# ADIOS: Storage and in situ I/O:
## Accelerating Scientific Knowledge Discovery with the Adaptable Input Output System

**Presenter**: S. Klasky[1,2,3],
**With contributions** from: N. Podhorszki[1] , A. Gainaru[1],M. Wolf[1], C. Atkins[8], William Godoy[1], Matthew Wolf[1] , Ruonan Wang[1], Chuck Atkins[8], Greg Eisenhauer[3], Junmin Gu[7], Philip Davis[5], W. Godoy[1], Jong Choi[1], Kai Germaschewski [10], Kevin Huck [11], Axel Huebl[7], Mark Kim[1], James Kress[1], Tahsin Kurc[1], Jeremy Logan[1], Kshitij Mehta[1], Franz Poeschel[8], Eric Suchyta[1], Keichi Takahashi, Lipeng Wan[1], Pradeep Subedi, Mark Ainsworth[19], Berk Geveci [8], Ben Whitney[1]

Science and Mathematics Division
[2] University of Tennessee, Knoxville, Department of Electrical Engineering and Computer Science
[3] Georgia Tech, School of Computer Science
[4] New Jersey Institute of Technology
[5] Rutgers University
[6] National Science Foundation, OAC
[7] Lawrence Berkeley National Laboratory
[8] Kitware
[9] Brown University
[10] University of New Hampshire
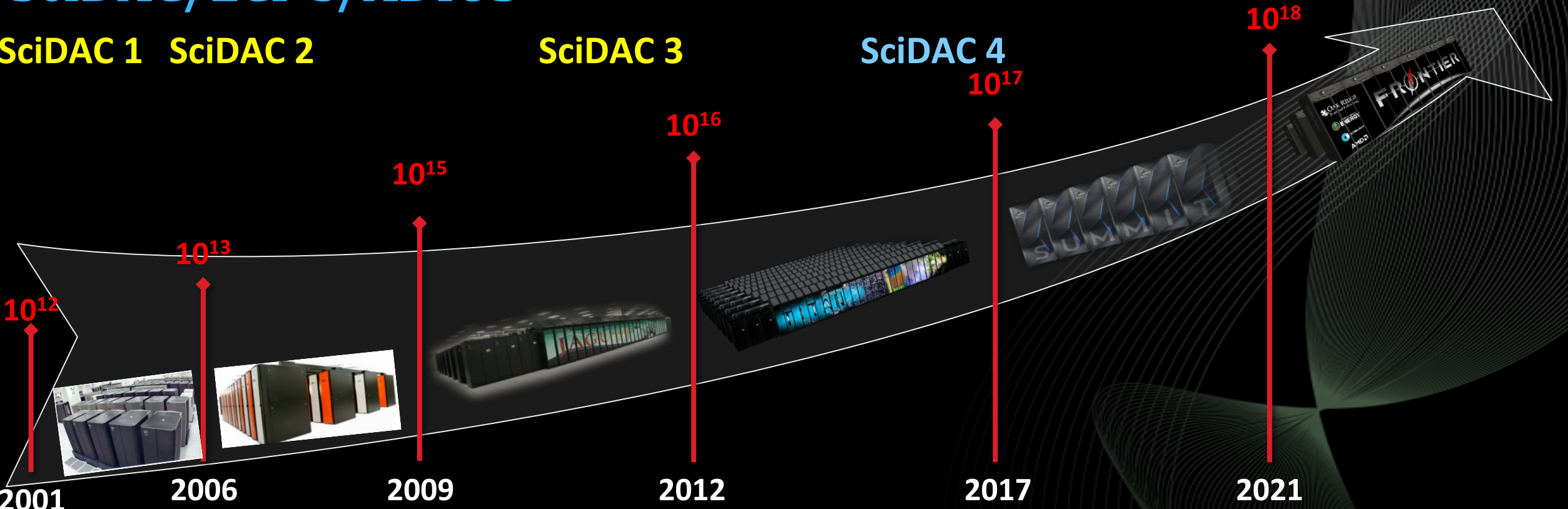[11] University of Oregon

[1] Oak Ridge National Laboratory, Computer

Community BOF, 3/30/2021

U.S. DEPARTMENT OF ENERGY | Office of Science

OAK RIDGE National Laboratory

# SciDAC/LCF's/ADIOS

**SciDAC 1**  **SciDAC 2**  **SciDAC 3**  **SciDAC 4**

$10^{18}$

$10^{17}$

$10^{16}$

$10^{15}$

$10^{13}$

$10^{12}$

**2001**    **2006**    **2009**    **2012**    **2017**    **2021**

**Seaborg:** 20 TF    **Jaguar:** 25 TF    **Jaguar 2:** 2.3 PF    **Titan:** 27 PF    **Summit:** 200 PF    **Frontier 1500 PF**

