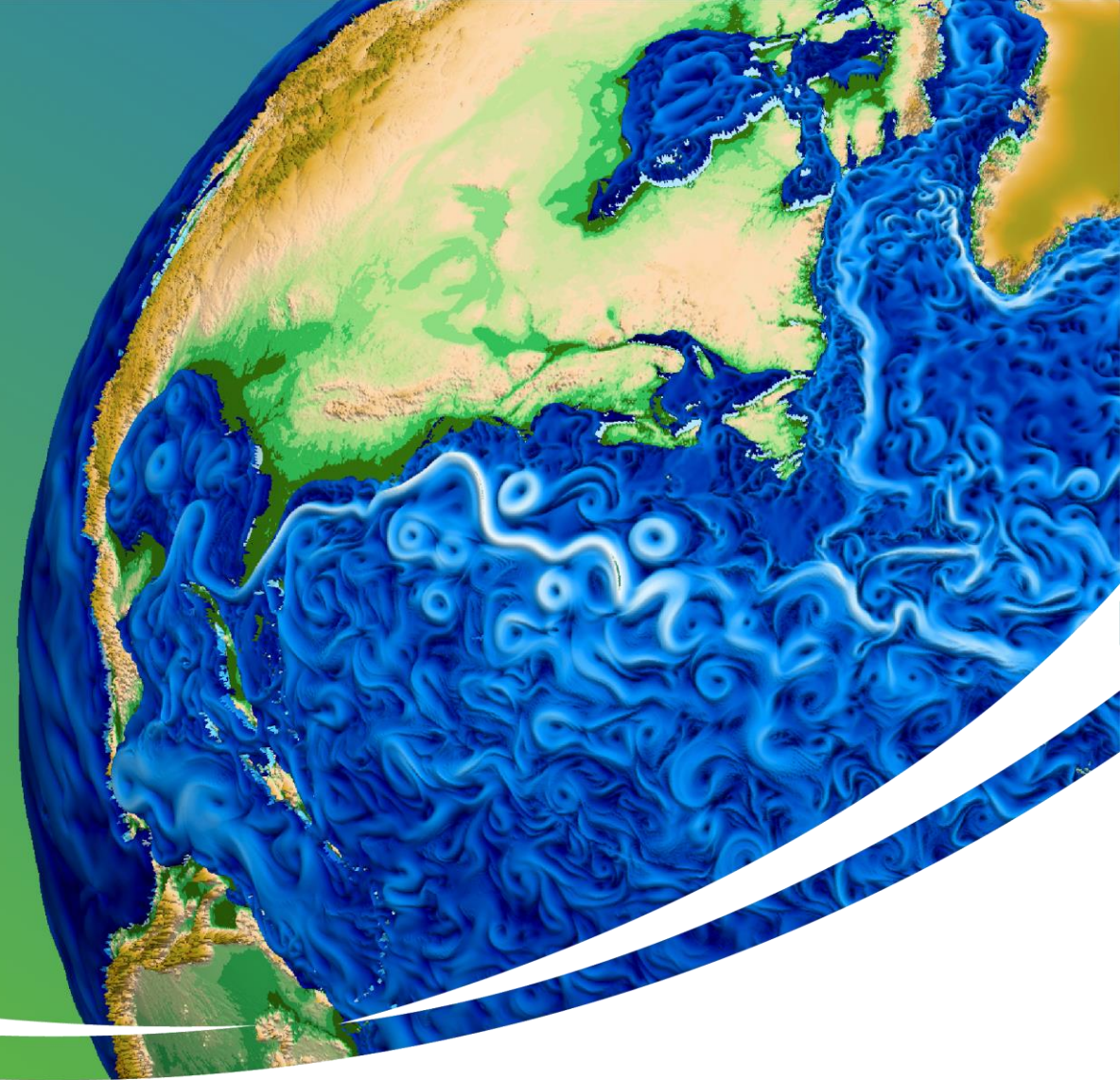


High Performance Computing for Climate

Sarat Sreepathi

Oak Ridge National Laboratory

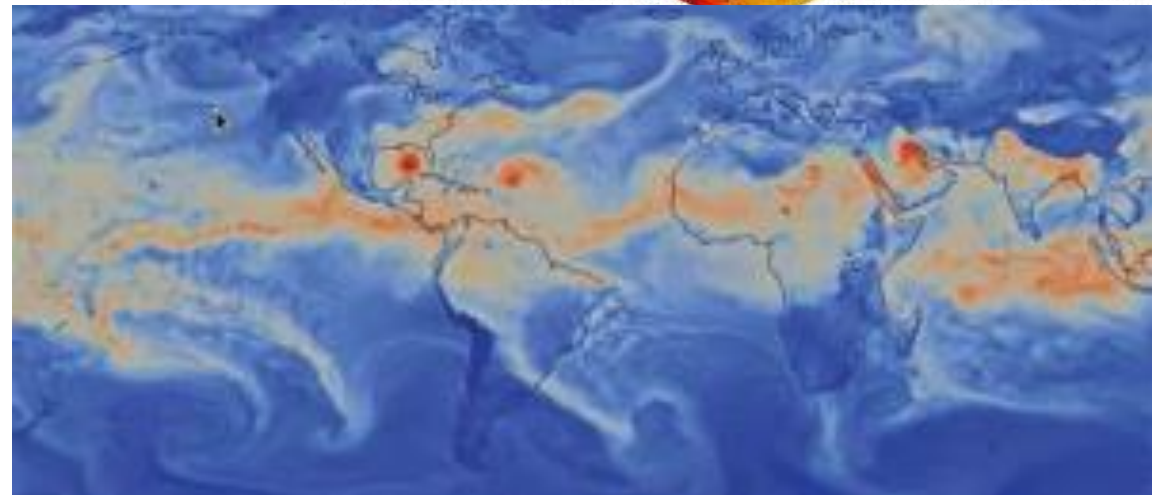
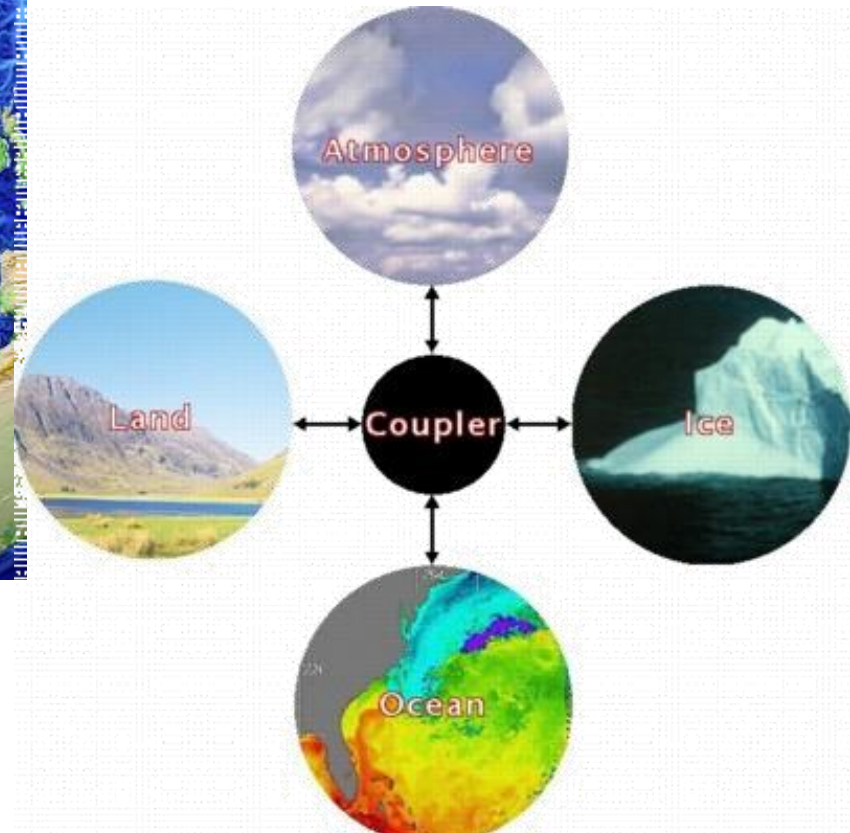
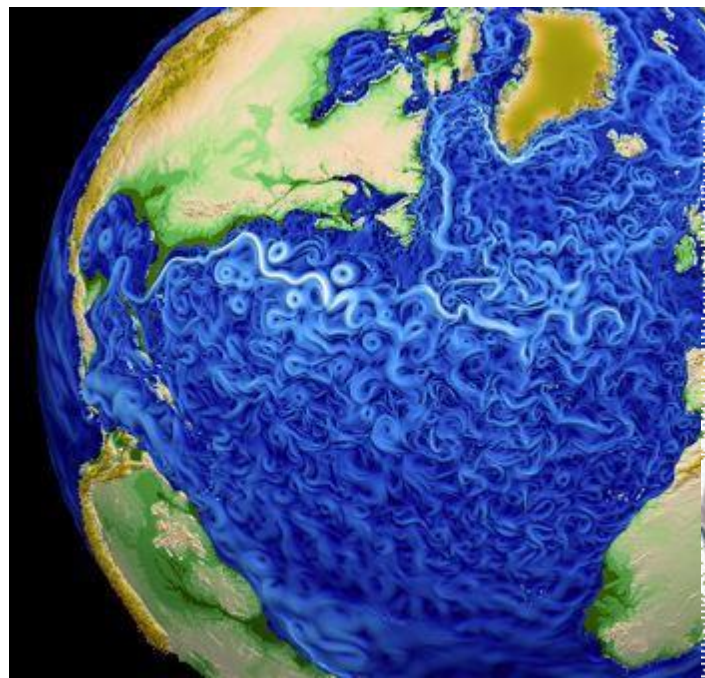
Performance Analytics for Computational Experiments



Energy Exascale Earth System Model (E3SM)

- Global Earth System Model
- Atmosphere, Land, Ocean, Ice, ... component models
- 8 DOE labs, 12 university partners, ... ~\$30+ M/year
- Development driven by DOE mission interests: Energy/water issues looking out 40 years
- **Key computational goal: Ensure E3SM effectively utilizes DOE exascale supercomputers**
- E3SM is open source / open development
 - Website: www.e3sm.org
 - Github: <https://github.com/E3SM-Project>

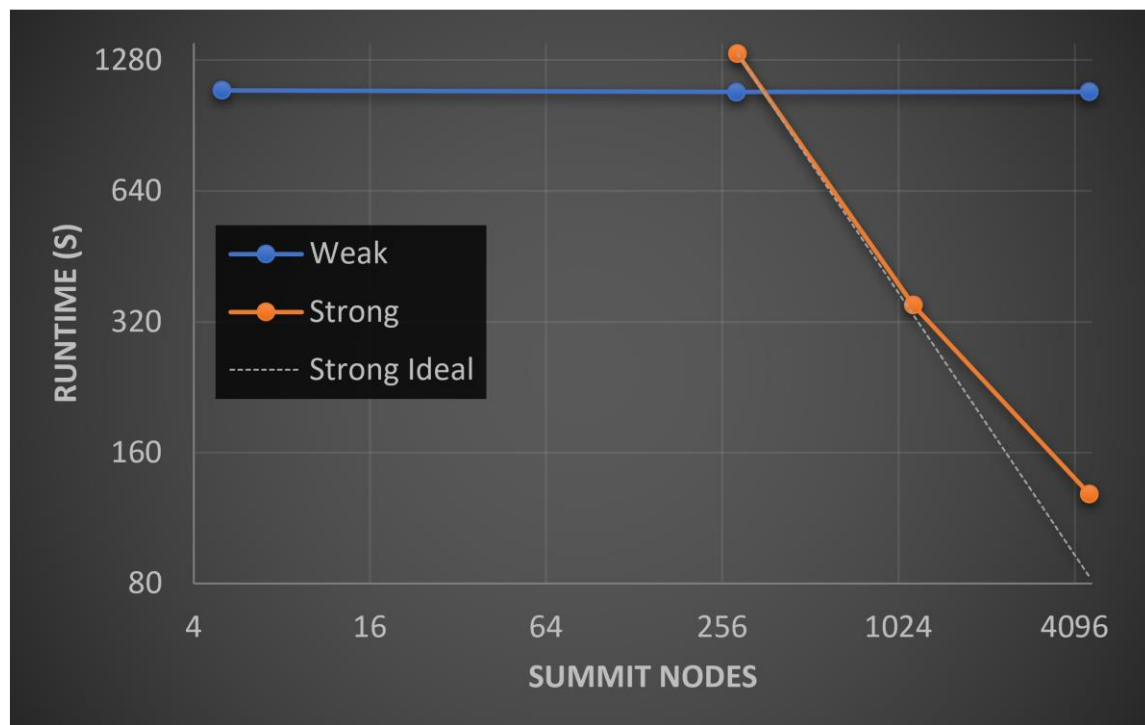
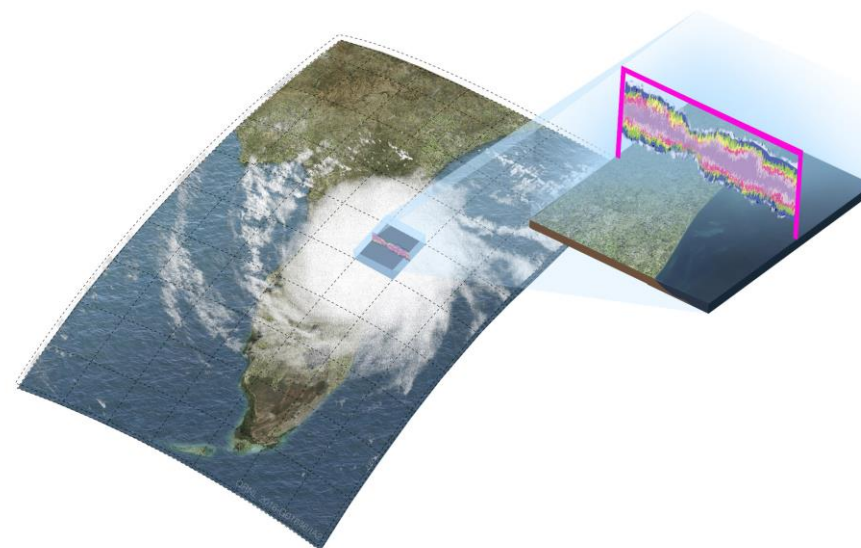
Mission: Use exascale computing to carry out high-resolution Earth system modeling of natural, managed and man-made systems, to answer pressing problems for the DOE.



E3SM-MMF Cloud Resolving Climate Model



Goal: Develop capability to assess regional impacts of climate change on the water cycle that directly affect the US economy such as agriculture and energy production.



