# AI, Machine Learning & Visualisation at UKAEA

**5th IFERC Workshop on GPU Fusion Applications**
**20 June 2024**

**S.Pamela, R.Akers, N.Amorisco, N.Bhatia, D.Brennand, J.Buchanan, N.Carey, E.Crovini, O.El-Zobaidi, S.Etches, V.Gopakumar, S.Jackson, E.Lewis, E.Ozturk, K.Pentland, C.Siddle, L.Zanisi, J.Brandstetter, M.Hoelzl, G.Huijsmans**
**and many others…**

**Foreword: Only a small fraction of all activities (and co-authors)**

**- The various teams**

**- Overview of main activities**

**- Neural-Parareal**
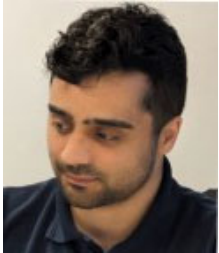
**- Foundation Models**

Foreword: Only a small fraction of all activities (and co-authors)
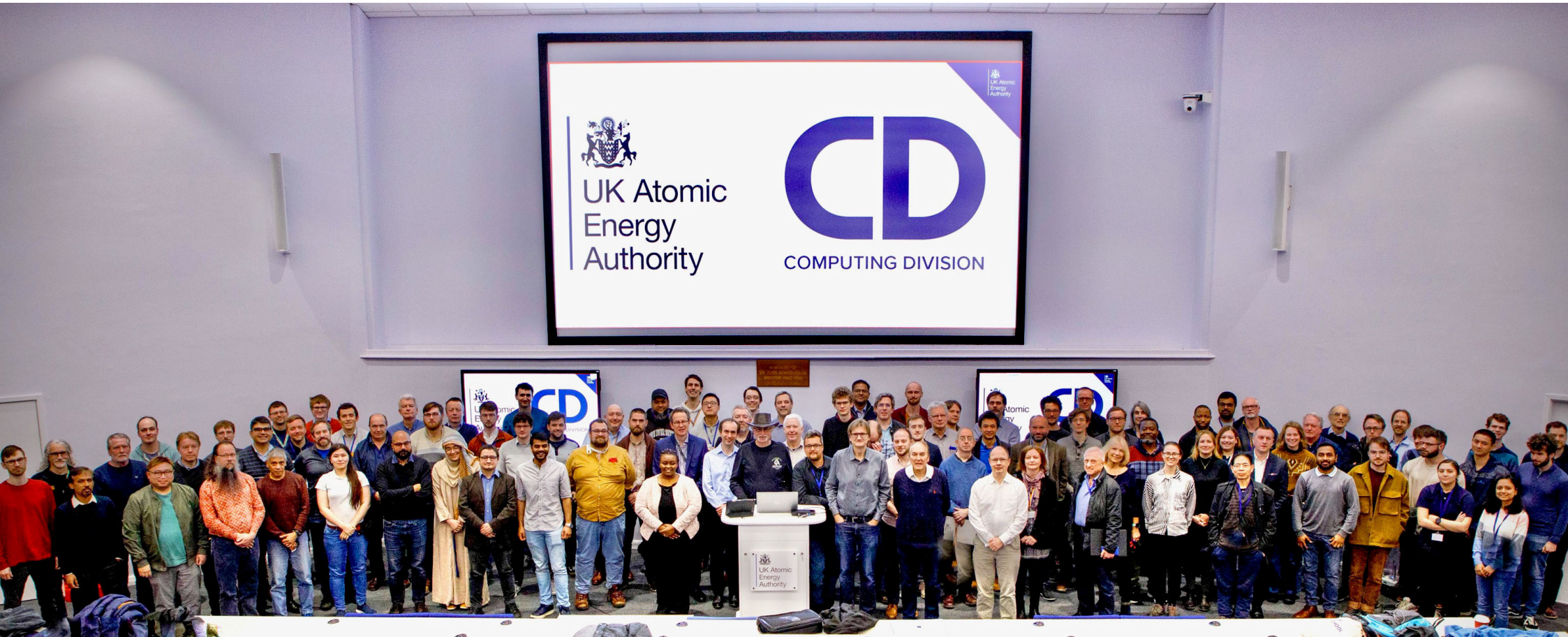
**- The various teams**

