# Redefining reproducibility for E3SM on multicore systems

Kate Evans, presenter

Plus ORNL project members:
Salil Mahajan (lead),
Matt Norman,
Joe Kennedy,
Min Xu

# E3SM v1 release, July 2018
# Development for v2 is ongoing

- Model Components, each have new features in development:
  - **Atmosphere:** cloud microphysics, aerosols, variable resolution, etc. (EAM)
  - **Land:** biogeochemistry, soil hydrology, land units (ELM)
  - **Ocean:** dycore solvers, coupling to ice (MPAS-O)
  - **Land-ice:** new components (MALI, BISICLES)
  - Etc.

- All components have code updates in anticipation of new computing architectures
  - Code refactoring (Fortran + OpenACC & C++/Kokkos most common)
  - Consideration of new algorithms that favor less local memory, data transfer, efficiency
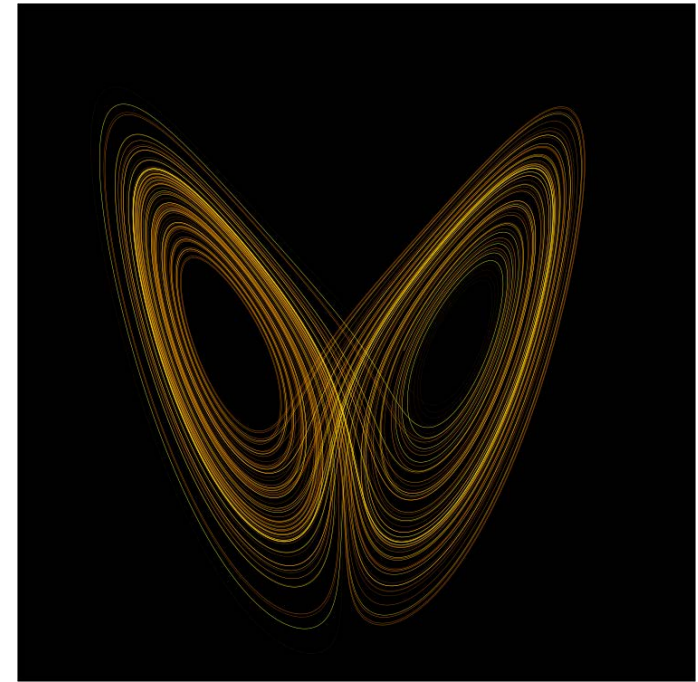
OAK RIDGE
National Laboratory

# There are categories of code updates:

- Changes that **do no**t affect the climate and **should be** bit-for-bit reproducible

  – E.g. Adding a new compset, inclusion of new output variables

- Changes that **do not** affect the climate and **will not be** bit-for-bit reproducible

  – E.g. code porting, GPU kernel, etc.
  – Climate statistics are the same

- Changes that **do** affect the climate and **will not be** bit-for-bit reproducible

  – E.g. New parameterizations modules, new tunings
  – Climate statistics are not the same

Goal: Test the null hypothesis that climate simulation is "similar".

**OAK RIDGE**
National Laboratory

3

# Motivation: Bit-for-bit is not achievable on target computing systems for E3SM

- Truncated Floating Point arithmetic:
  - Round-off differences
  - Non-associative:
    - $(-1 + 1) + 2^{-53} \neq -1 + (1 + 2^{-53})$
  - Optimizations, hybrid architectures, threading

- Climate models are chaotic and non-linear, so round-off differences grow quickly

- Goal: identify systematic bugs in a non-BFB reproducible environment **that allows for a reasonable development cycle**



Lorenz attractor
(*Source:en.wikipedia.org/wiki/Chaos_theory*)

OAK RIDGE
National Laboratory

Open slide master to edit

# Reproducibility Part 1: BFB

- ## E3SM Testing Suite:

  - * APT (auto promotion test (default length))
    * CME (compare mct and esmf interfaces (10 days))
    * ERB (branch/exact restart test)
    * ERH (hybrid/exact restart test)
    * ERI (hybrid/branch/exact restart test, default 3+19/10+9/5+4 days)
    * ERS (exact restart from startup, default 6 days + 5 days)
    * ERT (exact restart from startup, default 2 month + 1 month (ERS with info dbug = 1))
    * ICP (cice performance test)
    * LAR (long term archive test)
    * NCK (multi-instance validation vs single instance (default length))
    * NOC (multi-instance validation for single instance ocean (default length))
    * OCP (pop performance test)
    * P4A (production branch test b40.1850.track1.1deg.006 year 301)
    * PEA (single pe bfb test (default length))
    * PEM (pes counts mpi bfb test (seq tests; default length))
    * PET (openmp bfb test (seq tests; default length))
    * PFS (performance test setup)
    * PRS (pes counts hybrid (open-MP/MPI) restart bfb test from startup, default 6 days + 5 days)
    * SBN (smoke build-namelist test (just run preview_namelist and check_input_data))
    * SEQ (sequencing bfb test (10 day seq,conc tests))
    * SMS (smoke startup test (default length))
    * SSP (smoke CLM spinup test (only valid for CLM compsets with CLM45 and CN or BGC))

Um… what if its not BFB?

OAK RIDGE
National Laboratory

Open slide master to edit

# Reproducibility Part 2: Expert Opinion

- **Some years** of a control run
  - scientifically validated on a trusted machine

- **Some years** of the perturbed run

- **Expert opinion** from a subjective evaluation of plots, tables, etc.

- **Expensive, slow and subjective,** no quantitative standardized metric or cost function analysis.

- Although: simpler models had less complexity, fewer multiscale features



Careers "made" on showing the climate is "good enough" with new numerical dycores, packages, features (e.g. Evans et al 2013, 2014)

Open slide master to edit

# Isn't there a better way?

- Perturbation growth test (B. Singh, PNNL)
  - a la Rosinski and Williamson (1997)
  - Remove branching/bugs/RNG issues
  - Only one time step, analyze by process

- Time step convergence test (H. Wan, PNNL)
  - Fast; only requires several time steps of data
  - Cannot track errors outside the code where convergence is assessed

- Statistical consistency test (A. Baker, NCAR)
  - Needing only a few time steps for almost all testing
  - Assesses total code output
  - Hard to determine location in code, but being addressed with a code search strategy



Compared against Constance (Intel,O0)

Legend: Cascade(Intel-O0), Mira(XLF-O0), Mira(XLF-O3), YS(Intel-O3)

Y-axis: Maximum Error in the temperature field; X-axis: Process # during a time Step

Max T (K) difference evolution in various computing environments. Process indices shown on x-axis refer to different physics parameterizations and/or Fortran code modules executed within one model time step. (courtesy, B. Singh, ACME-SM project proposal)

OAK RIDGE
National Laboratory

Open slide master to edit

# Ensemble Based Multivariate Approach

- Closest to original "expert" method in terms of set up (climate modelers stay in their happy place)
  - This has pros and cons, but means the code is tested just as it runs

- Can also be used for scientific analysis
  - Already being used to analyze long term atmospheric patterns, model sensitivity and UQ for sensor networks

- Suites of statistical tests can be applied.
  - But which ones are best?

- Some tests provide the geographic location of outliers

**OAK RIDGE**
National Laboratory

Open slide master to edit

# Ensemble Based Multivariate Approach

*Goal: Accelerate and add rigor to the verification of E3SM for non-BFB changes*

- Approach:
  - Ensemble vs. ensemble
  - Short (~1 year) ensembles of control and perturbed runs

- Short Ensembles:
  - Quantify natural variability, span possible climate states
  - Better utilizes multicore machines (*Mahajan et al., 2017*)

- Leverage two sample equality of distribution tests:
  - e.g. cross-match test, energy test, kernel test
  - Distribution-free/non-parametric
  - Effective at high dimensions, low sample sizes
  - Used widely in other fields, e.g. genetics, image processing, etc.

