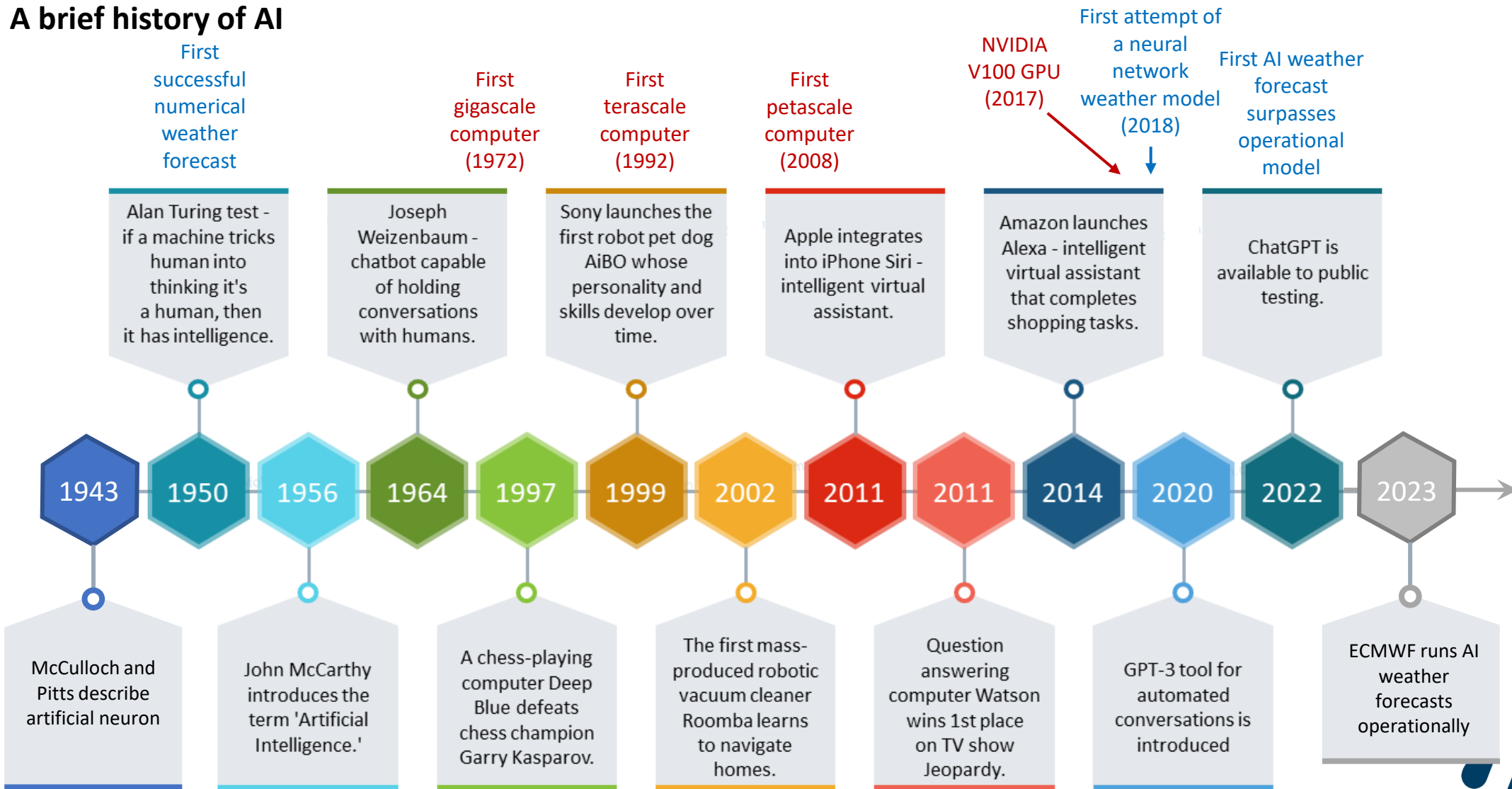


Is there an end to deep learning for weather and climate?

Martin Schultz, Jülich Supercomputing Center, Germany

A brief history of AI



Adapted from <https://www.infodiagram.com>

Review Article | Published: 02 September 2015

The quiet revolution of numerical weather prediction

[Peter Bauer](#) , [Alan Thorpe](#) & [Gilbert Brunet](#)

[Nature](#) **525**, 47–55 (2015) | [Cite this article](#)

48k Accesses | **1239** Citations | **1116** Altmetric | [Metrics](#)

1960-2010

Perspective | Published: 22 February 2021

The digital revolution of Earth-system science

[Peter Bauer](#) , [Peter D. Dueben](#), [Torsten Hoefler](#), [Tiago Quintino](#), [Thomas C. Schulthess](#) & [Nils P. Wedi](#)

[Nature Computational Science](#) **1**, 104–113 (2021) | [Cite this article](#)

18k Accesses | **94** Citations | **300** Altmetric | [Metrics](#)

2005-2025

FOURCASTNET: A GLOBAL DATA-DRIVEN HIGH-RESOLUTION WEATHER MODEL USING ADAPTIVE FOURIER NEURAL OPERATORS

... skillful medium-range forecasting

... A 3D High-Resolution System for Accurate Global Weather Forecast

The AI revolution in weather and climate modeling

Ashesh Chattopadhyay Rice University Houston, TX 77005	Mertcan Mardani NVIDIA Corporation Santa Clara, CA 95051	Thorsten Kurth NVIDIA Corporation Santa Clara, CA 95051
David Hall NVIDIA Corporation Santa Clara, CA 95051	Zongyi Li California Institute of Technology Pasadena, CA 91125	Kanyar Azizadehshel Purdue University West Lafayette, IN 47907
Pedram Hasanizadeh Rice University Houston, TX 77005	Karthik Kashnath NVIDIA Corporation Santa Clara, CA 95051	Animeshree Anandkumar California Institute of Technology Pasadena, CA 91125

... used weather simulator—called “GraphCast”—which outperforms the most accurate previous operational medium-range weather forecasting system in the world. GraphCast is an autoregressive model, based on graph neural networks, which we trained on historical data for Medium-Range Weather Forecasts (ECMWF’s ERA5 reanalysis), at 6-hour time intervals, of five surface variables and geopotential pressure levels, on a 0.25° latitude-longitude grid, which is the resolution at the equator. Our results show GraphCast is more accurate than the current operational forecasting system, HRES, on 90.0% of the 2760 variables. GraphCast also outperforms the most accurate previous

... a few deep neural networks with about 256 million parameters in total. The spatial resolution is comparable to the ECMWF Integrated Forecast Systems (IFS). More importantly, for the first time, an end-to-end numerical weather prediction (NWP) model in terms of accuracy (latitude-weighted global, specific humidity, wind speed, temperature, etc.) and in all time ranges (from one hour to one week). GraphCast is an autoregressive model, based on graph neural networks, which we trained on historical data for Medium-Range Weather Forecasts (ECMWF’s ERA5 reanalysis), at 6-hour time intervals, of five surface variables and geopotential pressure levels, on a 0.25° latitude-longitude grid, which is the resolution at the equator. Our results show GraphCast is more accurate than the current operational forecasting system, HRES, on 90.0% of the 2760 variables. GraphCast also outperforms the most accurate previous