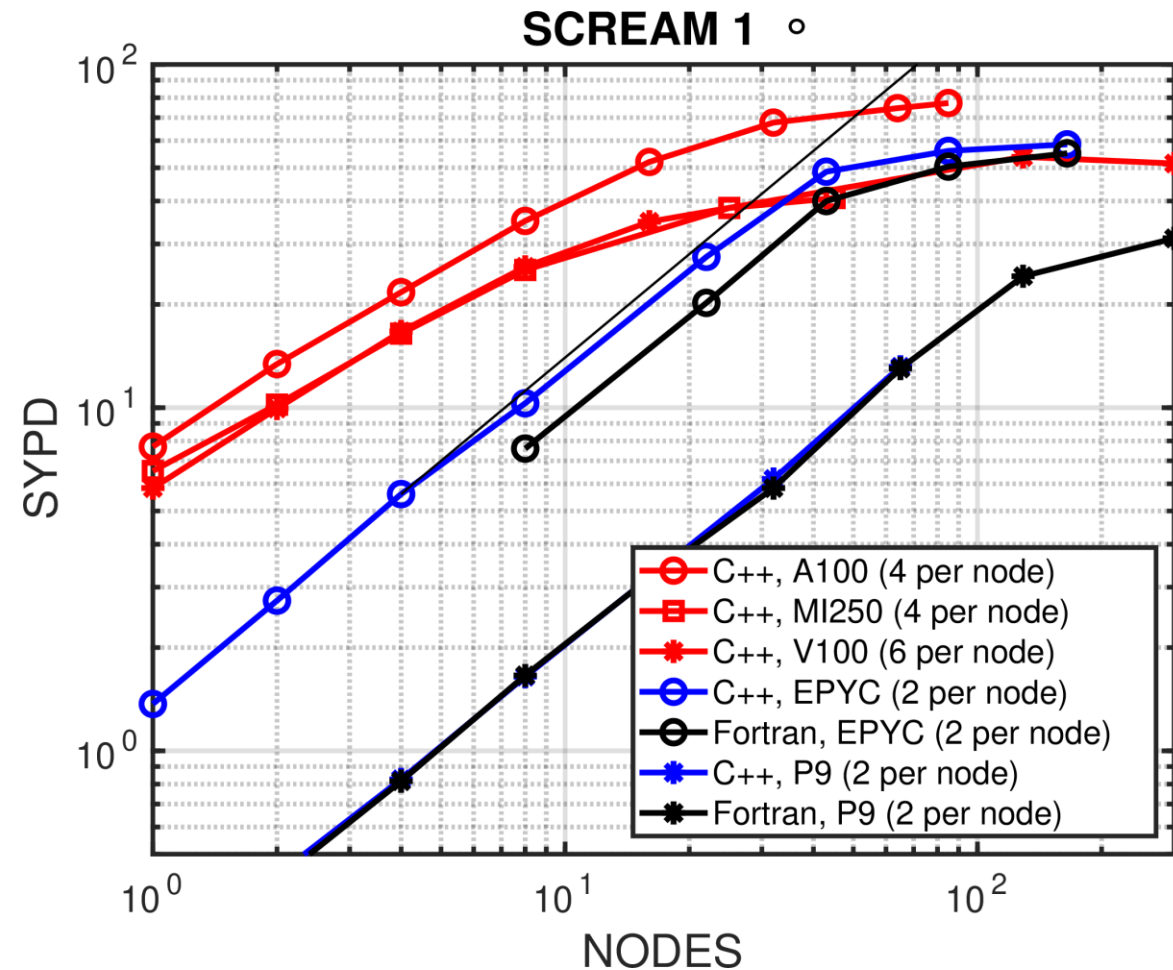
- Multiscale Modeling Framework (MMF) / Super-Parameterization
- Replaces traditional parameterizations with cloud resolving model within each grid cell of global climate model

Programming Models

- C++ with templates (Kokkos or YAKL)
 - Robust and well supported solution across most hardware
 - Requires minimal vendor support
- Fortran with OpenACC or OpenMP offload
 - Relies heavily on (lagging) vendor compiler support
 - Remains immature w.r.t. advanced Fortran features
 - Good performance requires major code refactoring
- Domain Specific Languages
 - Promising approach (e.g. GT4Py/GridTools, PSyclone)
 - Need additional investments to support algorithms & meshes in E3SM components
 - Most experience within DOE labs is with C++

E3SM's Atmosphere model (EAMXX in "SCREAM" configuration)
 1 degree resolution: 128 vertical levels, nonhydrostatic (NH) dycore, 10 tracers, P3/SHOC physics with prescribed aerosols, no convective parameterization

- Performance portability
 - IBM P9, AMD EYPC
 - NVIDIA V100, A100
 - AMD MI250
- CPU performance:
 - C++/Kokkos as fast or faster than Fortran
- GPU performance:
 - Large scaling range where GPU nodes are 4-10x faster than CPU nodes



Early Evaluation of Fugaku A64FX Architecture Using Climate Workloads

Sarat Sreepathi
Oak Ridge National Laboratory

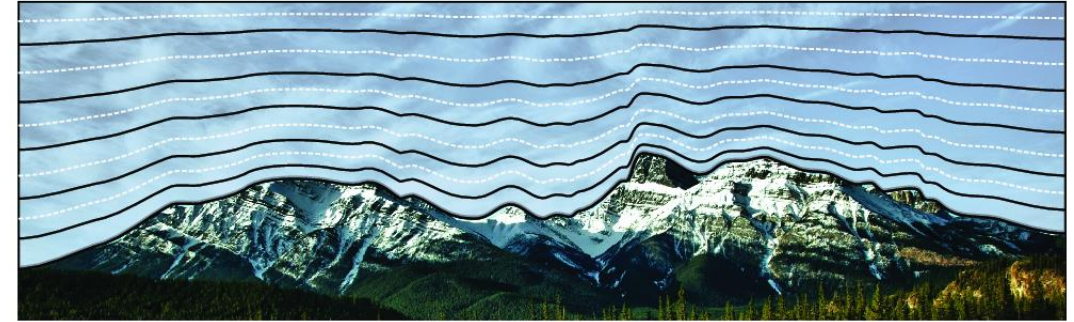
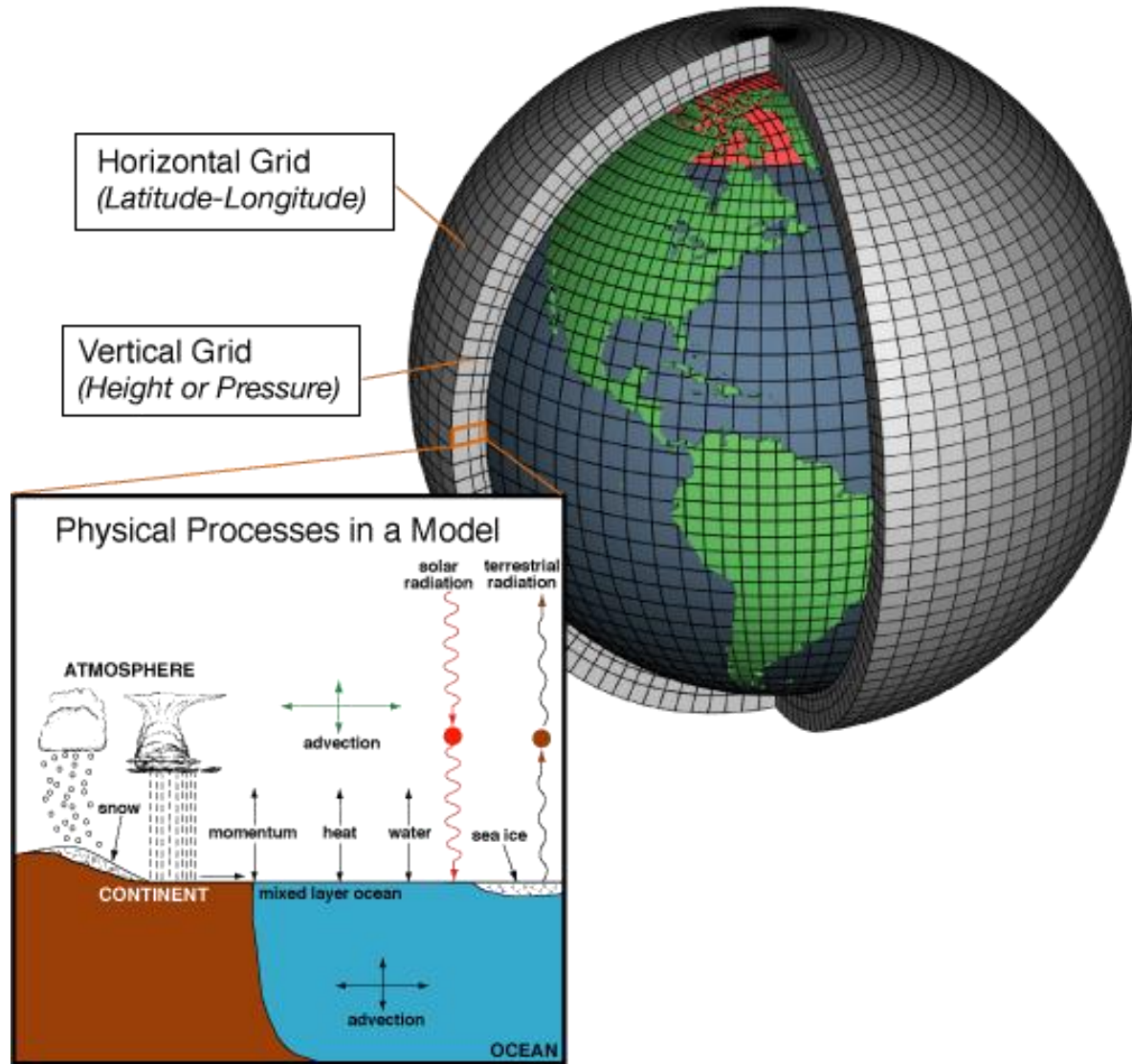
Mark Taylor
Sandia National Laboratories

Adapted from talk given at
EAHPC Workshop
IEEE Cluster 2021
September 7, 2021



EXASCALE COMPUTING PROJECT

Atmosphere Component



hydrostatic-pressure terrain-following coordinates

- Dynamical Core
 - Solves the Atmospheric Primitive Equations
 - Linear transport of 40 atmospheric species
 - 72 vertical levels – 0.8 km avg. spacing
 - Benchmark (two versions): Fortran (preqx) and C++ (preqx_kokkos)

Terrain following figure: D. Hall, CU Boulder

Source: http://celebrating200years.noaa.gov/breakthroughs/climate_model/welcome.html

Fugaku



HOME LISTS STATISTICS RESOURCES ABOUT MEDIA KIT

Home »RIKEN Center for Computational Science »
Supercomputer Fugaku - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu...

SUPERCOMPUTER FUGAKU - SUPERCOMPUTER FUGAKU, A64FX 48C 2.2GHZ, TOFU INTERCONNECT D

Site:	RIKEN Center for Computational Science
System URL:	https://www.r-ccs.riken.jp/en/fugaku/project
Manufacturer:	Fujitsu
Cores:	7,630,848
Memory:	5,087,232 GB
Processor:	A64FX 48C 2.2GHz
Interconnect:	Tofu interconnect D
Performance	
Linpack Performance (Rmax)	442,010 TFlop/s
Theoretical Peak (Rpeak)	537,212 TFlop/s
Nmax	21,288,960
HPCG [TFlop/s]	16,004.5
Power Consumption	
Power:	29,899.23 kW (Optimized: 26248.36 kW)
Power Measurement Level:	2

<https://www.top500.org/system/179807/>

- #2 on Top500
- RIKEN Center for Computational Science
- Key Characteristics of A64FX*
 - Arm 64-bit with 512-bit SVE (Scalable Vector Extensions)
 - High Bandwidth Memory
 - Low Power

*https://www.fujitsu.com/downloads/SUPER/a64fx/a64fx_datasheet_en.pdf



Architecture Comparison: Metrics

- Single node workload for understanding h/w trends (ca 2012+)
- Performance Efficiency metric: number of element remap timesteps **per second**

$$E_{perf} = \frac{N_e * N_t}{(prim_main_loop * num_devices)}$$

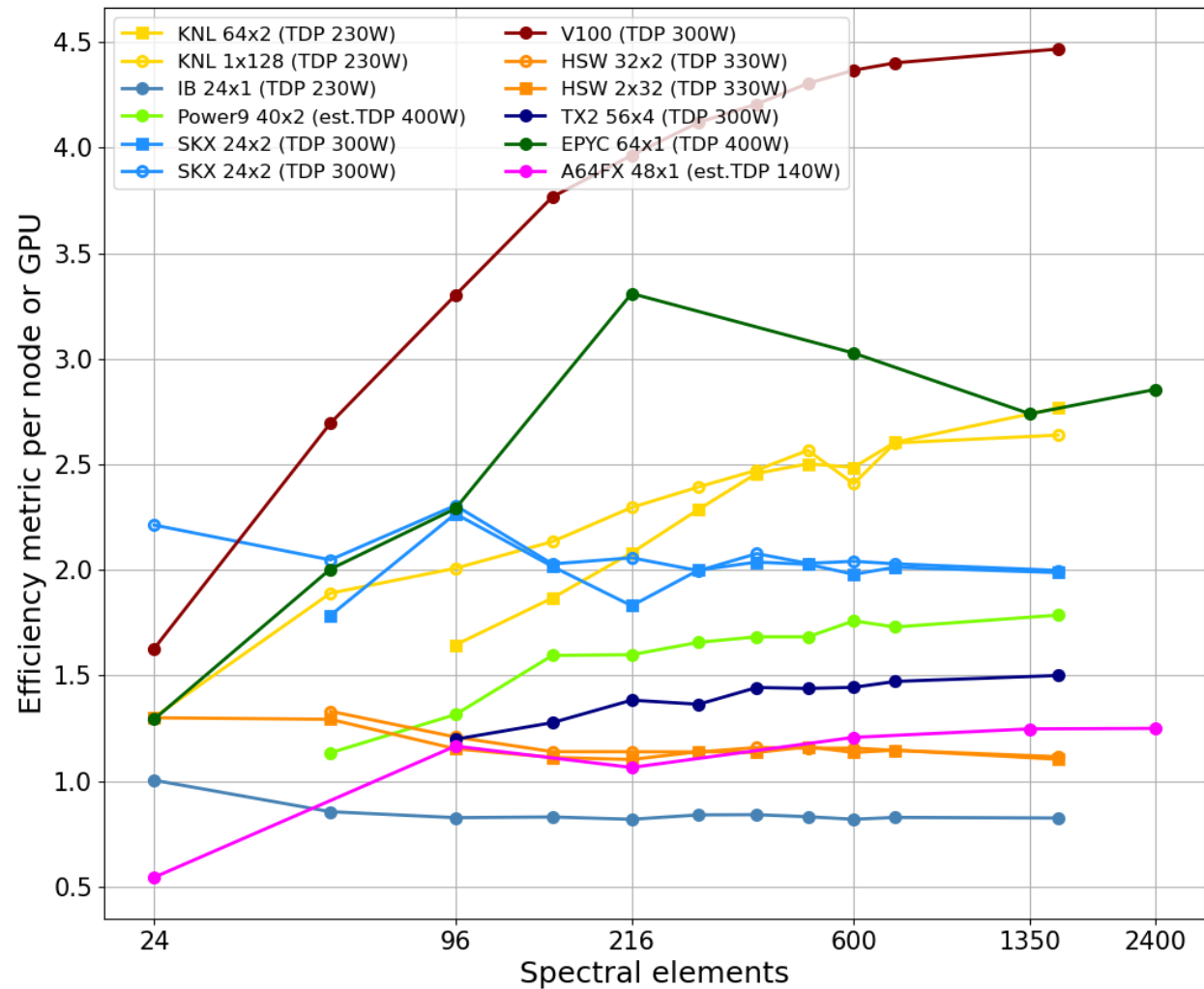
- N_e is the number of spectral elements
 - N_t is the number of remap timesteps (34 for the Fugaku experiments)
 - *prim_main_loop* is the main computation loop timer
 - *num_devices* is 1 for CPU nodes or the number of GPUs per node for GPU systems
- Power Efficiency metric: number of element remap timesteps **per Watt**

$$E_{power_tdp} = \frac{E_{perf}}{TDP}$$

- Thermal Design Power (TDP)

Architecture Comparison: Performance Efficiency

E3SM HOMME Dycore Benchmark: Cross-Architecture Comparison
(A64FX, EPYC results are preliminary)



Inform configurations where GPU systems can outperform CPU systems

Fugaku Node: Single A64FX socket
GNU Fortran + MPI (48 ranks)

