

UEF 2019

Quantifying uncertainties and confidence level in ATM simulations – scientific aspects*

C. Maurer¹, D. Arnold¹, M. Haselsteiner¹, J. Brioude², L. Haimberger³, and F. Weidle¹

¹ Central Institute for Meteorology and Geodynamics (ZAMG), Vienna, Austria

² Atmosphere and Cyclone Lab (LACy), University of La Réunion Island, France

³ Institute for Meteorology and Geophysics (IMGW), University of Vienna, Austria



ZAMG
Zentralanstalt für
Meteorologie und
Geodynamik

*work performed under the CTBTO awarded contract for “Provision of Software Engineering Services for Atmospheric Transport Modelling, Data Acquisition, Processing and Dissemination” under funding from the European Union Council Decision VII.

1. Comprehensive Nuclear Test-Ban Treaty (CTBT)

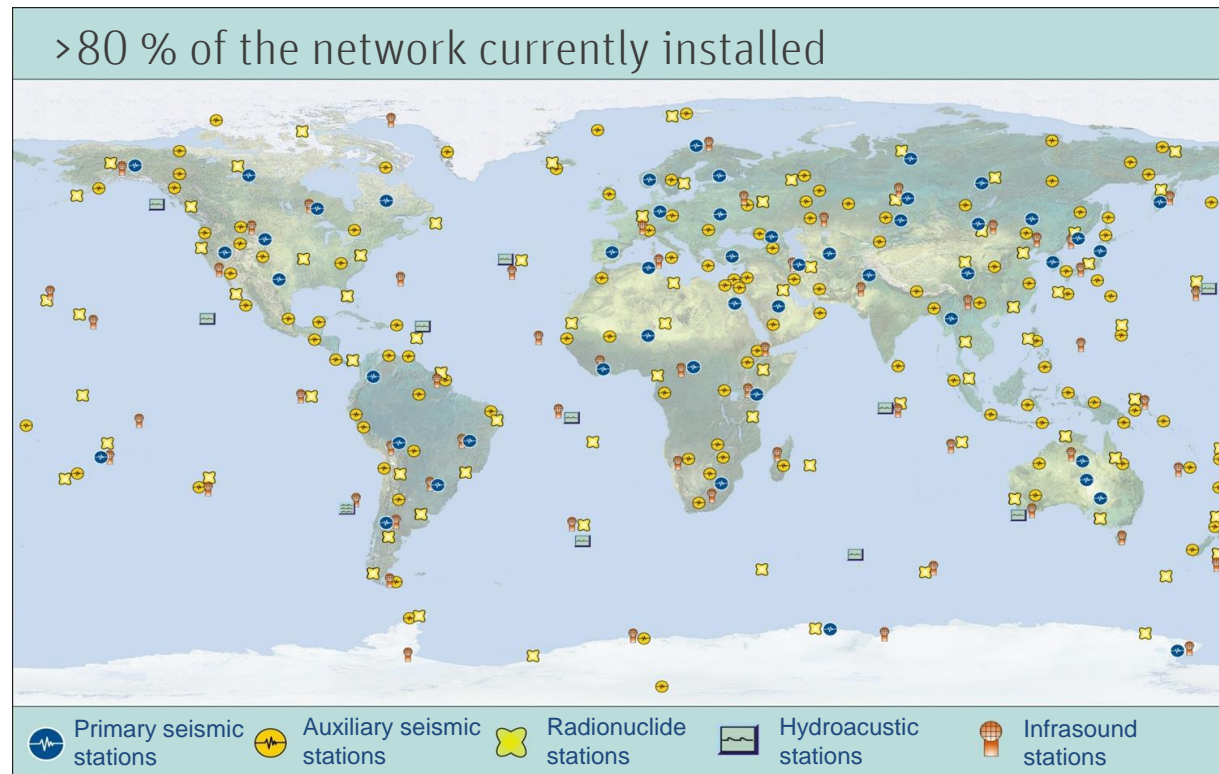
ARTICLE I - BASIC OBLIGATIONS

*Each State Party undertakes **not to carry out any nuclear weapon test explosion or any other nuclear explosion**, and to prohibit and prevent any such nuclear explosion at any place under its jurisdiction or control.*

Entry into force: As soon as 44 countries which operate nuclear reactors have ratified the Treaty (not just signed).

The International Monitoring System (IMS)

80 radionuclide stations foreseen, whereof 40 will be equipped with noble gas measuring devices



2. Aim



“CTBTO wants to develop tools to estimate uncertainties and confidence levels in (operational) ATM simulations using meteorological data from an Ensemble Prediction System (EPS). New tools will become part of CTBTO ATM pipeline in the second phase of this project.”

Tasks:

1. Perform a scientific literature review and define requirements.
2. Prepare a plan in consultation with the ATM team.
3. Define a set of products to display uncertainties and confidence levels in ATM simulations.
4. Prepare data to estimate uncertainties.
5. Run a prototype to generate products to display uncertainties and confidence levels.
6. Produce a short report including recommendations for the second phase of the project.

3.1. Literature review for EPS-based ATM



Main messages:

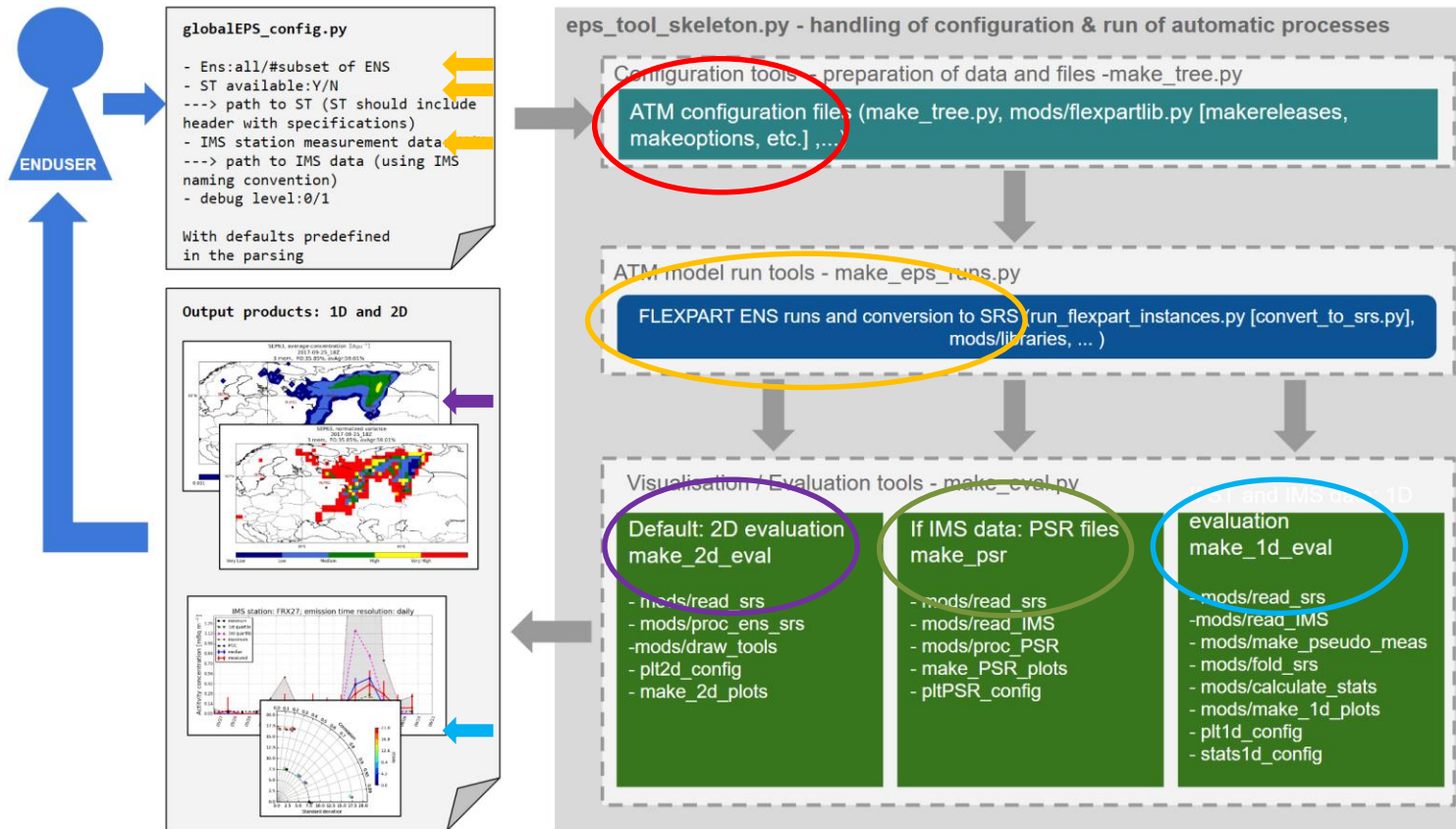
- If computationally affordable, ensemble dispersion modelling is the method of choice to take uncertainties of models and measurement data into account and create probabilistic forecasts.
- Authors tend to use the full ensemble (51 members - 50 plus control run - in case of ECMWF) or, to reduce resources, draw a (mostly randomly) subset. However, accuracy of forecast probabilities in certain scenarios clearly increases with the ensemble size.
- Clustering/finding representative members of EPS input for the purpose of ATM remains a research topic with no final solution. So far only special cases have been considered (e.g., high altitude transport of volcanic ash with strong horizontal forcing). Moreover, clustering is always confined to a certain region (e.g., Europe).
- No redundancy has to be expected due to the design of the NWP ensemble system – which may not be true for multi-model ensembles.

