

**XIII giornata della modellistica in aria(net)
Milano, 25/03/2026**



Machine Learning for Air Quality and Environmental Modelling

Overview of **ARIANET solutions**

Umberto Giuriato – Alessandro D'Ausilio

Is Machine Learning the Death-Of-Modelling?



ML refers to statistical techniques to automatically detect patterns in data by optimizing a performance metric (Murphy, 2022).

Unsupervised learning

Finding structure in unlabeled data

Meteorological-day clustering for faster long-term dispersion simulations

Road-segment classification for urban dispersion kernels and digital twins

Supervised learning

Learning from labeled data

Inverse modelling to rescale emission factors from observations

Mass-consistent concentration downscaling driven by geophysical proxies

Data fusion and concentration time-series forecasting

At *ARIANET*, we use machine learning **not to replace modelling**, but to make **traditional models more efficient** and to **integrate external data sources** into our modelling framework.



Unsupervised learning

Finding structure in unlabeled data

Meteorological-day clustering for fast long-term dispersion simulations

Road-segment classification for urban dispersion kernels and digital twins

Supervised learning

Learning from labeled data

Inverse modelling to rescale emission factors from observations

Mass-consistent concentration downscaling driven by geophysical proxies

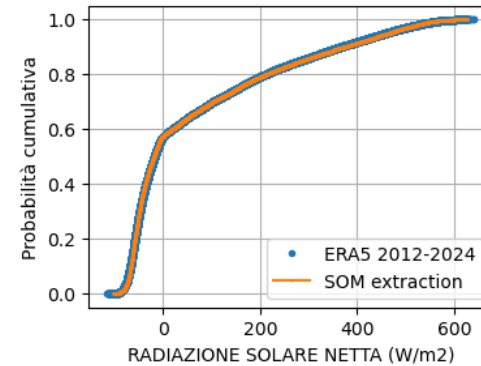
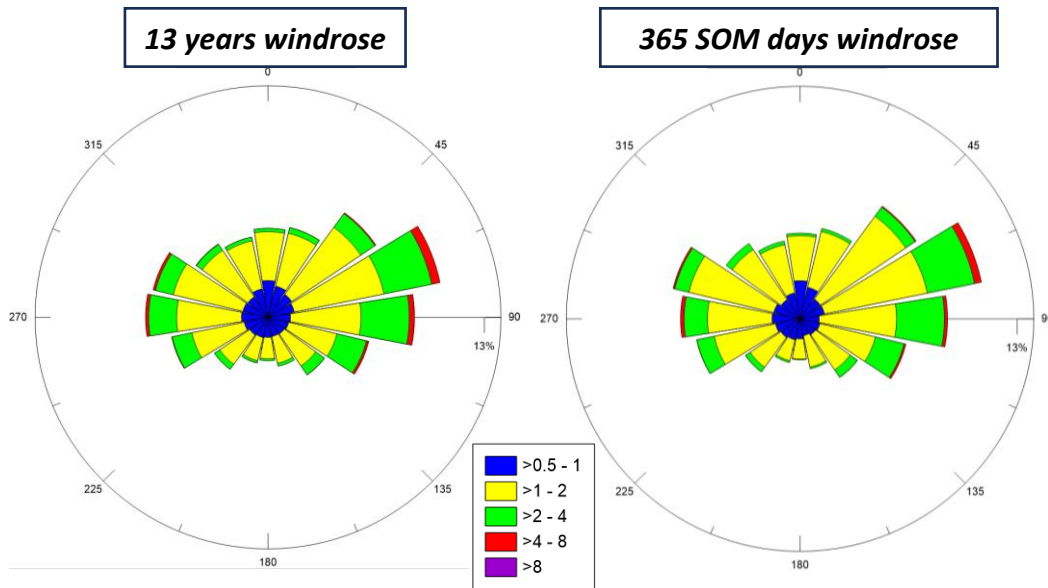
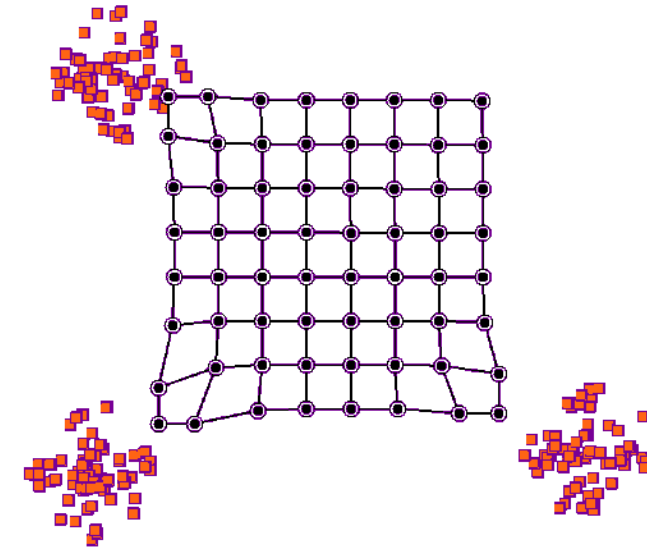
Data fusion and concentration time-series forecasting

Meteorological days clustering for fast long-term dispersion simulations

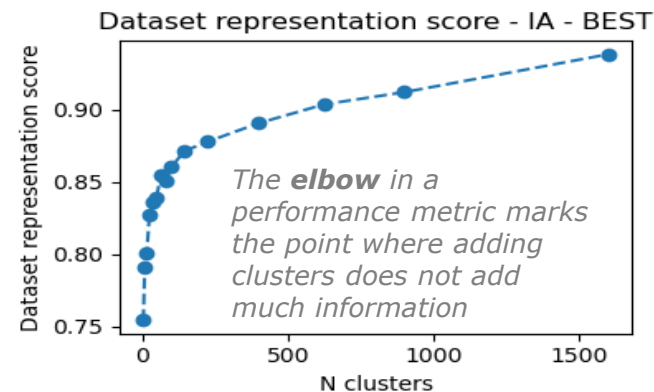
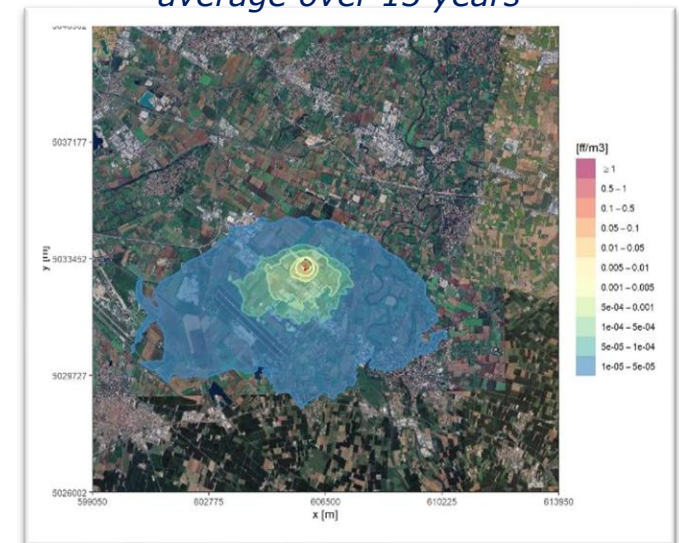
Need: Air quality assessments on long periods of time (e.g. in VIS for industrial plants or landfills). Need to fasten the generation of long-term statistics of concentration fields

Solution: The *Self Organizing Maps (SOMs)* algorithm automatically separates the days into groups and find a representant for each group, based on similarity of hourly meteorological variables (*wind speed, wind direction, cloud cover, net radiation*)

Workflow: Run dispersion (*SPRAY*) only for representative days, then reconstruct long-term statistics from their frequency



Output concentration average over 13 years



Some Studies delivered with SOMs #Ecologja
VIS Ecoeternit Landfill - VIS Tecnoinerti Landfill
VIS Acqua&Sole plant - AQ Assessment Cargill plant

Unsupervised learning

Finding structure in unlabeled data

Meteorological-day clustering for faster long-term dispersion simulations

Road-segment classification for urban dispersion kernels and digital twins

Supervised learning

Learning from labeled data

Inverse modelling to rescale emission factors from observations

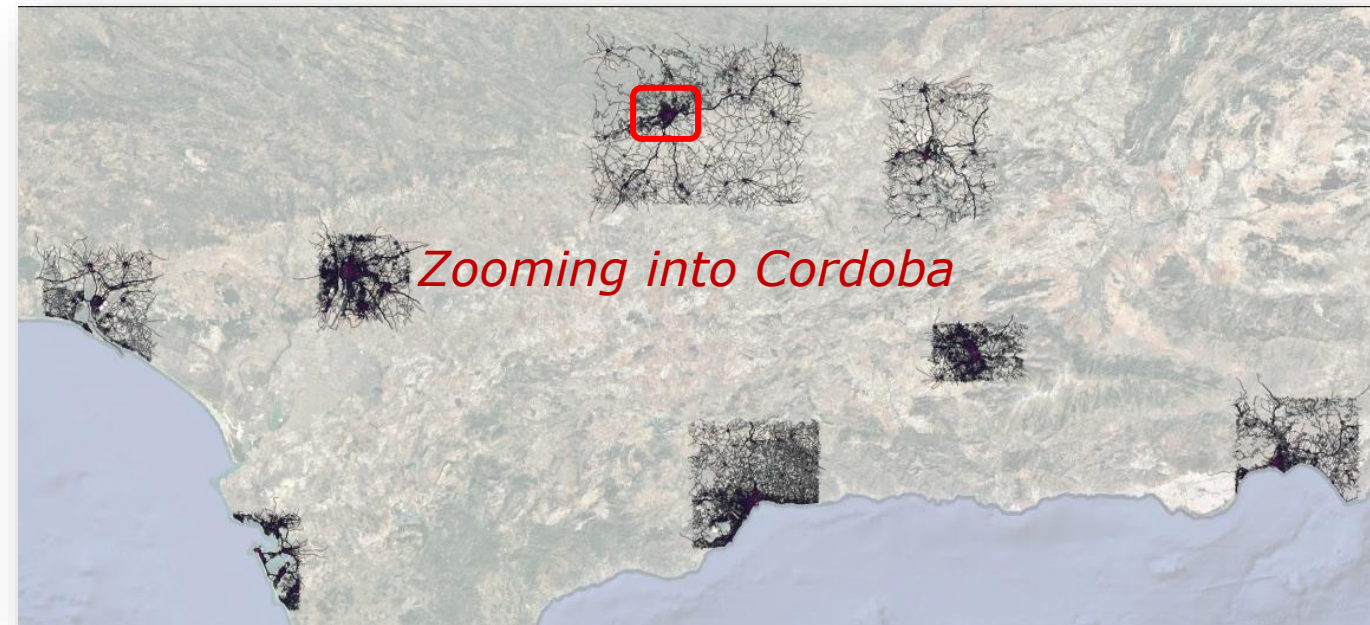
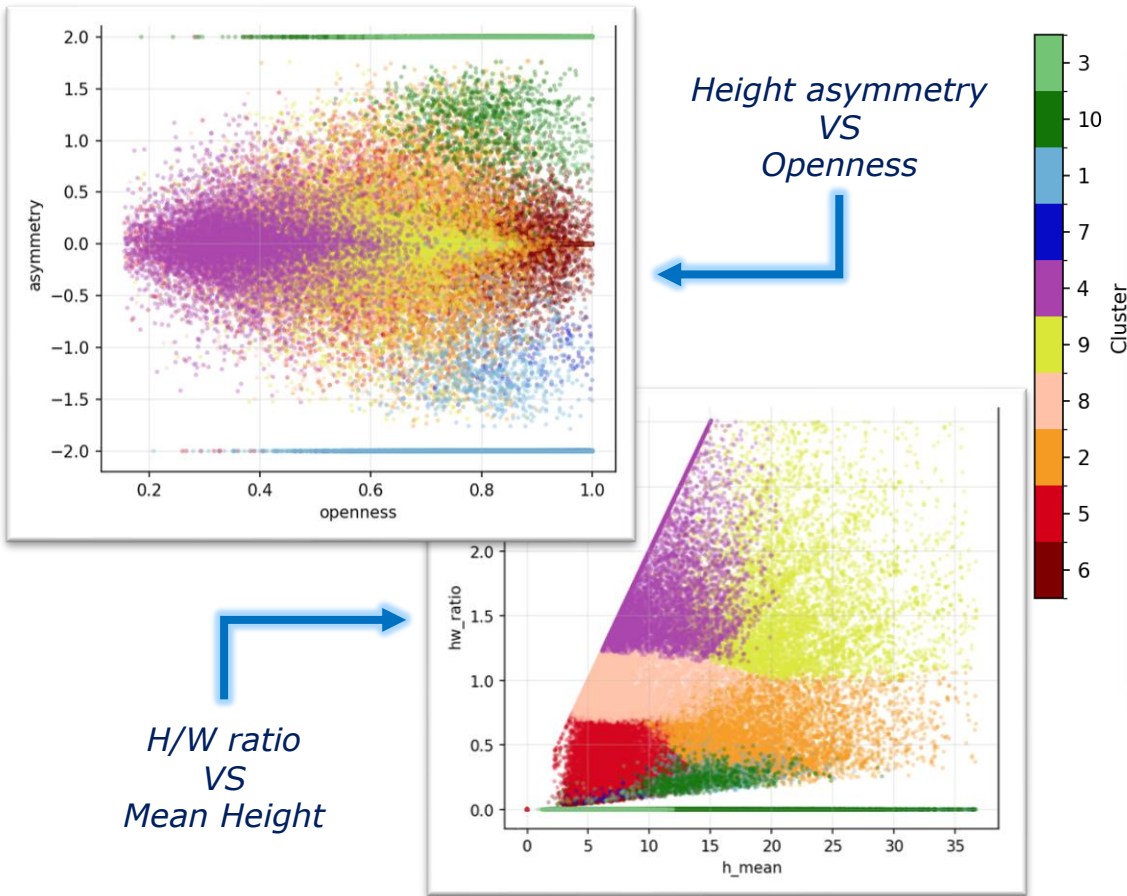
Mass-consistent concentration downscaling driven by geophysical proxies

Data fusion and concentration time-series forecasting

Unsupervised learning for roads clustering in an urban environment

Goal: Building representative dispersion kernels in urban environment to be employed in air quality systems

Task: We leverage clustering algorithms (eg *K-Means*, *Gaussian mixtures*, *SOMs*, *HDBSCAN*) to classify road segments in the features space that characterize urban environment: *H/W ratio*, *Height asymmetry*, *Openness* (1 - % built soil), ...