## ADIOS timeline

BP self-describing file format ADIOS 0.01

Research writing performance

Research reading performance

1.12 Burst buffer

Data Staging ADIOS 0.1

1.2 Data Staging

1.6 index, compression

1.10 Recovery/Ease of use

2.3 Summit optimizations

ADIOS 1.0 released

1.4 weak code coupling

1.8 Query /Indexing

1.11 compression

2.2 New framework in C++
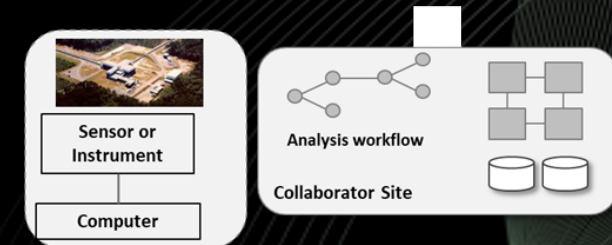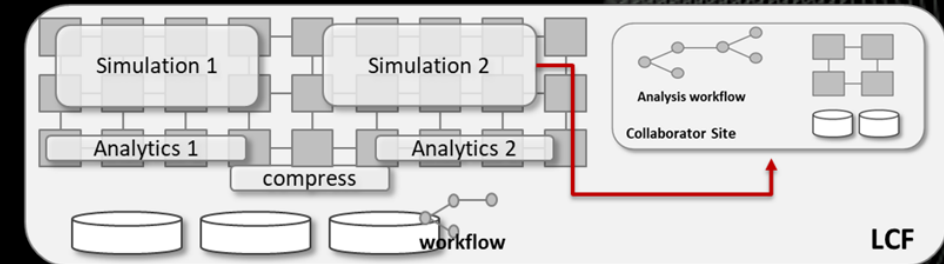
SciDAC funding in purple boxes

ECP funding

ASCR Research funding

OLCF funding

# ADIOS: High-Performance Publisher/Subscriber I/O framework

- An abstraction to allow for high-performance I/O to/from storage and for in situ processing

- Utilizes a publish/subscribe mechanism with self-describing data

- Optimized I/O engines for C/R, strong/loose in situ coupling WAN data streaming, and in-memory object storage

  - Fast Writing/Reading: BP4

  - DAOS optimizations: BP-DAOS

  - Write from GPU: BP-GPUDirect

  - Compatibility with HDF5: HDF5

  - Weak Code Coupling: SST

  - Tight Code Coupling: SSC

  - WAN streaming: DataMan

  - In memory object store: DataSpaces

  - Works with ECP reduction libs: MGARD, SZ, ZFP

- Typical for applications to achieve > 1 TB/s on Summit



3

# Sustainability: is a primary goal of the ADIOS project


Kitware

- ## Nightly testing
  - Testing on many different platforms

- ## Continuous Integration
  - Only allow tested code to be merged
  - Almost 2,000 tests for each commit

- ## Static and dynamic analysis reports
  - Compile-time and run-time analysis

- ## Code coverage
  - Level of testing

- ## External testing
  - Allow feedback from user projects

Nightly testing on target HPC platforms



coverage report

# ADIOS Approach: "How"

- I/O calls are of declarative nature in ADIOS
    - which process writes what: add a local array into a global space (virtually)
    - adios_close() indicates that the user is done declaring all pieces that go into the particular dataset in that timestep
- I/O strategy is separated from the user code
    - aggregation, number of sub-files, target file-system hacks, and final file format not expressed at the code level
- This allows users to choose the best method available on a system without modifying the source code
- This allows developers
    - to create a new method that's immediately available to applications
    - to push data to other applications, remote systems or cloud storage instead of a local filesystem

W. Godoy, N. Podhorszki, R. Wang, et al, ADIOS 2: The Adaptable Input Output System. A framework for High-Performance Data Management, SoftwareX, 2020.

