# Going Global

Taking code from **research**
to **operational** open ecosystem
for AI **weather forecasting**

Dr. Jesper Dramsch

Thanks to all contributors to AIFS and Anemoi

Rilwan Adewoyin, Mihai Alexe, Zied Ben Bouallègue, Matthew Chantry, Mariana Clare, Harrison Cook, Jesper Dramsch, Joffrey Dumont Le Brazidec, Rachel Furner, Vera Gahlen, Sara Hahner, Aaron Hopkinson, Gareth Jones, Simon Lang, Christian Lessig, Martin Leutbecher, Linus Magnusson, Michael Maier-Gerber, Gert Mertes, Gabriel Moldovan, Ana Prieto Nemesio, Cathal O'Brien, Florian Pinault, Ewan Pinnington, Jan Polster, Baudouin Raoult, Nina Raoult, Mario Santa Cruz, Jakob Schloer, Maria Luisa Taccari, Helen Theissen, Steffen Tietsche, Lorenzo Zampieri

Developed and used by meteorological centers across Europe.

AEMET, DWD, FMI, GeoSphere, KNMI, MET Norway, Meteo Swiss, Meteo France, RMI, & ECMWF

# Goals for this talk

| | | | |
|---|---|---|---|
| Growing software projects | Anticipating user needs | Learn about weather forecasting | AI Cautionary Tales |
| Upskilling From Coder to Software Architect | Planning for features you can't even know about | Growing a team to significant size | Have fun |

# Who am I?

Scientist for Machine Learning in Weather Forecasting

PhD in Machine Learning for Geoscience

Python, Data, AI Education

Maintain Pythondeadlin.es, ML.recipes, data-science-gui.de

Generally loud online

RADIOSONDES

SATELLITES

Humidity

800M
DAILY
OBSERVATIONS

Snow cover

AIRCRAFT

Trace gases and aerosols

RADARS

Pressure

Temperature

Precipitation

SYNOP

SHIP

BUOYS

Sea ice

Surface temperature

Waves

Wind

Vegetation

Sea surface temperature

Soil moisture

ECMWF

# ECMWF EARTH SYSTEM APPROACH

SUN

ATMOSPHERE

Turbulence

Solar radiation

Sea-ice atmosphere coupling

Sea-ice ocean coupling

OCEAN

Wind stress

Terrestrial radiation

Trace gases and aerosols

Evaporation

Human influences

Heat exchange

Precipitation

LAND

Land-atmosphere coupling

Forecast:
~6,400 CPUs in 30 Minutes

Data archive:
Over 1 Exabyte

What will machine learning for weather and climate predictions look like in 10 years from now?

Just Two Years Ago!

Machine learning will have no long-term effect

Machine learning will replace conventional models

Observation screening

Simple post-processing applications

Feature detection in model output

Bias correction in data assimilation

Emulation of parametrisation schemes

Learn model components from observations

Learn equations of motion

**The uncertainty range is still very large...**

Sets of training data from ERA5

Example of training loop

ML model

INPUT
OUTPUT

Checks accuracy against output

Corrects errors to improve accuracy

Each step of the training loop uses several ERA5 states

Predicts weather based on physical state of earth after learning from ERA5

Data assimilation | Initial conditions | Modelling | Predictions

Observation

Machine learning components in parts of the chain

Entirely machine learning

Training:
64 GPUs in 8 Days

Forecast:
1 GPU in 2 Minutes

# How did we end up here?

# Initial Commit

**ECMWF** EUROPEAN CENTRE FOR MEDIUM-RANGE WEATHER FORECASTS

# Typical research ML code

# It started out with a wish ✨



Configurations

Refactors

Inference Mode

# User Base

| Modifies Configs for Experimentation and Improvement of Anemoi Model | Modifies Codebase to implement new Features and Augment Anemoi Libraries | Runs the Anemoi Model in a common interface on reliable infrastructure |

**ECMWF**

# Keeping the User and Collaboration in Mind

**Researchers**

**Developers**

**Operations**

- Researchers
  - Quick switching of experiment values
  - Less interaction with core ML code
  - Low-key experiment tracking
- Developers
  - Modularity and Extensibility
  - Code quality
  - Separation of Concerns
- Operations
  - Minimal Dependencies
  - Consistent interfaces

# How do we facilitate collaboration?

Reading

Bonn

Bologna

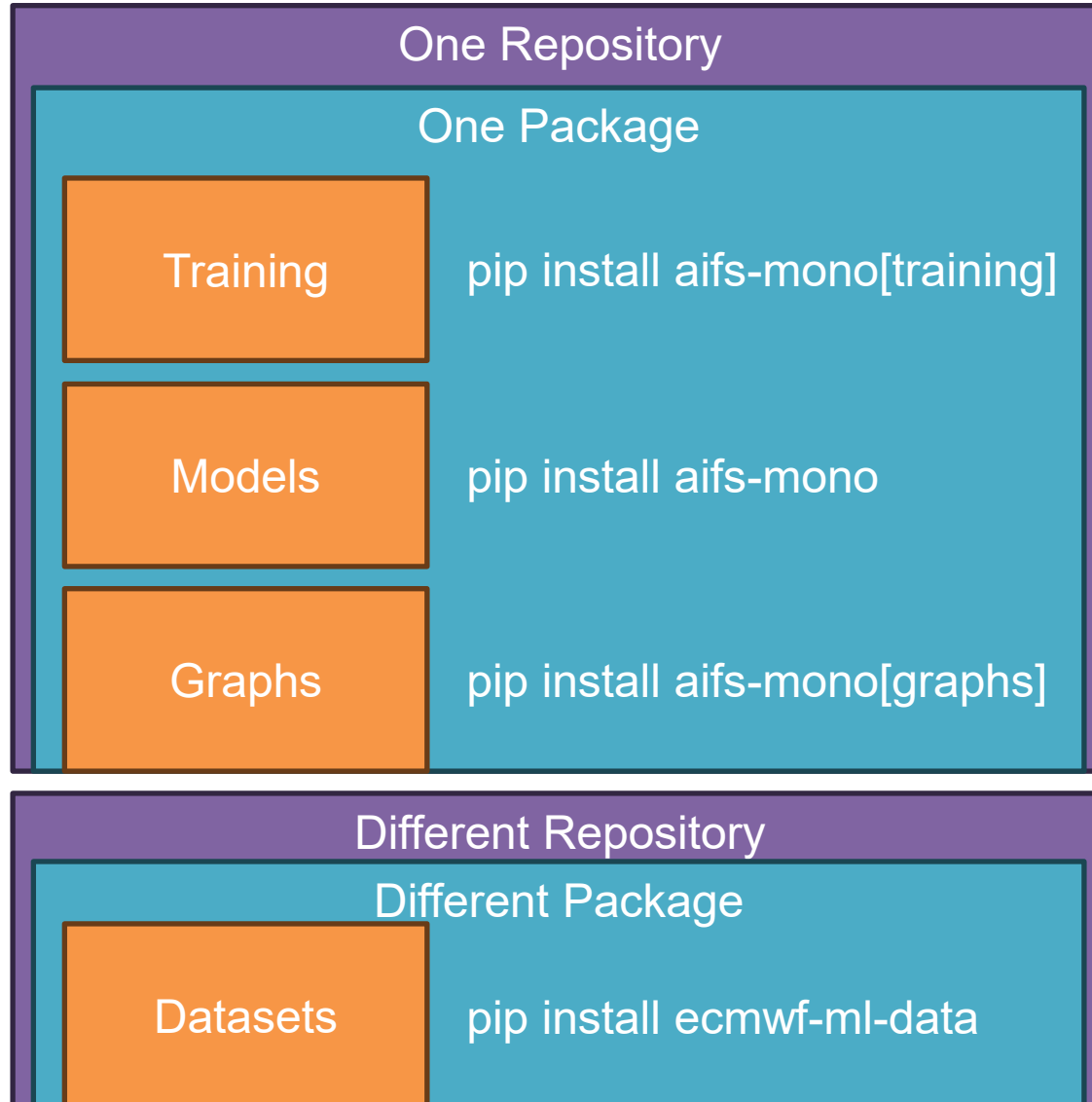# Modular and Extensible AIFS Trainer code