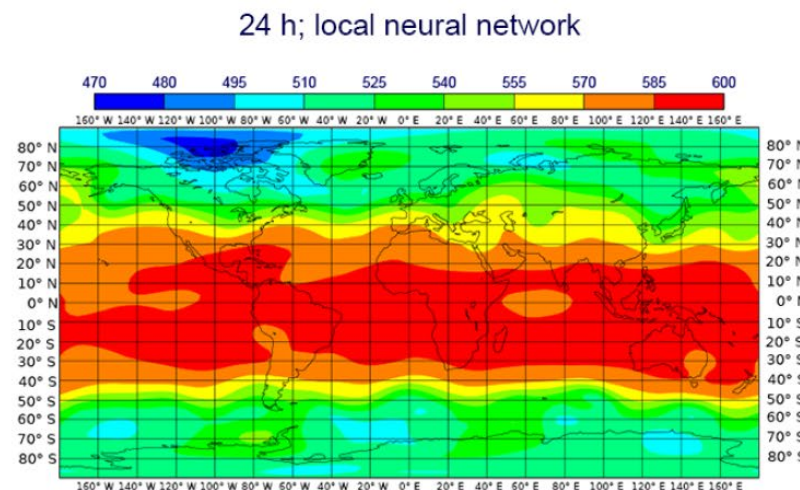
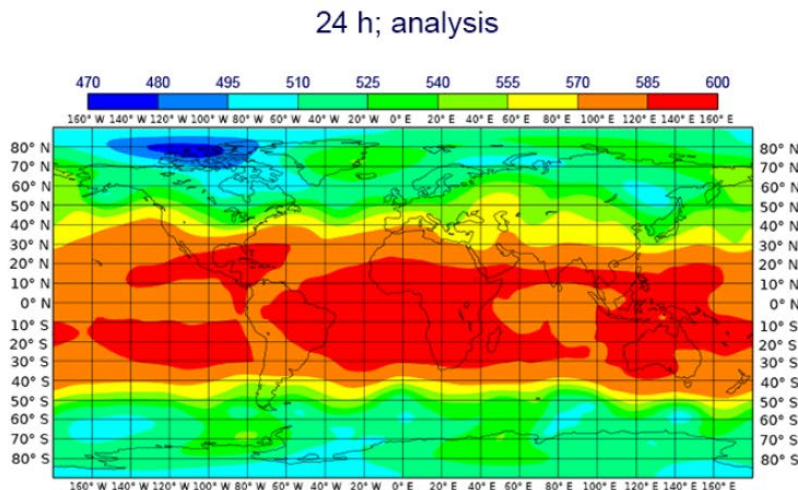
2022-

A few milestones on the path to AI weather forecasts

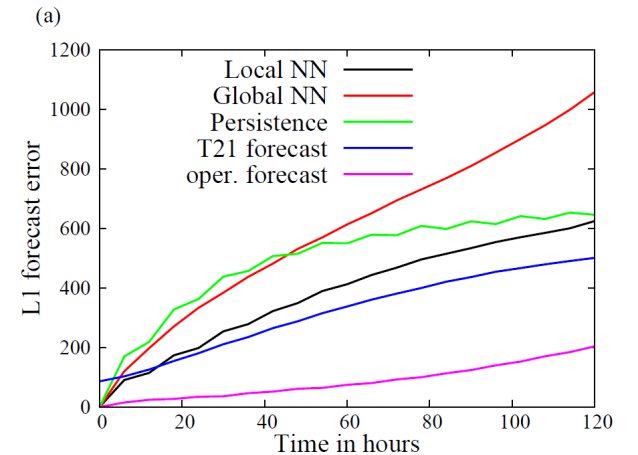
2018: First trial of a pure ML model: FNN with 4 layers

Input: hourly Z500, 1 pressure level, 1860 grid points (6° resolution), 67200 snap shots 2010—2017

Output: Z500 up to 120 hours ahead (autoregressive rollout)



Model error



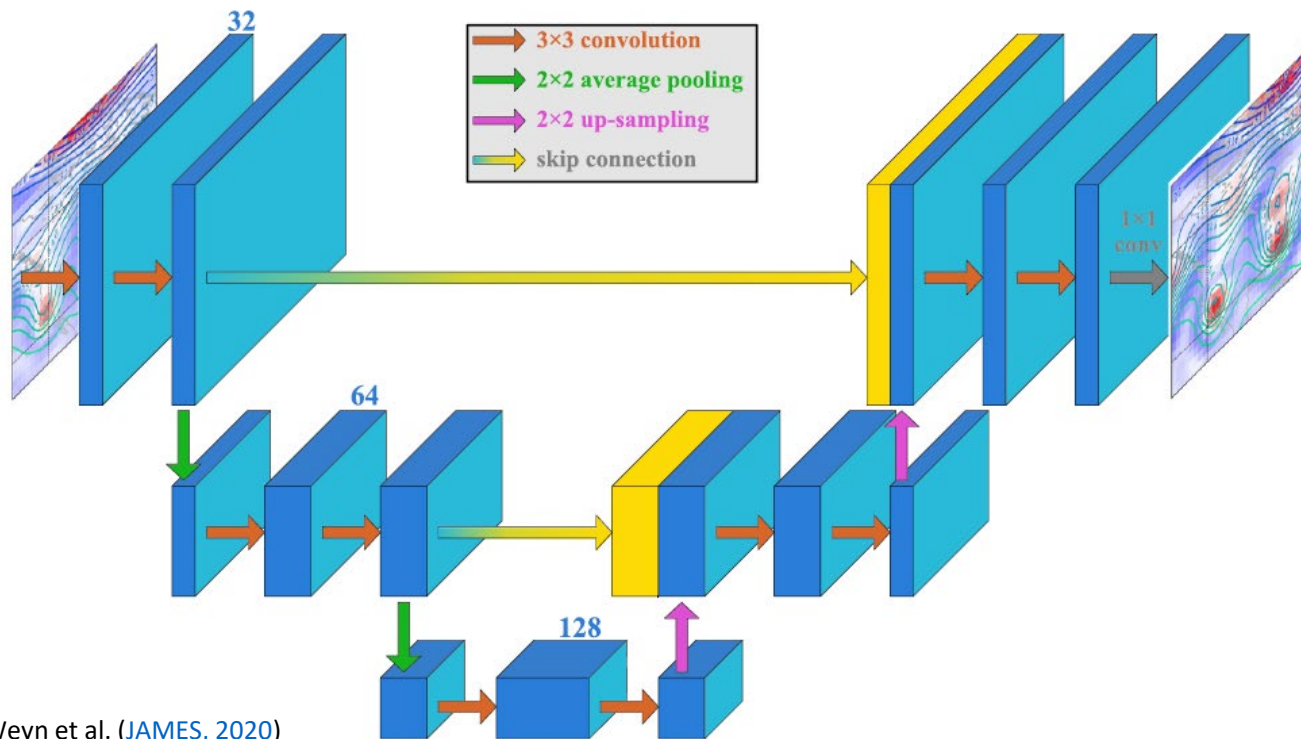
Düben and Bauer ([GMD, 2018](#))

A few milestones on the path to AI weather forecasts

2020: global forecasts up to 14 days with a U-net, 11 Conv2D layers

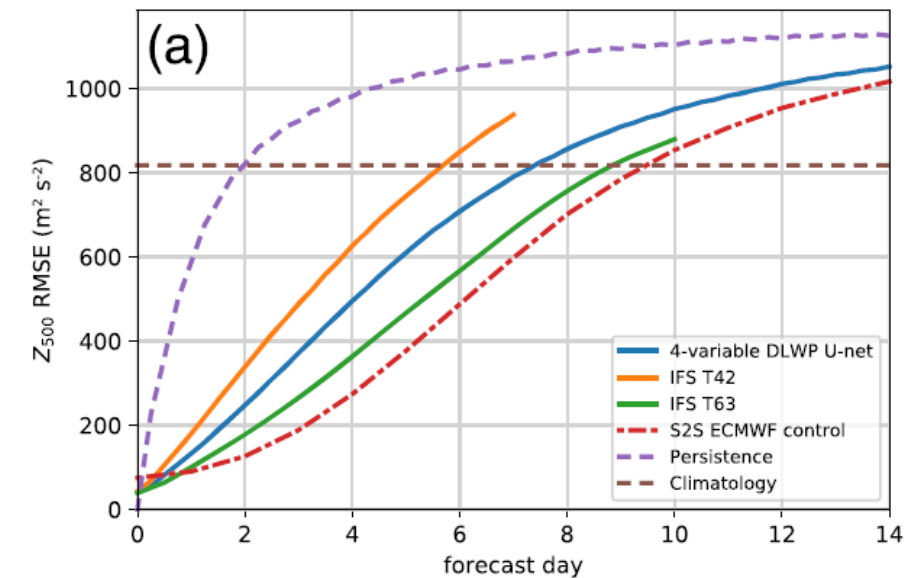
Input: 4 variables: Z_{500} , Z_{1000} , $\tau_{300-700}$, T_{2m} , 6-hourly data, 1917-2012 (100,000 samples), 2° horizontal resolution

Output: 4 variables, 6-hourly for t+6 and t+12



Weyn et al. ([JAMES, 2020](#))

