

AIFS

ECMWF's data driven forecast model

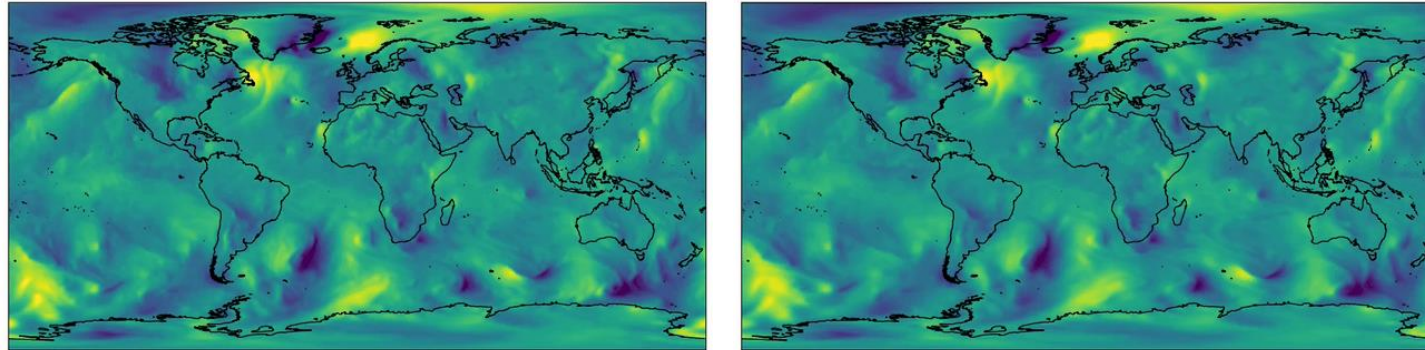
Simon Lang

Simon.lang@ecmwf.int

Mihai Alexe, Matthew Chantry, Jesper Dramsch, Florian Pinault,
Baudouin Raoult, Zied Ben Bouallegue, Linus Magnusson, Mariana
Clare, Peter Lean, Christian Lessig

AIFS

- > following Keisler 2022 and Lam et. al 2022
- GNN architecture: Interaction Networks (Battaglia et. al 2016)
- Graph representation, hidden multi-scale mesh, edge features

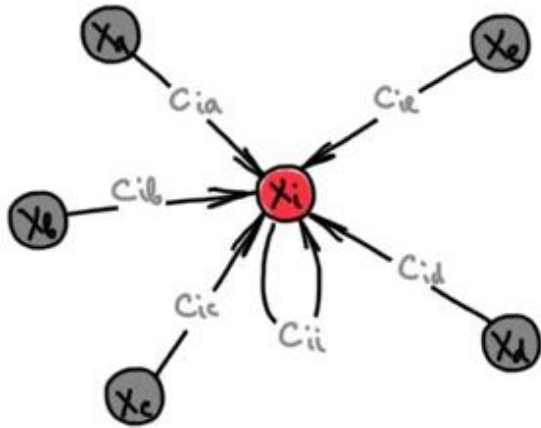


Why GNN : can handle arbitrary input grids, local grid refinement, changing grids etc. ; attractive for use in earth system science

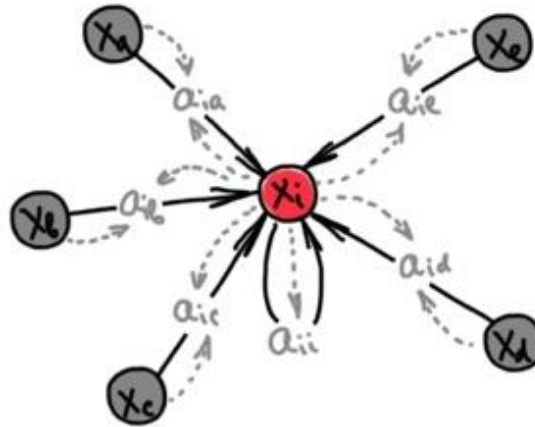
- built from flexible message-passing GNNs