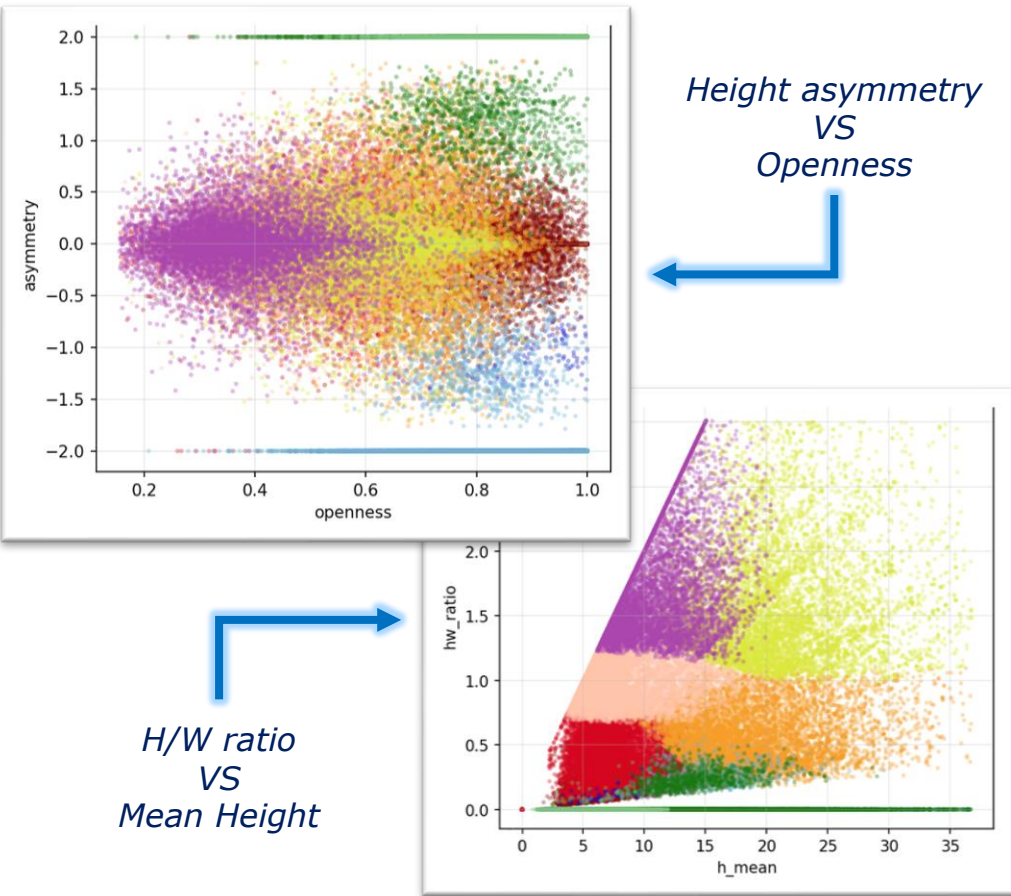
Some projects involved
 Andalusia Air Quality Forecasting system
 CALLIOPE - Digital Twin of the city of Taranto

Left asymmetric canyon Tall buildings	Shallow canyon Very low builtup	Shallow canyon Low buildings	Medium-shallow canyon Medium builtup
Left asymmetric canyon Short buildings	Medium canyon Short buildings	Deep canyon Low builtup	Deep canyon High builtup
	Right asymmetric canyon Short buildings	Right asymmetric canyon Short buildings	Right asymmetric canyon Tall buildings

Unsupervised learning for roads clustering in an urban environment

Goal: Building representative dispersion kernels in urban environment to be employed in air quality systems

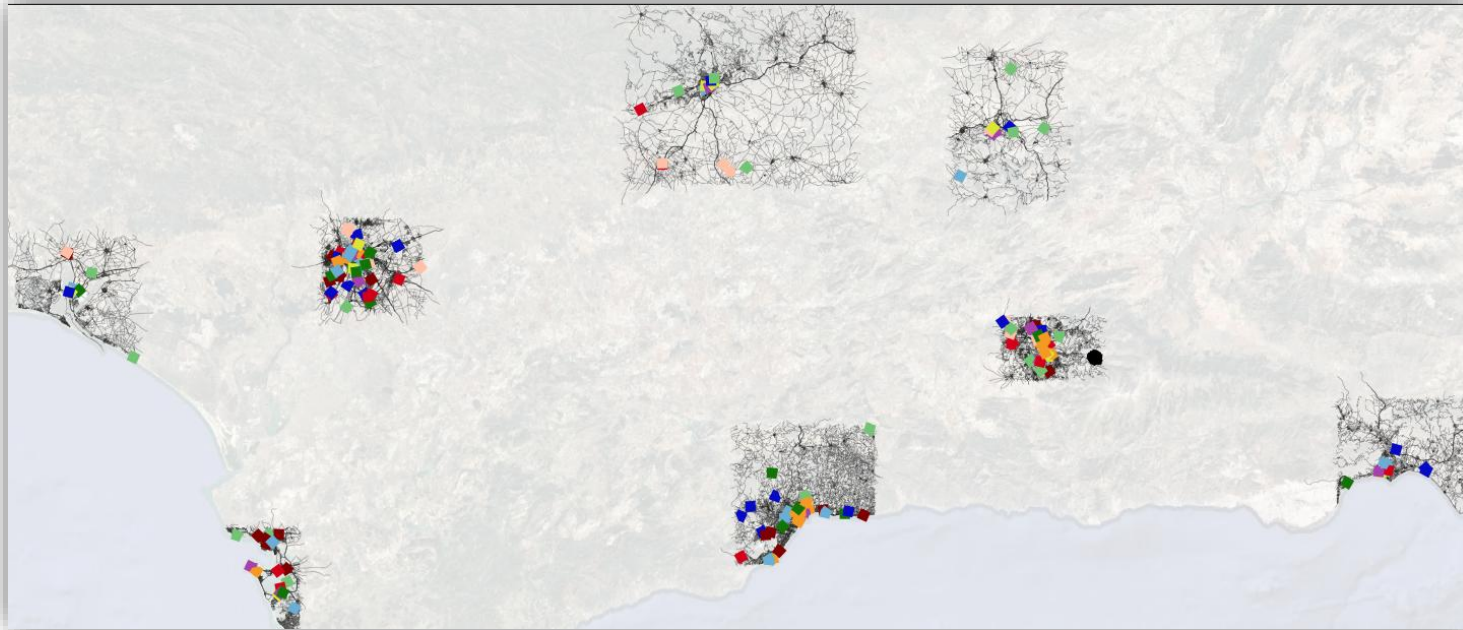
Task: We leverage clustering algorithms (eg *K-Means*, *Gaussian mixtures*, *SOMs*, *HDBSCAN*) to classify road segments in the features space that characterize urban environment: *H/W ratio*, *Height asymmetry*, *Openness* (1 - % built soil), ...



3	Shallow canyon Very low builtup	Shallow canyon Low buildings	Medium-shallow canyon Medium builtup
10	Medium canyon Short buildings	Deep canyon Low builtup	Deep canyon High builtup
1	Left asymmetric canyon Tall buildings	Right asymmetric canyon Short buildings	Right asymmetric canyon Tall buildings
7	Left asymmetric canyon Short buildings		
4			
9			
8			
2			
5			
6			

Some projects involved
Andalucia Air Quality Forecasting system
CALLIOPE – Digital Twin of the city of Taranto

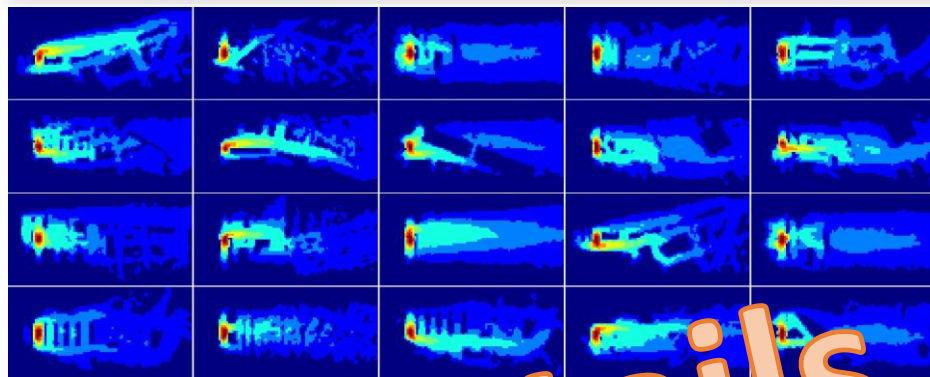
Unsupervised learning for roads clustering in an urban environment



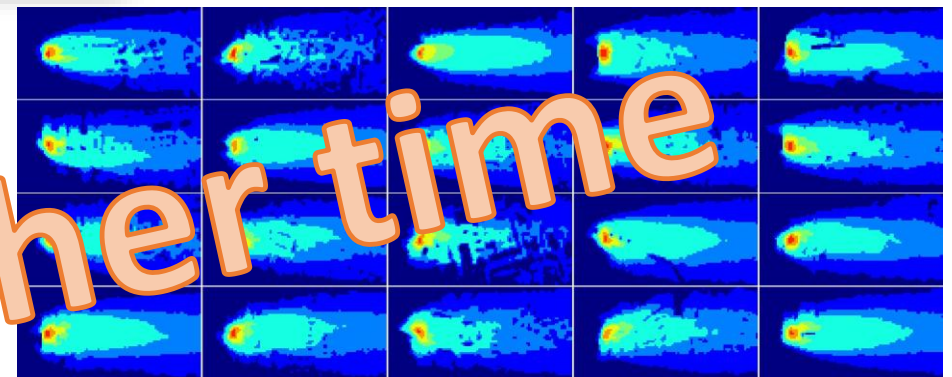
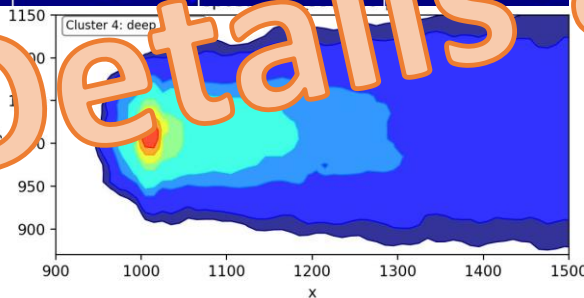
In each cluster we identify an *ensemble* of representative roads

For each ensemble, we build a set of dispersion kernels averaging *stationary PMSS concentrations*, in classified combinations of *Windspeed, Stability, Wind-road orientation*

Example of kernels building for two road types

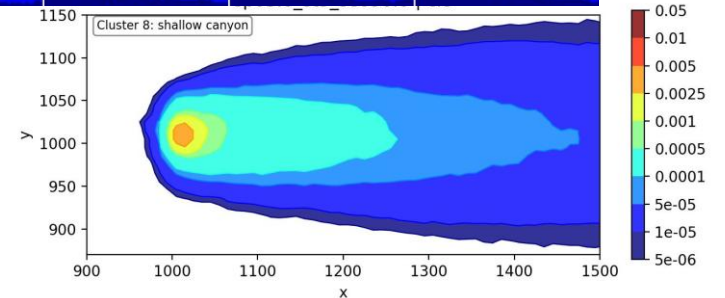


Deep canyon
Windspeed **3 m/s**
PGT **C**
Wind-road angle **90°**



Shallow canyon

Windspeed **3 m/s**
PGT **C**
Wind-road angle **90°**



Details another time

Unsupervised learning

Finding structure in unlabeled data

Meteorological-day clustering for faster long-term dispersion simulations

Road-segment classification for urban dispersion kernels and digital twins

Supervised learning

Learning from labeled data

Inverse modelling to rescale emission factors from observations

Mass-consistent concentration downscaling driven by geophysical proxies

Data fusion and concentration time-series forecasting

RemiScaler: Regularized inverse modelling for Emission reScaling

Goal: estimate source-specific emission rescaling factors by combining observed concentrations with modelled source-labelled concentrations.

Working assumptions

- Modelled concentrations are fully source-apportioned (typical of Lagrangian dispersion models).
- No dominant atmospheric chemistry breaks the approximate linearity between emissions and concentrations.