# Creating I/O abstractions to accelerate I/O to storage

- One change in the code or input file, to specify the engine

```
adios2::Engine writer = io.Open("analysis.bp",
adios2::Mode::Write);

writer.BeginStep()

writer.Put(varT, T.data());

writer.EndStep()

writer.Close()


adios2::Engine reader = io.Open("analysis.bp",
adios2::Mode::Read);

reader.BeginStep()

adios2::Variable<double> T =
reader.InquireVariable("Temperature");

std:vector<double> t;

reader.Get(varT, t);

reader.EndStep()

reader.Close()
```
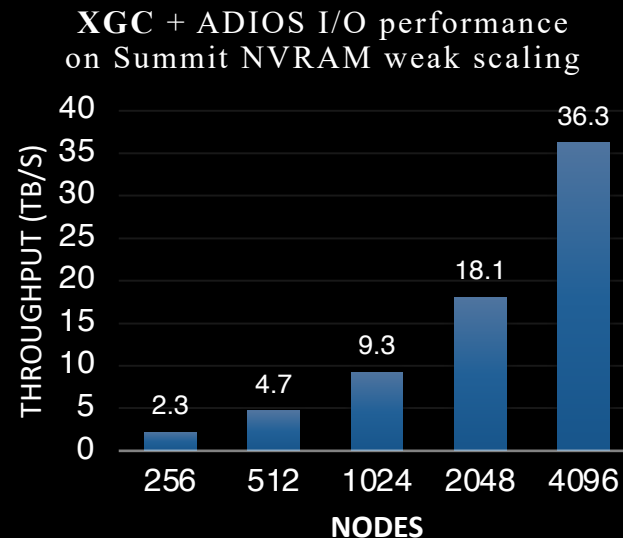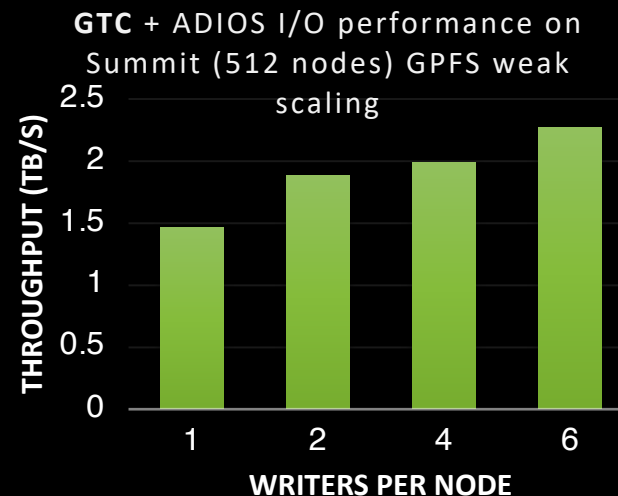
The APIs are identical for code coupling

**36.3 TB/s to NVRAM**

**XGC** + ADIOS I/O performance on Summit NVRAM weak scaling



Chart: THROUGHPUT (TB/S) vs NODES
- 256: 2.3
- 512: 4.7
- 1024: 9.3
- 2048: 18.1
- 4096: 36.3

**2.3 TB/s to GPFS**

**GTC** + ADIOS I/O performance on Summit (512 nodes) GPFS weak scaling



Chart: THROUGHPUT (TB/S) vs WRITERS PER NODE
- 1: ~1.45
- 2: ~1.9
- 4: ~2.0
- 6: ~2.3

# ADIOS performance results (measured by the app teams/not us)

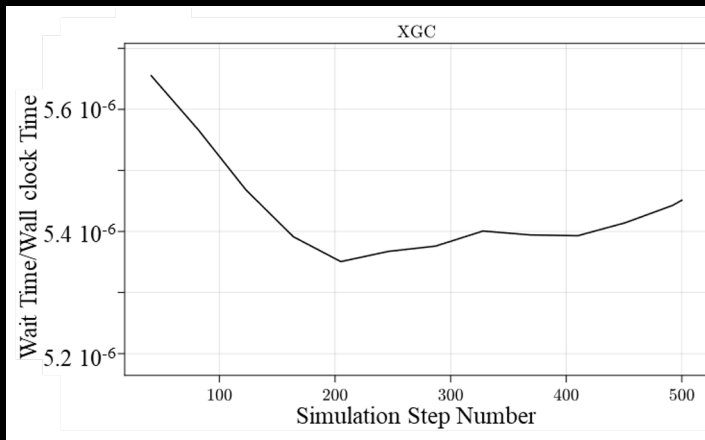## WDMApp

https://github.com/PrincetonUniversity/XGC-Devel

From the WDMApp annual ECP review

> XGC on 512 Summit nodes
> GENE on 6 Summit nodes
> $N_m$ = 9,640,480 vertices
> $N_p$ = 8,922 particles/vertex
> Timestep = 61.4 seconds.

XGC wait time during charge coupling



**Contact: Amitava Bhattacharjee (PPPL)**

## E3SM-MMF

https://github.com/E3SM-Project/scorpio/tree/master

ADIOS 2.x Port is integrated into master SCORPIO

SCREAM project evaluated it on their own and found 4-5x improvement in IO using ADIOS for TBs of data