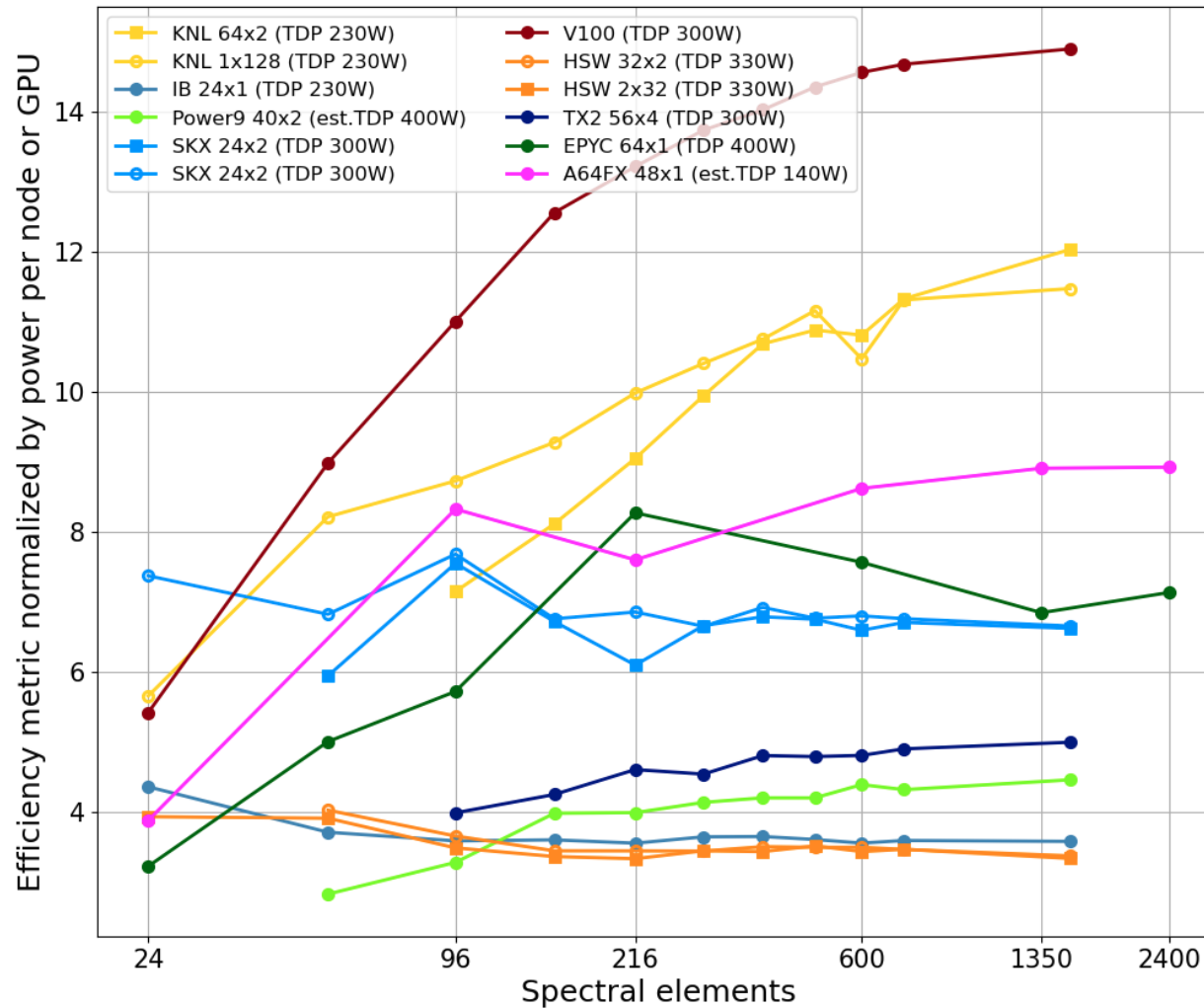
Note: Top Red (Volta V100), Pink (A64FX),
Orange (Dual-socket Haswell)
Higher is better

* Plot of the efficiency metric normalized by power consumption on various hardware architectures. The legend includes a short descriptor for each architecture along with the number of parallel processes times (x) the number of threads and includes TDP in parenthesis. Specifically, the labels map as follows: KNL (Intel® Knights Landing), IB (Intel®Ivy Bridge), SKX (Intel® Skylake), V100 (NVIDIA® Volta), HSW (Intel® Haswell), A64FX (Fujitsu® A64FX), Power9 (IBM® POWER9), TX2 (Marvell®ThunderX2), EPYC (AMD® EPYC).

Architecture Comparison: Power Efficiency

E3SM HOMME Dycore Benchmark: Cross-Architecture Comparison

(A64FX, EPYC results are preliminary)



A64FX: Promising performance/watt

Fugaku Node: Single A64FX socket
GNU Fortran + MPI (48 ranks)

Note: Top Red (Volta V100), Pink (A64FX), Yellow (KNL)
Higher is better

* Plot of the efficiency metric normalized by power consumption on various hardware architectures. The legend includes a short descriptor for each architecture along with the number of parallel processes times (x) the number of threads and includes TDP in parenthesis. Specifically, the labels map as follows: KNL (Intel® Knights Landing), IB (Intel® Ivy Bridge), SKX (Intel® Skylake), V100 (NVIDIA® Volta), HSW (Intel® Haswell), A64FX (Fujitsu® A64FX), Power9 (IBM® POWER9), TX2 (Marvell® ThunderX2), EPYC (AMD® EPYC).

Power and Performance tradeoffs

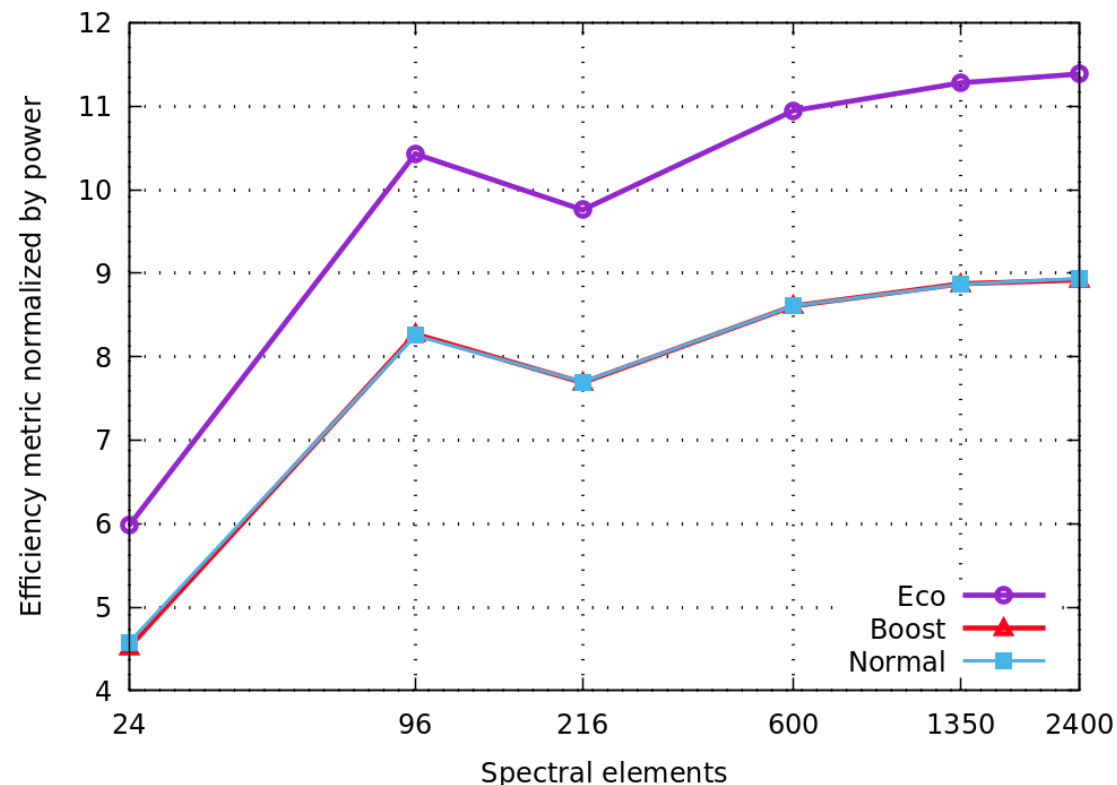
- Power Efficiency metric normalized by the measured power on the compute node

$$E_{power_measured} = \frac{E_{perf}}{measured_power}$$

- PowerAPI
- Three modes
 - Normal (2 GHz)
 - Boost (2.2 GHz)
 - **Eco** (2 GHz/`eco_state=2`)
- Fortran version with GNU

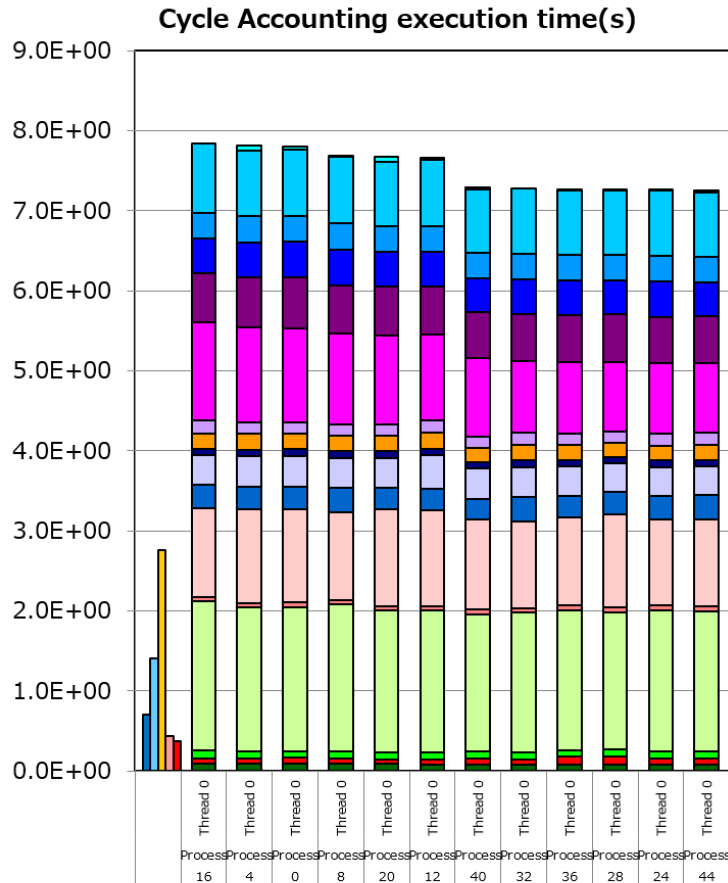
Higher is better

E3SM HOMME Dycore Benchmark: Power modes on Fugaku
(Single node: 48 MPI ranks)

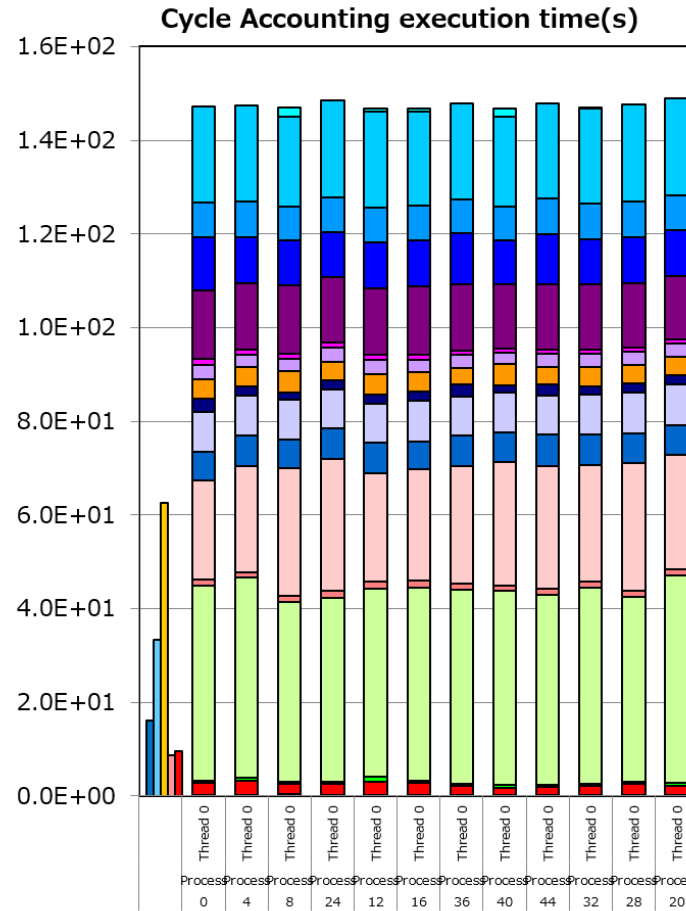


Performance Characterization: Instruction mix

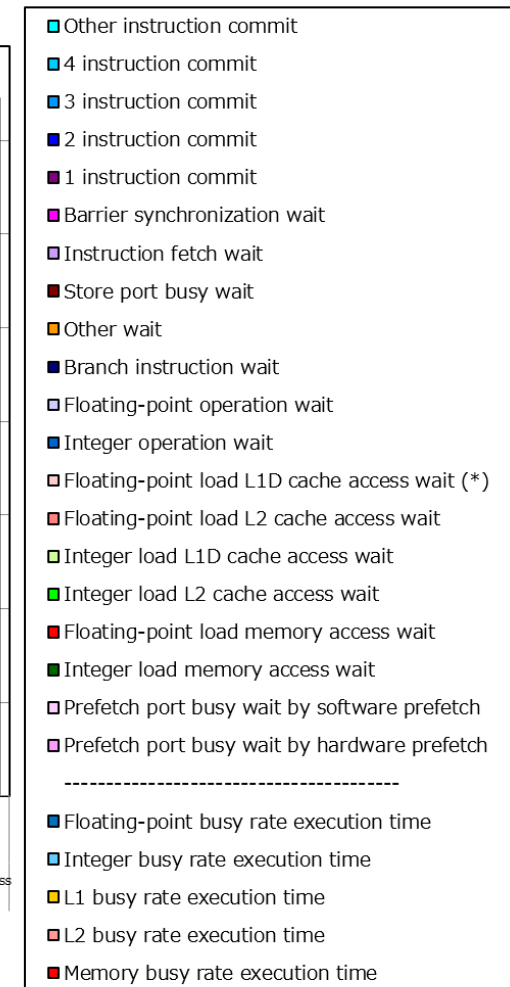
Spectral Elements: 96



Spectral Elements: 2400



20 Categories



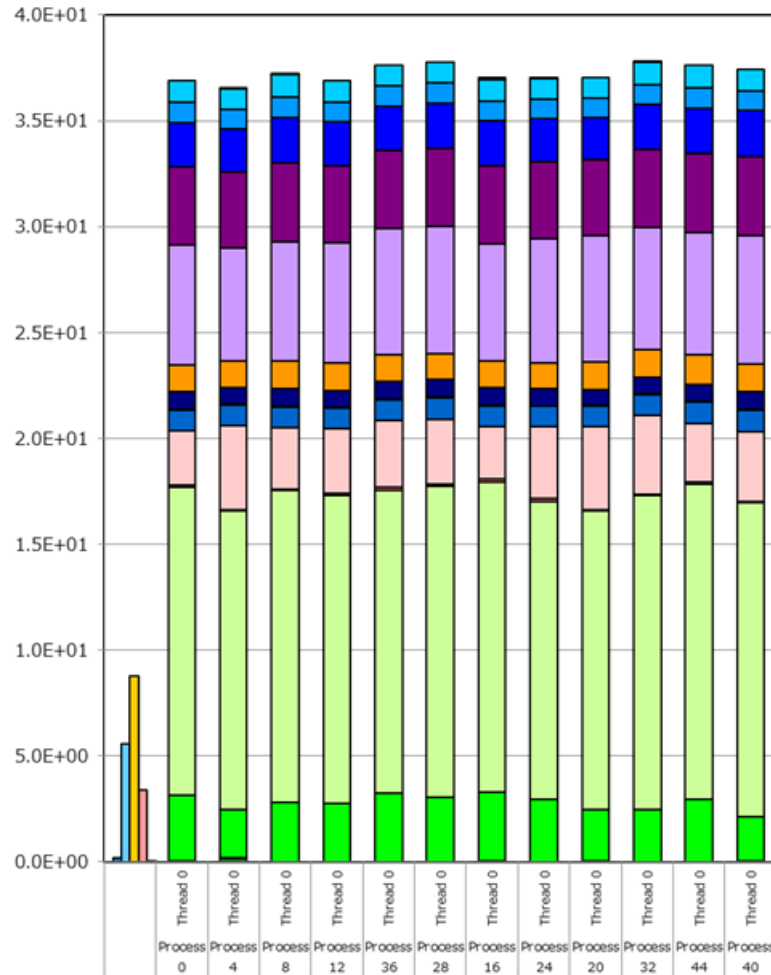
Significant fraction of runtime in the Integer Load L1D and Floating-point Load L1D cache access wait times