**OAK RIDGE**
National Laboratory

# Short Independent Simulation Ensemble (SISE)

$$T'_j = (1+x')T_j$$

$x'$ is uniform random number transformed to range from $(-10^{-14}, 10^{-14})$



L1-norm of the absolute differences for hourly 850mb T (Kelvin)
*(Courtesy: Matt Norman)*

OAK RIDGE
National Laboratory

# Short Independent Simulation Ensemble (SISE)

Problem to solve: Multivariate two sample equality of distribution testing for high dimension, low sample size

# Packing simulations together is **economical** relative to a Single Long Run (SLR)

- Single Long run:
  - Less work per core with large core counts
  - Increase in MPI communications
  - Smaller MPI messages -> Large MPI latency
  - MPI cost > 90%

- 100 1-yr SISE vs. 100-yr long run
  - 100x greater workload per node on the same nodes
  - Faster throughput, and easy to use large core counts
    - Significantly reduced relative MPI and PCI-e overheads
    - Higher priority (the cool kids queue) on leadership class machines (e.g. Titan, Cori, etc.)



CAM5''0.25°''Titan'

*Strong scaling of a single long run. Courtesy: Mark Taylor and more, circa ~2012*

OAK RIDGE
National Laboratory

# Short Independent Simulation Ensemble (SISE)

- Example: EAM (E3SM atmosphere spectral element) two degree component:
  - SLR (100 years): 1536 elements given 96 nodes, 16 elements per node, takes **weeks** to finish
  - SISE (100 1yr runs): 1536 elements given 48 nodes each, 32 elements per node (total nodes: 4800), takes **less than a day** to finish
  - Took a while to analyze for success, we kept finding new bugs!
    - Random number generator was not so random
    - Restart bug for submonthly configurations for 3D variables



Temperature at 850 mb, first 10 years, arrow: model restart month

Open slide master to edit

# Test: Equality of Distribution

Kolmogorov Smirnov (KS) testing framework:

- Null Hypothesis ($H_0$): Two simulation ensembles (SISE) represent the same climate state.

- Use **global annual means** of all standard model output variables (**158 variables**)

- $H_0$: A variable between the two SISEs belong to the same distribution.

- Test $H_0$ for each variable using a KS test.

- Test statistic ($t$): No. of variables that reject $H_0$ at a given confidence level (say 95%).

- $H_0$ rejected if $t > a$, where $a$ is some critical number for a significance level (Type I error rate).

- $a$ is empirically from an approximate **null distribution of t** derived using **resampling** techniques.



Illustration: KS test

OAK RIDGE
National Laboratory

# Significance Level (Type I Error rate): Resampling

- Simulations from the two ensembles of size *n* and *m* are pooled together.

- Simulations from the pool are then randomly assigned to one of two groups of sizes *n* and *m*.

- The *t-statistic* is then computed for the random drawing.

- Repeat

- If all possible random drawings are made, the null distribution of *t* is exact.
  - We conduct 500 drawings - approximate null distribution.

**OAK RIDGE**
National Laboratory

# Model Verification Using SISE: E3SM v1
## Known Climate Changing Perturbation

- Configuration: Active atmosphere & land, prescribed cyclical F2000 SSTs and sea-ice distribution (FC5)

- Spatial Resolution: ~500km at the equator (5 degrees), 30 vertical layers
- Machine Configuration: PGI compiler on Titan
- Ensembles: Machine-precision level random perturbations to the initial 3-D temperature field
  - 30 member SISE
  - $T'_j = (1+x')T_j$, $x'$ is random number transformed to range from (-10$^{-14}$, 10$^{-14}$)

- Turn a tuning parameter knob: zm_c0_ocn (control case: 0.007, modified: 0.045)

**OAK RIDGE**
National Laboratory

# KS Testing Framework Results

| Name | Description | Ens. Size |
|------|-------------|-----------|
| Default c0_ocn | Default model settings | 30 |
| Perturbed c0_ocn | Perturbed model parameter | 30 |

| Comparison | Test Statistic (t) | Critical No. | H0 Test |
|------------|-------------------|--------------|---------|
| Default vs. perturbed c0_ocn | 119 | 13 | Reject |

OAK RIDGE
National Laboratory

# Power Analysis: KS Testing Framework



Power Analysis of KS Testing Framework

Fewer ensembles mean less sensitivity. How well do we know how sensitive the world is to changes on forcing?