**Hierarchical Configs with Hydra**

**Tidier Separation**

**Modular Trainer Class**

# AIFS mono-package + External dependencies

## One Repository

### One Package

| | |
|---|---|
| **Training** | pip install aifs-mono[training] |
| **Models** | pip install aifs-mono |
| **Graphs** | pip install aifs-mono[graphs] |

## Different Repository

### Different Package

| | |
|---|---|
| **Datasets** | pip install ecmwf-ml-data |

- Pros
  - ▲ Convenient to develop
  - ▲ All code in one place
  - ▲ Quick to release

- Cons
  - ▼ Weird to install
  - ▼ Not all code in „ecosystem"
  - ▼ Complex
  - ▼ No unified testing or infrastructure

# Focusing on Configurability, Extensibility and Modularity

```
Model / GNN.yml
activation: GELU
num_channels: 512

model:
 _target_:
anemoi.models.models.encoder_processor_decoder.
AnemoiModelEncProcDec

processor:
 _target_:
anemoi.models.layers.processor.GNNProcessor
 _convert_: all
 activation: ${model.activation}
 trainable_size:
${model.trainable_parameters.hidden2hidden}
 sub_graph_edge_attributes:
${model.attributes.edges}
 num_layers: 16
 num_chunks: 2
 mlp_extra_layers: 0

encoder:
 _target_:
anemoi.models.layers.mapper.GNNForwardMapper
```

```
Model / GraphTransformer.yml
activation: GELU
num_channels: 1024

model:
 _target_:
anemoi.models.models.encoder_processor_decoder.
AnemoiModelEncProcDec

processor:
 _target_:
anemoi.models.layers.processor.GraphTransformerPr
ocessor
 _convert_: all
 activation: ${model.activation}
 trainable_size:
${model.trainable_parameters.hidden2hidden}
 sub_graph_edge_attributes:
${model.attributes.edges}
 num_layers: 16
 num_chunks: 2
 mlp_hidden_ratio: 4 # GraphTransformer
 num_heads: 16 # GraphTransformer

encoder:
```

```
Model / Transformer.yml
activation: GELU
num_channels: 1024

model:
 _target_:
anemoi.models.models.encoder_processor_decoder.
AnemoiModelEncProcDec

processor:
 _target_:
anemoi.models.layers.processor.TransformerProcesso
r
 _convert_: all
 activation: ${model.activation}
 num_layers: 16
 num_chunks: 2
 mlp_hidden_ratio: 4 # Transformer only
 num_heads: 16 # Transformer only
 window_size: 512
 dropout_p: 0.0

encoder:
 _target_:
anemoi.models.layers.mapper.GraphTransformerFor
```
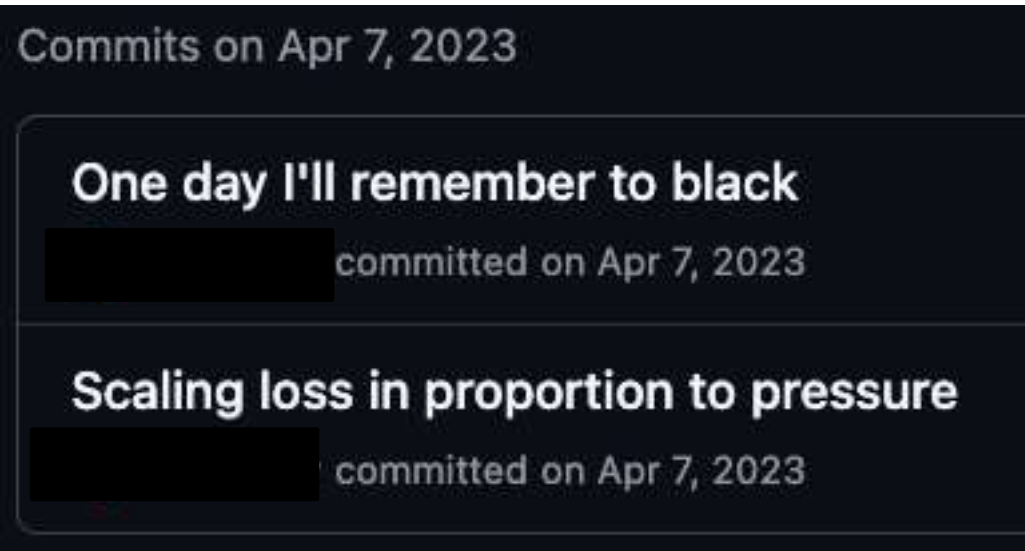
> anemoi-training train model=gnn      > anemoi-training train model=graphtransformer      > anemoi-training train model=transformer

