# Building Foundational Models for Environmental Modelling and Prediction

Christian Lessig, Ilaria Luise, Martin Schultz,
Michael Langguth, Alberto di Meglio et al.

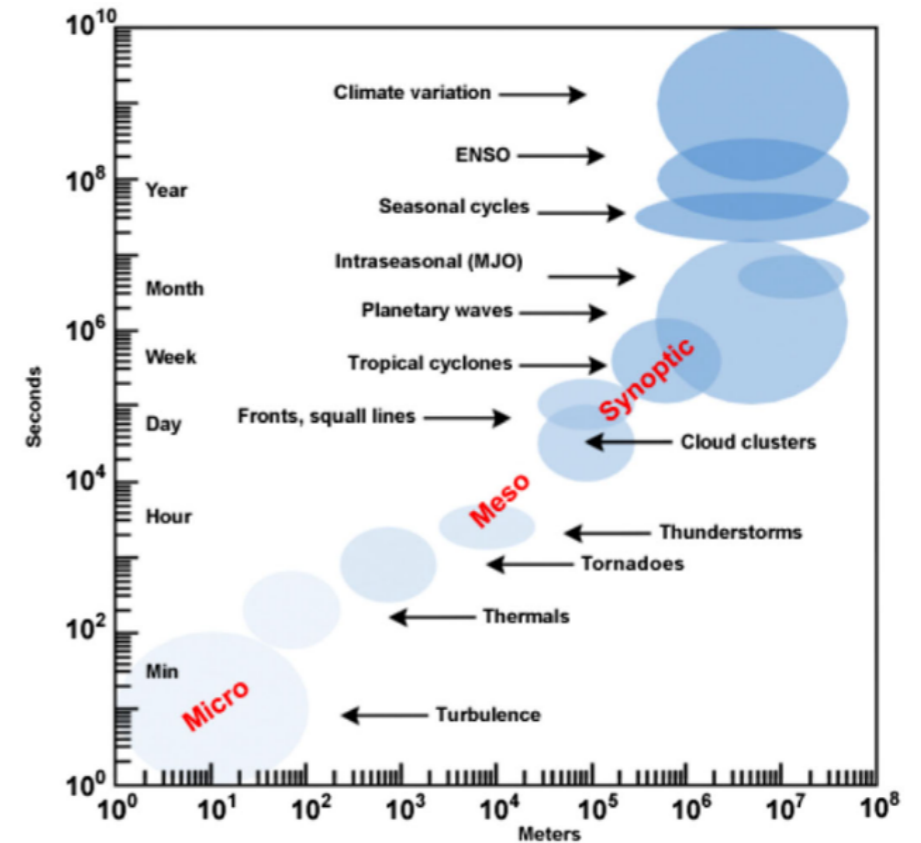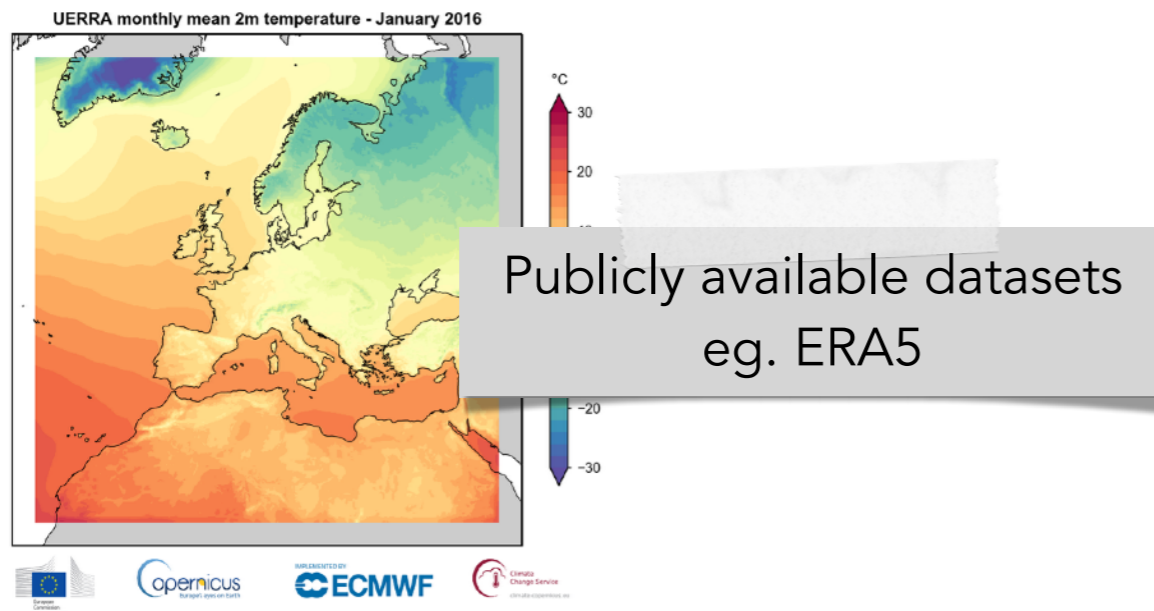ACAT 2024 | March 2024 - Stony Brook

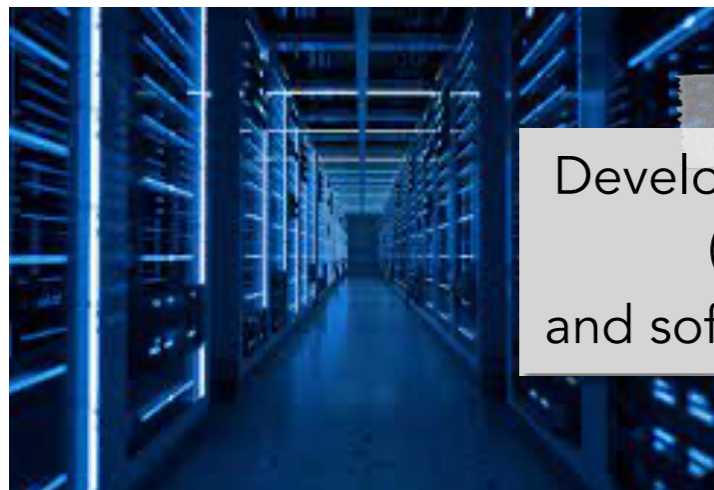# Motivation and scientific challenge

**Atmosphere:**

- Complex phenomena involving multiple scales
- No complete classical model to simulate the dynamics
- Very large amounts of **observational** data available



UERRA monthly mean 2m temperature - January 2016

Publicly available datasets
eg. ERA5



From V. M. Galfi, V. Lucarini, F. Ragone, and J. Wouters. Applications of large deviation theory in geophysical fluid dynamics and climate science. La Rivista del Nuovo Cimento, 44(6):291–363, 2021.

**We have hundreds of TB of available atmospheric observations.**

**Can we use the information in these datasets for the next generation of improved weather and climate models?**
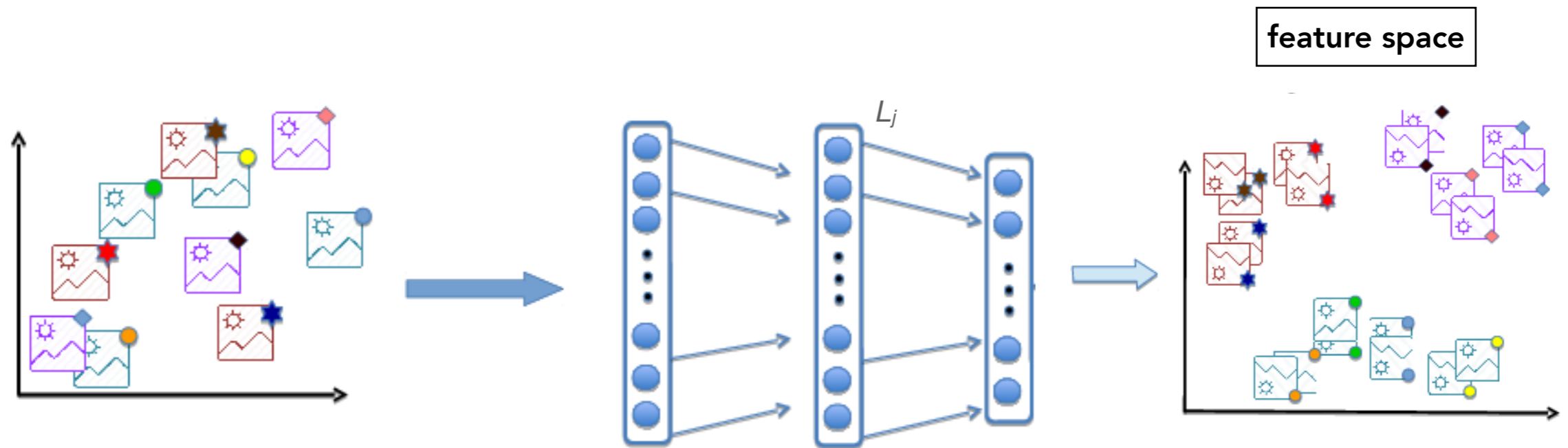
Developments in hardware
(GPU clusters)
and software (exa-scale ML)

