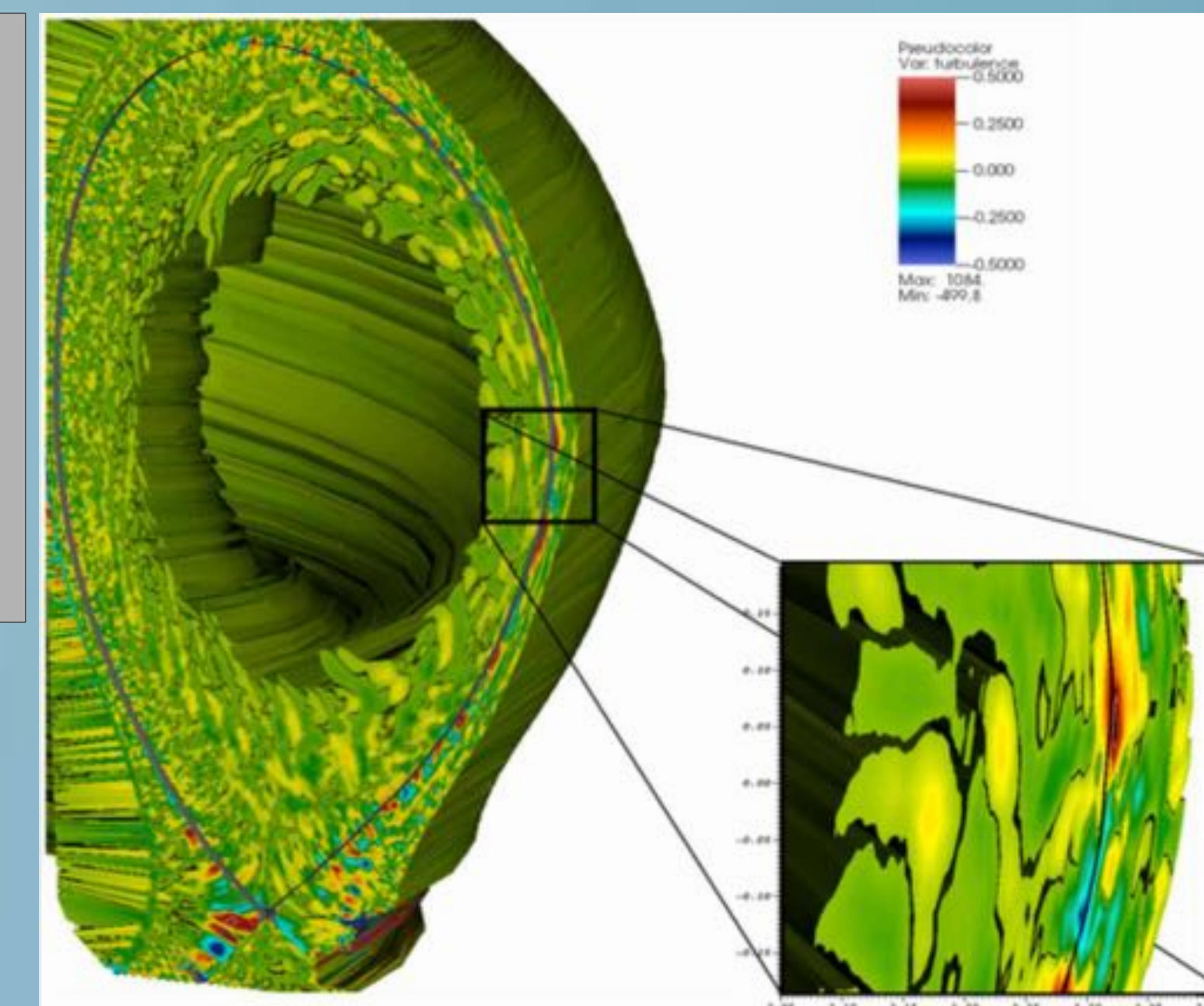# Data Management Challenges In HBPS

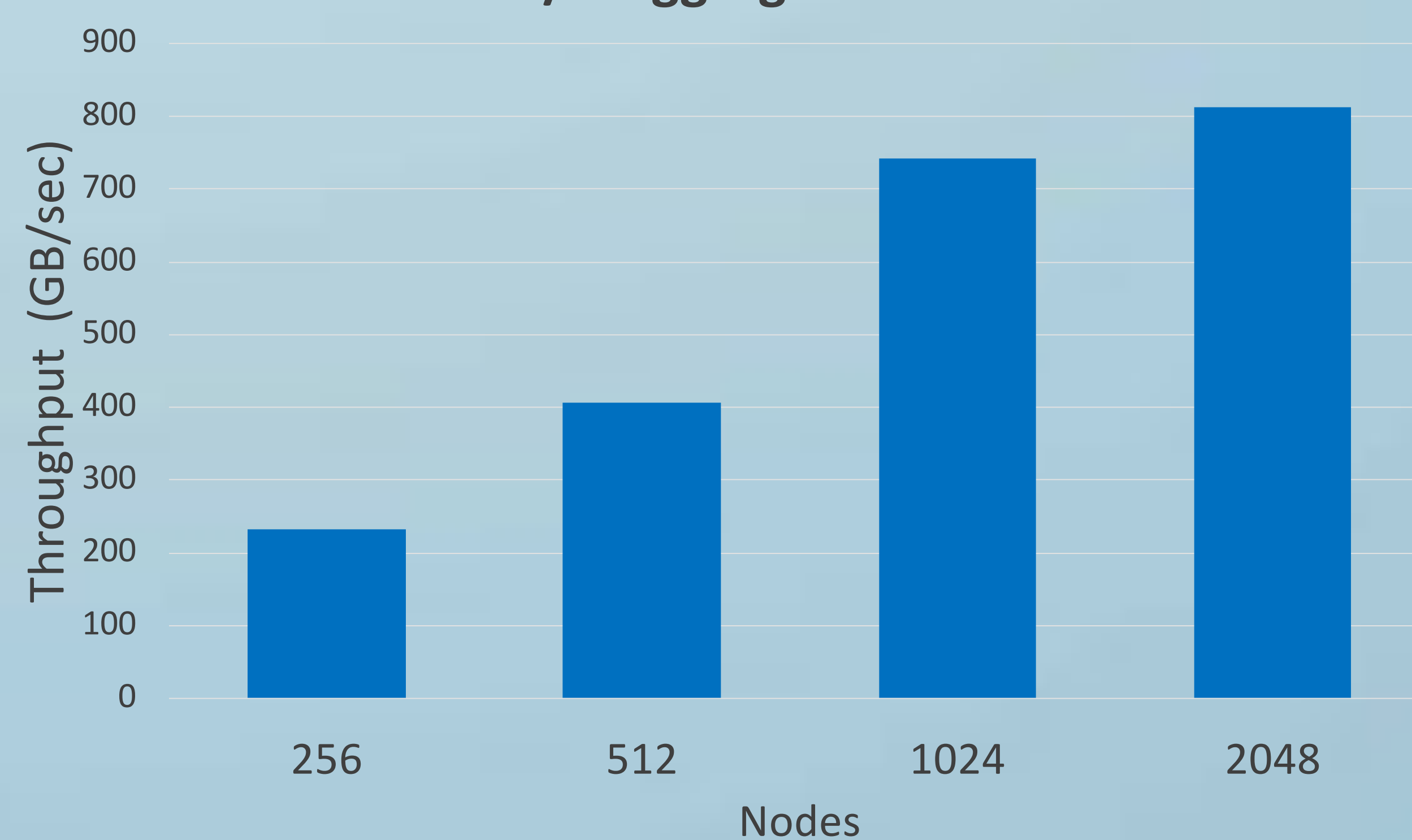Jong Youl Choi[1], Michael Churchill[2], Davide Curreli[3], Sonata Mae Valaitis[3], Robert Hager[2], Seung-Hoe Ku[2], E. D'Azevedo[1], Bill Hoffman[4], David Pugmire[1], Scott Klasky[1], C. S. Chang[3]