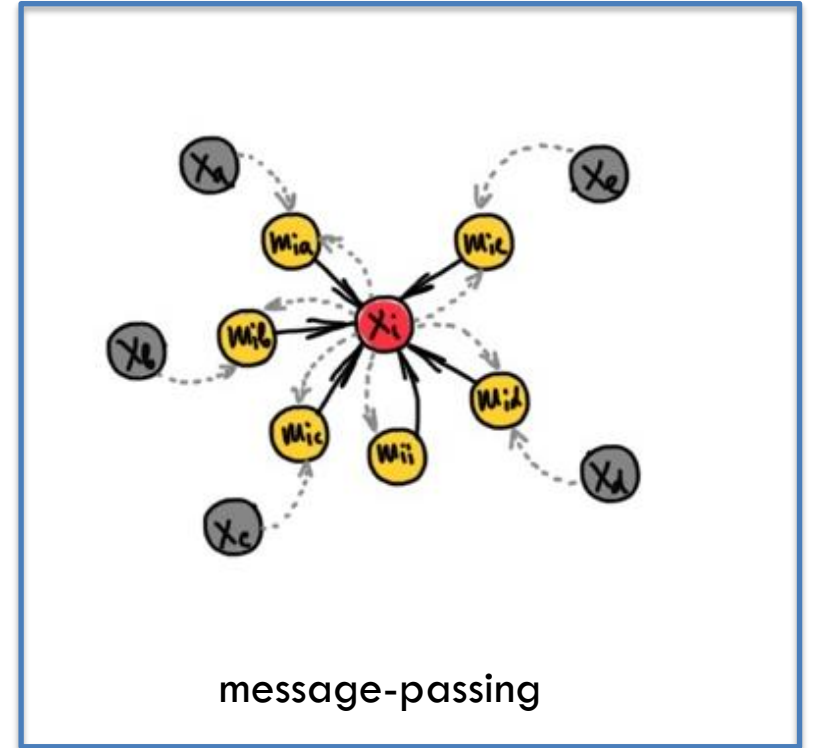
expressive power



convolutional



attentional

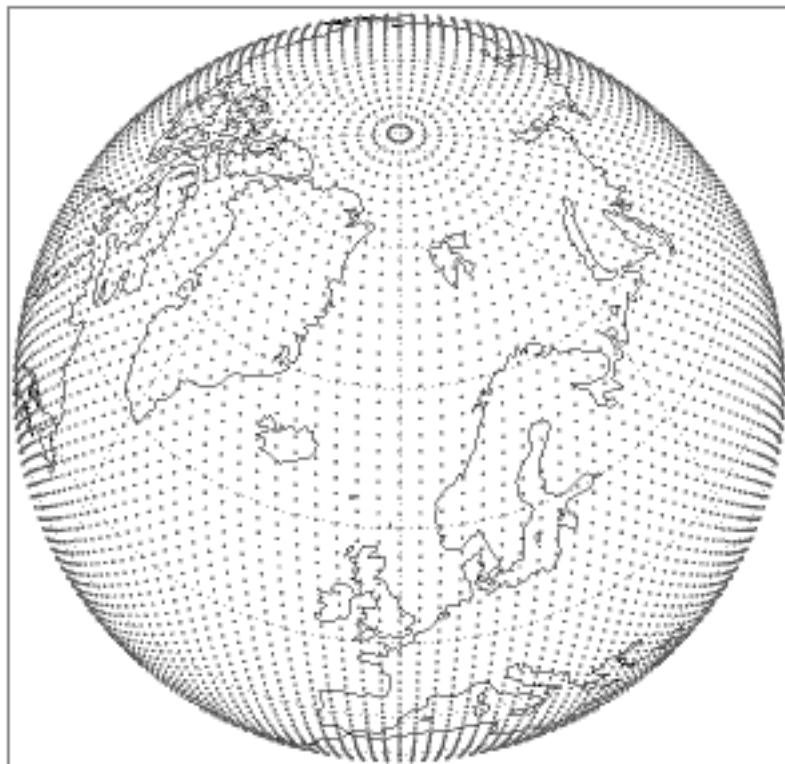


message-passing

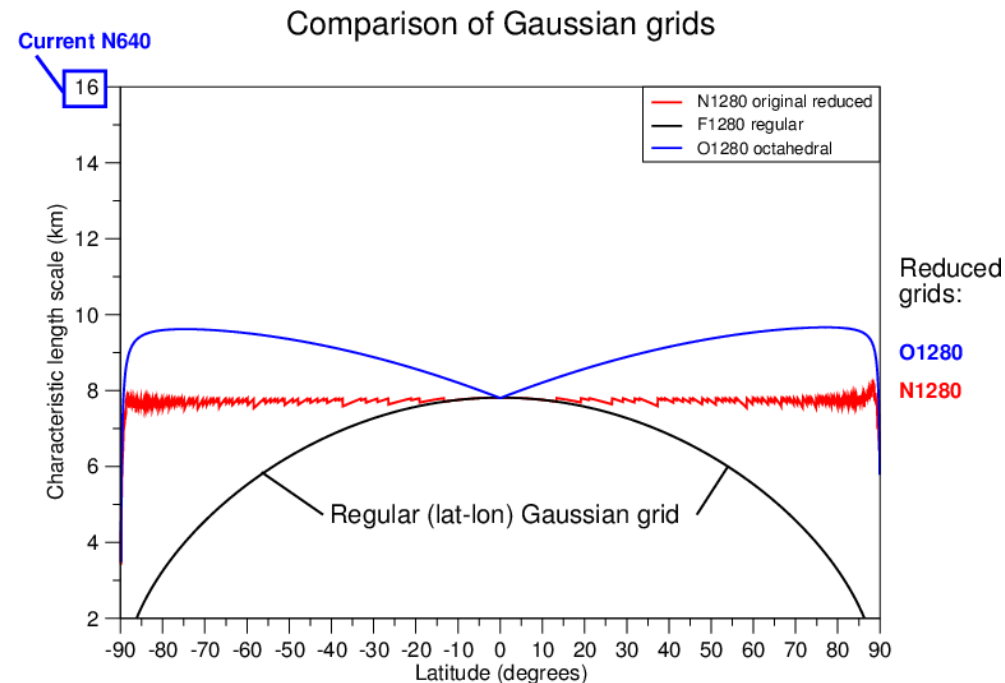
<https://arxiv.org/abs/2301.08210>

$$\mathbf{h}_u = \phi \left(\mathbf{x}_u, \bigoplus_{v \in \mathcal{N}_u} \psi(\mathbf{x}_u, \mathbf{x}_v) \right)$$

AIFS grid



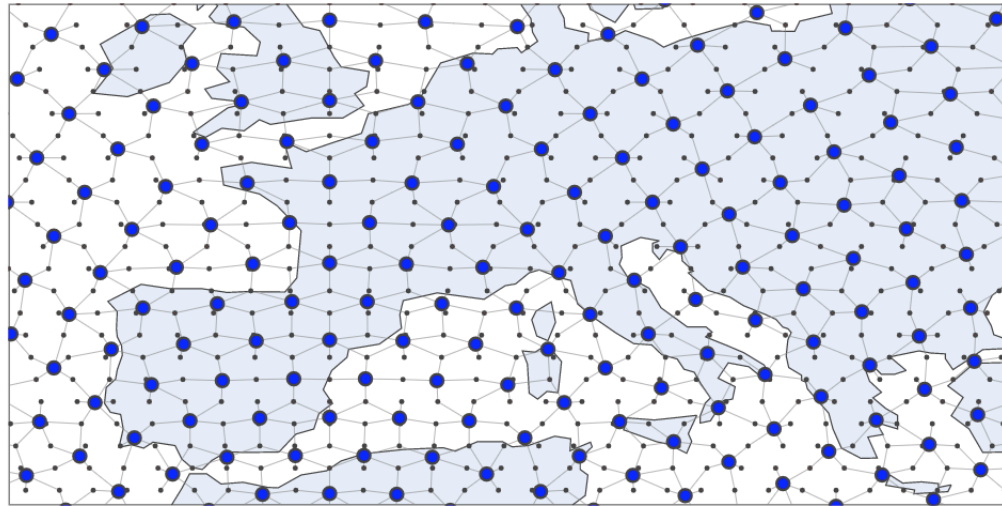
o96 / o160 / n320



- + (equi-)area weights
- + weighting along plevs (vertical)
- + per-variable weights in the loss

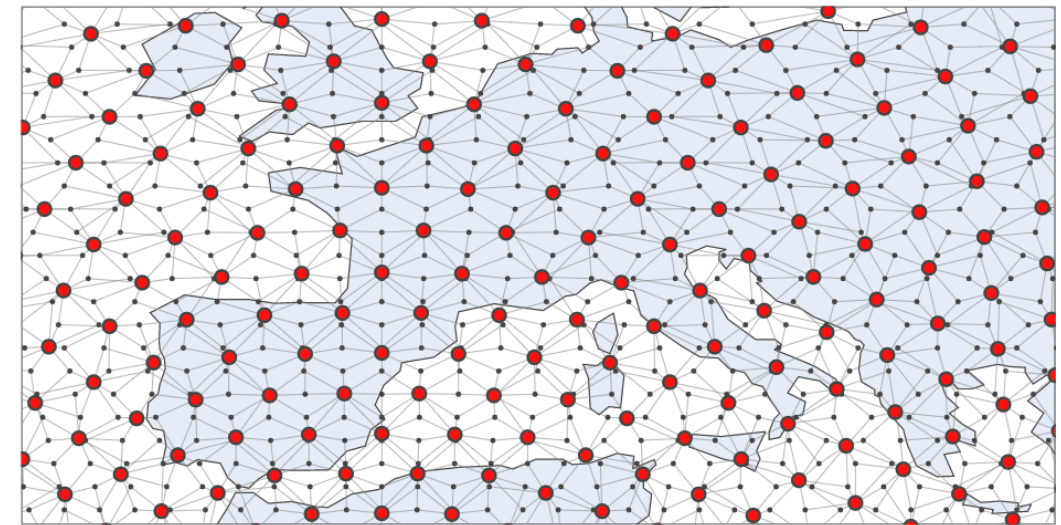
<https://confluence.ecmwf.int/display/FCST/Introducing+the+octahedral+reduced+Gaussian+grid>

1 x Encoder



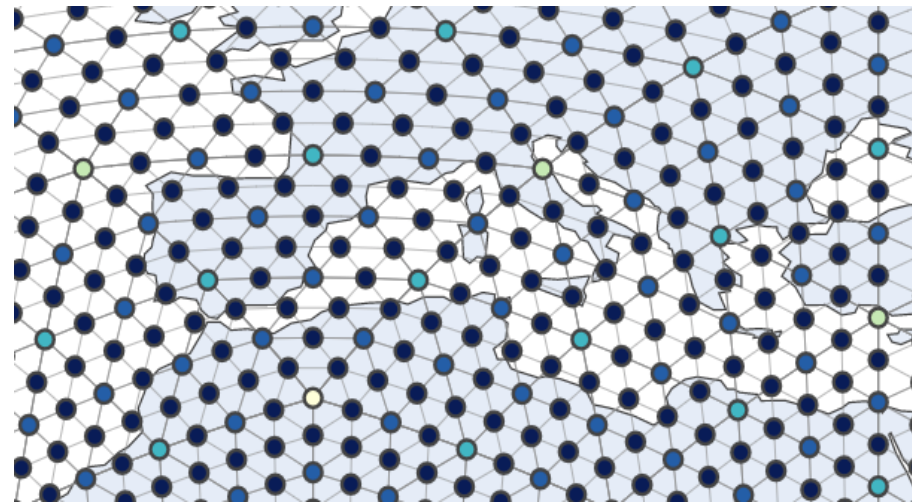
Skip-connection
(residual)