Oak Ridge National Laboratory

# Model Verification Using SISE: Compiler optimization choices with E3SM v0

- **Configuration**: Preindustrial, active atmosphere (CAM5) and land (CLM4)
- **Spatial Resolution**: 208km at the equator (2 degrees), 30 vertical layers
- **Machine Configuration**: PGI compiler on Titan
- **Ensembles**: Machine-precision level random perturbations to the initial 3-D temperature field

| Name | Description | Ens. Size |
|---|---|---|
| SLR | Long control simulation (100 years, -O2 optimization) | 1 |
| SISE-DEFAULT | Short 1-yr simulation ensemble with default (-O2) optimization | 65 |
| SISE-O1 | Short 1-yr simulation ensemble with -O1 optimization | 59 |
| SISE-FAST | Short 1-yr simulation ensemble with -fast optimization | 62 |
| SISE-LND-INIT | Short simulation ensemble with land initialized with states from 70 different years of the SLR | 70 |

**OAK RIDGE**
National Laboratory

# Compiler optimization choices

## KS Testing Framework Results

| Comparison | Test Statistic ($t$) | Critical Value ($\square$) | $H_0$ Test |
|---|---|---|---|
| SISE-DEFAULT vs. SISE-O1 | 1 (0.6%) | 17 | Accept $H_0$ |
| SISE-DEFAULT vs. SISE-FAST | 24 (15.2%) | 14 | Reject $H_0$ |
| SISE-O1 vs. SISE-FAST | 23 (14.6%) | 16 | Reject $H_0$ |

Aggressive compiler choices (SISE-FAST) with the PGI compiler on Titan **can** result in climate-changing simulations.

OAK RIDGE
National Laboratory

# Extended Verification and Validation for E3SM:

- Python based toolkit:
  - Runs control and perturbed ensembles
  - Post-processes model output
  - Conducts tests
  - Publishes results and auxiliary plots, tables

OAK RIDGE
National Laboratory

# Summary:

- Short runs and ensembles are the only viable path for model verification as model expense grows
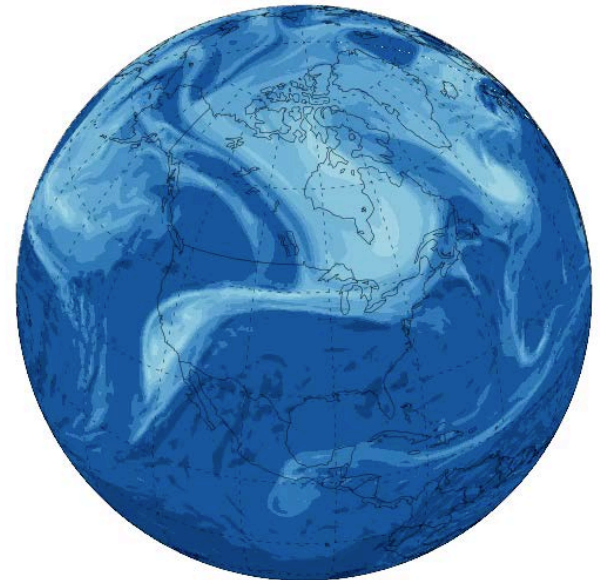
- A multivariate testing framework (EVE) is presented for climate reproducibility:

- We demonstrated this with known climate changing perturbations (and provided detection limits), choice of compiler optimization, and verifying how frozen the model was after months of software updates

- Future work:
  - Evaluate applicability of low-resolution results at high-resolution
  - Apply to shorter runs (monthly and daily vs. yearly)
  - Optimize multivariate tests, e.g. use different kernel functions, distance metrics



