

DINCAE: A neural network to reconstruct missing data in satellite images

Alexander Barth¹, Aida Alvera-Azcárate¹, Matjaz Licer², and Jean-Marie Beckers¹

Paper: <https://doi.org/10.5194/gmd-13-1609-2020>

Video: <https://youtu.be/MJaEncQv0eE>

¹GeoHydrodynamics and Environment Research (GHER), University of Liège, Liège, Belgium

²National Institute of Biology, Marine Biology Station, Piran, Slovenia
GHER, University of Liège, Belgium



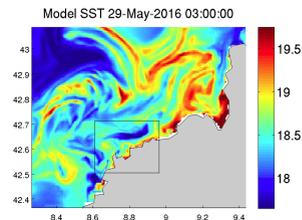
in situ
measurements



satellite
measurements



ocean models



Remote sensed data

- generally **more accurate than models**
- good **spatial coverage**

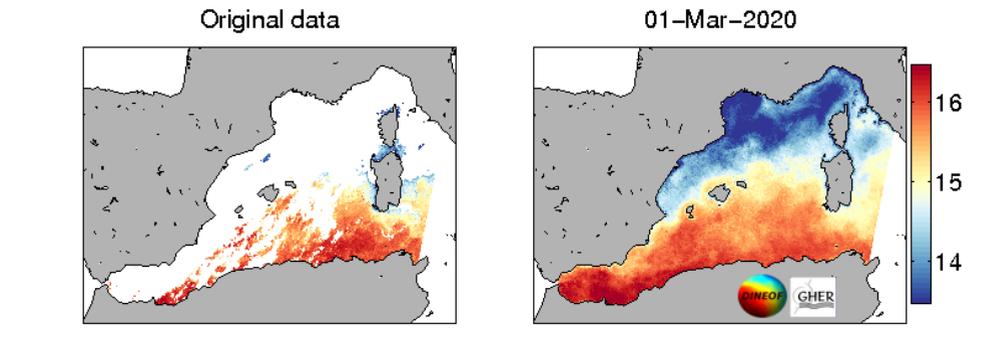
But:

- **gaps** due to e.g. clouds
- just the **surface**
- just **some parameters** (temperature, sea surface height, salinity, ...)
- no prediction (obviously)

Ideally:

- train a neural network on a large collection of full images
- reconstruct the missing data using the trained network

However: **only very few images do not have any clouds**

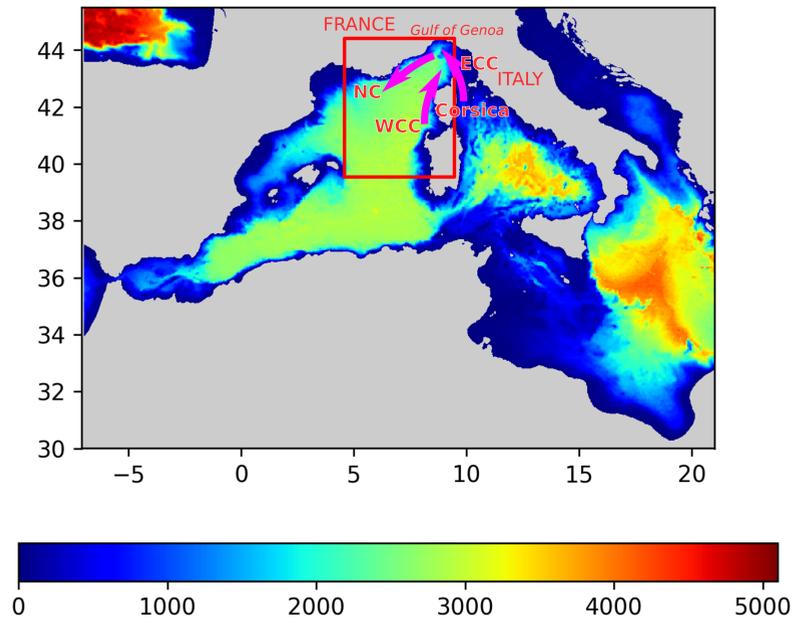


Application to SST

- Having **long time series** to train neural networks is quite important
- Advanced Very High Resolution Radiometer (AVHRR) dataset, from 1 April 1985 to 31 December 2009.
- One single image is composed by 112 x 112 grid points.
- **Cross-validation**: in the last 50 images we removed data according to the cloud mask of the first 50 images of the SST time series.
 - not used at all during either the training or the reconstruction phases
 - can be considered independent.
 - 106816 measurements have been withheld this way.