- Overview of main activities

- Neural-Parareal

- Foundation Models

# The UK Fusion Computing Lab

# STEP: Spherical Tokamak For Energy Production

**STEP: Deliver energy to the grid by 2040**

# Tokamak Design is an AI & Exascale Challenge

Cannot build 20 demonstration power plants
=> _Exascale and AI_ is needed to design & optimise STEP and future fusion power plants

- complex engineering
- in-silico design optimization
- model-based predictions with large uncertainty

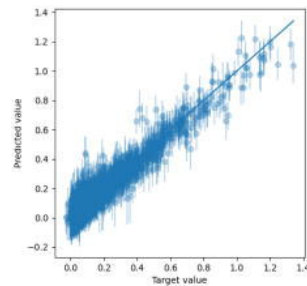Foreword: Only a small fraction of all activities (and co-authors)

- The various teams

- **Overview of main activities**

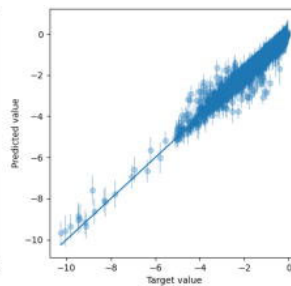- Neural-Parareal

- Foundation Models

**Will Hornsby, in collaboration with Digilab ltd**
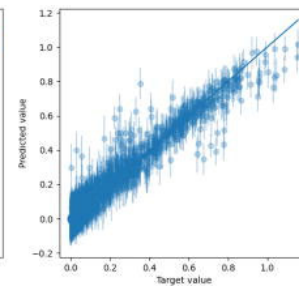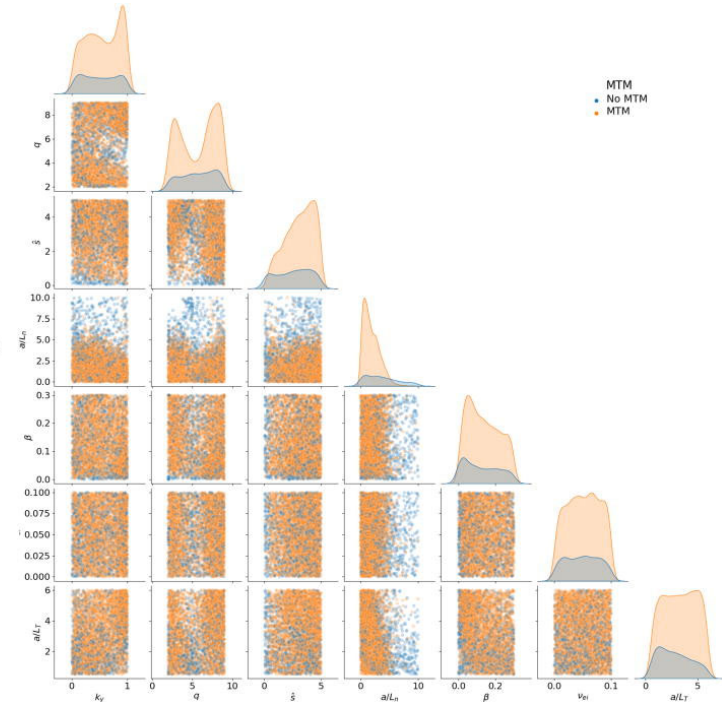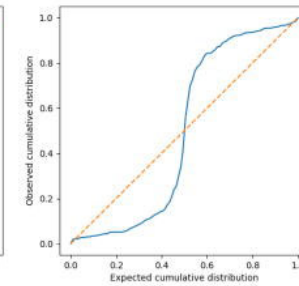  **Using Gaussian Processes to emulate MTM stability with GS2**
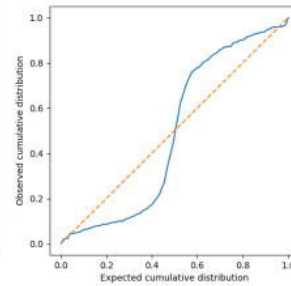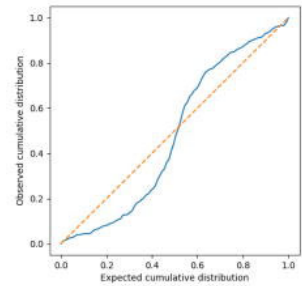  **W.Hornsby et al., *Phys. Plasmas 31, 012303 (2024)***
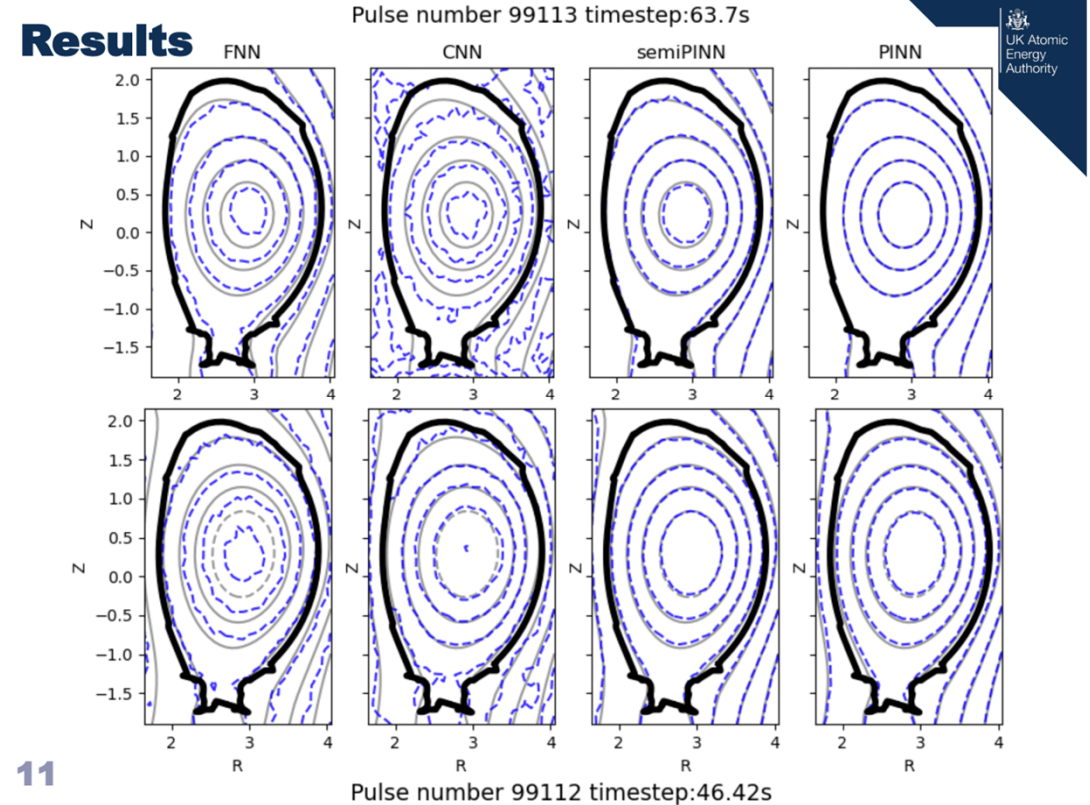