Left: pink section is Barrier synchronization wait

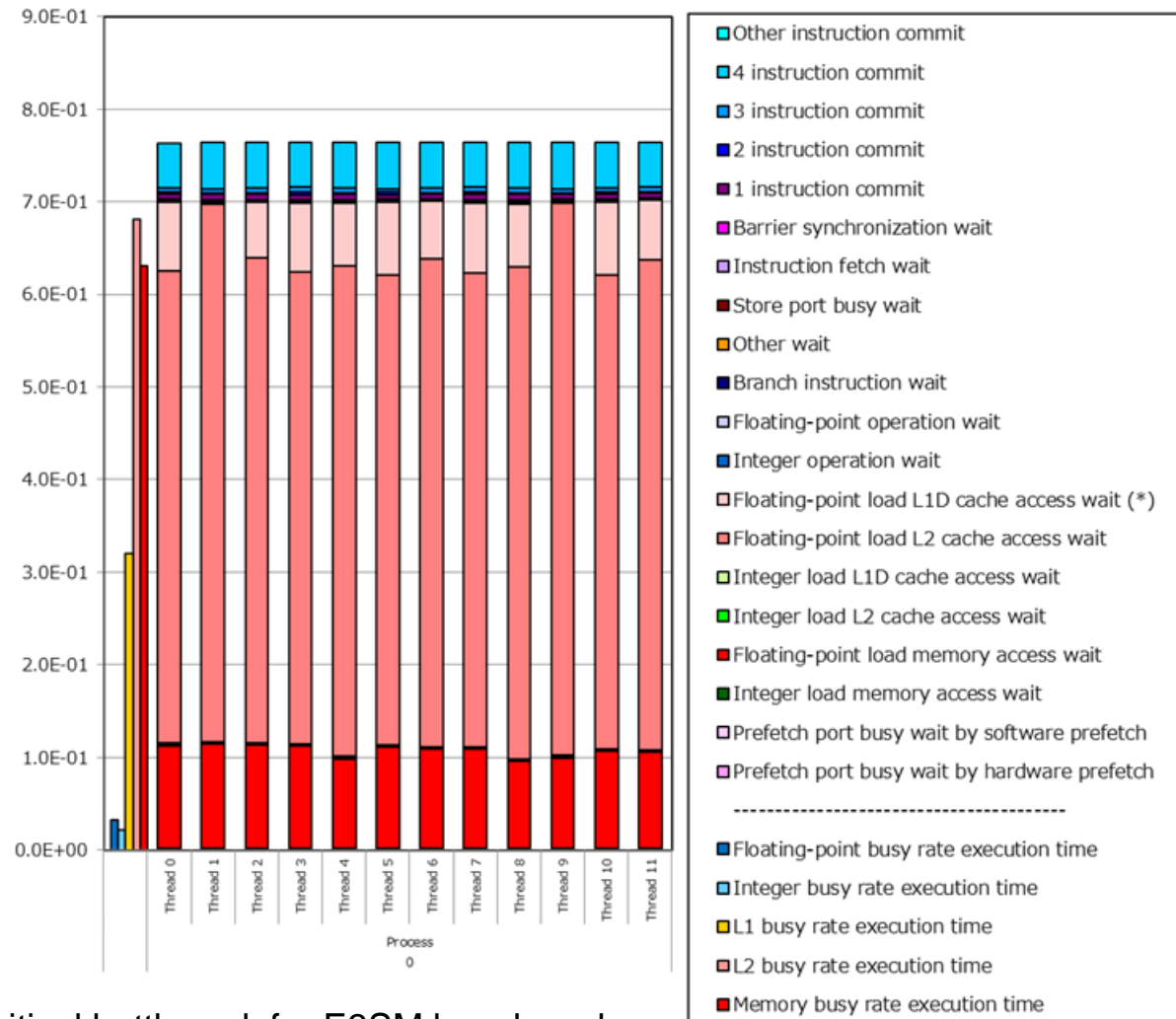
(*)Include wait time for integer L1D cache access

Instruction mix Comparison: Gradient sphere kernel vs. STREAM TRIAD

Gradient Sphere : Cycle Accounting execution time(s)



STREAM TRIAD: Cycle Accounting execution time(s)



Integer L1D cache access wait times critical bottleneck for E3SM benchmark
Mitigate high instruction latencies (INT: 5 cycles, FP: 8 cycles, SVE: 11 cycles)

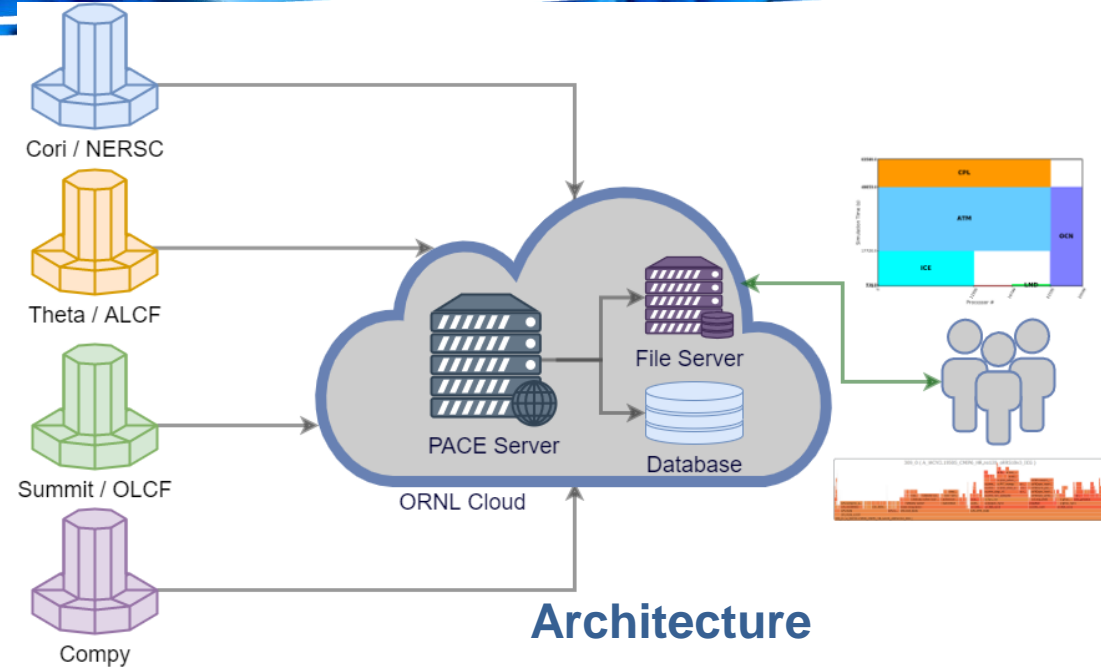
(*) Include wait time for integer L1D cache access

Performance Analytics for Computational Experiments



Summary

- Captures every E3SM experiment run on DOE supercomputers *automatically*
 - Performance Summary & Provenance
 - Facilitate performance research

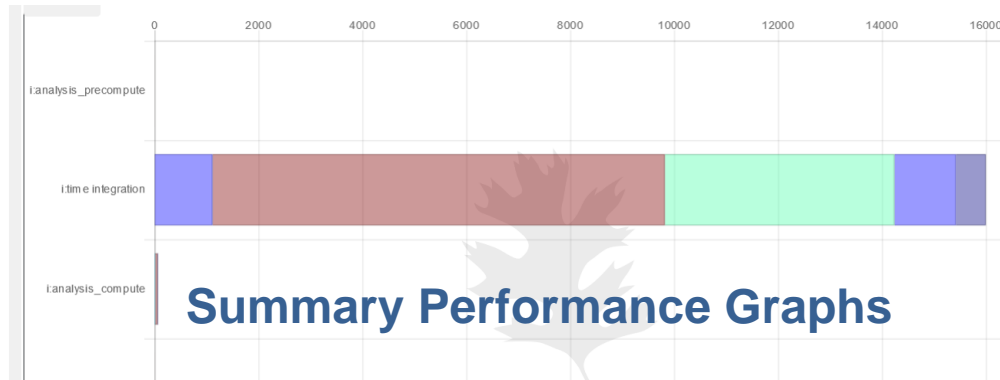
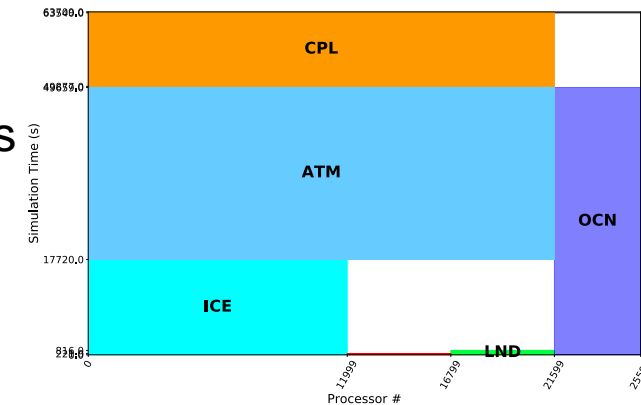


Architecture

Stats

- 130k experiments
- 3+ million input files
- 200+ users
- 14 platforms

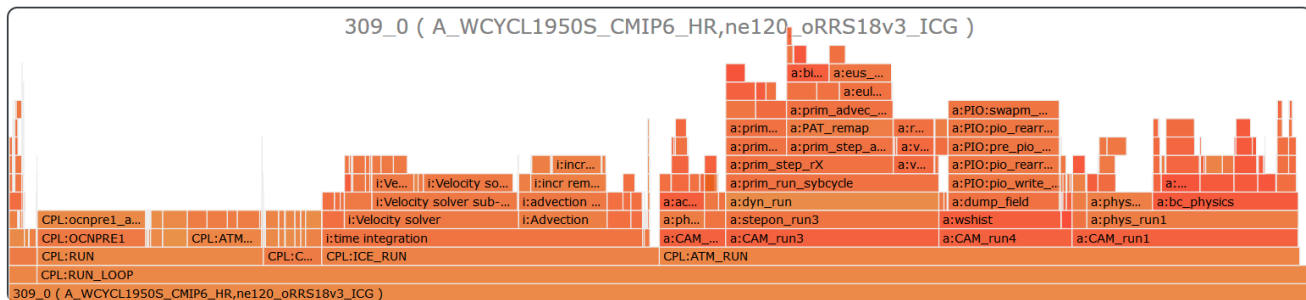
Load Balancing



Summary Performance Graphs

```

PIO:pio_get_var_id_double
PIO:pio_get_var_0d_text
PIO:pio_get_var_id_int
CPL:INIT
CPL:RUN_LOOP_BSTART
CPL:RUN_LOOP
CPL:CLOCK_ADVANCE
CPL:RUN
CPL:COMM
CPL:ICE_RUN
i:analysis_precompute
i:time integration
i:analysis_compute
i:analysis_restart
i:analysis_write
CPL:ATM_RUN
    
```

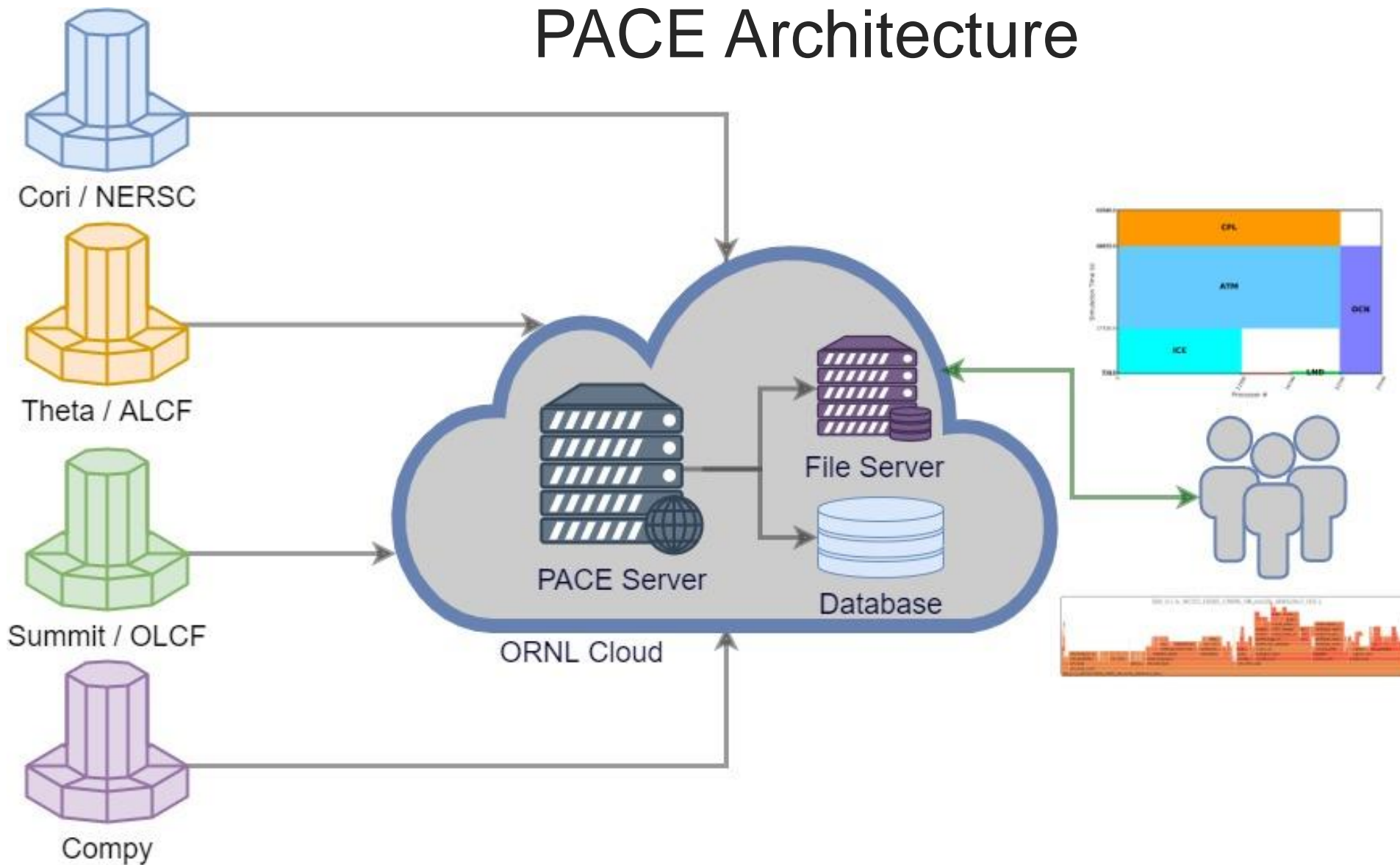


<https://pace.ornl.gov>

E3SM Performance Data

- Lightweight performance profiling by default
 - Utilizes General Purpose Timing Library (GPTL) timers
 - Mark start/stop at defined application phases
 - Aggregate statistics for parallel processes
 - Collect computation, communication and I/O performance data
 - Support for PAPI hardware counters
- Performance Archiving
 - Enabled on supported platforms at OLCF, ALCF, NERSC etc.
 - Archive performance data in project wide locations
 - Provenance data for context and reproducibility
 - System state and various logs