Inverse model

Ridge regression with non-negative constraints and *Tikhonov* regularization

$$\hat{\mathbf{k}} = \operatorname{argmin}_{\mathbf{k}} \left[\|\mathbf{W}^{1/2}(\mathbf{H}\mathbf{k} - \mathbf{c}_{\text{obs}})\|^2 + \lambda \|\mathbf{r} \odot (\mathbf{k} - \mathbf{1})\|^2 \right]$$

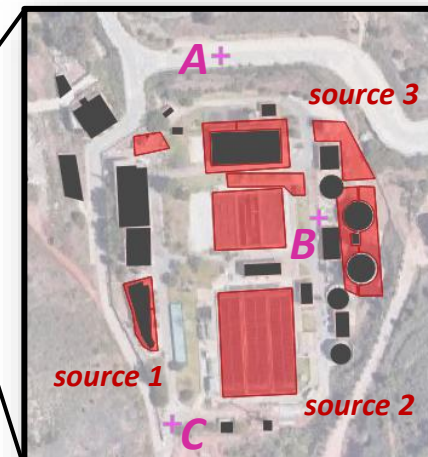
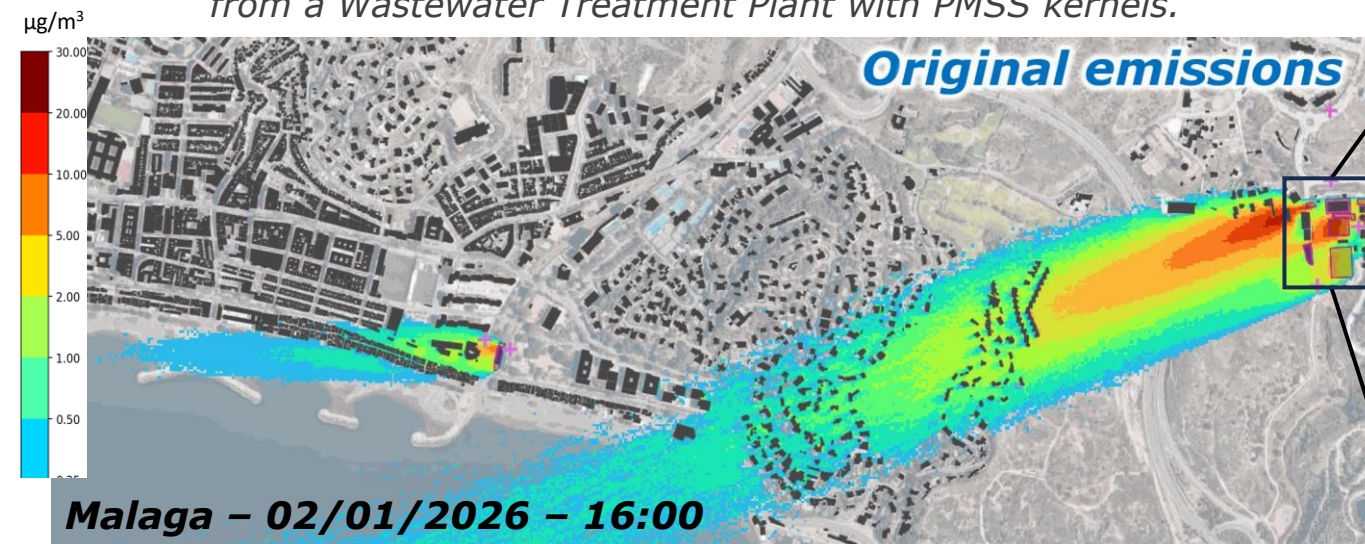
H: Modeled concentration transfer matrix

c_{obs}: Observed concentrations - **W:** Observations weights matrix

k: Emission scaling factors

λ : Ridge hyperparameter - **r:** Source inertias to rescaling

Example application *Forecasting system of H₂S dispersion from a Wastewater Treatment Plant with PMSS kernels.*



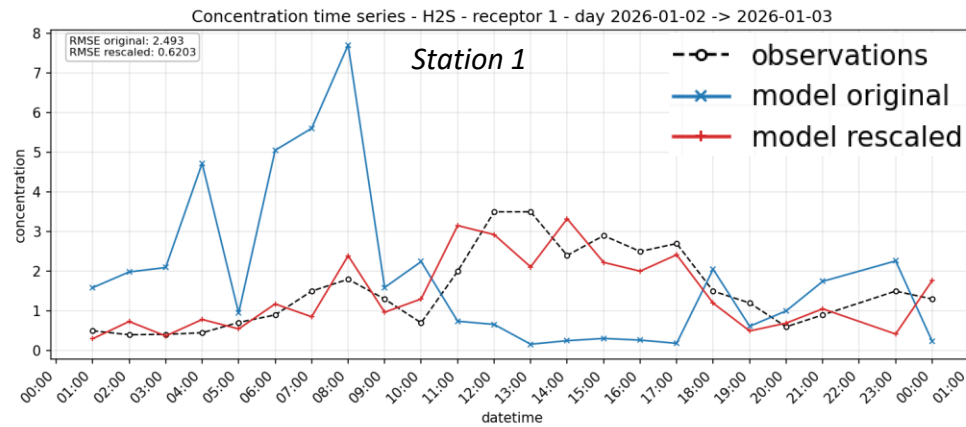
	<i>src</i> 1	<i>src</i> 2	...	<i>src</i> n	H matrix
	H ₁₁	H ₁₂	...	H _{1n}	Obs point A , Time 0
	H ₂₁	H ₂₂	...	H _{2n}	Obs point A , Time 1
	H ₃₁	H ₃₂	...	H _{3n}	Obs point B , Time 0
	H ₄₁	H ₄₂	...	H _{4n}	Obs point C , Time 0

	H _{m1}	H _{m2}	...	H _{mn}	Obs / time <i>m</i>

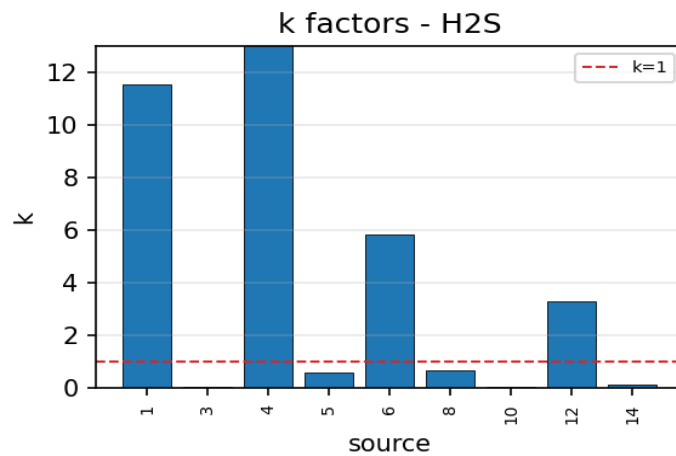
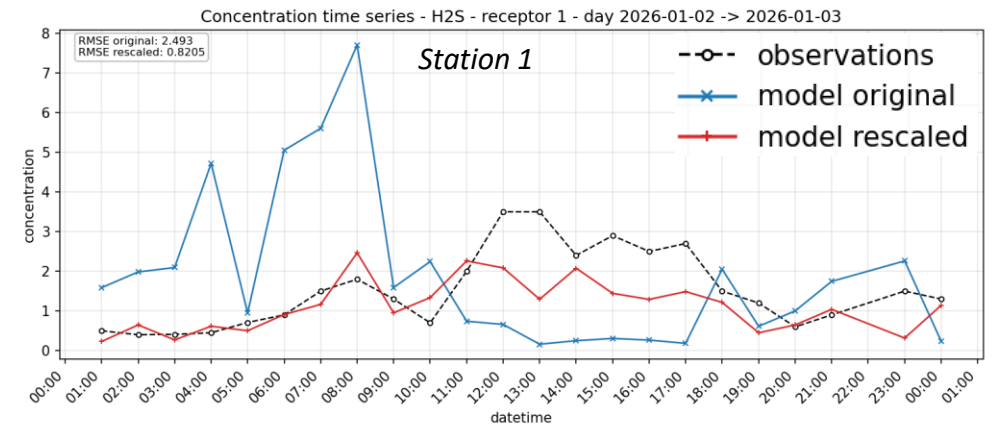
Why regularization is needed

The inverse problem is *ill-posed*: different emission combinations can reproduce the same concentrations*. Regularization stabilizes the solution and keeps scaling factors close to prior values.

Without regularization

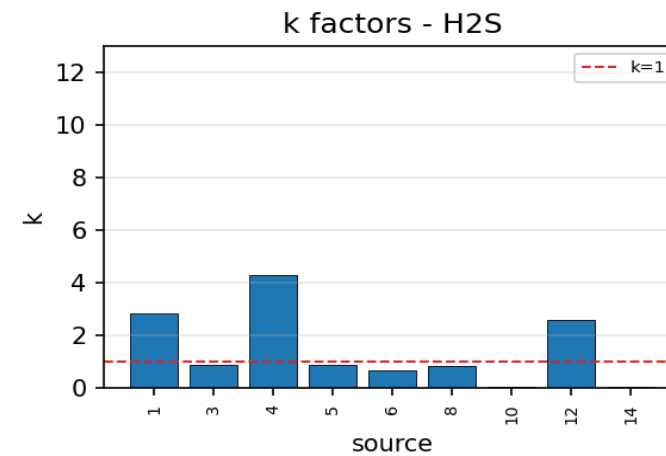


With regularization $\lambda=0,5$



RMSE original emi
2,49 $\mu\text{g}/\text{m}^3$

RMSE rescaled emi
0,62 $\mu\text{g}/\text{m}^3$



RMSE original emi
2,49 $\mu\text{g}/\text{m}^3$

RMSE rescaled emi
0,82 $\mu\text{g}/\text{m}^3$

* Practical lever: more observations in time help constrain the solution and reduce degeneracy between sources.

Unsupervised learning

Finding structure in unlabeled data

Meteorological-day clustering for faster long-term dispersion simulations

Road-segment classification for urban dispersion kernels and digital twins

Supervised learning

Learning from labeled data

Inverse modelling to rescale emission factors from observations

Mass-consistent concentration downscaling driven by geophysical proxies

Data fusion and concentration time-series forecasting

Mass-consistent concentration downscaling with Random Forest

Idea – Use coarse-grid concentration together with fine-grid proxies (*roads, land use, population, ...*) to **redistribute** concentration from 1 km to 100 m

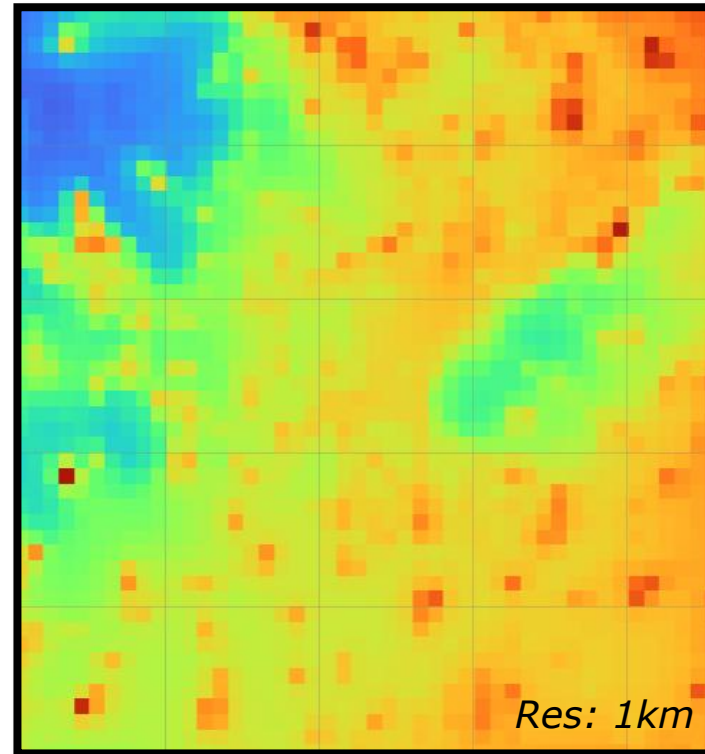
Training – A **Random Forest** is trained on concentration using geophysical proxies as features

$$C_{ML} = RF(\text{proxies})$$

Mass constraint – Within each 1 km FARM cell, the average of the 100m predictions matches the original coarse-cell concentration

$$\hat{C}_{ij}^{ML} = \frac{C_k^F N_i N_j}{\sum_{ij} C_{ij}^{ML}} C_{ij}^{ML}$$

Interpretation – ML does not model concentration at 100 m; it learns **redistribution weights** on the fine grid, while preserving the mass of the original FARM cell



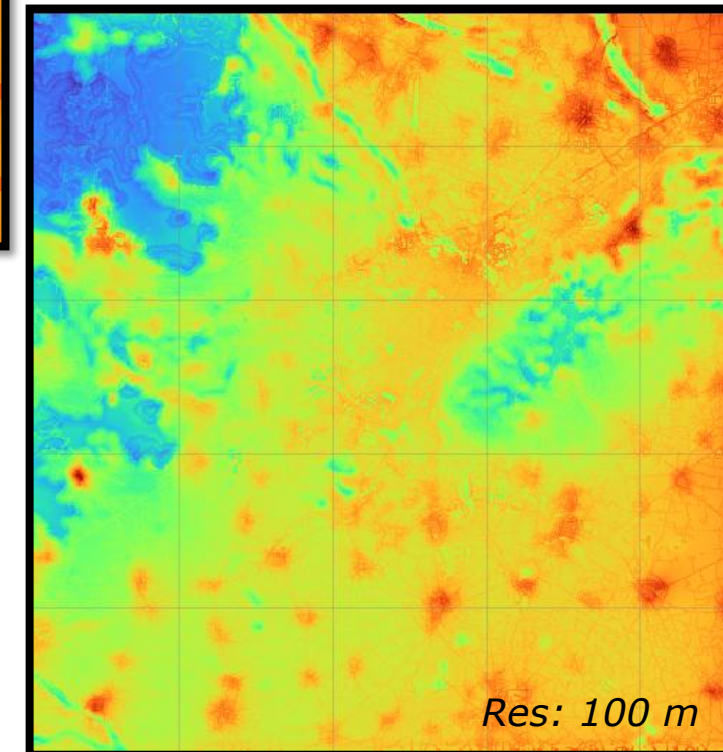
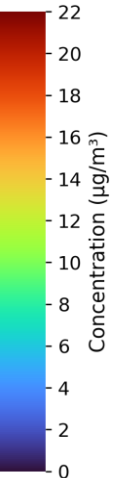
FARM cell



Mass-consistent downscaling



Annual average
PM10 concentration
over Turin



Some projects involved

MISTRAL Horizon: Emission scenarios platform on Rybnik
Andalucia Air Quality Forecasting system

Unsupervised learning

Finding structure in unlabeled data

Meteorological-day clustering for faster long-term dispersion simulations

Road-segment classification for urban dispersion kernels and digital twins

Supervised learning

Learning from labeled data

Inverse modelling to rescale emission factors from observations

Mass-consistent concentration downscaling driven by geophysical proxies

Data fusion and concentration time-series forecasting

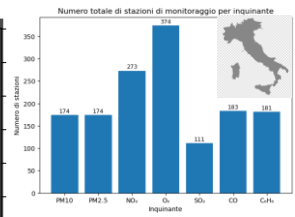
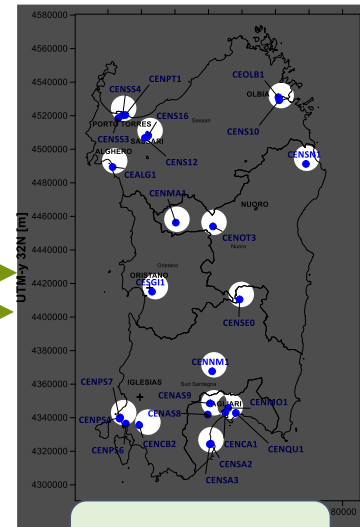
Mal-Air: Random Forest for Data Fusion

Osservazioni
Serie temporali orarie/giornaliere degli inquinanti
PM10, NO2, PM2.5

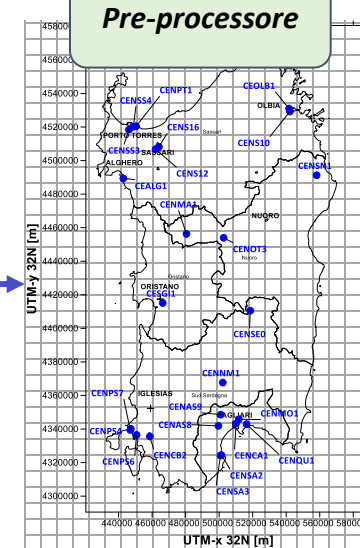
Predittori - Proxies
Spazio-temporali
FARM/CAMS (# km)
Leaf Area Index (LAI)
Meteorological (ERA5)

Temporali (homogeneous)
Periodic functions of Julian day, day of week

Spaziali (stazionari)
Length roads in buf.
Population density
Land use
...



