# Current Status of
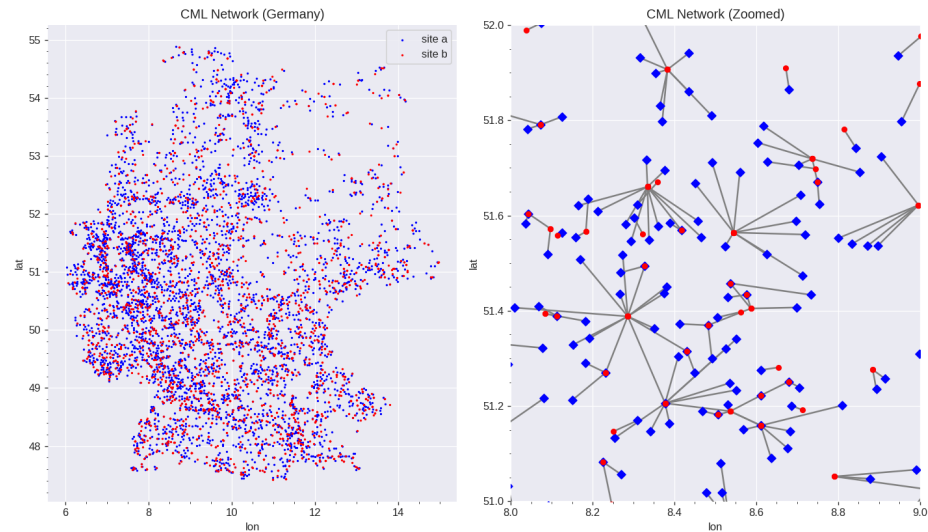# CML Data Assimilation
# and
# Targeted Covariance Inflation
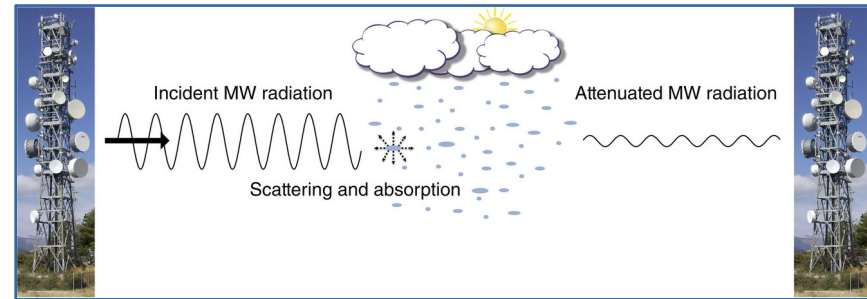
RealPEP

**Deutscher Wetterdienst**
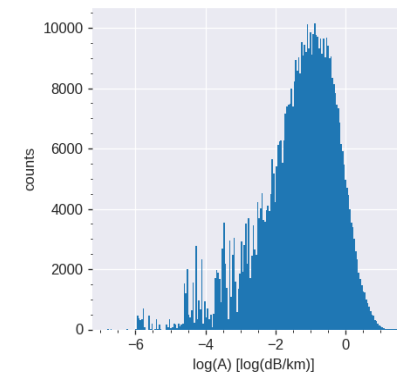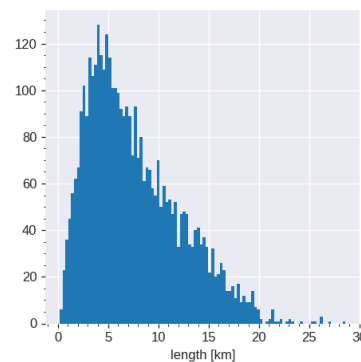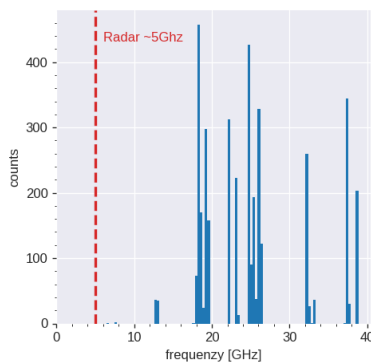**Wetter und Klima aus einer Hand**
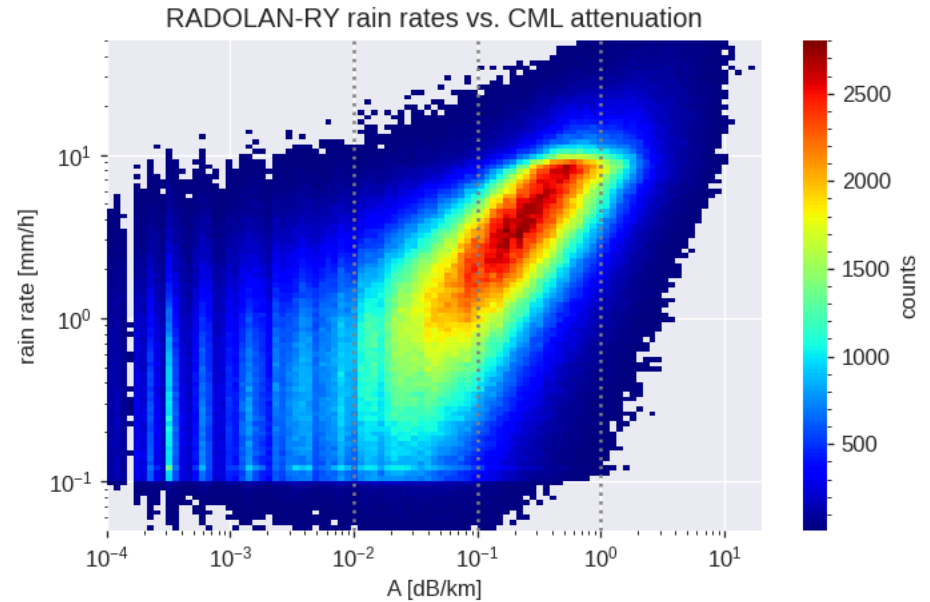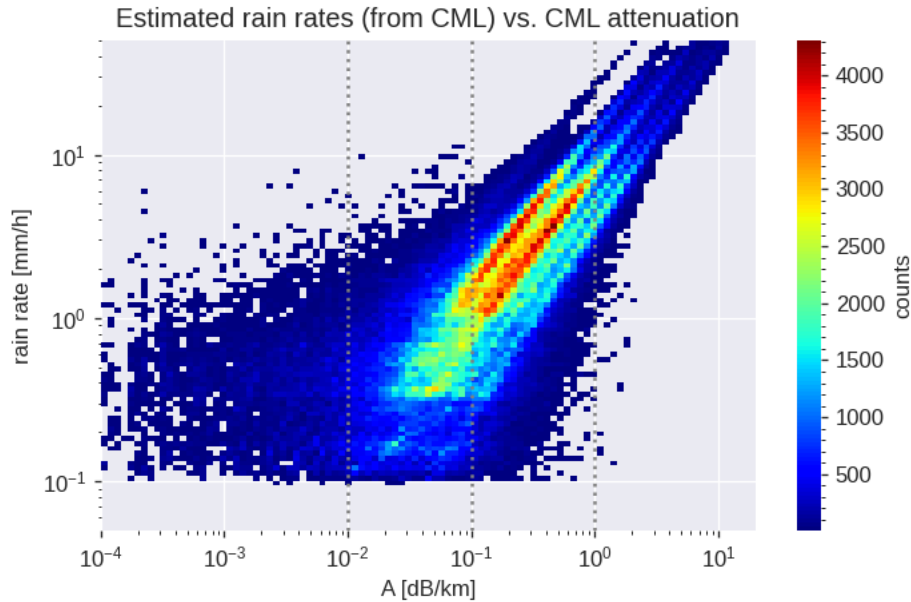
DWD