1 x Decoder



e.g.
Era5 n320 ~ 540 000 Nodes
Hidden mesh "h6" ~ 40 000 Nodes
~ 800 000 edges

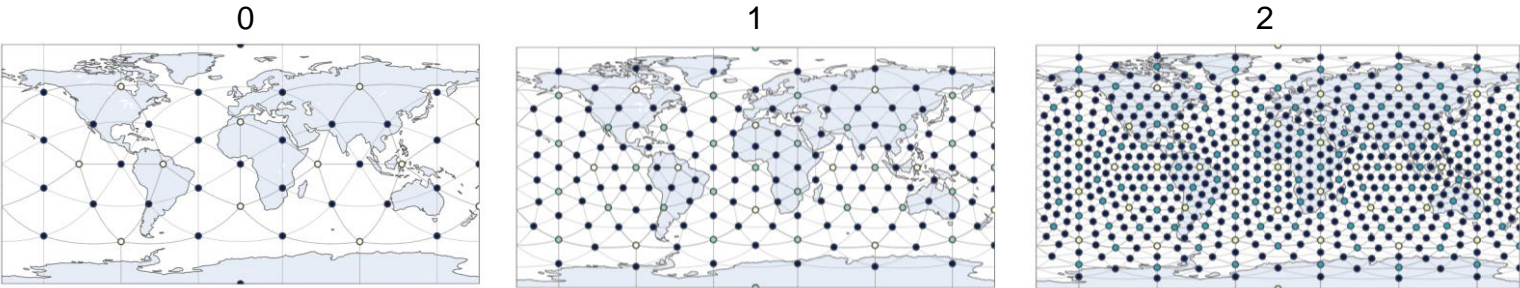
16 x Processor



e.g.
Hidden mesh "h6" ~ 40 000 Nodes
Era5 n320 ~ 540 000 Nodes
~ 1.6 million edges

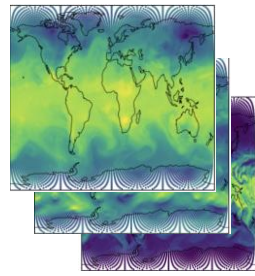
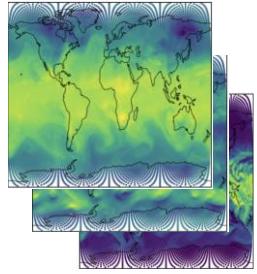
e.g.
Hidden mesh "h6" ~ 40 000 Nodes
Multi-scale interconnectivity
~ 320 000 edges

Simultaneous multi-level message passing



...

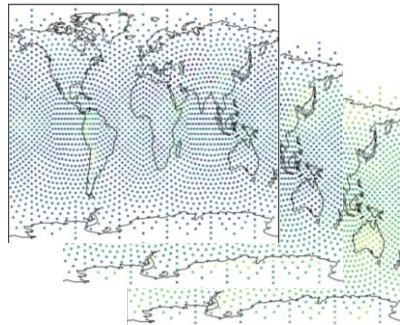
Atmospheric state:
 $X(t), X(t-6h)$



encoder

decoder

processor:
8-16 MP rounds



AIFS_{t -> t+6h}

previous
 $X(t)$

Prediction:
 $X(t+6h)$

$WMSE_{t+6h}$

AR predictions

AIFS_{t + 6h -> t+12h}

$WMSE_{t+12h}$

Aggregate
WMSE

AIFS current reference model

Model:

- O96 ERA5 grid, ~1-degree
- “Level 5” hidden grid, ~2-degree
- MLP dimension 768 (cf. 512 for GraphCast)
- 80 million parameters

Variables:

13 pressure levels – u, v, w, q, t, z
surface: 2t, 10u, 10v, 2d, sp, msl, sst

Training:

Step 1: 4 days on 16 GPUs to minimise errors for single 6h step

Step 2: 34 hours on 16 GPUs to minimise errors up to 3 days

Step 3: 4 hours on 16 GPUs minimising errors up to 3 days on operational analysis

Total ~6 days on 16 GPUs

Scaling up AIFS

*transformer = fully connected GNN

- Modest number of parameters, model size currently not much of an issue
- Data size is large, lots of grid points -> lots of nodes and edges (similar to very large sequence length in transformer based* models)

GPU memory is limited => 40 GB on ATOS, 64 GB on LUMI / Leonardo

For large (parameters) models there exist quasi out of the box libraries, this is not the case for big input data.

Memory saving options:

Do not keep everything in memory -> use re-computation in backward : large memory savings. We trade compute for memory as much as possible.

Multi GPU model : shard one model instance across multiple GPUs ; forward and backward pass different in terms of required communication

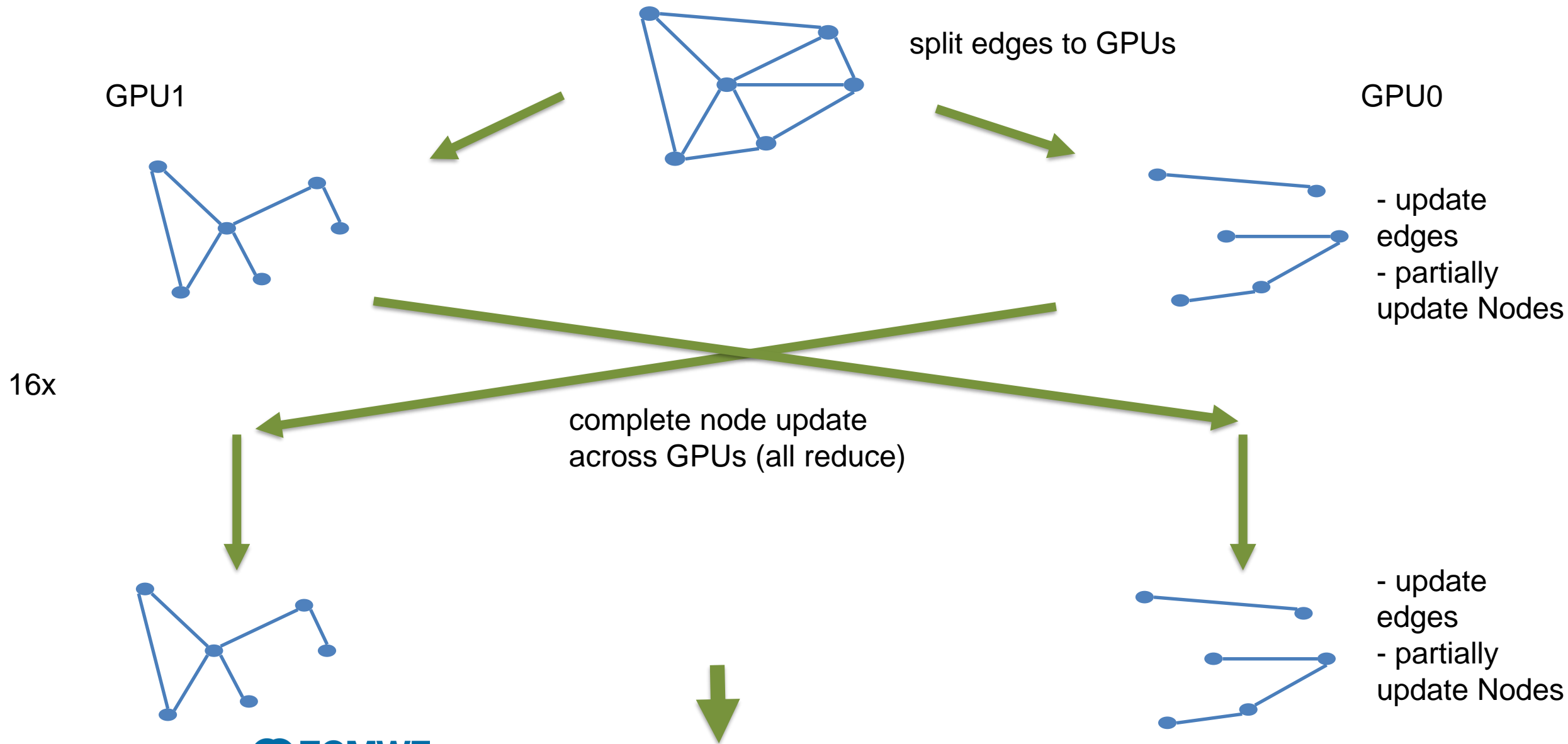
=> increases instantaneous available memory for more edges / nodes and also possible to split activations across GPUs