Allenamento (2024)



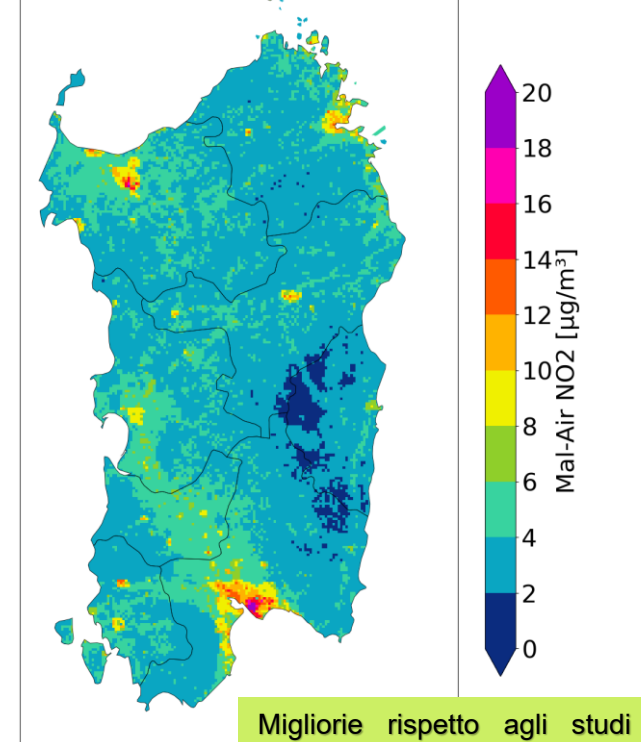
Pre-processore

Random Forest

Inferenza (2024)

7 modelli (PM10, PM2.5, NO2, O3, SO2, C6H6, CO)

Risoluzione target: 1km



Migliorie rispetto agli studi precedenti

- Proxy meteo ad alta risoluzione
- Prodotti satellitari
- Approccio multi predittore CAMS applicato a tutti gli inquinanti

Allenamento —————
Inferenza —————

Dati di Input
Modello / processo

Some projects involved:
Valutazione modellistica dello stato di Qualità dell'aria sulla Regione Sardegna per l'anno 2024; SPOTT; CALLIOPE

Classi di uso del suolo – Corine Land Cover 2018

Aree urbane (tessuto urbano continuo e discontinuo)

Aree industriali e commerciali

Aeroporti e infrastrutture aeroportuali

Altre superfici artificiali

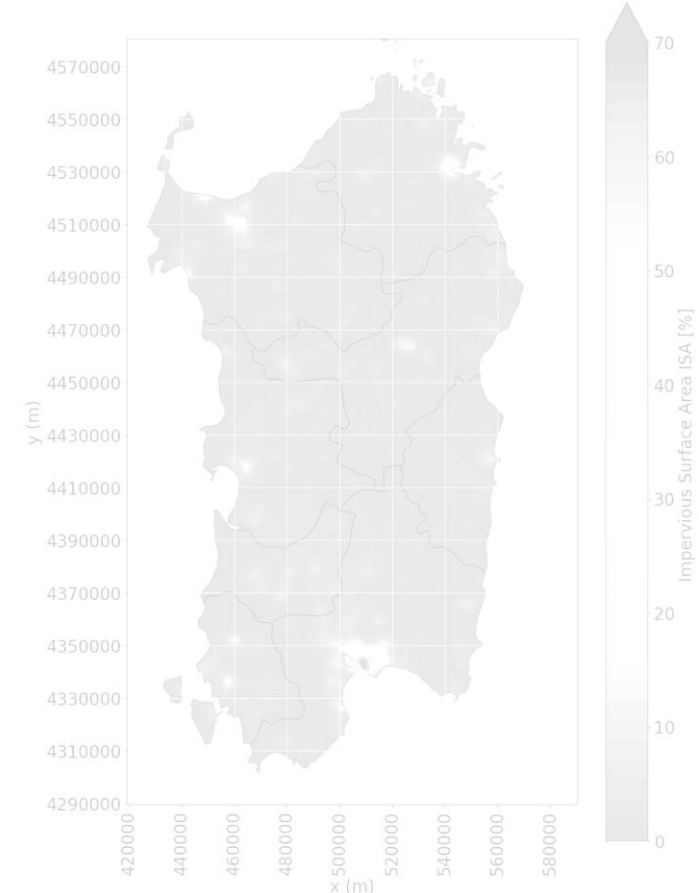
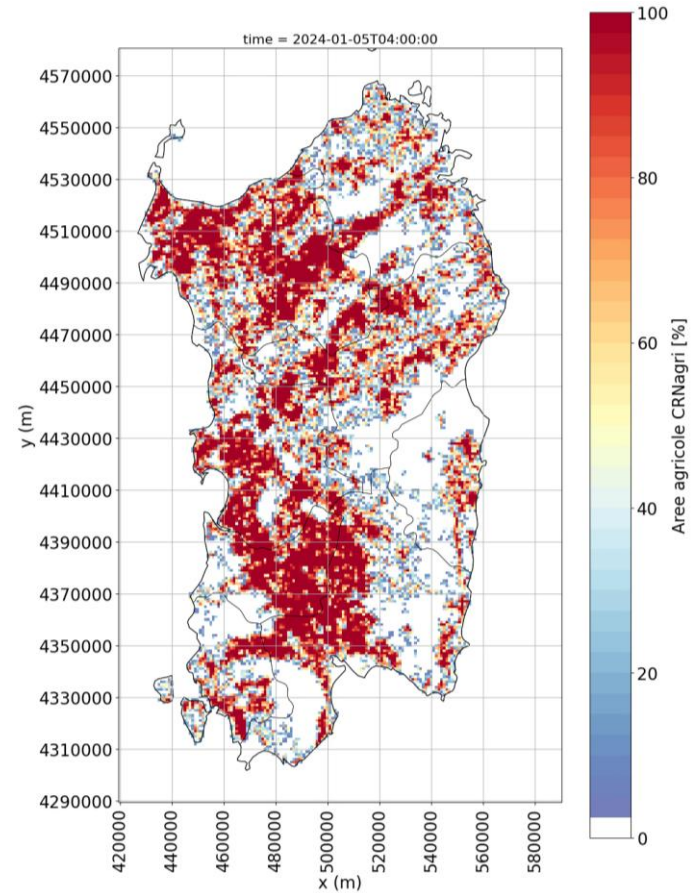
Aree agricole

Aree forestali

Prati, aree arbustive e superfici naturali aperte
(pianure brulle, spiagge naturali)

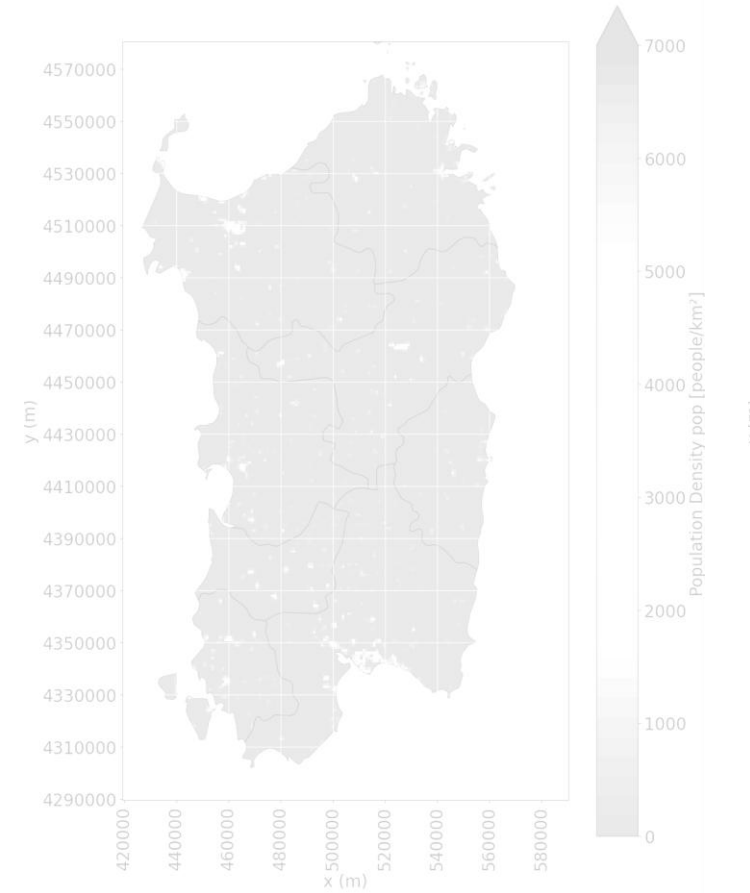
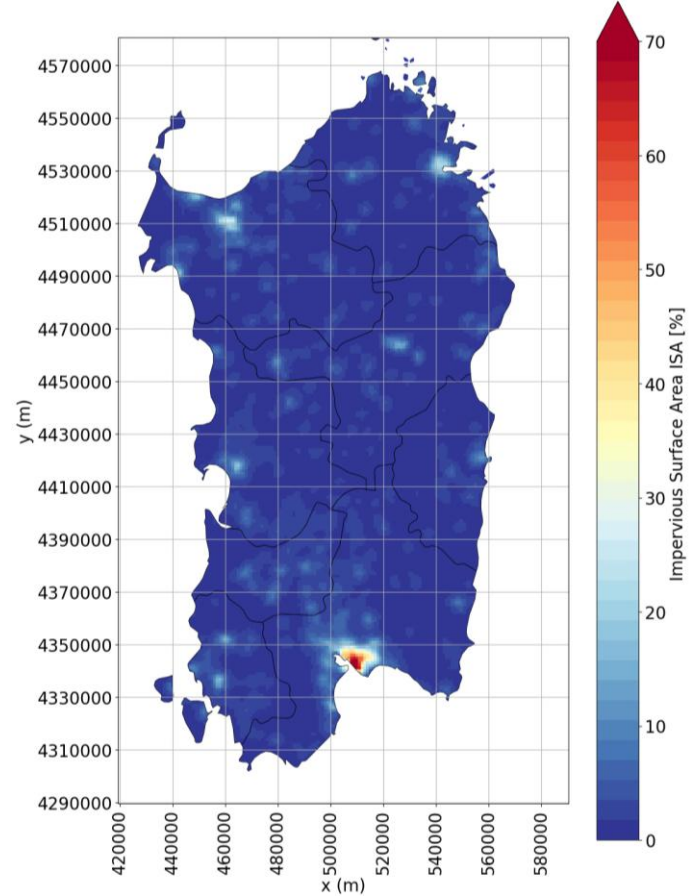
Risoluzione orig: 100 m

Predittore statico

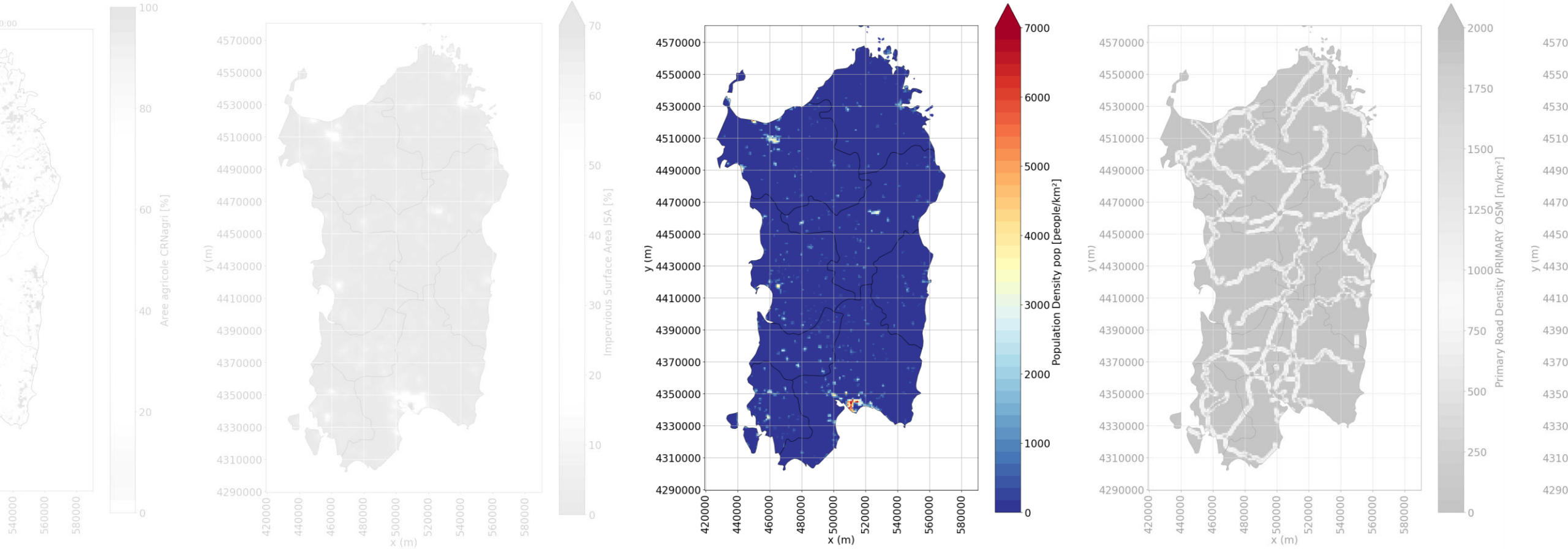


Inventario globale della distribuzione spaziale e della densità della superficie impermeabile costruita