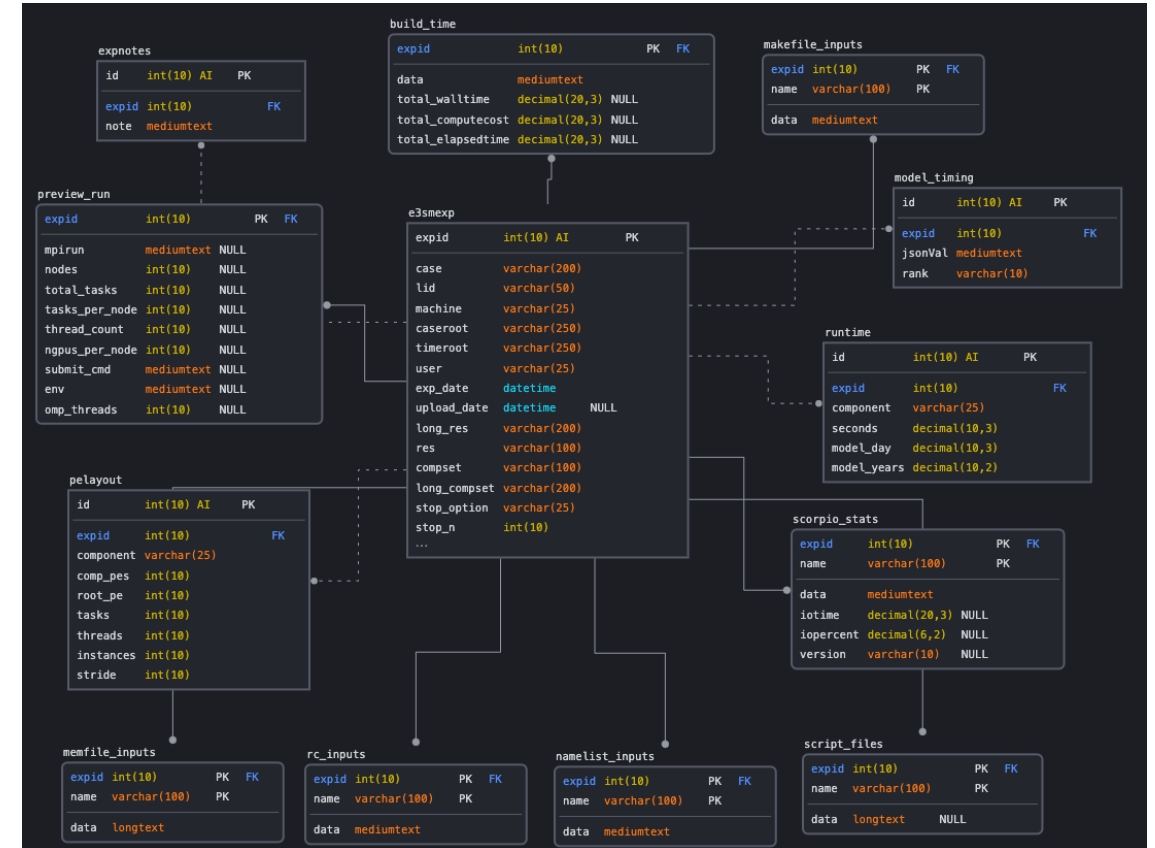
PACE Architecture



Technology Stack

- Infrastructure
 - ORNL Cloud (CADES)
 - OpenStack VM
- Nginx Web Server + Reverse Proxy
- Python-Flask middleware
 - Application Server
 - Process model inputs/timings
- Minio File Server
 - Object based storage for raw data
- MariaDB database
 - Structured and semi-structured data
 - Flexible Schema
- JavaScript
 - Frontend and visualization

Database Schema



Last but definitely not least:
Cybersecurity compliance at a DOE lab

Usage

- Search for existing experiment using case, compset, grid, user etc. (Autocomplete supported)
 - Sort by desired criterion
- Click on a row from search results to dive into specific experiment
- Experiment details page contains
 - Metadata: user, machine, date etc.
 - Provenance: Browse model inputs
 - Performance overview
 - Model, Component throughputs
 - Process layout diagram
 - Links to detailed performance graphs

Sort by ▾

Ascending Descending

ID	User	Machine	Compset	Res	Case	Total PEs	Run Length (days)	Throughput (sim_years/day)	Init time	ExpDate	Summary Charts
39596	mwu1	cori-knl	A_WCYCL20TRS_CMIP6	ne30_oECv3_ICG	20201019.DECKv1b_H3_...	32000	1825	3.15	351.632	2020-10-27 22:10:07	<input type="checkbox"/> Global Stats <input type="checkbox"/> Rank 0 <input type="button" value="More"/>
39438	xudo627	compy	ICLM45	CLMMOS_USRDAT	Amazon_Calibration_e...	400	12410	163.76	11.363	2020-10-27 21:16:40	<input type="checkbox"/> Global Stats <input type="checkbox"/> Rank 0
39586	ndk	cori-knl	F2010-SCREAM-HR-DYAM...	ne1024pg2_r0125_oRRS...	fne1024pg2tri.s32-o...	3145728	0	0.00	1256.788	2020-10-27 20:43:33	<input type="checkbox"/> Global Stats <input type="checkbox"/> Rank 0

Experiment Details

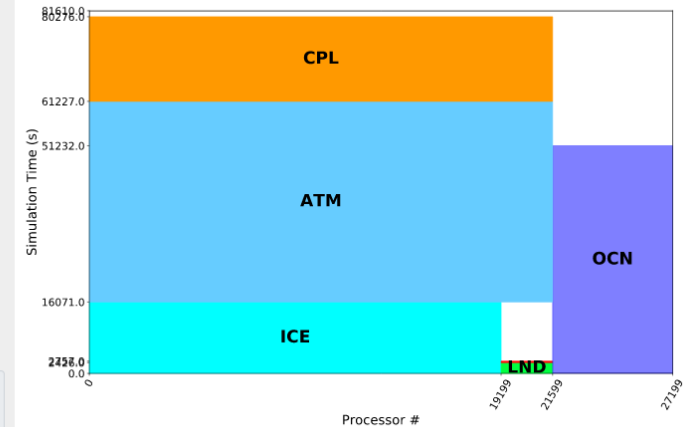
```

Id: 3262
User: azamatm
JobID: 312232.190213-030655
Machine: theta
Compset: A_WCYCL1950S_CMIP6_HR
Case: theta.20180906.branch_noCNT.A_WCYCL1950S_CMIP6_HR.ne120_oRRS18v3_ICG
Res: ne120_oRRS18v3_ICG
Version: v1.0-0-27-g2f3b0aec4
Date: 2019-02-14 02
Run_time: 81610.278 sec
Init time: 1454.002 sec
Final_time: 0.545 sec
Total PEs: 435200
Model cost: 930012.44 pe-hrs/sim_year
Run length: 242 days
Stop_n: 8
Stop Option: nmonths
MPI tasks/node: 32

Model Throughput: 0.70 SYPD

Graphs:
Summary Global statistics
Tree graph : Rank 0 Rank 19200 Rank 21600
Flame graph : Rank 0 Rank 19200 Rank 21600
Atm process distribution

Additional Notes
    
```



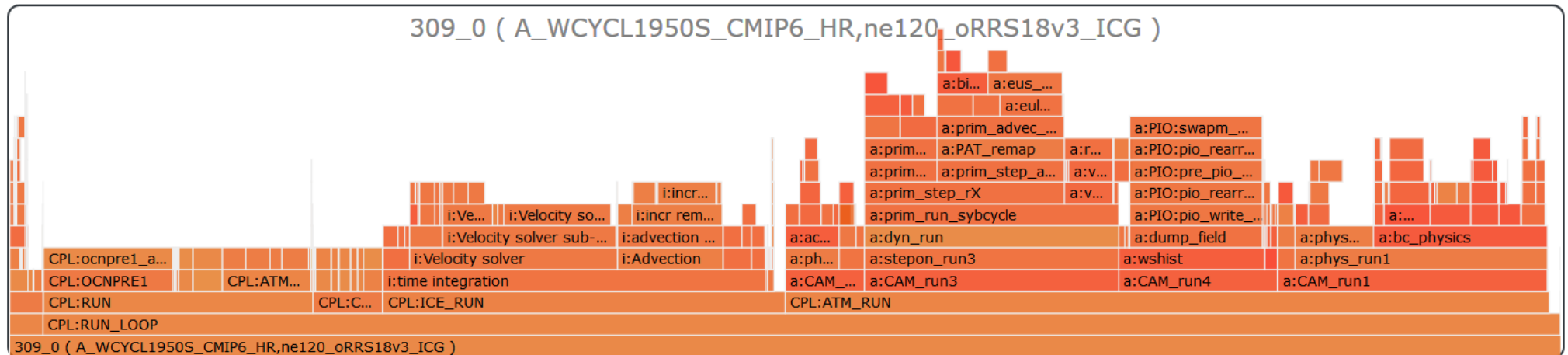
Tree Graph

Summarize time taken by model components
Recursively explore time taken by model sub-regions



Flame Graph

High-level overview of a parallel process execution time



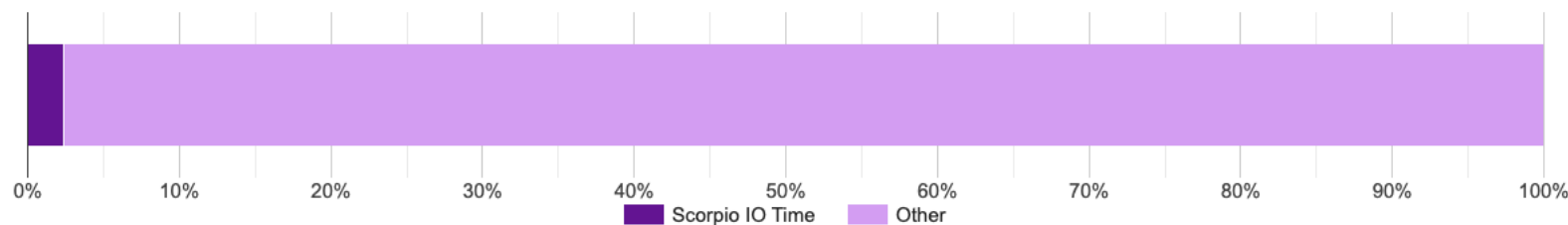
I/O Performance

Overall I/O Summary Statistics

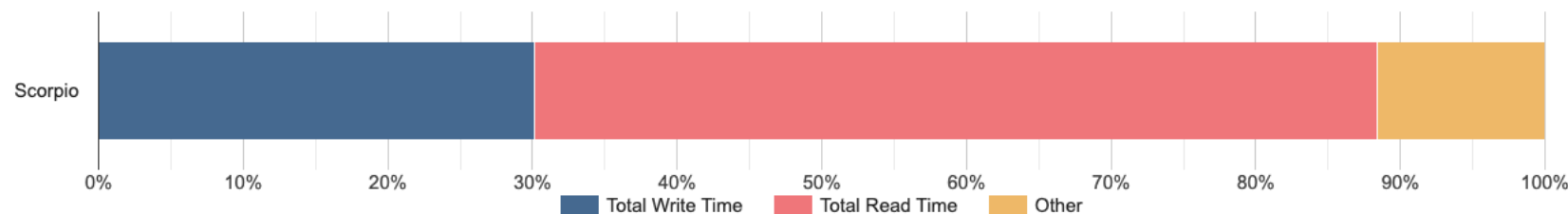
► [Collapse](#)

Total Run Time(s) ▲	Total IO Time(s) ▲	Total Write Time(s) ▲	Total Read Time(s) ▲	Avg Write Throughput(MB/s) ▲	Avg Read Throughput(MB/s) ▲	Total Write Size(GB) ▲	Total Read Size(GB) ▲	Na
140,791.72	3,285.97	991.81	1,911.73	1,673.52	246,476.79	1,740.44	494,084.65	Sc

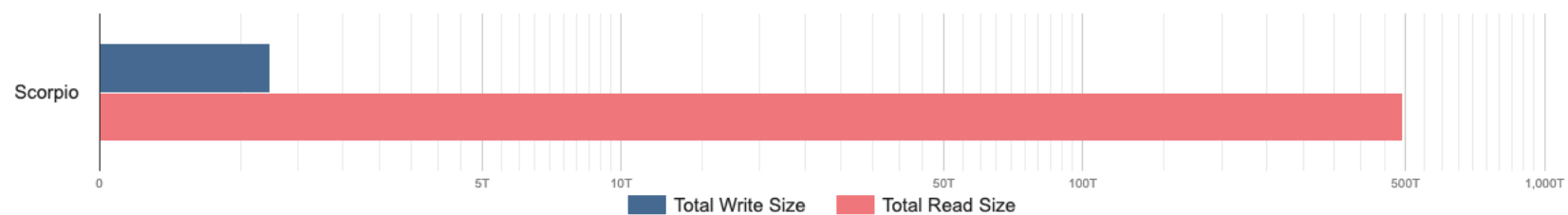
Scorpio IO time VS Model RunTime



Write VS Read Scorpio IO Time



Total Write Size VS Total Read Size



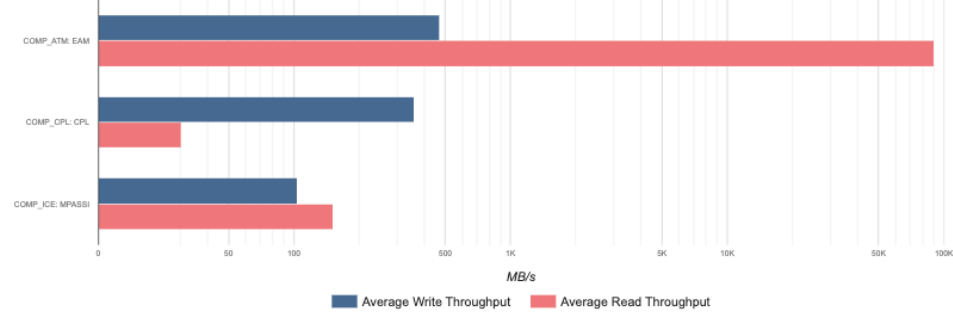
I/O Performance Details

Model Component I/O Statistics

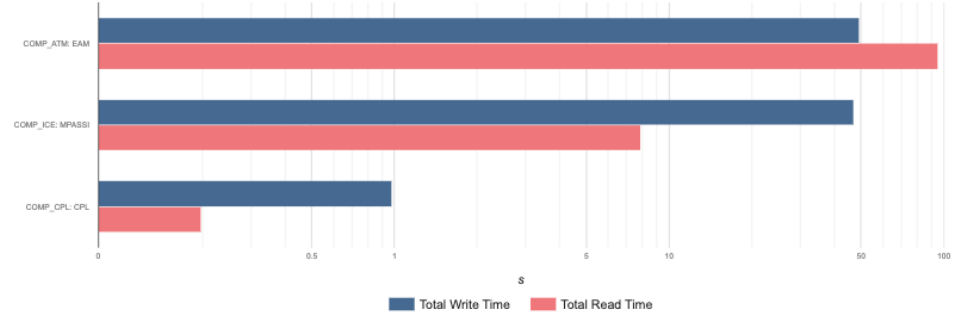
► Collapse

Total IO Time(s)	Total Write Time(s)	Total Read Time(s)	Avg Write Throughput(MB/s)	Avg Read Throughput(MB/s)	Total Write Size(GB)	Total Read Size(GB)	Name
147.65	49.34	95.06	470.36	90,044.44	24.33	8,975.11	COMP_ATM: EAM
56.87	46.79	7.85	103.50	152.12	5.08	1.25	COMP_ICE: MPASSI
1.29	0.98	0.20	359.36	30.09	0.37	0.01	COMP_CPL: CPL

