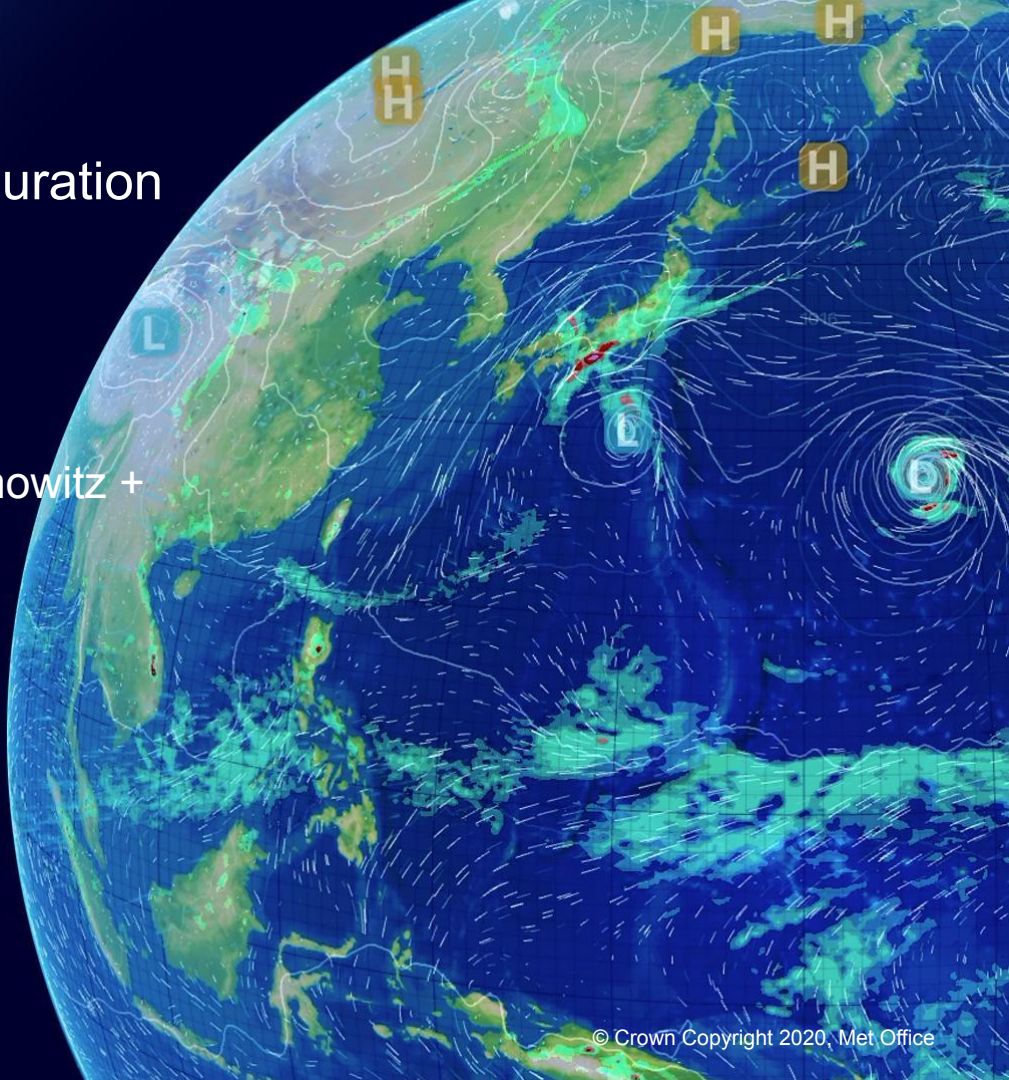


Developing the next standard configuration
for standalone JULES using a
benchmarking system based on
ModelEvaluation.org.

