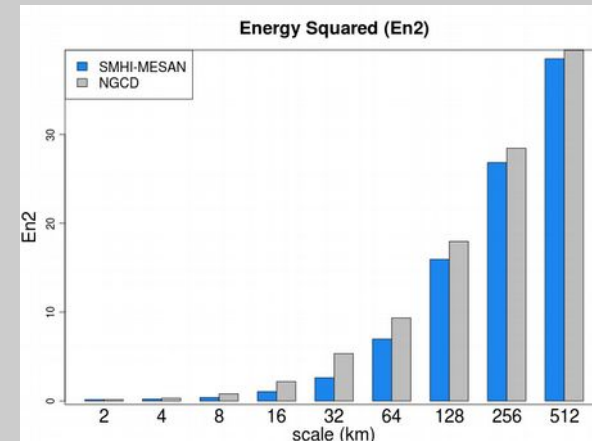
An intensity-scale skill score to assess the added value of enhanced resolution

B. Casati, A. Glazer, J. Milbrandt, and V. Fortin

http://presentations.copernicus.org/EMS2015-250_presentation.pdf



Case 1



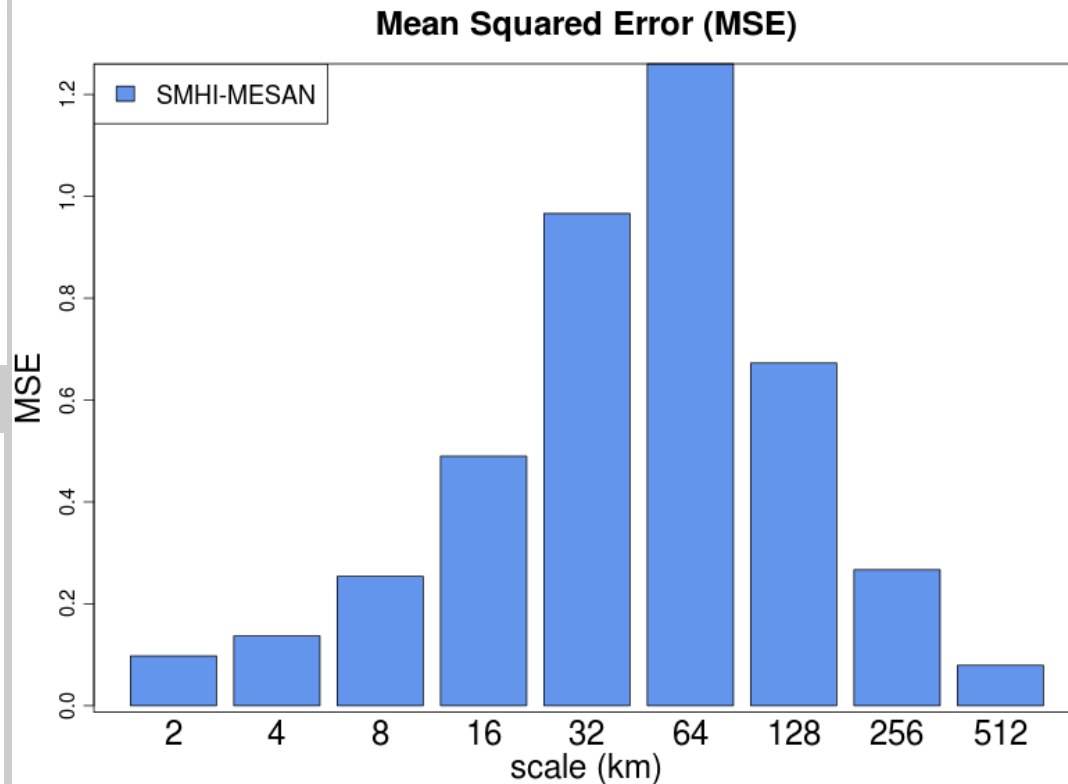
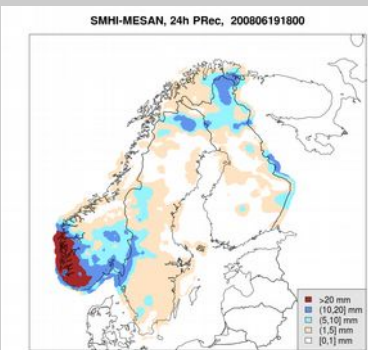
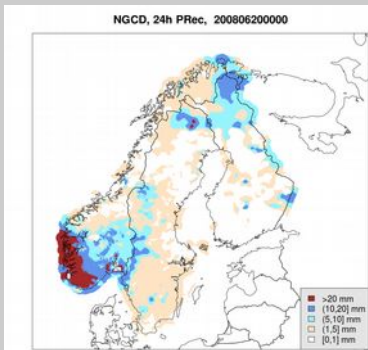
Case 2

case studies: RR for 1 day in summer 2008

Mean Squared Error(MSE) for each scale component

$$\text{MSE}(\mathbf{x}^{\text{rra}}, \mathbf{x}^{\text{ref}}) \equiv \text{En}^2(\mathbf{x}^{\text{rra}} - \mathbf{x}^{\text{ref}}) = \overline{(\mathbf{x}^{\text{rra}} - \mathbf{x}^{\text{ref}})^2}$$