Ilaria Luise, CERN - ilaria.luise@cern.ch

# Key ingredient: Representation learning

Representation learning:
- Learn a **task-independent representation** of the data in the **feature space** of the neural network
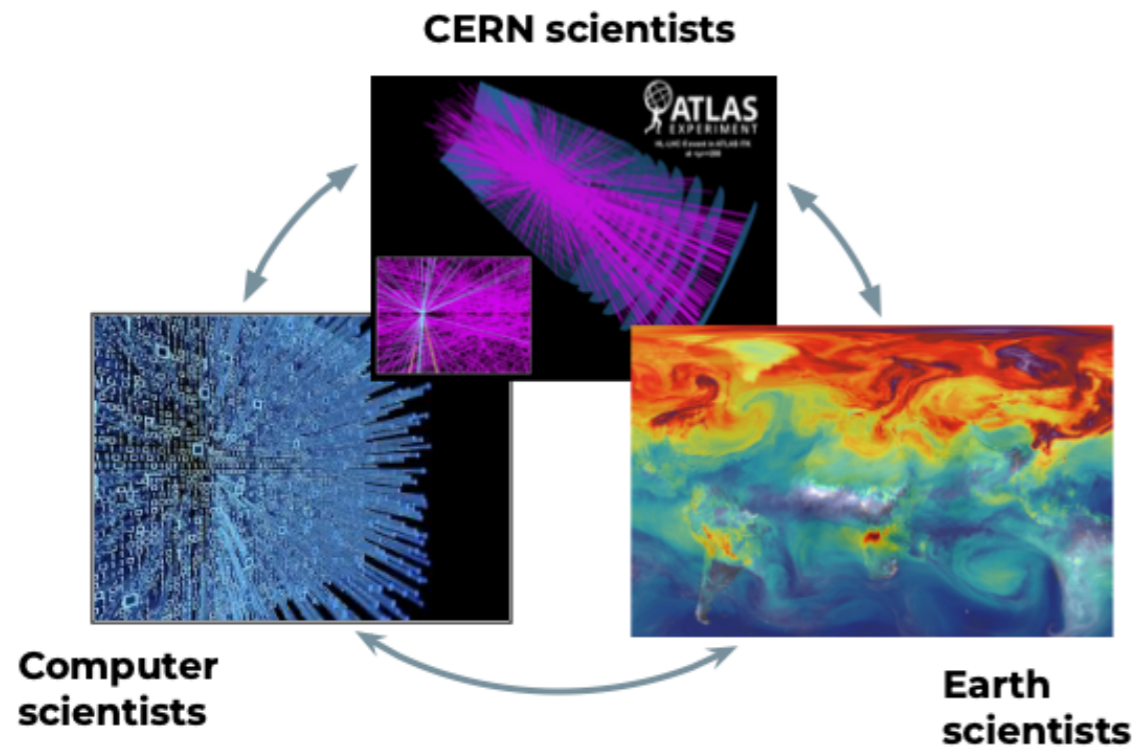


**feature space**

$$L_j : \mathbb{R}^n \to \mathbb{R}^m$$

Effective "encoding" of the data useful for many applications

**Data are clustered according to common properties in the "feature space"**

# AtmoRep: Introduction



**CERN scientists**

**Computer scientists**

**Earth scientists**

Solve common scientific challenge(s) in high-energy physics and weather/climate science using AI/ML

**Model complex, nonlinear phenomena and improve current simulations**

Access multi-scale dependencies of a given process

Earth science: eg. better understand convection phenomena

CERN: eg. particle-jet showers reconstruction

**Condense dataset information in a compact representation**

better handle the information in downstream applications.

eg. condense the info in a few GB rather than TB

**Explore potential of unsupervised learning for scientific applications**

Extract new information directly from data

eg. learn unknown correlation patterns

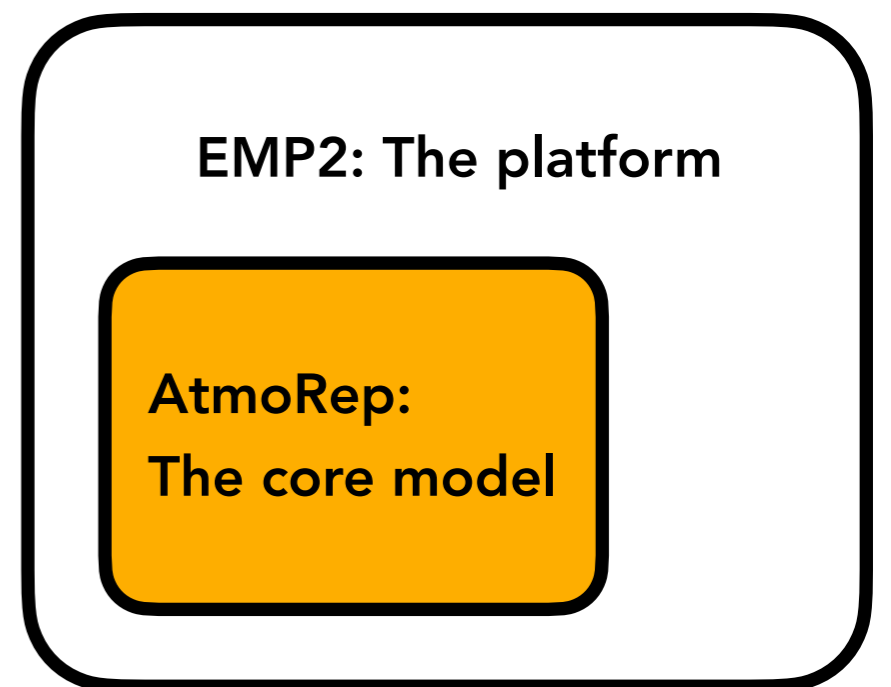Earth science: eg. early detection of extreme events

CERN: eg. anomaly detection

**Common Goal:**

**Develop a proof of concept of representation learning for scientific applications based on observations**

Ilaria Luise, CERN - ilaria.luise@cern.ch

# AtmoRep:
# A foundation model for the atmosphere

*The core model*
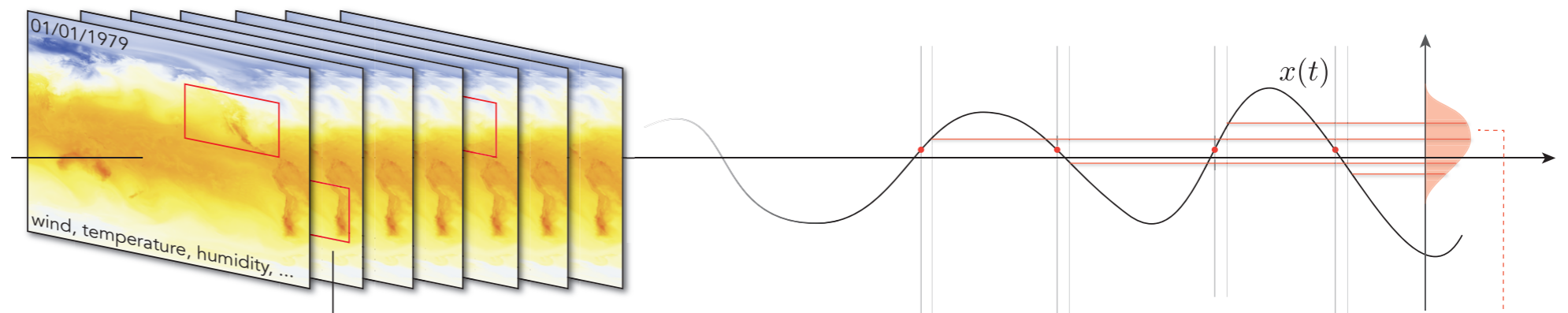
EMP2: The platform

AtmoRep:

The core model

# Key ingredient: What is a foundation model for us?

**The spatio-temporal (4D) evolution of a dynamical system can be summarised as**

*Probability of getting the state y given the initial state x and the auxiliary info α*  →  $p(y \mid x, \alpha)$  ←  *Auxiliary info: position, absolute time etc..*

$x(t)$



01/01/1979

wind, temperature, humidity, ...