- Make it easy to switch components with hydra

- Full tracking even on terminal overrides

- Easy to extend with new models and components

refactored Simon's new transformer
experiment tracking with Ana and Sara

# Why we use pre-commit hooks



Commits on Apr 7, 2023

One day I'll remember to black
committed on Apr 7, 2023

Scaling loss in proportion to pressure
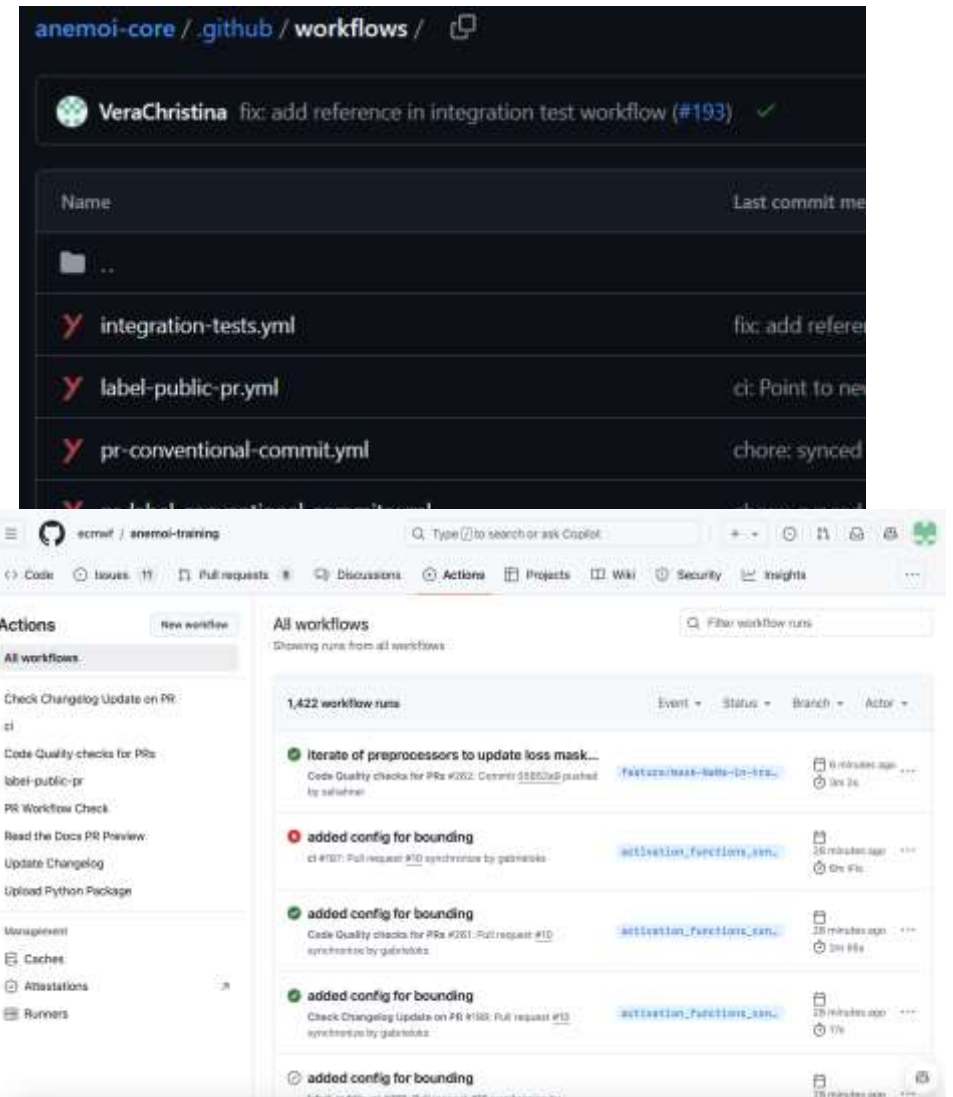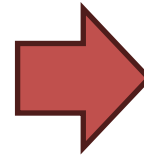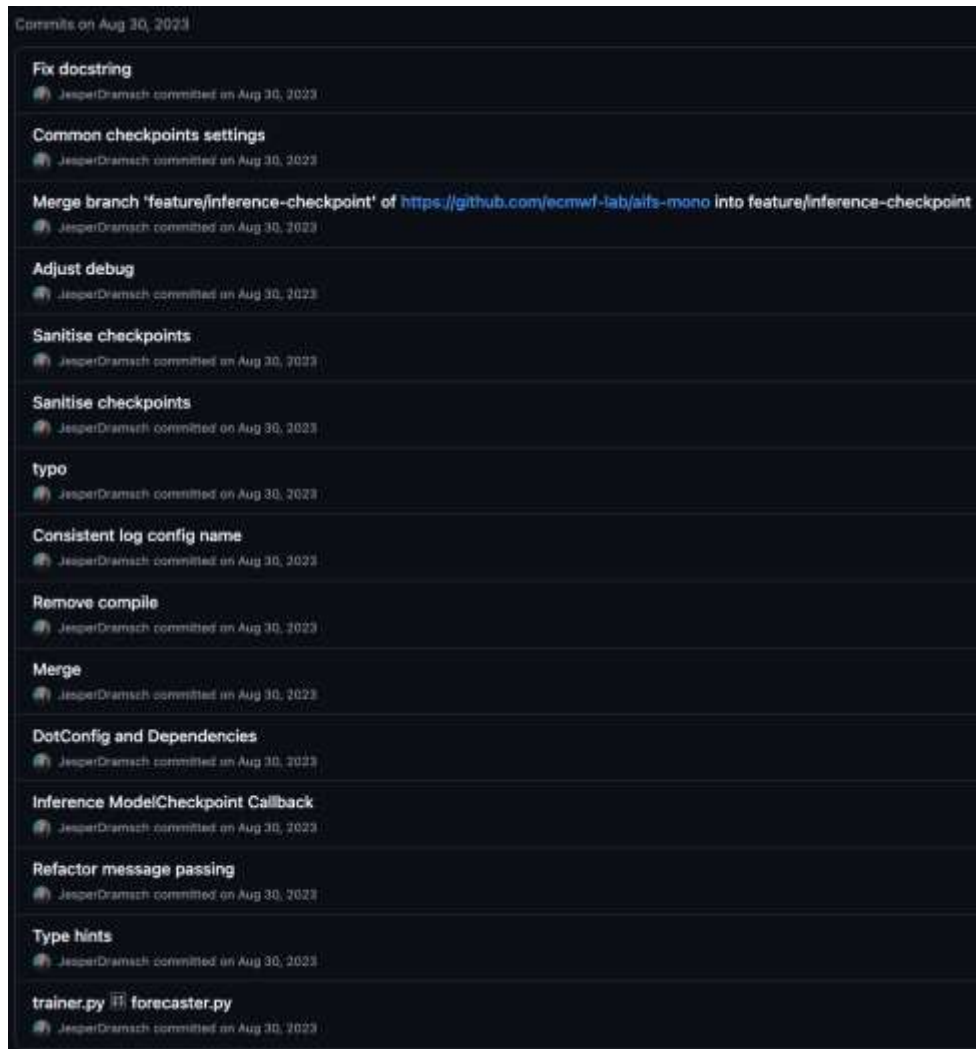committed on Apr 7, 2023

anemoi-training / .pre-commit-config.yaml

pre-commit-ci[bot] and gmertes [pre-commit.ci] pre-commit autoupdate (#177)
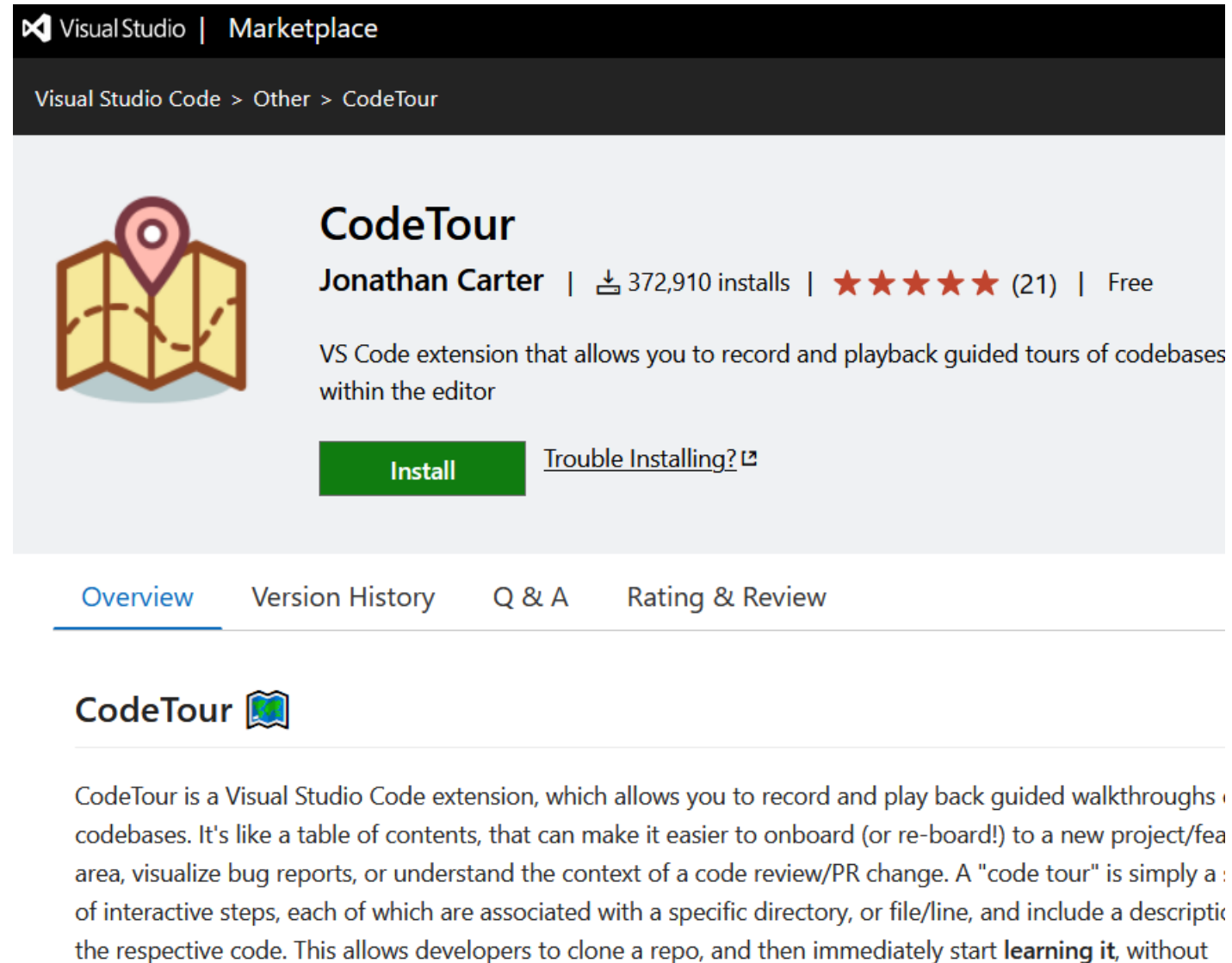
Code | Blame    78 lines (78 loc) · 2.54 KB

```
1    repos:
2    # Empty notebookds
3    - repo: local
4      hooks:
5      - id: clear-notebooks-output
6        name: clear-notebooks-output
7        files: tools/.*\.ipynb$
8        stages: [pre-commit]
9        language: python
10       entry: jupyter nbconvert --ClearOutputPreprocessor.enabled=True --inplace
11       additional_dependencies: [jupyter]
12   - repo: https://github.com/pre-commit/pre-commit-hooks
13     rev: v5.0.0
14     hooks:
15     - id: check-yaml # Check YAML files for syntax errors only
16       args: [--unsafe, --allow-multiple-documents]
17     - id: debug-statements # Check for debugger imports and py37+ breakpoint()
18     - id: end-of-file-fixer # Ensure files end in a newline
19     - id: trailing-whitespace # Trailing whitespace checker
20     - id: no-commit-to-branch # Prevent committing to main / master
```