K.Vobig

# CML Basics

- **Commercial microwave link (CML)** data successfully employed for the estimation of rain rates (QPE) ($\rightarrow$C. Chwala, P1)

- overall objective here ($\rightarrow$P3): **data assimilation** of CML data in **numerical weather prediction models** for **improving QPF**

  - able to contribute to bridging the gap between QPN and NWP?

  - (How much) does it improve QPF?

  - How does it compare to Radar data assimilation?

- in the following: discussion of **technical details** of CML data assimilation and presentation of **first results**

- CMLs are used to interconnect cell phone towers

- each CML consists of **sender** and **receiver**

- transmitted radiation gets **attenuated** by (e.g.) raindrops

- ~4000 CMLs in current dataset for June 2019
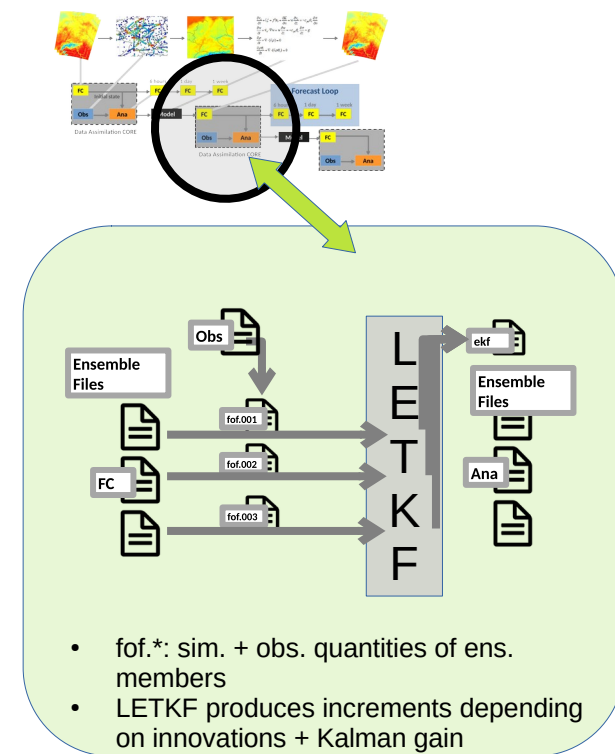
- temporal resolution 1min

- **CML frequency** significantly above DWD **Radar frequency** → different physics involved!