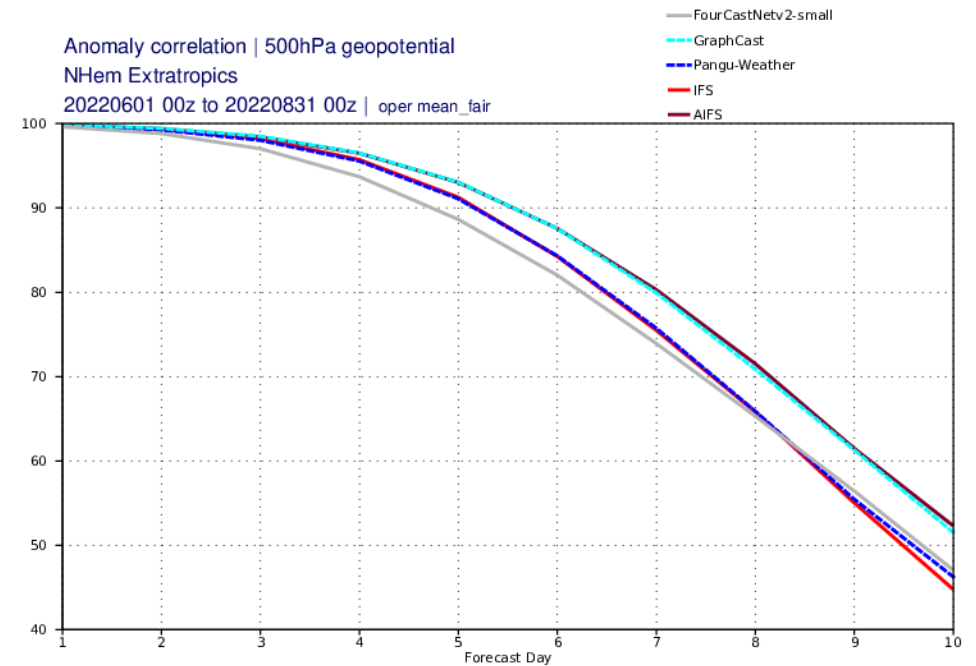
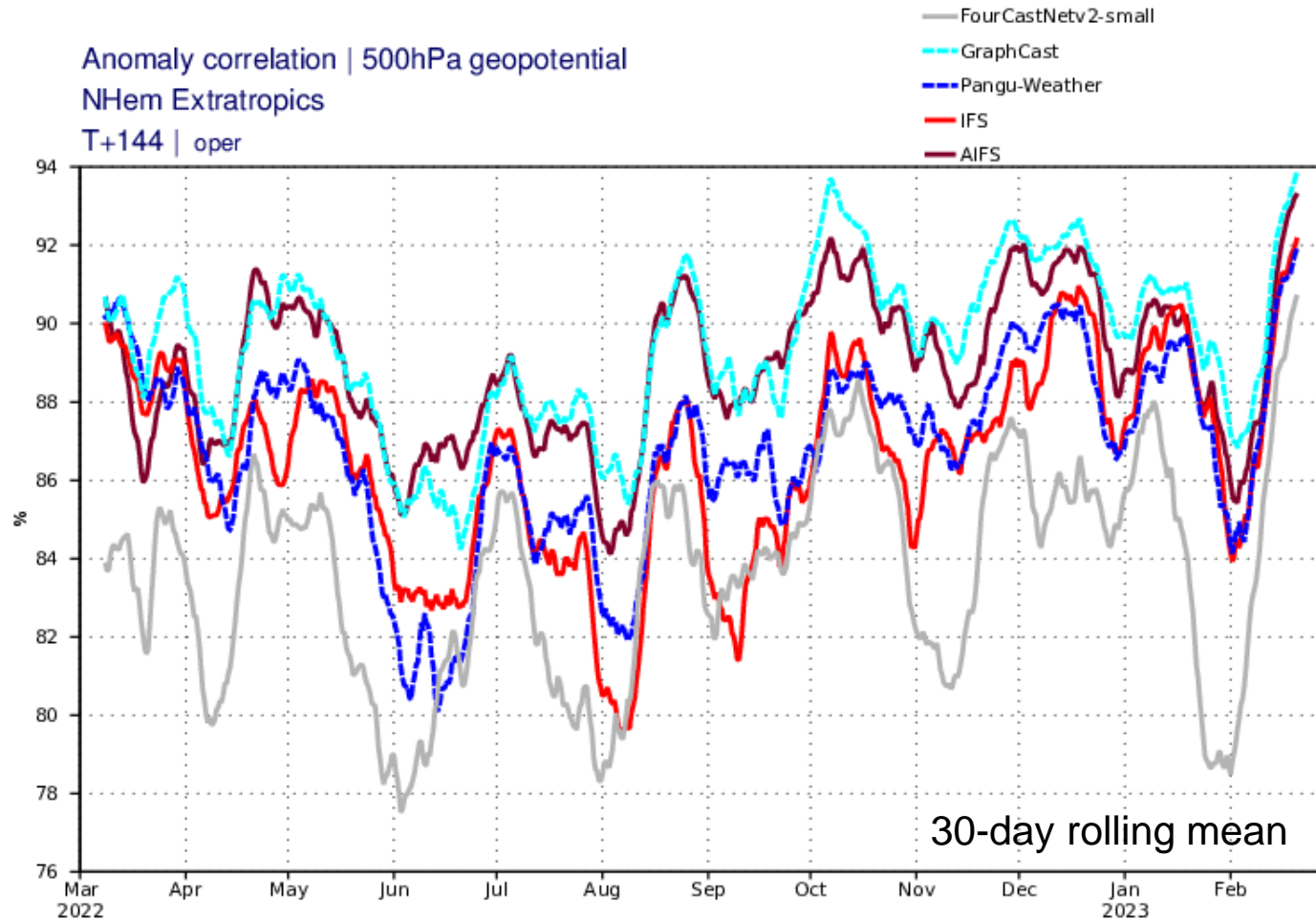
- Aggressive re-computing activations allows for rollouts of ~ 1. degree model on a single GPU
- Model parallel set-up makes it possible to go to native ERA5 resolution (and beyond?) and faster time to solution ; it will also allow to build larger models in the future

Model Parallel approach:

$$\text{edge_new} = f(x_0, x_1, \text{edge}) \quad x_new = f(x_0, \text{sum}(\text{edges_new}))$$



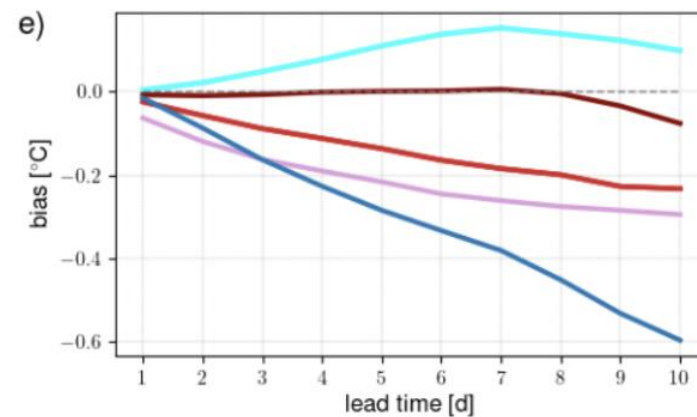
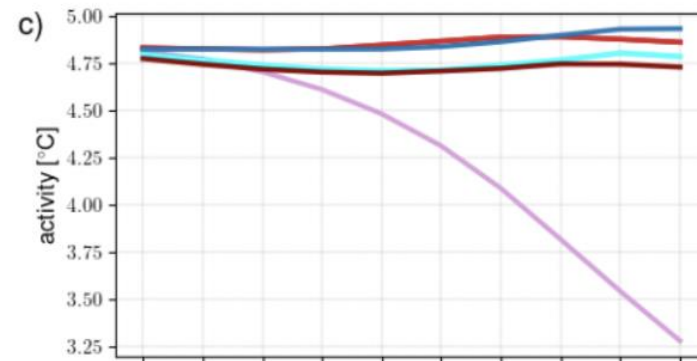
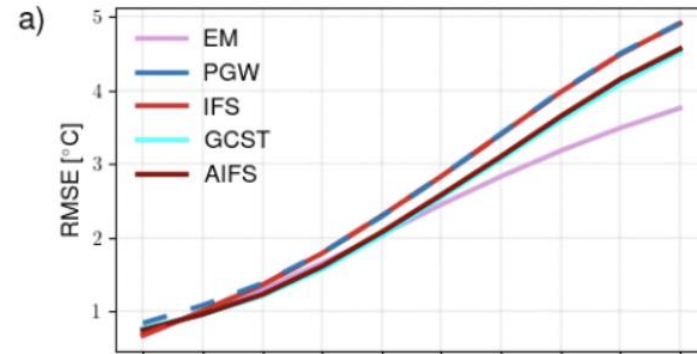
Headline score: anomaly correlation for Z500