New I-Case stresses IO
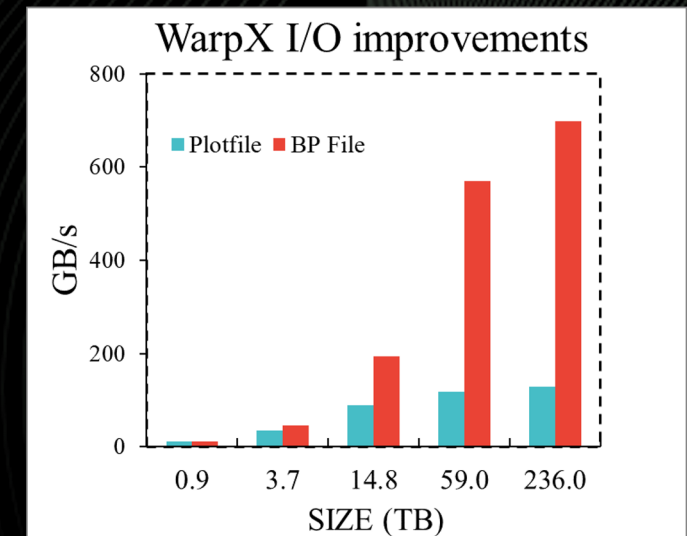
### I-Case benchmark on Summit



Simulating 1 day, 5 days and 10 days
Writing data every simulated hour
Run on Summit, 1344 MPI processes

**Contact: Mark Taylor (SNL), I-Case Peter Thornton (ORNL), SCREAM Peter Caldwell (LLNL)**

## WarpX

BP4 improved append performance for ADIOS and now applications can see the benefits of that

WarpX and in general, OpenPMD users can get high throughput



WarpX on Summit, weak scaling
6 GPUs, up to 256 nodes
ADIOS vs original AMReX Plot files

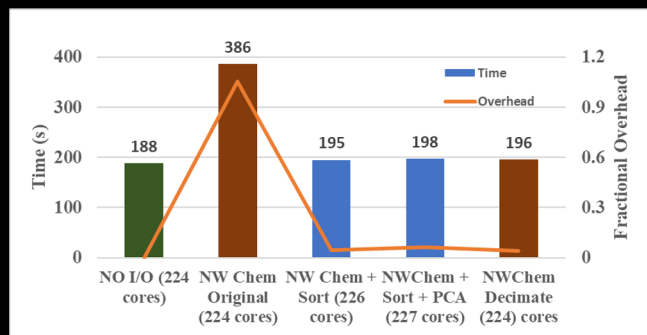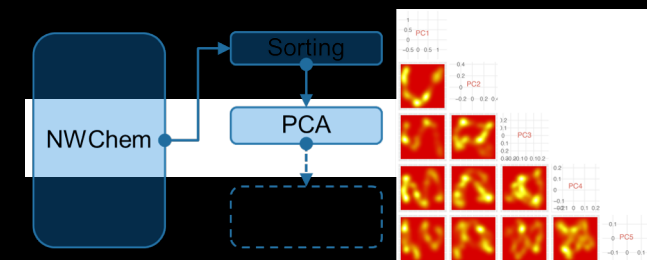**Contact: Jan-Luc Vay, Axel Huebl (LBNL)**

# ADIOS performance results (measured by the app teams/not us)

## NWChem

In situ sorting of atom trajectories can save 50% of runtime

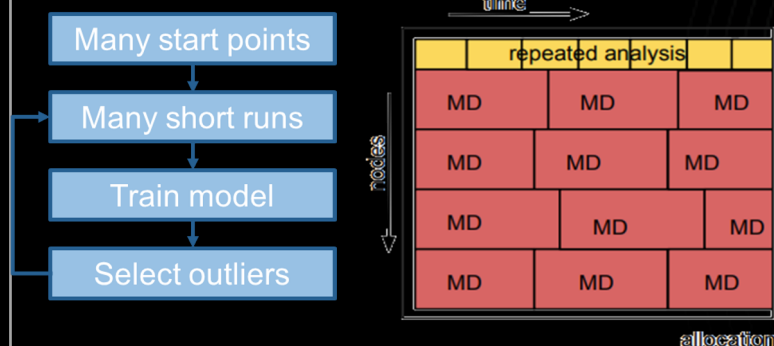Motion correction with PCA analysis (pbdR script) in situ



**Contact: Tjerk Straatsma (ORNL)**

## CANDLE/DeepDriveMD

CODAR collaboration for accelerating sampling of macromolecule potential energy surface via online coupling

Many concurrent MD runs + online training + inference (outlier search)

ADIOS for async collection of MD results to training allows for continuous simulation running and training
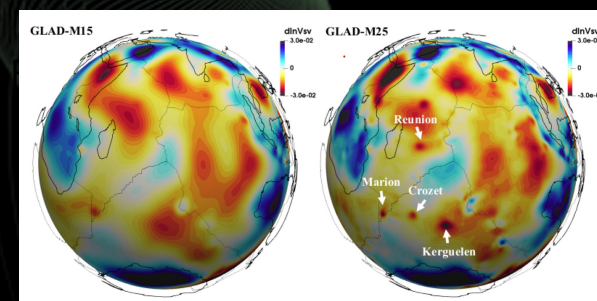


**Contact: Arvind Ramanathan, Igor Yakushin (ANL)**

## Specfem3D_globe

The Adaptable Seismic Data Format (ASDF) was developed that leverages the Adaptable I/O System (ADIOS) parallel library.