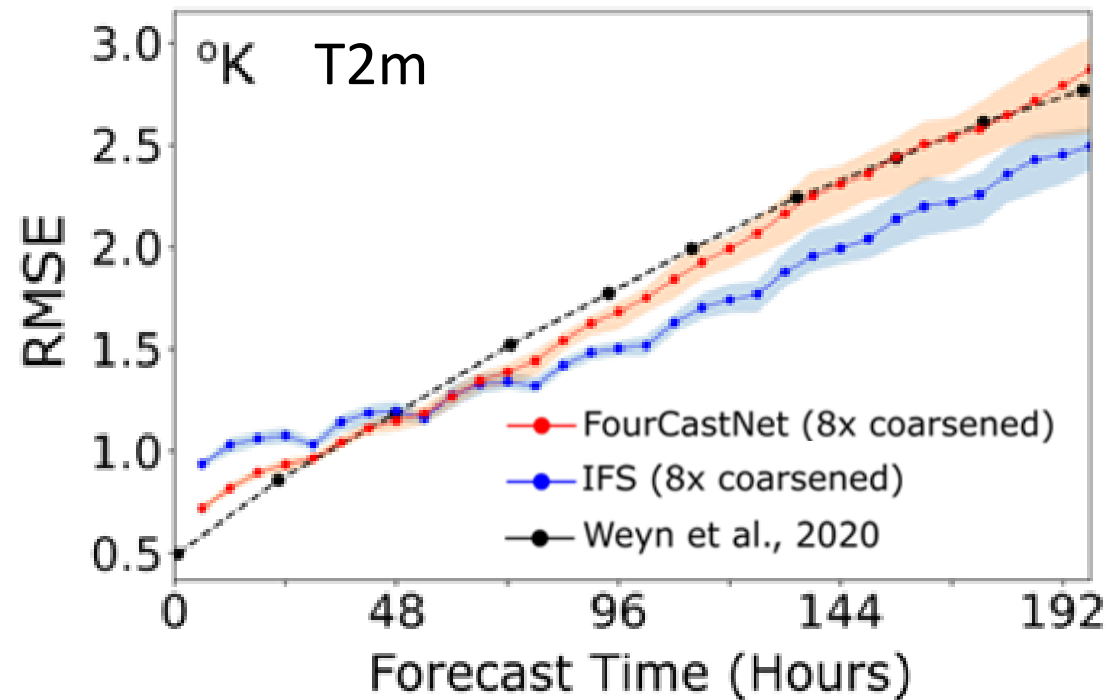
Model error (RMSE)



A few milestones on the path to AI weather forecasts

2022: Global forecasts at 0.25 ° resolution, 20 variables at 5 pressure levels, Transformer + Fourier Neural Operators

Model error

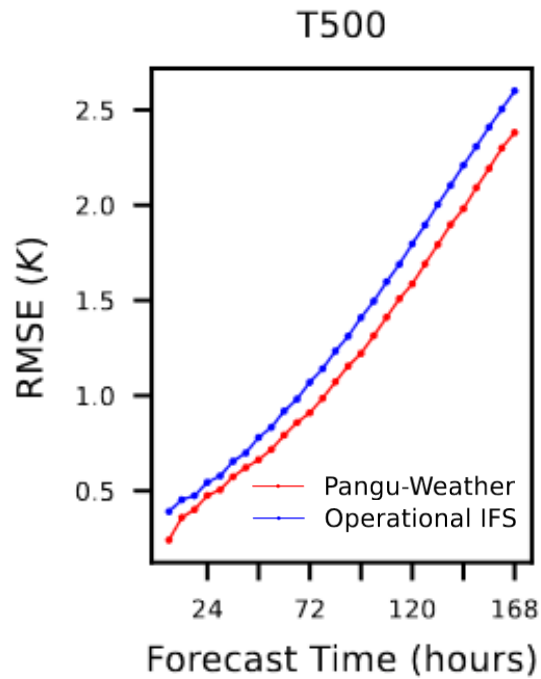


Pathak et al. (2022): Fourcastnet

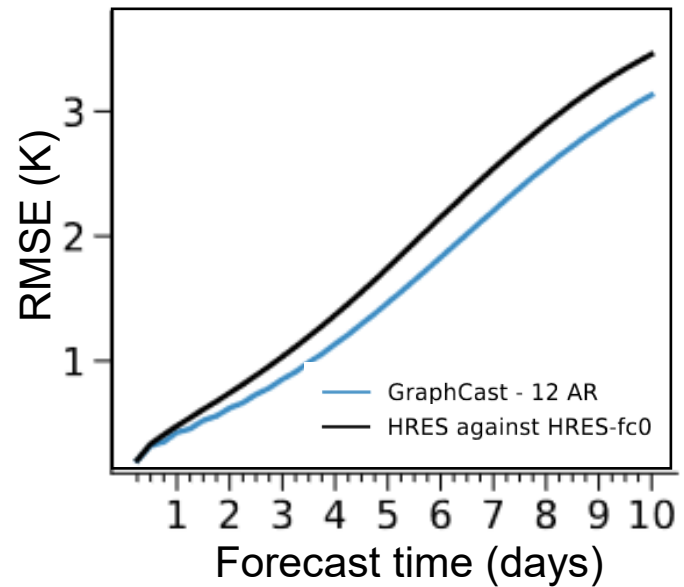
A few milestones on the path to AI weather forecasts