$x(t)$

Training

The distribution can be approximated by a large neural network
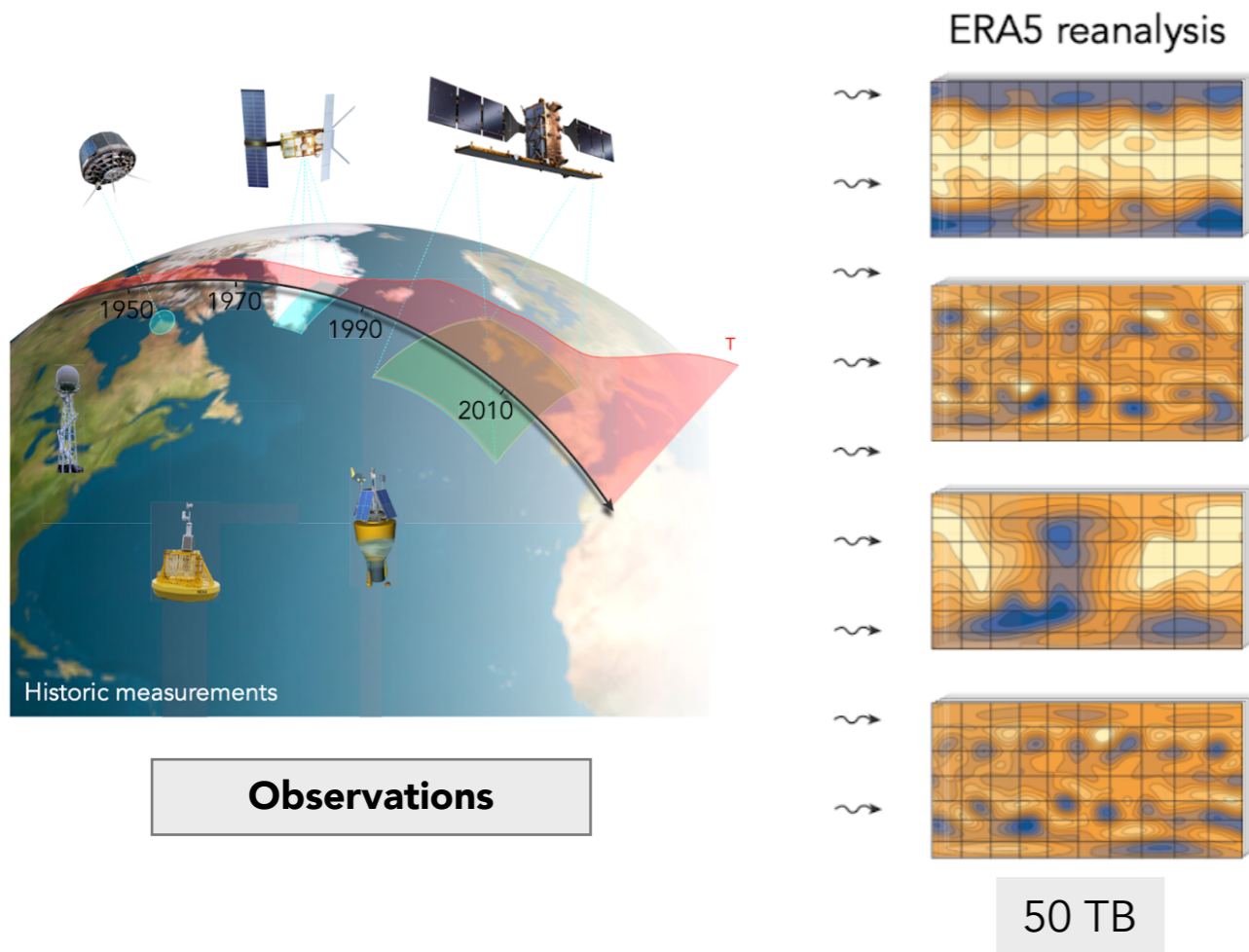
$$p(y \mid x, \alpha) \approx p_\theta(y \mid x, \alpha)$$

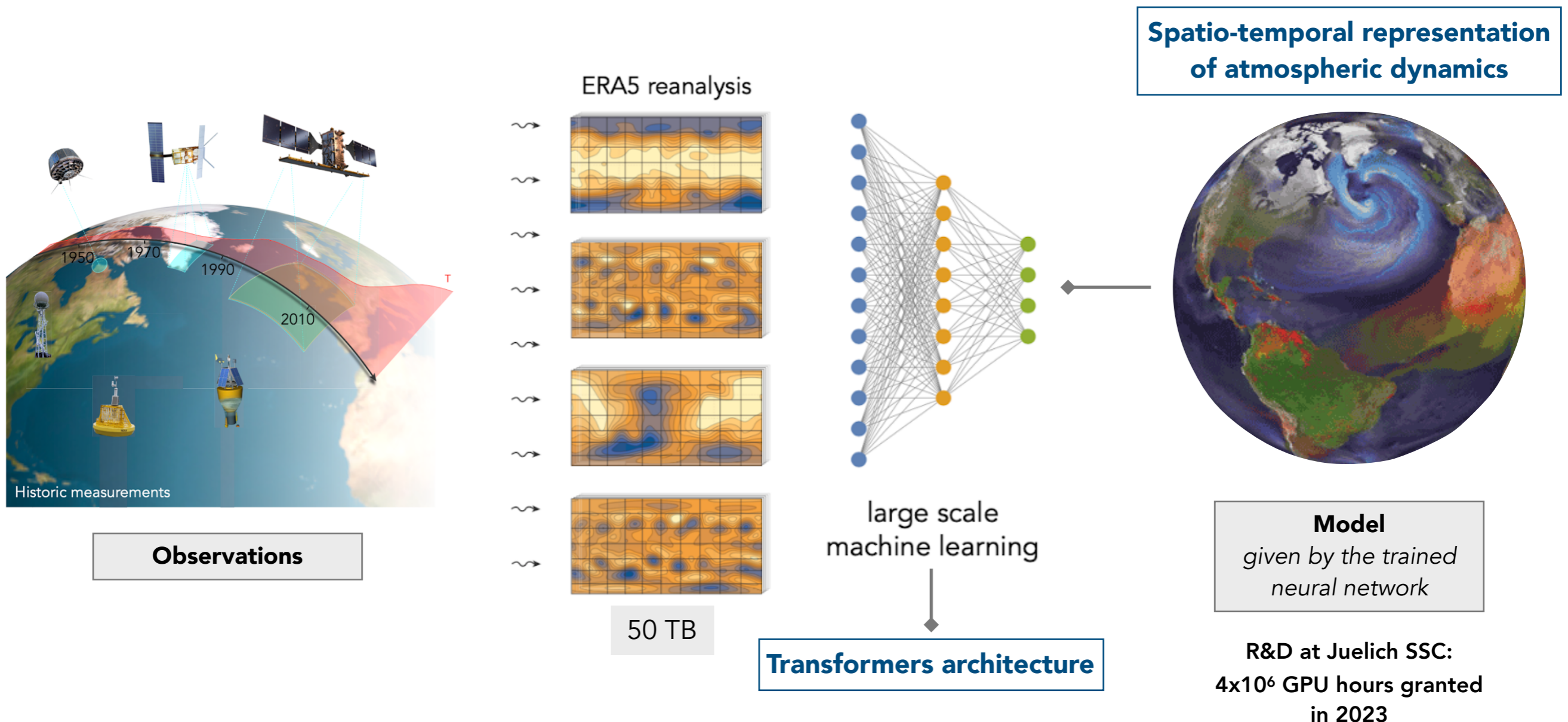$p_\theta(y \mid x, \alpha)$  →  **foundation model:** neural network that models data distribution for a specific domain

$t$

# The project in a nutshell

**A machine-learning based global environmental model trained on terabytes of observational data**
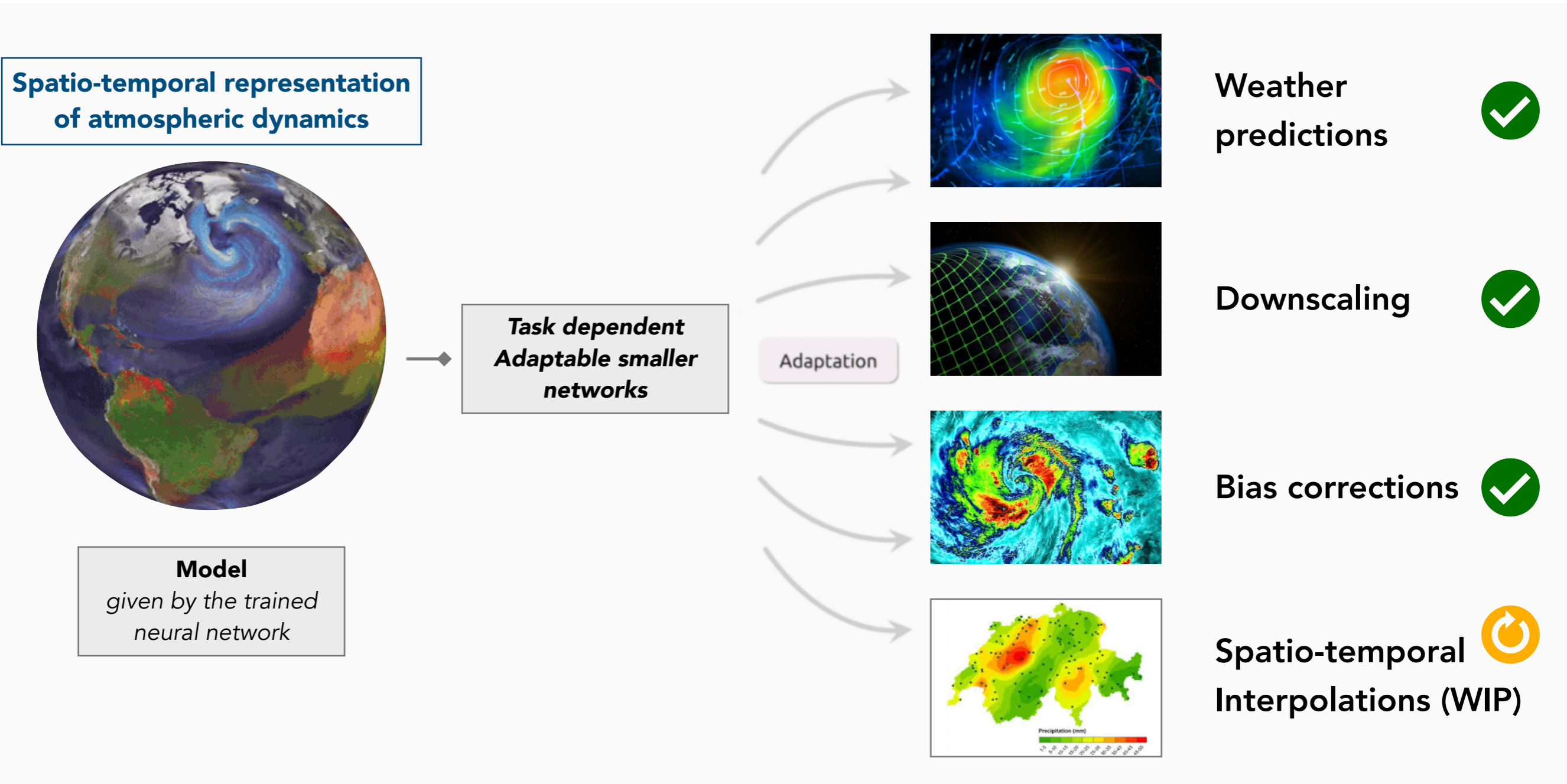
Ilaria Luise, CERN - ilaria.luise@cern.ch