It allows for recording, reproducing, and analyzing data on large-scale supercomputers

1.5 PB of data is produced in every workflow step, which is fully processed later in adjoint simulation
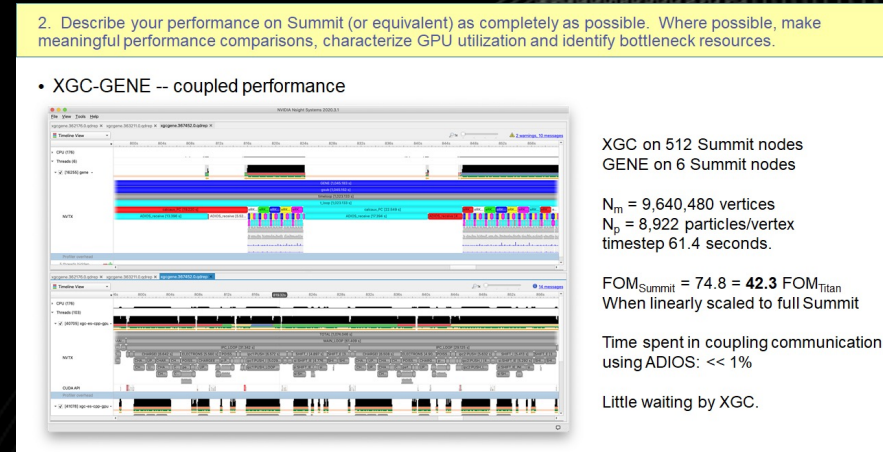https://www.olcf.ornl.gov/2019/07/05/tromp-titan/



Global adjoint tomography—model GLAD-M25, Geophysical Journal International, Volume 223, Issue 1, October 2020, Pages 1–21, https://doi.org/10.1093/gji/ggaa253
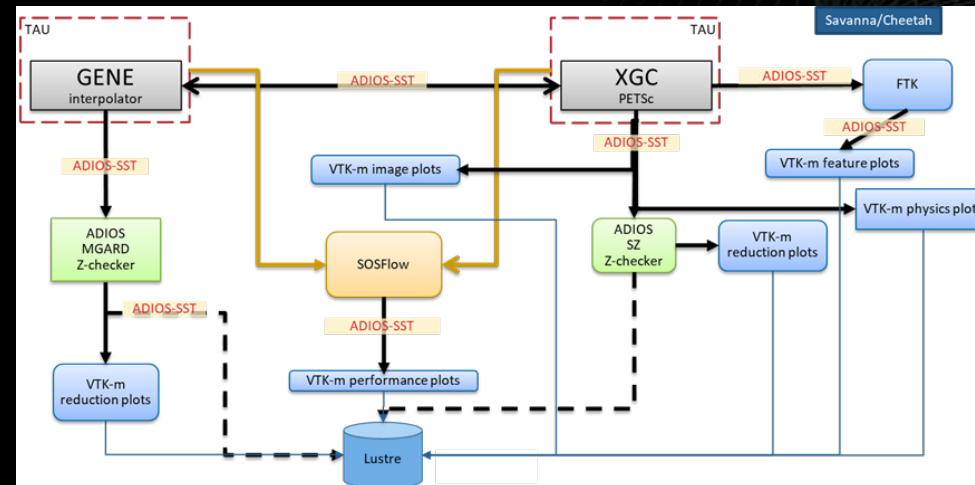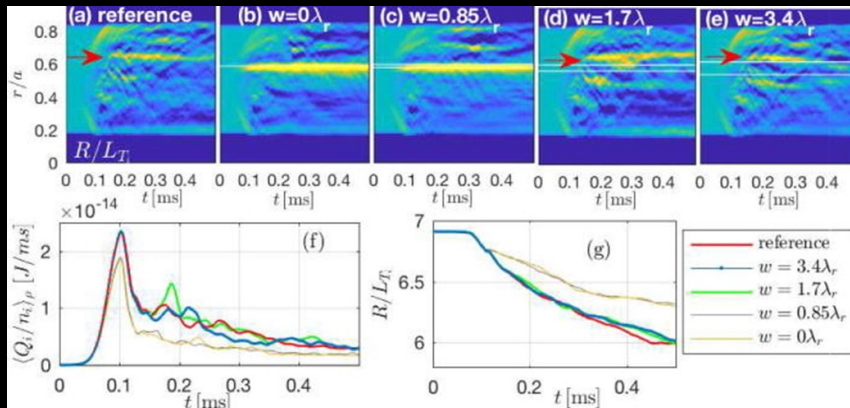
**Contact: Jeroen Tromp, Princeton University**