Predittore Statico



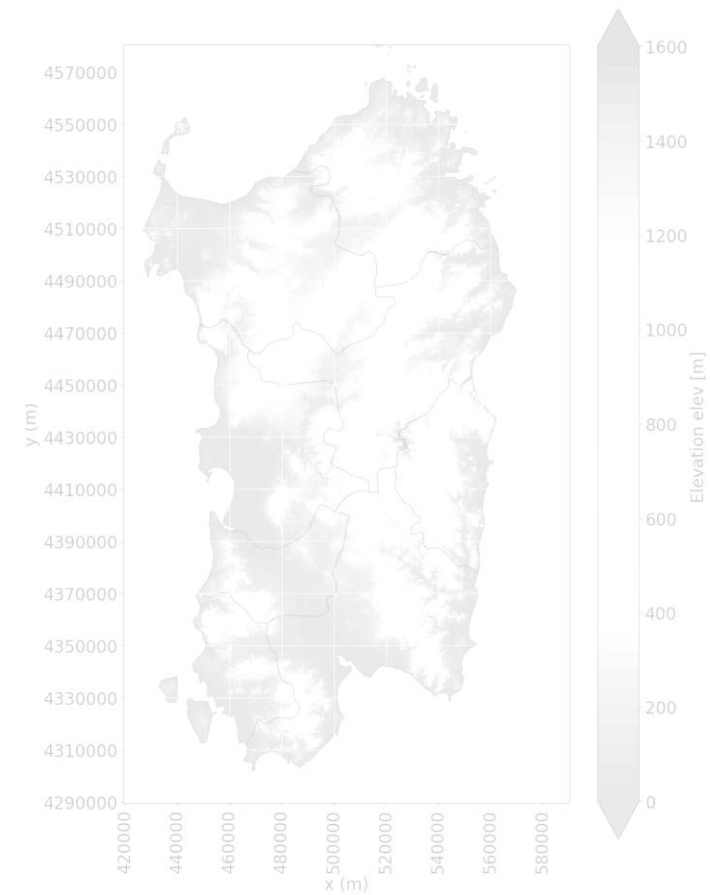
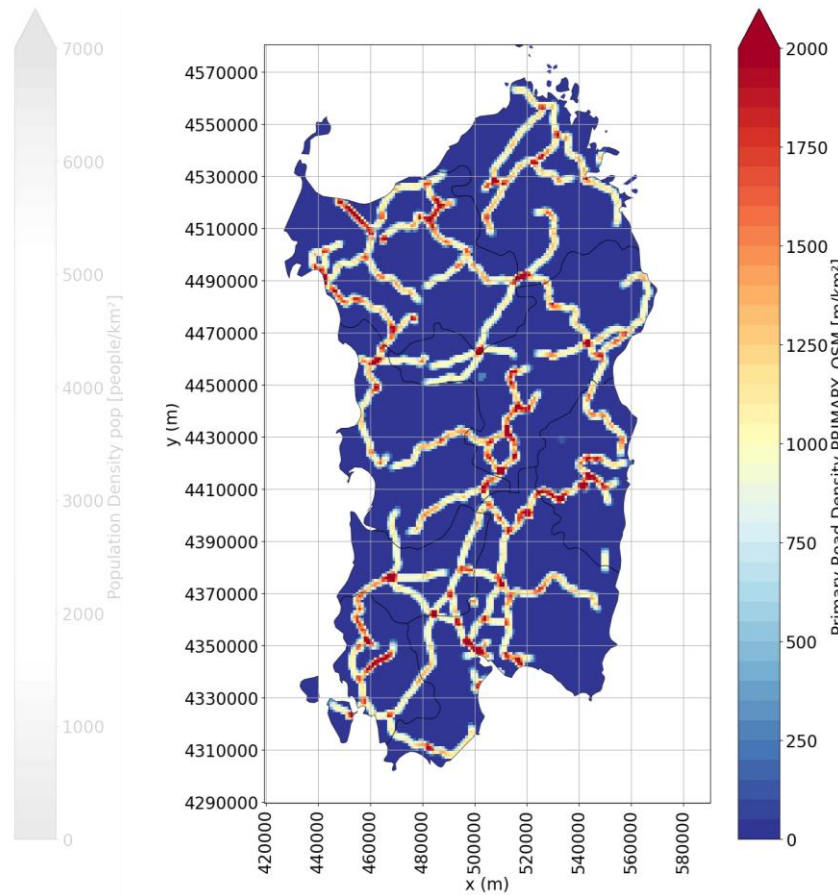
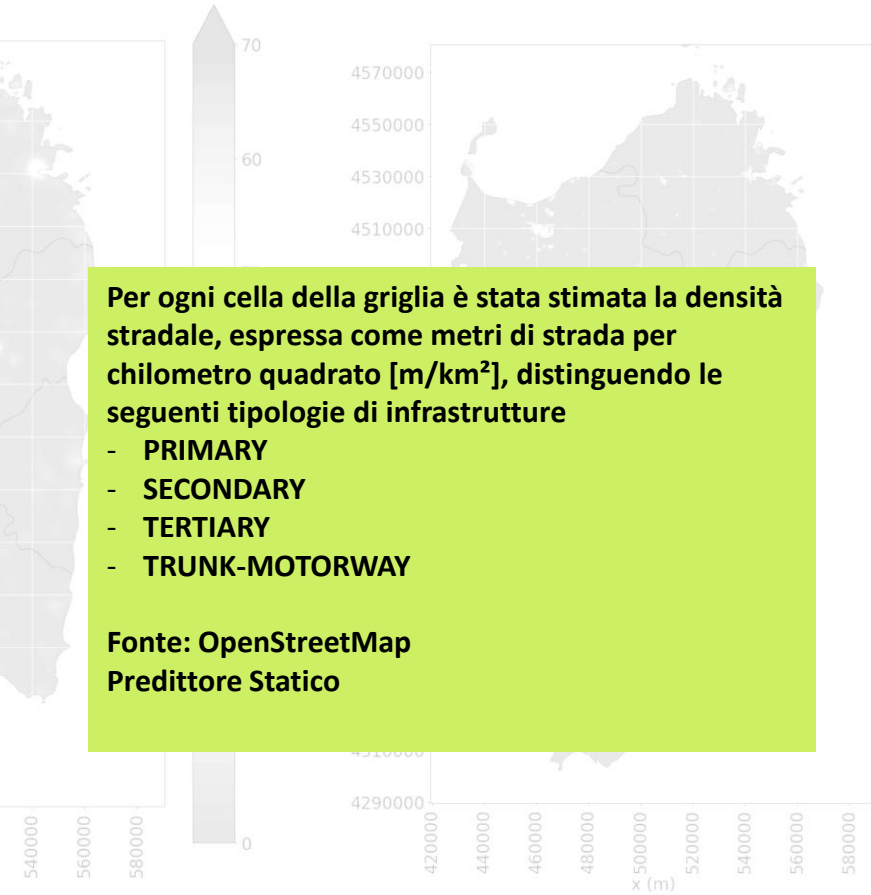
PREDITTORI



Per ogni cella della griglia è stata stimata la densità stradale, espressa come metri di strada per chilometro quadrato [m/km²], distinguendo le seguenti tipologie di infrastrutture

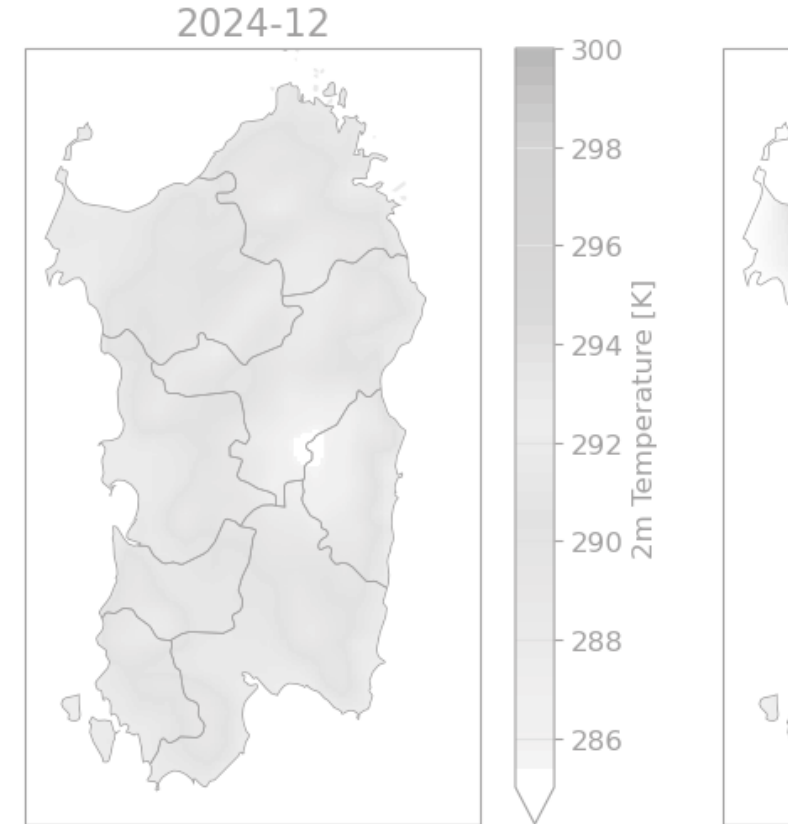
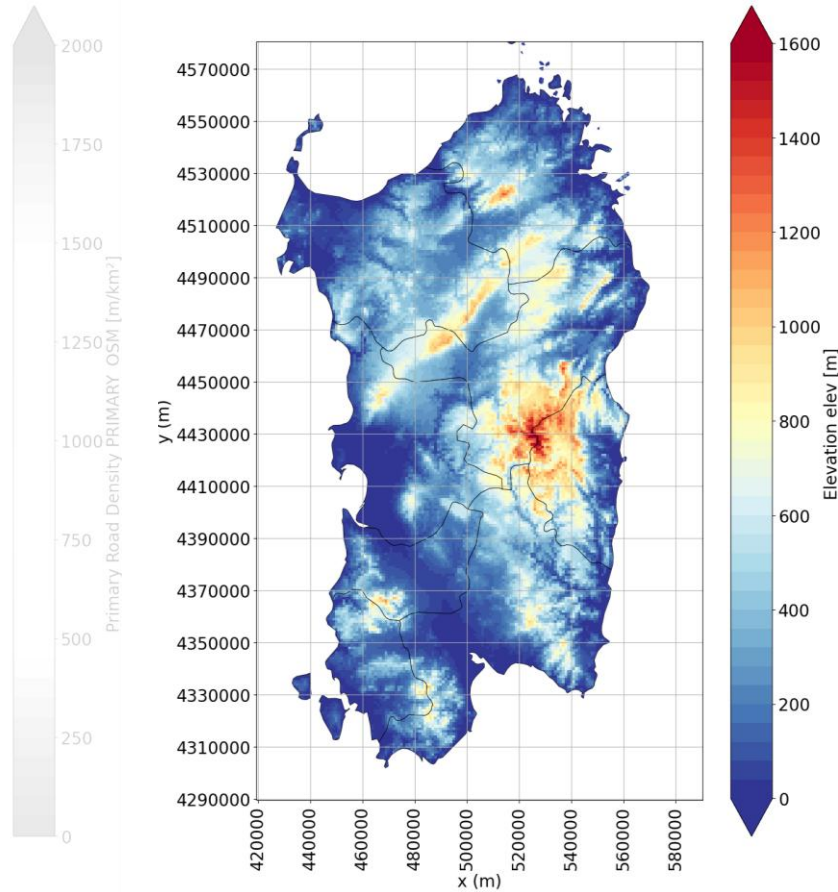
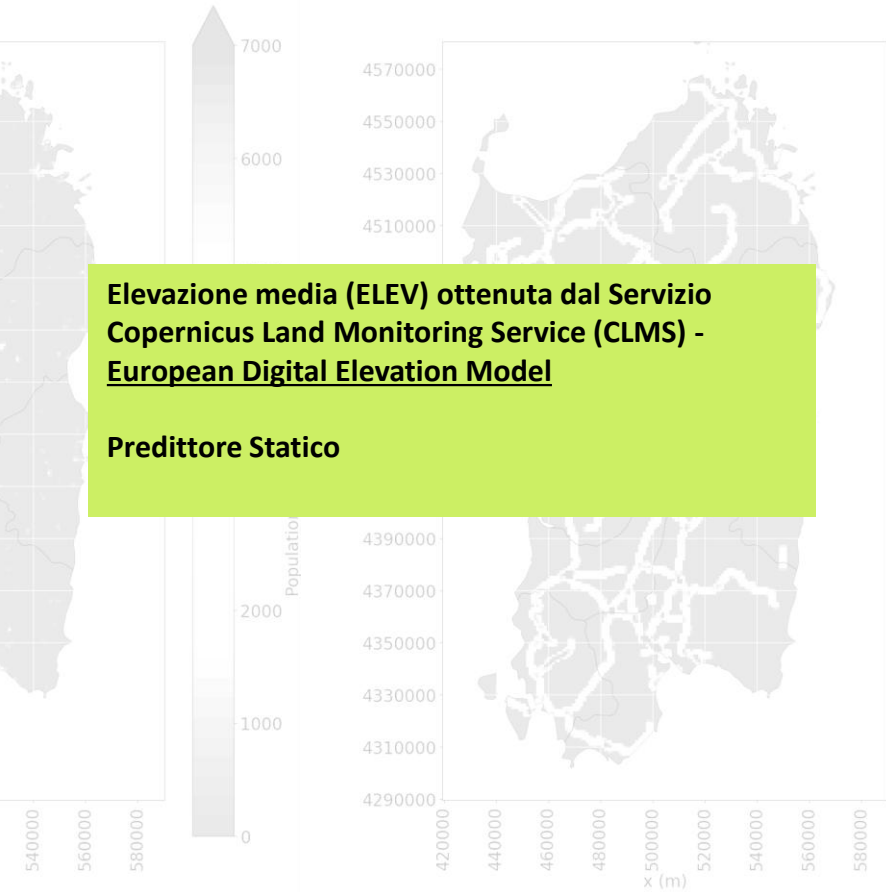
- PRIMARY
- SECONDARY
- TERTIARY
- TRUNK-MOTORWAY