# The project in a nutshell

**A machine-learning based global environmental model trained on terabytes of observational data**



Spatio-temporal representation of atmospheric dynamics

ERA5 reanalysis

Historic measurements

Observations

50 TB

large scale machine learning

Transformers architecture

**Model**
*given by the trained neural network*

R&D at Juelich SSC: $4 \times 10^6$ GPU hours granted in 2023

JÜLICH Forschungszentrum

8

Ilaria Luise, CERN - ilaria.luise@cern.ch

# Applications: one model for multiple purposes

**Use the learned representation to improve the state-of-the-art of specific weather & climate-related scientific applications**

# The dataset

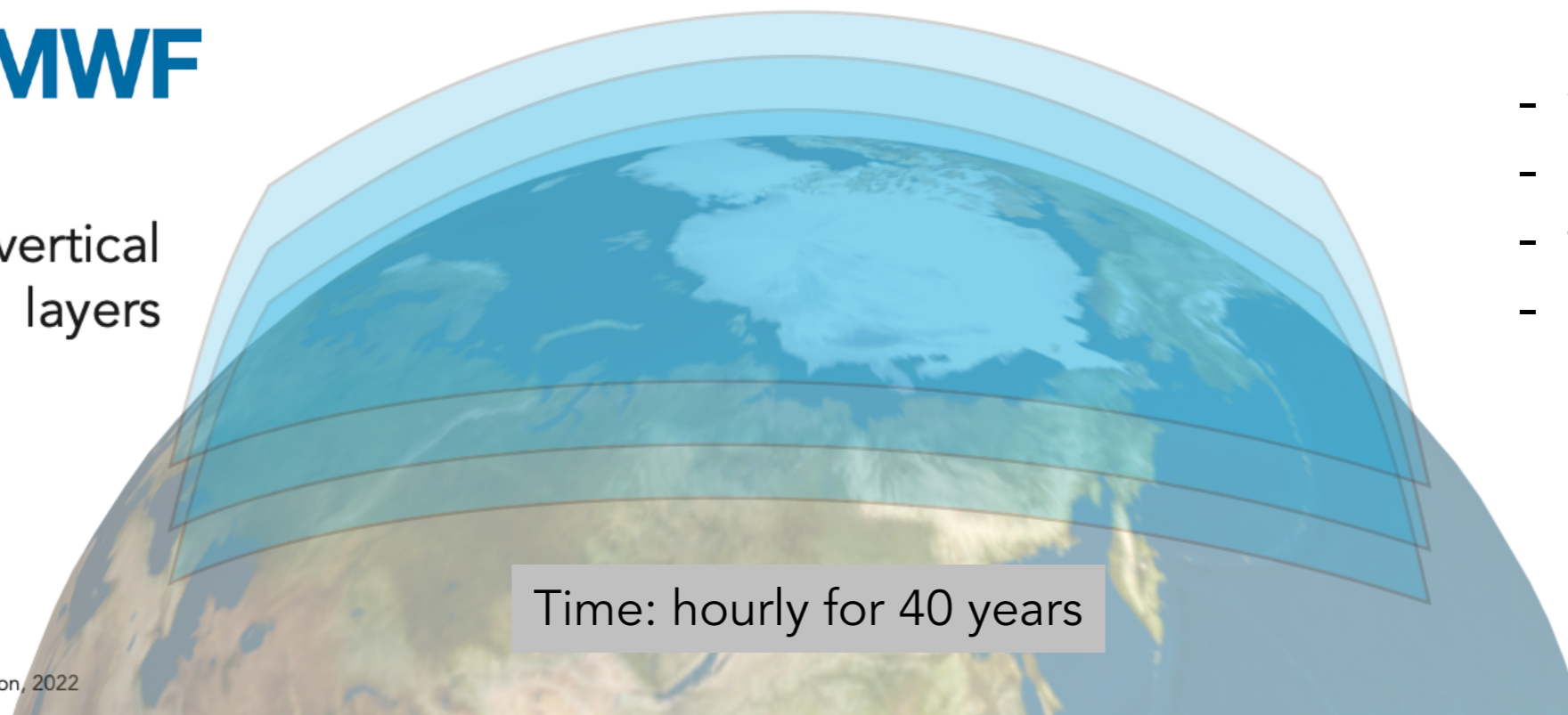**Publicly available pre-processed dataset of hourly spaced interpolated Earth observations: The <u>ERA5 reanalysis</u> from ECMWF**

**Subset of ERA5 reanalysis used at the moment for training:**

- Physical fields: vorticity, divergence (or wind velocity), vertical velocity, temperature, specific humidity, total precipitation
- Space: 721 x 1440 x 5 vertical layers
- Time: **randomly sample** over 24 time steps per day for 365 days for 40 years

721x1440 horizontal grid (0.25 degree)

137 vertical layers

- vorticity
- divergence
- temperature
- …

Time: hourly for 40 years

© Atmorep Collaboration, 2022

10

Ilaria Luise, CERN - ilaria.luise@cern.ch

# Intermezzo: The training protocol

**Use a variation of BERT masked language model from self supervised trainings in NLP**

Random sampling of neighbourhoods for training → stochastic gradient descent



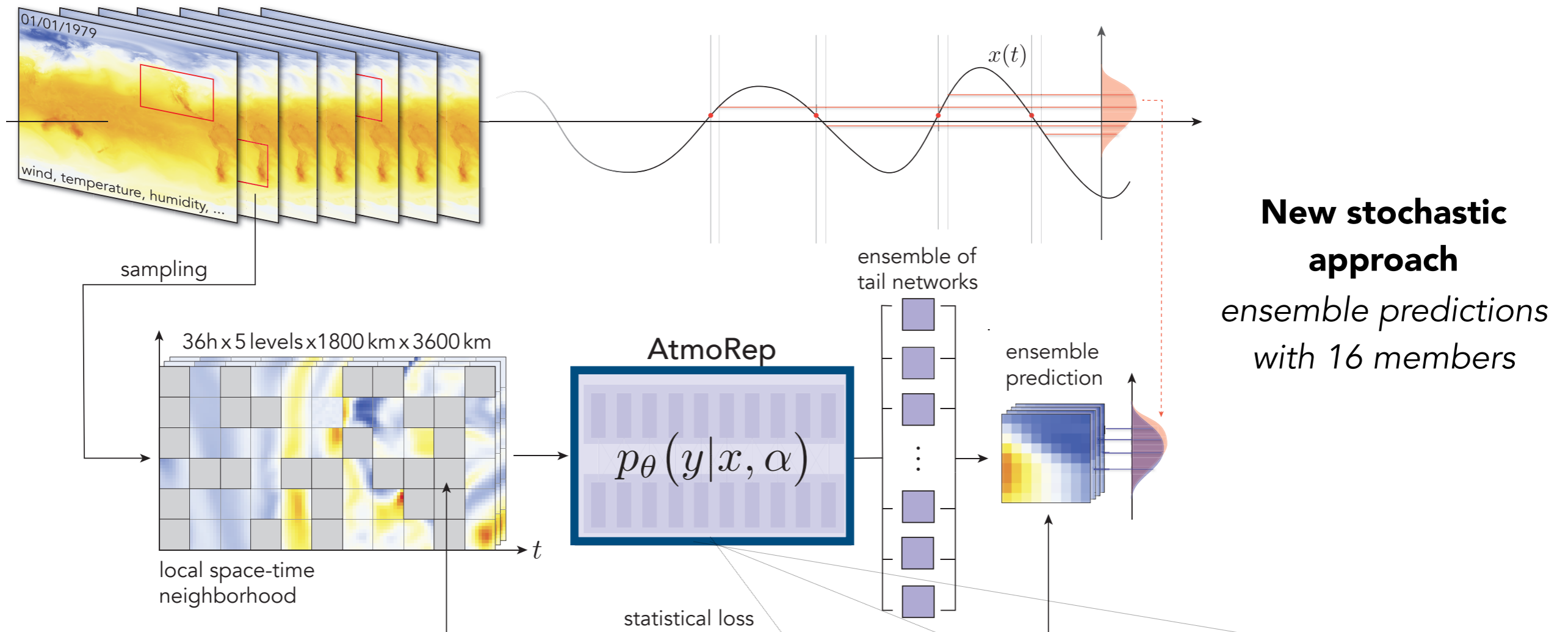**Split cube in small space-time regions (3D cubes) → tokens**
**Mask random tokens within the hyper-cube and try to predict them back**
visually: learn representation dynamics through interpolation