- use **path-integrated specific attenuation** for assimilation

  - referred to as A from now on

  - A [dB/km] = attenuation [dB] / distance [km]

  - direct **relationship of A with rain rate** (→ power law)
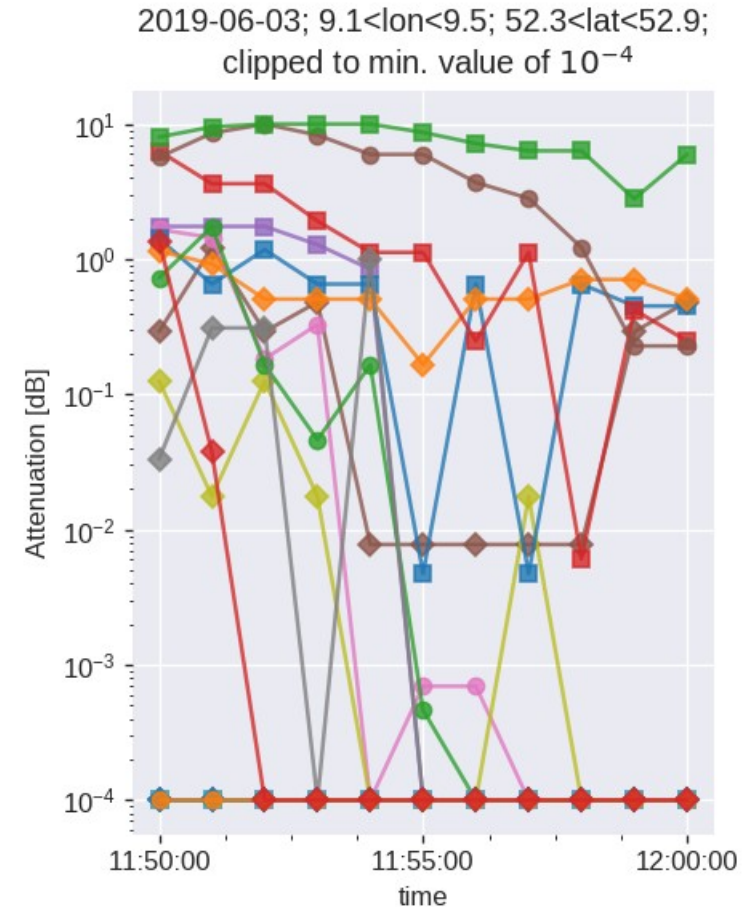
- most attenuations **very** small

Estimated rain rates (from CML) vs. CML attenuation

RADOLAN-RY rain rates vs. CML attenuation

- "linear" relationship (on double logarithmic scale)
  → hint at underlying power law

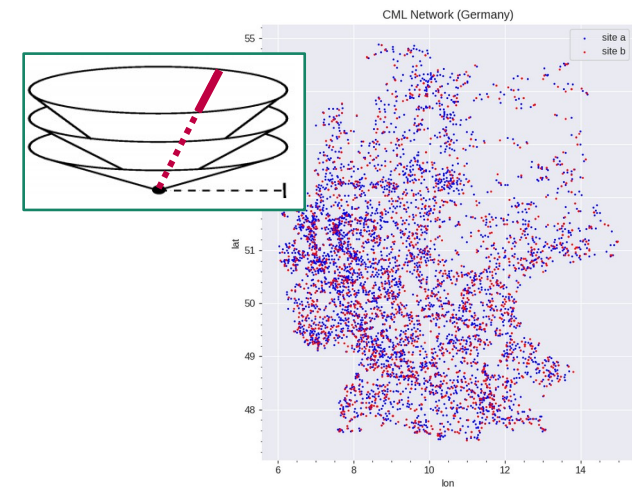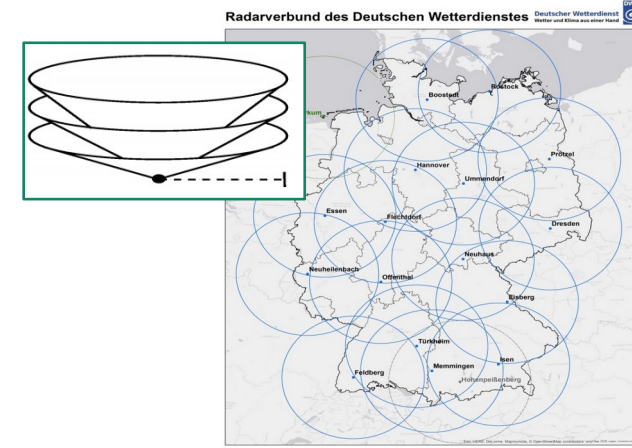- (very) noisy data for A < ~$10^{-2}$ dB/km

- for assimilating data **feedback/fof files** have to be generated

- each (ens.) fof file contains all data relevant to LETKF assimilation process (specific date)

- particularly, for each observation there has to be a **simulated model equivalent**

- **built automated system** for the construction of **CML feedback files**

  - includes all necessary data processing steps

  - implemented (mostly) in Python

  - integrated into new BACY experiment



- fof.*: sim. + obs. quantities of ens. members
- LETKF produces increments depending on innovations + Kalman gain

- **temporal superobbing/smoothing**:

  - for an assimilation at $t_0$ calculate the mean of all observations falling within a 10 min time window $[t_0 - 10 \text{ min}, t_0]$ for each CML