# 2.2.2.05 ADSE12-WDMApp: High-Fidelity Whole Device Modeling of Magnetically Confined Fusion Plasmas

PI: Amitava Bhattacharjee, PPPL,
C. S. Chang, PPPL

- Different physics solved in different physical regions of detector (spatial coupling)

- Core simulation: **GENE**
  Edge simulation: **XGC**
  Separate teams, **separate codes**

- Recently demonstrated first-ever successful kinetic coupling of this kind

- Data Generated by one coupled simulation is predicted to be > 10 PB/day on Summit



2. Describe your performance on Summit (or equivalent) as completely as possible. Where possible, make meaningful performance comparisons, characterize GPU utilization and identify bottleneck resources.

- XGC-GENE -- coupled performance

XGC on 512 Summit nodes
GENE on 6 Summit nodes

$N_m$ = 9,640,480 vertices
$N_p$ = 8,922 particles/vertex
timestep 61.4 seconds.

$FOM_{Summit}$ = 74.8 = **42.3** $FOM_{Titan}$
When linearly scaled to full Summit

Time spent in coupling communication using ADIOS: << 1%

Little waiting by XGC.

From FY21 WDMApp Review



Dominski, J., et al. "Spatial coupling of gyrokinetic simulations, a generalized scheme based on first-principles." *Physics of Plasmas* 28.2 (2021): 022301.
Merlo, G., et al. "First coupled GENE–XGC microturbulence simulations." *Physics of Plasmas* 28.1 (2021): 012303.
Cheng, Junyi, et al. "Spatial core-edge coupling of the particle-in-cell gyrokinetic codes GEM and XGC." *Physics of Plasmas* 27.12 (2020): 122510.

9

# Results: Seismic Tomography Workflow (PBs of data/run)

PI: Jeroen Tromp, Princeton

## Scientific Achievement

- Most detailed **3-D model of Earth**'s interior showing the entire globe from the surface to the core–mantle boundary, a depth of 1,800 miles
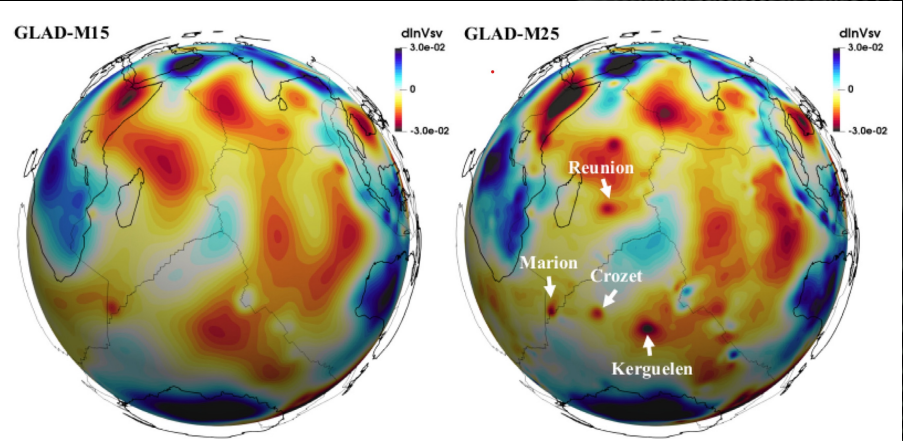
## Significance and Impact

- Updated (transversely isotropic) global seismic model GLAD-M25 where no approximations were used to simulate how seismic waves travel through the Earth. The data sizes required for processing are challenging even for leadership computer

- **7.5 PB** of data is produced in a single workflow step
  - which is fully processed later in another step.



Map views at 250 km depth of vertically polarized shear wave speed perturbations in GLAD-M15 (2017) and GLAD-M25 (2020) in the Indian Ocean. New features have emerged in GLAD-M25, such as the Reunion, Marion, Kerguelen, Maldives, Seychelles, Cocos and Crozet hotspots.