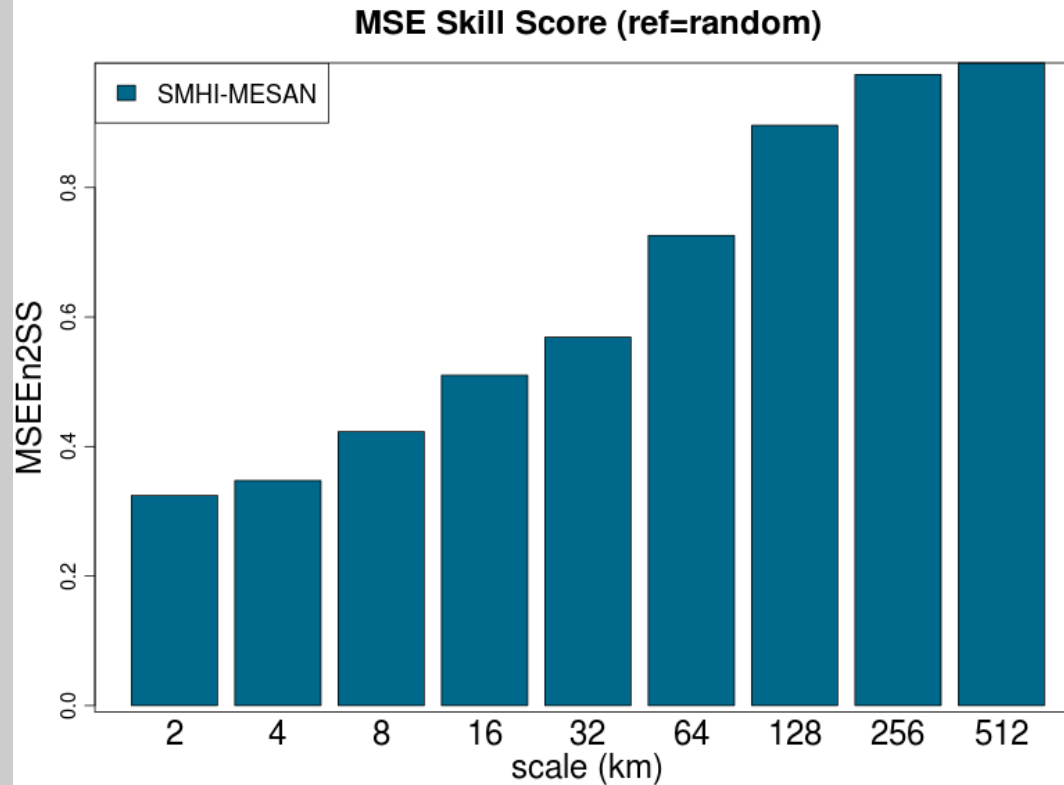
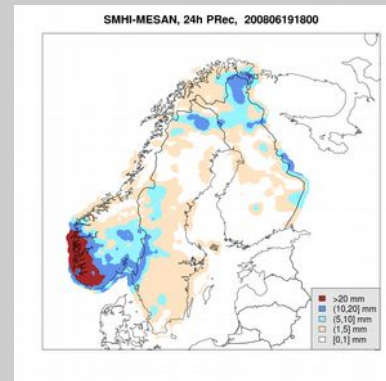
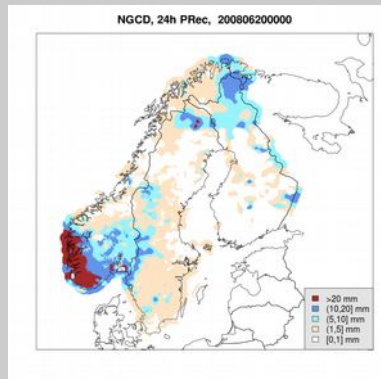
$$\text{MSE}_l(\mathbf{x}^{\text{rra}}, \mathbf{x}^{\text{ref}}) = \begin{cases} \langle \text{En}^2[W_l^m(\mathbf{x}^{\text{rra}}(t) - \mathbf{x}^{\text{ref}}(t))] \rangle & l = 1, \dots, L \\ \langle \text{En}^2[W_L^f(\mathbf{x}^{\text{rra}}(t) - \mathbf{x}^{\text{ref}}(t))] \rangle & l = L + 1 \end{cases}$$



