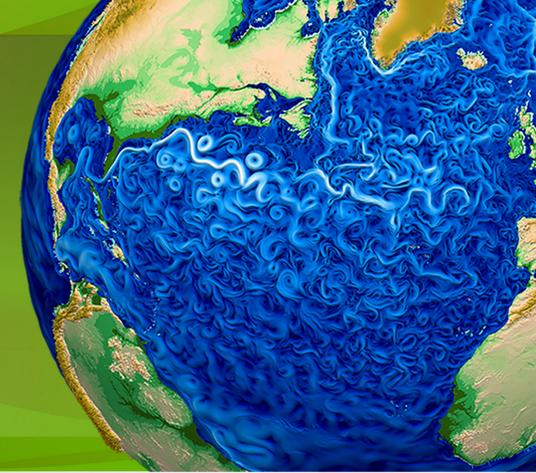


F: Climate reproducibility testing with EVE

Joseph H. Kennedy, Salil Mahajan, Katherine J. Evans, Peter Caldwell



Moving beyond bit-for-bit

Both internal (code) and external (machine) changes can affect a climate model's solution to a particular simulation

There are three types of changes:

1. Technical changes that continue to produce bit-for-bit **identical** solutions
2. Non-identical changes that produce a **statistically similar** solution
3. Changes that lead to a **different** solution

Only type 3 changes requires in-depth analysis of the changes, but there is no current capability to distinguish between type 2 and type 3 changes

We will enhance ACME's testing infrastructure to provide a robust climate reproducibility testing capability

Characteristics of successful testing

Concurrent to development

- Integrates into the development cycle
 - Useable, portable, flexible, extensible
 - Run frequently (easy) to continuously (scriptable)
- Minimal time to solution

Granular

- Functions → processes → components → model

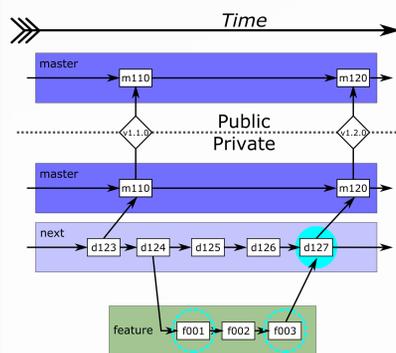
Informative

- Clear context
- Detailed analysis
- Appropriate metrics

Shareable

- Discussions across many institutions

ACME's development cycle



Developers' test suite

- Smaller set of tests; developers preferred machine
- Only tests feature changes
- Should exercise only feature relevant components

Integrators' test suite

- Larger suite of tests; multiple machines
- Tests integration of feature into next
- Should exercise whole model for unintended issues

Using EVE

```
(EVE-LIVVdev) fjk@pc0101123:~/Documents/Code/EVE/eve$ ./test -e extensions/crossmatch.json

EVE
Extended Verification and Validation for Earth System Models

Current run: 2017-06-02 14:11:12
User: fjk
OS Type: linux 4.4.0-78-generic
Machine: pc0101123

Previous output data found in output directory!
Backing up data to:
/home/fjk/Documents/Code/EVE/eve/vv_2017-06-02_20170602_125842

Beginning extensions test suite

Cross Match Test:
Null hypothesis: Accept
Critical value: 0.05
Test statistic: 0.8781241220279931

Extensions test suite complete

Done! Results can be seen in a web browser at:
/home/fjk/Documents/Code/EVE/eve/vv_2017-06-02/1index.html

(EVE-LIVVdev) fjk@pc0101123:~/Documents/Code/EVE/eve$
```

Above: EVE's command line interface was used to run the crossmatch test, one of the multivariate climate reproducibility tests that will be available. A general summary of the test results is displayed, and the path to the generated web output is given.

Right: EVE's web output showing detailed results of the crossmatch test, including a description of the test.

The screenshot shows the EVE web interface. At the top, it says 'EVE Home' and 'Extended Verification & Validation for Earth System Models'. Below that, there is a 'Multivariate crossmatch' section. The 'crossmatch' section contains a description of the test and a table of results. The table has columns for 'T', 'critical', 'h0', 'approval', 'set1', 'set2', 'a1', 'Ea1', 'Va1', and 'dev'. The results are as follows:

T	0.8781241220279931
critical	0.05
h0	accept
approval	0.7882201648686771
set1	1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35
set2	2,3,4,5,6,7,8,9,10,12,13,14,15,16,17,18,19,21,22,23,24,25,26,27,28,30,31,32,33,34,35
a1	19
Ea1	16.692307492307693
Va1	8.315581854043392
dev	0.80026084682832757

At the bottom of the page, there are logos for Oak Ridge National Laboratory, U.S. Department of Energy, and Office of Science. A note at the bottom states: 'EVE was developed by the CMDV Software project and funded by DOE BER under the CMDV program. Contact us on github'.

Integrating with ACME

Three main steps for each test:



1. Launch the test
 - Add a new ensemble test type to CIME
 - Strategy for each type of climate reproducibility test may be needed
2. Post process the test ensemble
 - Launch automatically when tests finish
 - Integrate with CIME and/or ACME post processing
3. Analyze the test results
 - Quickly tell if tests pass/fail
 - Detailed info on fail to help developers find bugs

Questions to be answered:

- When should they be run and by whom?
 - Developers' vs integrations' test suite
- Where are the tests applicable?
 - What types of issues can be identified
- What are the costs?
 - Both computational and personnel
- What do the developers need and want?
 - Successful tests are used