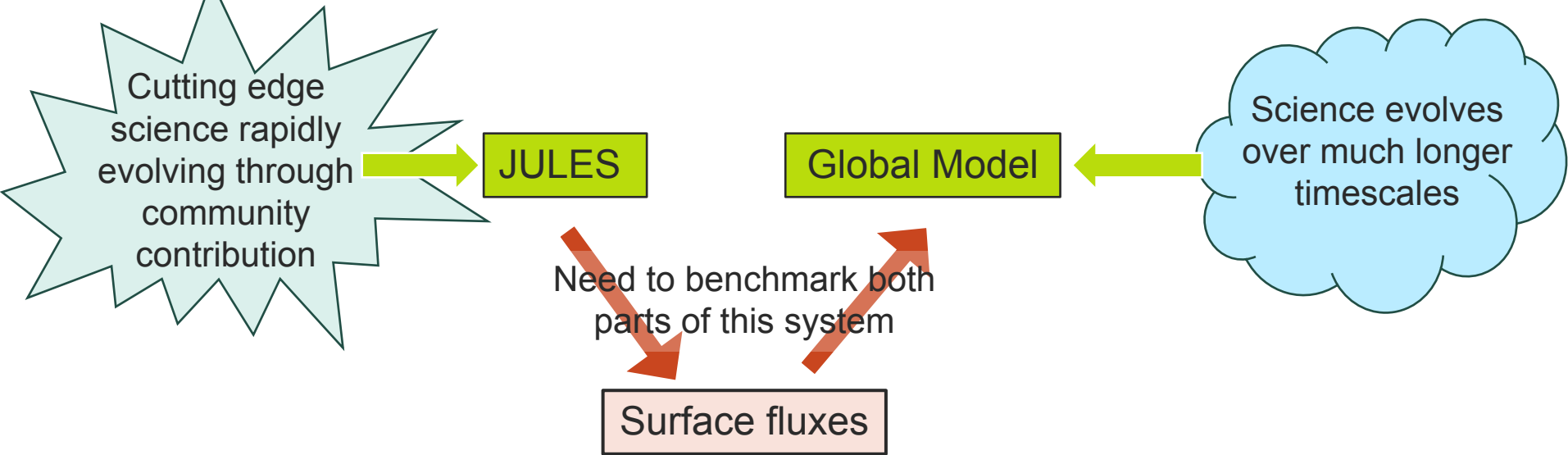
Heather Rumbold, Martin Best, Gab Abramowitz +
JULES Community