Study area

- the red rectangle: delimits the studied region
- color represents the bathymetry in meters
- the main currents: the Western Corsican Current (WCC), the Eastern Corsican Current (ECC) and the Northern Corsican Current (NC)



Bayes rule

- For Gaussian-distributed errors:
 - prior: $N(x^f, \sigma^f)$
 - observations: $N(y^o, \sigma^o)$
 - posterior: $N(x^a, \sigma^a)$
- Bayes:

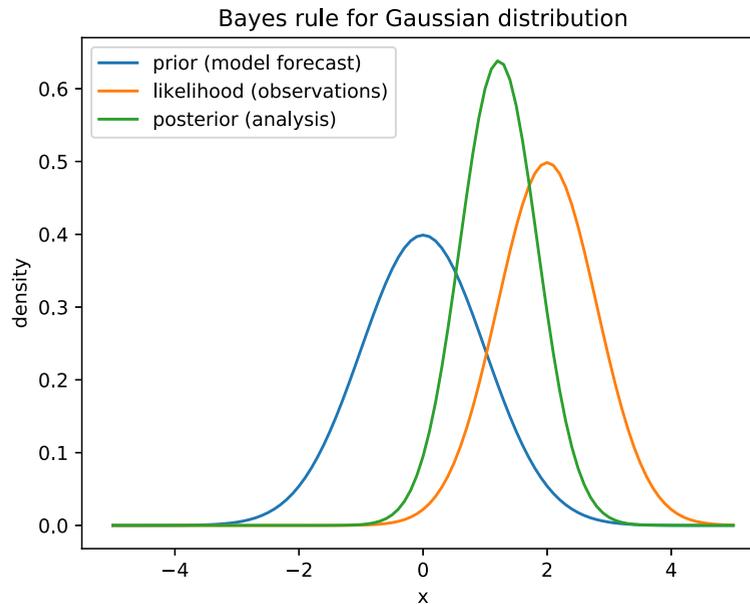
$$p(x|y^o) = \frac{p(x)p(y^o|x)}{p(y^o)}$$

- Mean and variance of posterior given by:

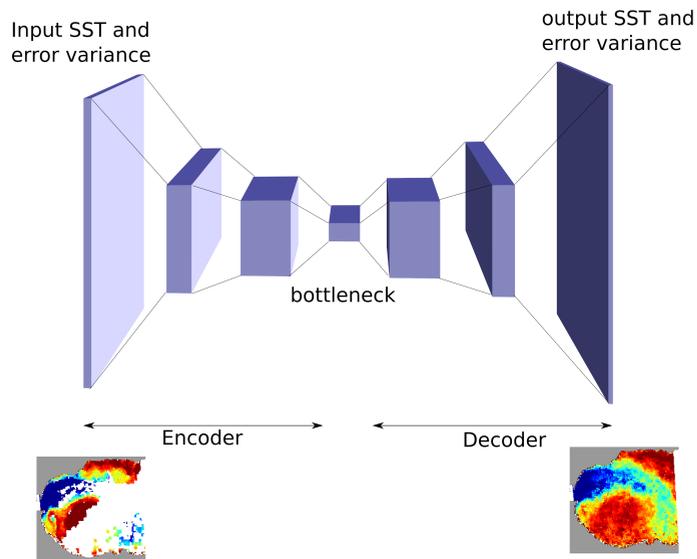
$$\sigma^{a-2} x^a = \sigma^{f-2} x^f + \sigma^{o-2} y^o$$

$$\sigma^{a-2} = \sigma^{f-2} + \sigma^{o-2}$$

- Concept of **information**: $\sigma^{f-2} x^f$ and $\sigma^{o-2} y^o$



Structure



Auto-Encoder: used to efficiently compress/decompress data, by extracting main patterns of variability