(a) Growth Rate.  (b) Mode Frequency.  (c) Electron Heat Flux.

**Nico Amorisco (UKAEA)**
**Steve Etches (UKAEA)**
**Emily Lewis (UCL PhD)**
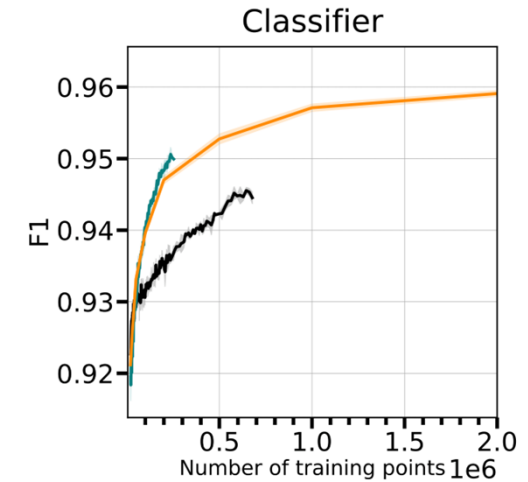**Omar El-Zobaidi (Placement)**



*N.Amorisco et al., "FreeGSNKE: A Python-based dynamic free-boundary toroidal plasma equilibrium solver", Phys. Plasmas 31, 042517 (2024)*