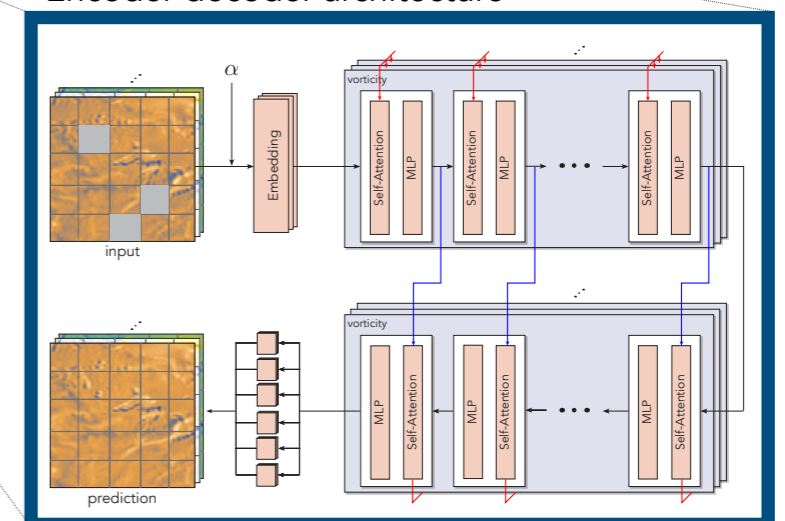*Default: 12 x 6 x 12 tokens with 3 x 9 x 9 grid points*

Ilaria Luise, CERN - ilaria.luise@cern.ch

# The network architecture

pre-processed historical observational record $x(t)$ (ERA5 reanalysis)



01/01/1979

wind, temperature, humidity, ...

$x(t)$

**New stochastic approach**
*ensemble predictions with 16 members*

sampling

36 h × 5 levels × 1800 km × 3600 km

AtmoRep

$$p_\theta\left(y\,|\,x,\alpha\right)$$

ensemble of tail networks

ensemble prediction

local space-time neighborhood

$t$

statistical loss

Encoder decoder architecture

**Approximate the 4-Dim PDF of the process using a Transformers-based network with 3.5 billion parameters**

Ilaria Luise, CERN - ilaria.luise@cern.ch

# Short term weather forecasting

*Zero-shot applications*

# Results: Target - ERA5

specific humidity, June 15th 2018 13:00 UTC

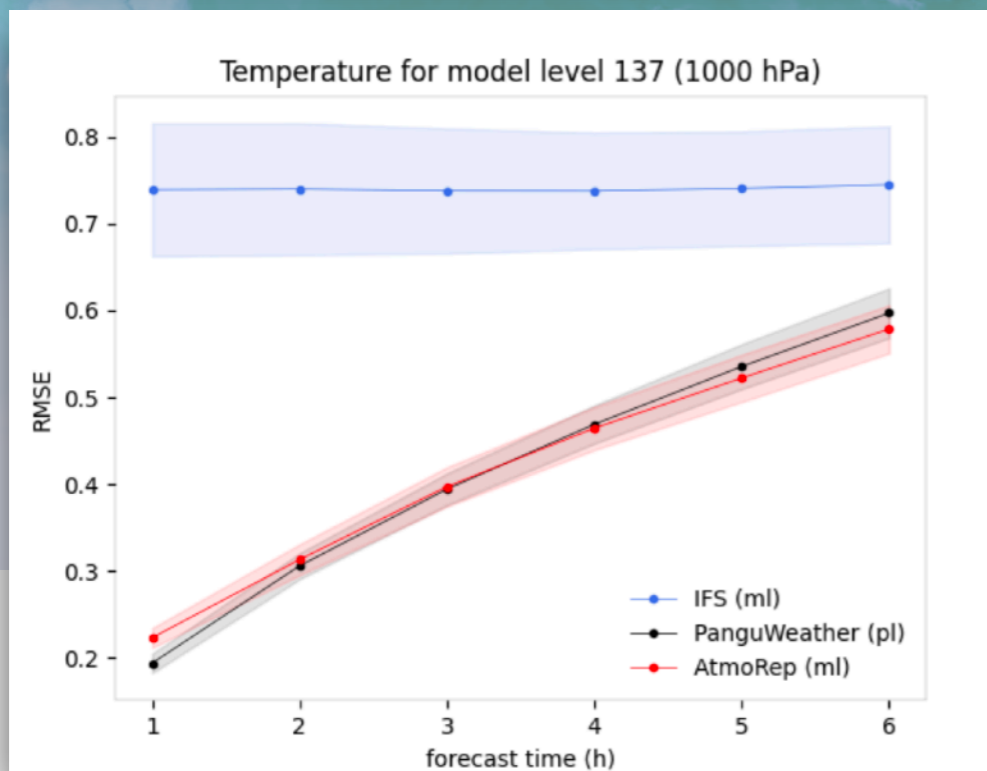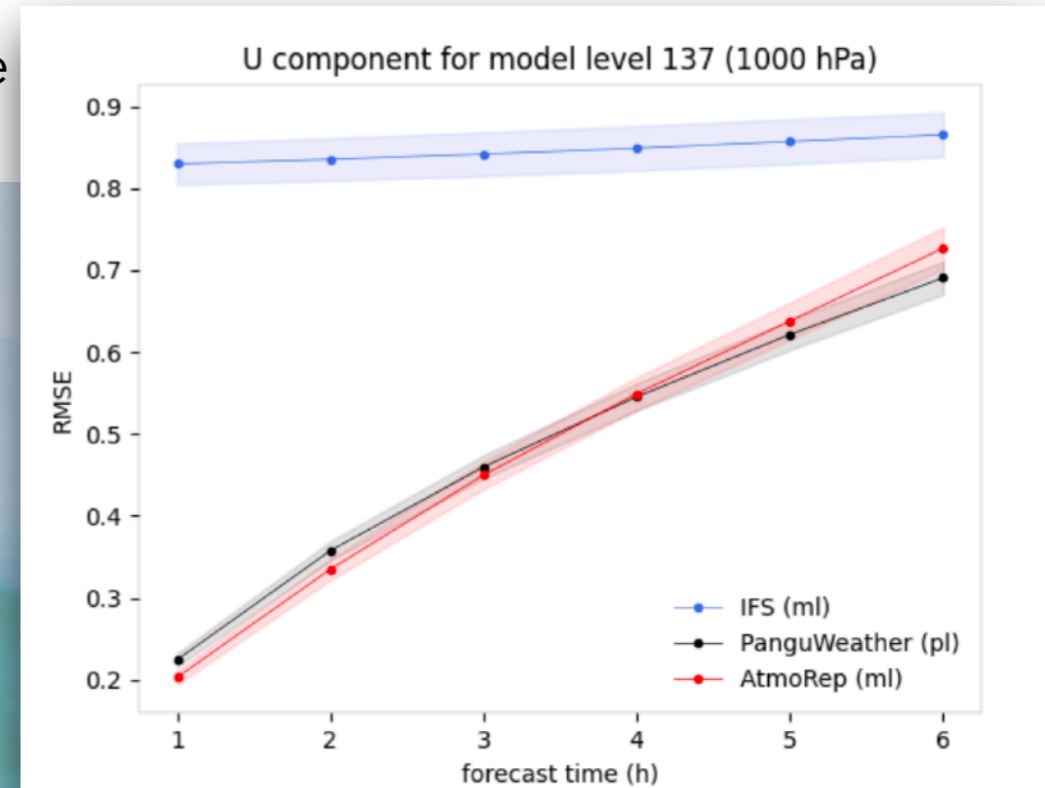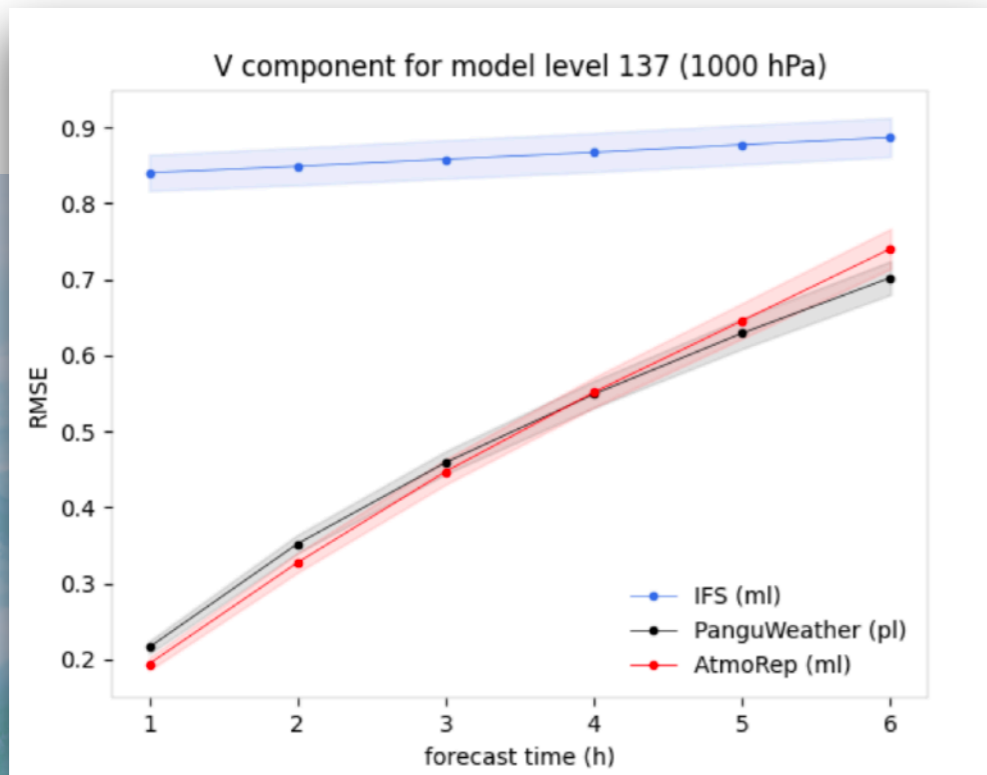# Results: Prediction - AtmoRep