[1]ORNL, [2]PPPL, [3]Univ. of Illinois Urbana-Champaign, [4]Kitware

**HBPS** High-fidelity Boundary Plasma Simulation

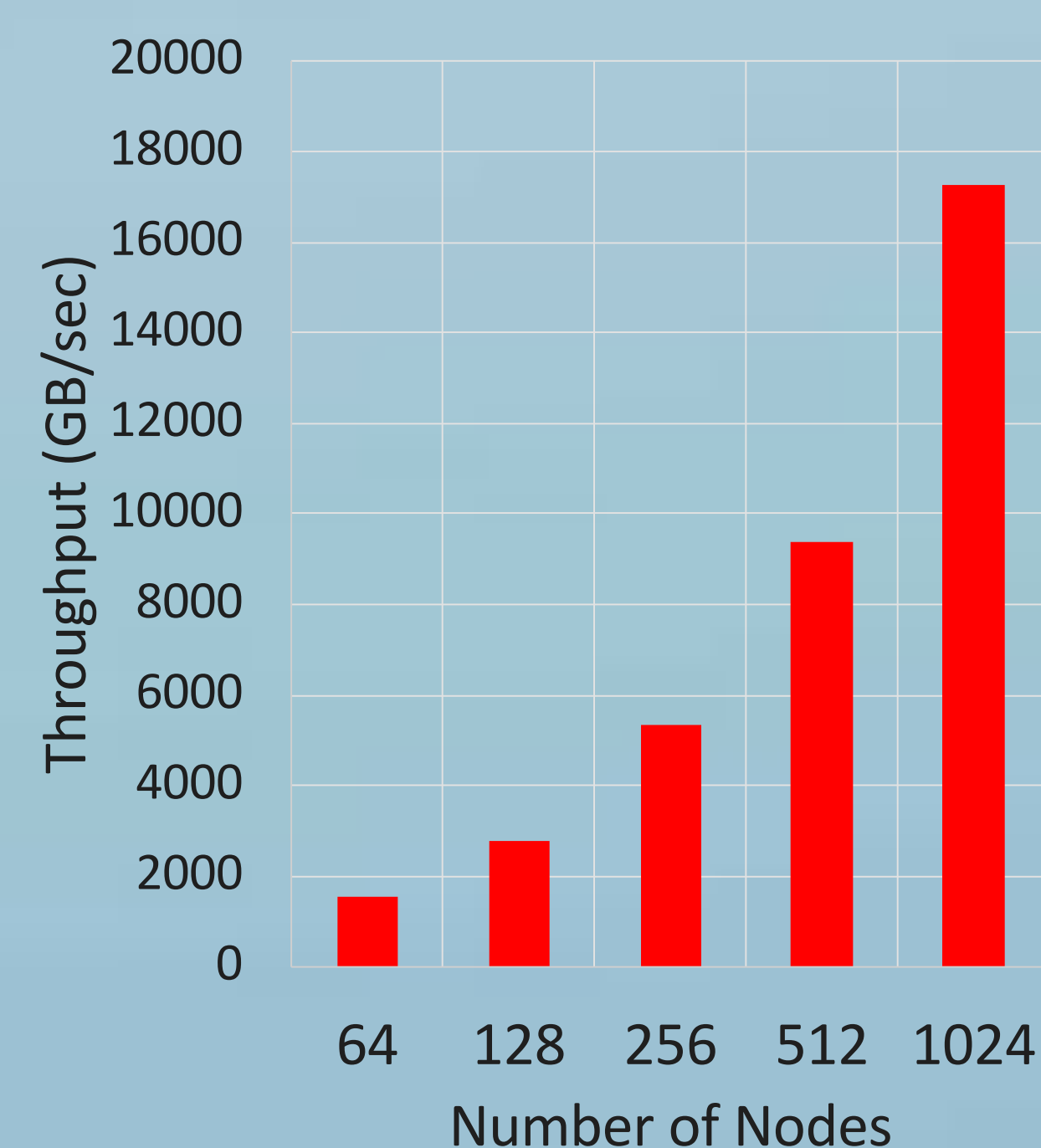

## XGC I/O Performance

We maintain cutting edge I/O performance for XGC on various file systems, including SSDs and NVMe, on Cori, Theta, and Summit.
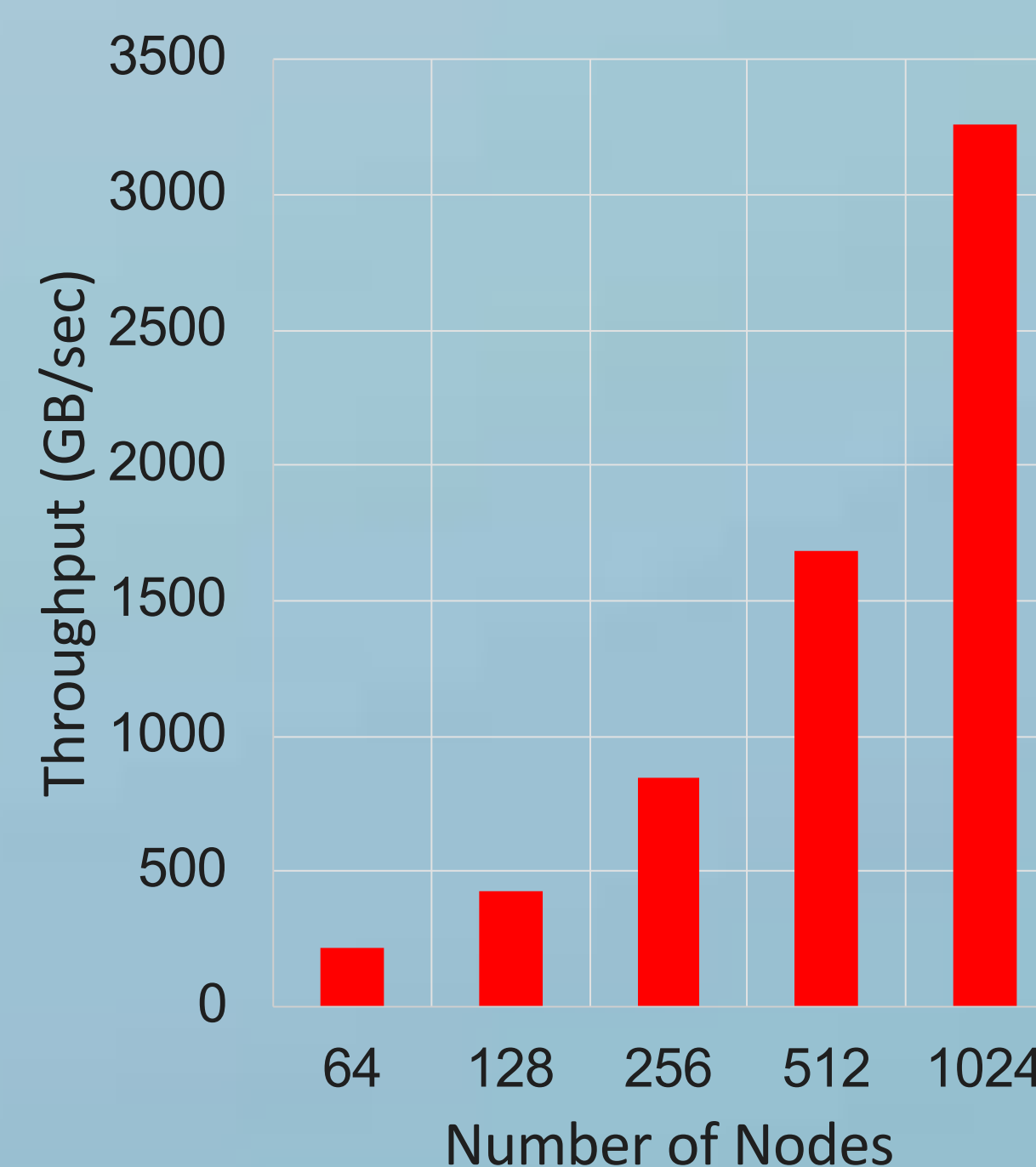
**XGC Checkpoint Writing on Summit GPFS with I/O aggregation**





Our team continues to innovate to take full advantage of the new memory and storage technologies, and to provide the highest levels of performance.