Seamless Global Modelling Workshop

Thursday June 5th, 2025



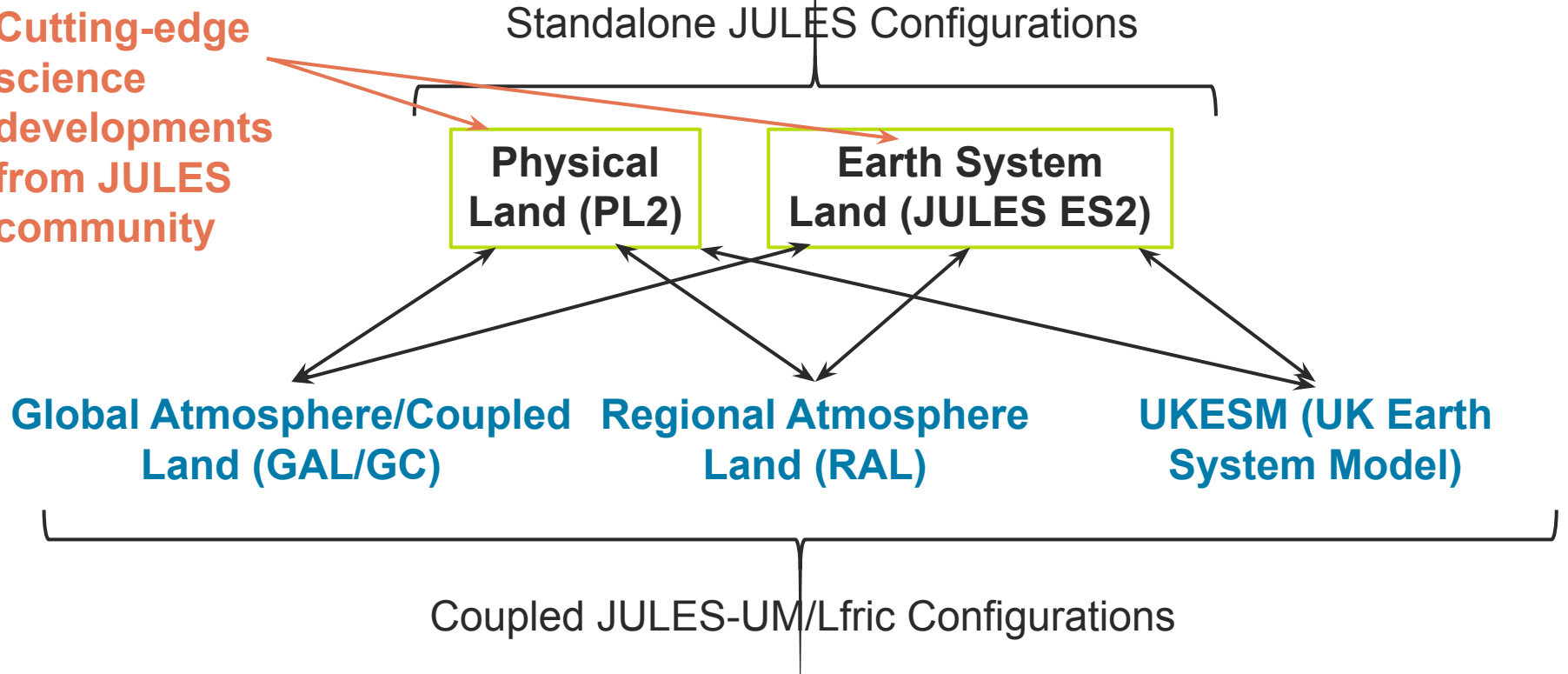
Motivation



- Standalone JULES configurations are evaluated and benchmarked by observations.
- This allows us to isolate the scheme and develop the best configurations for the land.
- However, compromises must be made when science is pulled through to global model configurations.
- Can we use a standalone benchmarking system to identify errors to improve the atmospheric model?

 **Met Office** **Configurations required for Standalone JULES**

Cutting-edge science developments from JULES community

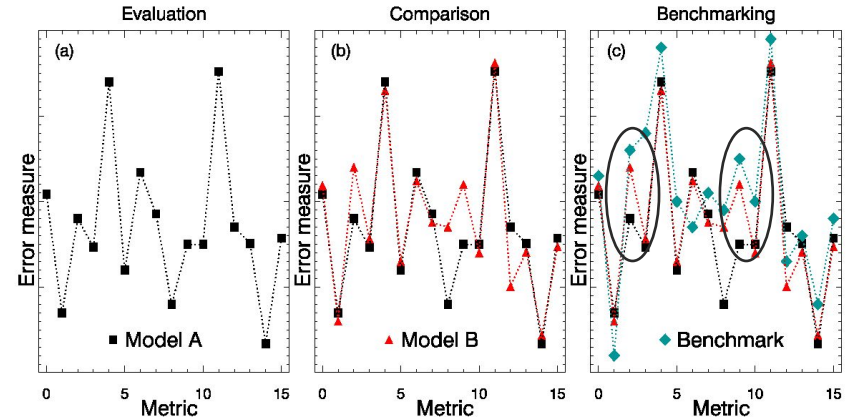


Aim: To create **standard configurations** that give the **best simulation** of the **physical and earth system environments**, demonstrated through **predefined benchmarks**.

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,^c 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.



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?

e.g. previous model configuration
 Does adding new science code improve JULES compared to the previous configuration?

- **Community web-based environment** for model evaluation and benchmarking.
- Users run their **land models offline** in single-site mode, forced by **locally observed sub-daily meteorology** (downloaded from ME.org), for a **wide variety of sites** across the globe.
- Model output is then **uploaded to the web application** and **analysed** using scripts that compare output with evaluation **data products, other models and empirical or ML benchmarks**.

Advantages:

- System shared by many international research groups, providing opportunities to benchmark against many other land surface models.
- Able to utilise standard testing experiments such as linear regression and machine learning models.
- Can utilise an extensive range of metrics and statistics.
- Use of consistent, reliable and robust, quality-controlled data for many sites.

The screenshot shows the ModelEvaluation.org website. At the top, there is a navigation bar with links for Home, Info, Data Sets, Experiments, Model Profiles, Model Outputs, and Analyses. A green button indicates the user is not logged in. Below the navigation bar, a welcome message is displayed. A central diagram illustrates the workflow: 'Your machine' (Run your model in your local environment) feeds into 'Upload your model output', which then goes to 'View evaluation' on the website. Conversely, 'Choose experiment' and 'Download driving data' are also shown as part of the process. To the right, a line graph titled 'Smoothed On: 14-day running mean. Obs - US-Me2_FLUXNET2015 Model - CABLE_FLUXNET2015' displays three data series: US-Me2_FLUXNET2015 (red), CABLE_FLUXNET2015 (blue), and CABLE_FLUXNET2015_obs (green). The graph shows seasonal fluctuations in gPP (gC/m²/day) over time. Statistics for the models are provided: US-Me2_FLUXNET2015 (Min = 1206, Max = 412), CABLE_FLUXNET2015 (Min = 1243, Max = 727), and CABLE_FLUXNET2015_obs (Min = 167.6, Max = 18.3). The score for the smooth is 0.563 and the score for all is 0.519.