specific humidity, June 15th 2018 13:00 UTC

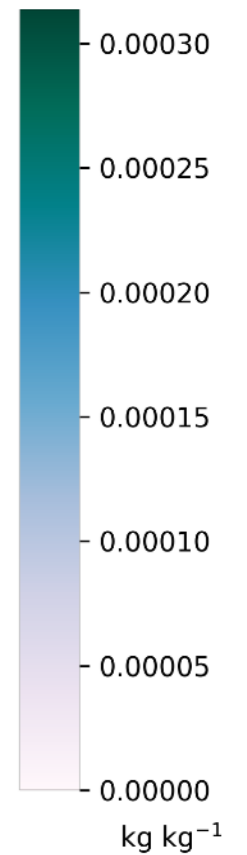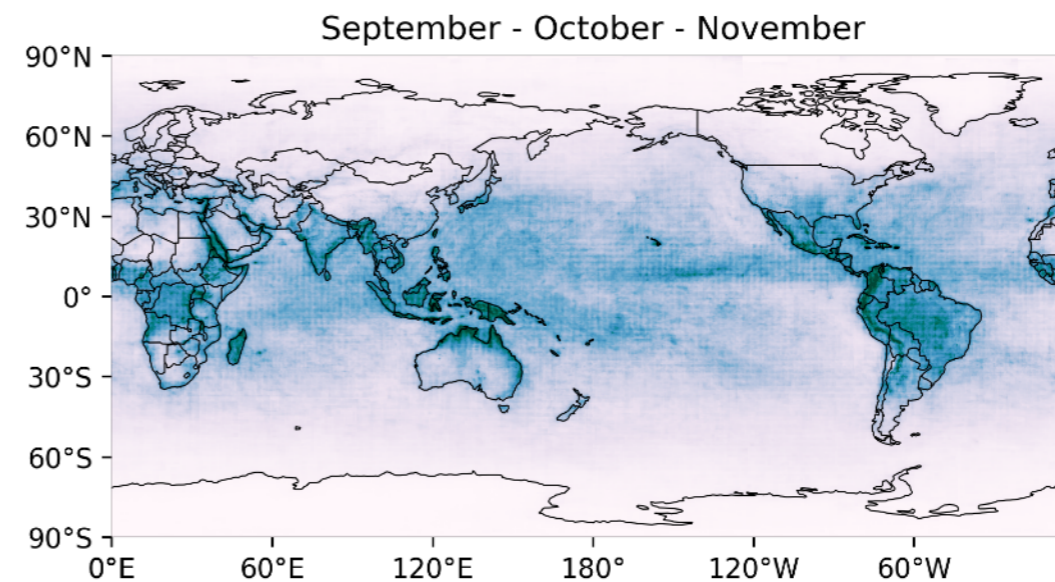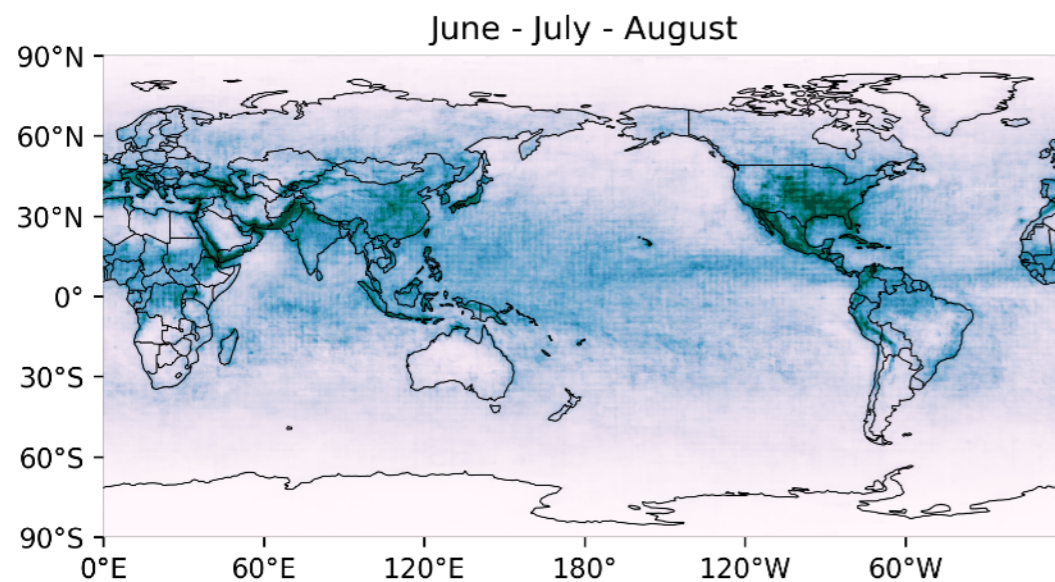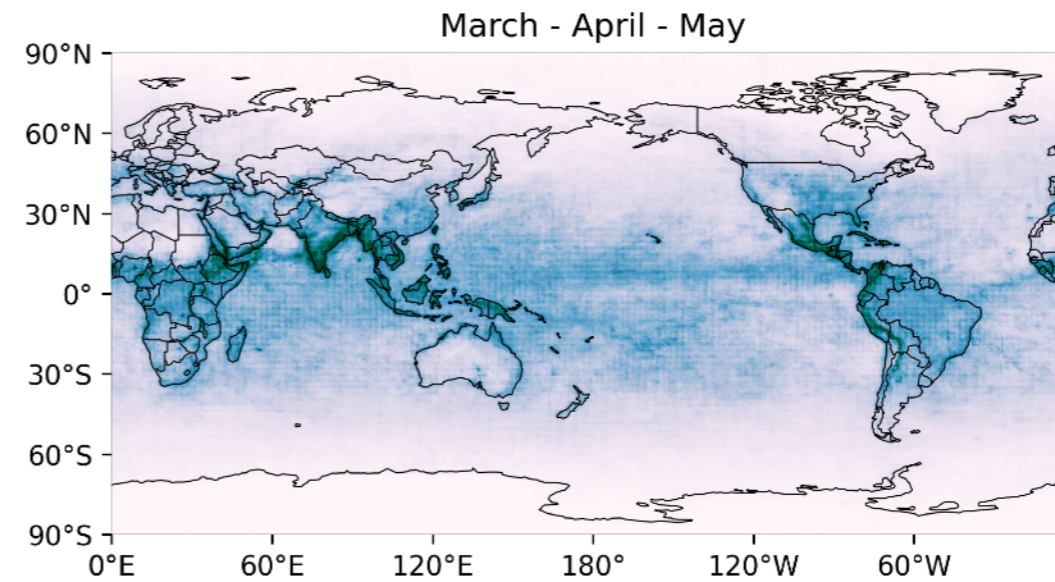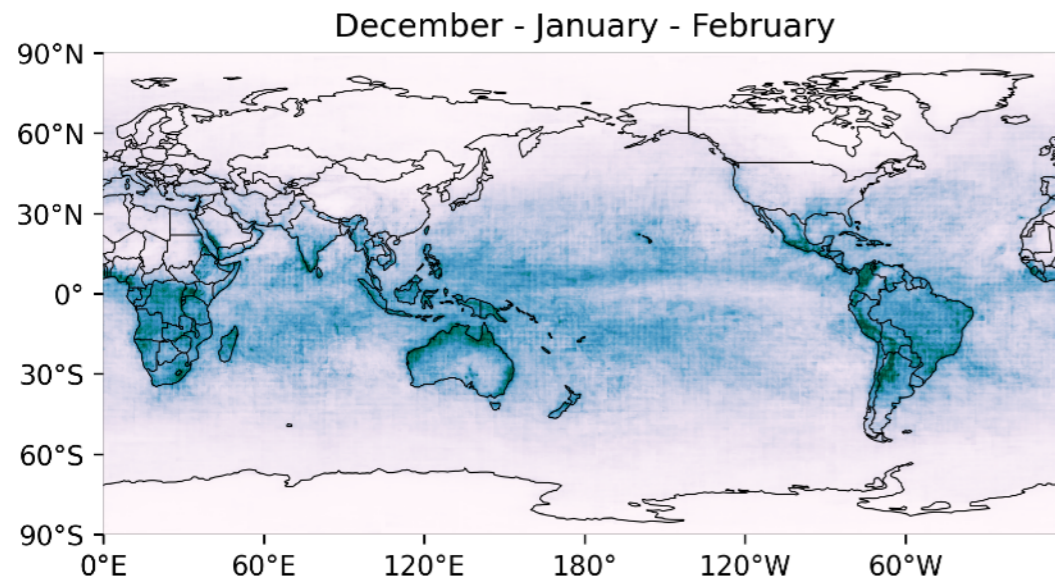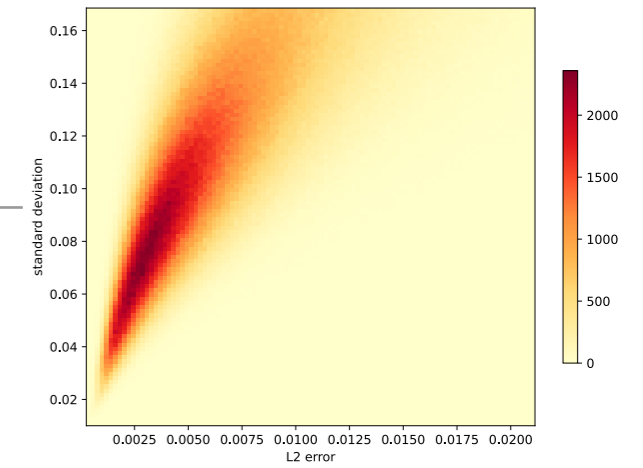# Results: Prediction - AtmoRep
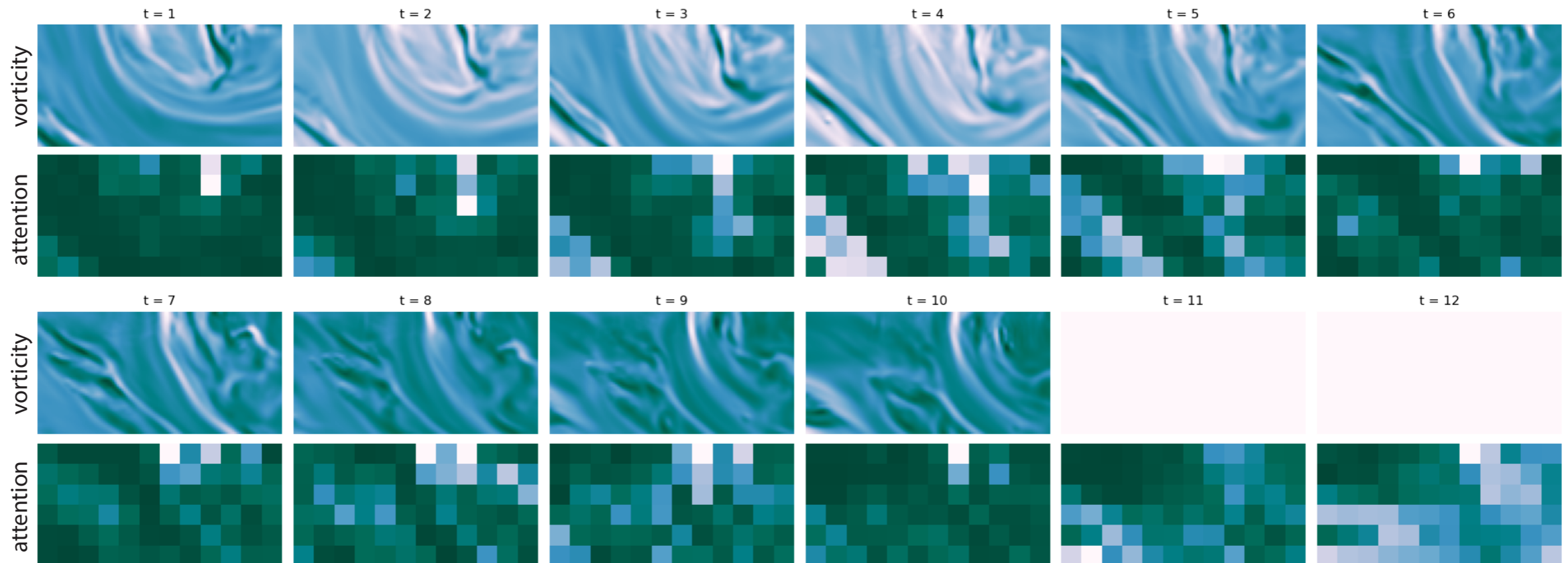
# Results: ensemble variability