**Emily Lewis, PhD at UCL:**
**Surrogate of plasma equilibrium**

**Lorenzo Zanisi (UKAEA)**
**Enrico Crovini (Imperial PhD)**
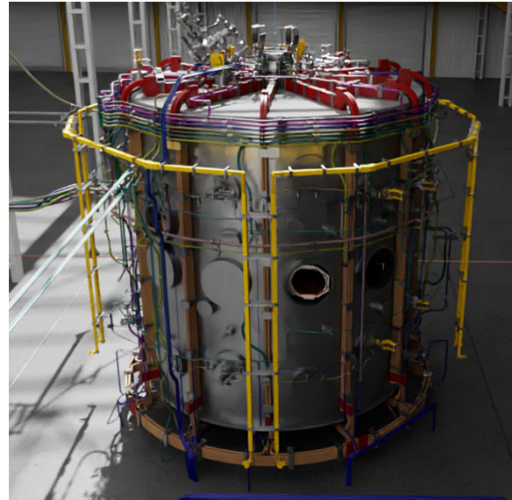**Theo Brown (UCL PhD)**
**Catherine Siddle (Grad-Scheme)**



**Active Learning for Qualikiz, L. Zanisi et al 2024 Nucl. Fusion 64 036022**

**E.Crovini et al. "Automatic JOREK calibration via batch Bayesian optimization", Physics of Plasmas 31, 063901 (2024)**
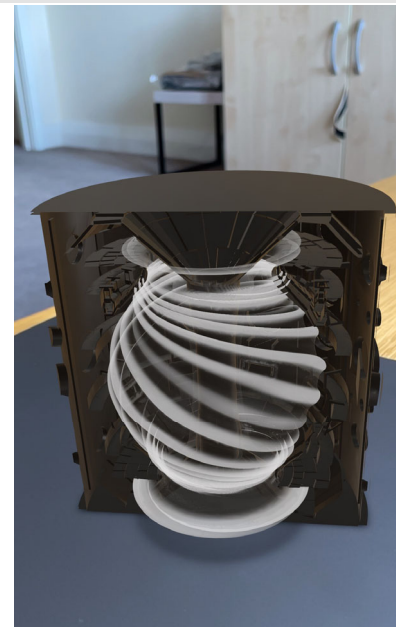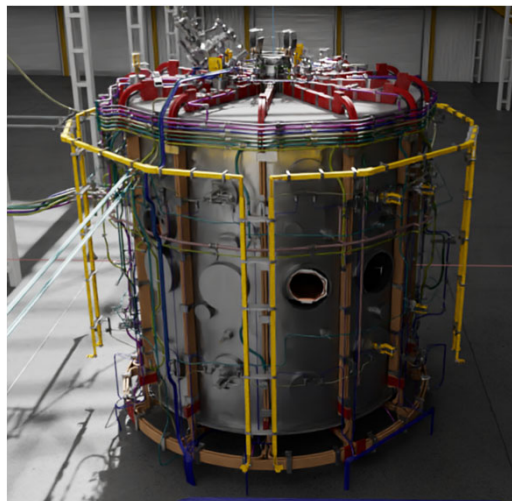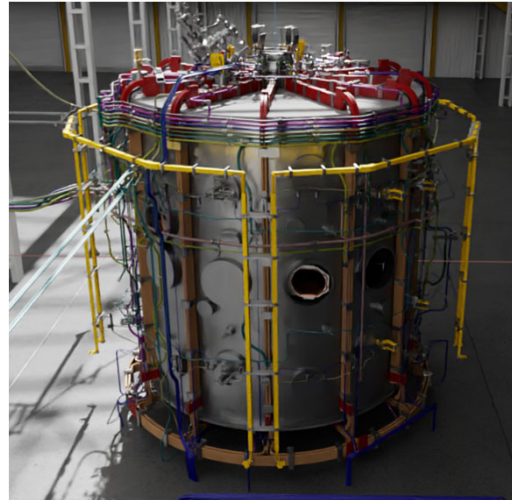
**Nitesh Bhatia (UKAEA)**
**Ekin Ozturk (Imperial PhD)**