CONTOUR FROM -.000001 TO .00002 BY .000001

Potential vorticity at 300mb, E3SM present day test run

**OAK RIDGE**
National Laboratory

Open slide master to edit

# We are hiring! If you have expertise* in one or more of the following I would like to talk to you

- A passion for coding for >petaflop systems

- Understanding of modeling the atmosphere
  - Dynamics
  - clouds

- Understanding of modeling the ocean

- Understanding software ecosystems and good habits

- Diversity in every dimension

\* i.e. ninjas

# Model Verification Using SISE:
# Frozen model configuration v0 vs. v1

- **Configuration**: F1850C5 compset (frozen after v0 bug-fixes, v0.4)
- **Spatial Resolution**: 208km at the equator (2 degrees), 30 vertical layers
- **Ensembles**: Machine-precision level random perturbations to the initial 3-D temperature field

- Goal: Evaluate if efforts towards exascale computing impact climate reproducibility:
  - New scientific features, code refactoring
  - CIME (Common Infrastructure for Modeling the Earth System) update
  - Compiler and Software library updates

| Name | Ens. Size | CIME | PGI | p-netcdf |
|---|---|---|---|---|
| **v0.4-2015** | 30 | 4.0 | 15.3 | 1.5.0 |
| **master** | 30 | 5.0 | 17.5 | 1.7.0 |
| **v0.4** | 27 | 4.0 | 17.5 | 1.7.0 |

OAK RIDGE
National Laboratory

# Frozen model configuration v0 vs. v1

| Comparison | Test Statistic (t) | Critical no. (α) | H0 Test |
|---|---|---|---|
| v0.4-2015 vs. master | 6 (3.6%) | 13 | Accept H0 |
| v0.4 vs. master | 8 (4.2%) | 13 | Accept H0 |
| v0.4-2015 vs. v0.4 | 5 (3%) | 13 | Accept H0 |

Software infrastructure updates are not climate changing.
Frozen model configuration reproducible!

OAK RIDGE
National Laboratory

# Short Ensembles: Scientific Utility



Control Case (1850S) — Perturbed Case (2000S)

Jan 0001   Jan 0002   Jan 0003   Jan 0004   Jan 0005   Jan 0081

Fast Response

**SST (2000S – 1850S)**          **Precipitation (2000S – 1850S)**

OAK RIDGE
National Laboratory

# Test for Extremes

- Distribution tests perform poorly on distribution with different tails
  - Known for univariate tests, unexplored for multivariate tests.

- Use Generalized Extreme Value (GEV) theory (*e.g. Mahajan et al. 2015, Evans et al. 2014*).
  - max./min. of a process belong to GEV distribution.
  - Analogous to central limit theorem
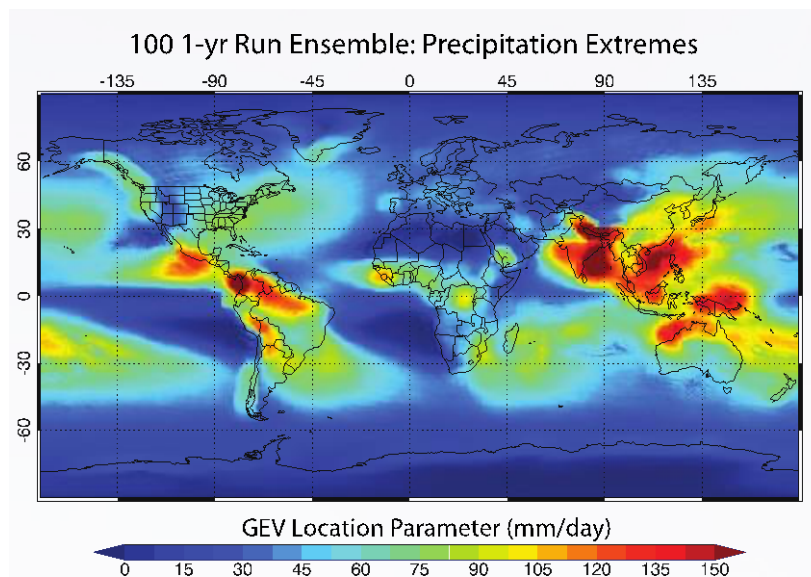  - GEV parameters normally distributed asymptotically



Generalized extreme value densities

All with $\mu = 0$, $\sigma = 1$. Asterisks mark support-endpoints

$$G(z) = \exp\left\{ -[1 + \xi(\frac{z - \mu}{\sigma})]^{-1/\xi} \right\}$$

$$z : 1 + \xi(z - \mu)/\sigma > 0$$

where $\mu$, $\sigma$ and $\xi$ represent the location, scale and shape parameter respectively.