| I/O System | Summit ORNL | Theta ANL | Cori NERSC |
|---|---|---|---|
| Locality | Node local | Node local | Remote Shared |
| System | Local filesystem | Local filesystem | Cray WARP |
| Capacity | 800 GB per node | 128 GB per node | 288 Server 50 TB limit per job |
| Parallel Filesystem | GPFS Lustre | Lustre | Lustre |

## XGC Software Process

Agile XGC development
- Incorporate a modern CMake build system
- Continuous Integration testing system
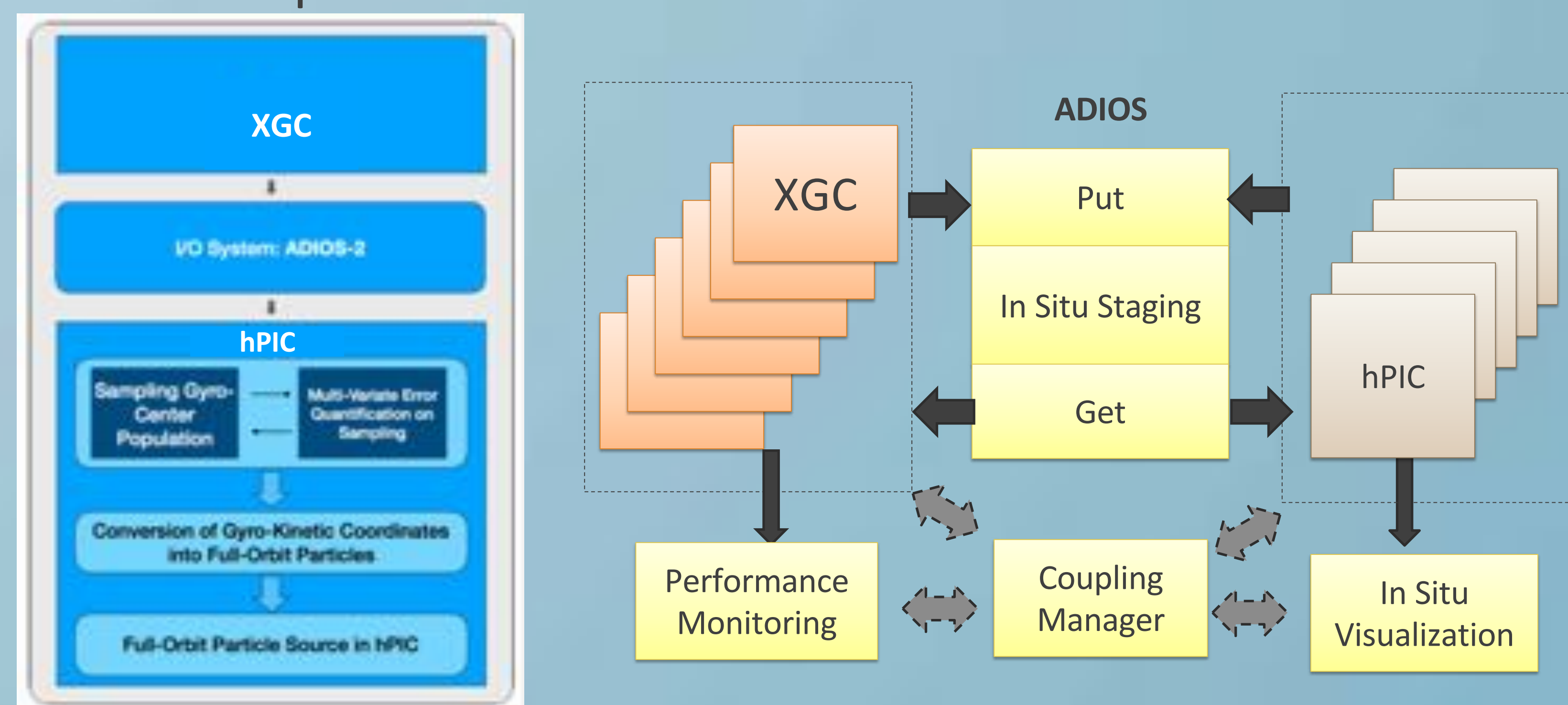- Git workflow incorporated with CI system
- Integrate CDash into github