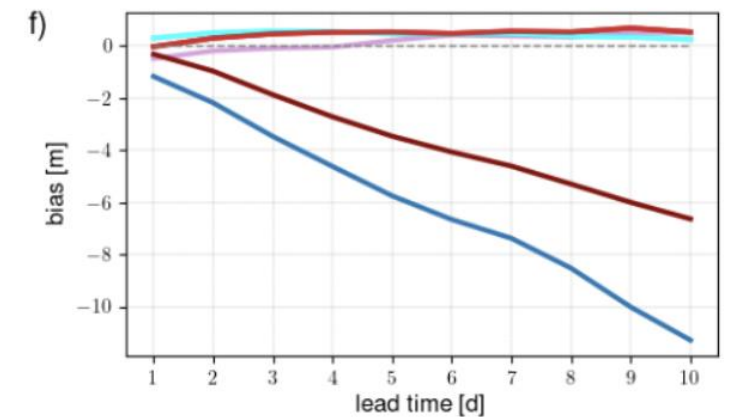
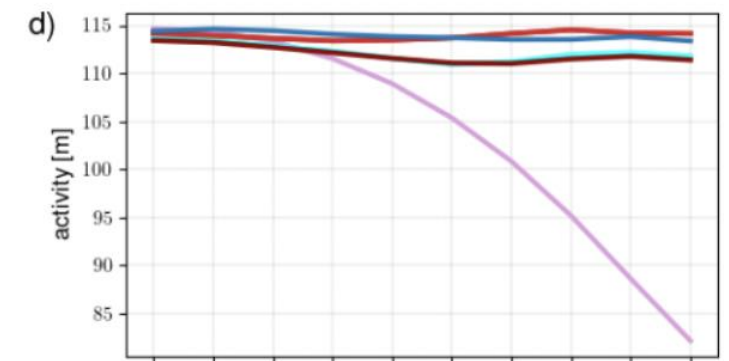
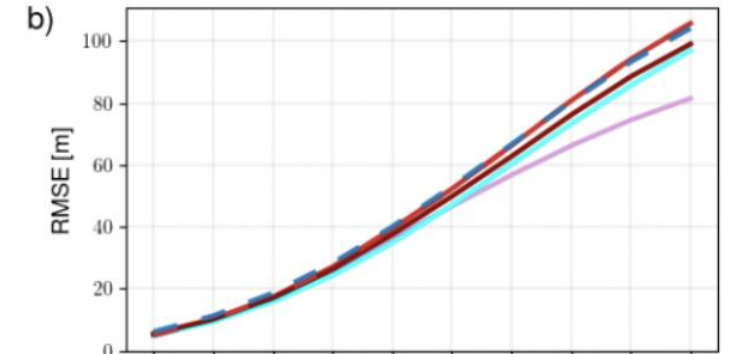
RMSE, activity, and bias

DJF 2022/2023
NHem Extratropics

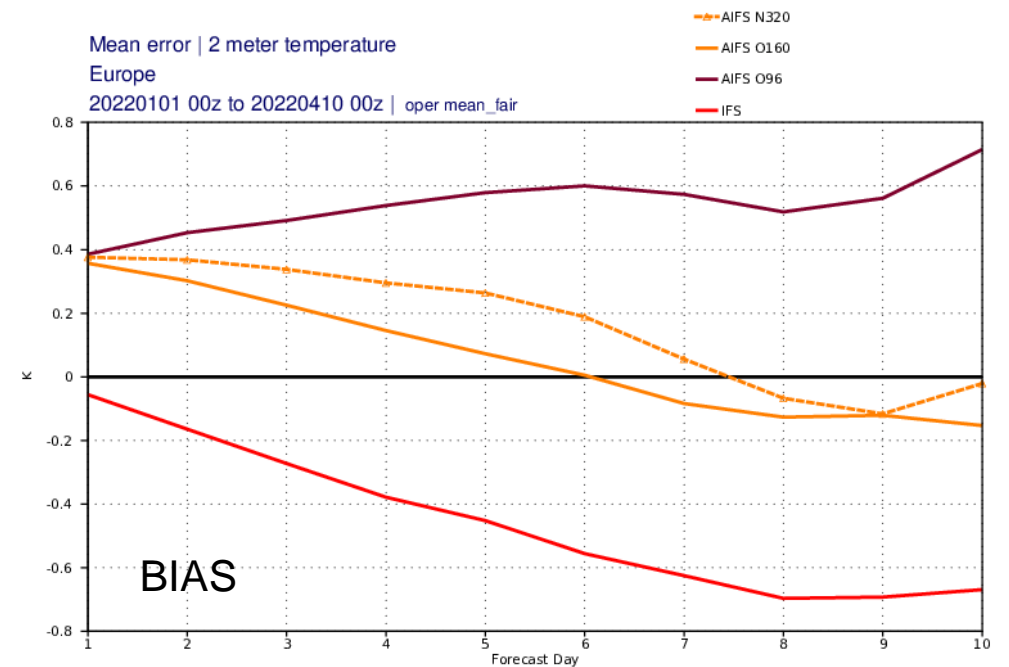
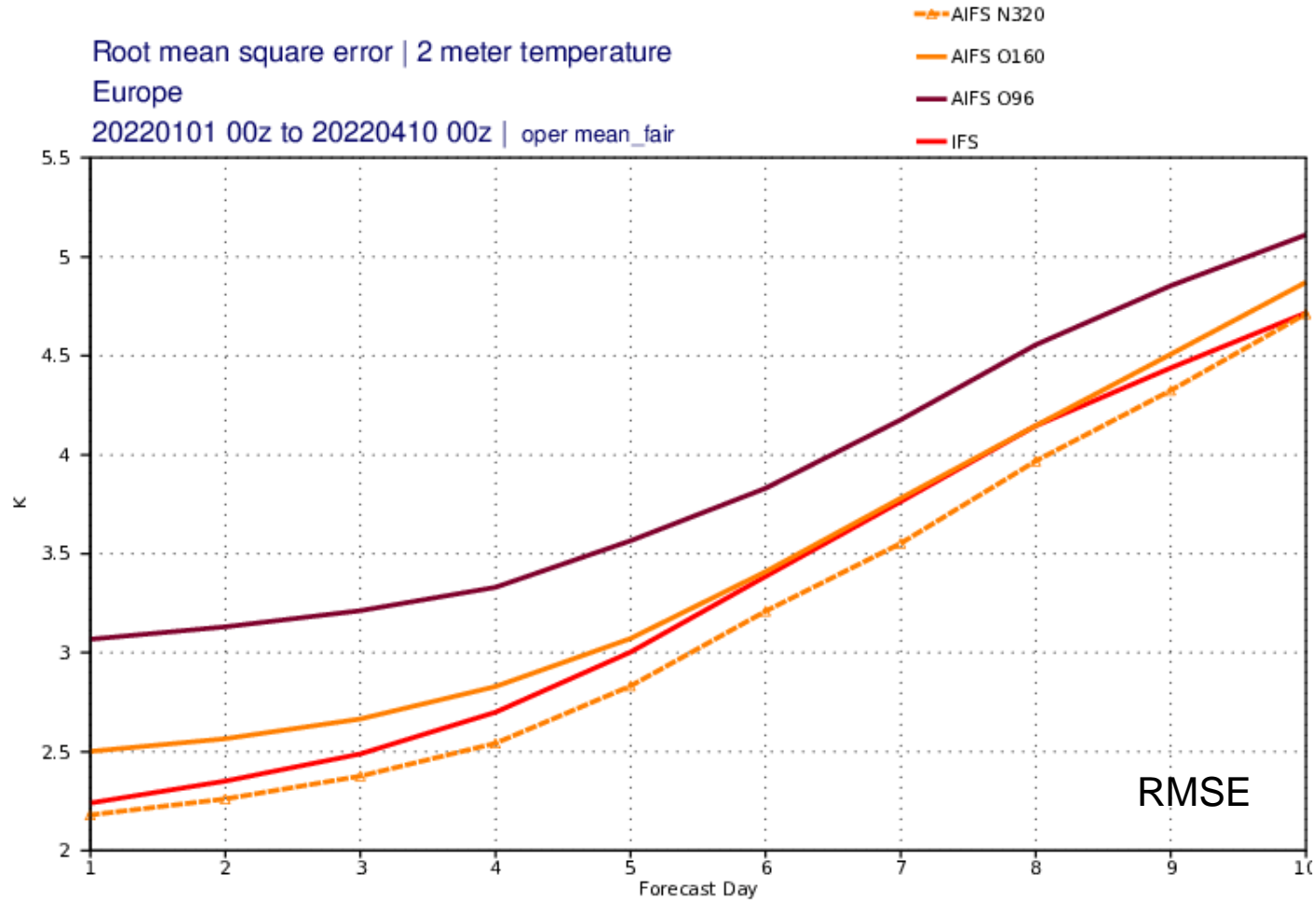
T850