2022/23: The breakthrough – DL models outperform IFS HRES

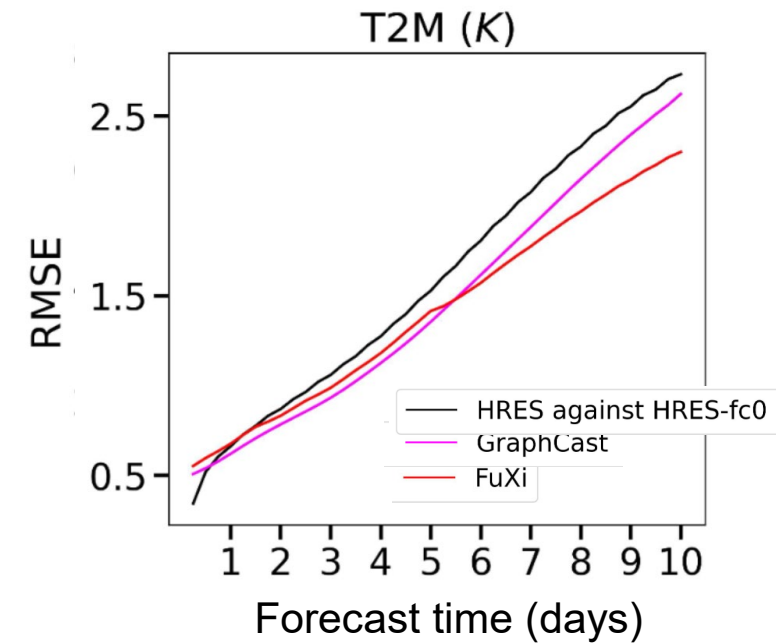
Model error



Bi et al. (2022): Pangu-Weather



Lam et al. (2022): GraphCast

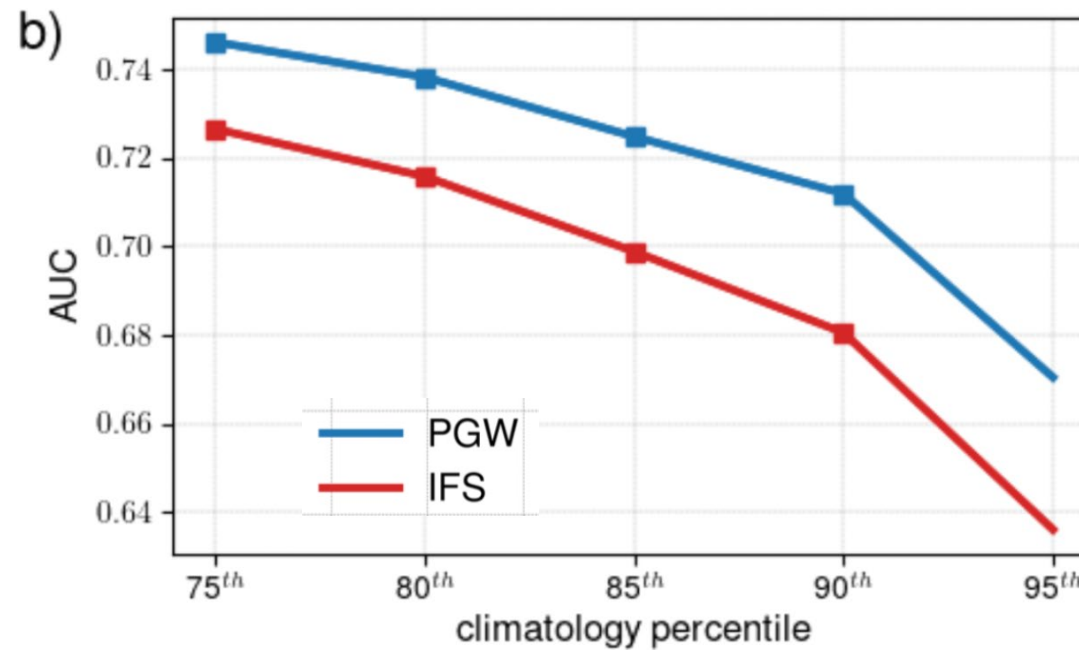


Chen et al. (2023): FuXi

Limitations of DL weather models

2022/23: The breakthrough – DL models outperform IFS HRES; **but are they really better?**

T2m extremes Europe, summer 2022

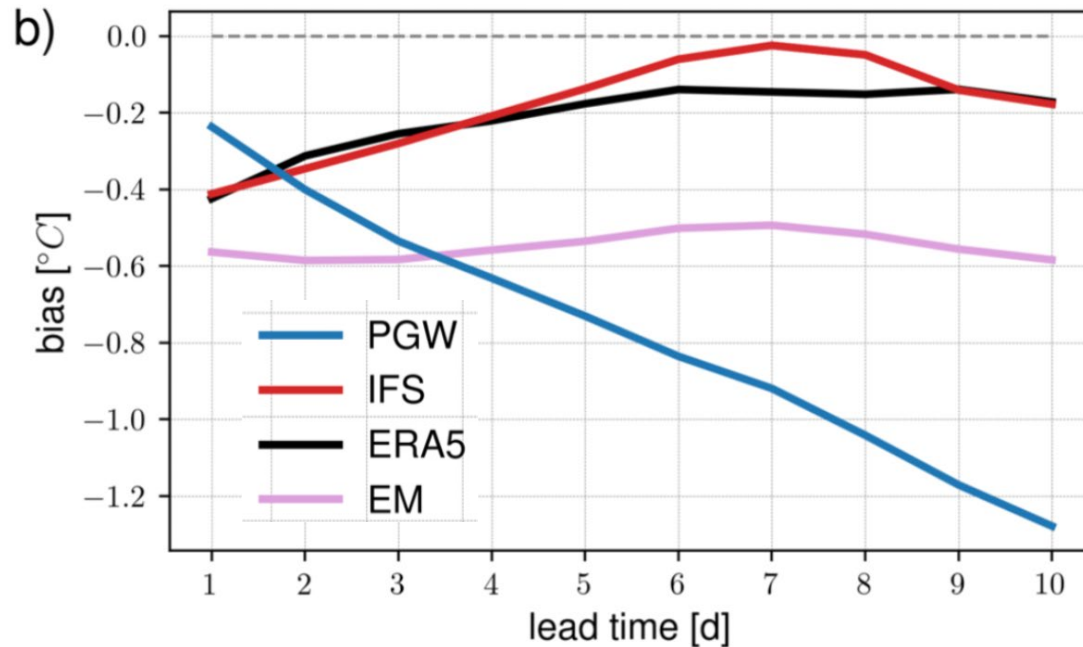


Ben-Bouallegue et al. ([Arxiv, 2023](#))

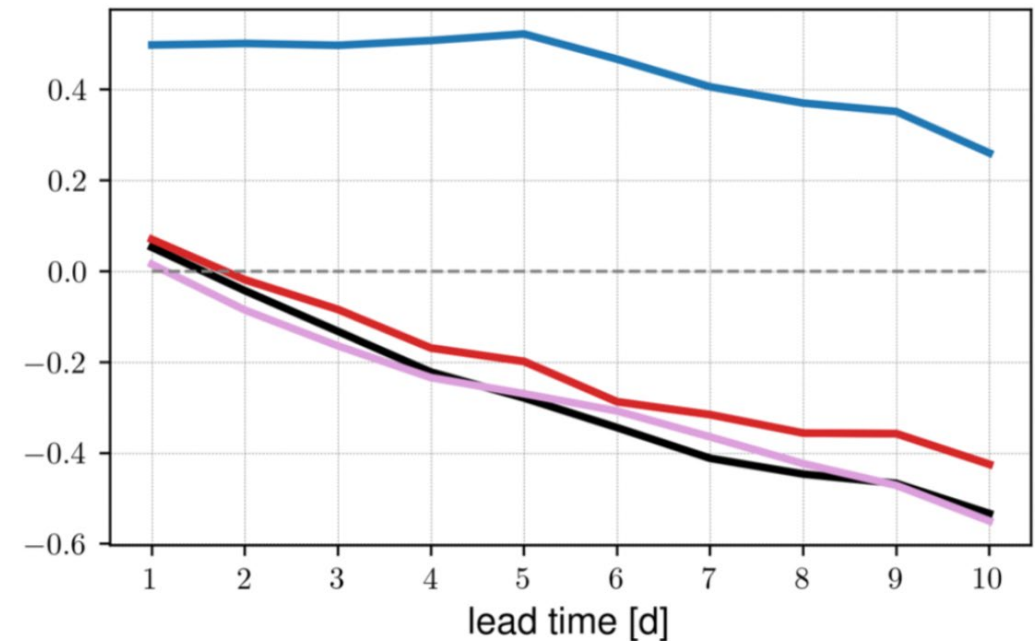
Limitations of DL weather models

2022/23: The breakthrough – DL models outperform IFS HRES; **but are they really better?**