100 1-yr Run Ensemble: Precipitation Extremes
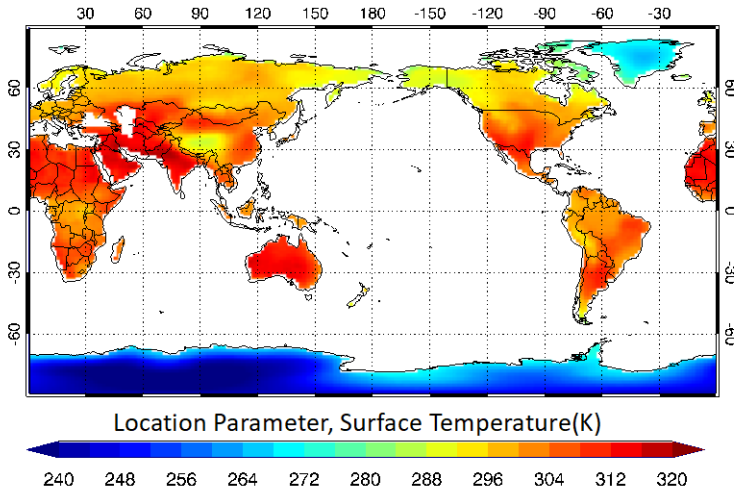
GEV Location Parameter (mm/day)

# Climate Extremes Test

- Null Hypothesis ($G_0$): Simulation of extremes of a variable between two SISE is statistically indistinguishable.

- Annual maxima for each grid point are fit to a GEV distribution.

- $G_0$: Extremes at each grid point are statistically indistinguishable

- Test statistic ($g$): No. of grid points that reject $G_0$

- $G_0$ rejected if $t > b$, where $b$ is some critical number, obtained using resampling techniques.
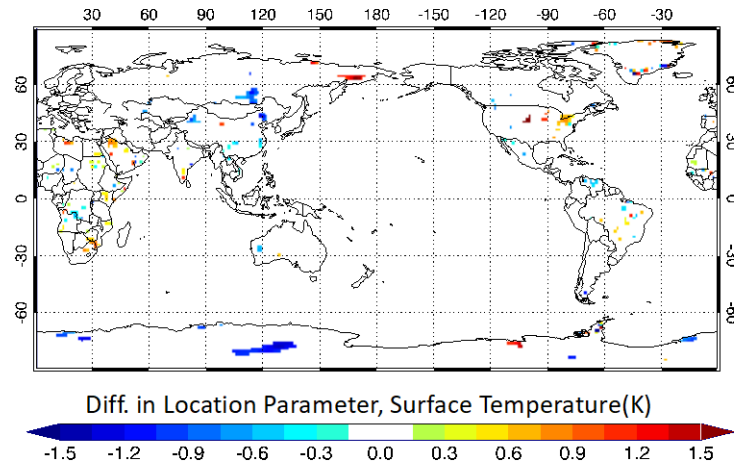
**OAK RIDGE**
National Laboratory

# Climate Extremes

a.
### Surface Temperature Extremes: Default

Location Parameter, Surface Temperature(K)

240  248  256  264  272  280  288  296  304  312  320

c.
### Precipitation Extremes: Default

Location Parameter, Precipitation Rate (mm/day)

0  6  12  18  24  30  36  42  48  54  60

b.
### Default – O1

Diff. in Location Parameter, Surface Temperature(K)

-1.5  -1.2  -0.9  -0.6  -0.3  0.0  0.3  0.6  0.9  1.2  1.5

d.
### Default – O1

Diff. in Location Parameter, Precipitation Rate (mm/day)