- Similarity to EOFs

Convolutional: works on subsets of data, i.e. trains on local features

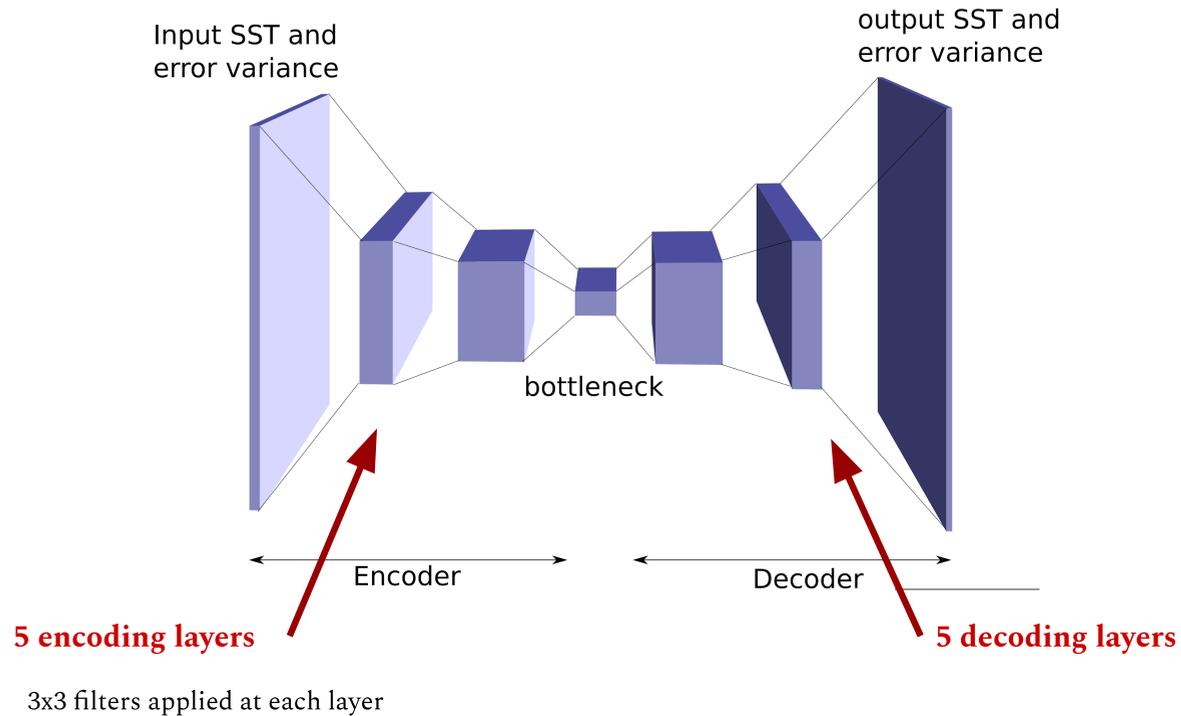
Missing data handled as data with different initial errors

- If missing, error variance (σ^2) tends to ∞

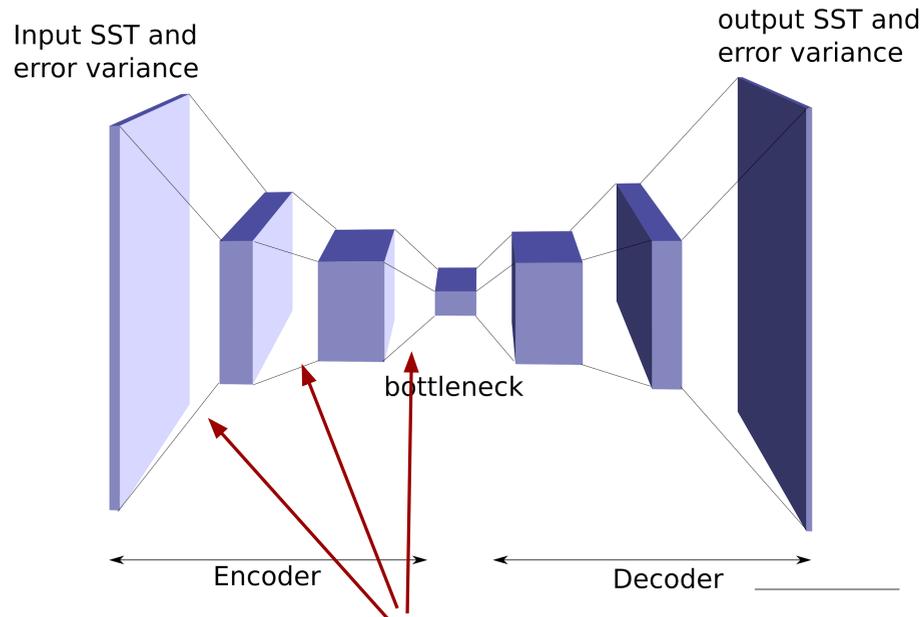
Input data:

- SST/ σ^2 (previous day, current day, following day)
- $1/\sigma^2$ (previous day, current day, following day)
- Longitude
- Latitude
- Time (cosine and sine of the year-day/365.25)