Hosted at Australian supercomputing facility (NCI) with HPC back-end capability

PLUMBER2 Benchmarking Workflow

PLUMBER2 dataset
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

```
{  
  "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  
  },  
}
```

Rose workflow
Run JULES for all
sites in json file

task	status	host	job-system	job-ID	T-submitt	T-start	T-finish	dT-hour	submit-message
NETO_SPCZ	running								
jules_FR-1q7	waiting			436649	10:45:12Z	10:45:12Z			job(0) started
jules_AD-Whr	waiting								
jules_FR-1q3	waiting								
jules_US-10r	waiting								
jules_ES-1M6	waiting								
jules_AD-1low	waiting								
jules_DE-1M5	waiting								
jules_CN-1Ch	waiting								
jules_US-1M2	waiting								
jules_ES-1q5	waiting								
jules_US-1M6	waiting								
jules_AD-1Ch	waiting								
jules_DE-1R1	waiting								
jules_IT-1Noe	waiting								
jules_US-1M0	waiting								
jules_CN-1Ch	waiting								

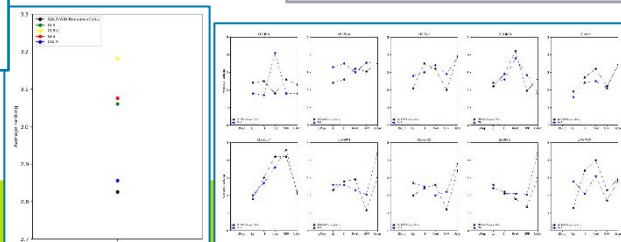
Calculate output
metrics*:
JULES vs. Obs for
every site, for LWup,
LE, H, Rnet, GPP &
Ustar, for each
configuration

Configuration/
benchmark are
ranked for each
metric, var, site

Perform the
average
ranking of
metrics for
every site and
variable

Benchmark configurations:
Standalone versions of
GL8, GL8.1, GL9 & GAL9

Plot output and
document ticket



* Metrics: RMSE, MBE, NME,
StdDevdiff, Correlation, 5thdiff, 95thdiff,
PDFOverlap, Skewness, kurtosis