T2m Europe, summer 2022

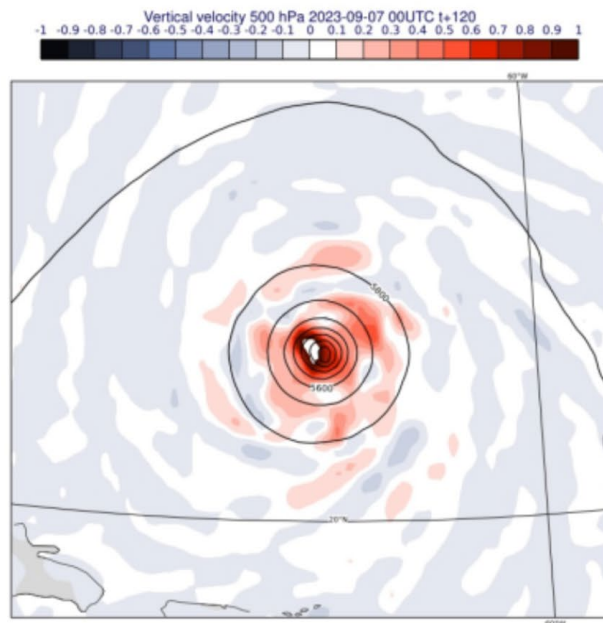


T2m Europe, winter 2022/23

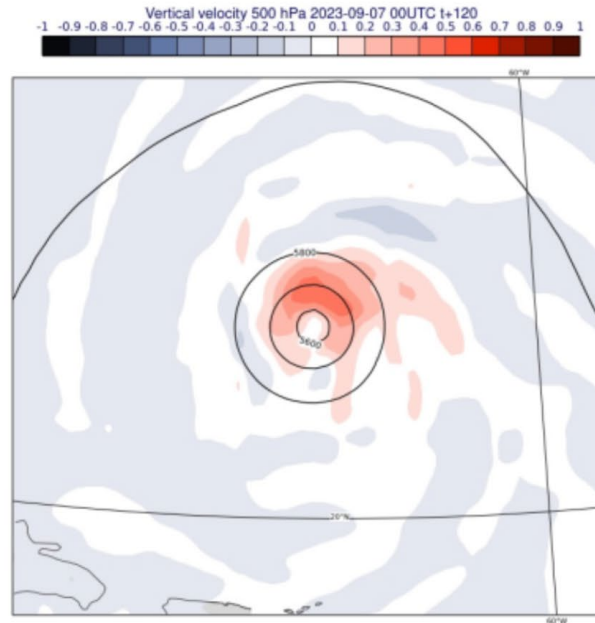


Limitations of DL weather models

IFS-HRES



ERA5



PanguWeather

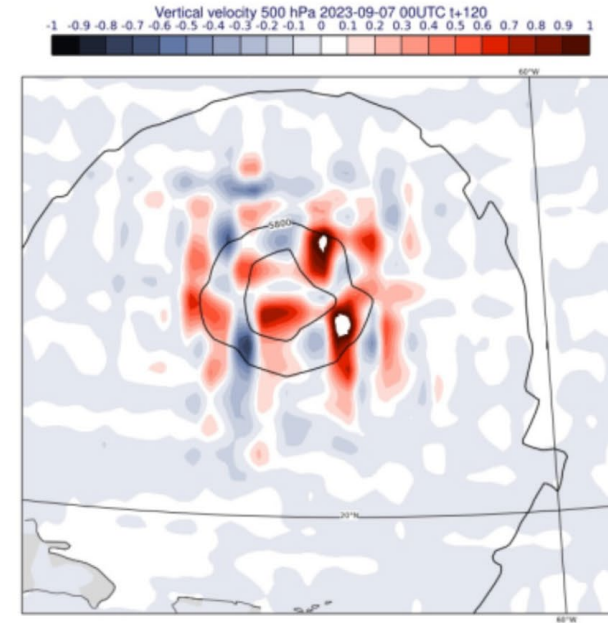


Figure 10: *Plots of t+120 hours forecast vertical velocities (shaded, units: m/s) at 500 hPa from:*

Bonavita, [2024](#)

DNN stability and robustness

„However, for the local networks that use the special treatment of the area around the pole, as discussed in the previous section, the forecast error diverges for the 7×7 and the 9×9 configuration. [...likely that this can be fixed...]“ (Düben and Bauer, 2018)

„Every one of the 4-week forecasts initialized twice weekly in the 2-year test set (210 total forecasts) was free from instabilities and the amplification of spurious perturbations.“ (Weyn et al., 2020)

„Moreover, the similarities in error growth of a data-driven forecast and a standard NWP forecast indicate similar sensitivities to chaos between ML-based and physically-based models.“ (Ben Bouallegue et al., 2023)

„More generally, the discussion above and the results presented here highlight one of the main challenges for the next generation of data-driven ML prediction models, namely, how to produce forecasts that are skilful and at the same time dynamically and physically consistent at all relevant spatial scales.“ (Bonavita, 2024)

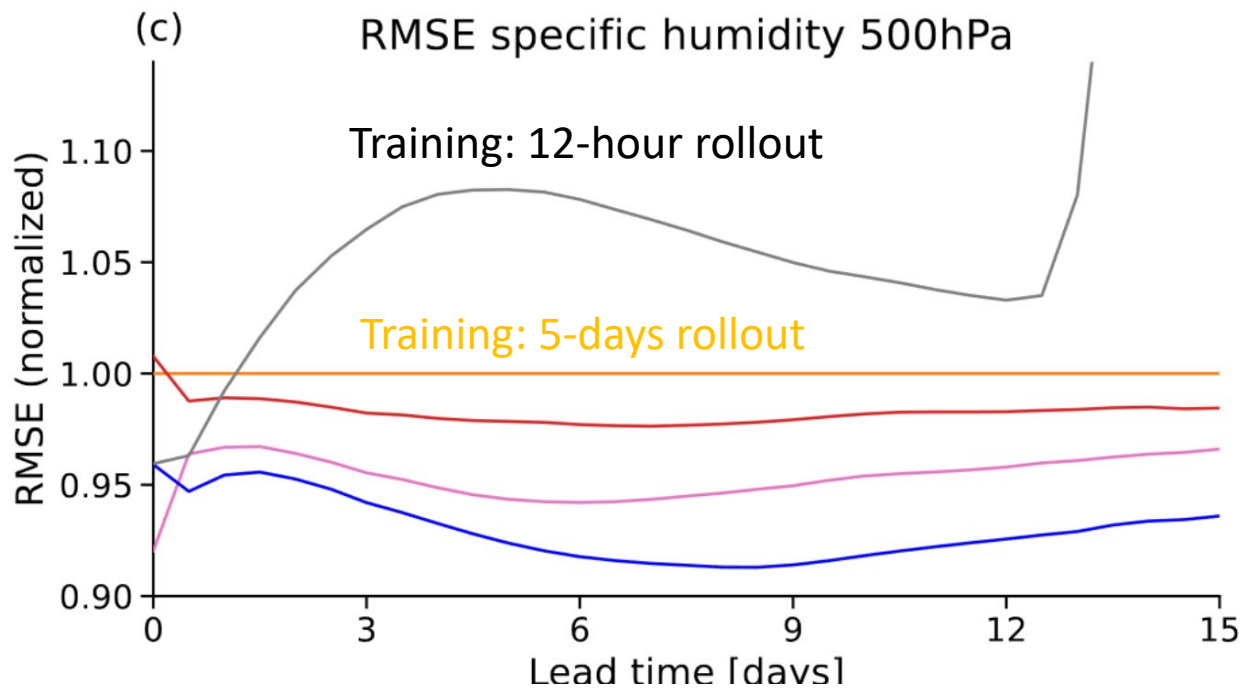
Long-term rollout

Weyn et al. (2020): up to 1 year

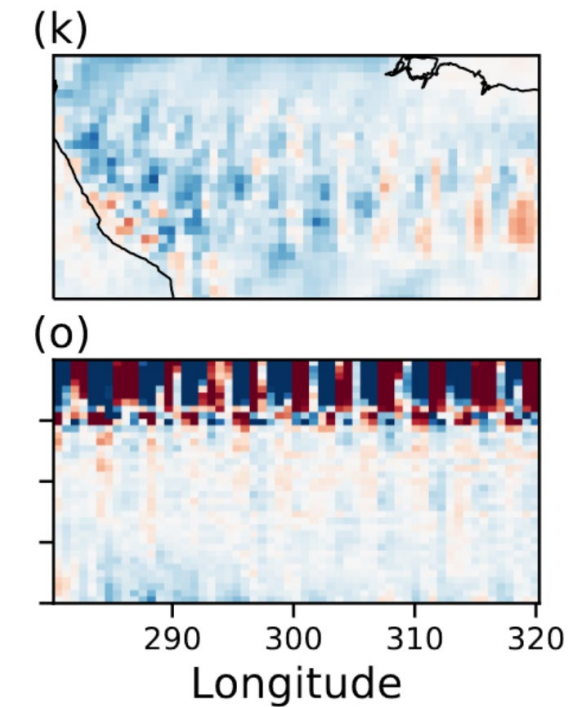
Watt-Meyer et al. (2023) [ACE]: up to 100 years

Kochkov et al. (2024) [NeuralGCM]: up to 40 years (22 out of 37 runs stable)