3.2 Literatur review for EPS-based ATM

Display and evaluation:

- *Agreement in Threshold Level, Agreement in Percentile Level, Overlap Plots, Maximum Concentration Plots* and *Box Plots* are most common.
- Many of the usual statistical scores suffer from important drawbacks. The main judgment of the model performance should not be based on the correlation coefficient. Also, statistics can be heavily influenced if some modelled values are near zero while nuisance sources cause the measured values to be at or just above a detection limit. To mitigate the deficiencies of independent statistical measures, a combined Rank measure should be considered.

4. System Prototype design

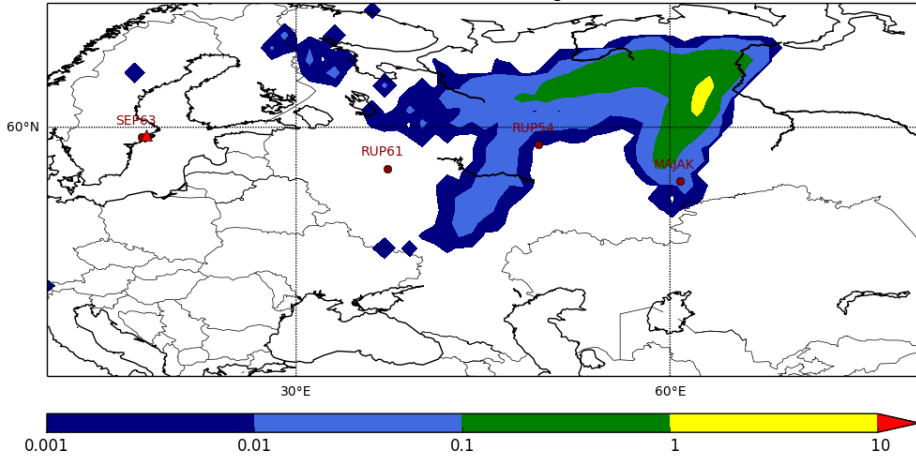


- Namelist file for configuration
- Launch atmospheric transport model FLEXPART runs based on (randomly selected or full) ECMWF-EPS input at selected locations, either the source location (fwd) or any (or all) of the IMS stations (bwd).
- Postprocess FLEXPART results.

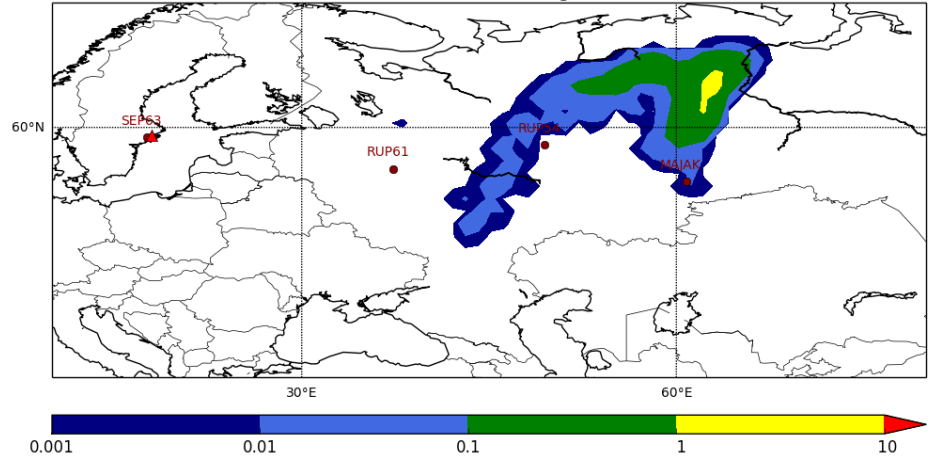
- Plot (ensemble) Source Receptor Sensitivity (SRS) fields/Fields of Regard (FOR).
- Calculate and plot the (ensemble) Probable Source Region (PSR)-fields if measurements are available and can be correlated with SRS values for each grid point and time step.
- Fold SRS fields with a source term if both the source term and measurements are available, plot resulting time series and perform statistical evaluation.

5.1. 2D-Products (SRS)

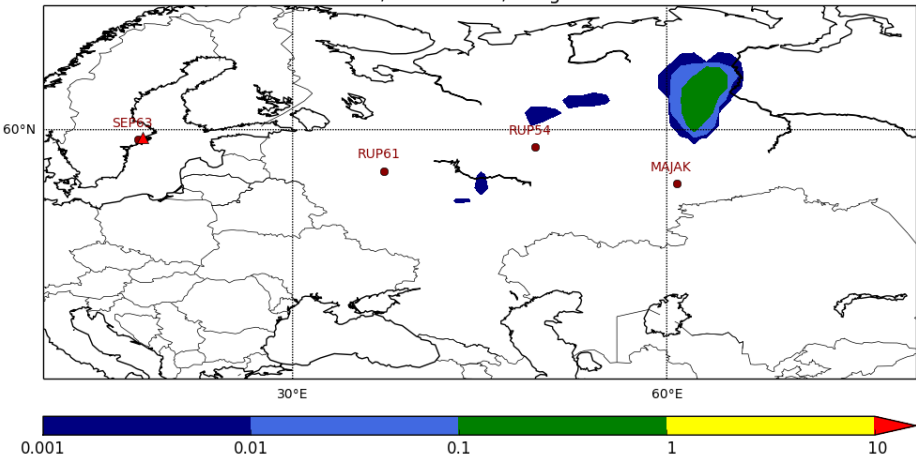
SEP63, average-concentration [Bq m^{-3}]
2017-09-25_18Z
3 mem, FO:35.85%, avAgr:59.01%



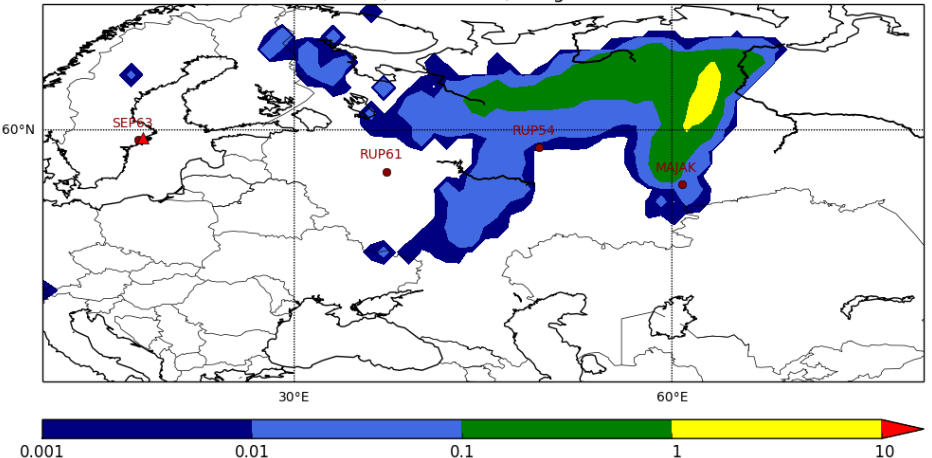
SEP63, median-concentration [Bq m^{-3}]
2017-09-25_18Z
3 mem, FO:35.85%, avAgr:59.01%



SEP63, min-concentration [Bq m^{-3}]
2017-09-25_18Z
3 mem, FO:35.85%, avAgr:59.01%

