



Introduction

Discussions of how climate models should be evaluated tend to rely on either philosophical arguments about the status of models as scientific tools, or on empirical arguments about how well runs from a given model match observational data. These lead to quantitative measures expressed in terms of model bias or forecast skill, and ensemble approaches where models are assessed according to the extent to which the ensemble brackets the observational data.

Such approaches focus the evaluation on *models* per se (or more specifically, on the simulation runs they produce), as if the models can be isolated from their context. Such approaches may overlook a number of important aspects of the use of climate models:

- the process by which models are selected and configured for a given scientific question.
- the process by which model outputs are selected, aggregated and interpreted by a community of expertise in climatology.
- the software fidelity of the models (i.e. whether the running code is actually doing what the modellers think it's doing).
- the (often convoluted) history that begat a given model, along with the modelling choices long embedded in the code.
- variability in the scientific maturity of different components within a coupled earth system model.

These omissions mean that quantitative approaches cannot assess whether a model produces the right results for the wrong reasons, or conversely, the wrong results for the right reasons (where, say the observational data is problematic, or the model is configured to be unlike the earth system for a specific reason).

Furthermore, quantitative skill scores only assess specific versions of models, configured for specific ensembles of runs; they cannot reliably make any statements about other configurations built from the same code.





Quality as Fitness for Purpose

Evaluation of climate models should not be about "the model", but about the relationship between a modelling system and the purposes to which it is put. More precisely, it's about the relationship between particular ways of building and configuring models and the ways in which the runs produced by those models are used.



- To provide inputs to assessments of the current state of climate science:
- To explore the consequences of a current theory;
- To test a hypothesis about the observational system (e.g. forward modeling);
- To test a hypothesis about the calculational system (e.g. to explore known weaknesses);
- To provide homogenized datasets (e.g. re-analysis);
- To conduct thought experiments about different climates;
- To act as a comparator when debugging another model;