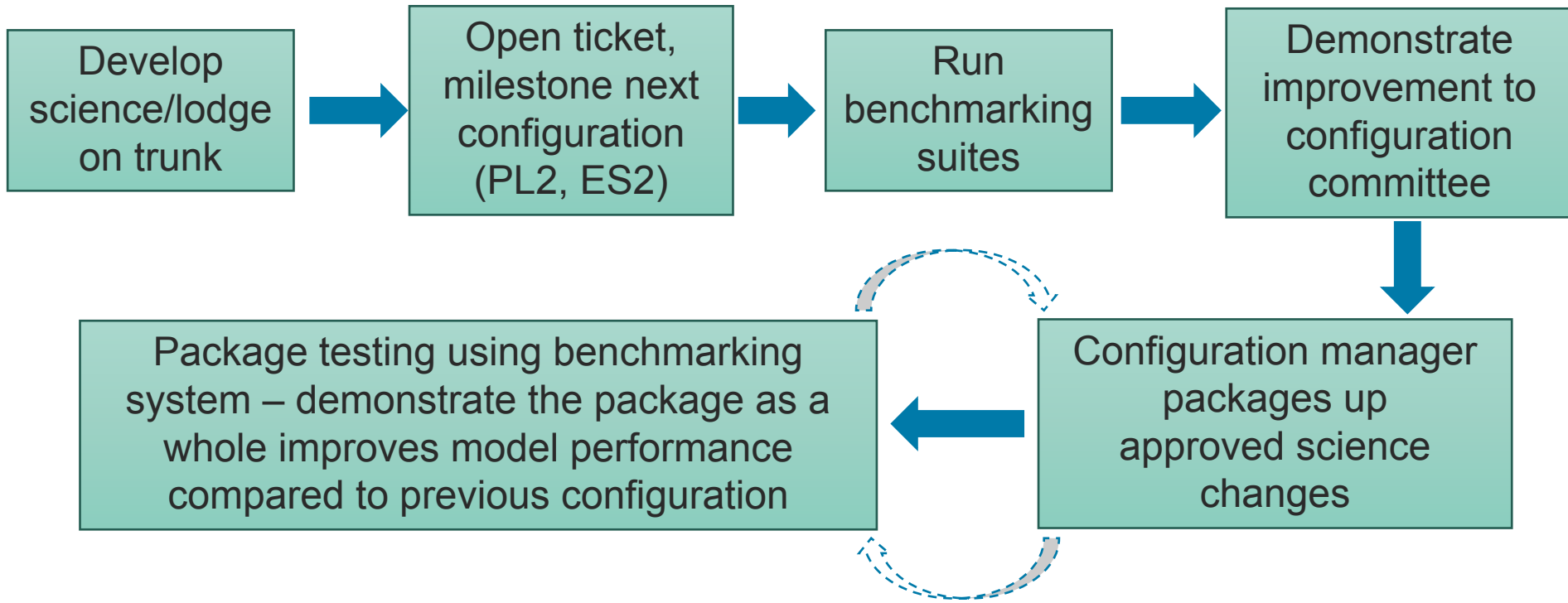
Analysis

- Compares JULES vs Obs for every **site, variable, model configuration and model benchmarks.**
- 10 Statistical metrics are calculated for each instance.

Sites	170 FLUXNET sites, global coverage, with site metadata AND evaluation data
Metrics	RMSE, MBE, NME, StdDevdiff, Correlation, 5thdiff, 95thdiff, PDFoverlap, Skewness, Kurtosis
Variables	LE, H, LWup, Rnet, GPP and momentum flux
New Model Configuration	PL2 package test
Benchmarks	GL8, GL8.1, GL9 & GAL9

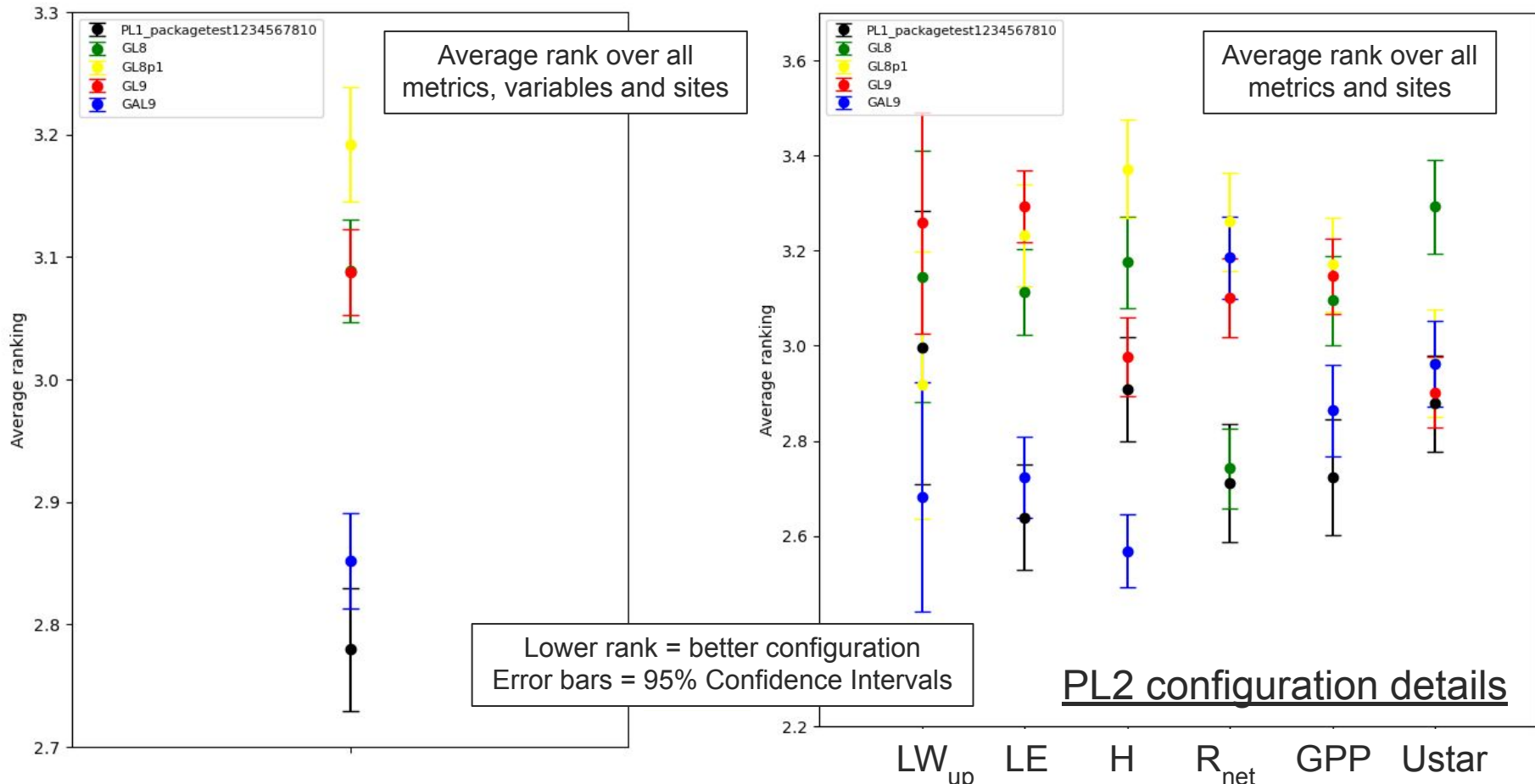
- Metrics values are ranked from 1 to 5 (1 being the best result)
- These rankings are averaged over all statistics and all sites to give an average ranking for all the variables separately.
- A final averaging is performed over all variables to give an overall ranking.