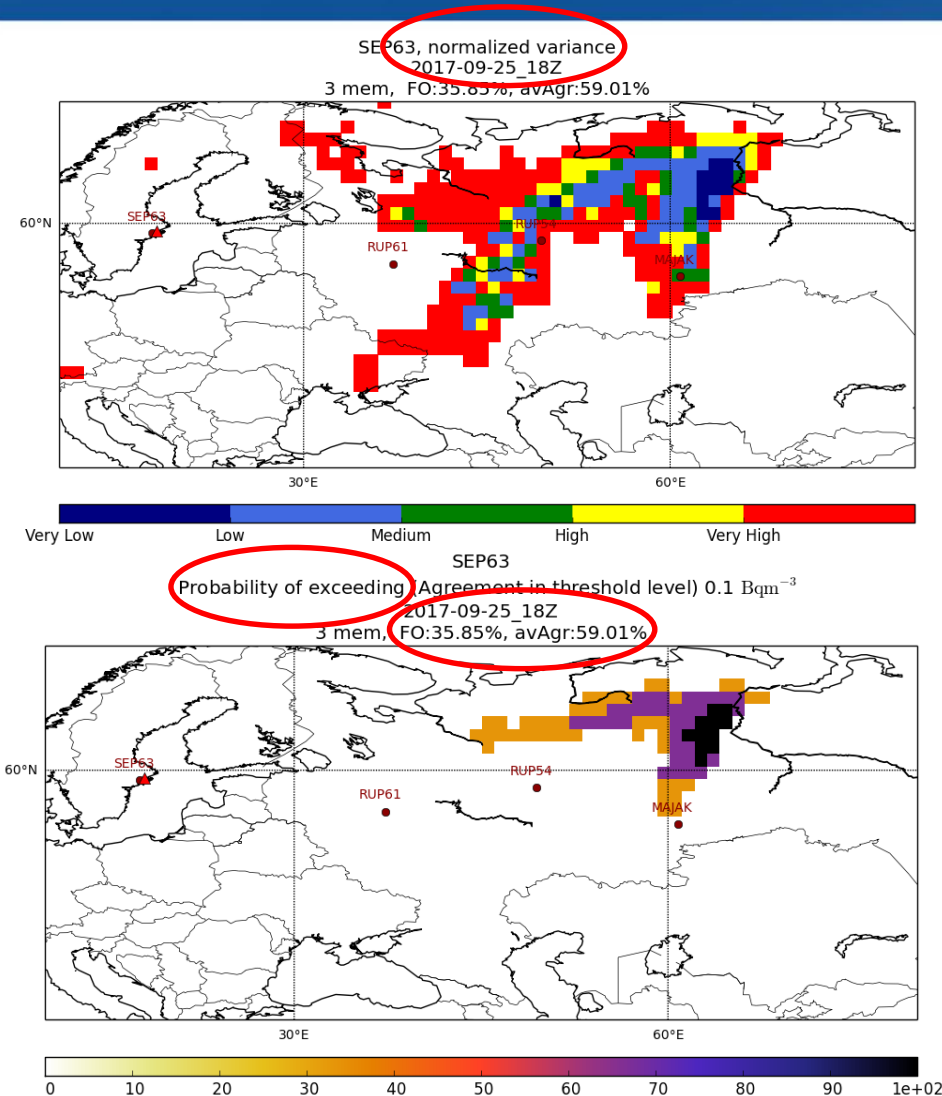


SEP63, max-concentration [Bq m^{-3}]
2017-09-25_18Z
3 mem, FO:35.85%, avAgr:59.01%



Prototype design based on the „Ruthenium event in Sept./Oct. 2017

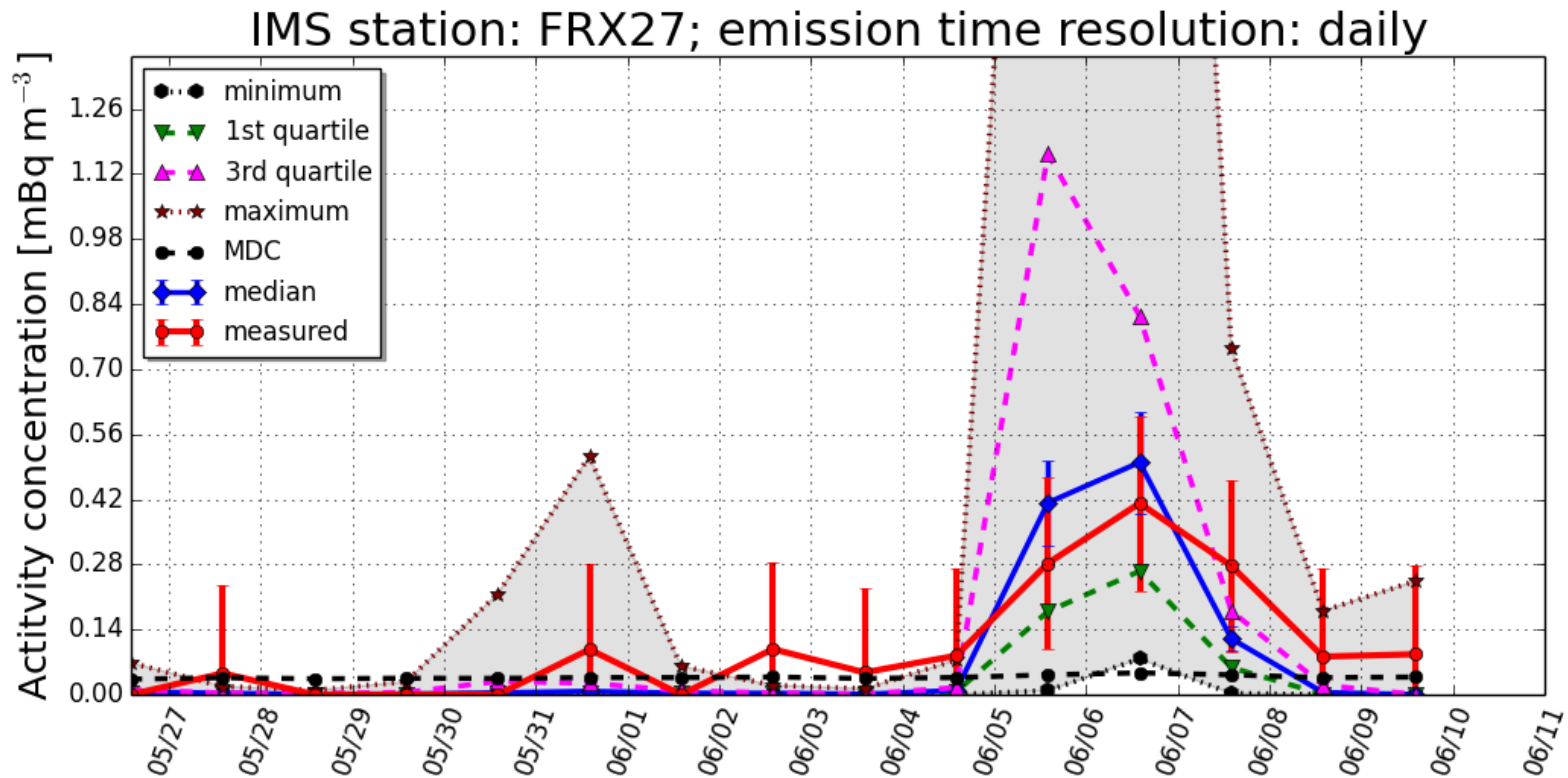
5.2. 2D-Products (SRS)



- Both **median** and **mean** are used according to literature in the context of ensemble dispersion modelling. However, **median** is preferred.
- Minimum** and **maximum** indicate extreme events, which may not be visible using the median display.
- The **normalized variance** (normalized with the squared mean value, maximum value N-1) provides a **rough estimate of uncertainty**.
- The **probability of exceedance (Agreement in Threshold Level)** includes a threshold definition in order to **focus on values above a certain level**.