Structure



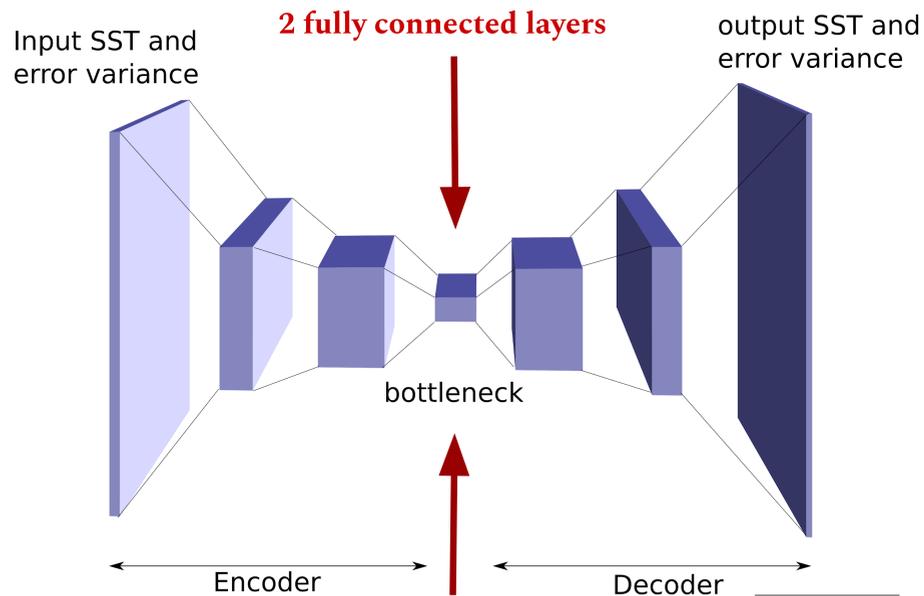
Structure



Average pooling layers

Reduce size by retaining the average value on 2x2 boxes

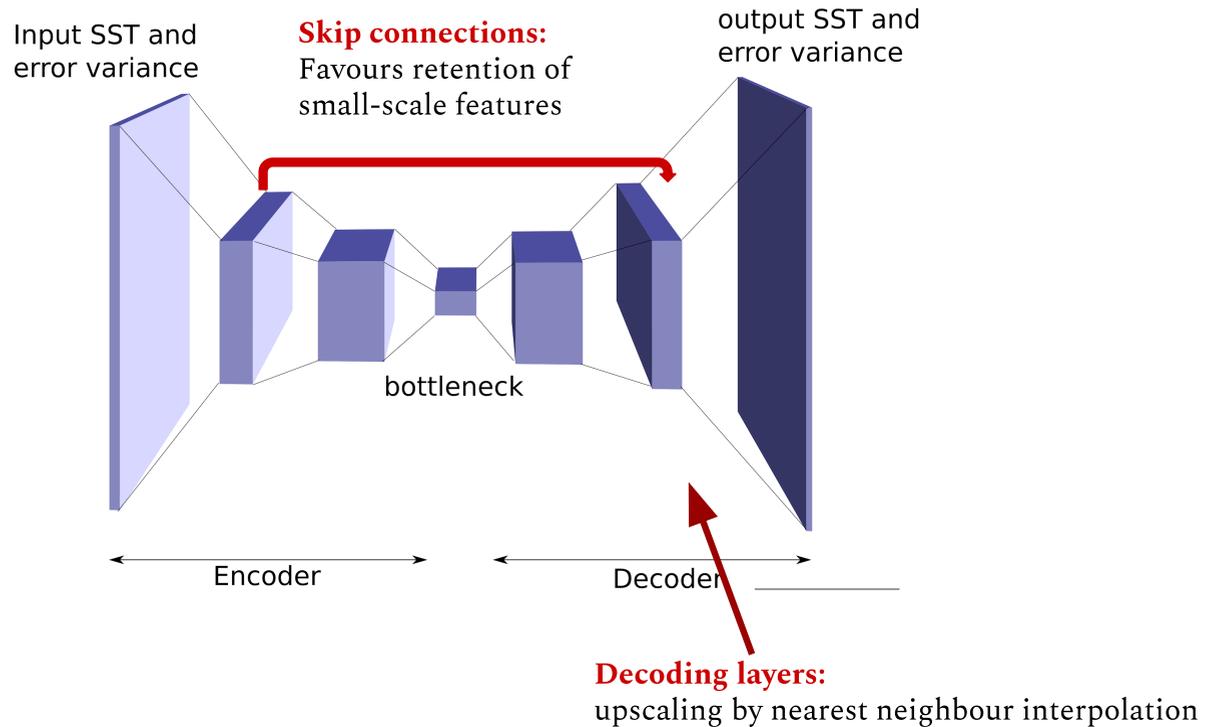
Structure



+ 2 drop-out layers

Take out 30% of neurons (pixels) to avoid overfitting

Structure



Training

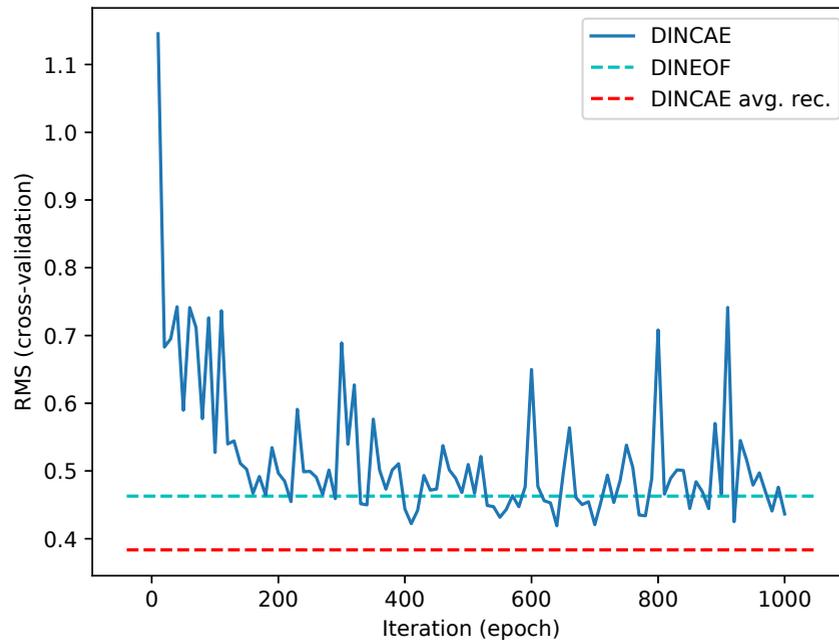
- Partitioned into so-called mini-batches of 50 images
- The entire dataset are used multiple times (epoch)
- For every input image, more data points were masked (in addition to the cross-validation) by using a randomly chosen cloud mask during training (data set **augmentation**).
- The output of the neural network (for every single grid point i, j) is a Gaussian probability distribution function characterized by a mean \hat{y}_{ij} and a standard deviation $\hat{\sigma}_{ij}$.