Case study on temperature instability NeuralGCM; after 139 days



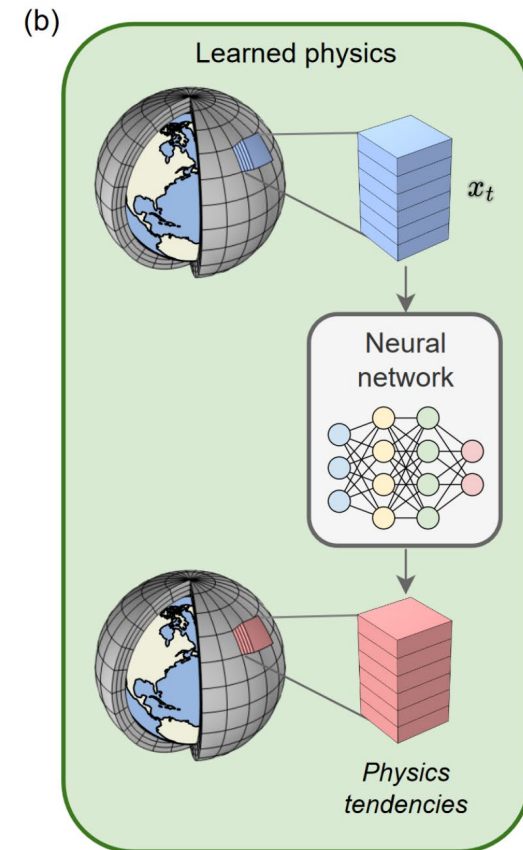
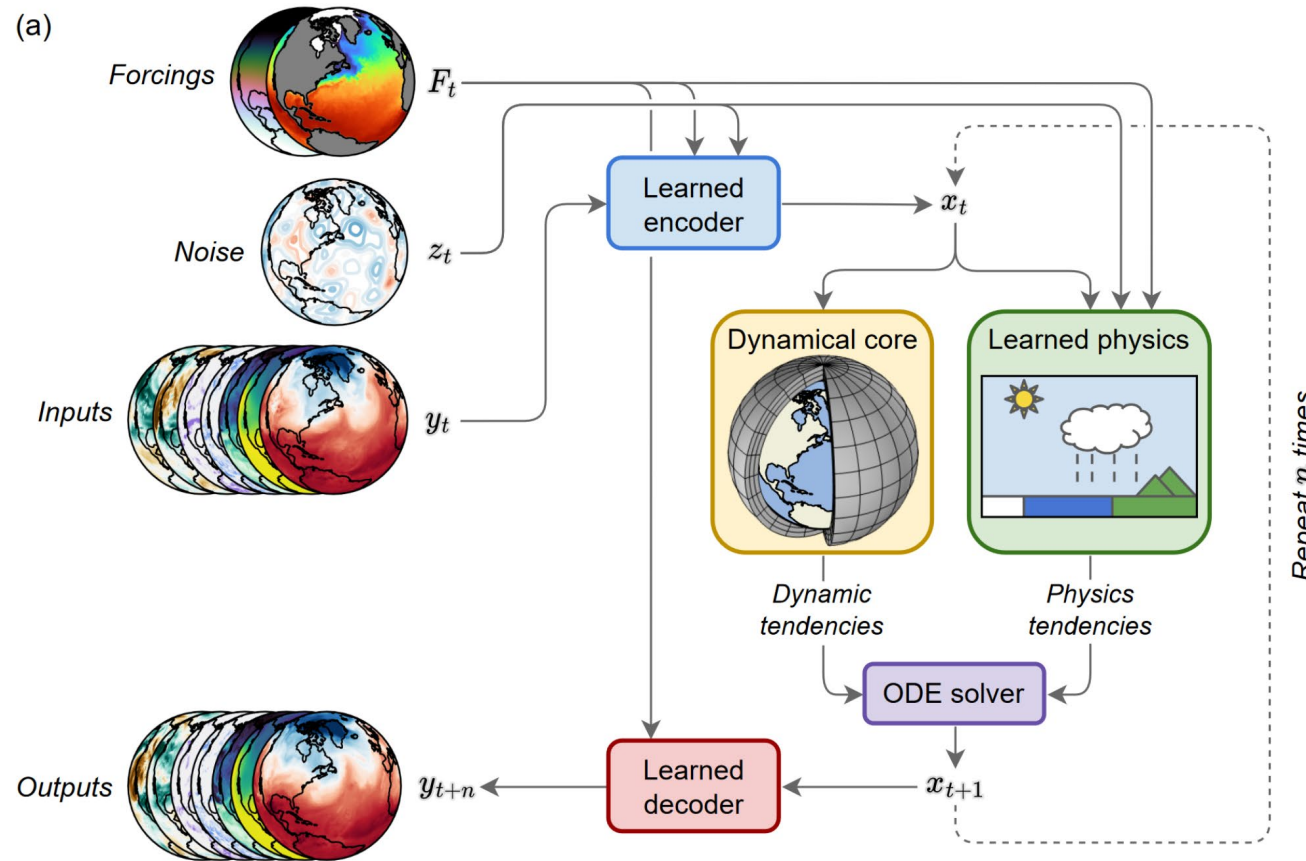
Kochkov et al. ([Arxiv, 2024](#))



Kochkov et al. ([Arxiv, 2024](#))

Climate time scales; NeuralGCM

High quality deterministic forecasts up to 10 days, probabilistic forecasts up to 15 days, and stable „weather“ on century time scales

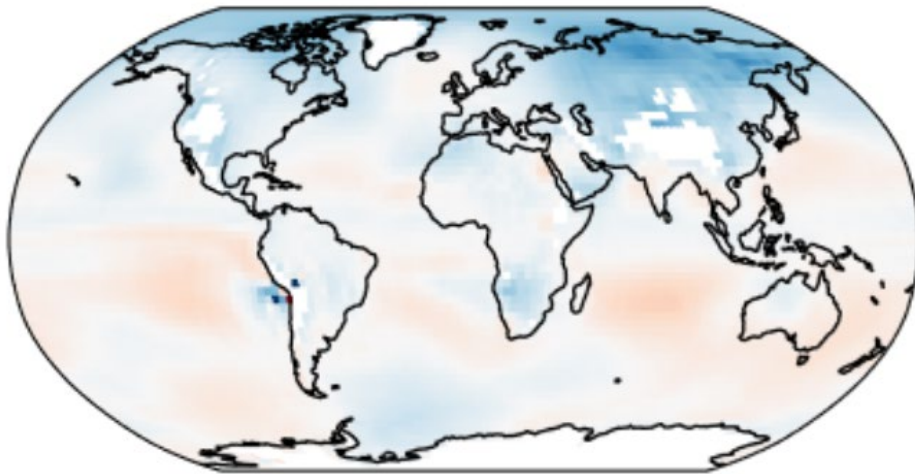


Climate time scales: Neural GCM achieves reduced bias in climate predictions

850 hPa temperature bias averaged 1981-2014

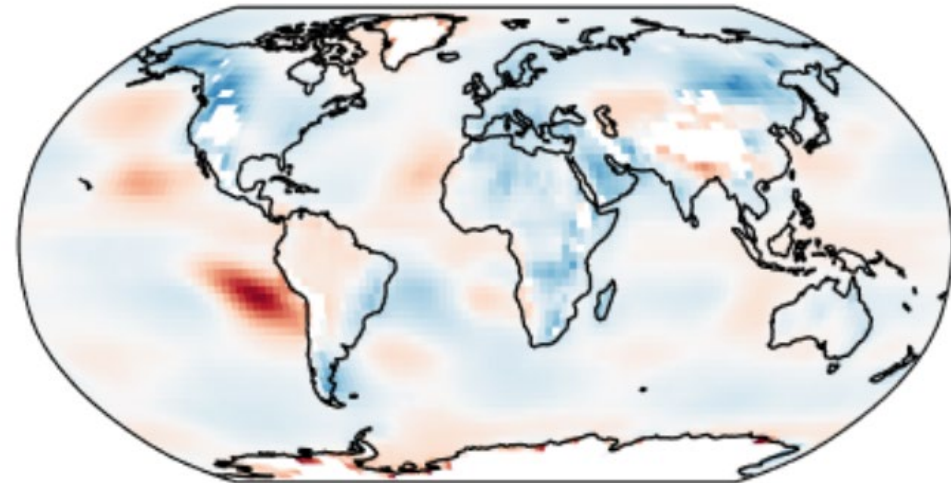
Worst of 22 Neural GCM simulations

RMSE=0.440 K



Best of 22 CMIP6 GCM simulations

GFDL-AM4 RMSE=0.475 K



Kochkov et al. ([Arxiv, 2024](#))