„Figure of Overlap“ and „Average Agreement“ as indicator of consistency between ensemble members

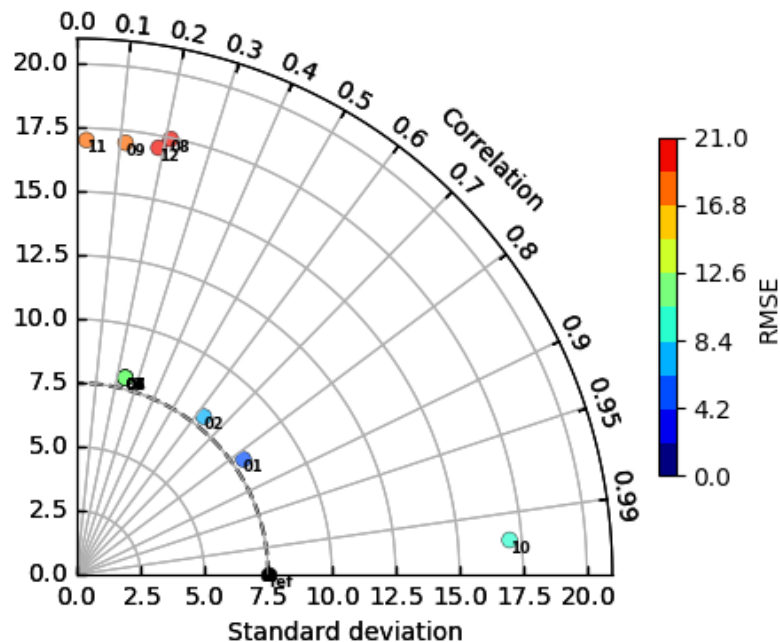
6. 1D-Products (concentration time series) & numerical statistics



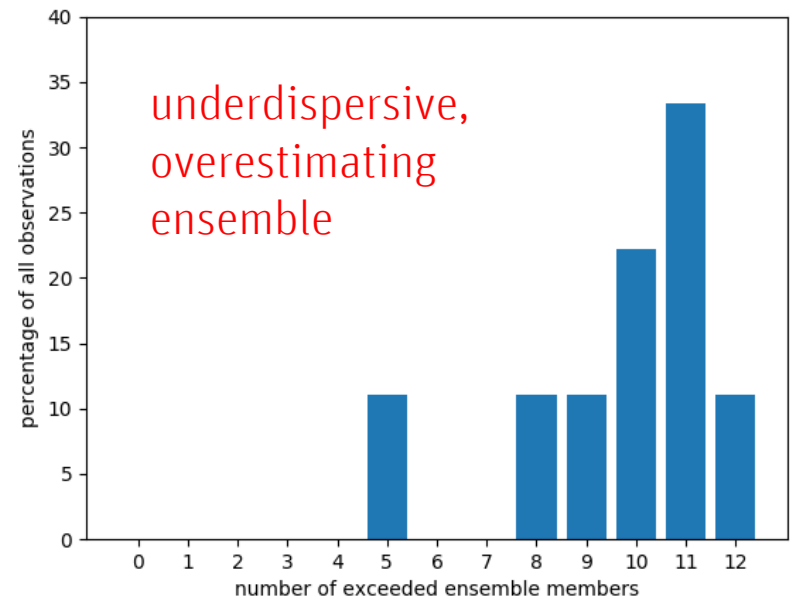
Statistics:

- R^2 , the Fractional Bias, the Kolmogorov-Smirnov Parameter and the Accuracy (based on MDC) will be given individually and combined into a Rank measure.
- Brier Score based on MDC will be calculated.
- An alternative skill score S (personal communication with P. Seibert) will be tested.

7. Graphical statistical evaluation



Taylor Diagram: Data points represent the time series of a member or an ensemble statistic at a measurement station.

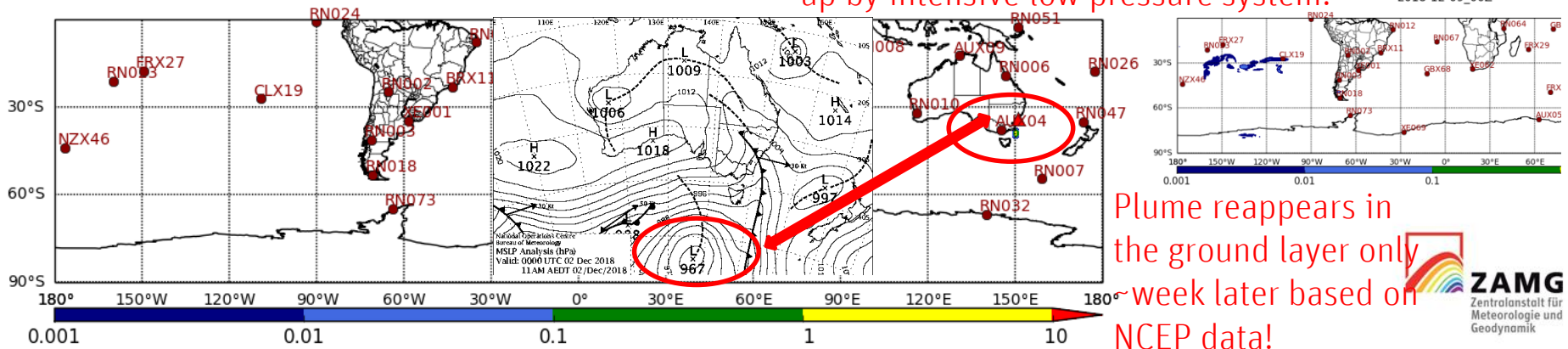


Talagrand Diagram: Evaluating the structure of the ensemble. Ideally one wants a flat diagram.

8. Test Cases (real & synthetic)

1. Real case (11.3.-25.3.2011) - Fukushima NPP accident: Real IMS Cs-137 and Xe-133 measurements and source term from *Stohl et al.* (2012).
2. Synthetic cases (1.12.2018-15.12.2018), producing pseudo-measurements based on FLEXPART and NCEP-GFS data:
 1. Hypothetical ANSTO puff release (1E15 Bq Xe-133, Dec., 1st, 00:00-01:00 UTC, summer time conditions)
 2. Hypothetical DPRK-test site puff release (1E15 Bq Xe-133, Dec., 1st, 00:00-01:00 UTC, winter time conditions)

Xe-133, activity-concentration [Bq m^{-3}] Hypothetical ANSTO plume gets sucked up by intensive low pressure system!



9. ECMWF EPS data

- Update of flex_extract_v7.0.3 (www.flexpart.eu) for the current purpose.
- **Re-caluculation of the full ensemble** (12-hourly short-term forecasts) for the real test case (from 2011). **Perturbed forecasts** on model level fields are available only 30 days after creation!
- Using the **Ensemble Analysis** product (available at 00/06/12/18 UTC, *stream=ELDA*,) for the synthetic test cases: 25 members, **observations and model tendencies perturbed with stochastic noise**. Another 25 symmetric members are added by calculating the differences between each member and the mean.
- Combining analysis (6- or 12-hourly) and short term-forecasts (up to 9 hours).
- Perform **test cases both with full ensemble and a 10 member subset**.
- ERA5-10-members ensemble is no alternative, due to not being real-time data.

11.1. PSR fields for the synthetic DPRK test case

*Can we confine the PSR
at known release time
better than with the
operational run only
based on the metrics as
introduced above?*

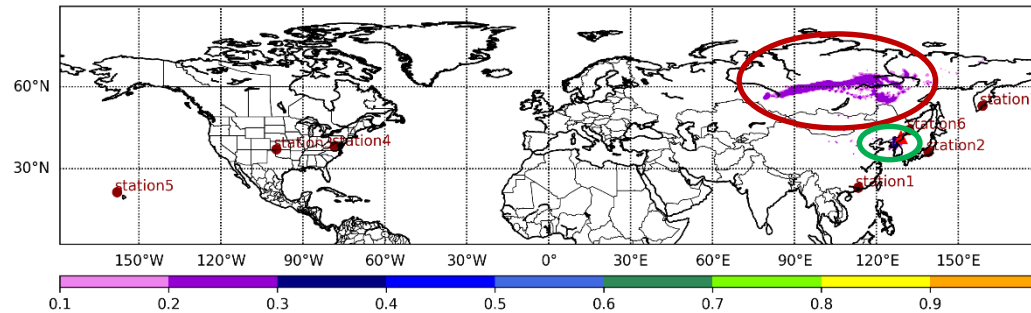
Improvement +

Wrong source
location

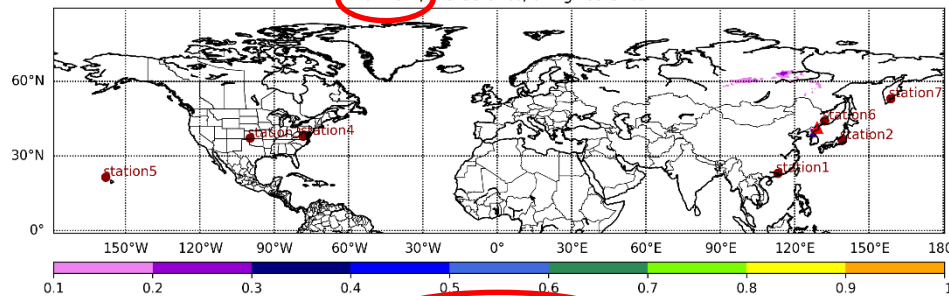
Correct
source
location

Deterioration —

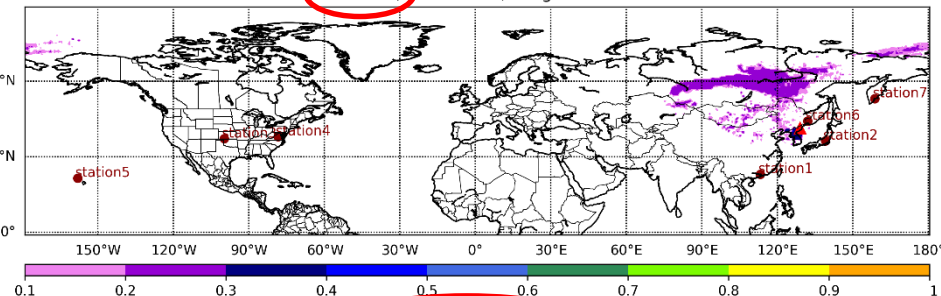
DPRK, correlation
2018-12-01_00Z
hRes run of synthetic case