# Why we do squash commits now 👀

# Do Code Tours!

- Shortcuts learning (multi-week -> 1 hour)

- Puts design decisions in context

- Gives New Hires space to ask questions

- Can be tailored to specific needs

- Establishes collaborative aspects

- Optionally: use tools

Visual Studio | Marketplace

Visual Studio Code > Other > CodeTour

## CodeTour

Jonathan Carter | ⤓ 372,910 installs | ★★★★★ (21) | Free

VS Code extension that allows you to record and playback guided tours of codebases within the editor

**Install**    Trouble Installing? ↗

Overview    Version History    Q & A    Rating & Review

## CodeTour 🗺️

CodeTour is a Visual Studio Code extension, which allows you to record and play back guided walkthroughs codebases. It's like a table of contents, that can make it easier to onboard (or re-board!) to a new project/fea area, visualize bug reports, or understand the context of a code review/PR change. A "code tour" is simply a of interactive steps, each of which are associated with a specific directory, or file/line, and include a descriptio the respective code. This allows developers to clone a repo, and then immediately start **learning it**, without

# Growing beyond a single git repo

# The Vision of AIFS becoming the Anemoi Ecosystem
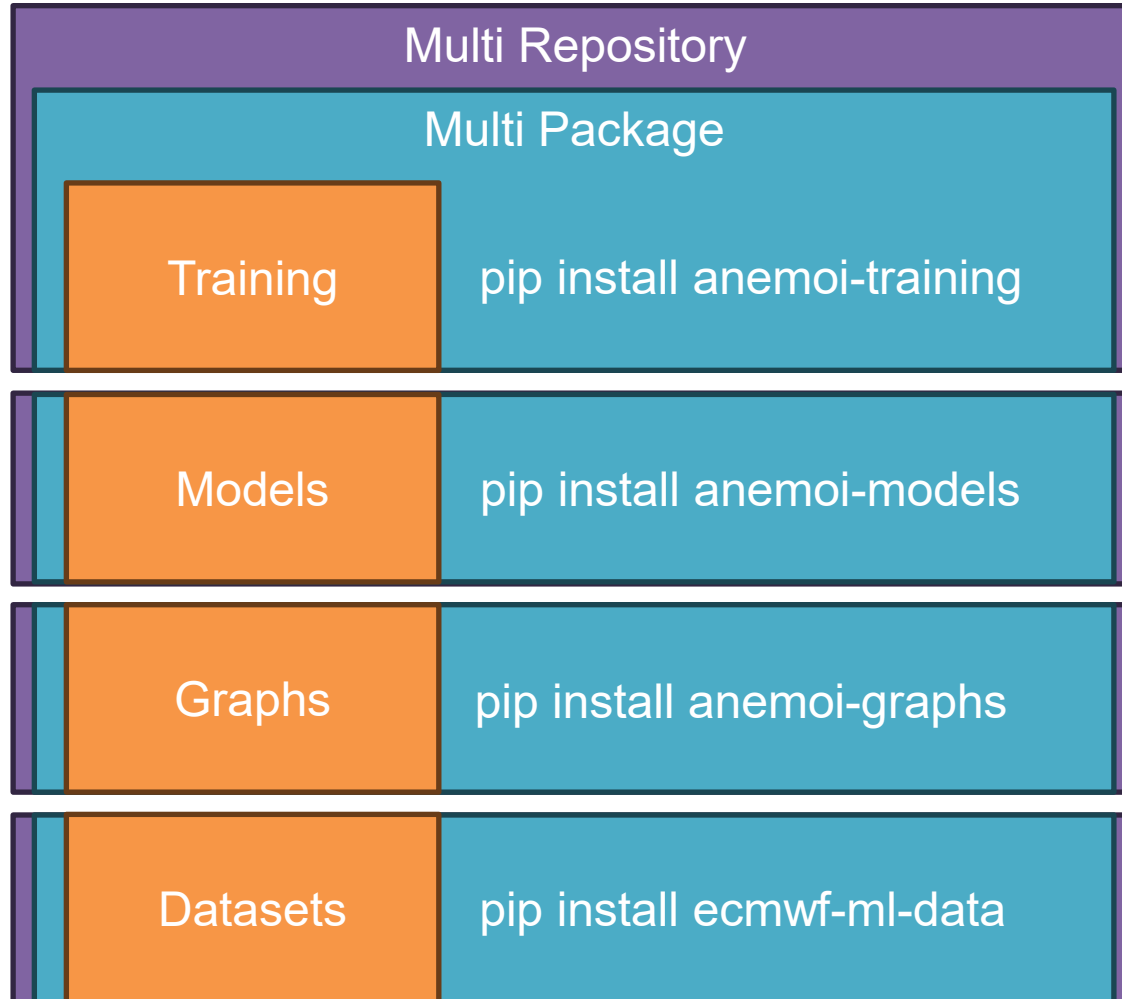
# Common interface and modularity

# Anemoi Ecosystem (multi repo)

| Multi Repository | |
|---|---|
| **Multi Package** | |
| Training | pip install anemoi-training |
| Models | pip install anemoi-models |
| Graphs | pip install anemoi-graphs |
| Datasets | pip install ecmwf-ml-data |

- Pros
  - ▲ Separation of concerns
  - ▲ Neat repo structure at root
  - ▲ Intuitive to Release individually
  - ▲ Complexity contained within repo
  - ▲ Easy to delegate responsibility
  - ▲ Different styles of collaboration possible

- Cons
  - ▼ PRs across repos for changes
    - ▼ Working on „develop" can break
    - ▼ Out of sync
    - ▼ Synced Releases difficult
  - ▼ CI/CD hard to set up and maintain
  - ▼ Dependencies are hell to manage
  - ▼ End-to-end tests difficult to set up
  - ▼ Could silo knowledge in teams/repos
  - ▼ Can be difficult „keeping up"

# Alternative: Anemoi Ecosystem (mono repo)

| One Repository | |
|---|---|
| **Multi Package** | |
| Training | pip install anemoi-training |
| Models | pip install anemoi-models |
| Graphs | pip install anemoi-graphs |
| Datasets | pip install ecmwf-ml-data |