  - **smooths** out **erratic fluctuations** of attenuations

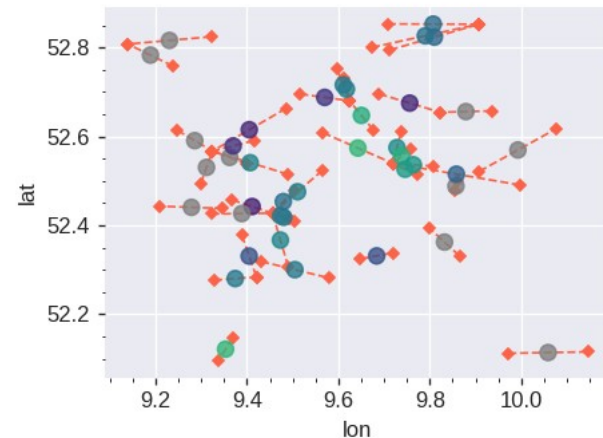- **outlook**: also perform spatial thinning and/or superobbing



2019-06-03; 9.1<lon<9.5; 52.3<lat<52.9; clipped to min. value of $10^{-4}$

- using the **Radar forward operator EMVORADO** in offline mode for **simulating attenuations**, i.e., calculating relevant model equivalents

- important **differences** Radar vs. CML:

    - **Radar**: 17 stations, many azimuths, few elevations, frequency ~5 GHz

    - **CML**: ~4000 "stations"/sender, individual azimuth/elevation (only one per station) and frequency within 10 – 40 GHz



Radarverbund des Deutschen Wetterdienstes



CML Network (Germany)

- two **main inputs** for EMVORADO (many other config. options):
  - ◆ ICON **model fields** (regular grid) for **hydrom.** qr, qg, qv, ...
  - ◆ **auto-generated namelist** with information **for each** CML
    - CML sender is interpreted as Radar station
    - lat/lon/level of "station", azimuth/elev. of ray, frequency, ...
- extract **path-integrated one-way attenuation** from output
- perform EMVORADO run for **each ensemble member**
- current **limitations**:
  - ◆ single EMVORADO run not able to simulate all (~4000) CMLs
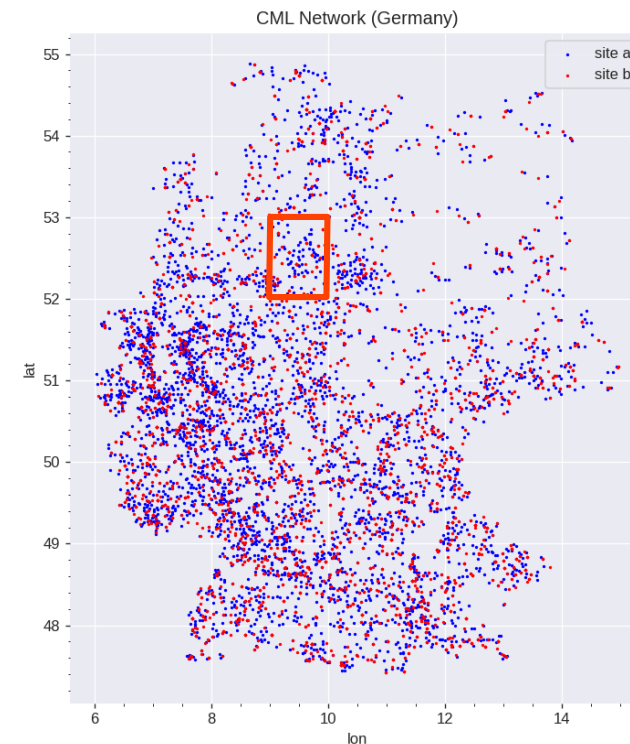  - ◆ simulation does not include water vapor attenuation

- **collect** processed **observed** and **simulated** data for specific **assim. date**

- use **halfway** lat/lon/level of each CML in feedback files

- CML data currently assimilated as SYNOP observation (*obstype*) and using an experimental *codetype* and *varno*

- write all data into feedback (netcdf) file

# CML Case Study I

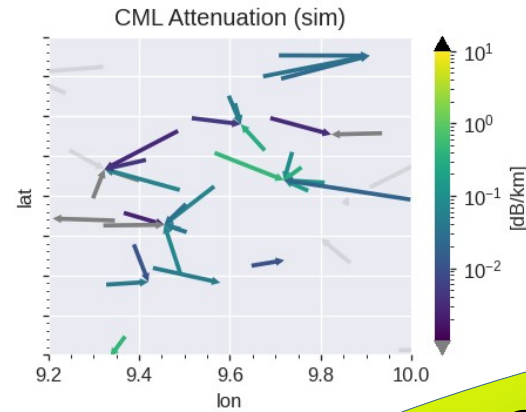- perform assimilation on 2019-06-03 at 12:00:00

- only use CMLs within region
  9.2° < lon < 10° and 52.1° < lat < 52.9°

  ◆ evades EMVORADO limitation

  ◆ 40 CMLs within this region

- **only CML** data is set to **active** here!



CML Network (Germany)

consistency/ plausibility check : ✔

- representation of corresponding "ekf" file (LETKF output)

- shaded background
  → special assimilation state

result: system works (technically)

# CML: Case Study II

- perform **BACY cycle** for **2 days** (on 2019-06-03/04)

- only use **CMLs within region** 8° < lon < 10° and 51° < lat < 52°

- 185 CMLs within this region

- **CML and CONV** data set to active

- vertical localization off



CML Network (Germany)

# CML Case Study: Obs. Err. Stat.

- humidity and temperature stat. (AIREP/TEMP) looked "okay"

- BUT: CML data itself is pulled in wrong direction. Possible causes:

  - localization ($\rightarrow$ correlations)

  - observation error (here: 20%)

- next steps: look at effects of CML data assimilation more closely

  - performing single "core-more runs": single assimilation followed by an ICON model run

  - study LETKF output, increments, and model dynamics (under parameter changes)



error [obs - ana]

RMSE=0.1422; ME=-0.0030; #=8525
RMSE=0.1030; ME=-0.0032; #=8514

cml-off
cml-on

error [obs - fg]

RMSE=0.1250; ME=-0.0036; #=8524
RMSE=0.1385; ME=-0.0058; #=8512

cml-off
cml-on

- two options for observation error

  - "err-rel": relative 20% (not really realistic!)