Bluriness

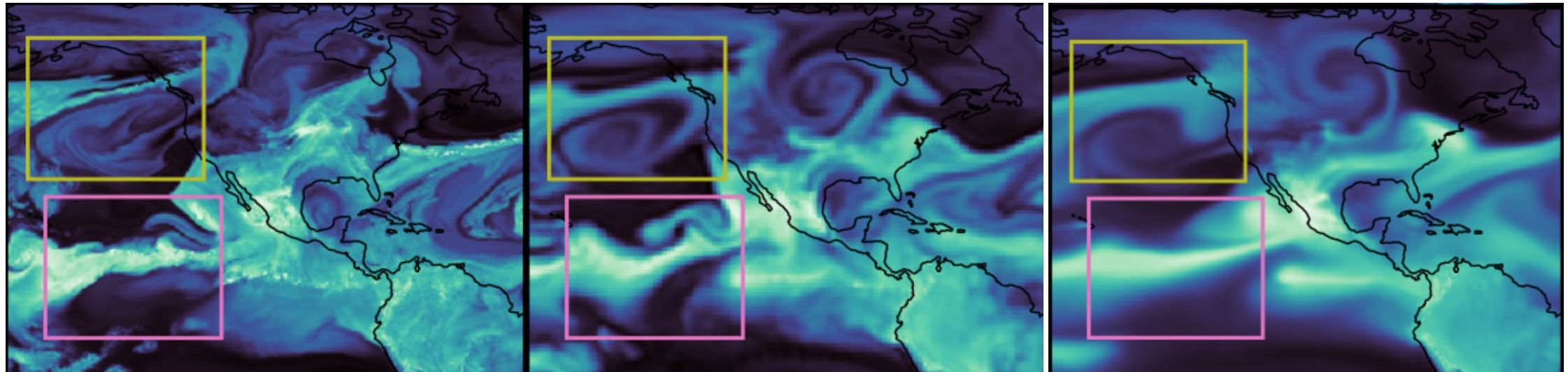
„The data-driven forecast appears smoother than the operational IFS forecast but the level of smoothness does not seem to increase with the forecast lead time, as we might expect when training toward RMSE.“ (Ben Bouallegue et al., 2023)

Specific humidity, 700 hPa, t+10 days

IFS HRES (0.1°)

NeuralGCM (0.7°)

Graphcast (0.25°)



Kochkov et al. ([Arxiv, 2024](#))

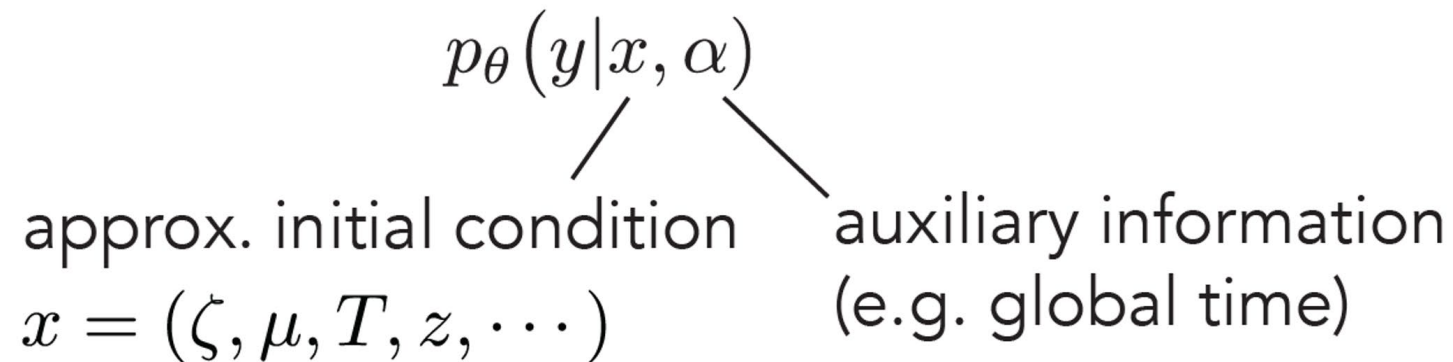
Ensemble modelling

Classical approach: perturb initial conditions and parametrisations (also adopted by FuXi)

Deep learning: Use generative models

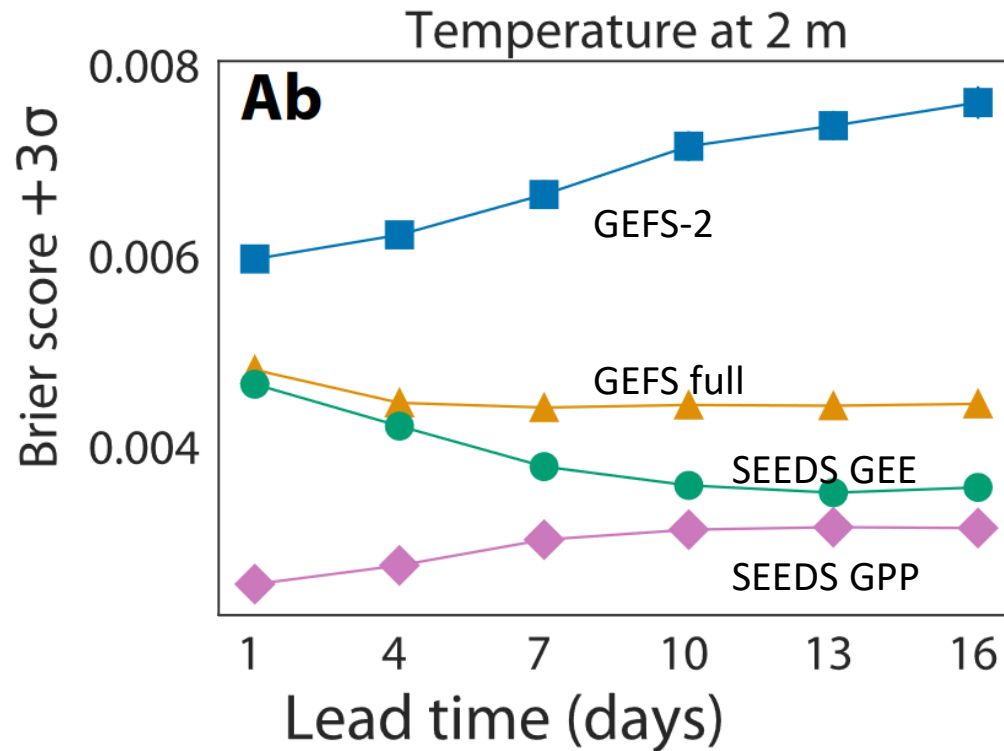
Example: AtmoRep (Lessig et al., [2023](#)) → talk by Ilaria Luise tomorrow

Inherent probabilistic formulation (and probabilistic loss)



Ensemble modelling

Google's Scalable Ensemble Envelope Diffusion Sampler (SEEDS)



See also:

J. Leinonen, U. Hamann, D. Nerini, U. Germann, G. Franch, Latent diffusion models for generative precipitation nowcasting with accurate uncertainty quantification. arXiv:2304.12891 [physics.ao-ph] (25 April 2023).

Z. Gao, X. Shi, B. Han, H. Wang, X. Jin, D. Maddix, Y. Zhu, M. Li, Y. Wang, PreDiff: Precipitation nowcasting with latent diffusion models. arXiv:2307.10422 [cs.LG] (28 December 2023).

H. Addison, E. Kendon, S. Ravuri, L. Aitchison, P. Watson, Machine learning emulation of a local-scale UK climate model, in *NeurIPS 2022 Workshop on Tackling Climate Change with Machine Learning* (Climate Change AI, 2022); www.climatechange.ai/papers/neurips2022/21/paper.pdf.

S. Bassetti, B. Hutchinson, C. Tebaldi, B. Kravitz, DiffESM: Conditional emulation of Earth system models with diffusion models. arXiv:2304.11699 [physics.ao-ph] (23 April 2023).

Challenges ahead

- Direct use of observations (work in progress)
- Climate scenarios
- Multi-scale models (has been demonstrated → SEEDS)
- Upper atmosphere and tracer transport (work initiated)
- Earth system modeling: ocean, sea ice, land, biogeochemical cycles, atmospheric chemistry

Upper atmosphere more difficult? Or simply overlooked?