$$J(\hat{y}_{ij}, \hat{\sigma}_{ij}) = \frac{1}{2N} \sum_{ij} \left[\left(\frac{y_{ij} - \hat{y}_{ij}}{\hat{\sigma}_{ij}} \right)^2 + \log(\hat{\sigma}_{ij}^2) + 2 \log(\sqrt{2\pi}) \right]$$

- The first term: mean square error, but scaled by the estimated error standard deviation.
- The second term: penalizes any over-estimation of the error standard deviation.

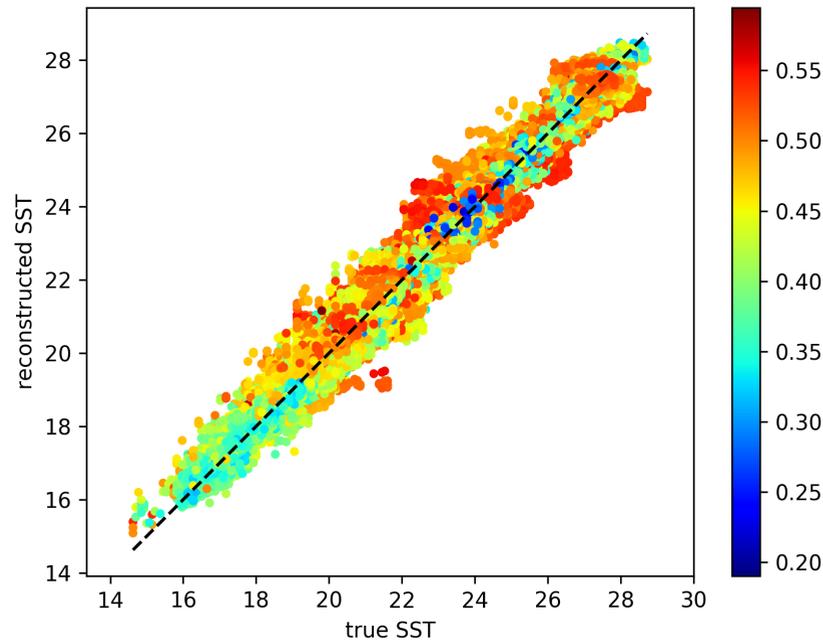
Results

- RMS difference with cross-validation dataset as a function of iteration.
- The solid blue line represents the DINCAE reconstruction at different steps of the iterative minimization algorithm.
- The dashed cyan line is the DINEOF reconstruction and the dashed red line is the **average DINCAE reconstruction** between epochs 200 and 1000.