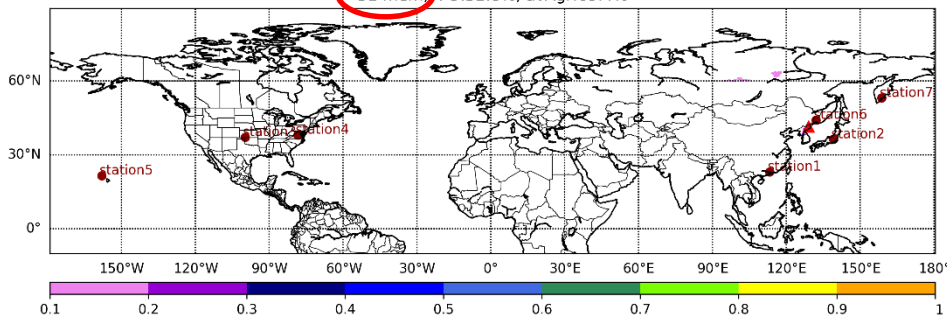
DPRK, min-correlation
2018-12-01_00Z
10 mem, FO:33.37%, avAgr:69.32%



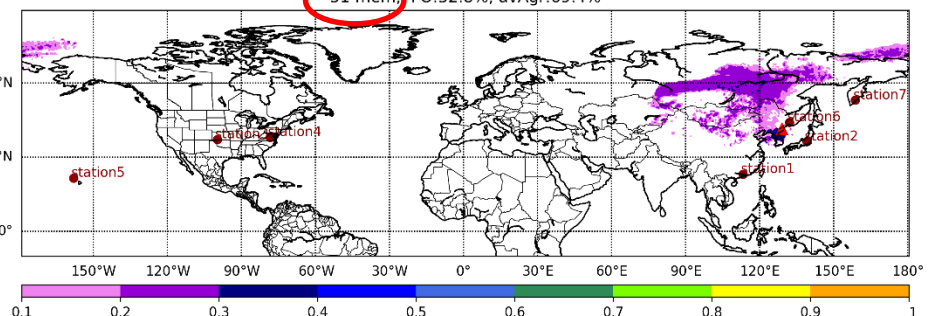
DPRK, max-correlation
2018-12-01_00Z
10 mem, FO:33.37%, avAgr:69.32%



DPRK, min-correlation
2018-12-01_00Z
51 mem, FO:32.8%, avAgr:69.4%

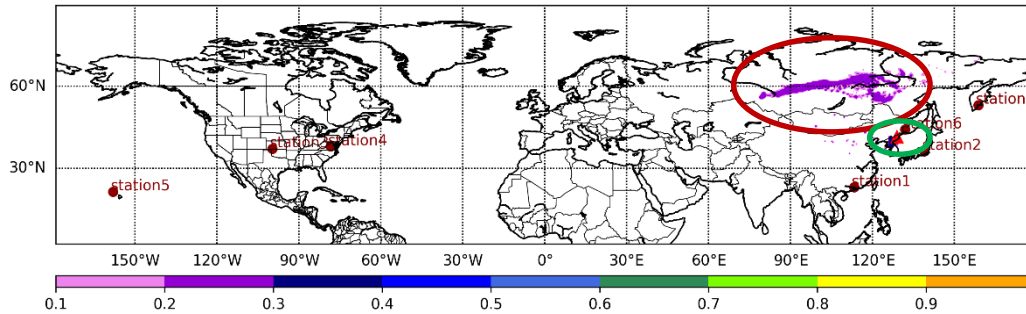


DPRK, max-correlation
2018-12-01_00Z
51 mem, FO:32.8%, avAgr:69.4%



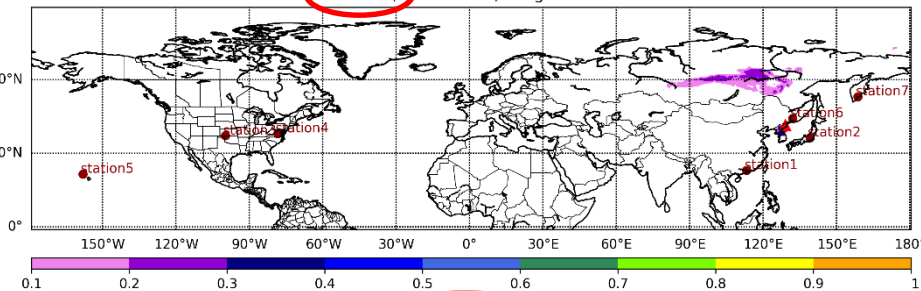
11.2. PSR fields for the synthetic DPRK test case

DPRK, correlation
2018-12-01_00Z
hRes run of synthetic case

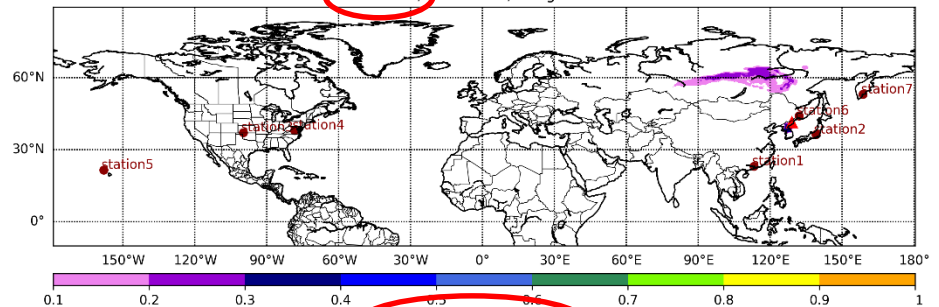


Wrong source
location
Correct
source
location

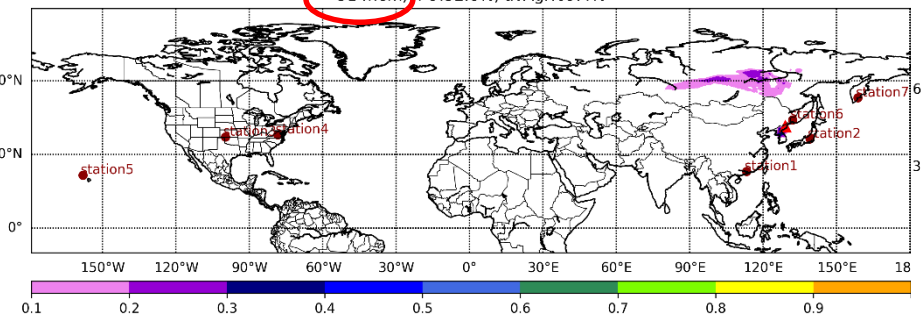
DPRK, average-correlation
2018-12-01_00Z
10 mem, FO:33.37%, avAgr:69.32%



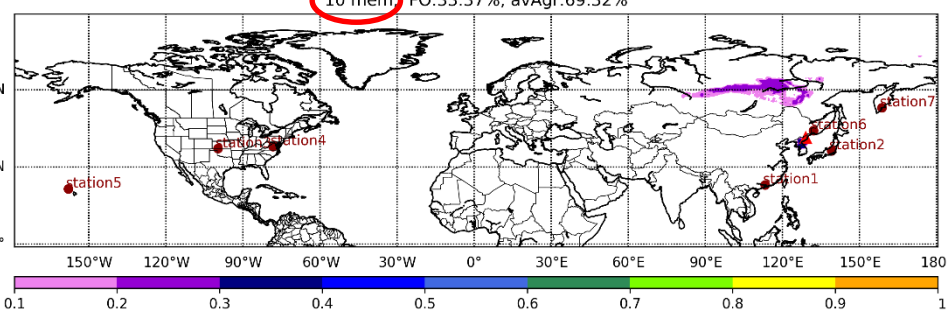
DPRK, median-correlation
2018-12-01_00Z
51 mem, FO:32.8%, avAgr:69.4%