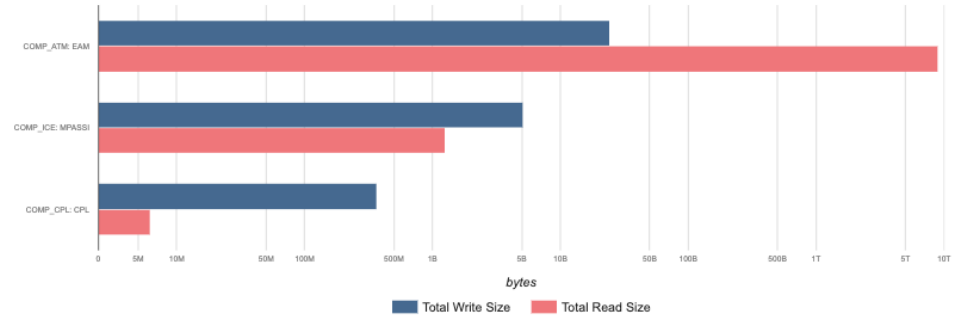
Avg Write Throughput VS Read Throughput



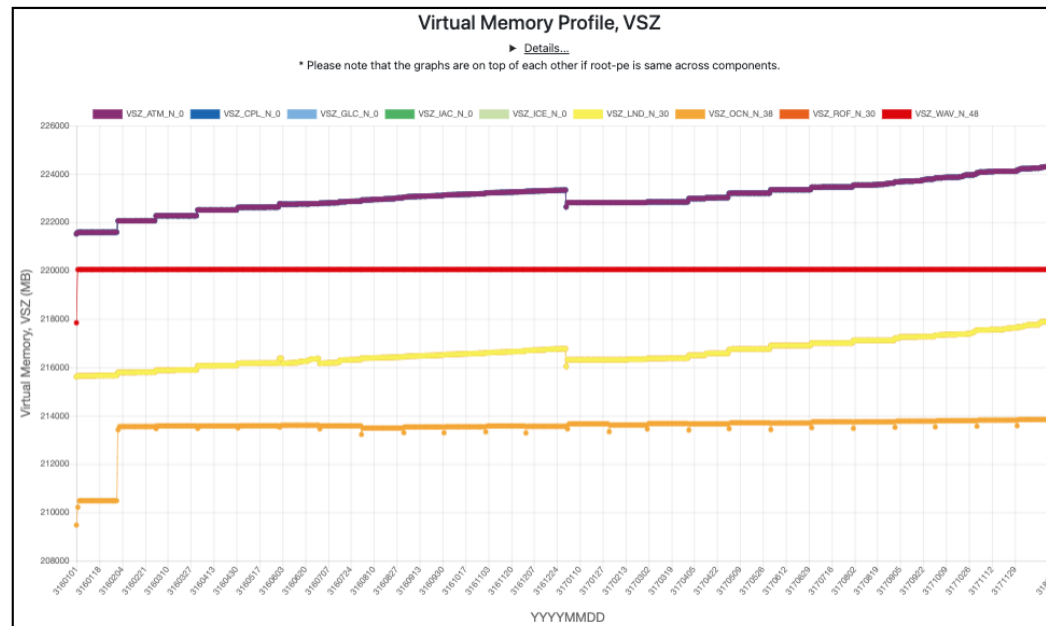
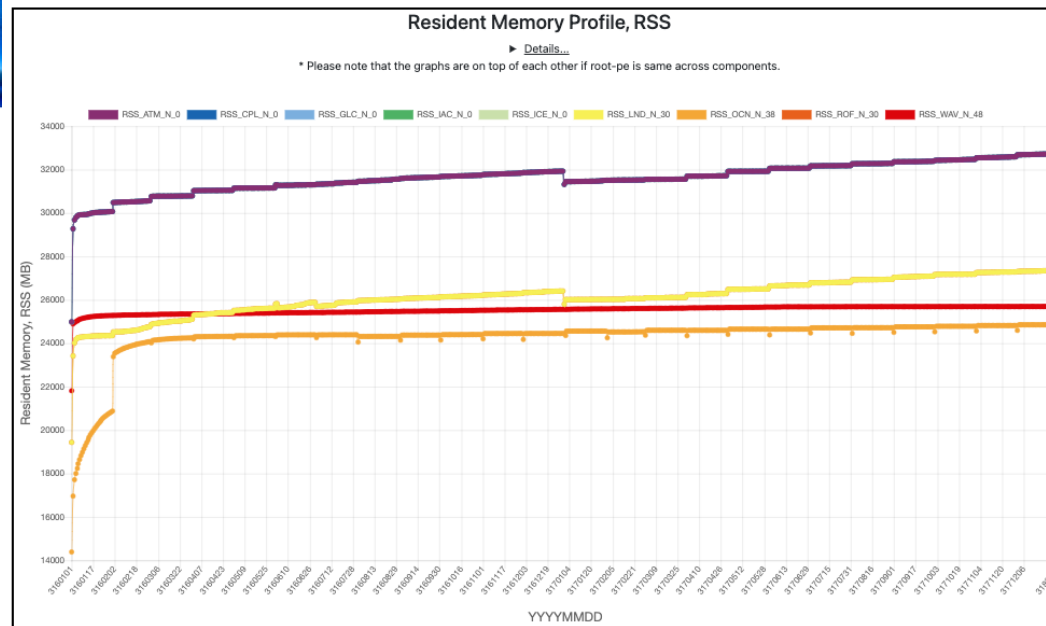
Total Write Time VS Read Time



Total Write Size VS Read Size



Memory Profiles



Build Profiles

Build Times

Wall Time for build: 2039.520 sec.

Total Compute Cost of build: 2423.561 sec.

Note: This is the total cost associated with compilation across components.

Typically, the wall time is lower due to parallel builds.

► [Details](#)

Time (s) ▼	File Name ▲
2,039.520000	Total_Build
139.131729	CMakeFiles/Ind.dir/___/elm/src/data_types/VegetationDataType.F90.o
63.014300	CMakeFiles/ocn.dir/___/core_ocean/driver/mpas_ocn_core_interface.f90.o
60.025588	CMakeFiles/ice.dir/___/core_seaice/model_forward/mpas_seaice_core_interface.f90.o
52.475535	CMakeFiles/atm.dir/___/eam/src/physics/cosp2/local/cosp.F90.o
51.702540	CMakeFiles/Ind.dir/___/elm/src/data_types/ColumnDataType.F90.o
51.515716	CMakeFiles/Ind.dir/___/elm/src/biogeochem/CNCarbonFluxType.F90.o
35.808562	CMakeFiles/atm.dir/___/eam/src/physics/clubb/mt95.f90.o
27.578397	CMakeFiles/Ind.dir/___/elm/src/external_models/fates/main/FatesHistoryInterfaceMod.F90.o
25.040277	CMakeFiles/rof.dir/___/mosart/src/riverroute/RtmMod.F90.o

Ongoing and Future

Assistant

- Simulation planning
- Process layouts
- Data analytics
- Anomaly detection
- Allocation reports



Steve The Minion – from Pixabay
<https://pixabay.com/photos/minions-banana-steve-the-minion-2552584/>

Wizard

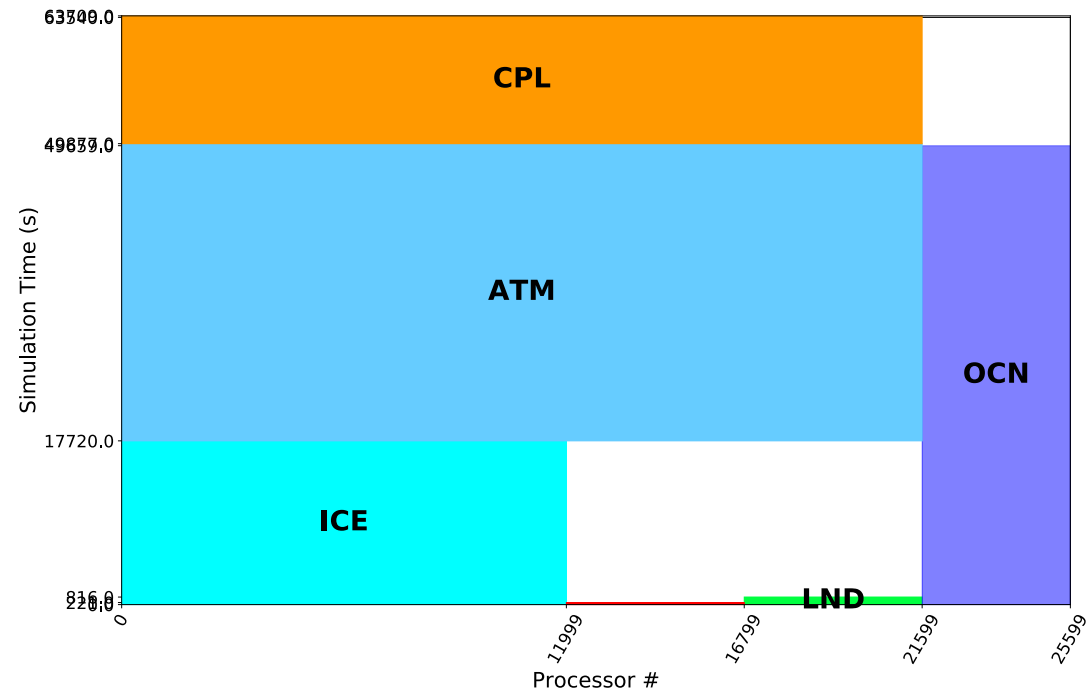
- Recommend optimizations
- Optimal resource allocations
- Machine Learning
- Communication optimization
- Active monitoring and reporting



Dennis Jarvis from Halifax, Canada / CC BY-SA
(<https://creativecommons.org/licenses/by-sa/2.0>)

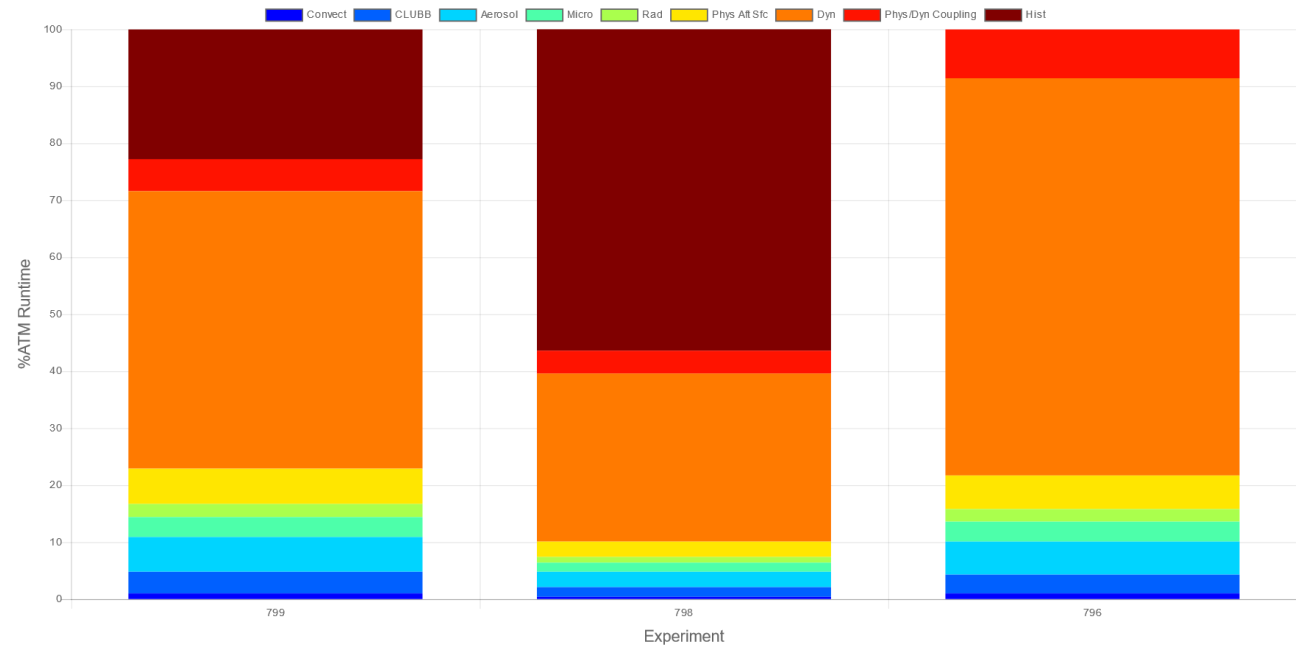
Performance Research Directions

Resource Allocation and Load Balancing



MPI Task Mapping

Targeted Optimization



Atmosphere model time distribution

EarthInsights: Parallel Clustering of Large Earth Science Datasets on the Summit Supercomputer

Sarat Sreepathi¹, Jitendra Kumar¹, Forrest M. Hoffman¹,
Richard T. Mills², Vamsi Sripathi³, William W. Hargrove⁴

¹Oak Ridge National Laboratory

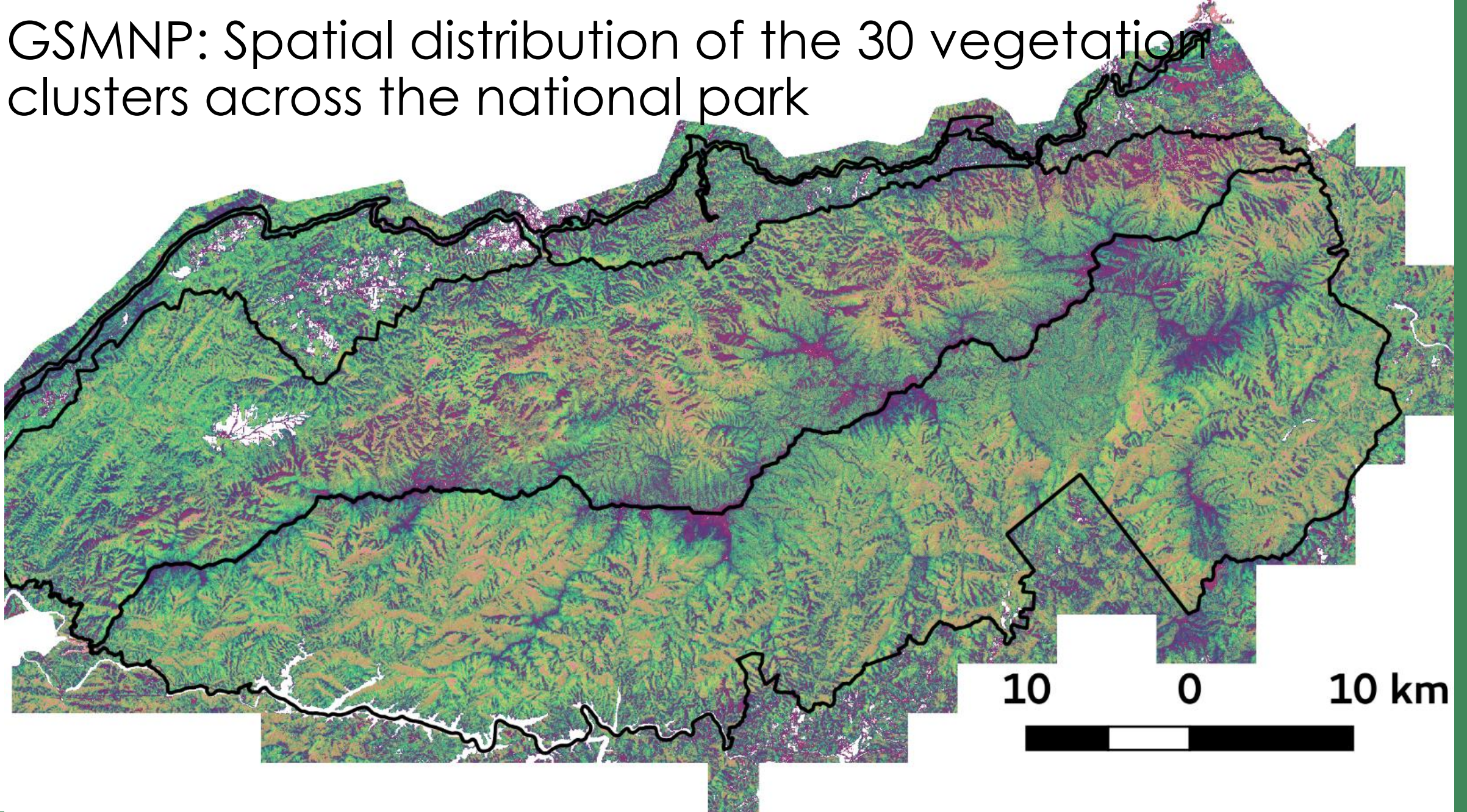
²Argonne National Laboratory

³Intel Corporation

⁴USDA Forest Service

ORNL is managed by UT-Battelle, LLC for the US Department of Energy

GSMNP: Spatial distribution of the 30 vegetation clusters across the national park