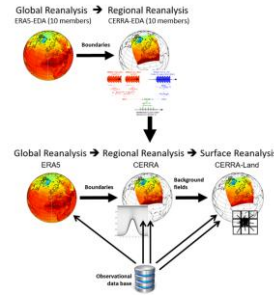
Fonte: OpenStreetMap
Predittore Statico



Elevazione media (ELEV) ottenuta dal Servizio Copernicus Land Monitoring Service (CLMS) - European Digital Elevation Model

Predittore Statico

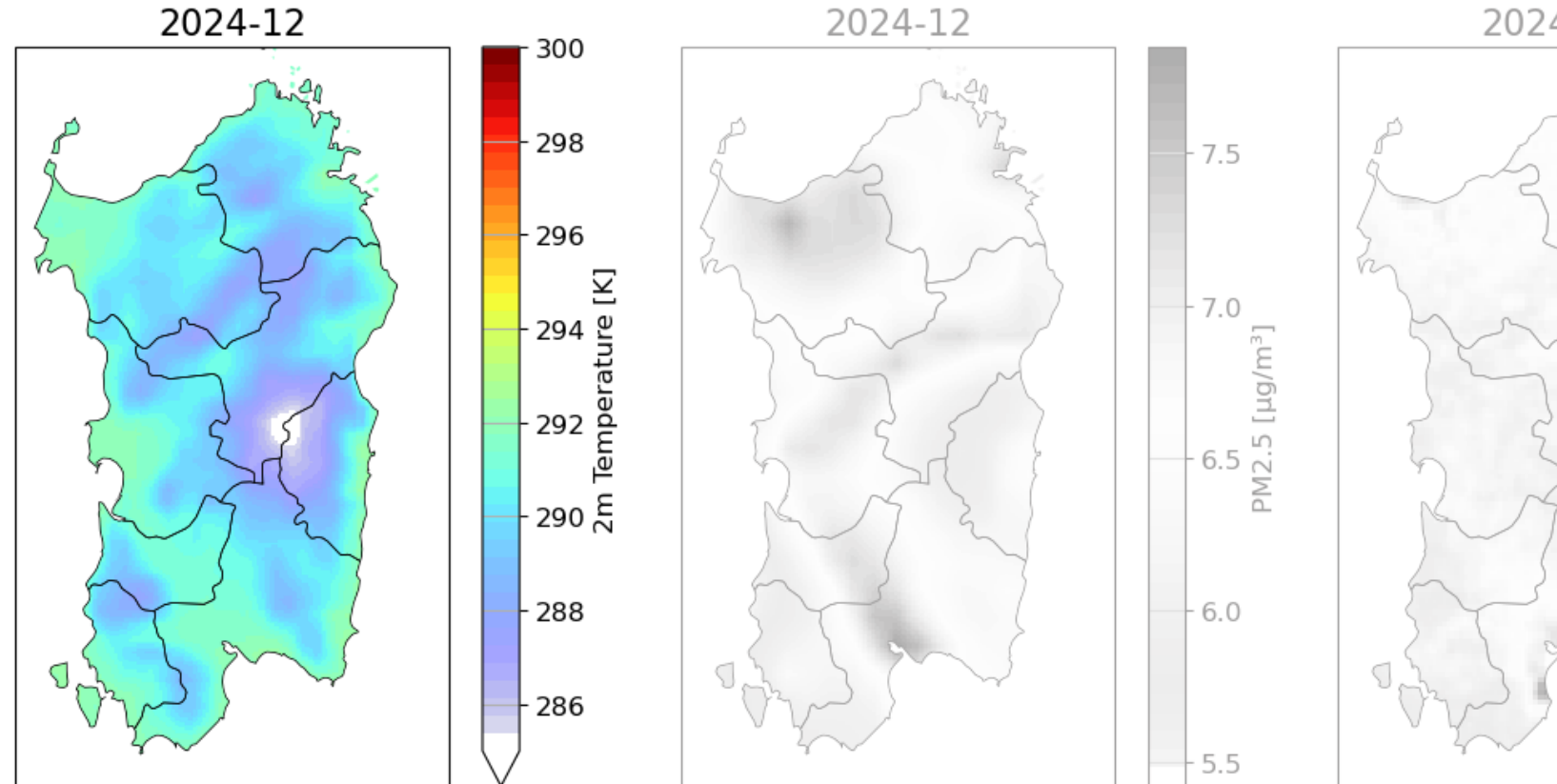




CERRA – [Copernicus regional reanalysis for Europe \(CERRA\)](#) | Copernicus

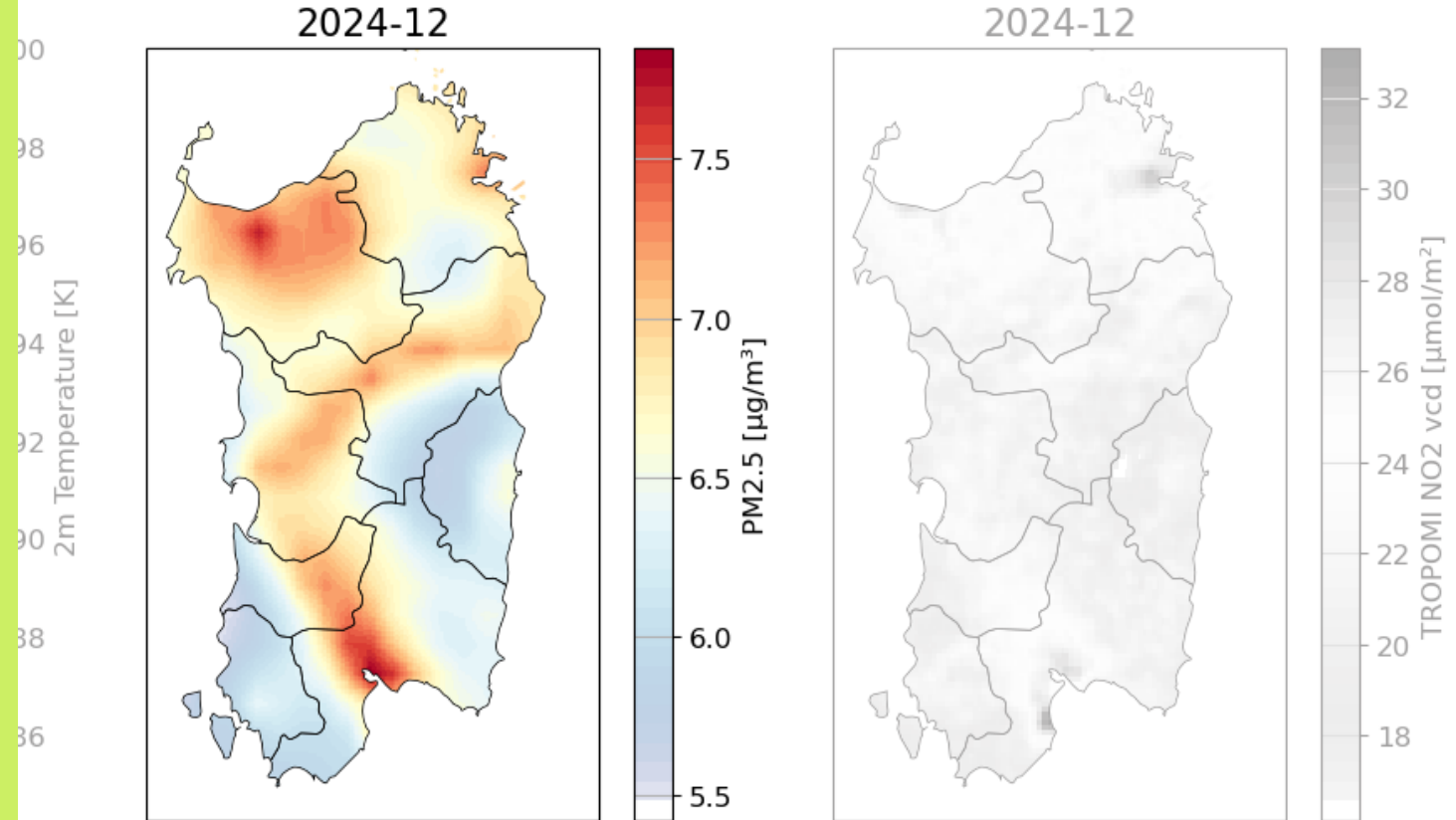
Dataset sviluppato con il sistema numerico **HARMONIE-ALADIN**, forzato ai bordi da **ERA5**

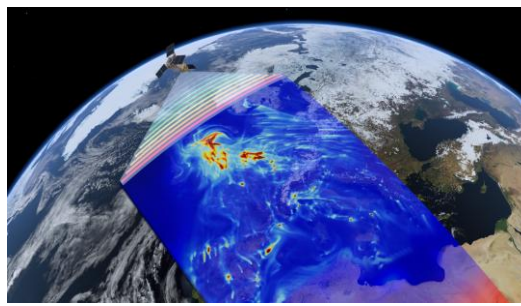
- **Configurazione utilizzata nello studio**
- **Prodotto:** *CERRA sub-daily regional reanalysis on single levels*
- **Risoluzione spaziale:** ~5.5 km
- **Risoluzione temporale:** 3 ore
- **Periodo:** fino a novembre 2024 (*dicembre 2024: campi ERA5*)
- **Variabili impiegate**
- Vento a 10 m (velocità e direzione, espressa come componenti seno/coseno), t2m



Campi di concentrazione **orari** di:
PM10
PM2.5
NO₂
O₃
SO₂
CO
DUST (PM10)

Risoluzione spaziale **~10 km**
PM10 e PM2.5 aggregati su **base giornaliera**
Altri inquinanti utilizzati su **base oraria**
Utilizzo **congiunto di più inquinanti** per tenere conto delle **interazioni chimico-fisiche in atmosfera**