DPRK, average-correlation
2018-12-01_00Z
51 mem, FO:32.8%, avAgr:69.4%

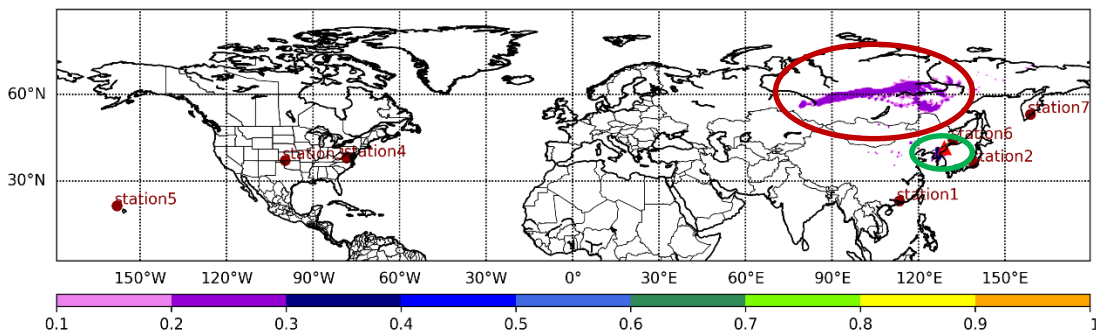


DPRK, median-correlation
2018-12-01_00Z
10 mem, FO:33.37%, avAgr:69.32%



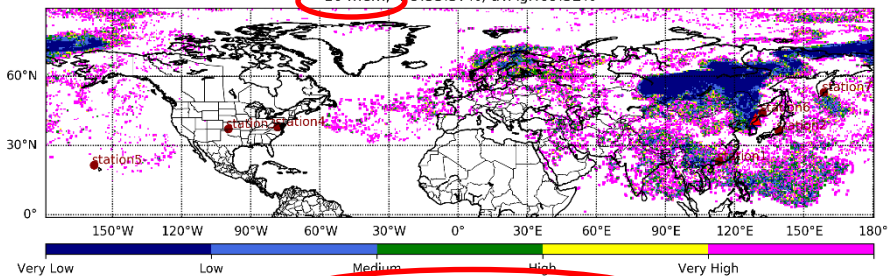
11.3. PSR fields for the synthetic DPRK test case

DPRK, correlation
2018-12-01_00Z
hRes run of synthetic case

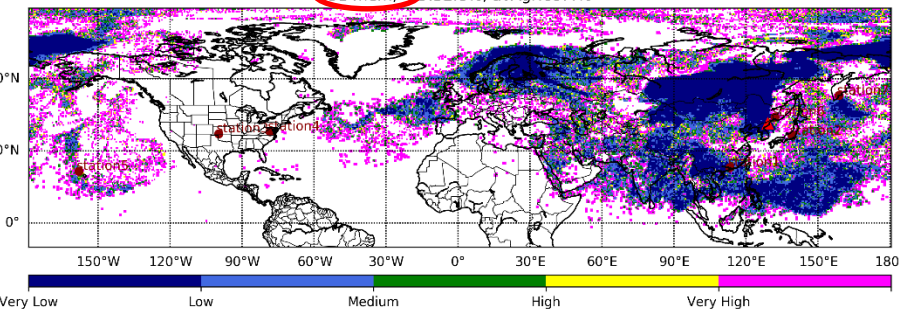


Wrong source
location
Correct
source
location

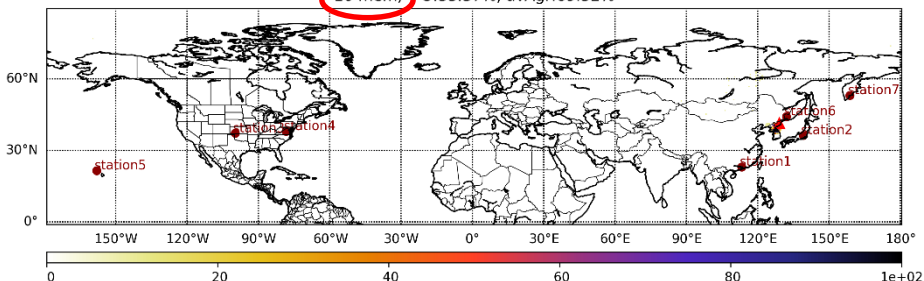
DPRK, normalized variance of correlation
2018-12-01_00Z
10 mem, FO:33.37%, avAgr:69.32%



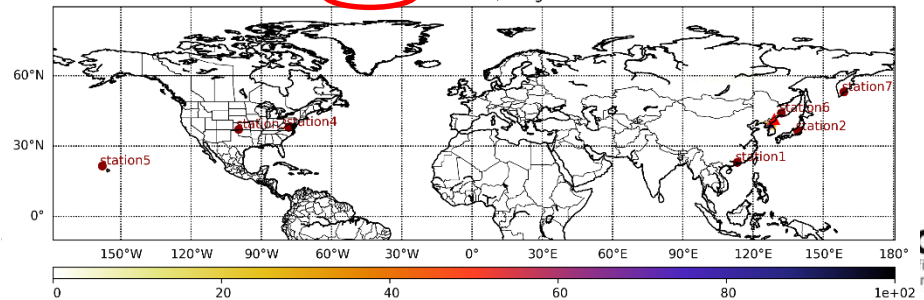
DPRK, normalized variance of correlation
2018-12-01_00Z
51 mem, FO:32.8%, avAgr:69.4%



DPRK corr
Probability of exceeding (Agreement in threshold level) 0.3
2018-12-01_00Z
10 mem, FO:33.37%, avAgr:69.32%

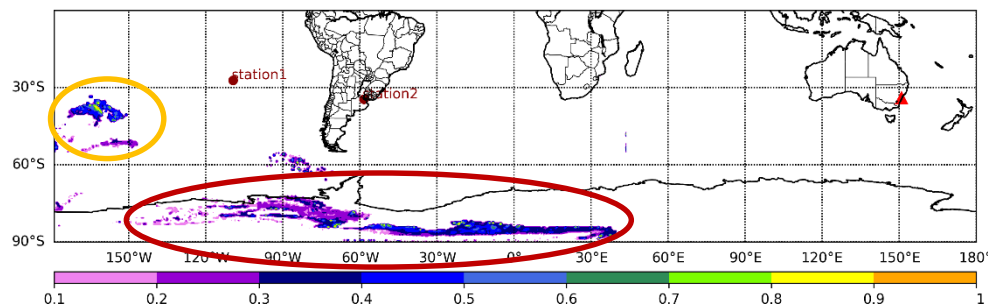


DPRK corr
Probability of exceeding (Agreement in threshold level) 0.3
2018-12-01_00Z
51 mem, FO:32.8%, avAgr:69.4%



12.1. PSR fields for the synthetic ANSTO test case

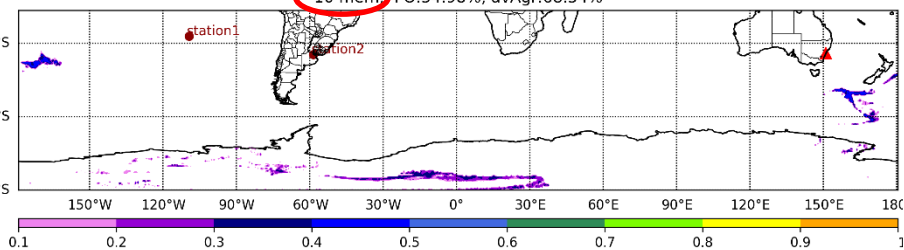
ANSTO, correlation
2018-12-01_00Z
hRes run of synthetic case