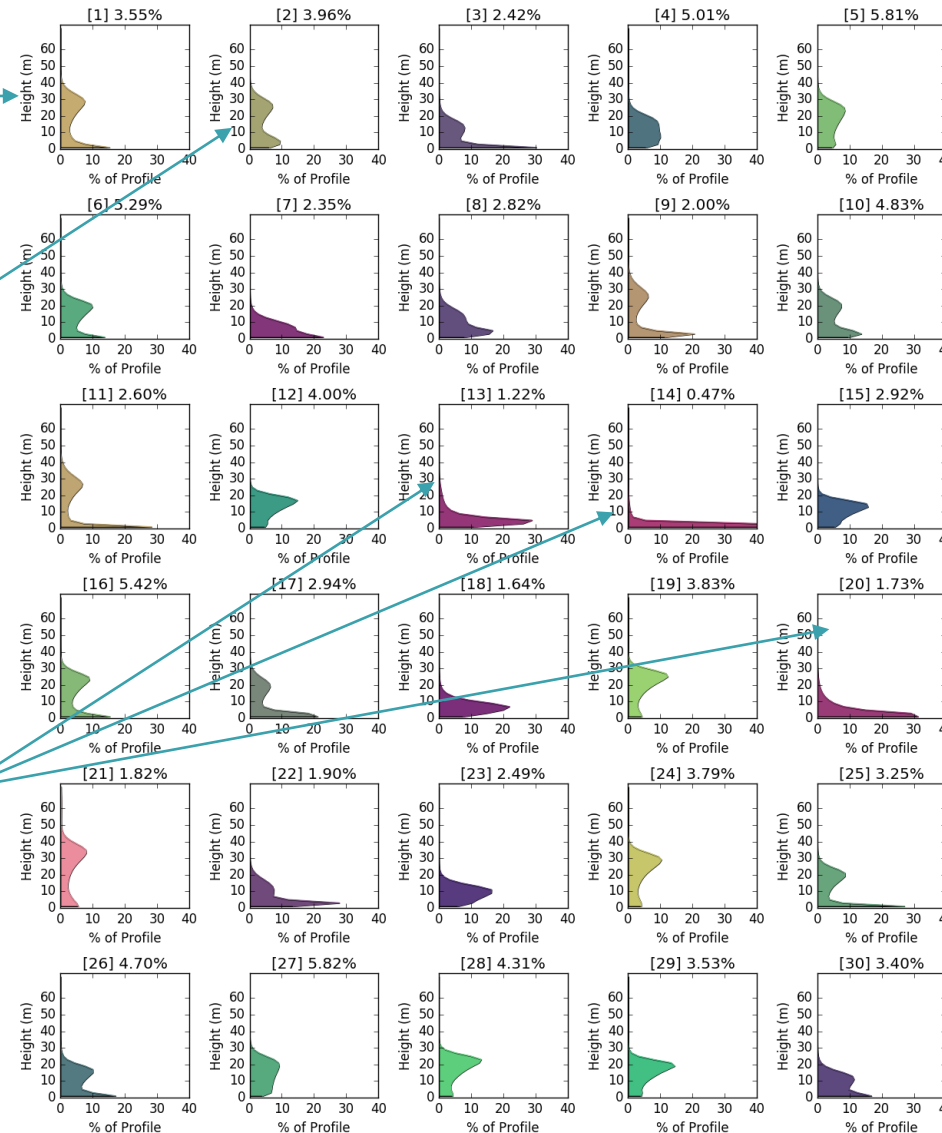


GSMNP: 30 representative vertical structures (cluster centroids) identified

tall forests with low understory vegetation

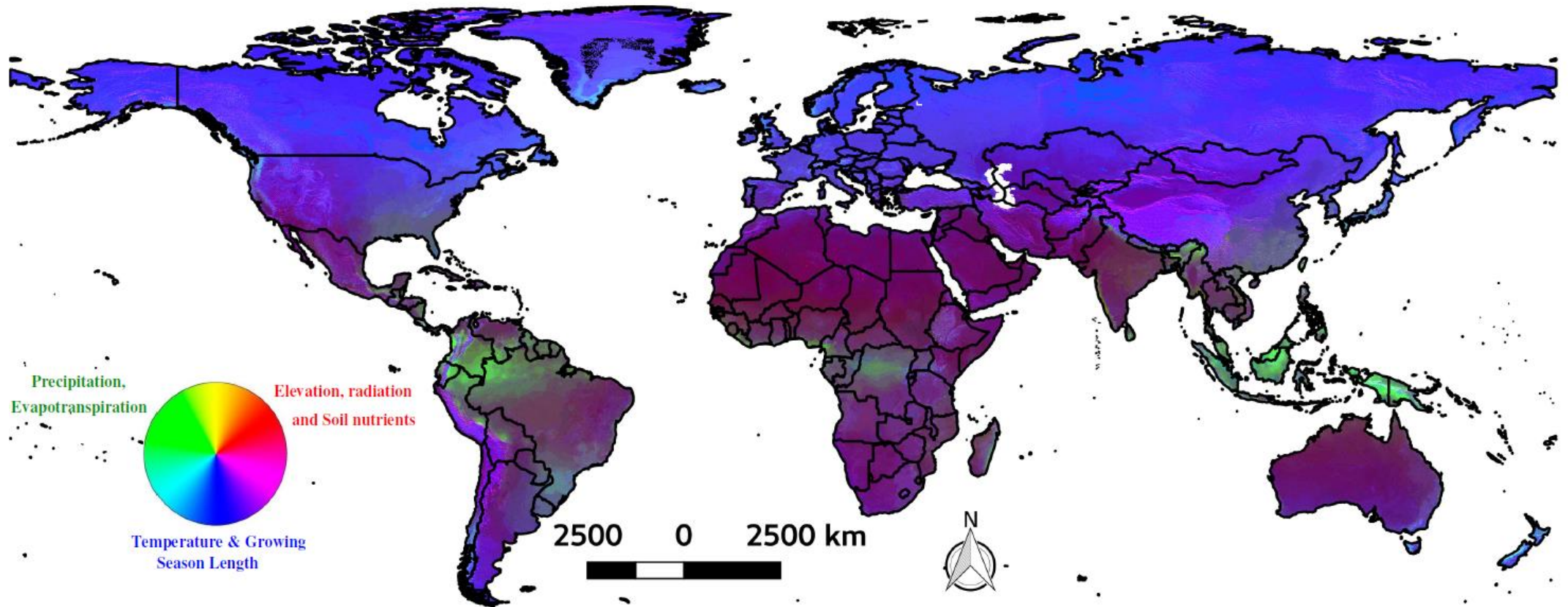
forests with slightly lower mean height with dense understory vegetation

low height grasslands and heath balds that are small in area but distinct landscape type

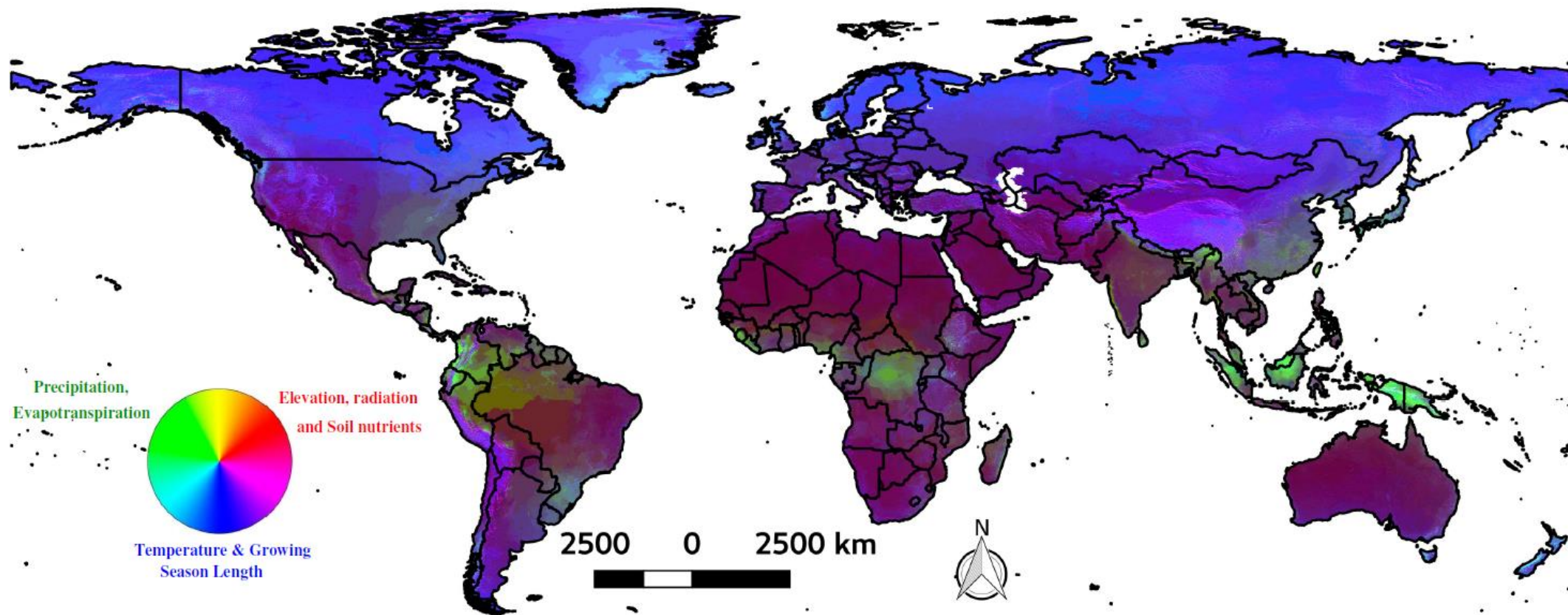


Global Climate Regimes: 1000 clusters

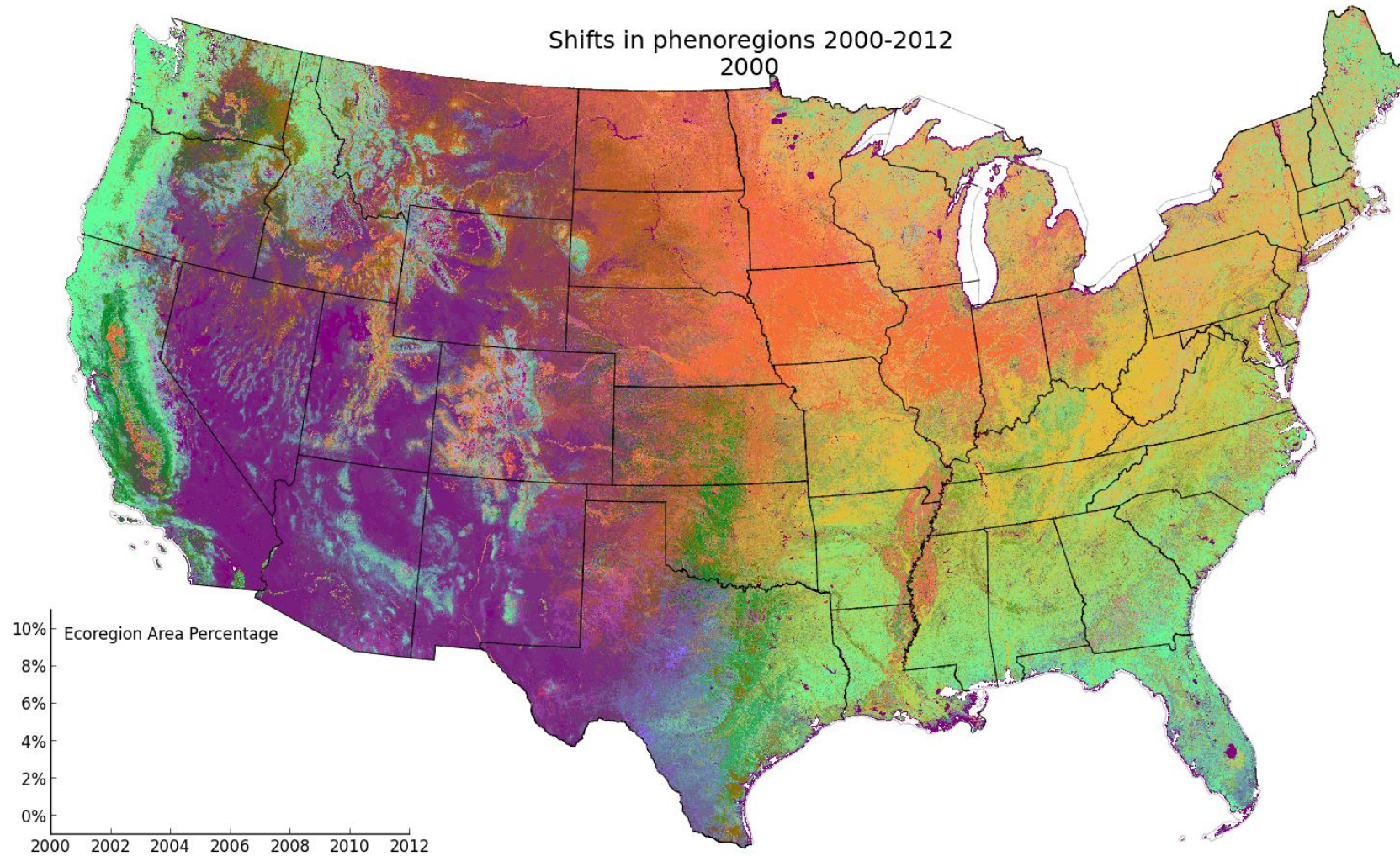
Contemporary using Similarity color scheme



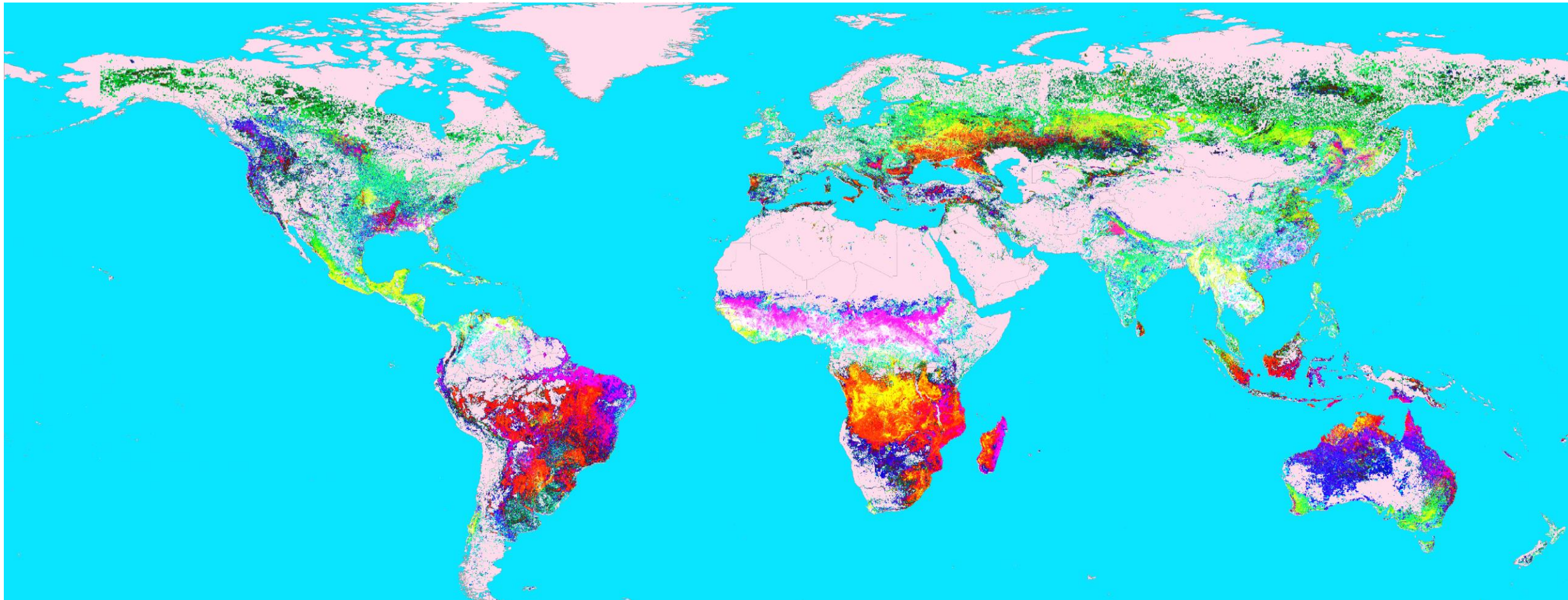
Global Climate Regimes: 1000 clusters 2100 using Similarity color scheme