case studies: RR for 1 day in summer 2008

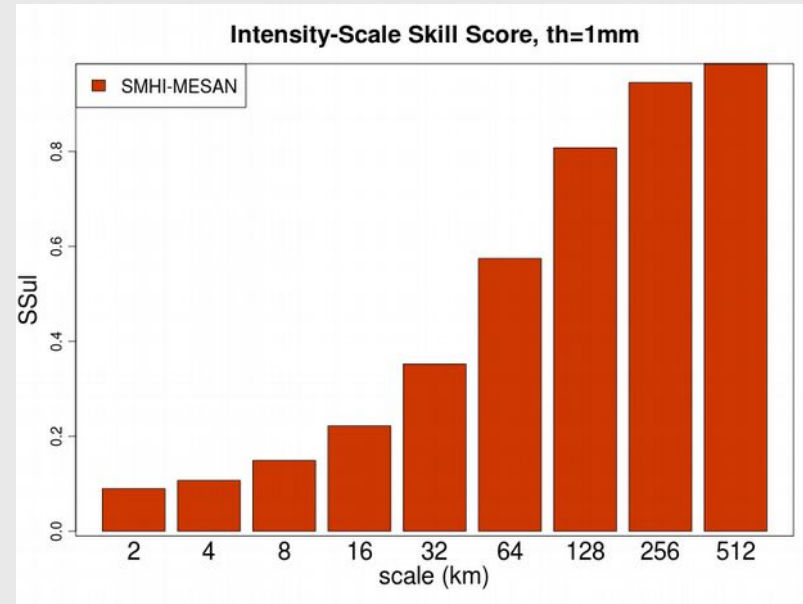
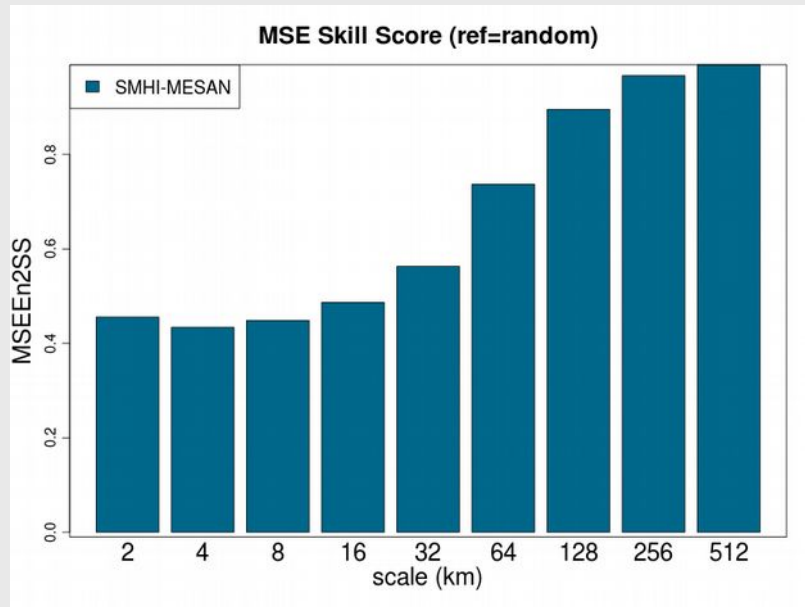
MSE skill-score for each scale component

$$SS_l \equiv \frac{\text{MSE}_l(\mathbf{x}^{\text{rra}}, \mathbf{x}^{\text{ref}}) - [\text{MSE}_l(\mathbf{x}^{\text{rra}}, \mathbf{x}^{\text{ref}})]_{\text{random}}}{[\text{MSE}_l(\mathbf{x}^{\text{rra}}, \mathbf{x}^{\text{ref}})]_{\text{best}} - [\text{MSE}_l(\mathbf{x}^{\text{rra}}, \mathbf{x}^{\text{ref}})]_{\text{random}}}$$



Scale-Separation verification technique

aggregation over JJA 2008



Conclusions. Wavelet based scale-separation MSE skill-score and scale-separation statistics are:

- informative on RRA bias, error and skill on different scales
- suitable for comparing models with different resolutions
- RRA performances for specific intensity events (thresholded binary fields)

UERRA-EVA: Evaluation software tools

<https://github.com/UERRA-EVA>

- a. Assess the added value of enhanced resolution in RRA:
wavelet-based scale-separation MSE skill score
→ *to be included in EVA_gridobs*
- b. Scale-decomposition of the Brier score for for the verification of probabilistic RRA...

EVA_gridobs

UERRA common evaluation procedure: assessing uncertainties in reanalysis by evaluation against gridded observational datasets

Reanalysis and gridded observations are assumed to be: on the same coordinate reference system and grid; same temporal aggregation.

Tip: use fimex for regridding (<https://wiki.met.no/fimex/start>)

List of Skill-Scores/Tests

Developed and tested for daily precipitation

- ☒ Probability Density Functions (PDFs) related skill-scores (PDFs are approximated by comparing discrete histograms):
 - ☒ difference between modes: $\text{mode}(\text{reanalysis}) - \text{mode}(\text{observation})$
 - ☒ relative precipitation biases = $(\text{mode}(\text{reanalysis}) / \text{mode}(\text{obsevation}) - 1) * 100\%$
 - ☒ overlapping skill-score (see [1], Eq.(1))
- ☒ two-sample Kolmogorov-Smirnov (K-S) test, or Smirnov test (see [2], Eqs (5.17-18))
- ☐ Fractional skill-score
- ☐ Optical-flow

[1] Mayer, S. et al. (2015). Identifying added value in high-resolution climate simulations over Scandinavia. Tellus A

[2] Wilks, D. S. (2011). Statistical methods in the atmospheric sciences (Vol. 100). Academic press.

Thanks for your attention!



cristianl@met.no