  - "err-mixed": absolute 0.1 dB/km + relative 20% (more realistic?)

- interesting: large region with missing spread ( →tci?)

**state**
1: active
5: passive
7: rejected

top view

bacy.cml/iodir_coremore.testing/,20190603130000,QV  vs.  bacy.cml/iodir_coremore.testing.new-obs-error/,20190603130000,QV

"obs-rel"  "obs-mixed"  diff ["left" - "right"]

side view

bacy.cml/iodir_coremore.testing/,20190603130000,QV  vs.  bacy.cml/iodir_coremore.testing.new-obs-error/,20190603130000,QV

"obs-rel"  "obs-mixed"  diff ["left" - "right"]

- model field increments for QV from LETKF
- reduced 3D fields to 2D fields via mean along dim. height/y
- clear difference between choices for obs. err.

bacy.cml-ref/iodir_coremore.testing/,obs,20190603130000,20190603130000 vs. bacy.cml-ref/iodir_coremore.testing/,sim,20190603130000,20190603130000

bacy.cml/iodir_coremore.testing/,sim,20190603130000,20190603133000 vs. bacy.cml-ref/iodir_coremore.testing/,sim,20190603130000,20190603133000

- discrepancies between obs./sim REFL at assimilation time

- clear impact of CML data assim. after 30 minutes

- first version for assimilating CML data (integrated into BACY)

- first assimilation results seem plausible

- performed first BACY cycles comparing "CONV" vs "CONV+CML"

- next steps:

  - further study the detailed effects of CML data assimilation (as already begun via single "core-more" exp.)

  - single-obs. experiments (great for studying correlations)

  - study impact of parameters like obs. error, localization, ...

  - general quality control, spatial thinning/superobbing, bias correction

# TCl Basics

- even for **large discrepancies** between obs./sim. REFL LETKF might give **small increments** due to **small ensemble spread**

- targeted covariance inflation (TCI) approach:

  - **check conditions** (missing spread, large enough obs., ...)

  - apply suitable model: each ensemble member gets individual **"virtual" simulated REFL** leading to an increased spread
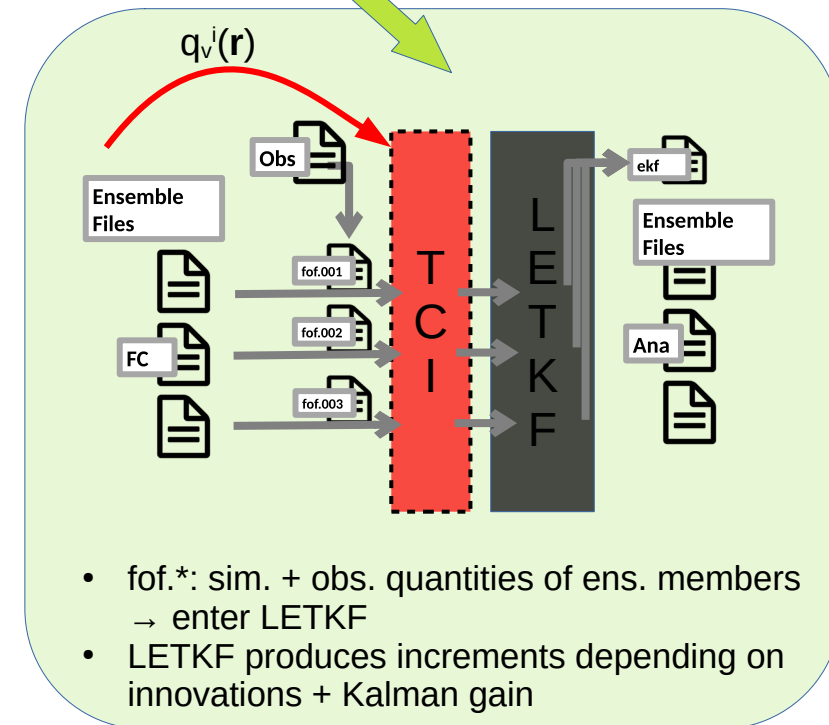
- implemented via **pre-processing of feedback (fof) files** before entering the LETKF

- apply TCI algorithm and **alter simulated Z** in feedback files

- each member processed separately

- use altered feedback files as input for LETKF

$q_v^i(\mathbf{r})$

- fof.*: sim. + obs. quantities of ens. members → enter LETKF
- LETKF produces increments depending on innovations + Kalman gain