Development Process for standalone JULES configurations



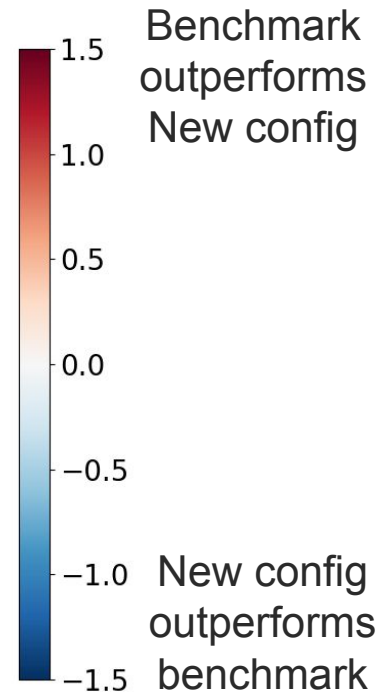
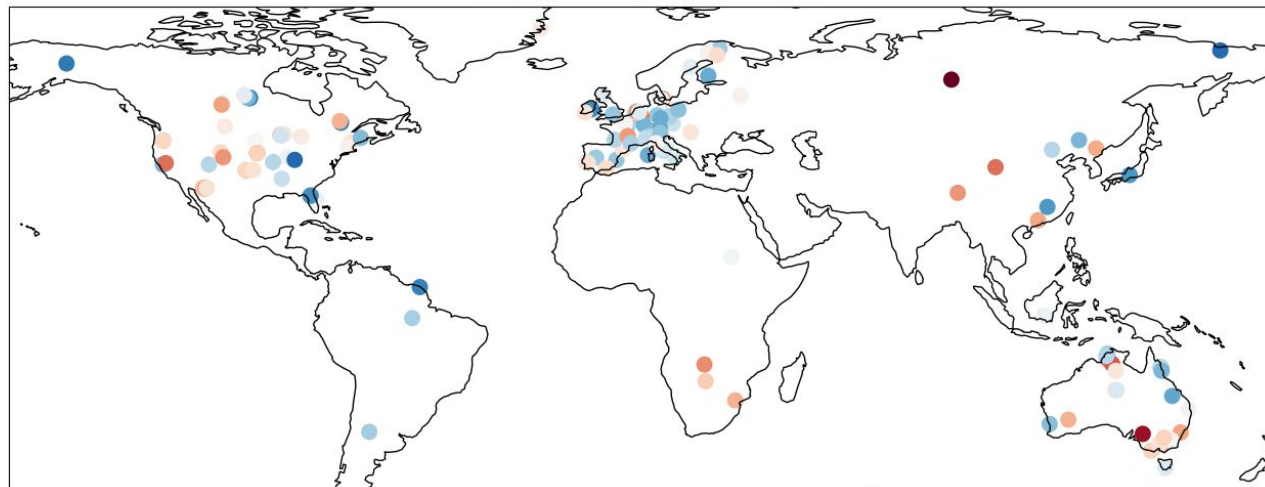
- Same process, as used by regional and global modelling, for the JULES community
- First cycle has been completed... configurations are now ready to release

Physical Land Package: PLUMBER2 benchmarking output



Met Office Average rank over all metrics and variables

Average rank difference over all variables between the model and latest benchmark (GAL9):
model minus benchmark