Z500

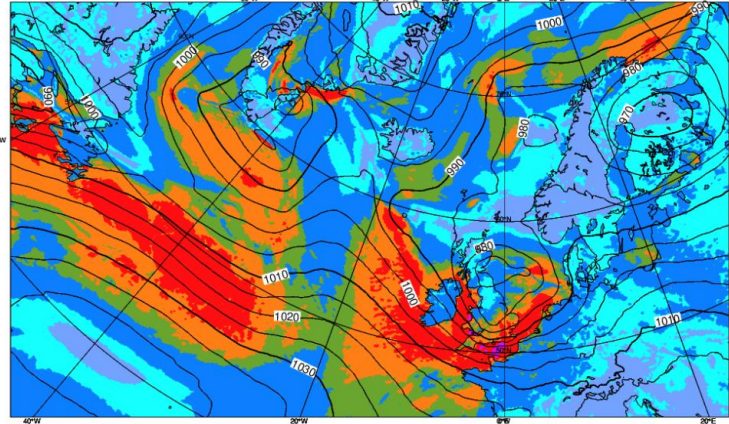
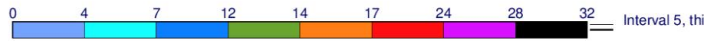


Verification against SYNOP observations



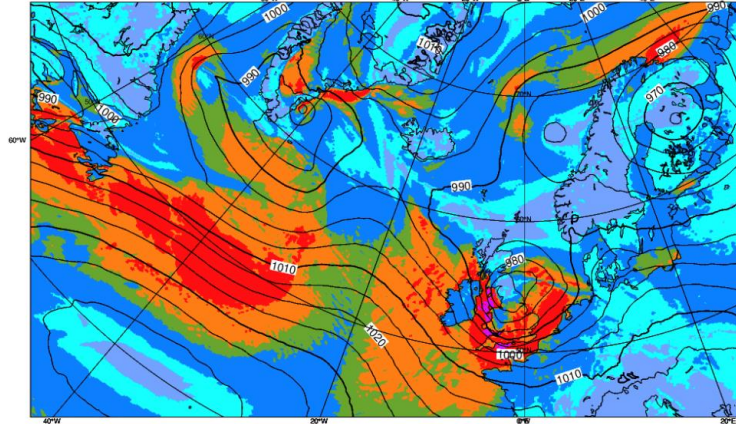
Analysis

MSLP+WS 2022021612 Step: 0
AN



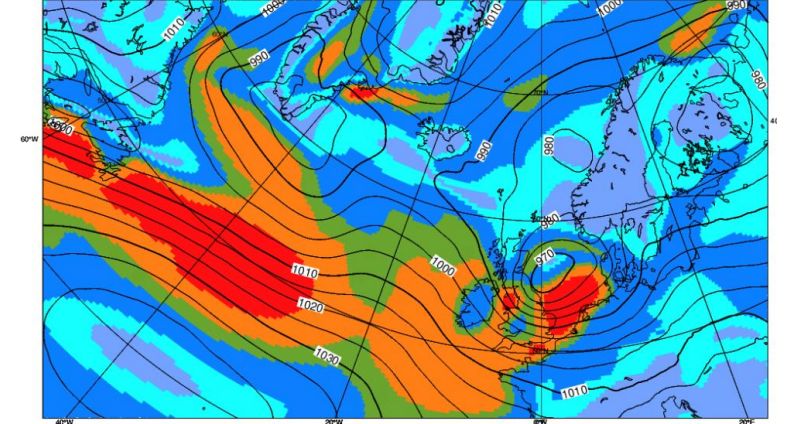
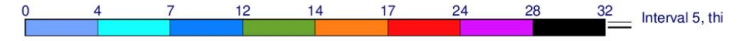
HRES

MSLP+WS 2022021600 Step: 60
HRES



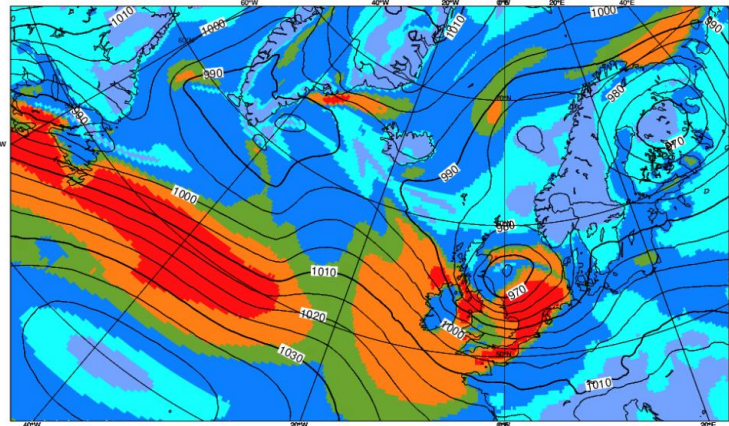
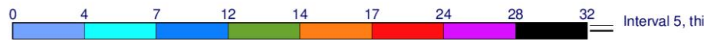
Fourcastnet

MSLP+WS 2022021600 Step: 60
i51c



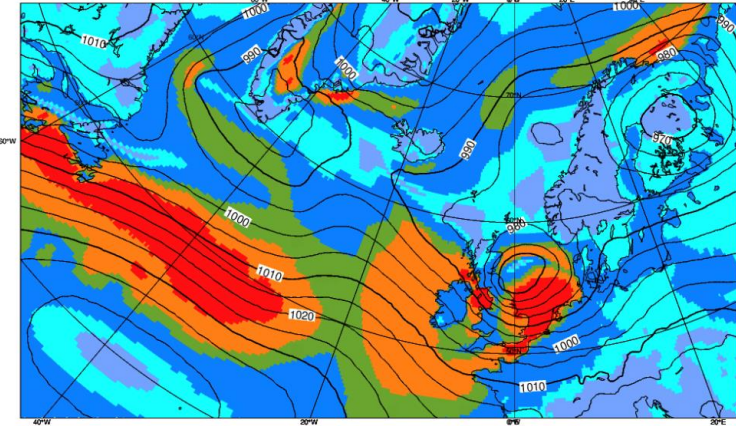
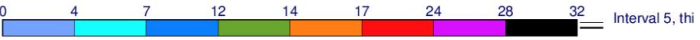
PanguWeather