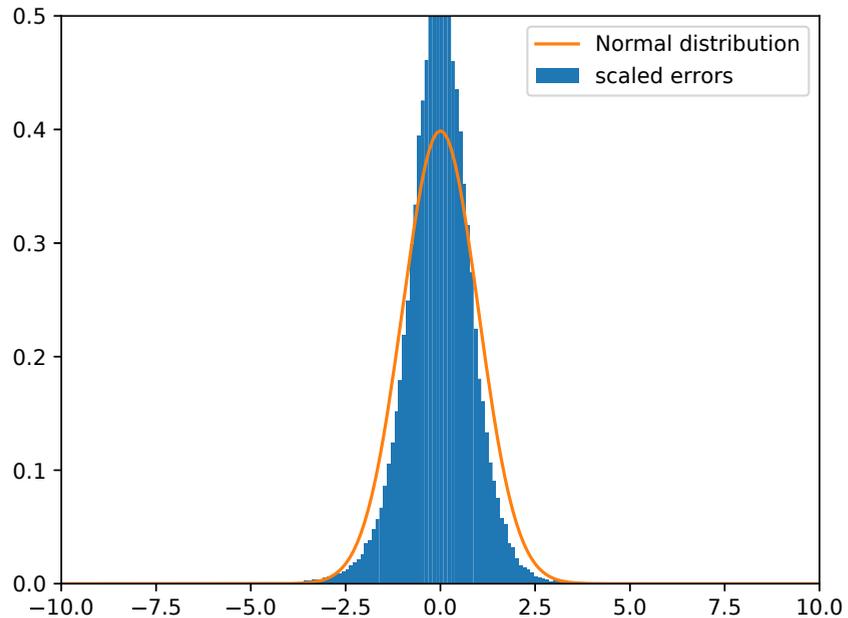
Error estimation

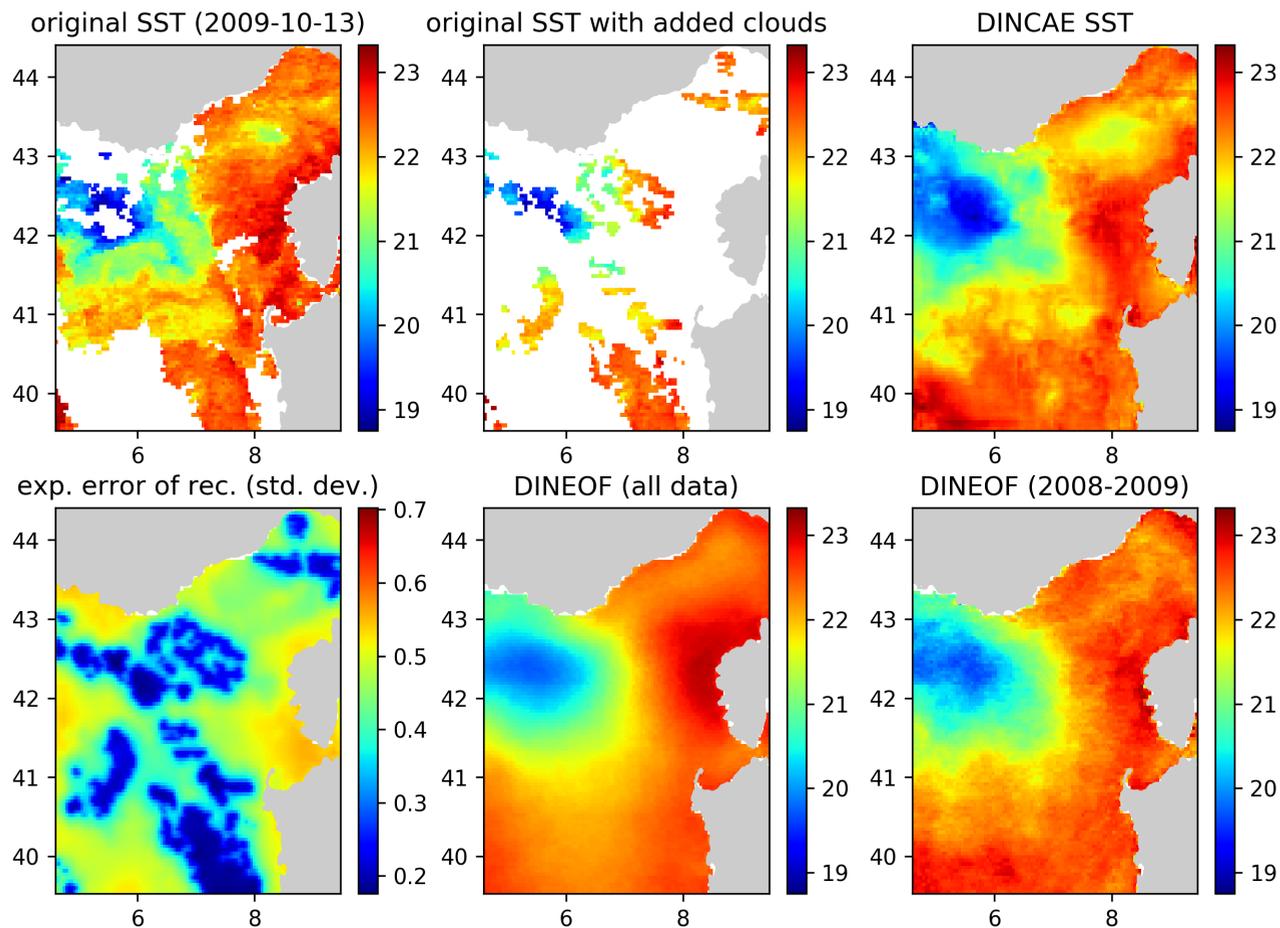
- scatter plot of the true SST (withheld during cross-validation) and the corresponding reconstructed SST.
- The color represents the estimated expected error standard deviation of the reconstruction.
- Low error values are expected to be closer to the dashed line.
- Reconstructed and cross-validation SST tend to cluster relatively well around the ideal dashed line.
- Typically the lower expected errors are found more often near the dashed line than at the edge of the cluster of points.