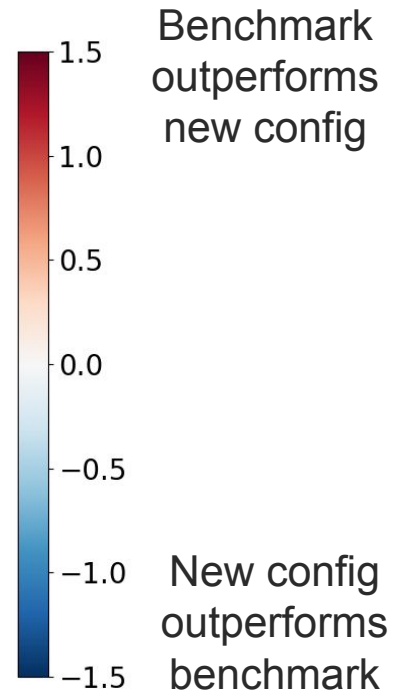
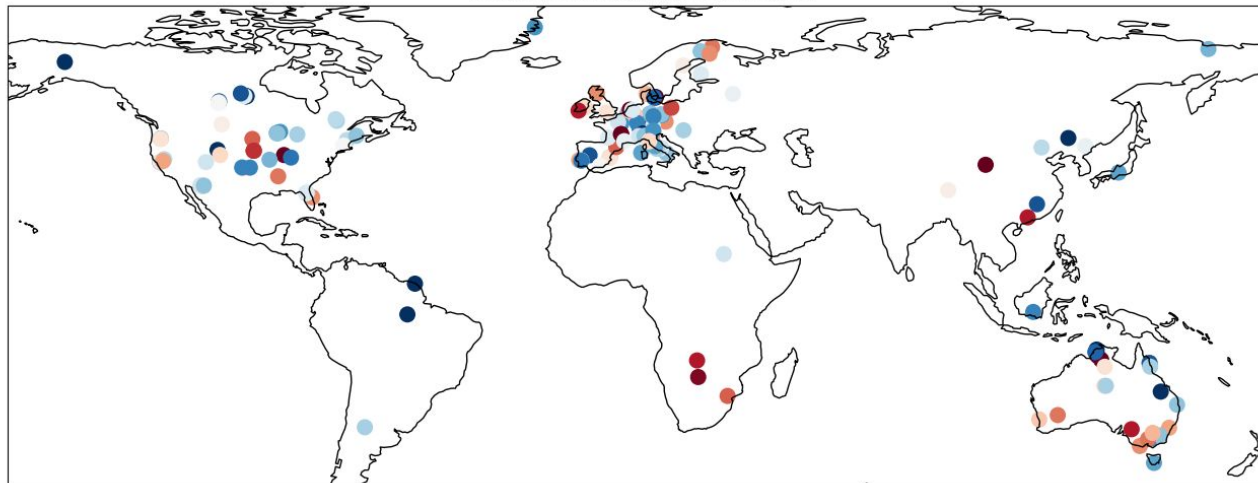
Model: JULES PL2. Benchmark: GAL9

If model rank is $>$ benchmark rank, then benchmark is doing better than model

Blue = better performing model i.e. PL2 configuration

Met Office Average rank over all metrics only: Latent heat flux

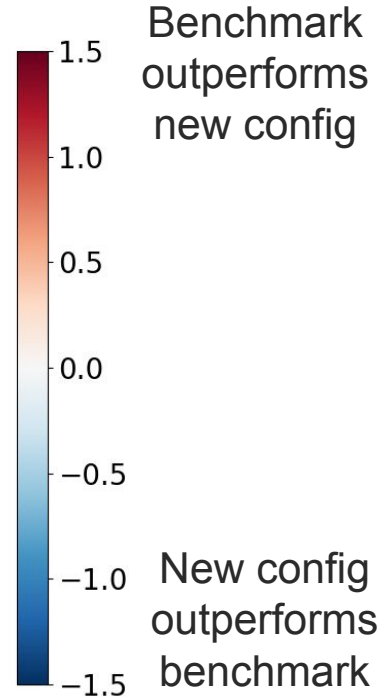
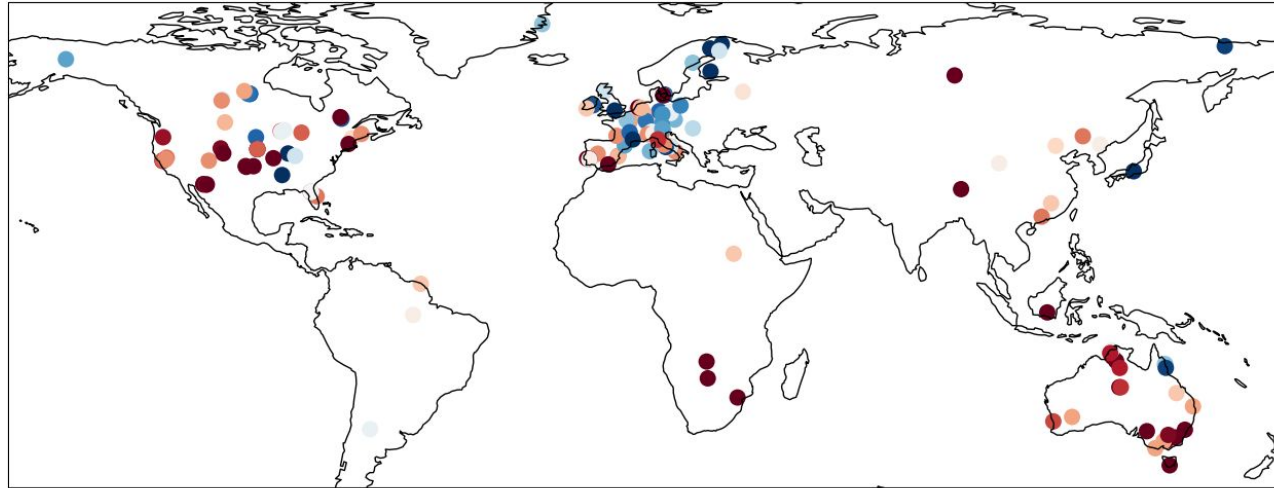
Average rank difference for latent heat flux between the model and latest benchmark (GAL9):
model minus benchmark