## Coupling Workflows

**The Fusion HPBS project is focusing on researching multi-way coupling science to study multi-scale/multi-physics.**
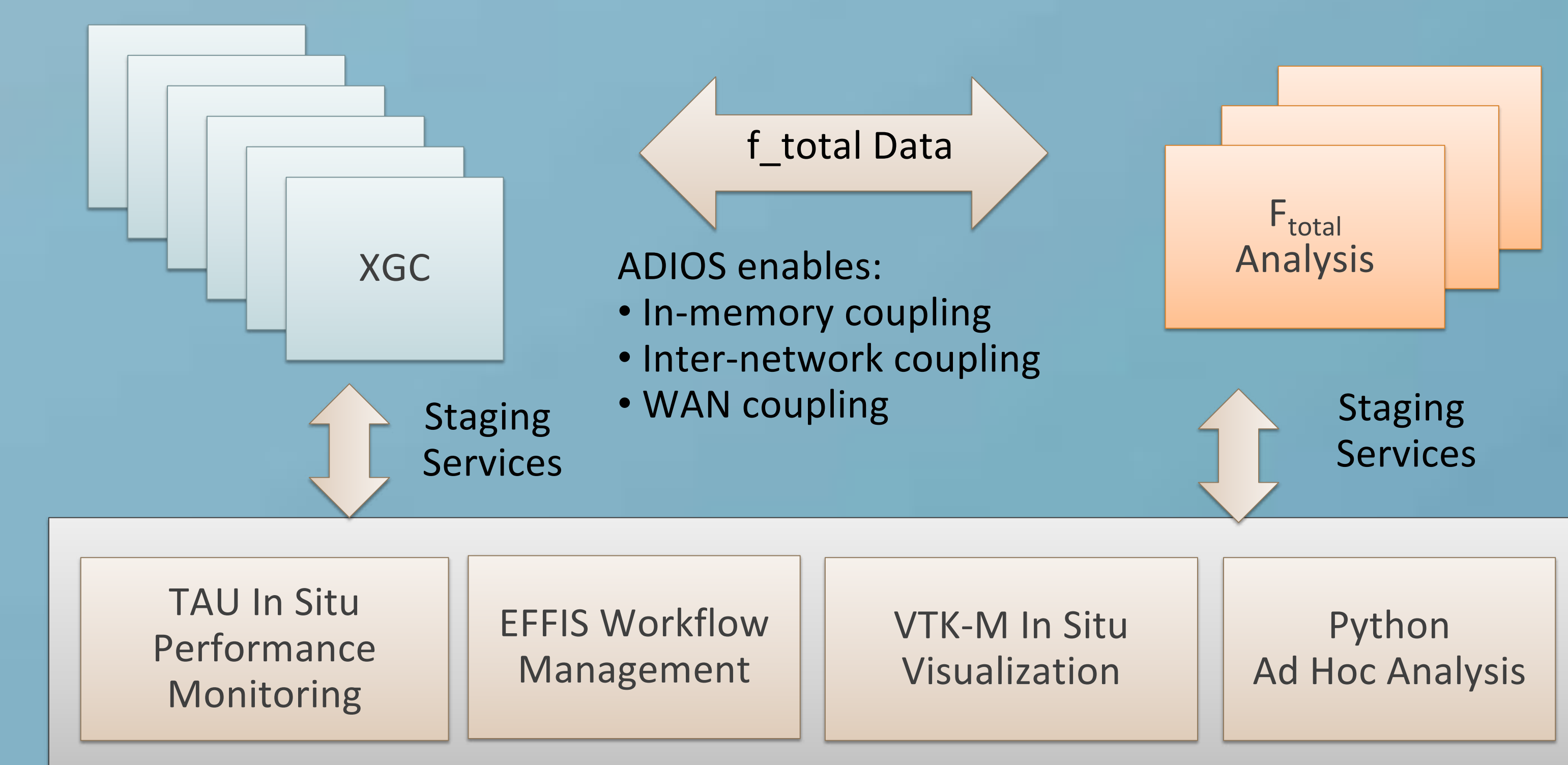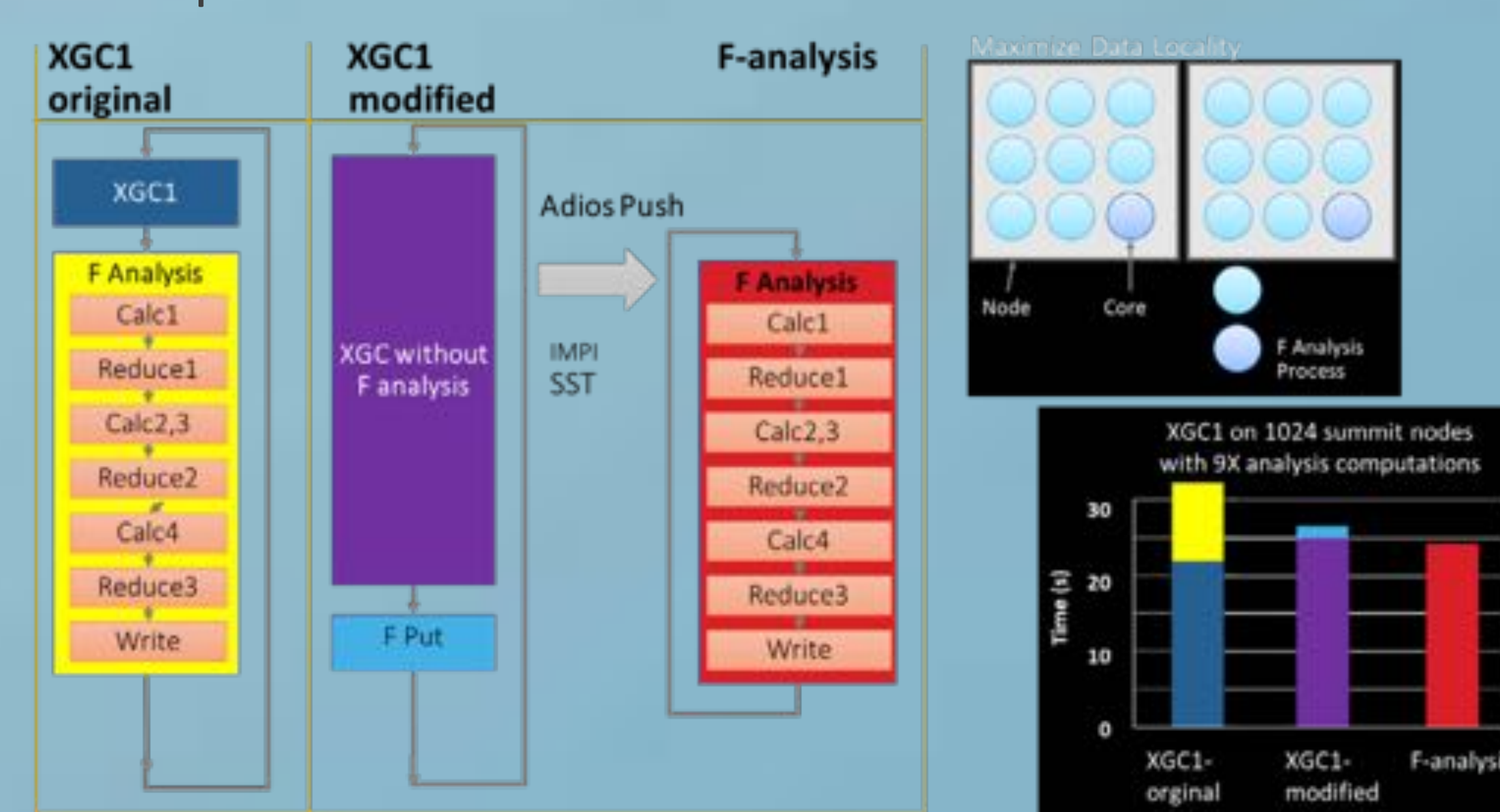
### 1) XGC and hPIC
- Plasma-material-interaction hPIC code coupled into XGC
- hPIC code has 6D marker particles, while XGC has 5D marker particles