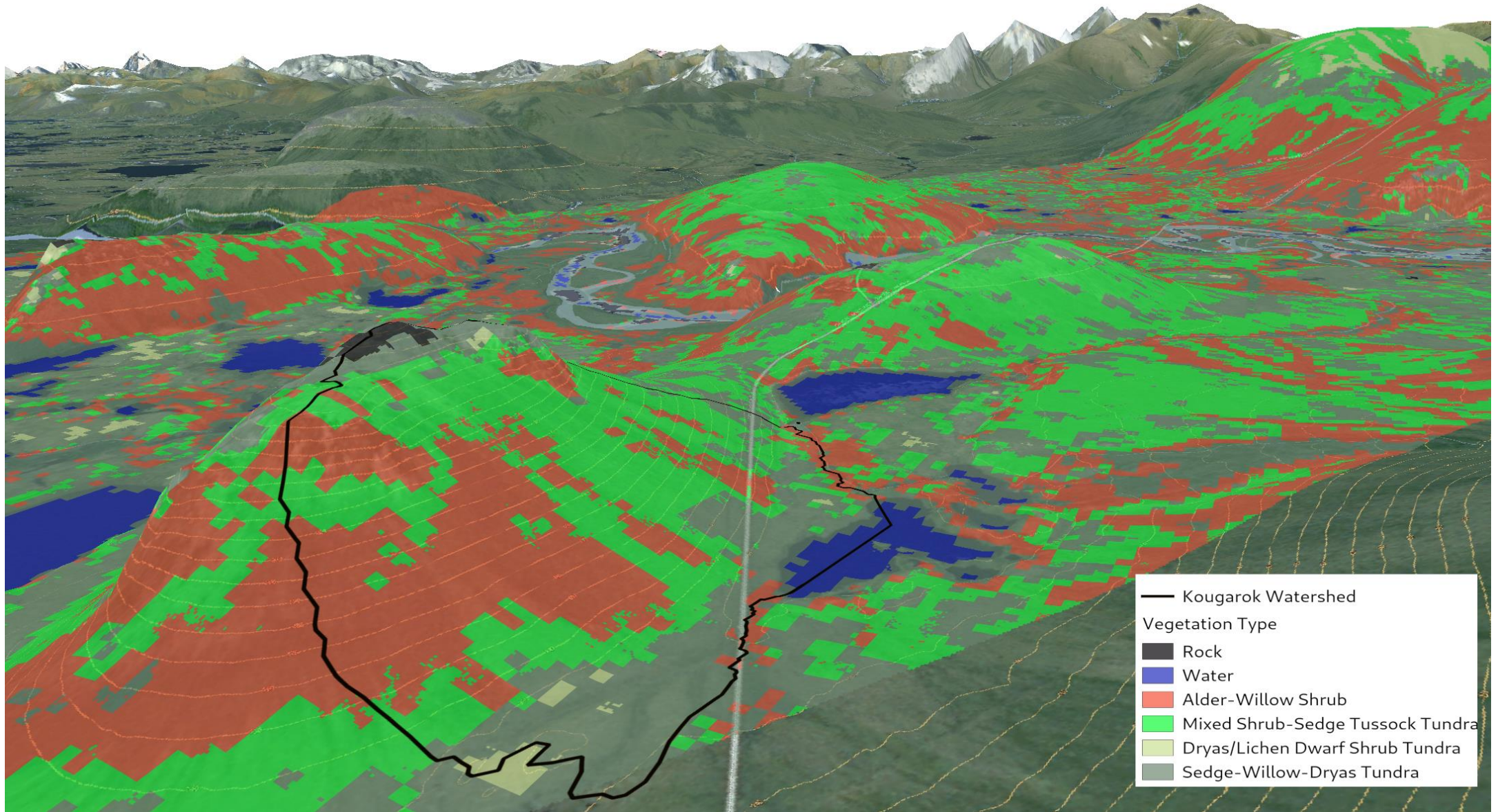
CONUS dynamic phenoregions



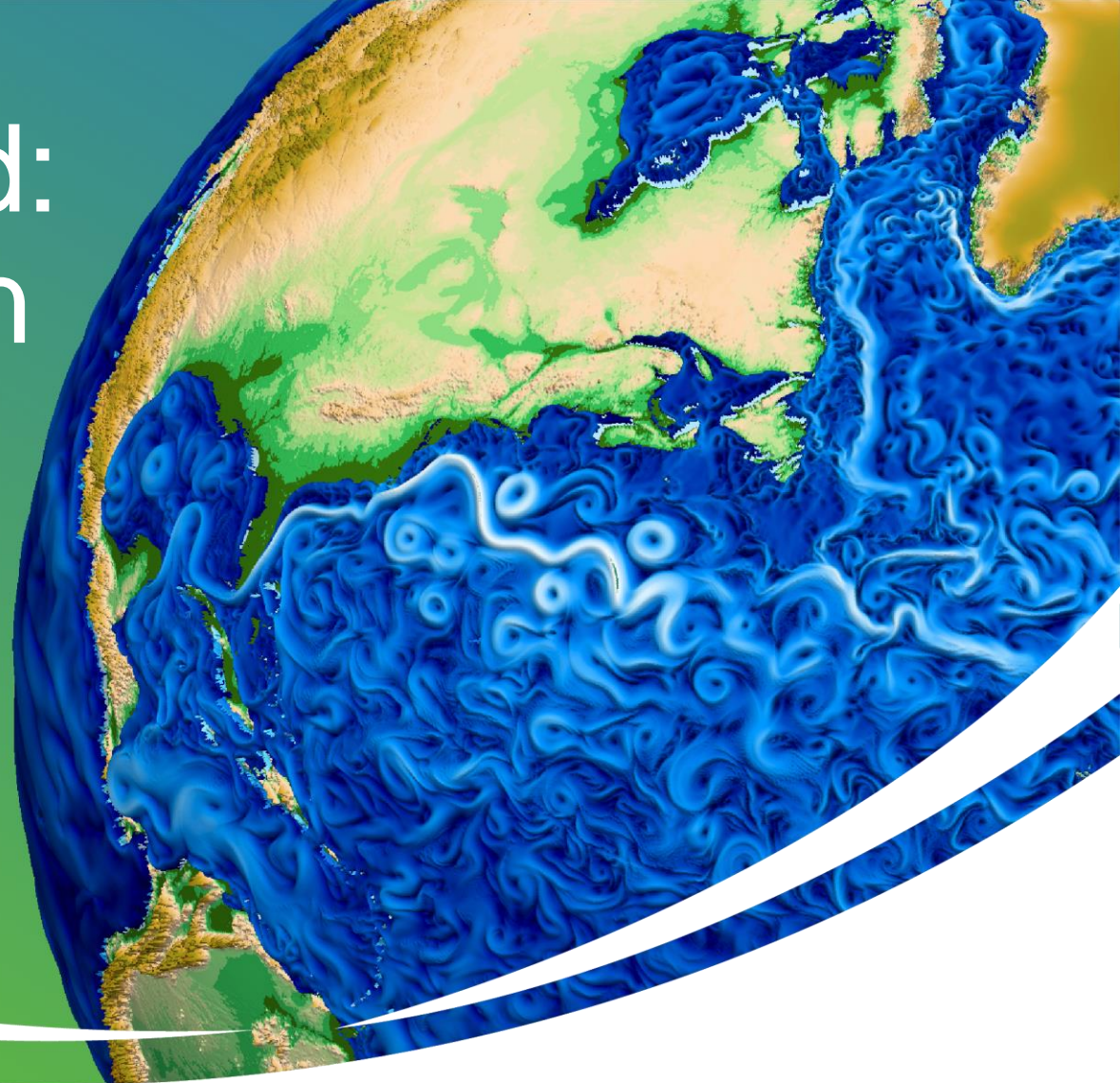
Global Fire Regimes



Arctic: High-resolution vegetation mapping

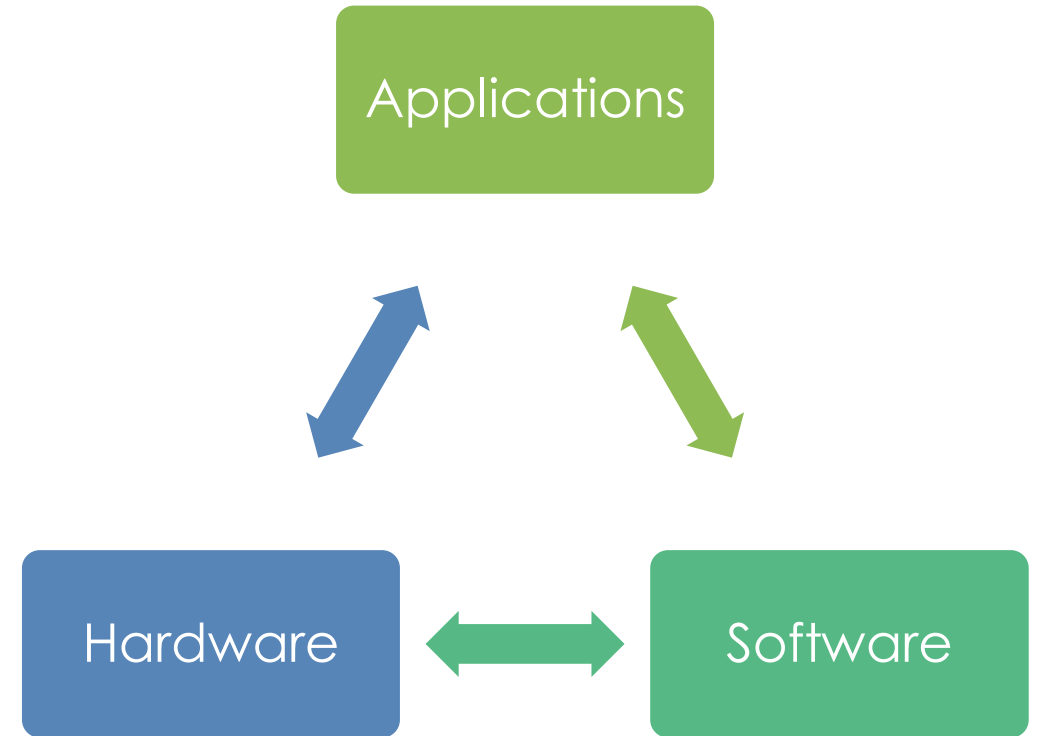


Exascale and Beyond: Application Co-design



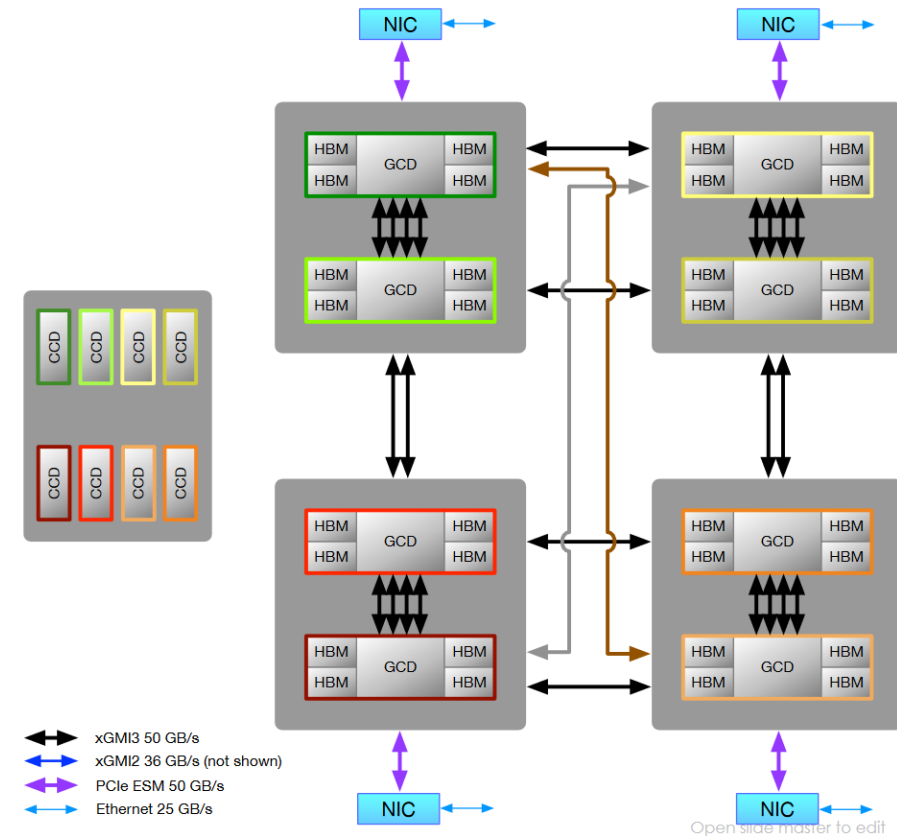
Co-design

- Feedback loop between applications, system s/w and computer architecture
- Application requirements inform (influence?) hardware design
- Technology choices and constraints guide problem formulation and design of algorithms.



Today's Exascale: Frontier

- 1.1 Exaflops (FP64 – HPL Benchmark)
- 29 MW
- 4,000 ft²
- 9,408 nodes
- Node
 - 4 AMD MI250X GPUs/node
 - Equiv. **8 logical GPU*s/node**
 - 512 GiB DDR4 (CPU) + **512 GiB HBM2e (GPU)**
 - **GPU Mem B/W: 8x 1,635 GB/s (13,080 GB/s Total)**
 - 1 AMD Trento CPU (64 cores)
- GPUs directly connected to high-speed interconnect
- **Aurora: Still under NDA**



Frontier Compute Node Architecture
1 CPU, 8 GPU*s

One cabinet of Frontier (24 ft²) has higher HPL than all of Titan (4,500 ft²) while using lower power (309 kW vs. 7 MW)



DOE Thinking

- [DOE RFI](#) – Summer 2022
 - Computing vendors and system integrators
 - Next generation supercomputers for 2025-2030 timeframe
- 10-20 FP64 exaflops in 2025 (8x from 2022)
- 100+ FP64 exaflops in 2030 (64x from 2022)
- 20-60 MW
- 4000 ft² (+ 50% more option)



Contract Opportunity
General Information

Follow

**Request for Information - Advanced
Computing Ecosystems (Dept. of
Energy)**

Optional

- Upgradability: Every 1-2 years
- Emerging accelerators (Quantum,...) if feasible
- Hybrid: On-prem + Cloud



New Golden Age for Computer Architecture

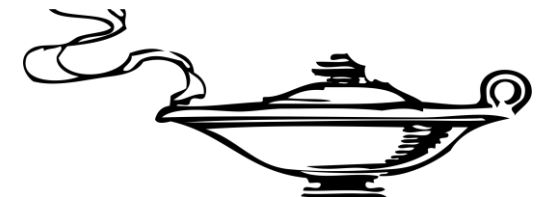
- Increasing heterogeneity
- Hybrid chips (APU/XPUs)
- Divergence of AI and HPC
- Open Source Hardware
 - [DARPA Electronic Resurgence Initiative \(ERI\)](#)
- Numerous [semiconductor startups](#)
- Chiplet-based System-on-Chips
- Widespread HBM
- 3D stacking
- Low-power ARM (A64FX, Grace,...)
- RISC-V
- Processing in memory
- Silicon Photonics, Optical Interconnects
- Moore's Law, Dennard Scaling
- Quantum: Optimization problems
 - Noisy Intermediate-Scale Quantum (NISQ): Practically useful?
- Neuromorphic: No clear fit
 - Spiking Neural Networks: Perhaps wavefront computations



Planning under uncertainty: A perspective

- Compute Architectures and Science: Friends or Frenemies?
 - Creativity for effective science
- High-end scientific computing: Leading vs. following
 - Cultivate and nurture vendor relationships
 - *Strategize ahead and influence vs. starting after general availability of an architecture*
- Co-design: Key application kernels and mini-apps
 - Impact on hardware: Skepticism warranted
 - Ray of hope (software/compiler)
- Changing Economics of Hardware Design
 - Fugaku (\$1B incl. R&D), Frontier \$600M procurement
- Wish: Imagine *relatively affordable* custom chips
 - Opinion: Better bet than fusion-powered quantum computers

(Domain-specific architectures)



Acknowledgments



EXASCALE COMPUTING PROJECT



U.S. Department of Agriculture, U.S. Forest Service,
Eastern Forest Environmental Threat Assessment Center.

©RIKEN
画像の無断使用・無断転載を禁じます

This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725.

Contact:
sarat@ornl.gov