□ More points where benchmark is outperforming model

Met Office Average rank over all metrics only: Sensible heat flux

Average rank difference for sensible heat flux between the model and latest benchmark (GAL9):
model minus benchmark



- Many more points where benchmark is outperforming model
- Highlights the need for a more rigorous system which explores the data space in more detail...

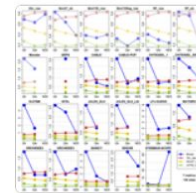
Uses Machine Learning models as our benchmarks - i.e. use ML to predict fluxes using meteorological variables as predictors

See seminar by Gab Abramowitz: "Machine Learning for Benchmarking Land Models", 21st May, [here](#)

- Simulations of latent, sensible heat, carbon fluxes at 154 flux tower sites.
- 20 different international land models.
- Linear regression, cluster and regression, LSTM, random forest as benchmarks – all out of sample
- Land model provided with observed vegetation type, reference and canopy height, LAI and 30-60min meteorological forcing. No calibration.
- How predictable are the surface fluxes? How good are existing mechanistic models at prediction? Under which conditions/environments do they excel, or we know they can improve?
- Results are broadly consistent: **Simple out-of-sample empirical models, including linear regression, comfortably outperforming mechanistic land models.**

On the predictability of turbulent fluxes from land: PLUMBER2 MIP experimental description and preliminary results

Gab Abramowitz , Anna Ukkola, Sanaa Hobeichi, Jon Cranko Page, Mathew Lipson, Martin G. De Kauwe, Samuel Green, Claire Brenner, Jonathan Frame, Grey Nearing, Martyn Clark, Martin Best, Peter Anthoni, Gabriele Arduini, Souhail Boussetta, Silvia Caldararu, Kyeungwoo Cho, Matthias Cuntz, David Fairbairn, Craig R. Ferguson, Hyungjun Kim, Yeonjoo Kim, Jürgen Knauer, David Lawrence, Xiangzhong Luo, Sergey Malyshev, Tomoko Nitta, Jerome Ogee, Keith Oleson, Catherine Ottlé, Phillipe Peylin, Patricia de Rosnay, Heather Rumbold, Bob Su, Nicolas Vuichard, Anthony P. Walker, Xiaoni Wang-Faivre, Yunfei Wang, and Yijian Zeng



Conclusions

- The **land component** can now be **formally assessed in isolation**, so we have a clearer **understanding** of the science processes.
- We now have a **benchmarking system** that runs very quickly giving us the ability to **assess sites and variables** individually for separate **individual science changes**
- **Pull through** of new science in the global model configurations will allow us to understand the **impacts of land science changes on the global model outputs**.
 - **Compromises** will still be needed to account for biases in the atmospheric model.
 - However, we are now able to repeat local assessment of the impacts of such compromises.
- **PLUMBER2 MIP** has provided lots more information which can be used to explore the **where**, **how** and **why** our land models aren't performing well and identify where it can be improved.
 - Provide the **most accurate fluxes** back to the global model.
- **Benchmarking** is critical to **target challenges in land modelling** and therefore **accelerate improvements** in the science across all scales.

Thank you for listening

Any Questions?

heather.rumbold@metoffice.gov.uk