MSLP+WS 2022021600 Step: 60
PanguWeather



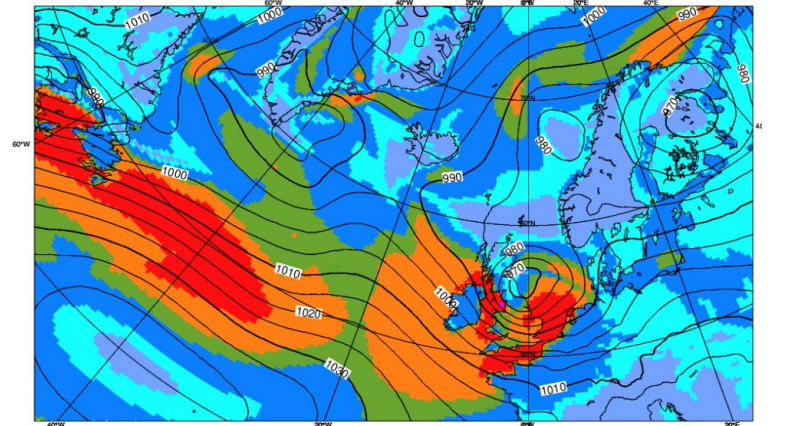
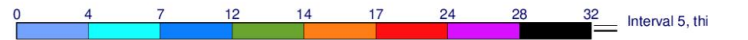
Graphcast

MSLP+WS 2022021600 Step: 60
i51d



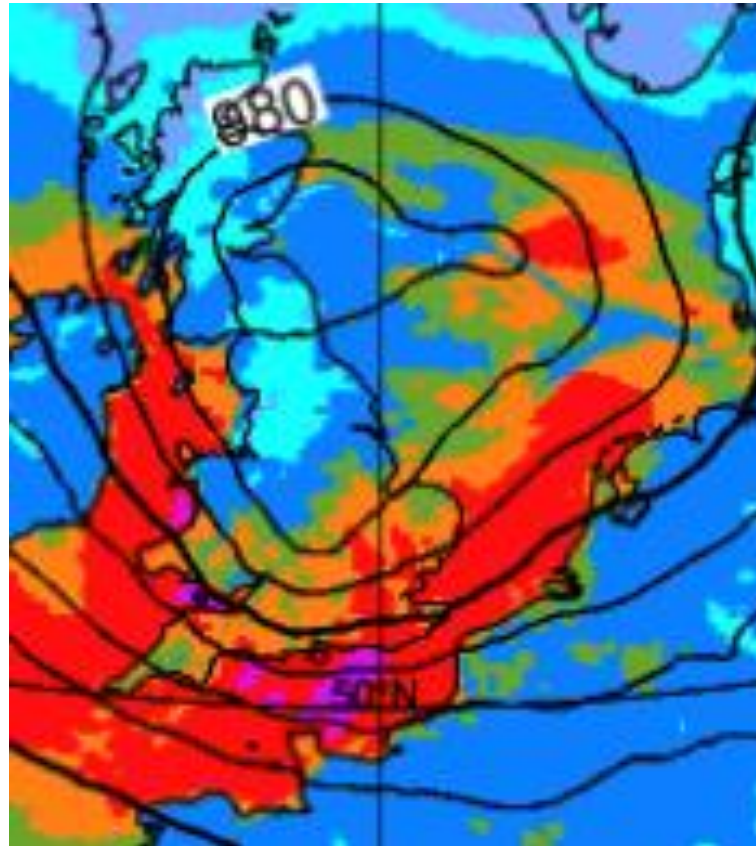
AIFS n320

MSLP+WS 2022021600 Step: 60
i5e6

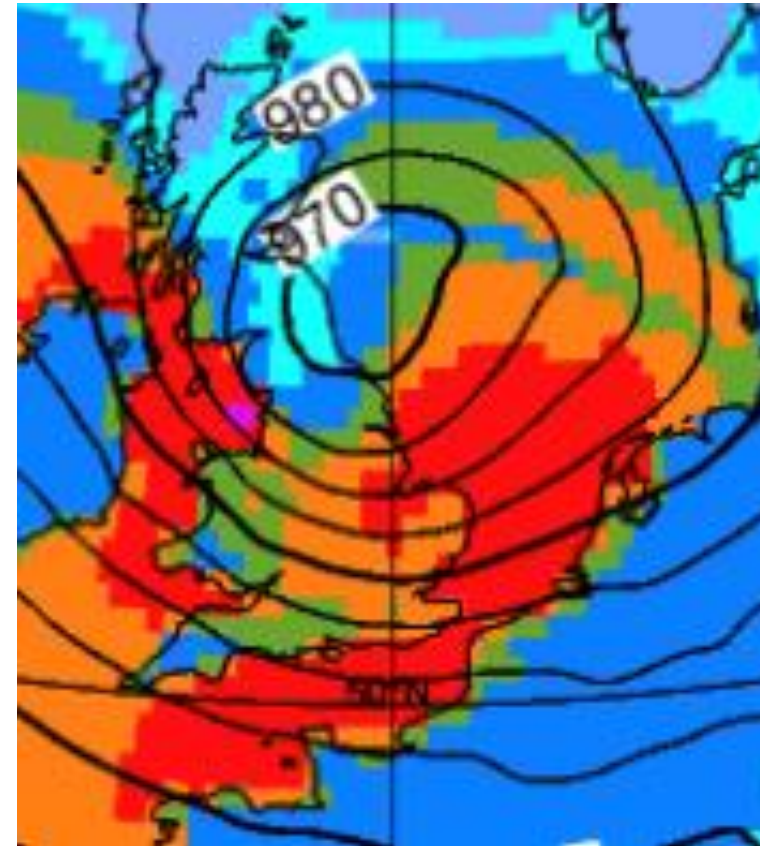


2022-02-16 00z + 60h (Storm Eunice over UK)