-15  -12  -9  -6  -3  0  3  6  9  12  15

to edit

# Climate Extremes

| Comparison | Variable | Test statistic ($g$) | Critical value ($\beta$) | $G_0$ Test |
|---|---|---|---|---|
| SISE-DEFAULT vs. SISE-O1 | Precipitation Rate | 5.1% | 6.5% | Accept $G_0$ |
| | Surface Temperature | 5.0% | 9.6% | Accept $G_0$ |
| SISE-DEFAULT vs. SISE-FAST | Precipitation Rate | 4.7% | 6.3% | Accept $G_0$ |
| | Surface Temperature | 3.6% | 9.6 % | Accept $G_0$ |
| SISE-O1 vs. SISE-FAST | Precipitation Rate | 5.2% | 6.5% | Accept $G_0$ |
| | Surface Temperature | 10.3% | 9.8% | Reject $G_0$ |

- All SISE simulations are identical to each other in terms of their simulation of climate extremes.
- The result is in contrast to the result of the KS-testing framework.
- It suggests that either optimization choices do not effect climate extremes, or
- Climate extremes are not a good metric to evaluate answer changes that might effect the simulation of the climate, with 60 ensemble members.

**OAK RIDGE**
National Laboratory

Open slide master to edit

# Single Long Run (SLR) vs. SISE

- SLR is clearly distinct from the SISE-DEFAULT

## KS Testing Framework Results

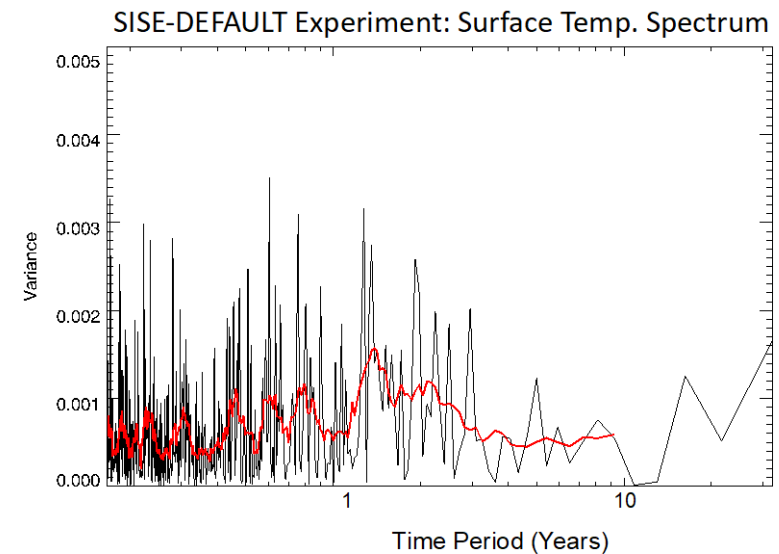| Comparison | Test Statistic ($t$) | Critical Value ($\square$) | $H_0$ Test Result |
|---|---|---|---|
| SLR vs. SISE-DEFAULT | 80 (50.6 %) | 15 | Reject $H_0$ |
| SLR vs. SISE-LND-INIT | 74 (48 %) | 13 | Reject $H_0$ |

OAK RIDGE
National Laboratory

# SLR vs. SISE

- Atmospheric models show that free atmospheric-only internal variability can include variability on longer time-scales (*e.g. James and James, 1989, Lorenz, 1990, Held, 1993, Marshall and Molteni, 1993*).