### 2) XGC and F-analysis coupling

In XGC and F analysis coupling, we move the F computation to a dedicated analysis code. XGC asynchronously offloads those computations via ADIOS and improves computational performance





ADIOS enables:
- In-memory coupling
- Inter-network coupling
- WAN coupling

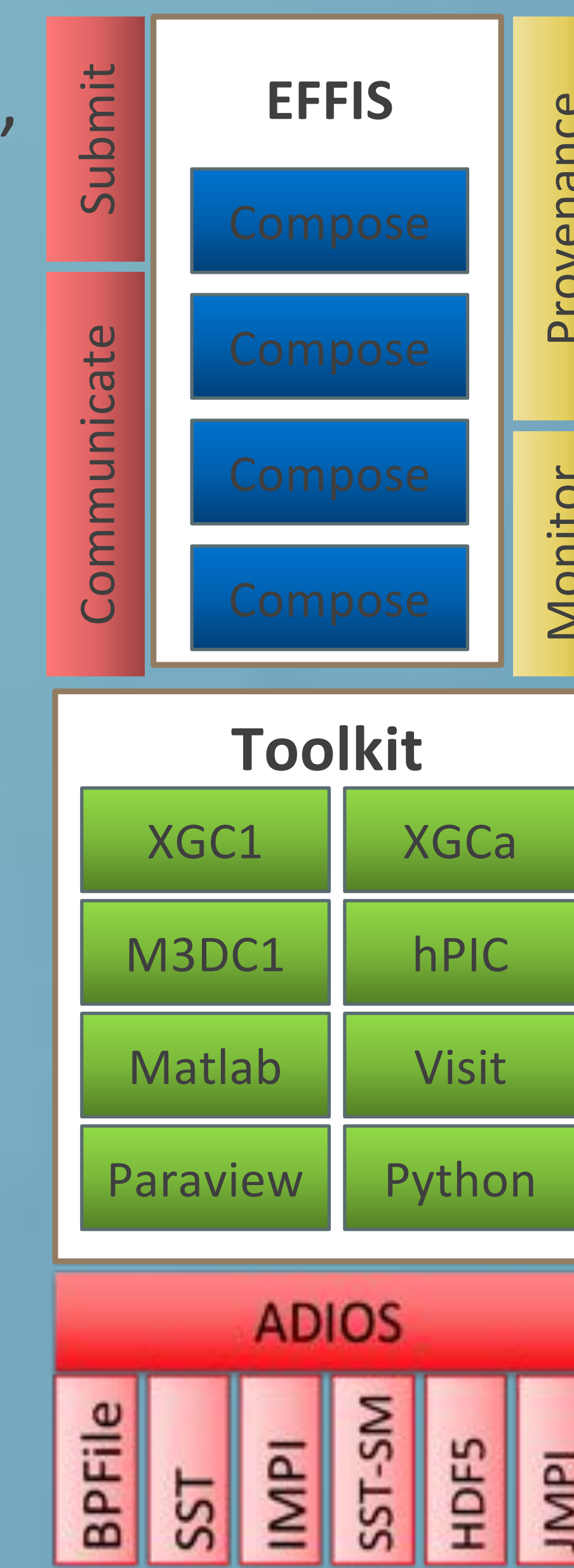| TAU In Situ Performance Monitoring | EFFIS Workflow Management | VTK-M In Situ Visualization | Python Ad Hoc Analysis |
|---|---|---|---|

### Research Details
a) To improve movement performance and flexibility, HBPS integrated with ADIOS for data management.
b) Developing multi-way coupling science cases to study multi-scale/multi-physics scenarios.
c) Exploiting data locality to improve performance
   - XGC computes 5D $f$ and electromagnetic field
   - Hand-off computational reduction of physics from XGC
   - Analysis code consumes in-memory f data

## EFFIS

**EFFIS** is an integrated platform of services to compose, launch, monitor, and control coupled applications.

EFFIS can simplify the complexity of composing, running, and monitoring applications on HPC systems.

We integrate **HBPS** with **EFFIS** to "easily" compose coupled HBPS workflows on HPC Resources (Cori, Theta, and Summit).