**Nitesh Bhatia (UKAEA)**
**Ekin Ozturk (Imperial PhD)**



https://niteshbhatia008.github.io/nb-webxr-viewer/

**Nitesh Bhatia (UKAEA)**
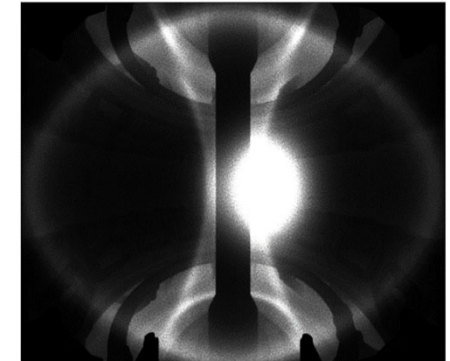**Ekin Ozturk (Imperial PhD)**
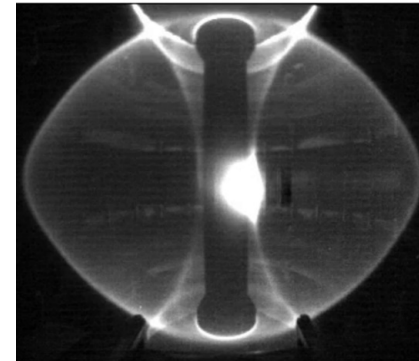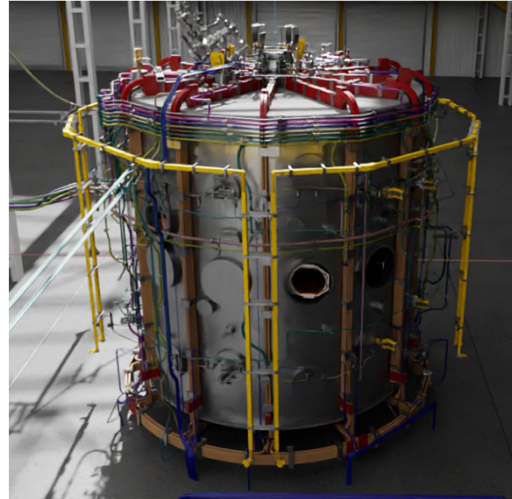
https://niteshbhatia008.github.io/nb-webxr-viewer/

**Nitesh Bhatia (UKAEA)**
**Ekin Ozturk (Imperial PhD)**

https://niteshbhatia008.github.io/nb-webxr-viewer/

**Vignesh Gopakumar (UKAEA)**
**Naomi Carey (UKAEA Apprenticeship)**
**Daniel Brennand (UKAEA Apprenticeship)**

**STORM/BOUT++**



**JOREK**

*V.Gopakumar et al. 2024 Nucl. Fusion 64 056025*
*N.Carey et al., IAEA-FEC 2023*

# AI activities at UKAEA: Foundation Models

Samuel Jackson (UKAEA)
Vignesh Gopakumar (UKAEA)
Naomi Carey (UKAEA Apprenticeship)
Lorenzo Zanisi (UKAEA)
Nathan Cummings (UKAEA)
Johannes Brandstetter (JKU)



**Main collaborations:**

- Linz
- Turing Institute
- IBM

Foreword: Only a small fraction of all activities (and co-authors)

- The various teams

- Overview of main activities

**- Neural-Parareal**

- Foundation Models

**Split time domain into parallel windows** [ J.-L. Lions et al., Comptes Rendus de l'Académie des Sciences, Série I. 332 (7): 661–668 (2015) ]

> fast approximation with *Coarse-Solver*
> correction using Fine-Solver

**Basics of Parareal: split time domain into parallel windows**
    **fast approximation with *Coarse-Solver***
    **correction using Fine-Solver**
    **better coarse solver => few iterations => high speedup**

**Basics of Parareal: split time domain into parallel windows**

       **fast approximation with *Coarse-Solver***

       **correction using Fine-Solver**

       **better coarse solver => few iterations => high speedup**

           **=> use Machine Learning surrogates**



**Top:**       **ground truth**
**Bottom:**    **neural PDE solver**

# Neural-Parareal: JOREK Demonstration

**JOREK: non-linear MHD solver for tokamak plasmas [ _jorek.eu_ ]**
**Comes with simplified models (basically NS in toroidal geometry, and Hasegawa-Wakatani)**
**Full simulations address things like plasma-edge filamentation, or disruption (loss of plasma control)**

**Blobs with 3-variables model (Navier Stokes in torus)**

- $\rho$, $T$, $\Phi$ (stream function)
- Plus 1 auxiliary variable: vorticity $w = \nabla^2 \Phi$