- **current model(s)**: based on simple linear regression
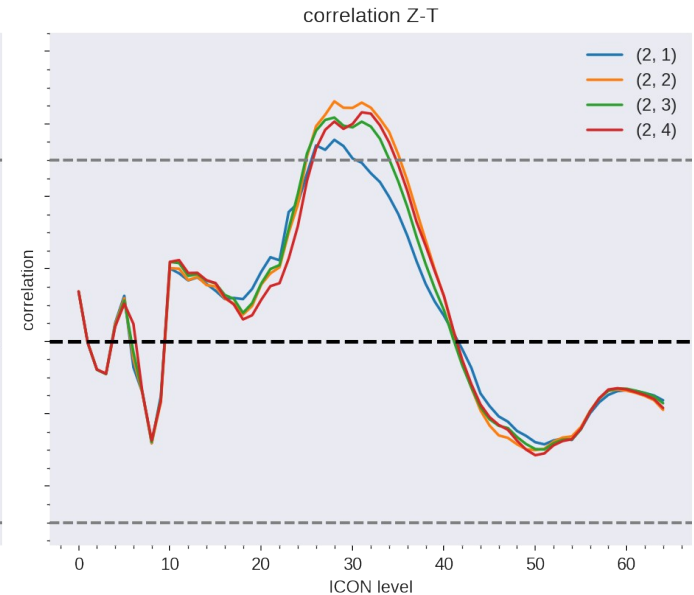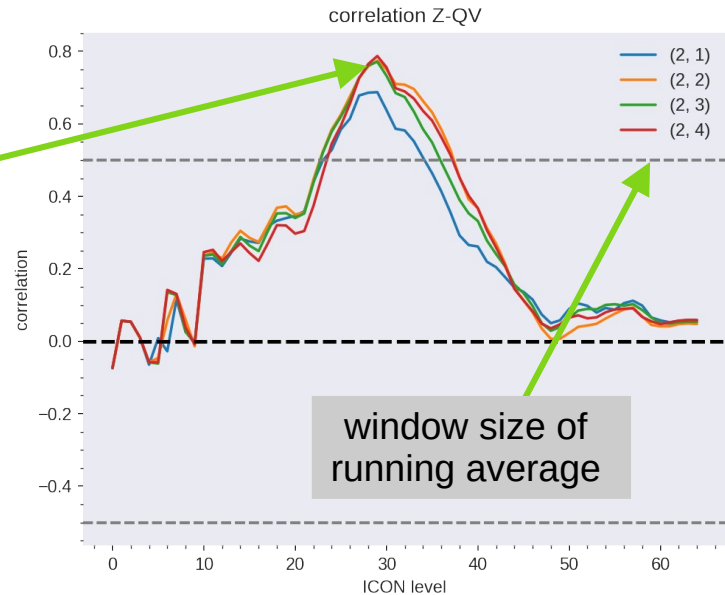  - $M_{h,h'}$ : $\delta Z_i(x,y,h,t) = \alpha * \delta qv_i(x,y,h',t)$
  - $\delta Z_i$, $\delta qv_i$: ensemble perturbations for Z, qv of i-th member
  - h, h': categorical/discrete heights
- **overall idea**:
  - spread of qv "imprinted" onto spread of Z
  - assim. "favors" members with more humidity
  - additional increments for humidity qv are produced
  - model (hopefully) generates qr/qs/qg → EMVORADO sim. REFL

adate:20190603140000, leadtime:20, ensemble_slice:slice(1, None, None)

- idea: training data should be **representative for convective events**

- built simple algorithm for the **detection of new cells**

  - employs **time series** of (binned) Radar data

  - gives **area and maximum position** of REFL $(x_0,y_0)$ of newly emerged cell at time $t_0$

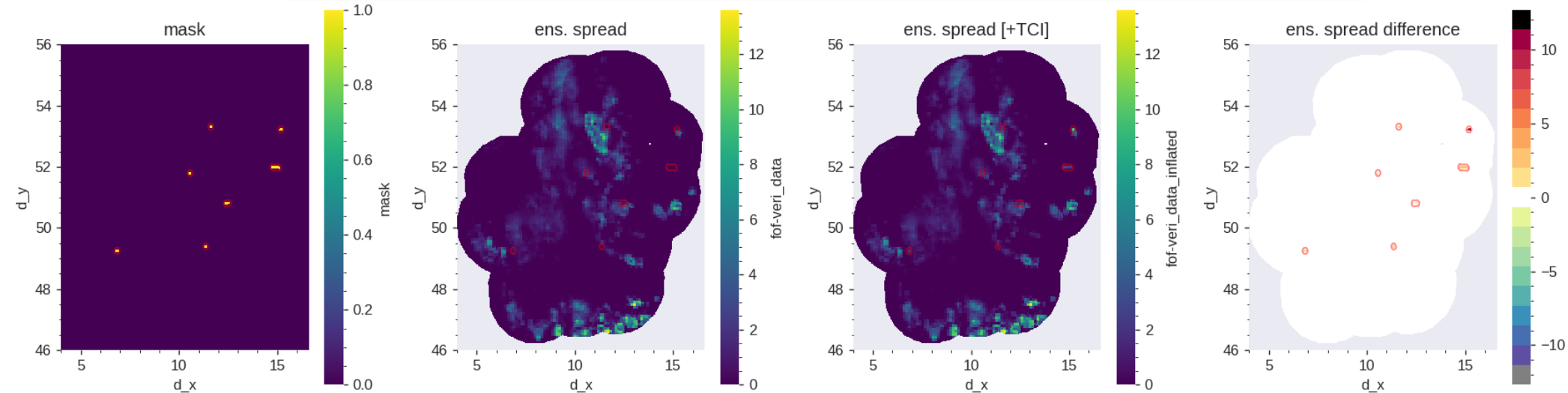- single instance (for training of model $M_{h,h'}$):  $\delta Z_i(x_0,y_0,h,t_0)$, $\delta qv_i(x_0,y_0,h',t_0)$

correlation Z-QV

correlation Z-T

window size of running average

- fixed height of REFLs of h="3000m-4000m"