EFFIS's using a python-like interface can allow "easy" integration to visualization tools (Visit, Python notebooks)



```
run:
  xgc:
    processes: 1024
    processes-per-node: 32
    path: xgc-build/xgc1-es
    groups:
      diagnosis.1d:
        plot:
          psi-plot:
            x: psi
            y: i_gc_density_1d
```
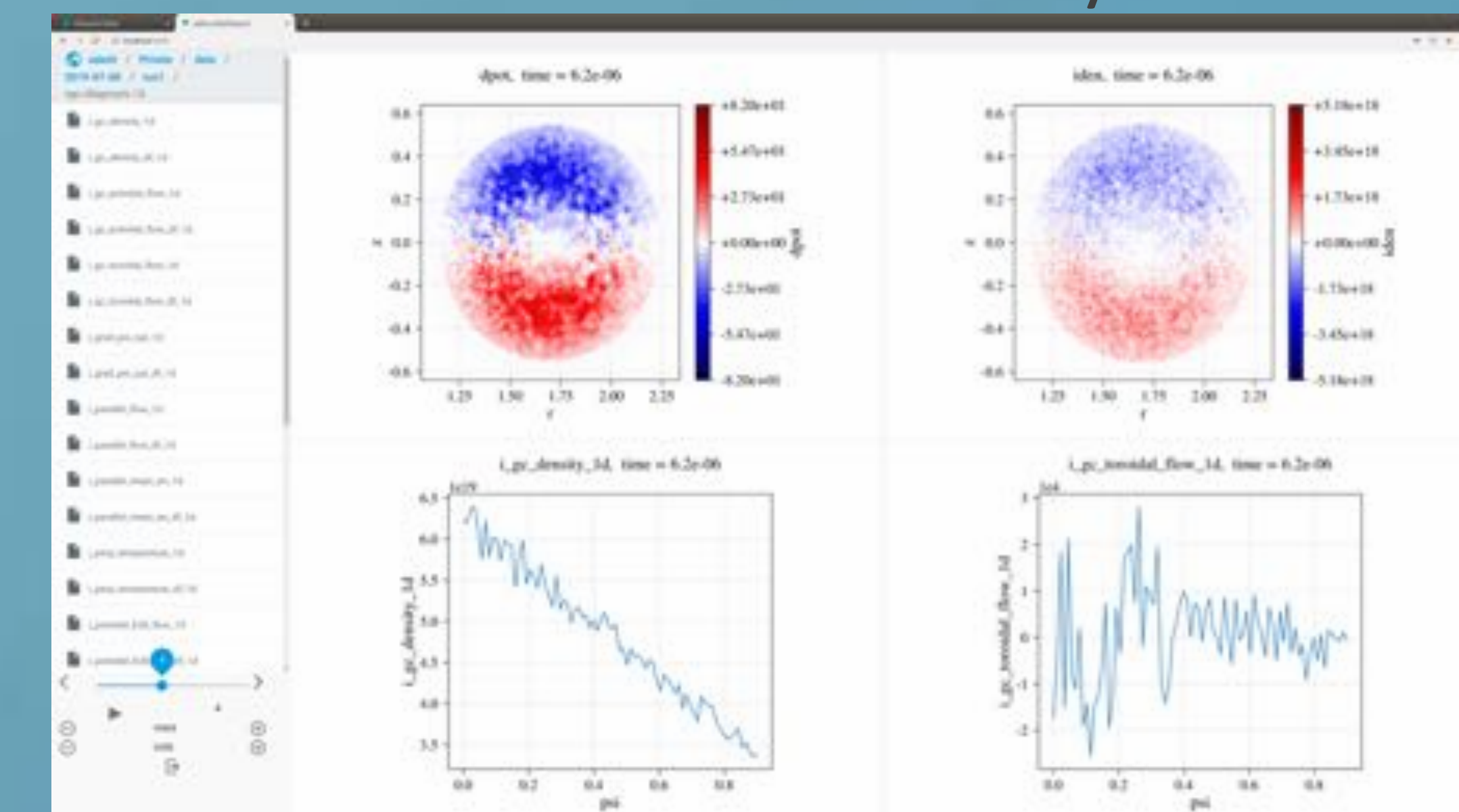
Example of EFFIS specification file. XGC run is configured to run with analysis application.

```
//@effis-begin reader_io-->"ConcentrationData", step=kstep;
adios2::IO reader_io = ad.DeclareIO("SimulationOutput");
adios2::IO writer_io = ad.DeclareIO("PDFAnalysisOutput");
```

Example of EFFIS instrumentation in XGC code using simple @effis pragmas.

EFFIS is integration with HBPS can provide:
- High Performant I/O for multiple codes
- Process placement (node sharing, co-location of codes on a node, etc.)
- Online dashboard functionality



- Services for concurrent analysis/visualization
- Run archival (e.g. long-term tape storage)
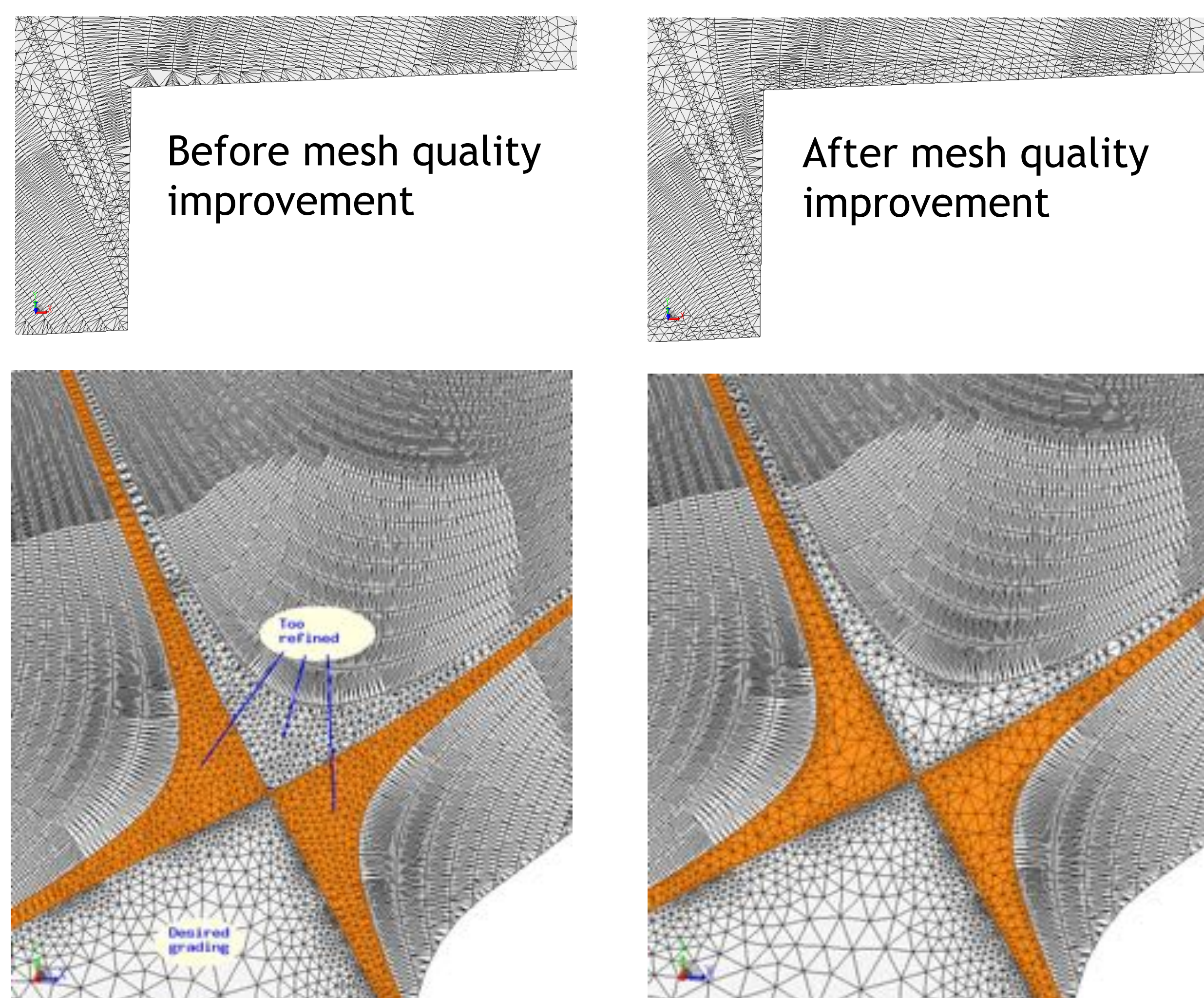- Source code association with runs