**Radially motion due to** $\nabla p$ **and toroidal geometry**

**Hotter blobs move faster**

**Blobs with Reduced-MHD model**
**Used extensively for fusion**

- 4 variables: $\rho$, $T$, $\Phi$, $\psi$ (magnetic potential)
- Plus 2 auxiliary variables:
    - Vorticity $w = \nabla^2\Phi$
    - Current $j = \nabla^2\psi$

**Effectively 6 variables**
**Behaviour is quite different**

Top:     electrostatic model
Bottom: electromagnetic model (RMHD)

**Blobs with Reduced-MHD model**
**Used extensively for fusion**

- **4 variables: $\rho$, $T$, $\Phi$, $\psi$ (magnetic potential)**
- **Plus 2 auxiliary variables:**
    - **Vorticity $w = \nabla^2 \Phi$**
    - **Current $j = \nabla^2 \psi$**

**Effectively 6 variables**
**Behaviour is quite different**

      **=> Blobs create their own internal current**
      **=> which in turn affects velocity**

**In theory, Parareal is "non-intrusive"**
**In practice, it requires a lot of work with the code's i/o**
**For FEM code, even more complex due to projection between resolutions**

**Except that a Neural PDE solver requires several input timesteps**

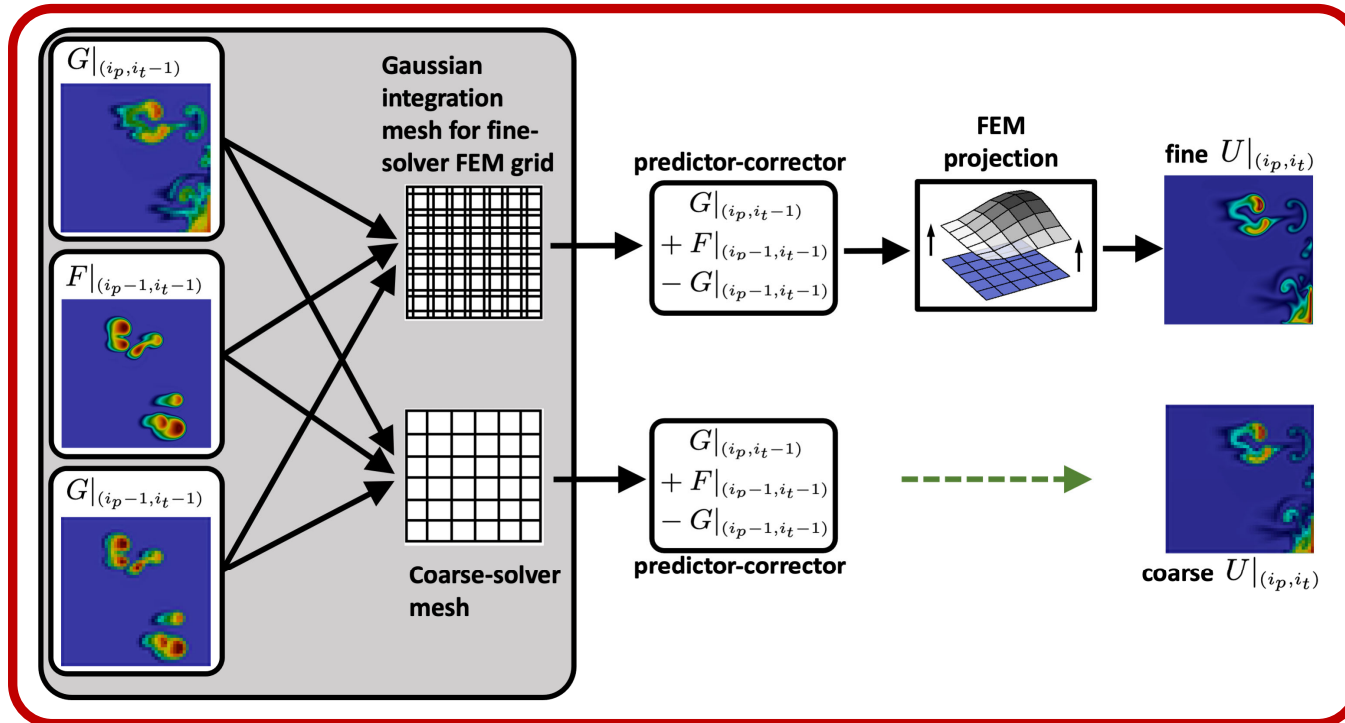   **=> Need to apply the predictor-corrector to many timesteps**

