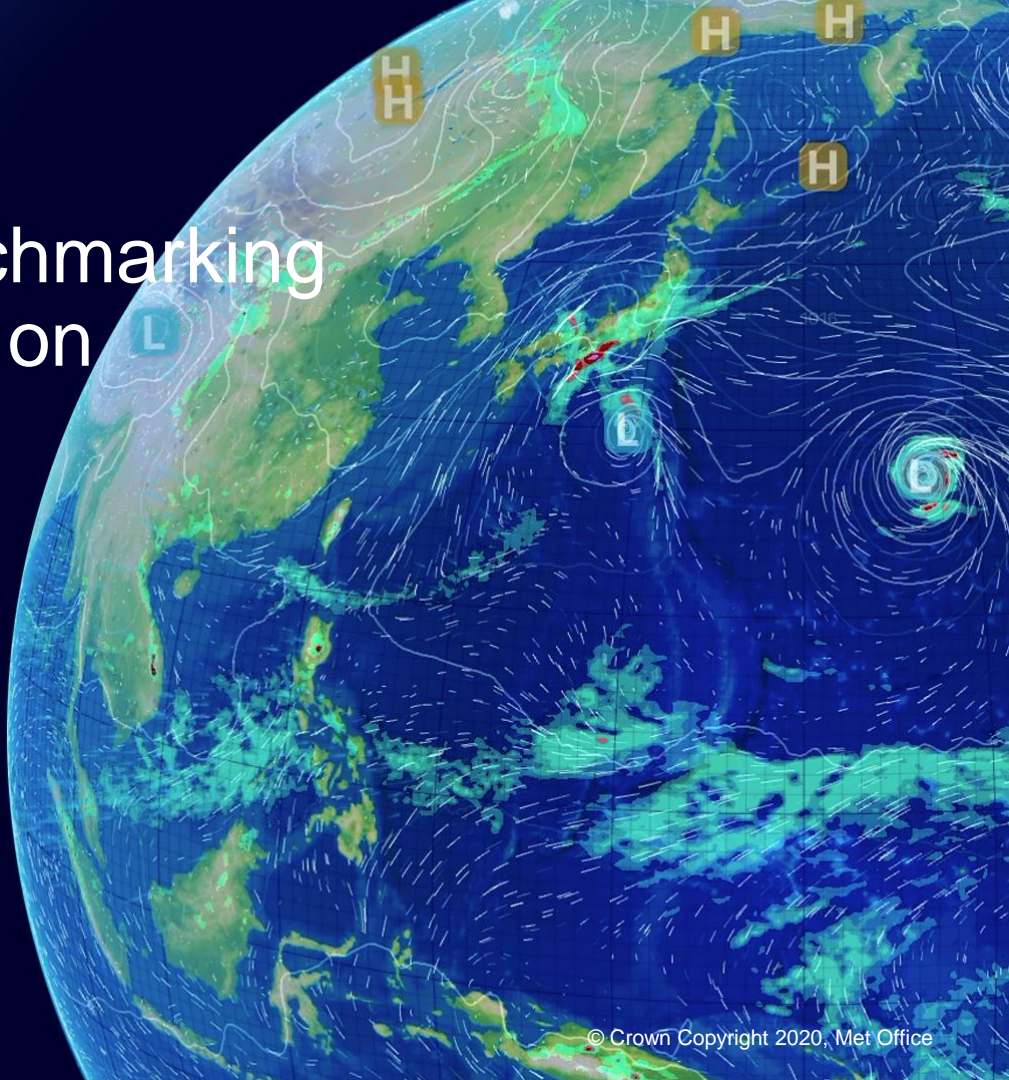# Developing the next benchmarking system for JULES based on ModelEvaluation.org