- h' determined through **maximum of correlation** → h'=30

- resulting model:

  - $\delta Z_i(h="3000-4000m") = 10^4\ dBZ * \delta qv_i(h'=30)$

- currently working on "new" TCI based on **machine learning**

  - goal: ultra-short **prediction of newly emerging REFL** and its magnitude ("rough" estimate)

  - learn ICON model dynamics for convection

  - not living within ensemble pert. space!

- **predictors**: qv, T at several heights (+spatial mean/std)

- **target: temporal derivative of REFL** ΔZ (initially vertically integrated qr must be zero →no rain!)

- employed ML algorithms: KNN, Decision Tree

- much **more flexible approach** (→ apply to CML data?)

# TCI Case Study

- set up two bacy experiment with/without application of TCI

- period: 2019-06-03 → 2019-06-10

- TCI is applied hourly at every assimilation step

  - TCI based on **simple linear model** (as shown previously)

  - TCI **applied to ALL radar data** over complete model domain

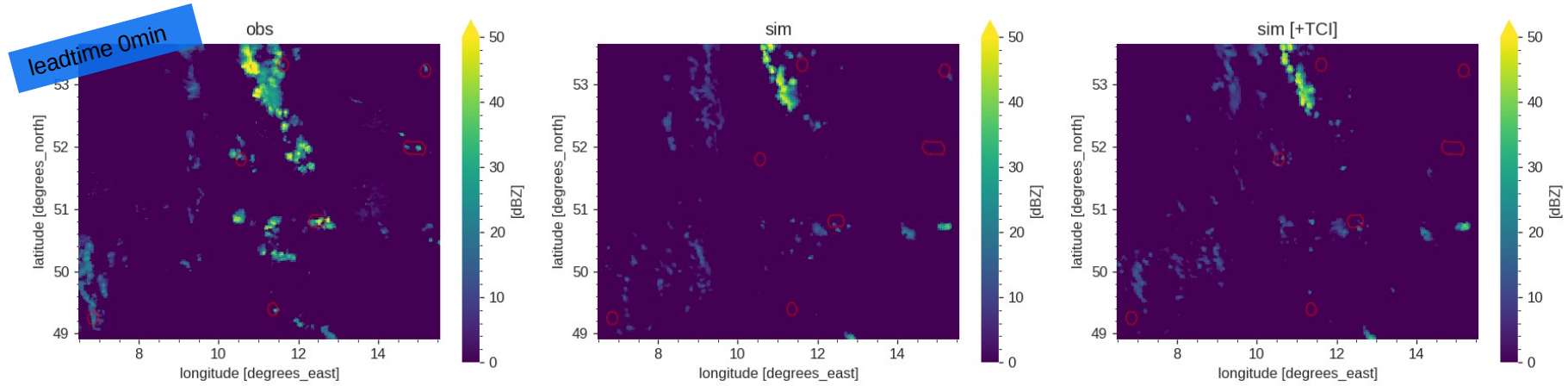- Initiate main forecast runs every 6h (max. leadtime 6h)

TCI monitoring for assimilation at 2019-06-05 15UTC; mean over stations/elevations
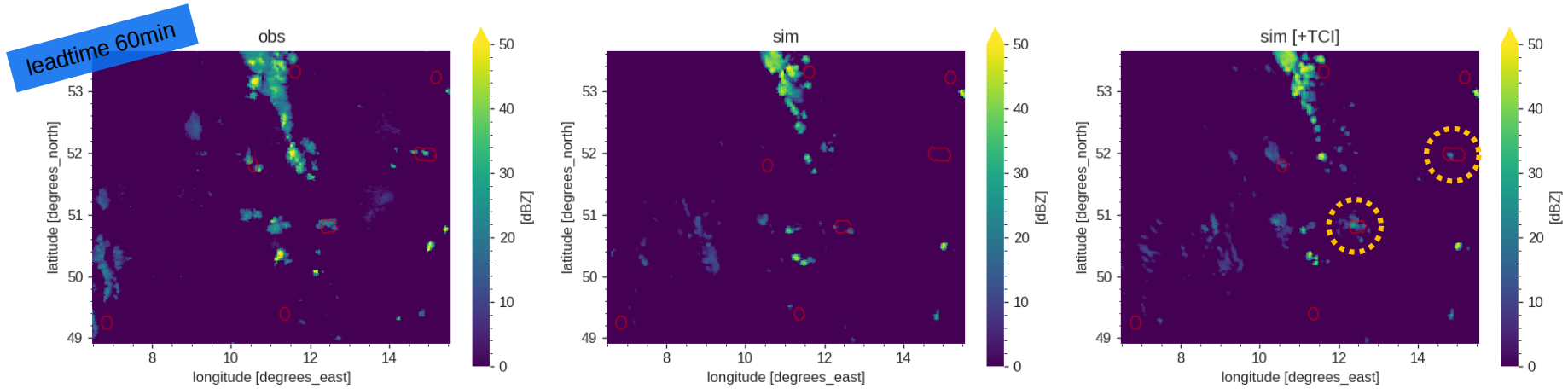
- "mask" shows if TCI got applied (also indicated with red contours)