   **=> Even more i/o**

   **=> Projections are costly => needs to be parallelised (otherwise can easily dominate workflow)**

   **=> end up with a lot of extra data!**

# Neural-Parareal: Parareal Demonstration

**It works!**

**Looking at difference of last timestep with ground truth
(top is ground truth, bottom is parareal evolution)**

**Top:**        ground truth
**Bottom:**       Parareal evolution (last timestep)

# Neural-Parareal: SSIM Measure

Using SSIM (Structural Similarity Index Measure)
Better than MSE for generic structures of blobs
"SSIM = 1" means 100% accuracy

No matter how bad your coarse solver, Parareal will
always converge to SSIM=1 at final iteration

**Top:** Parareal evolution (last time-step)
**Bottom:** Corresponding SSIM evolution

"fake" coarse solver == JOREK itself, but with controlled difference
- **Exactly same physics model**
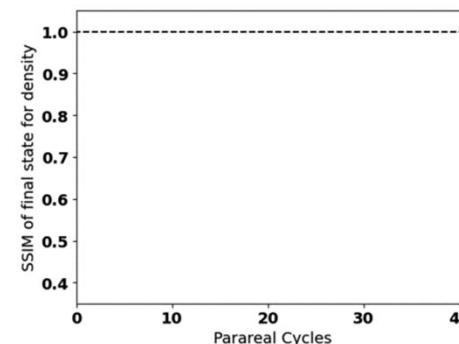- **Lower spatial resolution (half)**
- **Higher diffusion**
    - **=> Diff x 30 is a bad coarse solver**
    - **=> Diff x 3 is a good coarse solver**

**Neural Coarse solver gives better performance than Diff x 3**
        **=> it works really well**

Since Parareal simulations produce a lot of data, train NN as we launch more simulations

# Neural-Parareal: Main Result

**Ran 5 batches of 20 simulations**
**Speed-up increases significantly and quickly**



Speedup efficiency with an SSIM requirement of 0.7

Speedup efficiency with an SSIM requirement of 0.8

# GPU usage and scale

- **Moving away from CPU simulations**
  - **=> Many fusion codes still require CPU's**
  - **=> Need to move towards GPU codes (NekRS, CGYRO, XGC, OpenMC)**
  - **=> Strong collaboration with US-NL (Exascale Computing Project ECP)**



*Andy Davis & Nitesh Bhatia*
*NekRS running on Frontier*

- **Moving away from CPU simulations**
    - => **Many fusion codes still require CPU's**
    - => **Need to move towards GPU codes (NekRS, CGYRO, XGC, OpenMC)**
    - => **Strong collaboration with US-NL (Exascale Computing Project ECP)**

- **Transformers and Foundation models**
    - => **Most ML workflows are not GPU intensive**
    - => **Transformers are**
    - => **In visit at Linz (JKU) this week to learn from world experts**
    - => **Collaborations with JKU, Turing Institute, IBM, IAEA**
    - => **Strength of Transformers: they scale very well on GPUs**
    - => **Challenge of Transformers: data-hungry**



arXiv:2405.13063v2 [physics.ao-ph] 28 May 2024

## AURORA: A FOUNDATION MODEL OF THE ATMOSPHERE

Cristian Bodnar[*,1], Wessel P. Bruinsma[*,1], Ana Lucic[*,1], Megan Stanley[*,1], Johannes Brandstetter[3,†], Patrick Garvan[1], Maik Riechert[1], Jonathan Weyn[2], Haiyu Dong[2], Anna Vaughan[4], Jayesh K. Gupta[5,†], Kit Tambiratnam[2], Alex Archibald[4], Elizabeth Heider[1], Max Welling[6,†], Richard E. Turner[1,4], and Paris Perdikaris[1]

[1]Microsoft Research AI for Science
[2]Microsoft Corporation  [3]JKU Linz  [4]University of Cambridge  [5]Poly Corporation  [6]University of Amsterdam