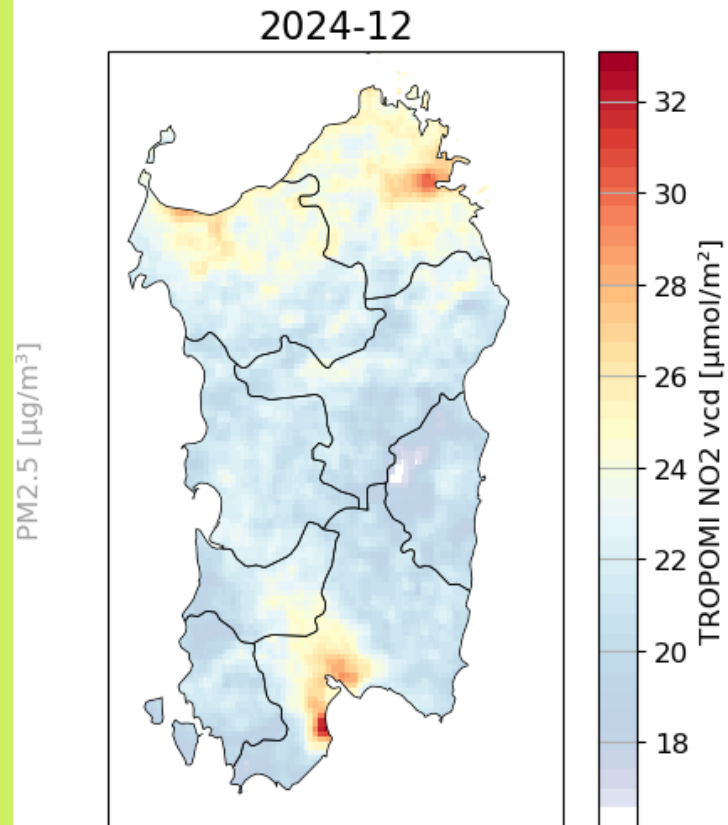




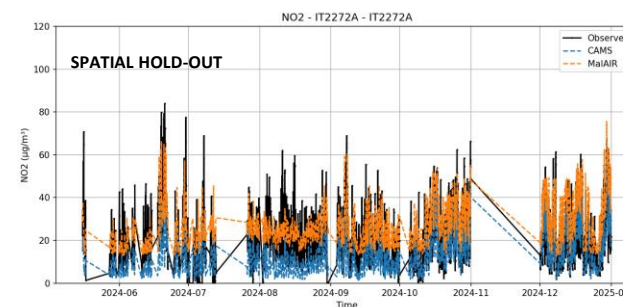
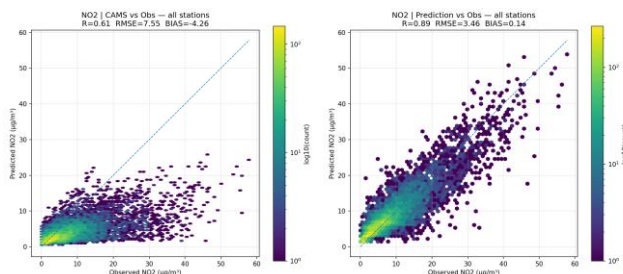
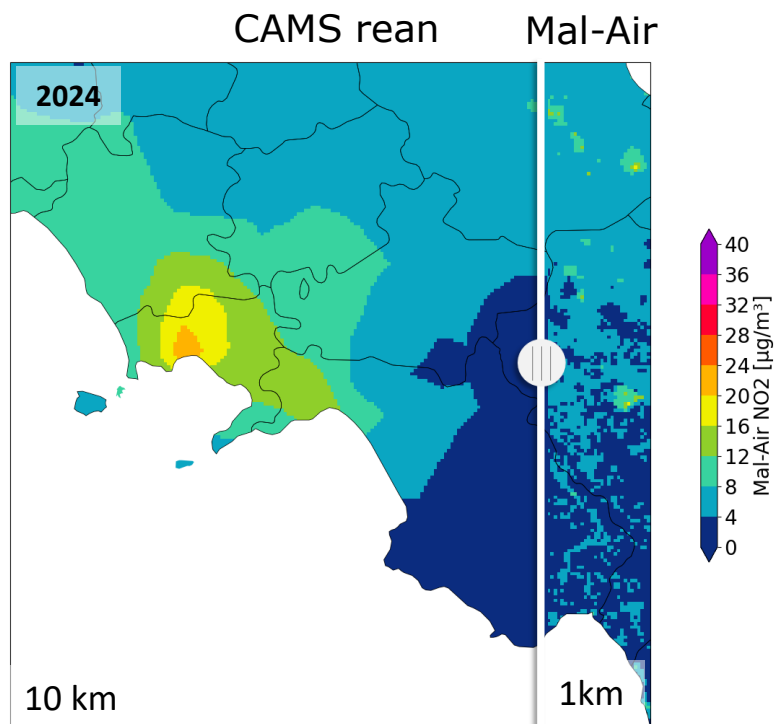
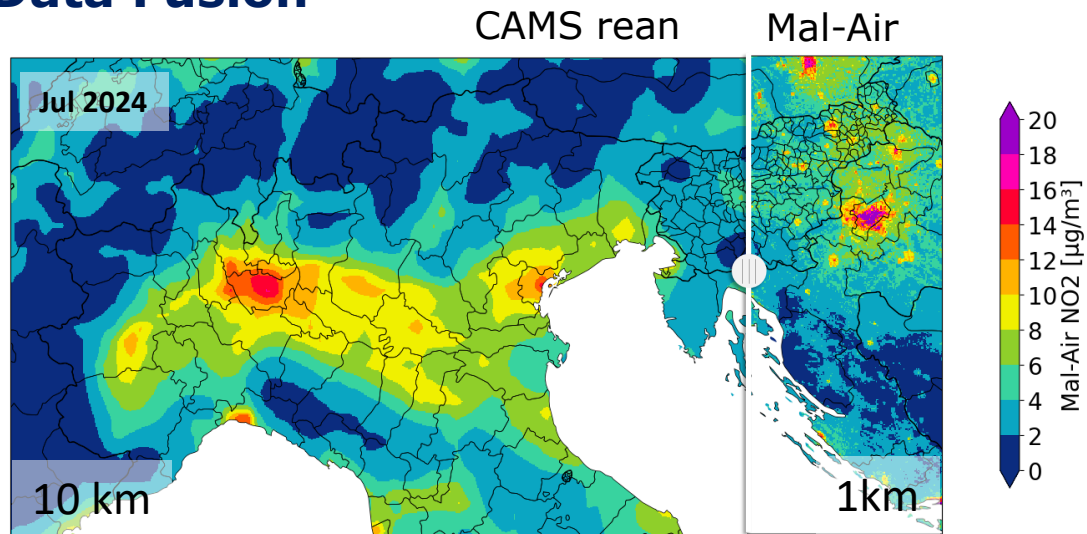
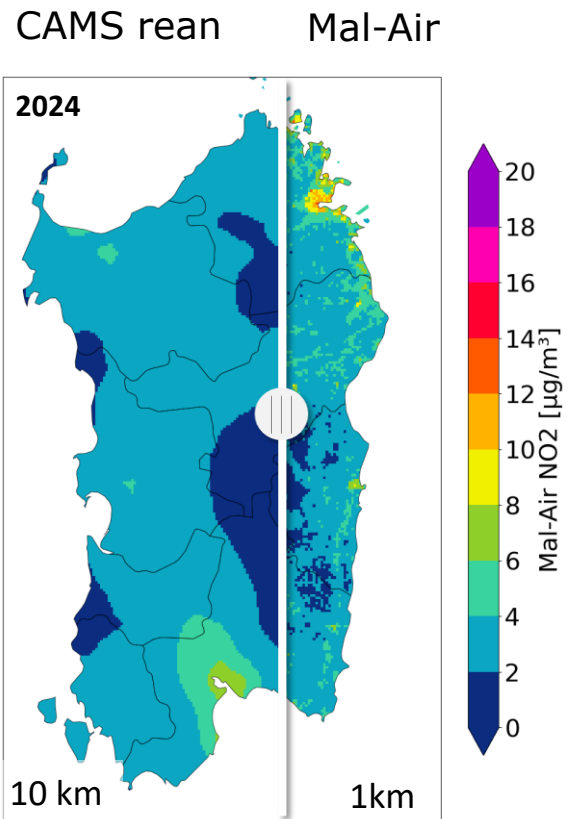
NO₂ troposferico da Sentinel-5P / TROPOMI (Copernicus)

- Colonna troposferica di NO₂ ($\mu\text{mol}/\text{m}^2$), **media mensile – anno 2024**
- **Risoluzione spaziale** $\sim 5.5 \times 3.5 \text{ km}^2$ (al nadir)
- Evidenza i principali **hotspot emissivi**:
 - aree urbane e industriali
 - principali assi viari
- Media calcolata su **tutte le osservazioni mensili** disponibili sulla regione
- **Orario di passaggio fisso** del satellite ($\sim 11:00-15:00$ locali):
 - rappresenta la **condizione tipica diurna**
 - **non** una media sull'intero ciclo giornaliero
- Dati disponibili dal **Copernicus Data Space**

Predittore dinamico

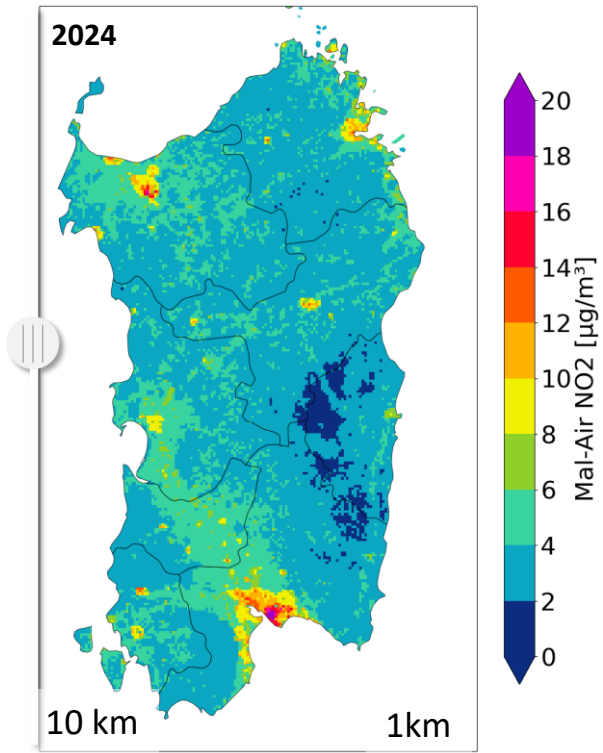


Mal-Air: Random Forest for Data Fusion

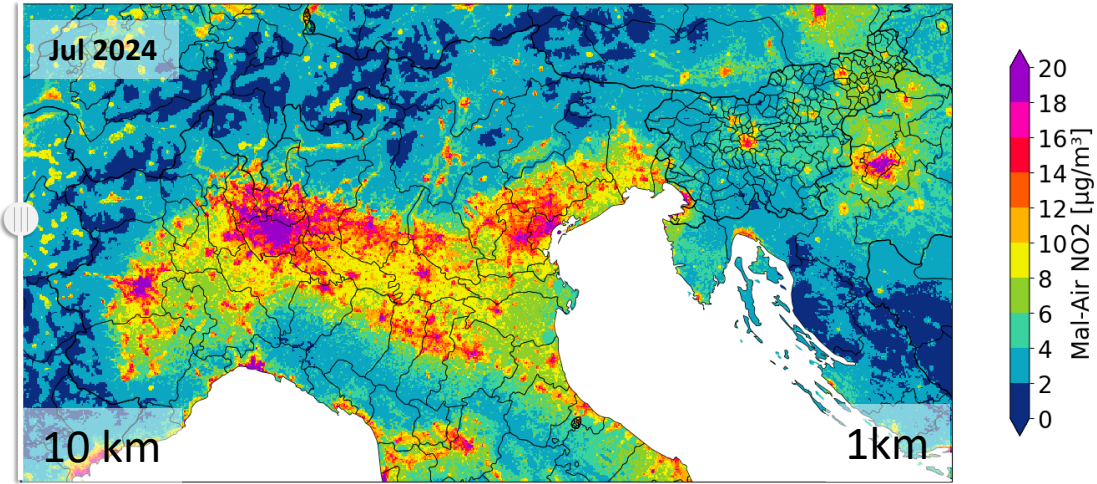


Mal-Air: Random Forest for Data Fusion

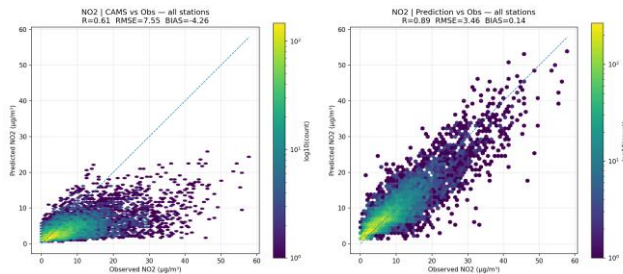
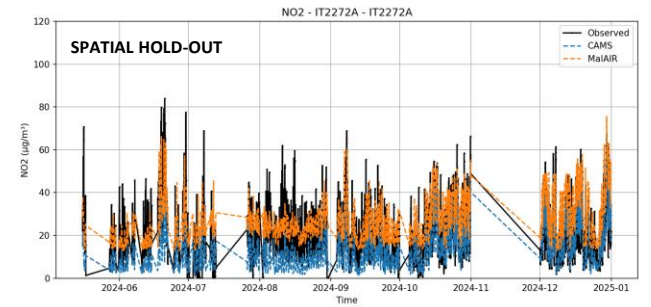
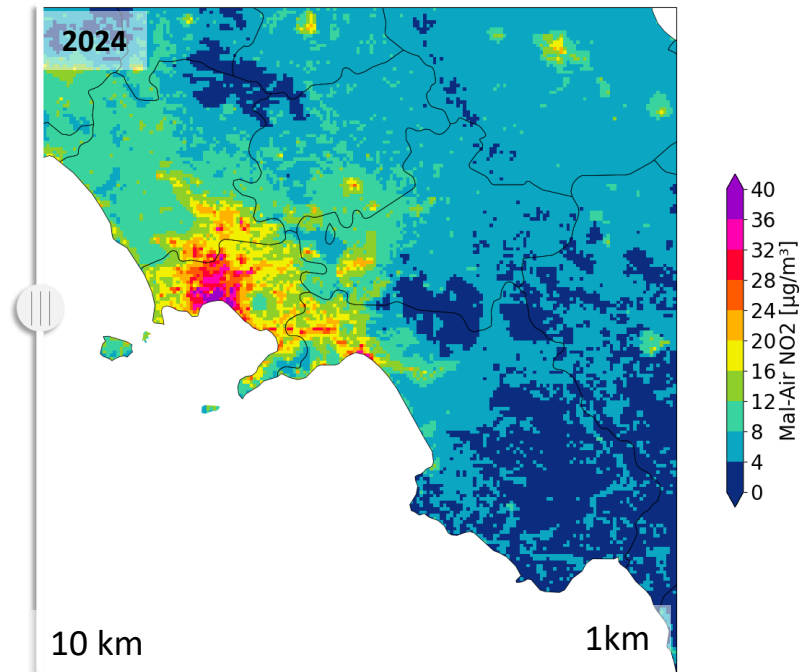
Mal-Air