- This low frequency variability is not captured by SISE.
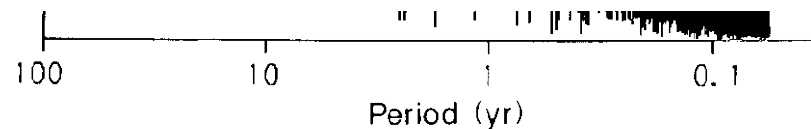
# Atmospheric Low-frequency Variability



FIG. 2 Spectrum of the time variations of two selected spherical harmonic coefficients during the last 96 years of the model integrations. *a*, The $Y_1^0$ temperature coefficients, averaged with respect to pressure, giving the

temperature contrast. The forcing was applied to this spherical harmonic component to simulate the annual cycle, which introduced a sharp peak at a period of one year. We might have expected baroclinic instability to lead directly to a large variance on the 1–7 day timescale but, after the annual cycle, ultra-low frequencies had the largest amplitudes in the spectrum. The maximum was for a period of 12 years. Figure 2*b*, shows the spectrum of the $Y_1^0$ vorticity coefficient which measures the solid-body rotation component of the atmospheric motion relative to the Earth. This coefficient was not forced directly, but varied as the temperature variations set the model atmosphere into motion. The seasonal cycle was expressed by a strong peak

*James and James, Nature, 1989*

Oak Ridge National Laboratory

Open slide master to edit

# Multivariate Cross-Match Test

- $n$ 1-yr control runs (~C)

- $m$ 1-yr modified runs (~M)

- Coarse grained: global annual means

- Multivariate vector for each run (size ~130)

- Pool vectors, $N = n+m$

- Pair vectors based on min. Mahalanobis distance

- $H_0$: C = M

- Test-statistic ($T$):
  - No. of pairs with one control run and one perturbed run

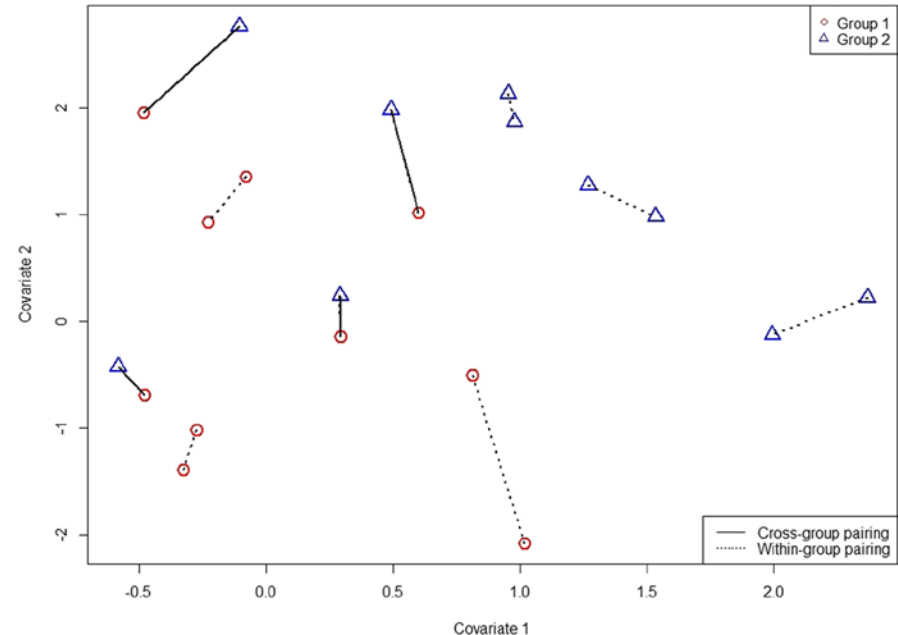- Test the null hypothesis using the exact null distribution



Illustration of cross matching for a bivariate case with $n = m = 10$. *(Ruth, 2014)*

Open slide master to edit

# Cross-Match Test

- Null distribution of T-statistic:

$$P(T = a_1) = \frac{2^{a_1}(N/2)!}{\binom{N}{n}\left(\frac{n-a_1}{2}\right)!\, a_1!\, \left(\frac{m-a_1}{2}\right)!}$$

  – i.e. when both samples belong to the same population

  – where $a_1$ is the no. of pairs with one control and one perturbed vector

  – Based on simple combinatorial arguments, thus exact
    - Analogous to the probability of drawing one red and one green ball
  – For e.g. for $n = m = 9$, $P(a_1 \leq 1) = 0.0259$

Open slide master to edit