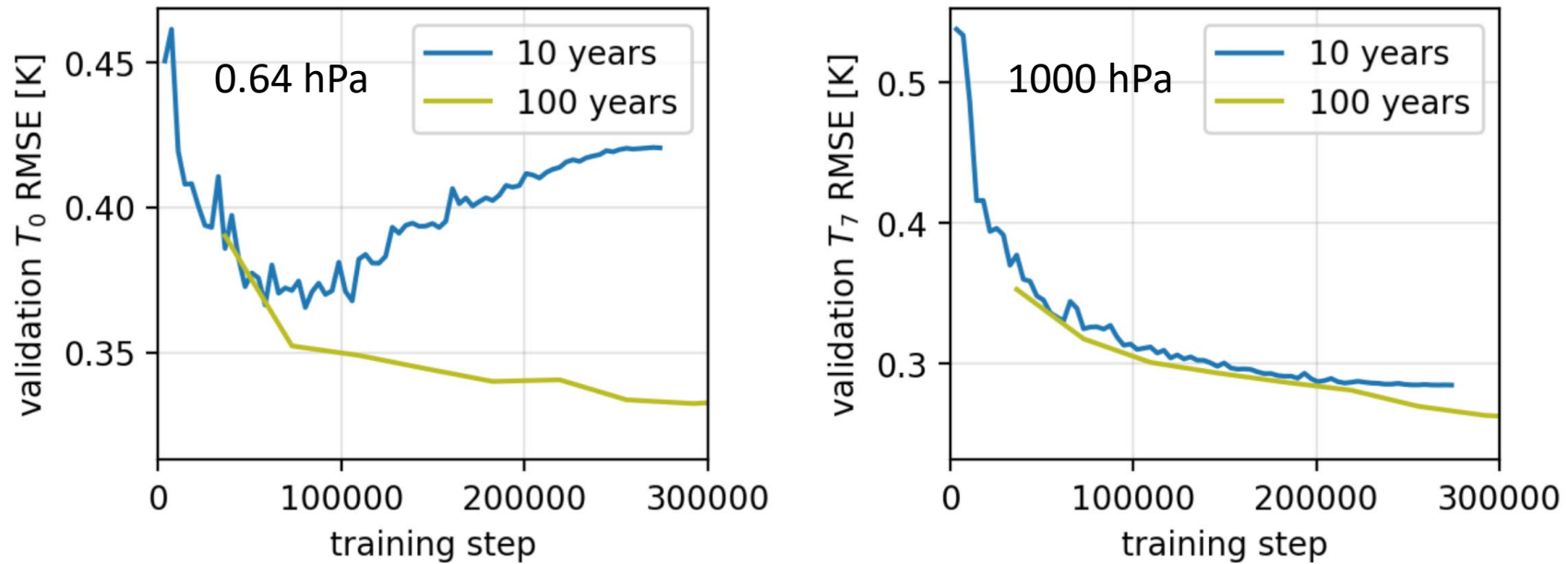
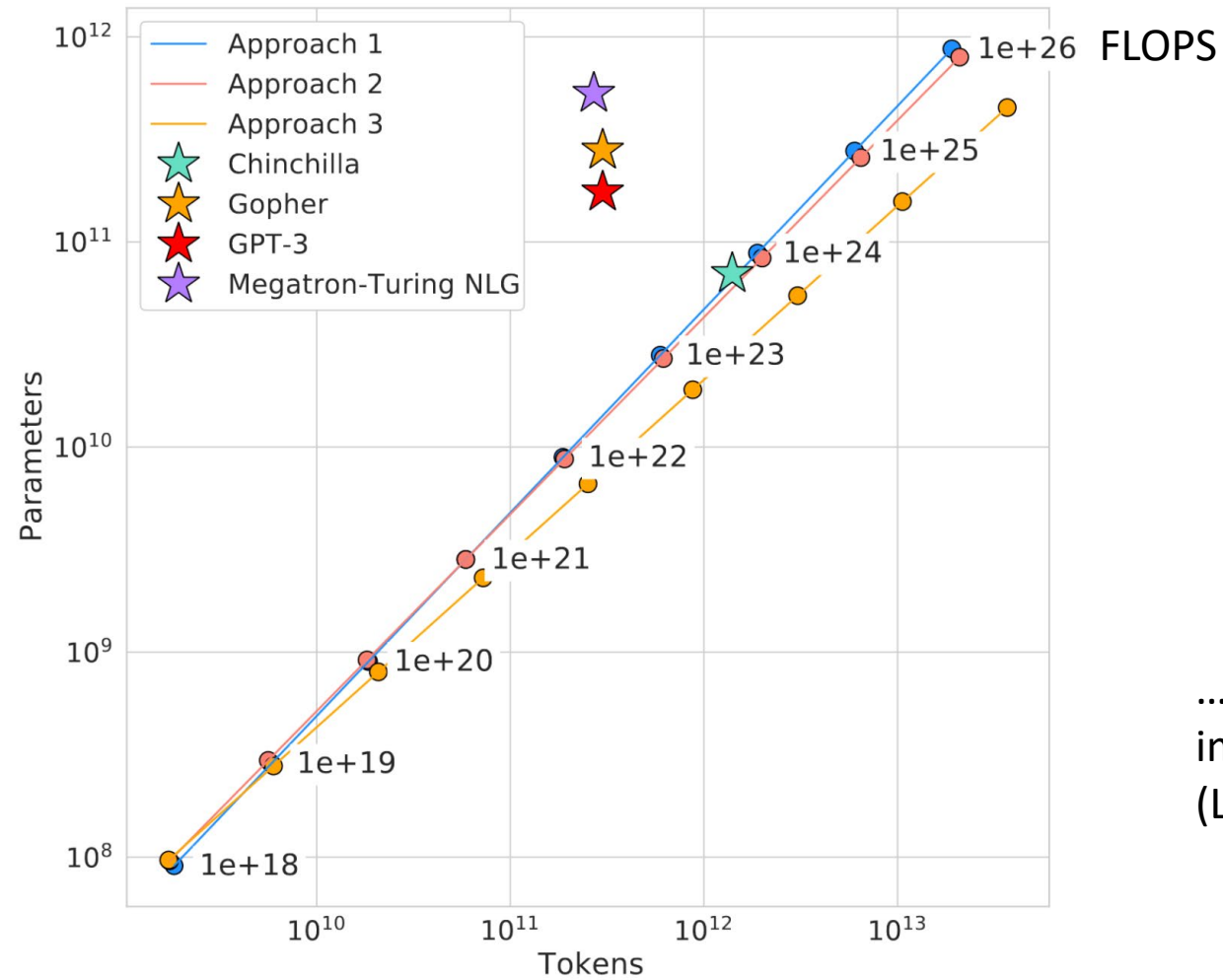


Figure 15: The validation RMSE of (left) T_0 and (right) T_7 for models trained on (blue) a 10-year dataset and (olive) a 100-year dataset.

Computational costs

„The computational cost of the model is negligible; it has a throughput of 256 ensemble members (at 2° resolution) per 3 min on Google Cloud TPUv3-32 instances and can easily scale to higher throughput by deploying more accelerators. [...] Training [of the 114 mio parameter model] takes slightly less than 18 hours on a 2 × 2 × 4 TPUv4 cluster.“ (Carver et al., 2024)

Scaling laws?



„Chinchilla“ optimal training; Hoffmann et al. ([Arxiv, 2021](https://arxiv.org/abs/2021.06.08))

... but: additional training improves performance!
(Llama3; Karpathy posts on X)

Training strategies (here: large language models)

Deciding on a model architecture

Deciding on a model parallelism strategy

Deciding on the model size

Scaling laws

Trade-off of large language model sizes

Issues and questions related to tensor precision

What to choose between fp32, fp16, bf16

Mixed-precisions for optimizers, weights, specific modules

How to finetune and integrate a model trained in a precision in another precision

Selecting training hyper-parameters and model initializations

Learning rate and learning rate schedules

Questions on batch size

Maximizing throughput

Avoiding, recovering from and understanding instabilities

Detecting instabilities early

Training tips to reduce instabilities

Issues with data and data processing

Debugging software and hardware failures

Tips on what metrics to follow during the training

Conclusions

- AI models will become the standard tools for weather forecasting at all scales
- AI models can produce excellent deterministic forecasts and quantify uncertainties well
- Some models (all?) exhibit some physical inconsistencies; these can likely be healed
- Larger models exhibit good robustness and (limited?) capabilities for extrapolation
- Tendency towards very large foundation models less clear than in language area (scaling laws?)
- Incorporating Earth system feedbacks on all time scales is probably the largest challenge ahead

Large-Scale Deep Learning for the Earth System

Bonn, Germany, 29 – 30 August 2024

[Registration](#) open until 31 July (student rate until 15 June)

Supported by