- Pros
  - Single PR for change across packages
    - Easier Refactors too
  - Easy CI including end-to-end tests
  - Easiest security scanning
  - Consistent coding standards
  - Common state of all „main"s
- Cons
  - „main" will move even faster
  - Complex release cycle
  - Different working styles might clash
  - Large repository size can make git slow
  - Might need „mono-repo tools"
  - Branching strategy complex
  - Easy to break „everything" accidentally

# Releasing Anemoi Training to the World

# Going Global

# Global Collaboration



Legend:
- Co-operation Agreements
- ECMWF Fellows
- Space Agencies
- South-East European Multi-Hazard Early Warning Advisory System
- Support for training and access to forecasts (World Bank)

Map labels:
ESA, EUMETSAT, ALADIN/HIRLAM, EC, WMO, CLRATP, CTBTO, HUNGARY, UKRAINE, REPUBLIC OF MOLDOVA, ROMANIA, SERBIA, BULGARIA, NORTH MACEDONIA, TURKEY, GREECE, LEBANON, JORDAN, KAZAKHSTAN, UZBEKISTAN, KYRGYZSTAN, TAJIKISTAN, TURKMENISTAN, CMA, JAXA, JMA, RIMES, CYPRUS, ISRAEL, SLOVENIA, CROATIA, BOSNIA & HERZEGOVINA, MONTENEGRO, ALBANIA, ACMAD, US NCAR, US NWS/NOAA, NASA, ESO, INPE

Fellows:
Patrick Eriksson, Sándor Baran, Daniela Jacob, Louise Nuijens, Hannah Cloke, Maria-Helena Ramos, Christian Grams, Daniela Domeisen, Marc Bocquet, Heini Wernli, Gabriele Pfister

# A Growing Ecosystem

# Anemoi Ecosystem ("Partial" mono-repo)

**Multi Repository**

**Multi Package**

| Training | pip install anemoi-training |

| Models | pip install anemoi-models |

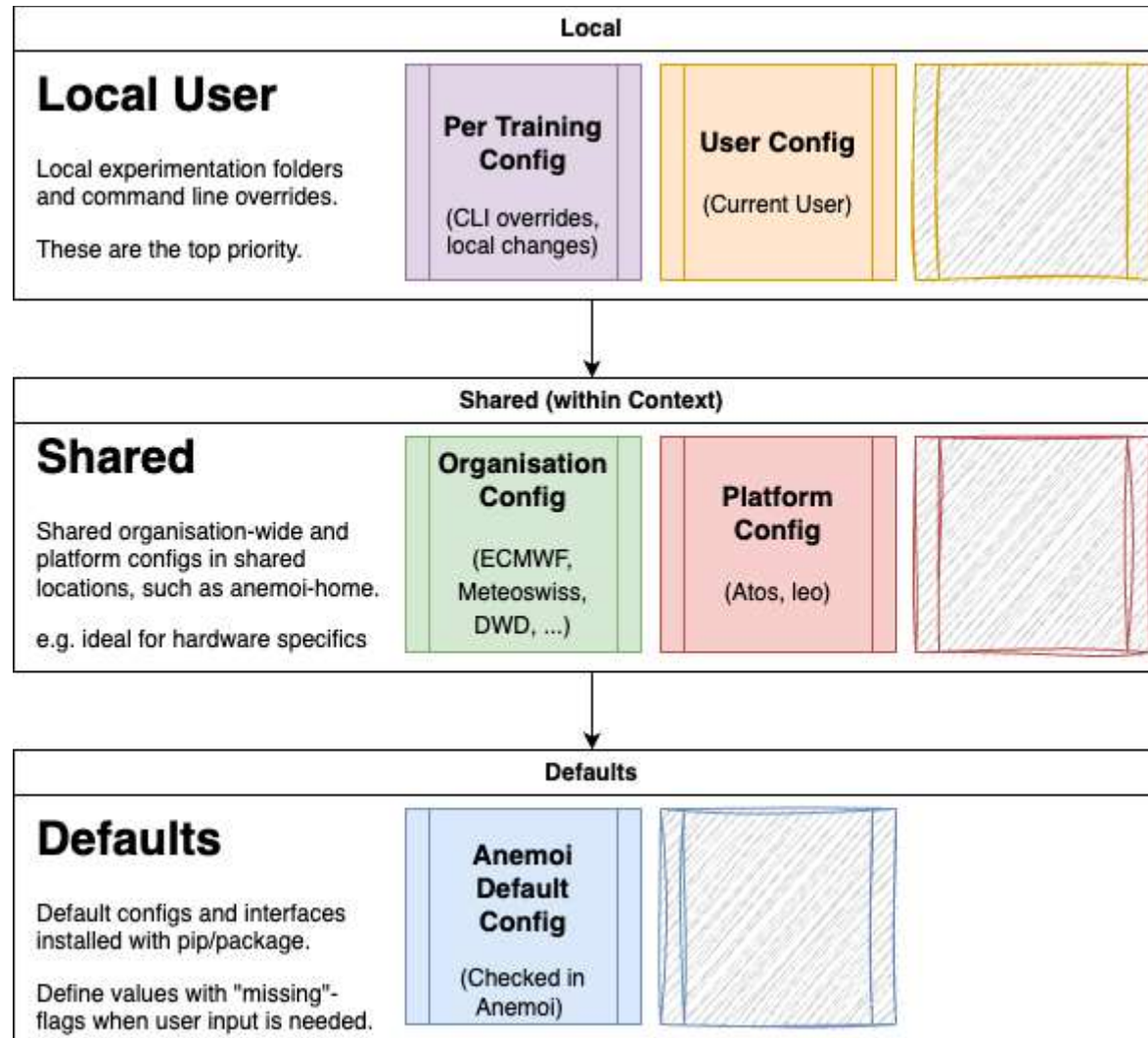| Graphs | pip install anemoi-graphs |

| Datasets | pip install anemoi-datasets |

- Pros
  - Some Separation of concerns
  - Puts tightly coupled code together
  - Different styles of collaboration possible
  - Simpler testing
  - Simpler configuration management

- Cons
  - Solves some CI but not „end-to-end"
  - Inconsistent workflows across anemoi
  - Some complex release workflows
  - Risk of creating artificial boundaries
  - Could complicate some dependencies

# How AI failed me here



```
1  #!/bin/bash
2
3  merge_repo() {
4      local repo_url="https://github.com/ecmwf/anemoi-$1"
5      local target_subtree=$1
6      local ref="${2-develop}"   # Default to 'develop'
7
8      echo "Merging repository from $repo_url into $target_subtree"
9
10     git subtree add --prefix "$1" "$repo_url" "$ref"
11
12     # add remotes for local exploration
13     git remote add "$target_subtree" "$repo_url"
14     git fetch "$target_subtree"
15
16     echo "Successfully merged $target_subtree"
17 }
18
19 # Example usage:
20 # ./merge_repos.sh
21
22 # Replace these with your actual repository names
23 merge_repo graphs
24 merge_repo models
25 merge_repo training
```