IFS

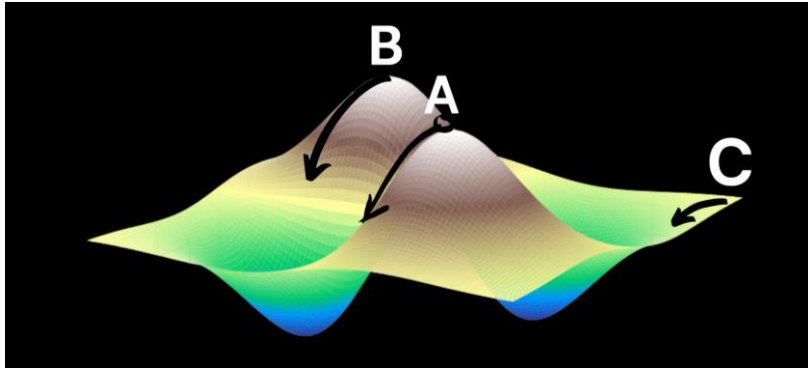


AIFS

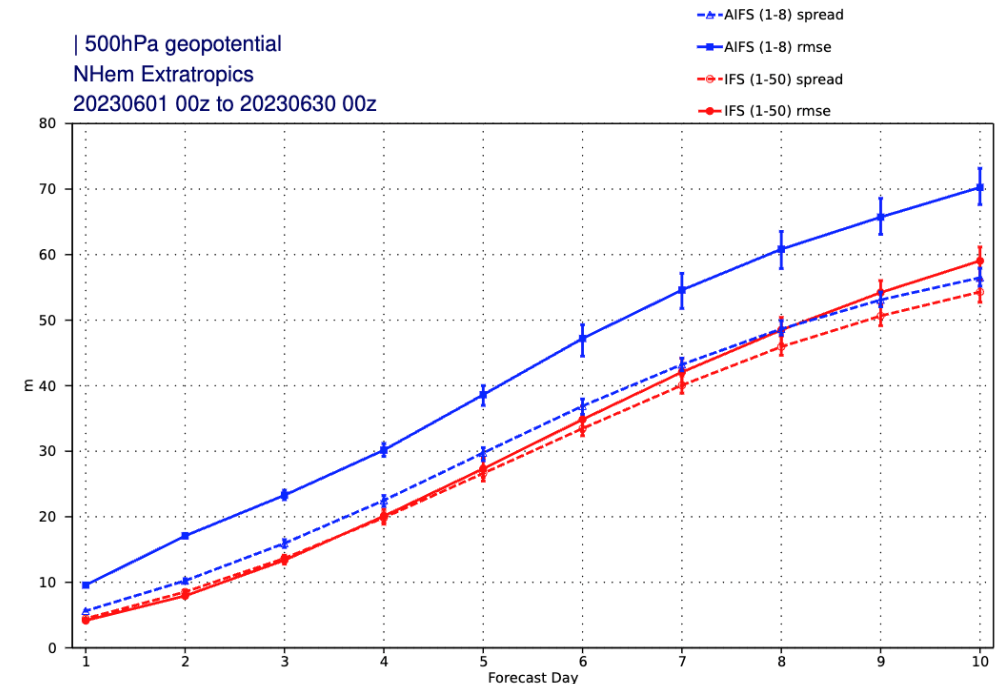


Towards Ensembles -> Preliminary ideas to capture uncertainty

1. Use the **ensemble initial conditions** to initialise the AI model
2. Using the randomness of the optimisation process to find different quasi-optimum solutions of the NN weights that can be used to build a **multi-model ensemble**



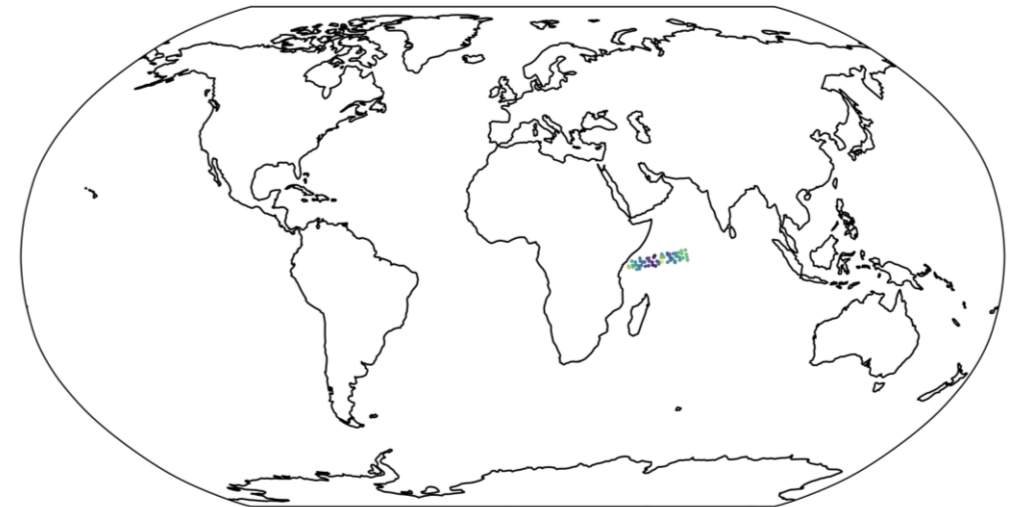
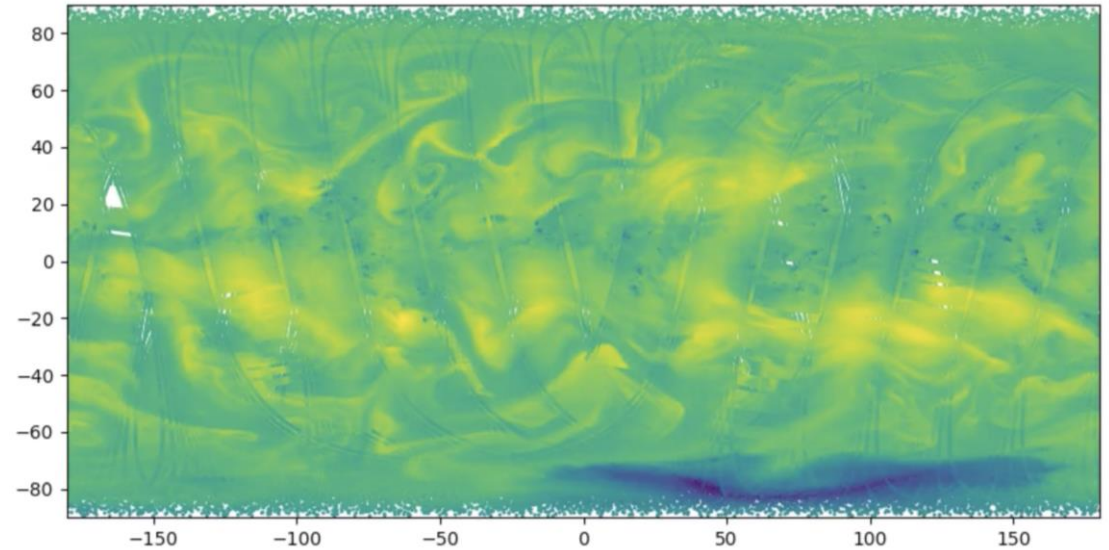
3. Train to minimise **probabilistic scores** like kernel CRPS or optimise distributions
4. **Generative models** where NN is trained on real samples and then generates new samples. Examples include Generative Adversarial Networks and Diffusion models



Learning from observations: challenges

Multitude of scattered and noisy data needs to be fused into coherent representation.

- Data is sparse in space and time
- Heterogeneity of data sources
 - Different vertical levels, ...
 - Many sources need to be combined
- Measured values are not standard physical fields and differ by source
 - Of interest are quantities like T2m, U10m, V10m, ...
 - Observation operators if necessary
- How to handle biases and quality control?



AI-Models Plugins for FOSS Data-Driven NWP

2211.02556 Pangu-Weather: A 3D High-Resolution Model
Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen, Xiaotao Gu, Qi Tian

2202.11214 FourCastNet: A Global Data-driven High-resolution Neural Operators
Jaideep Pathak, Shashank Subramanian, Peter Harrington, Sanjeev Raja, Ashes Zongyi Li, Kamyar Azizzadenesheli, Pedram Hassanzadeh, Karthik Kashinath, Zongyi Li, Kamyar Azizzadenesheli, Pedram Hassanzadeh, Karthik Kashinath, Zongyi Li

2212.12794 GraphCast: Learning skillful medium-range weather forecasting with a global deep learning model
Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wilmshausen, Wei Hua Hu, Alexander Merose, Stephan Hoyer, George Holm, B Battaglia

2306.03838 Spherical Fourier Neural Operators: Learning Stable Dynamics on the Sphere
Boris Bonev, Thorsten Kurth, Christian Hundt, Jaideep Pathak, Maximilian Baust, Karthik Kashinath, Anima Anandkumar

> pip install ai-models-panguweather

> ai-models panguweather

ONNX for model weights

> pip install ai-models-fourcastnet

> ai-models fourcastnet

PyTorch for code and model weights

> pip install ai-models-graphcast

> ai-models graphcast

Jax for code and model weights

> pip install ai-models-fourcastnetv2

> ai-models fourcastnetv2

PyTorch for code and model weights

prepmi

- **prepmi** is the companion tool to **ai-models**
- Uses ecFlow
- It allows to run inferences over many years
- Archives all outputs in the MARS archive in research mode
- It feeds into ECMWF's scores database so that models can be evaluated
- It allows users to run development code as well
- It can create ensembles using various combinations of models, inputs, ...

Summary

- We have an accurate baseline model
- Model can be scaled up to high-resolution
- Extensive supporting software infrastructure
- GNN can support arbitrary grids, including high resolution over only parts of the globe
- Using this basis we can explore
 - Different methods for constructing reliable ensemble forecasts
 - Using observations either with or replacing analysis.

Questions?