Heather Rumbold, Martin Best,
Gab Abramowitz, Adrian Lock
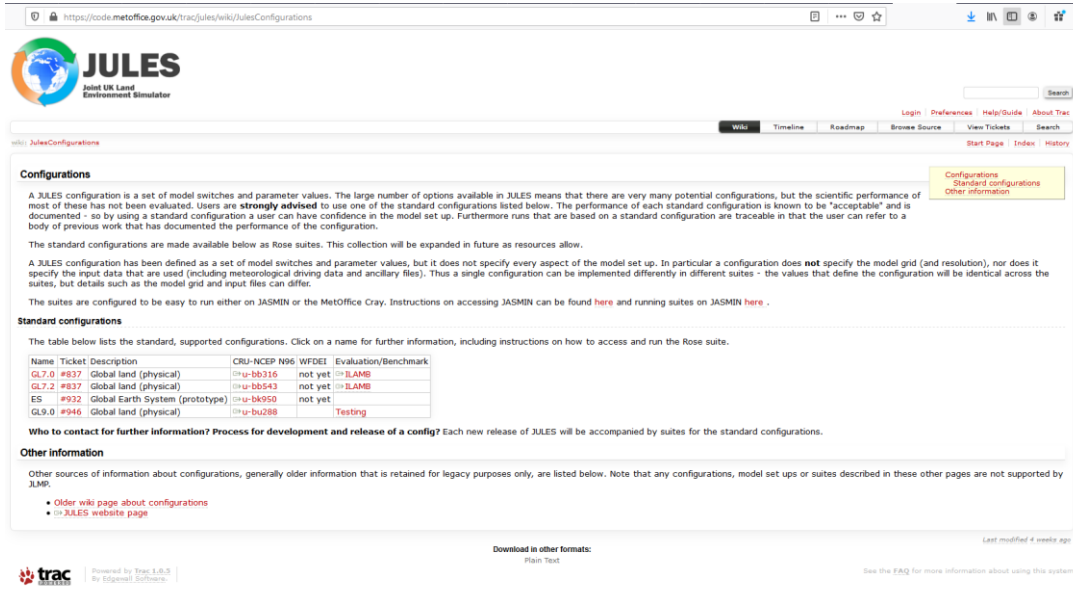
# Met Office

# Configuration Manager for the Global Land

My role…

- Maintain the standalone physical land model configuration versions on both the Met Office and NERC systems.
    - GL9 standalone is being finalised on Jasmin, ready for use shortly
- Build and maintain the comprehensive benchmarking system that will be used to assess new components for future configurations.
    - Generating a new benchmarking tool using ModelEvaluation.org for use along side existing tools

# JULES Standard Configuration

https://code.metoffice.gov.uk/trac/jules/wiki/JulesConfigurations



- Set of model and ancillary generation switches and parameter values

- GL/RL/ES are all required for the coupled system.
- Standalone only requires configurations for:
  1. Physical Land (weather/climate)
  2. Earth System

# JULES Standard Configuration

A configuration is not:

- Driving data

- The resolution

- The ancillary files (but can include the data sources)

- Application specific

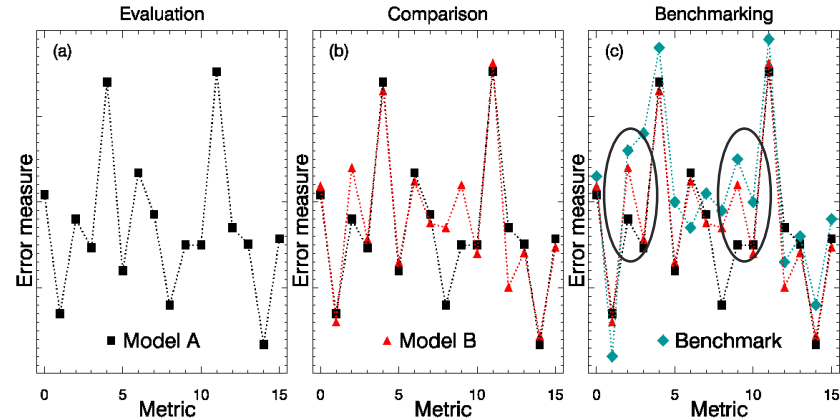Best combination of settings to give the best description of the physical environment

How do we know when we have a better description of the environment?

# What is benchmarking?

**The Plumbing of Land Surface Models: Benchmarking Model Performance**

M. J. Best,[a] G. Abramowitz,[b] H. R. Johnson,[a] A. J. Pitman,[b] G. Balsamo,[c] A. Boone,[d]
M. Cuntz,[e] B. Decharme,[d] P. A. Dirmeyer,[f] J. Dong,[g] M. Ek,[g] Z. Guo,[f] V. Haverd,[h]
B. J. J. van den Hurk,[i] G. S. Nearing,[j] B. Pak,[k] C. Peters-Lidard,[j]
J. A. Santanello Jr.,[j] L. Stevens,[k] and N. Vuichard[l]

(2015) Journal of Hydrometeorology, 16, 1425-1442.