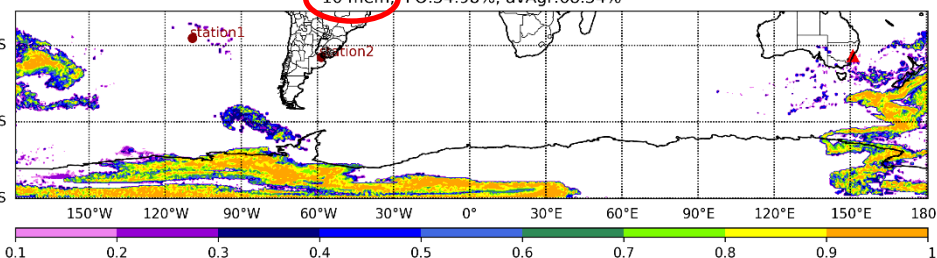
Wrong source
location
Approximate
source
location



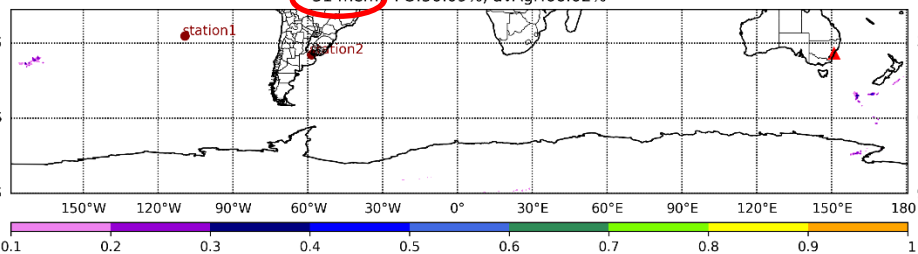
ANSTO, min-correlation
2018-12-01_00Z
10 mem FO:54.98%, avAgr:68.54%



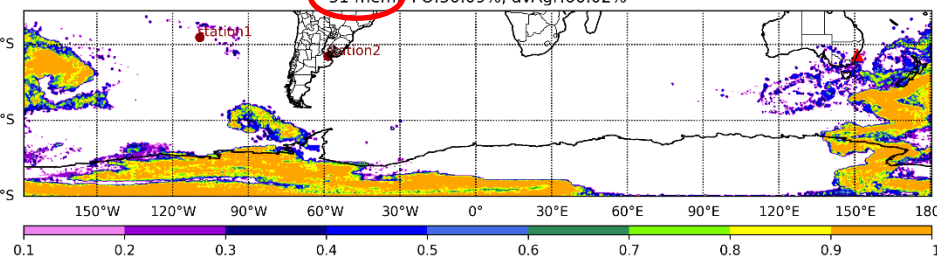
ANSTO, max-correlation
2018-12-01_00Z
10 mem FO:54.98%, avAgr:68.54%



ANSTO, min-correlation
2018-12-01_00Z
51 mem FO:50.09%, avAgr:66.02%

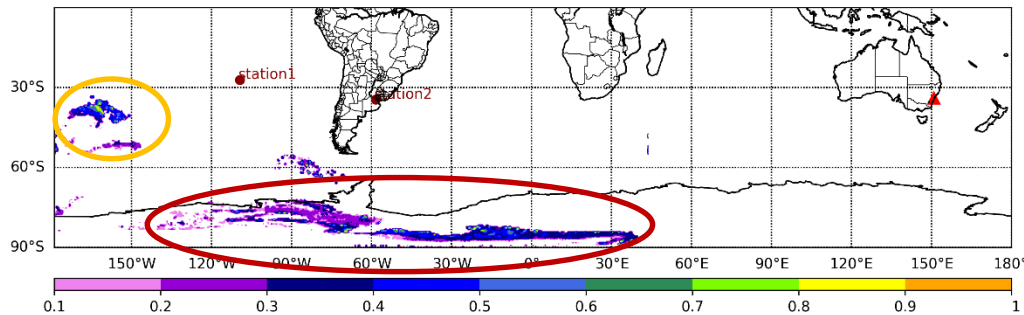


ANSTO, max-correlation
2018-12-01_00Z
51 mem FO:50.09%, avAgr:66.02%



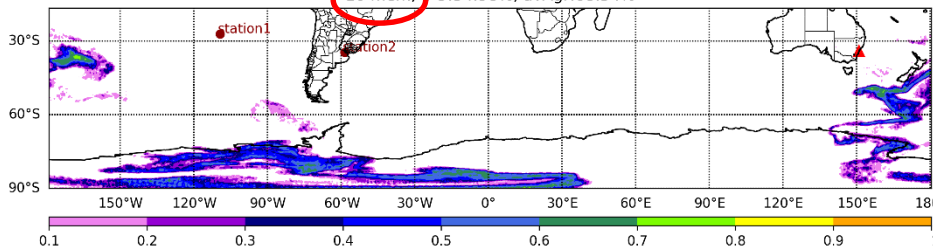
12.2. PSR fields for the synthetic ANSTO test case

ANSTO, correlation
2018-12-01 00Z
hRes run of synthetic case

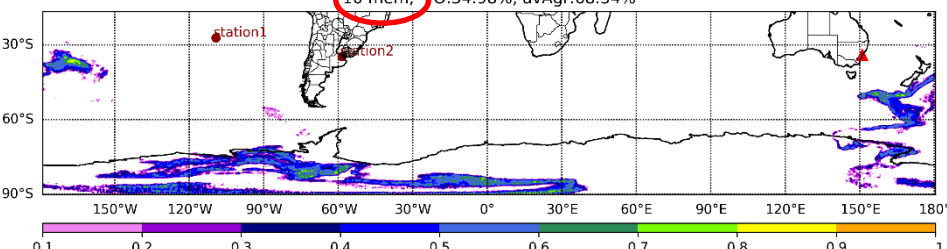


Wrong source
location
Approximate
source
location

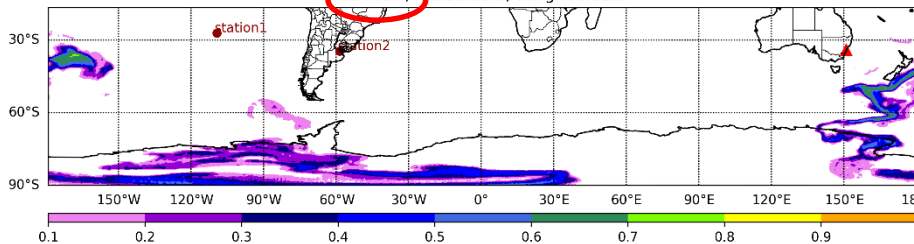
ANSTO, average-correlation
2018-12-01 00Z
10 mem, FO:54.98%, avAgr:68.54%



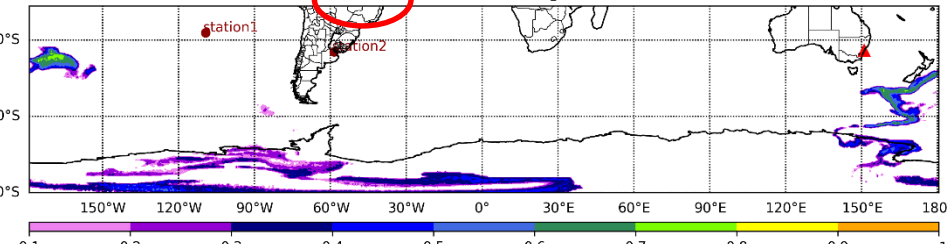
ANSTO, median-correlation
2018-12-01 00Z
10 mem, FO:54.98%, avAgr:68.54%



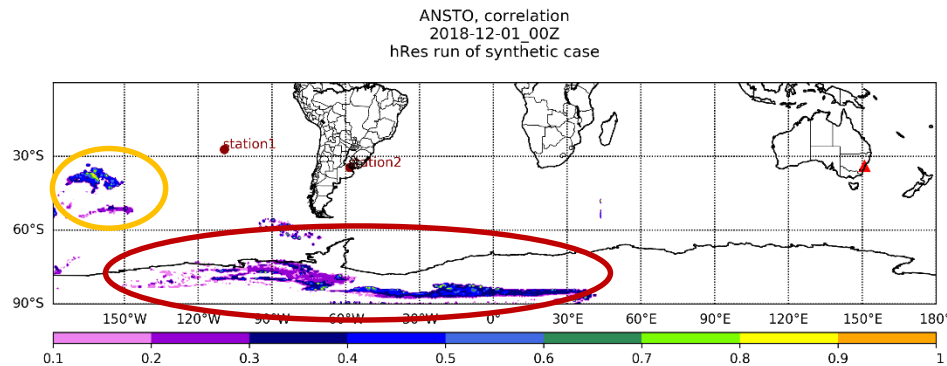
ANSTO, average-correlation
2018-12-01 00Z
51 mem, FO:50.09%, avAgr:66.02%