## Improvement by appending steps

- 3200 nodes ensemble run, 19200 GPUs
- 50 tasks at once
- 5.2 TB per task in 133 steps
- 260 TB total per 50 tasks
- 7.5 PB per 1500 tasks (total run)

| 50 tasks, 133 steps, 3200 nodes | Time |
|---|---|
| No I/O | 94s |
| BP3, one file per step | 235s |
| BP4 one dataset per job 133x reduction in # of files | 156s |



Wenjie Lei, Youyi Ruan, Ebru Bozdağ, Daniel Peter, Matthieu Lefebvre, Dimitri Komatitsch, Jeroen Tromp, Judith Hill, Norbert Podhorszki, David Pugmire **Global adjoint tomography—model GLAD-M25**, Geophysical Journal International, Volume 223, Issue 1, October 2020, Pages 1–21, https://doi.org/10.1093/gji/ggaa253

# FES Highlight: Established capability for near-real time networked analysis of big KSTAR data at NERSC (PPPL, ORNL, ESnet, NERSC, KSTAR, KISTI)

## Objectives

- Research and develop a streaming workflow framework, to enable near-real-time streaming analysis of KSTAR data on a US HPC
- Allow the framework to adopt ML/AI algorithms to enable adaptive near-real-time analysis on large data streams

## Impact

- Created a framework to enable US fusion researchers to have broader and faster access to the KSTAR data, enabling
  - Faster analysis of data
  - Faster and autonomous utilization of ML/AI algorithms for incoming data
  - More informed steering of experiment
  - Quicker utilization of US HPC for KSTAR collaboration



ECEI data

ADIOS DataMan

Quick analysis in ~10 minutes

NERSC

time

## Accomplishments

- Created end-to-end Python framework DELTA, streams data using ADIOS DataMan over WAN (at rates > 4 Gbps), asynchronously processes on multiple workers with MPI multi-threading
- Applied to KSTAR streaming data to NERSC Cori. Reduces time for an ECEi analysis from 12 hours on single-process to 10 minutes on 6 Cori nodes.
- Implemented deep convolutional neural networks for working with multi-scale fusion data, e.g. ECEi, for recognizing events of interest.[2]
- On-going: improve "adaptive" nature of data stream: adaptive compression at KSTAR source

Churchill RM, Klasky et al. A Framework for International Collaboration on ITER Using Large-Scale Data Transfer to Enable Near-Real-Time Analysis. Fusion Science and Technology. 2021 Feb 17;77(2):98-108
[2]R.M. Churchill, NeurIPS 2019

https://e3sm.org/pio2-adios-performance-improvement/

# SKA

Wang, Ruonan, et al. "Processing full-scale square kilometre array data on the summit supercomputer." *2020 SC20: International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*. IEEE Computer Society, 2020.

2020 Gordon Bell Nominee

- Change to ADIOS I/O: Total simulation time reduced from 9.5 hours to 6.1 hours on 1024 nodes on Summit

# Results: LAMMPS

PI: Steve Plimpton, Sandia

Results from ECP EXAALT Q4/FY19 milestone report (for 2.2.1.04 EXAALT ADSE04-54)
Summit 512 nodes
12B atoms, 5 TB