# Extend Configs to Institutional levels
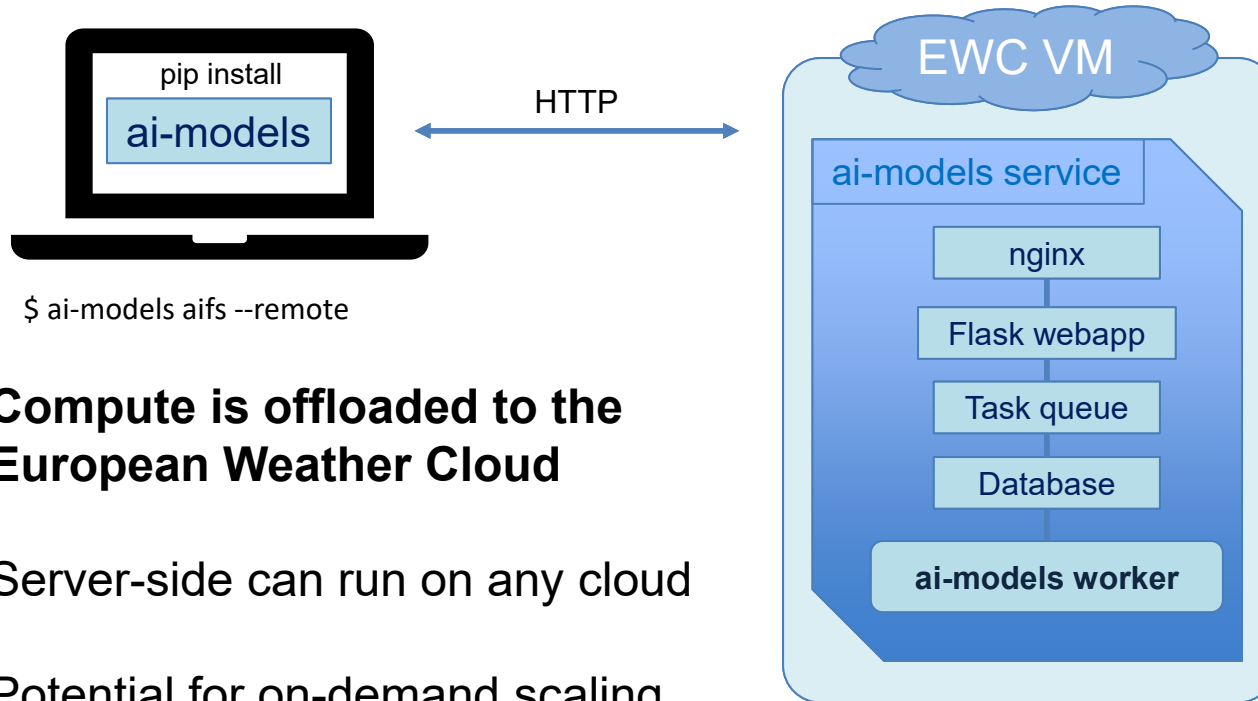
# Automate the Release process

# External and Global Collaboration

- Kick things off with a "hackathon"

  – You can even do a Code Tour at a hackathon!

- Maintain roadmaps / Kanban boards

- Code Reviews of PRs are essential

  – Consider dual reviews:

    • For code quality

    • For scientific validity / business case

- Maintain configurability, modularity and extensibility

- Automate what you can to keep morale high

- Don't be afraid to jump on calls

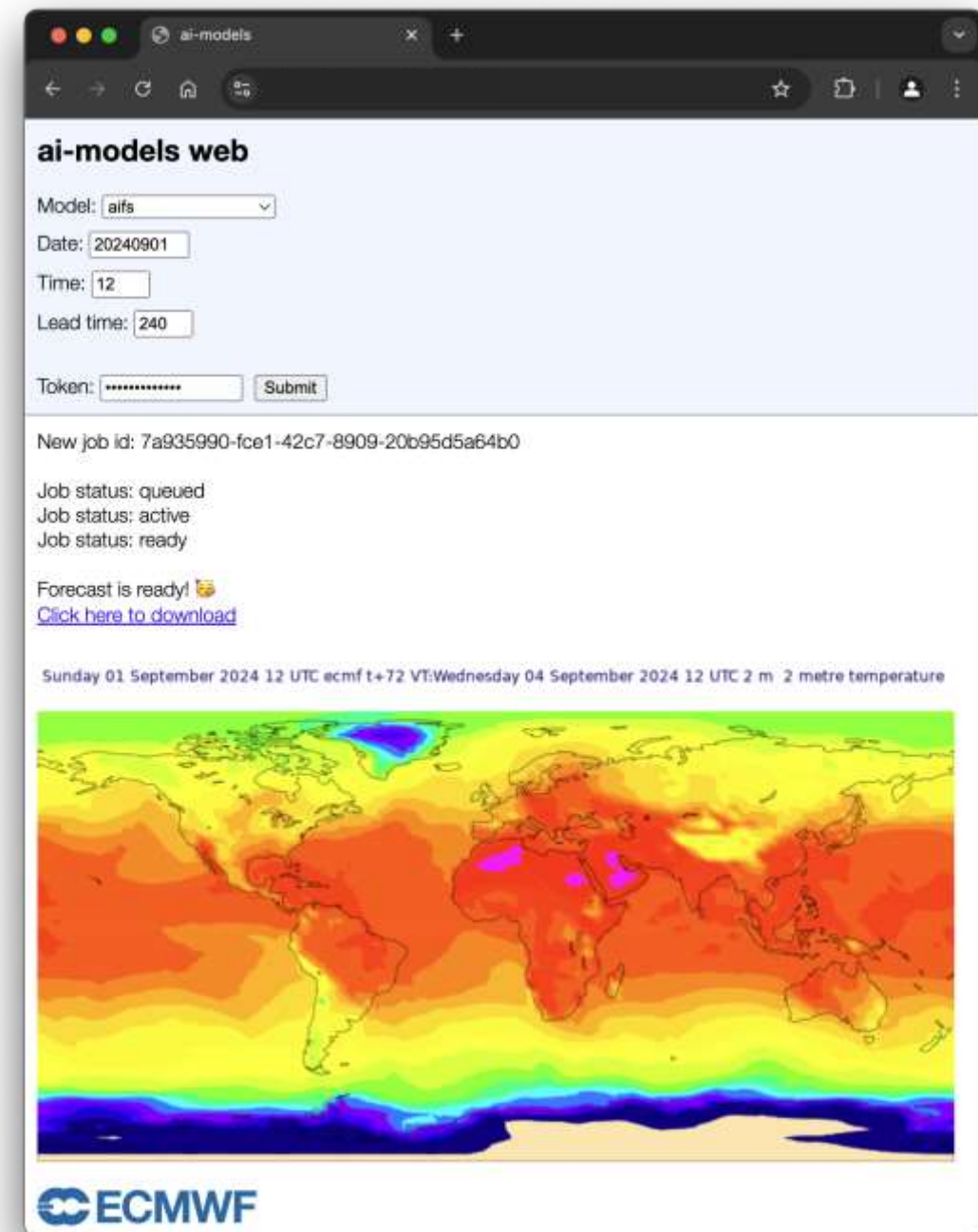**ECMWF**  EUROPEAN CENTRE FOR MEDIUM-RANGE WEATHER FORECASTS

# Success stories!

# Running AIFS anywhere



```
pip install
ai-models
```
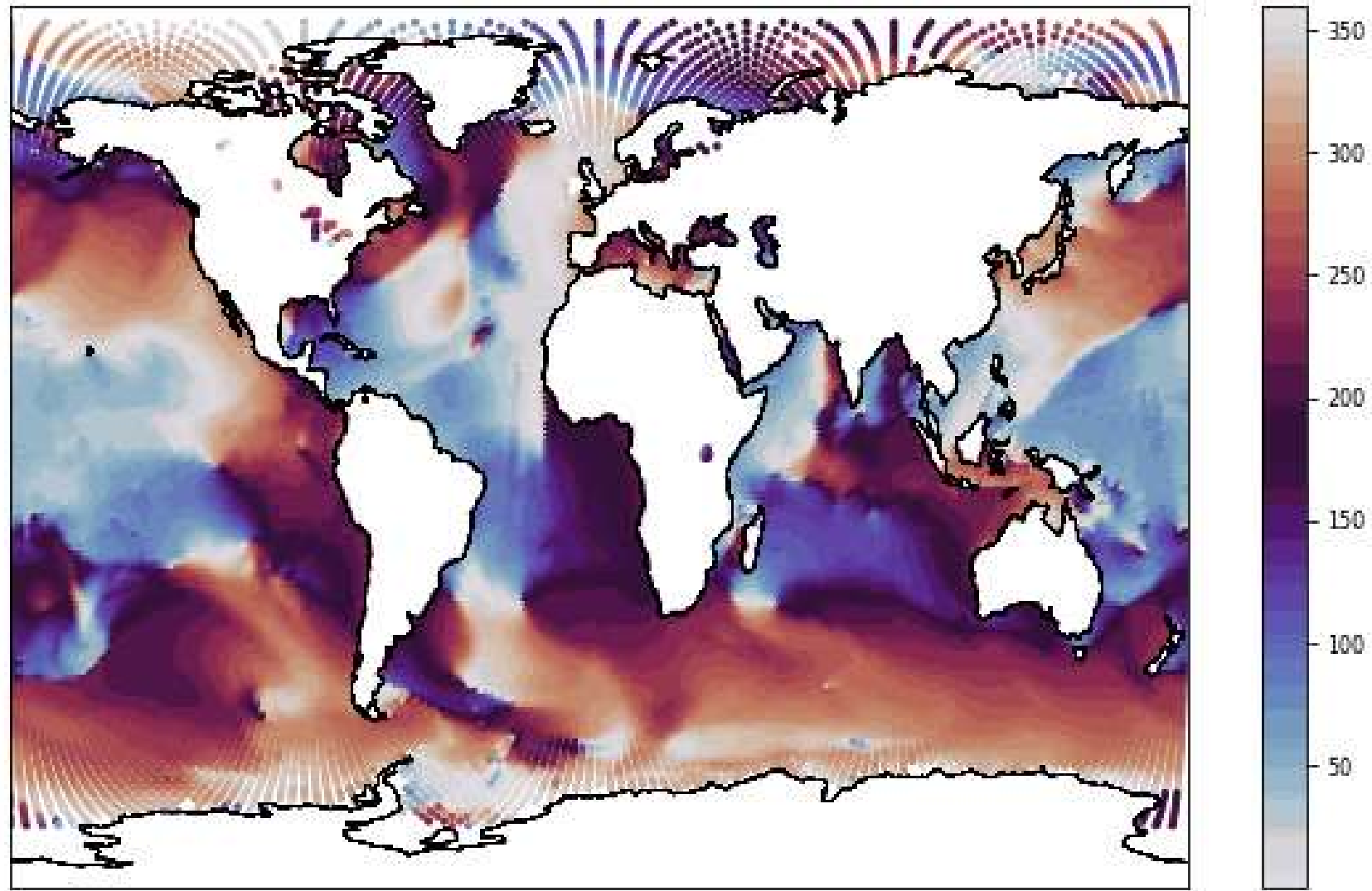
$ ai-models aifs --remote

- **Compute is offloaded to the European Weather Cloud**

- Server-side can run on any cloud

- Potential for on-demand scaling

- Forecast-in-a-box PoC

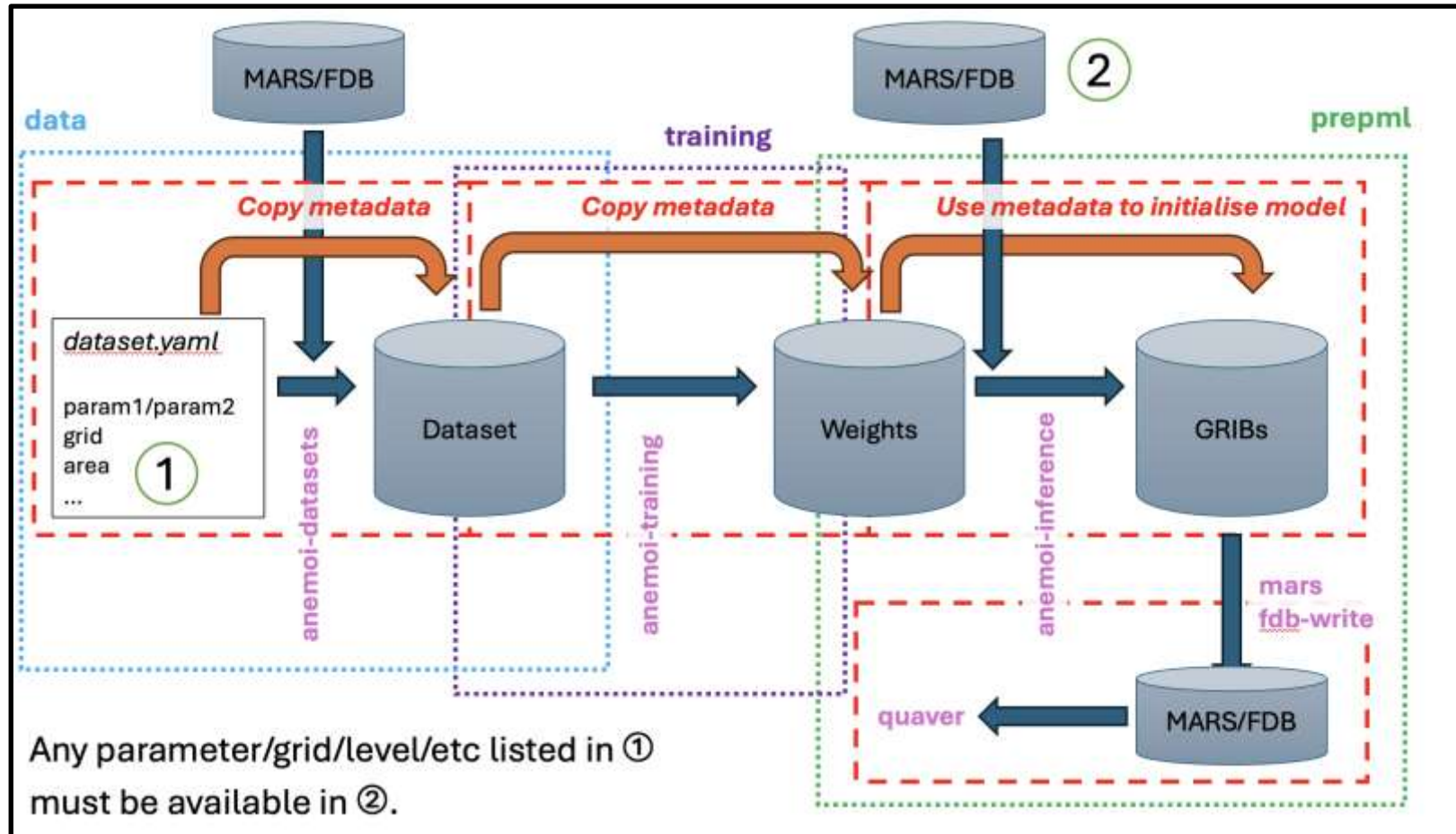**Provided continuity for AIFS and the other ML models during 2 GPU maintenance windows this year**



EWC VM

ai-models service
- nginx
- Flask webapp
- Task queue
- Database
- **ai-models worker**