Mal-Air



Mal-Air



Concentration time series forecasting

Local Forecasting Model

BASELINE

NaiveMean
NaiveSeasonal
NaiveDrift
NaiveMovingAverage

STATISTICAL

FFT
Kalman filter
ARIMA
VARIMA
Prophet
AutoARIMA

REGRESSION

XGBoost
LinearRegression
RandomForest
LightGBM
CatBoost

Global Forecasting Model

FOUNDATION

TimesFM-Google
Chronos-AWS

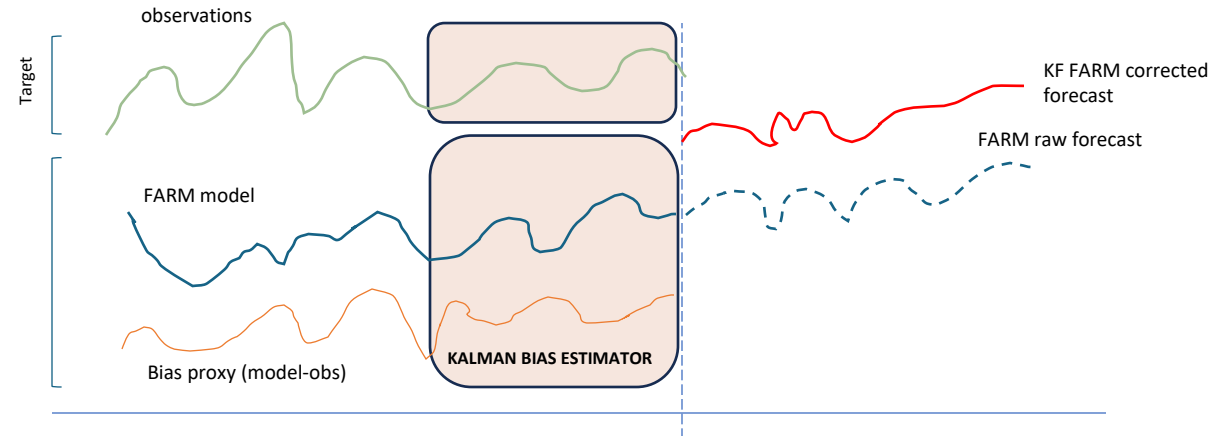
DEEPLARNING

NBEATS
TiDEModel
RNN-GRU
RNN-LSTM

Concentration time series forecasting with Neural Networks

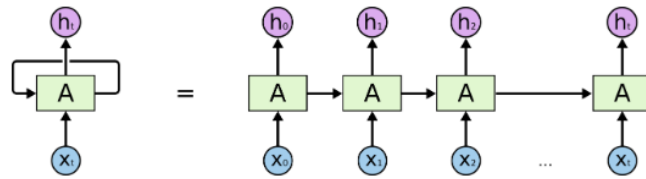
Kalman Filtering

linear, sequential, Gaussian assumptions

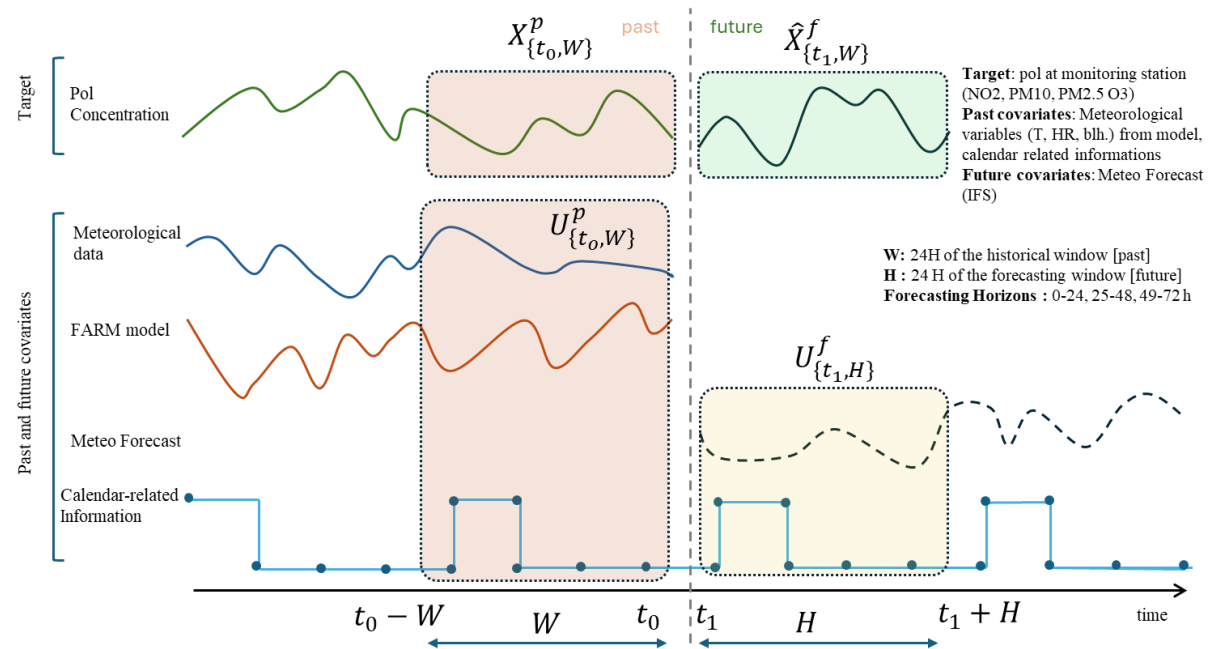
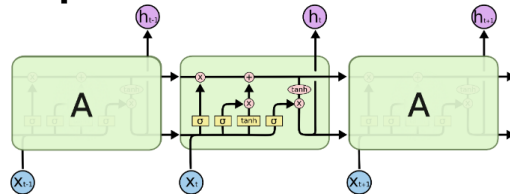


Recurrent Neural Networks

nonlinear, data-driven, multi-feature

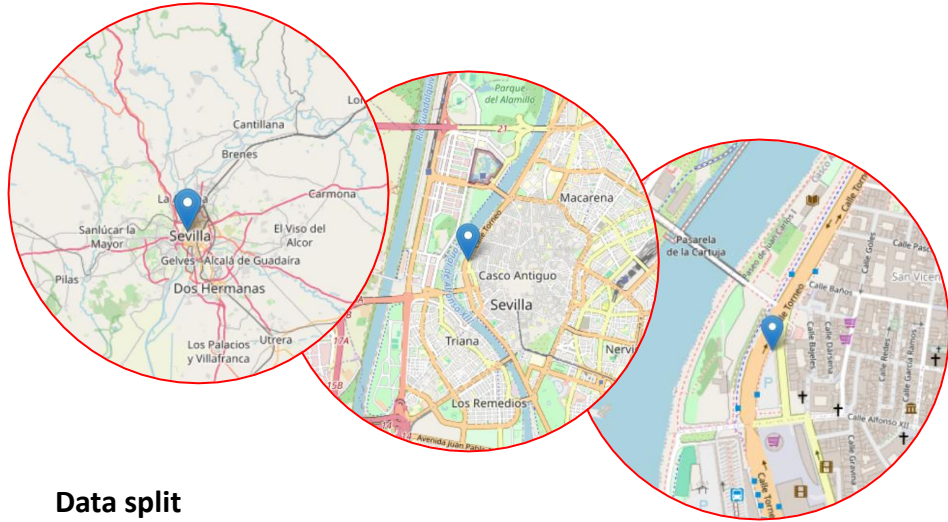


The problem of Long-Term Dependencies - LSTM



Recurrent Neural Network - example

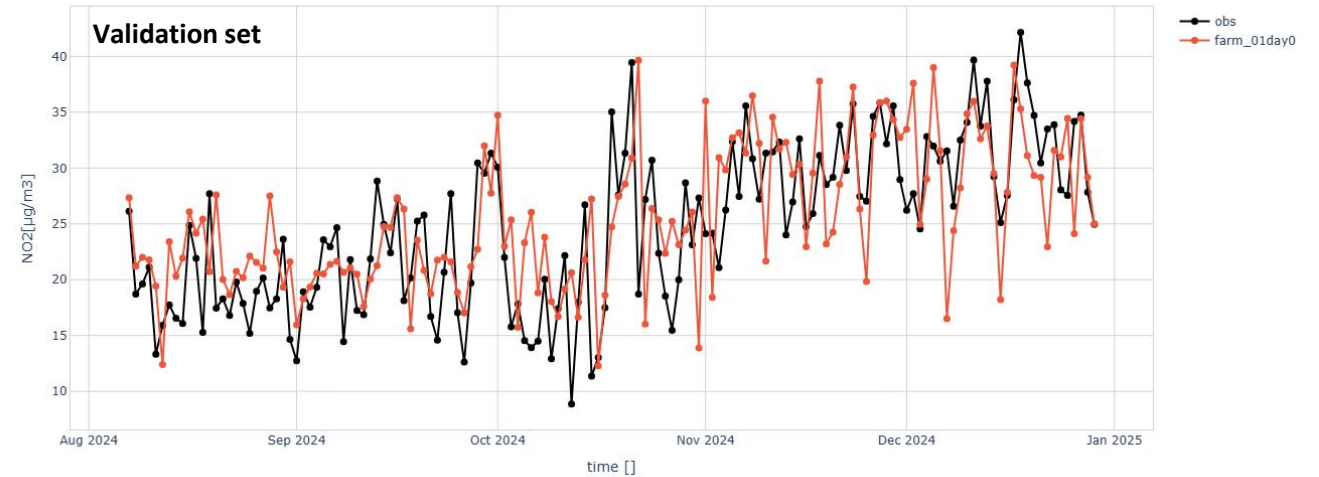
ES0890 – TORNEO Sevilla [Traffic]



Model : **BlockRNNModel-LSTM**

Loss : **MSE**

Time series - pol: NO2 Station: ES0890A Agg: da tstart: 2024-08-07 tend: 2024-12-31

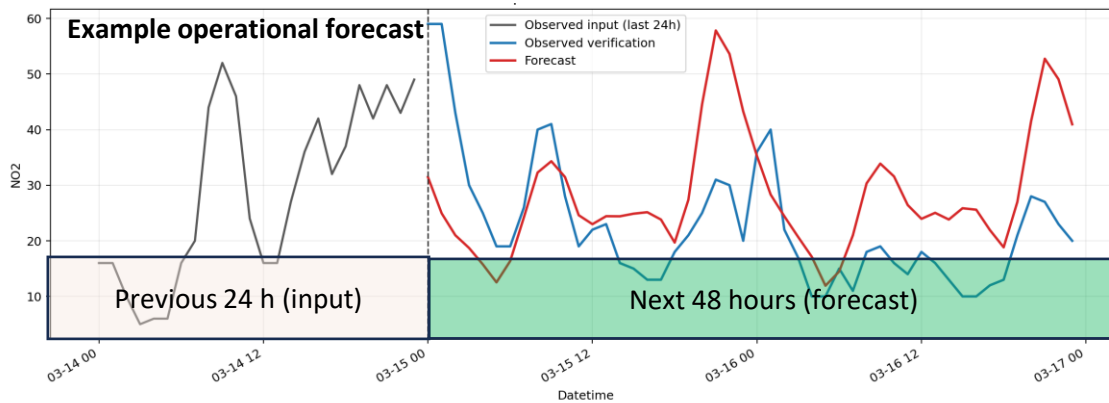


Data split

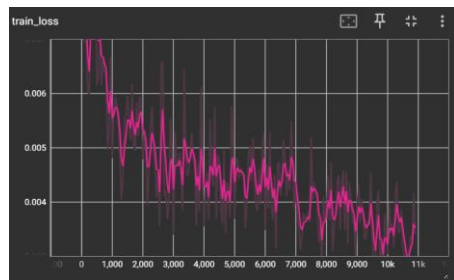
- Training : 2023 -> Aug 2024
- Validation : Aug 2024 -> Dec 2024
- Test : 2025 -> Jan 2026

Input

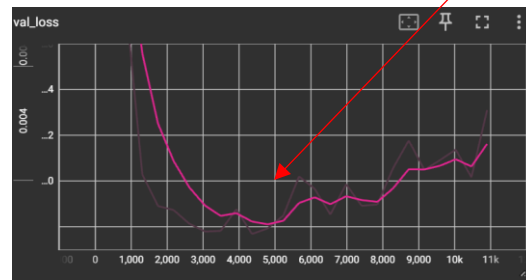
- Future covariates (+48h): CAMS reg fc (NO2, PM10)
- Past covariates (-24h): CAMS reg fc (NO2, PM10); EEA obs
- Calendar related info: hours_x, hours_y



Loss on training set

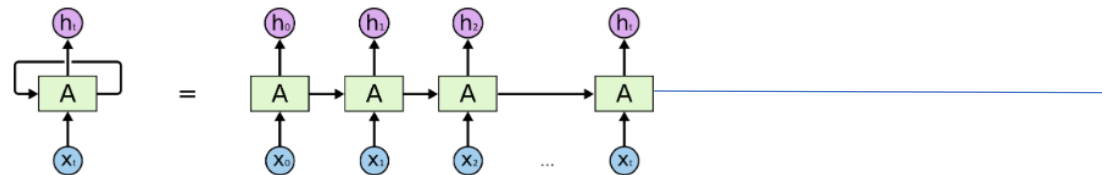
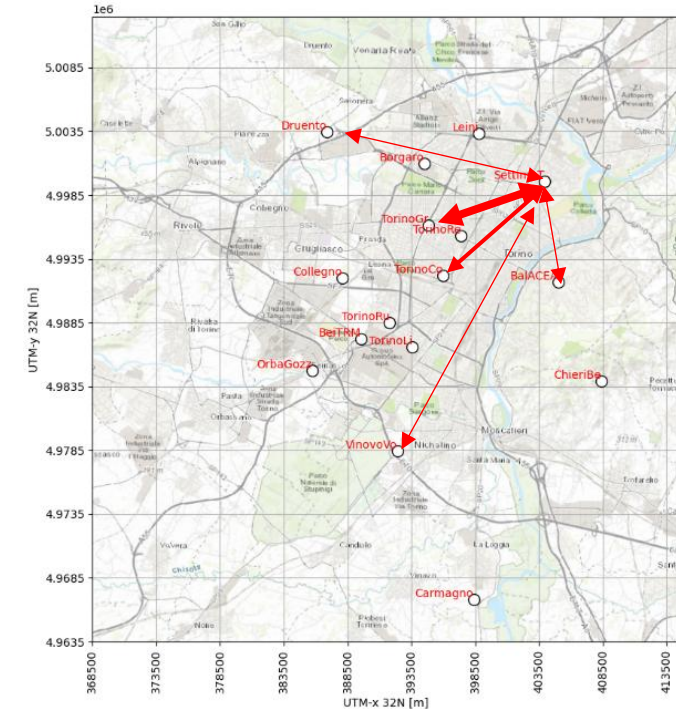
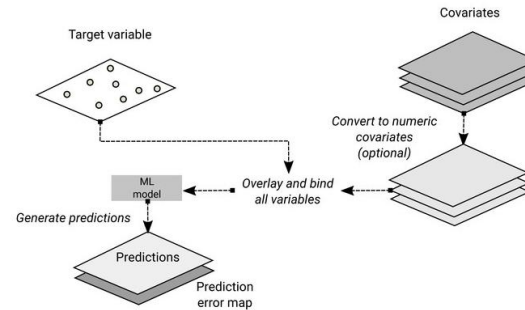
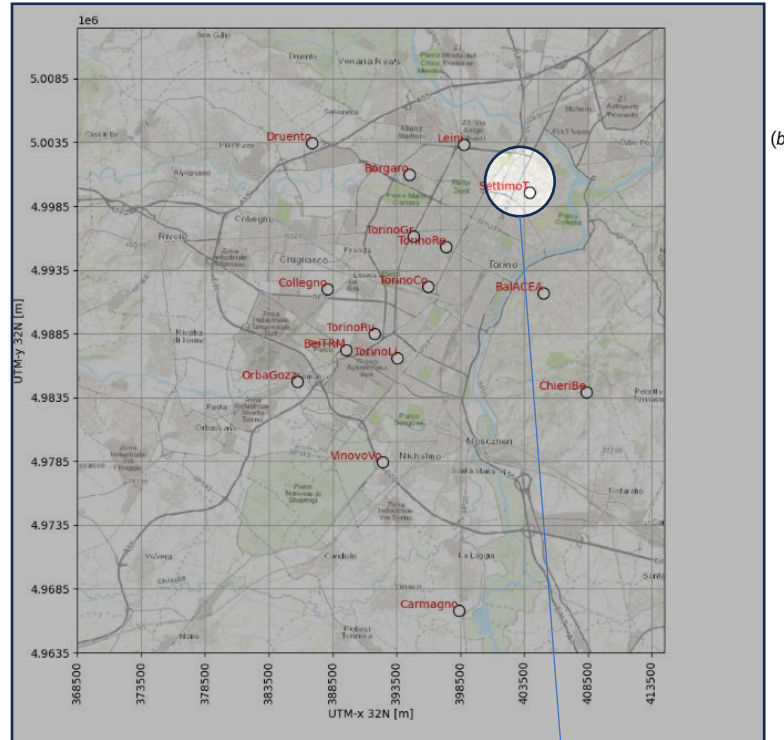


Loss on validation set



Recurrent Neural Network + Data fusion

Supervised learning (Random Forest)



Predictors accounting for geographical proximity
(and generic distance-based relations)
between observations
(to mimic spatial correlation used in Kriging
methods)