https://github.com/lammps/lammps/tree/master/src/USER-ADIOS

- USER-ADIOS package in LAMMPS for dump commands
  - dump atom/adios
  - dump custom/adios
- Output goes into an I/O stream
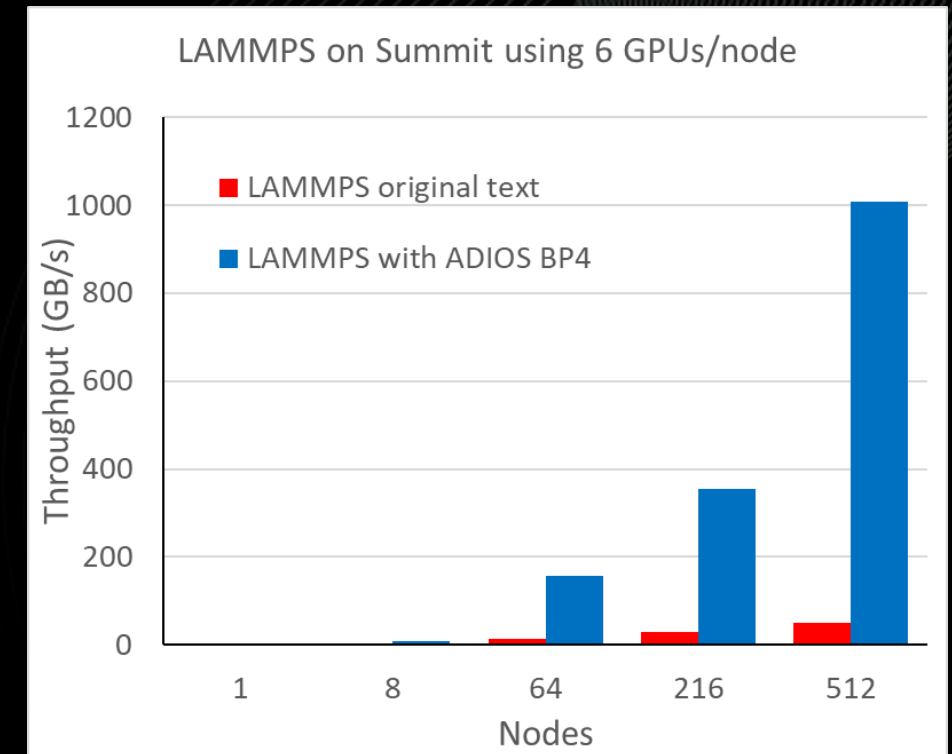  - BP4 file by default
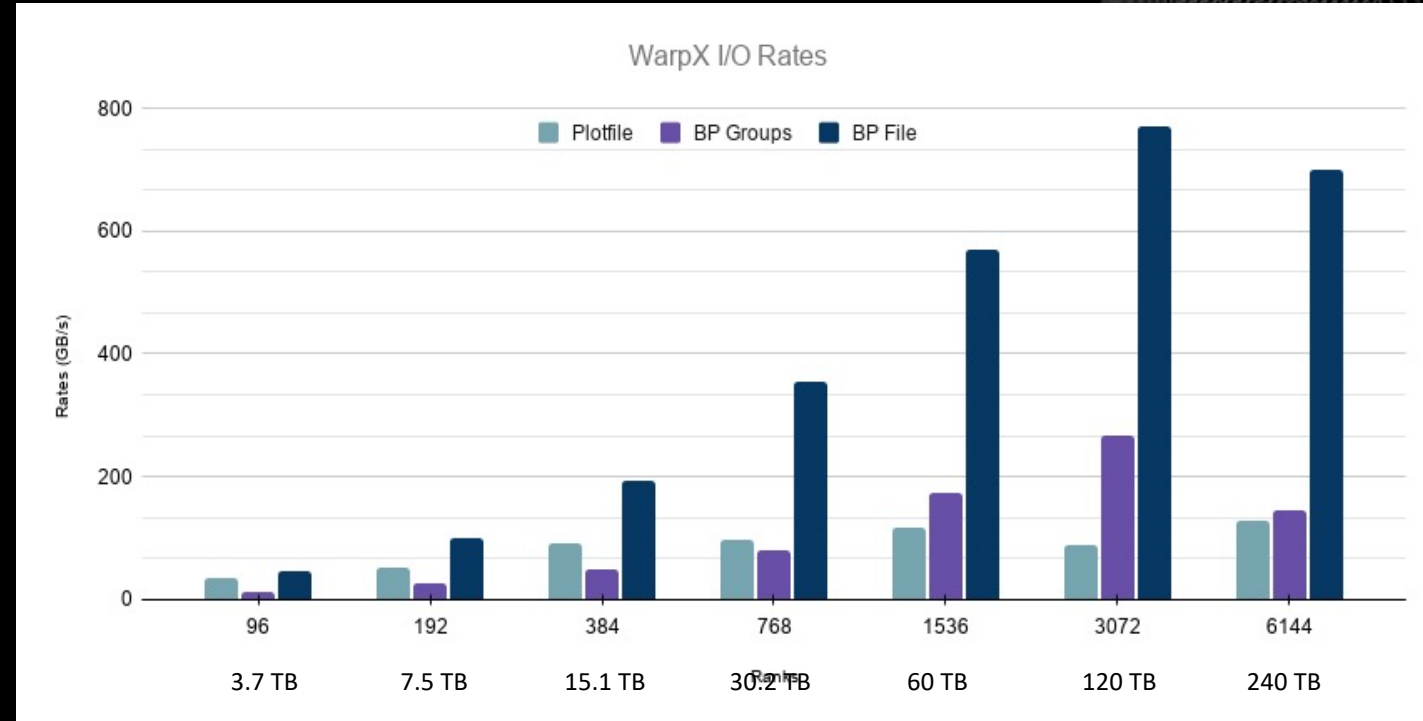  - Can use staging engines
- Concurrent reading is enabled



LAMMPS on Summit using 6 GPUs/node

Legend:
- LAMMPS original text (red)
- LAMMPS with ADIOS BP4 (blue)

Y-axis: Throughput (GB/s)
X-axis: Nodes (1, 8, 64, 216, 512)

# Results: WarpX

- BPFile:
  - BP4: Use one file for all outputs.
- BPGroups:
  - BP3: Use one file / timestep
- Plot:
  - The AMReX plot file.
- One way to improve the I/O performance, is to use one ADIOS file for all time steps



Summit, 6 GPUs, 6 cores per node, up to 1024 nodes

# Writing performance is great but what about reading?

- Codes such as the WarpX code, which uses AMReX can take advantage of ADIOS-BP4 for "fast" writing

- The challenge is reading

- Development of a clustering algorithm for WarpX/AMReX data for fast writing/reading performance



Lipeng Wan et al.: "Data Layout Strategies for Parallel I/O: The Good, The Bad and The Ugly", submitted to TPDS journal 2021/March