**specific humidity: standard deviation of the ensemble**
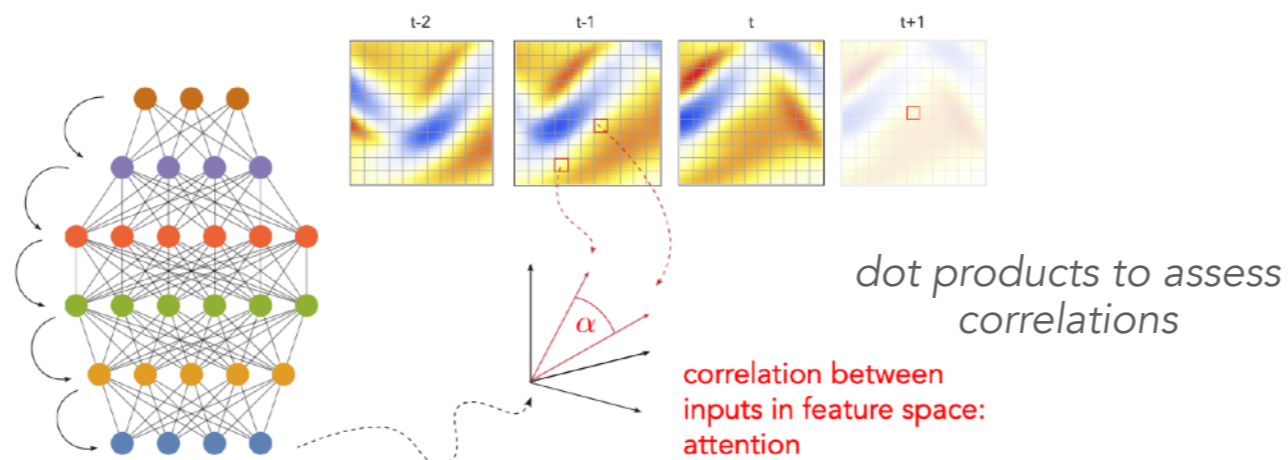
# Attention maps and interpretability

**Inspect the self-attention mechanism:**

**can we identify physics phenomena (e.g. hurricane formation) before they are even created?**



Attention:



*dot products to assess correlations*

correlation between inputs in feature space: attention

# Fine tuning on real data

*Data driven precipitation corrections & downscaling*

# Downscaling

ERA5

temperature, ml=137

COSMO-REA6

temperature, 2 m

Use COSMO REA6 data as **target** for the **loss minimisation**

loss function

COSMO ⟷ AtmoRep

target    prediction

# Downscaling

**Use the COSMO REA6 dataset (6 km resolution vs ~32 km in ERA5) to create a downscaled version of AtmoRep**



Comparison with a competing
AI-based model for downscaling:

# Bias corrections

**Precipitation rates are known to be suboptimal in ERA5**
**Use RADKLIM radar data to fine-tune the precipitation rates in AtmoRep**

Use *Radklim* data as *target*
  for the loss minimisation
  (just for total precip)



target
Radklim

prediction
AtmoRep

# Bias corrections: Results

**Precipitation rates are known to be suboptimal in ERA5**
**Use RADKLIM radar data to fine-tune the precipitation rates in AtmoRep**