ANSTO, median-correlation
2018-12-01 00Z
51 mem, FO:50.09%, avAgr:66.02%



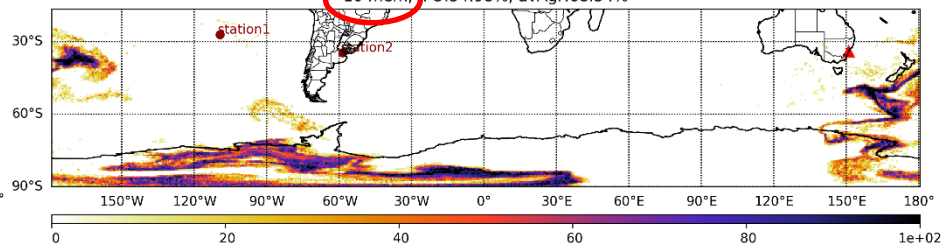
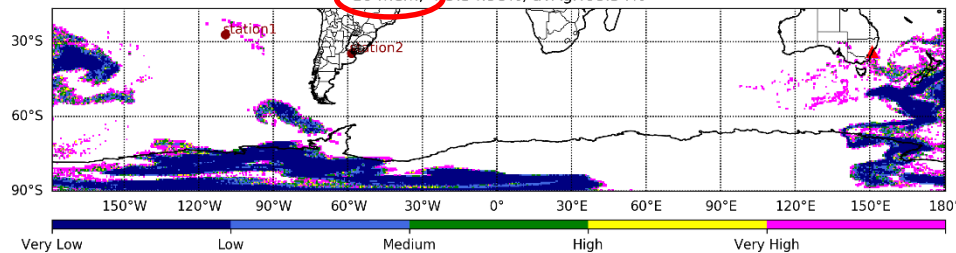
12.3. PSR fields for the synthetic ANSTO test case



Wrong source
location
Approximate
source
location

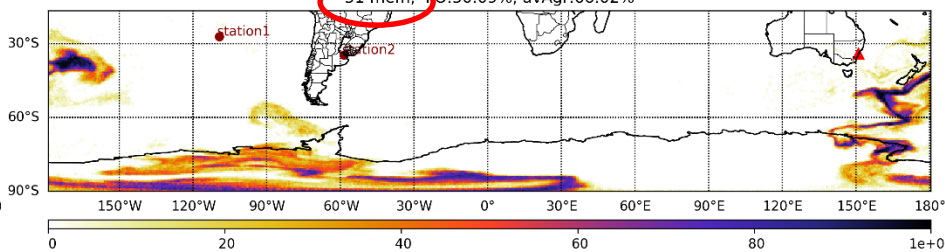
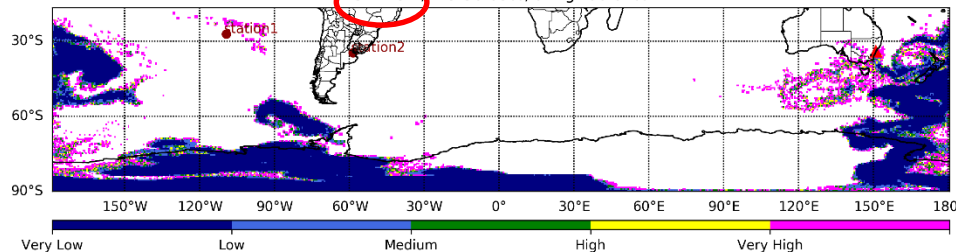
ANSTO, normalized variance of correlation
2018-12-01 00Z
10 mem, FO:54.98%, avAgr:68.54%

ANSTO
Probability of exceeding (Agreement in threshold level) 0.3
2018-12-01 00Z
10 mem, FO:54.98%, avAgr:68.54%



ANSTO, normalized variance of correlation
2018-12-01 00Z
51 mem, FO:50.09%, avAgr:66.02%

ANSTO
Probability of exceeding (Agreement in threshold level) 0.3
2018-12-01 00Z
51 mem, FO:50.09%, avAgr:66.02%



13.1. Preliminary conclusions

- The **best PSR ensemble metric** turns out to be the **ensemble minimum** as only those **parts of the PSR fields are kept** which are **at** (DPRK test case) **or near** (ANSTO case) to the **actual source** and **no threshold** needs to be defined. The **performance increases** when **switching from a 10 member to the full ensemble**. Whereas when dealing with **activity concentrations** the "**worst case**" (**ensemble maximum**) is extremely **relevant**, it is **not so** for PSR fields.
- **Probability of exceedance** plots applied to PSR fields **suffer from using a fixed threshold**. This is **unlike the situation with activity concentration values**, where a generally useful threshold is easier to define. In essence **thresholds would need to be redefined for every case** based on those values **deduced from the operational run**. For the DPRK test case this prerequisite was met by chance.

13.2. Preliminary conclusions

- **Average and median results** are **very similar**. They show a **small improvement** in terms of being able to locate the source region **for the DPRK test case**.
- The concept and the **usability of the "Normalized Variance"** when applied to PSR fields **needs to be re-visited**. Probably a **threshold for confining the region of interest to PSR values** above a certain value would be **necessary**. Variances for very tiny PSR values are not of interest.

THANK YOU FOR YOUR ATTENTION!

I. Additional material

have been skipped for an improper reason; personal communication). For the calculation of the case specific so-called *Agreement of Model p* with all others, $cAgreement_p$ from the case specific rank value of model p, $cRNK_p$, these authors preferred to give the percentage of the maximum $cRNK$ value (3.0) while excluding the trivial auto-correlation result as follows:

$$cAgreement_p = \frac{100}{3(N-1)} \sum_{i=1}^N \epsilon_{ip} cRNK_p \quad (4)$$

with N, total number of experiments' participants and $\epsilon_{ip} = 1$ for $i \neq p$ and 0 for $i = p$.

The belonging case and model specific anomaly, $cAnomaly_p$, to the case specific across participants' average agreement (cAV) is calculated as:

$$cAnomaly_p = cAgreement_p - cAV \quad (5)$$

with

$$cAV = \frac{1}{N} \sum_{p=1}^N cAgreement_p \quad (6)$$

II. Additional material



In order to assess probabilistic capabilities of the ensembles the *Brier Score* as function of time has been calculated for different thresholds in Galmarini et al. [16]. It is defined as the mean square error of a probability forecast:

$$BS = 1/n \sum_{i=1}^n (F_i - O_i)^2 \quad (13)$$

where n is the number of forecasts, F_i is the forecast probability on occasion i , O_i is the observation (0 or 1) on occasion i . The score weights larger errors more than smaller ones. De Meutter et al. [5] calculate the simulated probability as

$$F_i = \frac{t + 2/3}{n + 4/3} \quad (14)$$

where t is the number of members that simulate a value above a threshold and n is the total number of ensemble members n .

III. Additional material

the division by $M = 0$ is kept the $NMSE$ adopts infinity. Further, Seibert [34] states that a high value of the bias implies a high value of the RMSE, though the two data sets might just be shifted by some offset and otherwise agree quite well. Therefore, the definition of a third quantity, the bias-corrected RMSE, BC_RMSE , makes sense:

$$BC_RMSE = \frac{1}{n} \sqrt{\sum_i^N [(M_i - \bar{M}) - (O_i - \bar{O})]^2} \quad (10)$$

According to Sachs [33] (p. 128), this can also be calculated as:

$$BC_RMSE = \sqrt{RMSE^2 - B^2} \quad (11)$$

is then the BC_RMSE . According to Taylor [42] model standard deviation, BC_RMSE and R can be combined into a single skill score S_r .

$$S_r = 2(1 + R) \left(\frac{\sigma_m}{\sigma_o} + \frac{\sigma_o}{\sigma_m} \right)^2 \quad (20)$$

with σ_m and σ_o being the standard deviations of predictions and observations.

The correlation contribution becomes important for large values of BC_RMSE . If one wants to include also the relative bias FB into the skill score, Seibert [34] suggests:

$$S_b = \frac{1}{1 + \alpha FB^2} \quad (21)$$

A value of $\alpha = 10$ appears to give a relationship fulfilling Seibert's [34] subjective idea about such a skill score. Finally, both skill scores can be combined into a total skill score S :

$$S = \alpha S_r + (1 - \alpha S_b) \quad (22)$$

The value of α is rather arbitrary and would depend on the application. An additive and not multiplicative combination of the two scores is suggested because a model that has skill either in reproducing the mean or in reproducing the patterns should be attributed some total skill; the product of the two scores would be zero with one of the factors being zero. In any case, data sets with strongly non-normal distributions might better be transformed before applying any of the measures.