![Three panel figure showing Evaluation (a), Comparison (b), and Benchmarking (c) plots of Error measure versus Metric, with Model A (black squares), Model B (red triangles), and Benchmark (teal diamonds)](figure)

➢ Model outputs are compared to a predefined benchmark

➢ 3 types of benchmark:

1. Is it better than another model?

2. Is it fit for a particular application?

3. Can it effectively utilise available information?

"Ultimate" benchmark – model to be within the observational error

# What will benchmarking do for JULES?

- JLMP – Require a single configuration which generates the best simulation of JULES as a whole system
  - Is the new JULES configuration better the previous model configuration? (i.e. no 1)
  - E.g. Does adding X piece of new science code improve JULES compared to the previous configuration?
  - Old configuration version will become the benchmark

- JULES community are aiming for 2 or 3 – Best science for a specific area
  - E.g. Can the new configuration capture specific impacts (e.g. the river flow or snow depth) better than the old configuration?
  - E.g. If supplied with better inputs (e.g. high resolution veg ancillaries) it should be expected to perform better than a configuration without this.

**Met Office**

# A new benchmarking suite

Coming soon… upload automated within the suite

PLUMBER2 data 170 sites from FLUXNET2015, FLUXNET La Thuile & OzFlux + canopy height, LAI reference height & IGBP vegetation + HWSD soils

Python script → convert jules input variables into json file

Rose suite Run JULES for all sites in json file

Upload data to modelevaluation.org

Benchmarking output

Perform analysis

```
{
    "AR-SLu": {
        "data_start": "2010-01-01 00:00:00",
        "data_end": "2011-01-01 00:00:00",
        "data_period": 1800,
        "drive_file": "/data/users/hashton/PLUMBER2/met_f
        "latitude": -33.4648,
        "longitude": -66.4598,
        "spinup_start": "2010-01-01 00:00:00",
        "spinup_end": "2011-01-01 00:00:00",
        "main_run_start": "2010-01-01 00:00:00",
        "main_run_end": "2011-01-01 00:00:00",
        "timestep_len": 1800,
        "z1_tq_in": 11.0,
        "z1_uv_in": 11.0
    },
```

# Met Office
# ModelEvaluation.org

# Single site vs Observations

# Multi site analysis

- Variable breakdown…
- JULES LE beats the linear regression models
- H does not, however:

quantile value =
(highest rank – JULES model)/
(lowest rank – highest rank)

- H isn't as bad as it looks!
- Overall JULES is as good as or better than the benchmarks