Ilaria Luise, CERN - ilaria.luise@cern.ch

# Environmental modelling and prediction platform
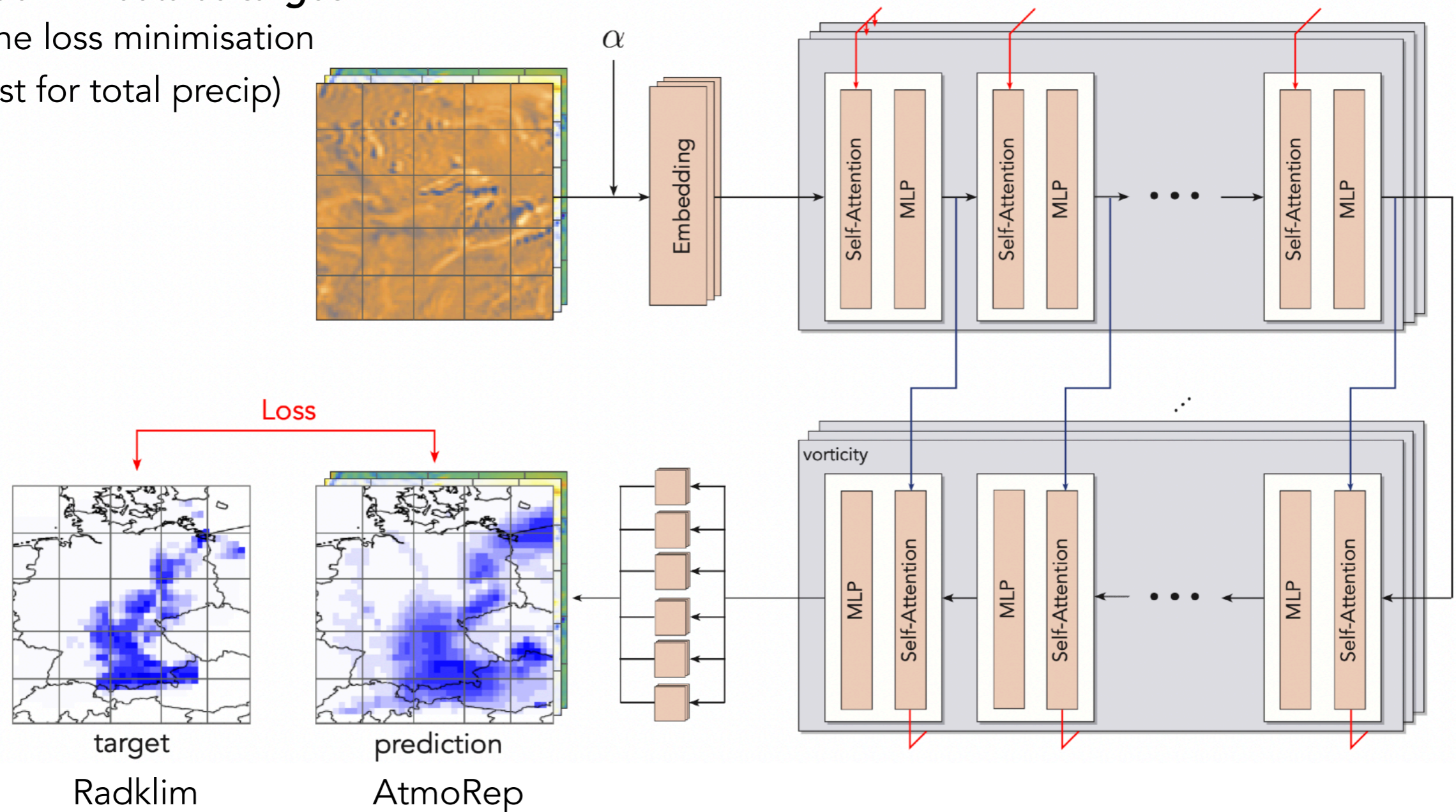
*The user-friendly platform*



EMP2: The platform

AtmoRep:

The core model

# Generalisation from HPC centres to clouds

**Future: develop the API & the user interface**

*Challenges: How foundation models interface with Digital Twin architectures?*
*Can we use foundation models as backbone for several digital twins*
*(e.g. core model = common atmospheric representation)?*

*Challenging part:*
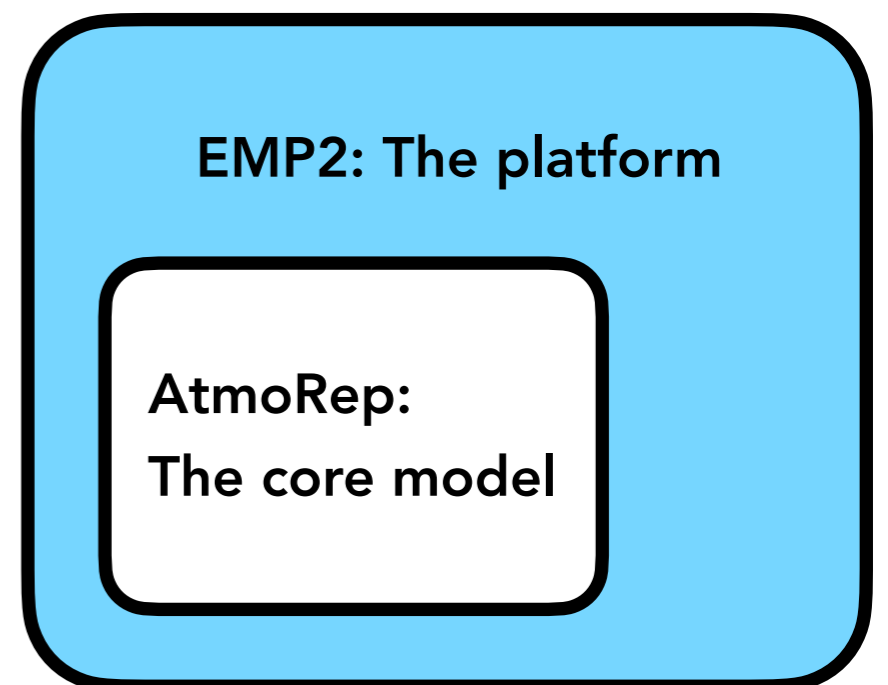*Close collaboration with the*
*members of InterTwin & CS4OD*
*projects at CERN*

*Final product:*
*Prototype of a user oriented platform*
*for environmental applications*



Onboard

(New CS4OD)

interTwin

**Goal: test EMP² within a digital twin existing architecture.**

*EMP² will be implemented as one of the use cases to test the Digital Twin architecture developed*
*through the InterTwin Project*

# Conclusions

**AtmoRep: First prototype of a multi-purpose model for Earth system applications**

The **model is available and testable** on the current applications:
nowcasting, downscaling, temporal interpolation and precipitation corrections.

**More infos:**
- Code is available on GitHub: link
- More infos on the website: www.atmorep.org
- **Pre-print on ArXiv: link**

**.. and some long term plans:**
- How to integrate "raw" observations?
- Coupled atmosphere+ocean system?

1950  1970  1990  2010

Historic measurements

Ilaria Luise, CERN - ilaria.luise@cern.ch

**Backup**

Technology

Operations

Collect

Aware

Compute

Respond

Visualize

Predict

Artificial Intelligence

Big Data Analytics

Visualization Tools

# Key ingredients: statistical loss

**? x N**

**Inspired by
cross entropy loss in
classification problems**

*Vorticity*

**x N**

tail network

loss

observation

*statistical interpretation:
measure the difference between the
pdf of the ML classification model
and the predicted distribution*

**Ensemble of tail networks
generates N predictions for
each pixel**

Ensemble
Prediction

**Interpret as probability
distribution for each pixel**
**(assumed Gaussian)**

**Loss: Minimize the difference between the mean of the distribution and the true value**

# Ensemble variability



$p(y|x;\alpha)$

$x(t)$

# A full set of possible applications

Many application areas for machine learning across ECMWF

| | | | |
|---|---|---|---|
| | Hydrology models and identification of human influence | Gravity wave drag emulation | ecPoint precipitation post-processing |
| Automated quality control of analysis | Clustering and categorisation for predictability and model errors | Hybrid ML/conventional land-surface models | Pollution anomaly detection and pollution downscaling | Post-processing for extreme predictions in Copernicus |
| Anomaly detection | Radiation emulation | ML preconditioning | Learn forecasts from observations | Feature detection |

Observations → Data assimilation → Numerical weather forecasts → Post-processing and dissemination

High-performance and (big) data processing infrastructure

| | | | |
|---|---|---|---|
| ML preconditioning for 4D-Var | Bias correction and model learning | Optimise data access | CliMetLab | Precipitation downscaling |
| Use land surface observations and data assimilation | Bias correction and gap filling for aerosol observations | Anomaly detection and workflow for IT infrastructure | | Fuse observations and predictions for sea ice products |
| Mapping of non-Gaussian to Gaussian distributions for data assimilation | | Optimise chillers for HPC | | Ensemble post-processing |

# EMP² vs ClimaX: differences & similarities

**ClimaX:**
## A foundation model for weather and climate

Tung Nguyen[*1,3], Johannes Brandstetter[2], Ashish Kapoor[1],
Jayesh K. Gupta[†1,], and Aditya Grover[†1,3]
[1]Microsoft Autonomous Systems and Robotics Research, [2]Microsoft Research AI4Science, [3]UCLA

6 February 2023!

| ClimaX | EMP² | |
|---|---|---|
| **Both are foundation models based on visual transformers!** | | |
| **Investigating similar downstream applications** | | |
| Trained on a randomised forecasting objective<br><sub>Goal: reconstruct states in the future</sub> | BERT-style training adapted to scientific data<br><sub>reconstruct *masked tokens* within a random hypercube</sub> | The model is less "forecasting oriented" |
| using ERA5 on pressure level variables | using ERA5 on model level variables | This is what ECMWF uses: an eye into the integration within their systems. |
| **deterministic predictions** | **stochastic predictions**<br><sub>Model uncertainty quantification through newly defined statistical loss</sub> | |
| **single transformer**<br><sub>Concatenation of fields in the variable aggregation step</sub> | **stack of transformers**<br><sub>one transformer for each field, coupled with cross attention.</sub> | Plug and play approach: new fields can be easily integrated. |
| private company | public research | |

Ilaria Luise, CERN - ilaria.luise@cern.ch