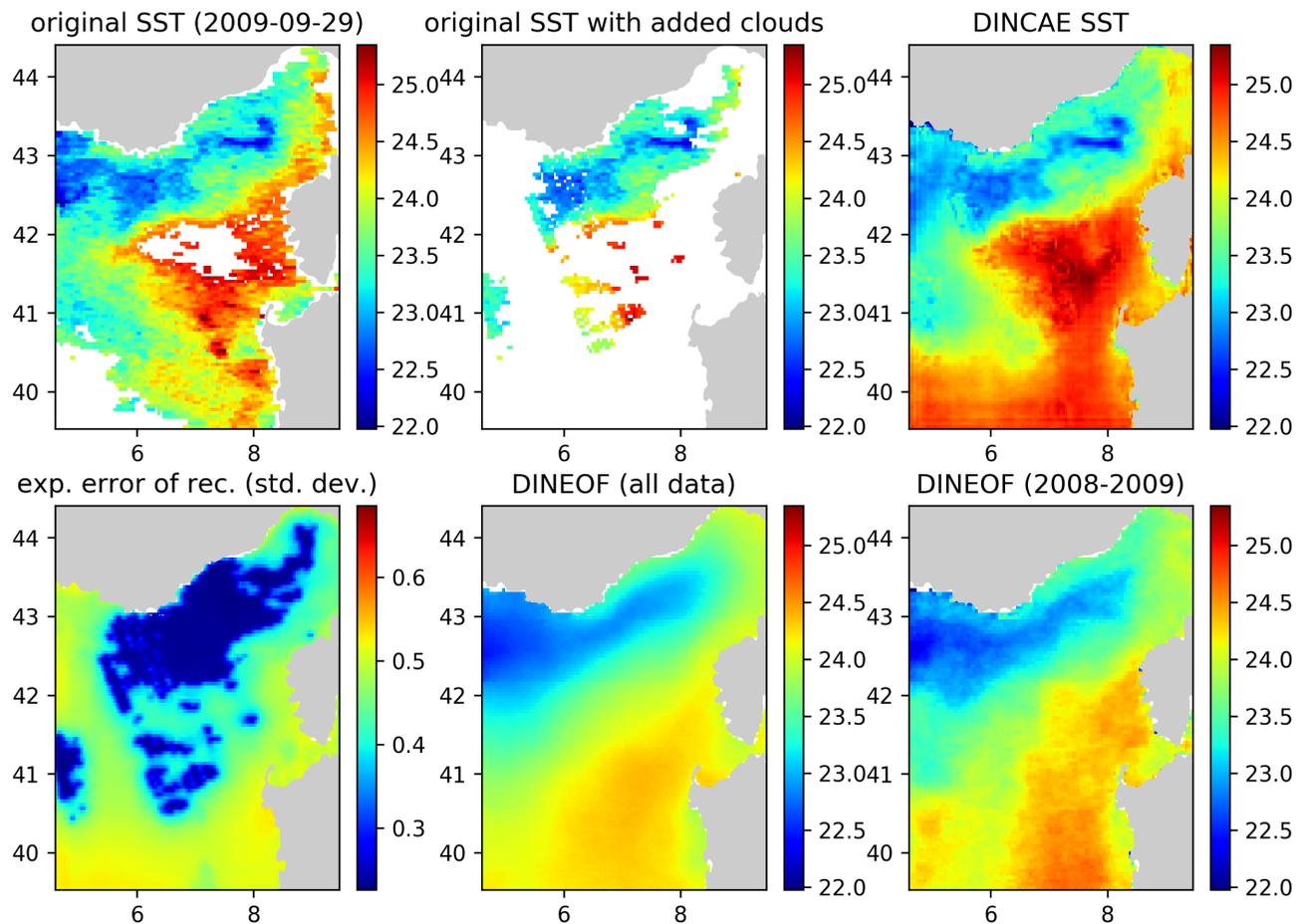
Error estimation

- Scaled errors are computed as the difference between the reconstructed SST and the actual measured SST (withheld during cross-validation) divided by the expected standard deviation error.





Example reconstruction with DINCAE and DINEOF for the date 2009-10-13



Example reconstruction with some artefacts for the date 2009-09-29

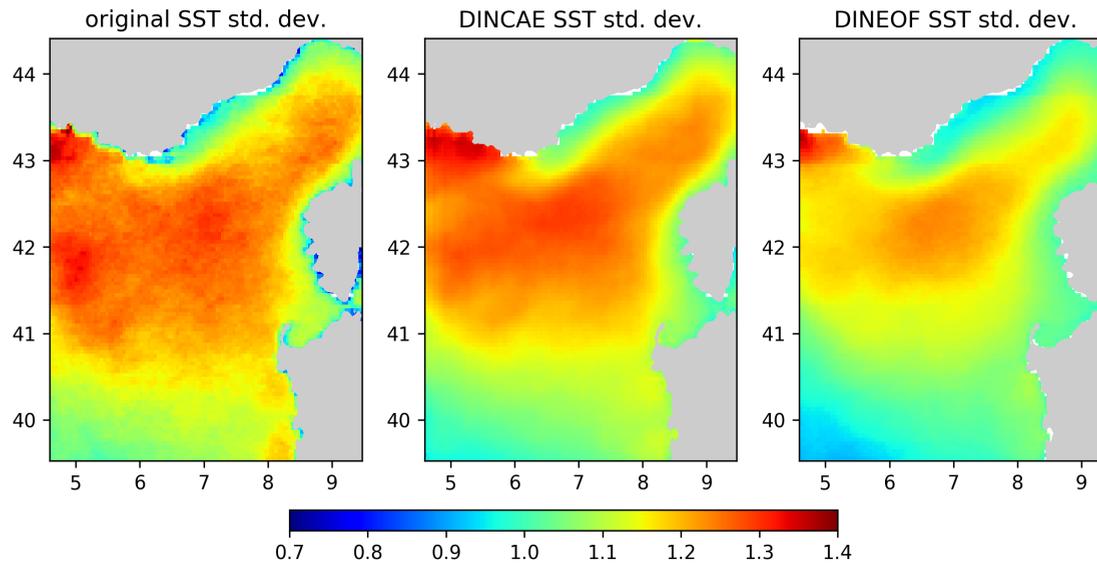
Independent validation

Comparison with the independent cross-validation data and the dependent data used for training (in °C)

	RMS	CRMS	bias
DINEOF	1.1676	1.1102	-0.3616
DINCAE	1.1362	1.0879	-0.3278

Comparison with the World Ocean Database for SST grid points covered by clouds

Variability



Standard deviation computed around the seasonal average

Conclusions

- **Practical way to handle missing data** in satellite images for neural networks
- Measured data **divided by its expected error variance**. Missing data are thus treated as data with an infinitely large error variance.
- The cost function of the neural network is chosen such that the network provides the reconstruction but also the **confidence of the reconstruction error**
- Reconstruction method **DINCAE compared favourably** to the widely used DINEOF reconstruction method which is based on a truncated EOF analysis
- The expected error for the reconstruction reflects well the **areas covered by the satellite measurements** as well as the areas with more intrinsic variability (like meanders of the Northern Current). The expected error predicted by the neural network provides a good indication of the accuracy of the reconstruction.
- The variability of the DINCAE reconstruction **matched the variability of the original data** relatively well.
- Code: <https://github.com/gher-ulg/DINCAE>
- Paper: Barth, A., Alvera-Azcárate, A., Licer, M., and Beckers, J.-M.: [DINCAE 1.0: a convolutional neural network with error estimates to reconstruct sea surface temperature satellite observations](#), Geosci. Model Dev., 2020.