# Performance Enhancements of XGC

E. D'Azevedo[1], A, Scheinberg[2], M. Shephard[3], P. Worley[4], S. Sreepathi[1], B. MacKie-Mason[5], T. Willians[5], and the SciDAC HBPS XGC Team

1. Oak Ridge National Laboratory , 2. Princeton Plasma Physics Laboratory, 3. Rensselaer Polytechnic Institute, 4. PHWorley Consunting, 5. Argonne National Laboratory

**HBPS** High-fidelity Boundary Plasma Simulation

## XGC Meshing

- Improved mesh quality in areas where flux curves interact with reactor wall
- Improved matched mesh gradation at x-point
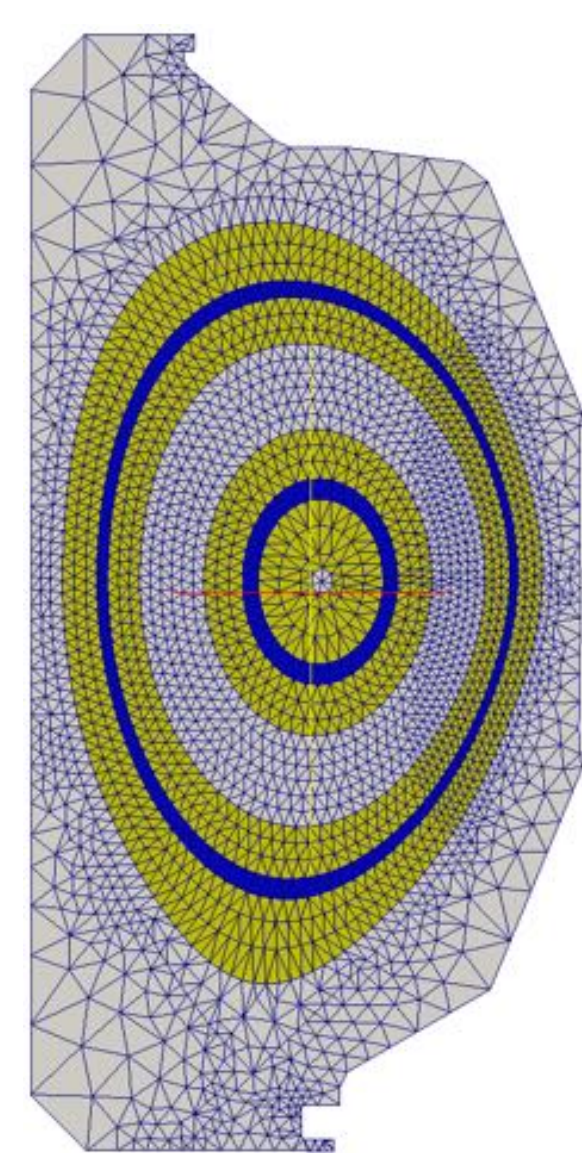- Reordering of mesh data for better memory access during XGC simulations

Before mesh quality improvement

After mesh quality improvement

Improved mesh gradation at X-point

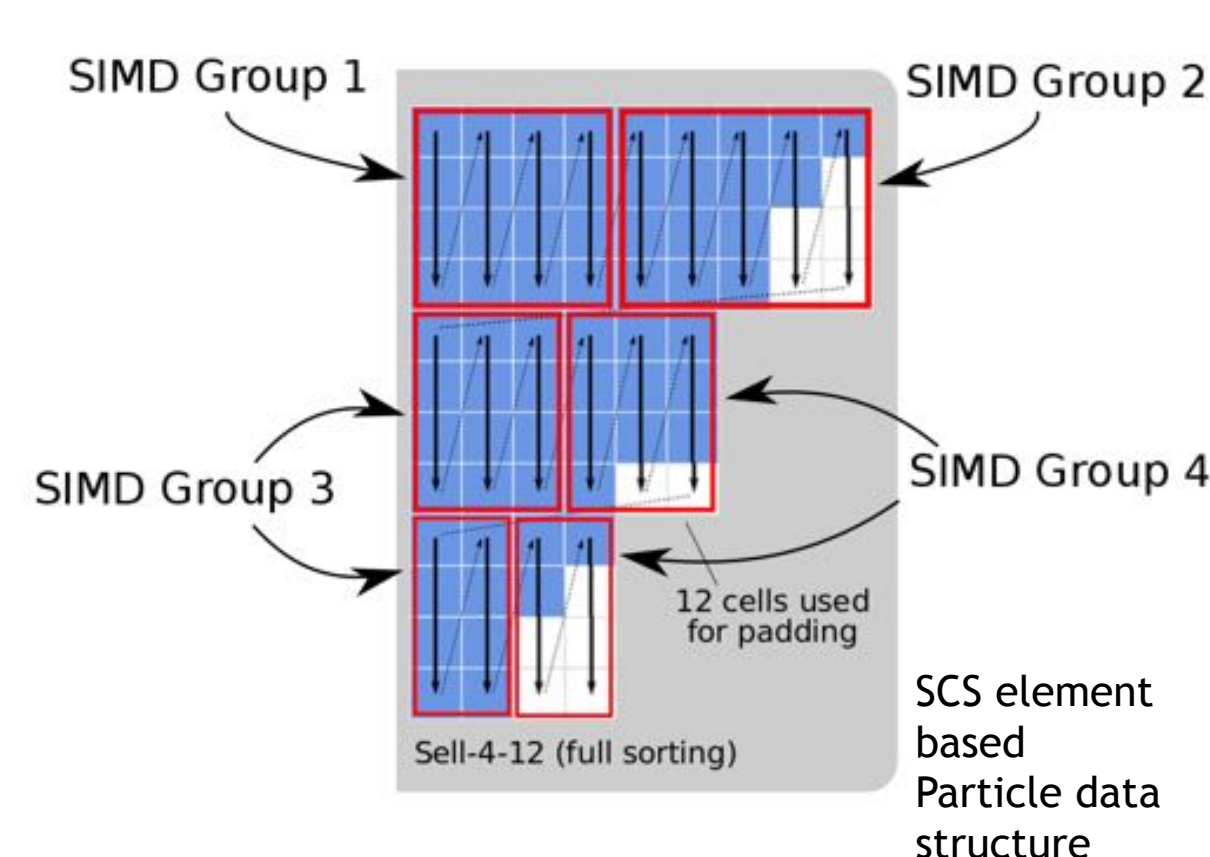## XGC based on Parallel Unstructured Mesh PIC (PUMIpic)