**Breakdown by variable**

Average metric quantile over sites and metrics

**Details**

- Model
- Emp2lin
- Emp1lin

**Metric quantile av. over:**

**8 metrics**
RMSE, MBE, NME
SDdiff, correlation, fifthdiff
ninetyfifthdiff, PDFoverlap

**170 sites**
AR-SLu - PLUMBER2
AT-Neu - PLUMBER2
AU-ASM - PLUMBER2
AU-Cow - PLUMBER2
AU-Cpr - PLUMBER2
AU-Ctr - PLUMBER2
AU-Cum - PLUMBER2
AU-DaP - PLUMBER2
AU-DaS - PLUMBER2
AU-Dry - PLUMBER2
AU-Emr - PLUMBER2
AU-GWW - PLUMBER2
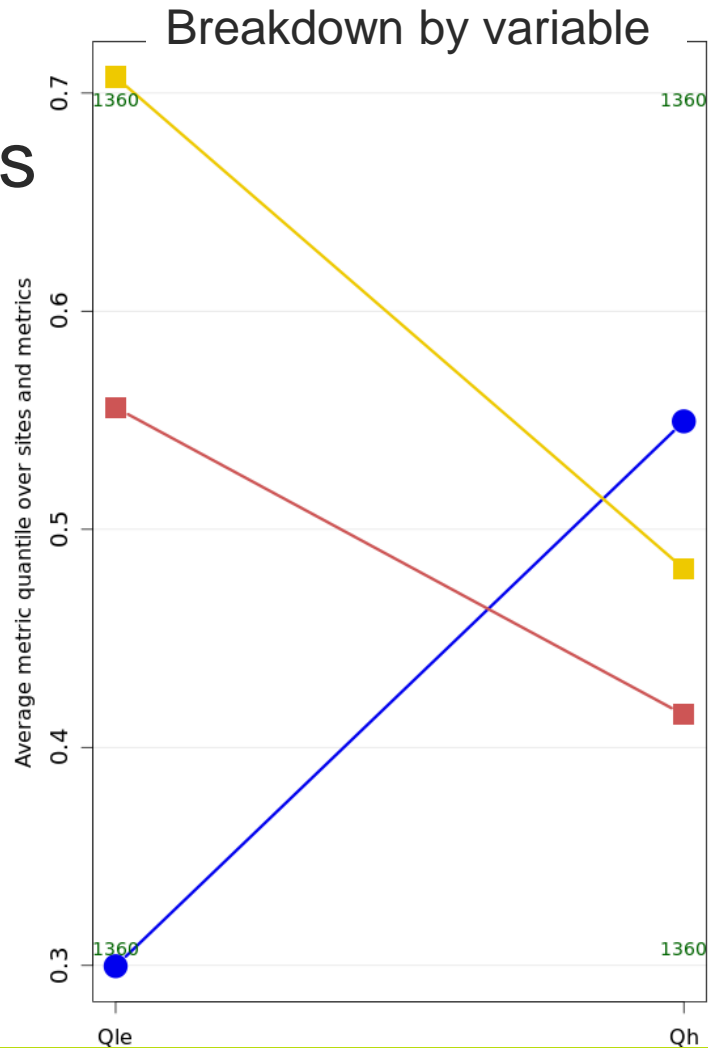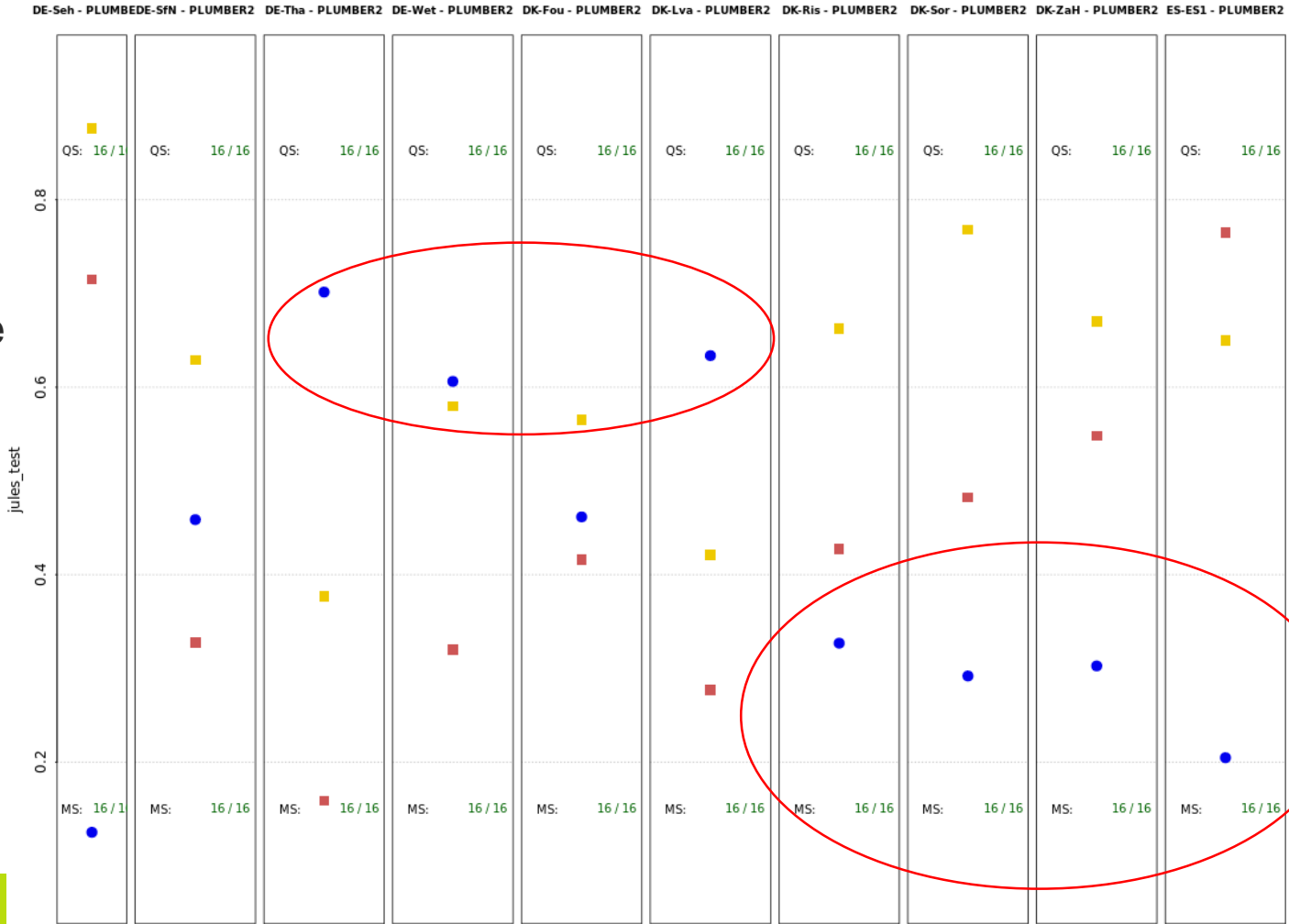AU-Gin - PLUMBER2
AU-How - PLUMBER2
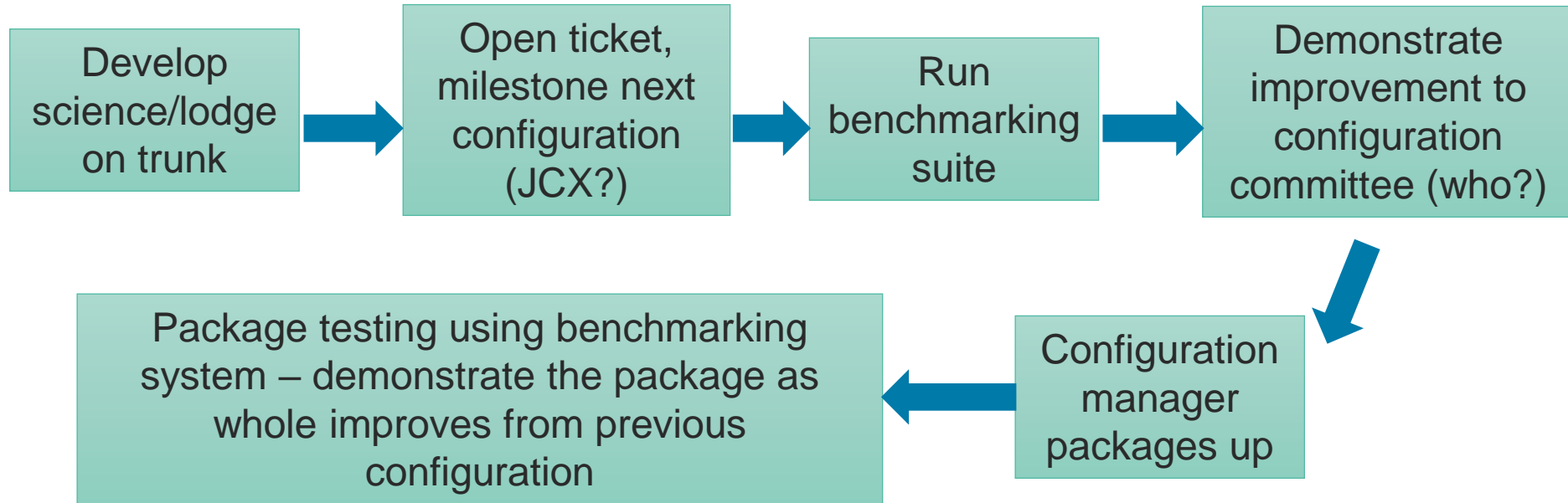AU-Lit - PLUMBER2
AU-Otw - PLUMBER2

# Met Office

# Development Process for standalone Physical Land configurations (work in progress)

| Develop science/lodge on trunk | → | Open ticket, milestone next configuration (JCX?) | → | Run benchmarking suite | → | Demonstrate improvement to configuration committee (who?) |

Package testing using benchmarking system – demonstrate the package as whole improves from previous configuration ← Configuration manager packages up
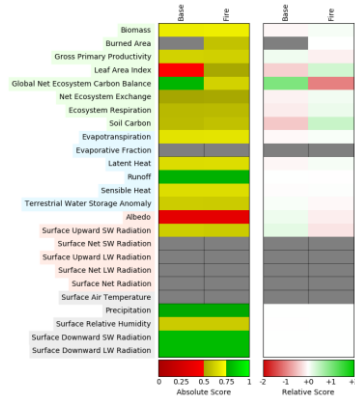
Benchmarking is an important part of this process!!

# How does this fit in with other tools going forward?



Physical Land



Earth System

Plus others….?



AutoAssess &
Validation Notes

Is the new JULES configuration better the previous model configuration?