- **main conditions** for specific obs.: vanishing ensemble spread/mean/det (+running average), sizeable observed REFL (Z>20), REFL height between 3000m and 4000m (all elevations)

dbzcmp @1.5; leadtime=0min

dbzcmp @1.5; leadtime=60min

- **reduced negative impact on humidity** stat. (for AIREP/TEMP) w.r.t. previous TCI implementations

- T/RH/WIND/REFL stat. for AIREP/RADAR unobtrusive

TIME SERIES PLOT for
period: 20190603 to 20190609
0S,6H 0M 0S,12H 0M 0S,18H 0M 0S UTC + (0S to 6H 0M 0S)
Forecast IDs: TCI1, TCI0
members: 0
method: fss
score: fss
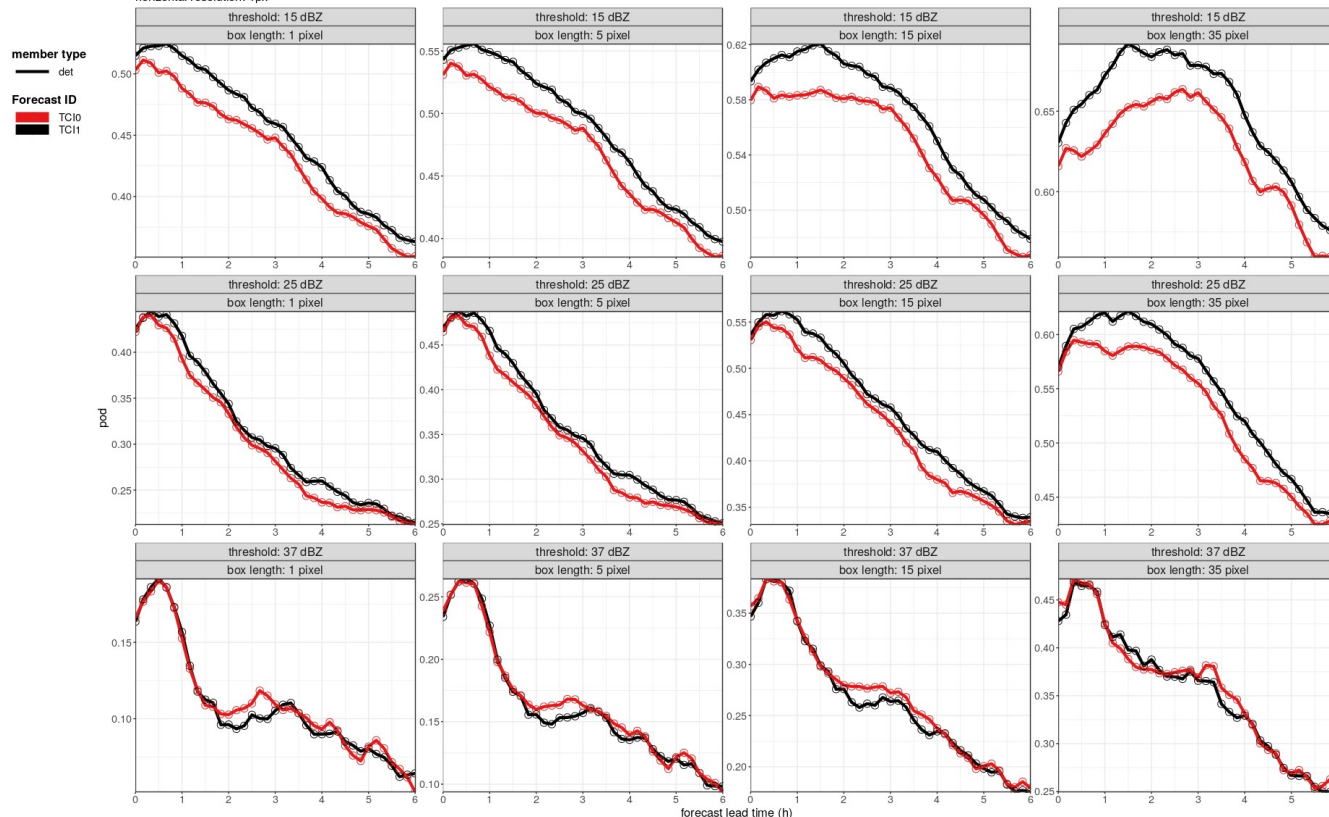#valid files: 23532/24864
horizontal resolution: 1px =

- Fractional Skill Scores (FSS) for dbzcmp from main forecast runs

- clear **positive impact** even after longer leadtimes!

- Probability of Detection (POD) for dbzcmp from main forecast runs
- clear **positive impact** even after longer leadtimes!

- overall, TCI results are promising

  - production of **"new" REFL cells** (consistent with observations)

  - positive impact on **fractional skill score** (w.r.t. dbzcmp)

  - **obs. err. stat.** results are unobtrusive (i.e. not too negative)

- next:

  - further studies necessary (verification of **longer time periods**)

  - continue work on **ML-based TCI**

Thank you for your attention!