# Adding in Earth System Components

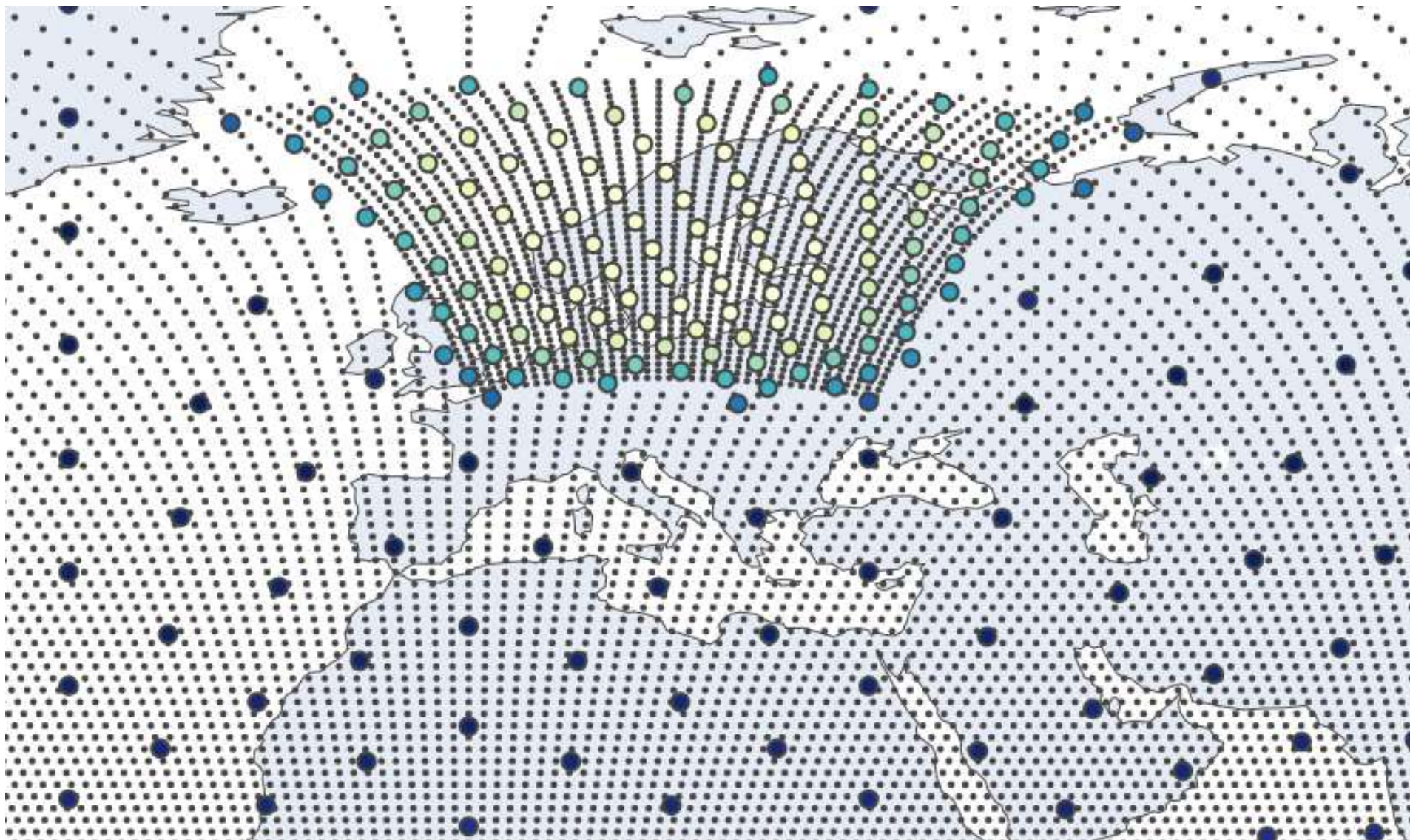the lower the better

Mean wave direction (14 day forecast)

# Tracking metadata through the entire ecosystem

All this is possible thanks to tracking the metadata.

# Enabling many types of graphs

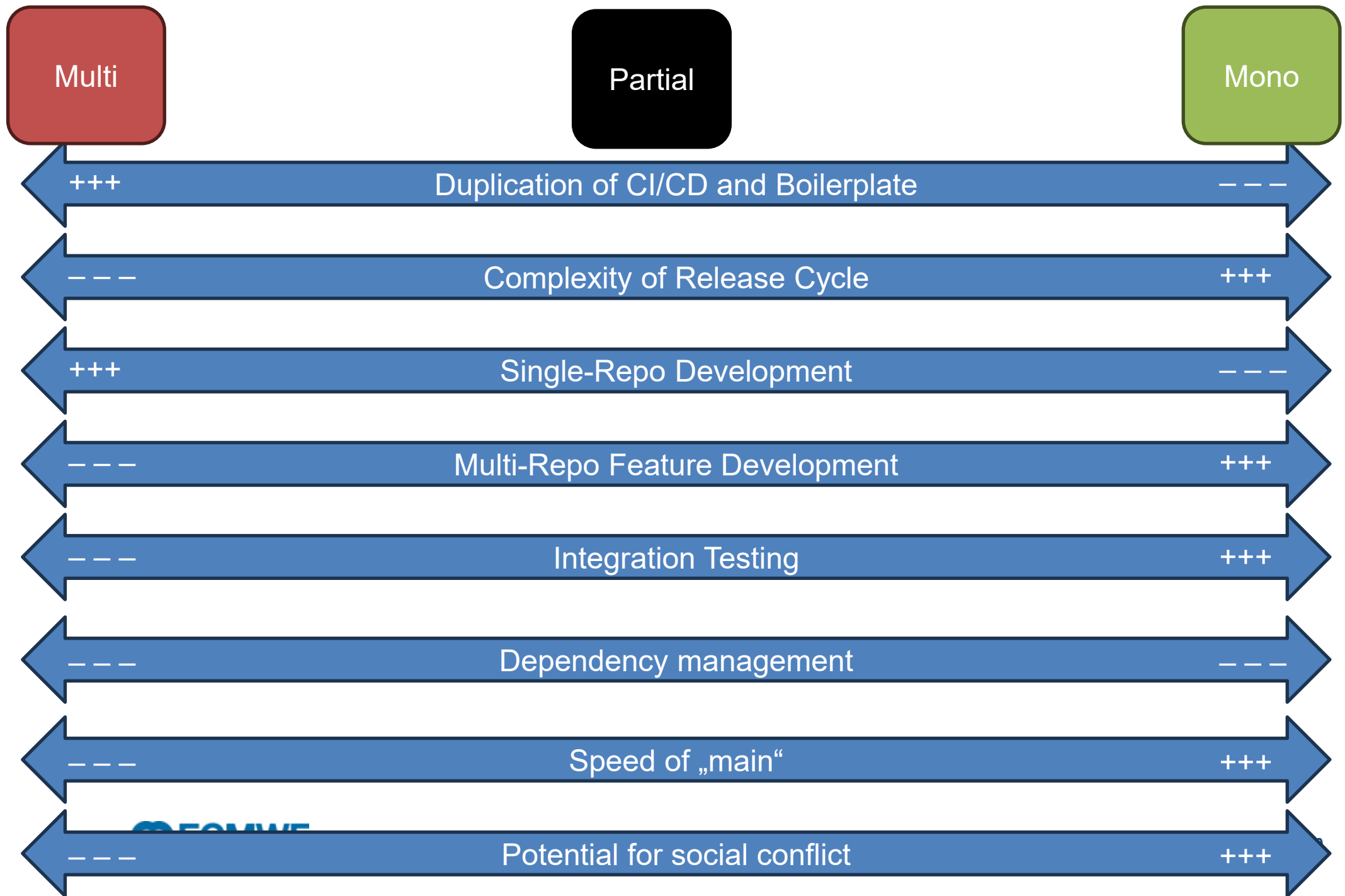# Enabling many types of graphs (courtesy Met Norway)



2023/02/06 18Z
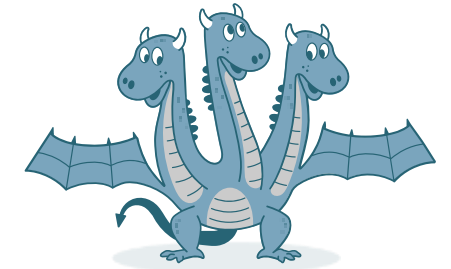
10m wind speed (m/s)

# Takeaways

# Learnings for growing software projects

- Get yourself a software architect or two

- Do Code Tours and in-person events

- Anticipate user needs and extensibility requirements

- Make decisions that are easy to reverse

- Make design choices that are easy to extend

- Don't just let AI agents run through your code

- Repo structures are difficult; use tools to make it easier

- Keep dependencies light, they are hell after all

- Collaborate and listen!

# Core Tech Stack