[*]Equal contribution  [†]Work done while at Microsoft Research

### ABSTRACT

Deep learning foundation models are revolutionizing many facets of science by leveraging vast amounts of data to learn general-purpose representations that can be adapted to tackle diverse downstream tasks. Foundation models hold the promise to also transform our ability to model our planet and its subsystems by exploiting the vast expanse of Earth system data. Here we introduce Aurora, a large-scale foundation model of the atmosphere trained on over a million hours of diverse weather and climate data. Aurora leverages the strengths of the foundation modelling approach to produce operational forecasts for a wide variety of atmospheric prediction problems, including those with limited training data, heterogeneous variables, and extreme events. In under a minute, Aurora produces 5-day global air pollution predictions and 10-day high-resolution weather forecasts that outperform state-of-art classical simulation tools and the best specialized deep learning models. Taken together, these results indicate that foundation models can transform environmental forecasting.

## 1 Introduction

Deep learning foundation models have revolutionised various scientific domains, such as protein structure prediction (Abramson et al., 2024), drug discovery (Chithrananda et al., 2020), computer vision (Betker et al., 2023), and natural language processing (OpenAI, 2024). The key tenets of foundation models include *pretraining*, where a single large-scale neural network learns to capture intricate patterns and structure from a large corpus of diverse data; and *fine-tuning*, which allows the model to leverage its learned representations to excel at new tasks with limited training data (Bommasani et al., 2021; Brown et al., 2020).

The Earth system is a complex and interconnected network of subsystems, such as the atmosphere, oceans, land, and ice, which constantly interact in intricate ways. In a rapidly changing climate, accurate understanding of these subsystems becomes increasingly important. We envision that foundation models can revolutionise our ability to model subsystems of the Earth, and eventually the whole Earth.

Amongst the Earth's subsystems, the atmosphere stands out as particularly data-rich (Reichstein et al., 2019; Bauer et al., 2015) and therefore constitutes ripe ground for pretraining a foundation model. Classical atmospheric simulation approaches, such as numerical weather prediction (NWP), are costly and unable to exploit this wealth of data (Bauer et al., 2015). Recent deep learning approaches are cheaper, more flexible, and have shown great promise in specific prediction tasks with abundant training data (Lam et al., 2023; Bi et al., 2023; Chen et al., 2023a,b; Han et al., 2024; Kochkov et al., 2024; Lessig et al., 2023; Pathak et al., 2022; Bonev et al., 2023; Andrychowicz et al., 2023; Ham et al., 2019; Nguyen et al., 2023a,b). However, these methods struggle when atmospheric training data are scarce (Chantry et al., 2021) or heterogeneous (Reichstein et al., 2019), and they lack robustness in predicting extremes (Charlton-Perez et al., 2024). By learning generalizable representations from vast amounts of diverse data, foundation models have been able to overcome analogous challenges in other domains (Zhai et al., 2022; Radford et al., 2021; Bommasani et al., 2021; Nguyen et al., 2023a).

# GPU usage and scale

- **Moving away from CPU simulations**
  - **=> Many fusion codes still require CPU's**
  - **=> Need to move towards GPU codes (NekRS, CGYRO, XGC, OpenMC)**
  - **=> Strong collaboration with US-NL (Exascale Computing Project ECP)**

- **Transformers and Foundation models**
  - **=> Most ML workflows are not GPU intensive**
  - **=> Transformers are**
  - **=> In visit at Linz (JKU) this week to learn from world experts**
  - **=> Collaborations with JKU, Turing Institute, IBM, IAEA**
  - **=> Strength of Transformers: they scale very well on GPUs**
  - **=> Challenge of Transformers: data-hungry**

- **Several GPU clusters => portability is key**
  - **Leonardo (CINECA) Nvidia**
  - **CSD3 (Cambridge) Nvidia**
  - **Dawn (Cambridge) Intel**
  - **Isembard-AI (Bristol) Nvidia**
  - **LUMI (Finland) AMD**
  - **Frontier (USA) AMD**

# Conclusion

- **AI & Machine Learning at UKAEA has ramped up over last 3 years**

- **Several projects running in collaboration with internal/external partners**

- **Integration into larger framework (eg. Neural-Parareal, JINTRAC)**

- **Currently on visit at Linz with Johannes Brandstetter**
    - **=> Starting Transformers and Foundation Models**
    - **=> will increase GPU usage a lot in near future**


**Thank you for your attention!**