PUMIpic – Components to support PIC operations on distributed unstructured meshes (2D and 3D)

- Mesh centric – no independent particle structure
- Distributed mesh with overlaps (PICparts)
- Particle migration and load balancing between pushes
- Adjacency-based particle containment determination
- Focused on structures for execution on GPUs
  - Omega GPU ready mesh topology being integrated
  - Particles stored by element in new SCS data structure
- Test shows on-par performance using less memory

Two PICparts

| ptcls (Ki) | no sorting time (s) | full sorting time (s) |
|---|---|---|
| 128 | 2.298661 | 3.642041 |
| 256 | 2.895464 | 3.415048 |
| 512 | 3.79263 | 3.851178 |
| 1024 | 4.972283 | 4.090044 |
| 2048 | 7.089673 | 4.389198 |
| 4096 | 11.578984 | 4.799475 |

SIMD Group 1
SIMD Group 2
SIMD Group 3
SIMD Group 4
12 cells used for padding
Scs1-4-12 (full sorting)
SCS element based Particle data structure

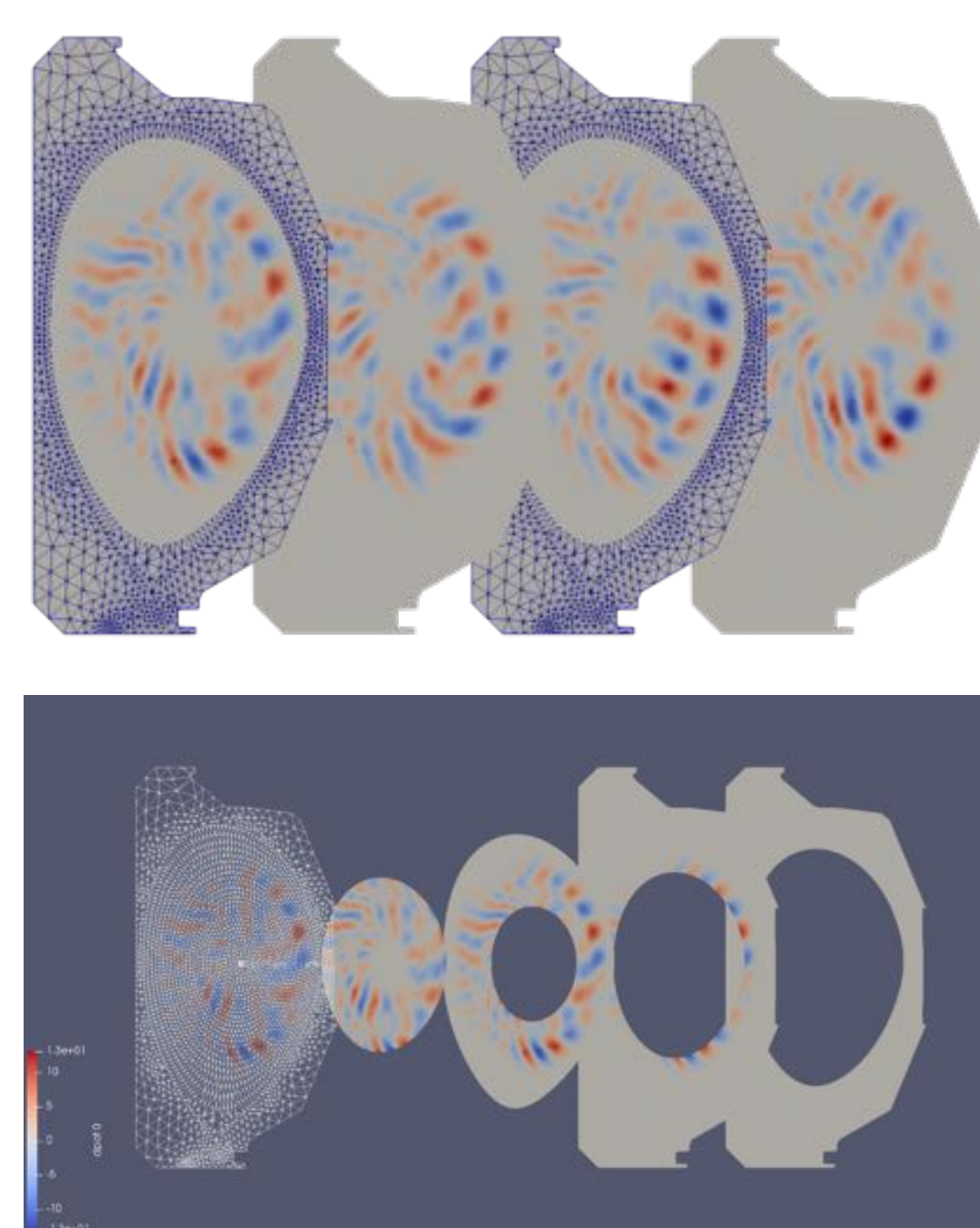**Implementing XGC physics and Numerics with PUMIpic:**
- Since all core data structures are changed code, code being rewritten in C++
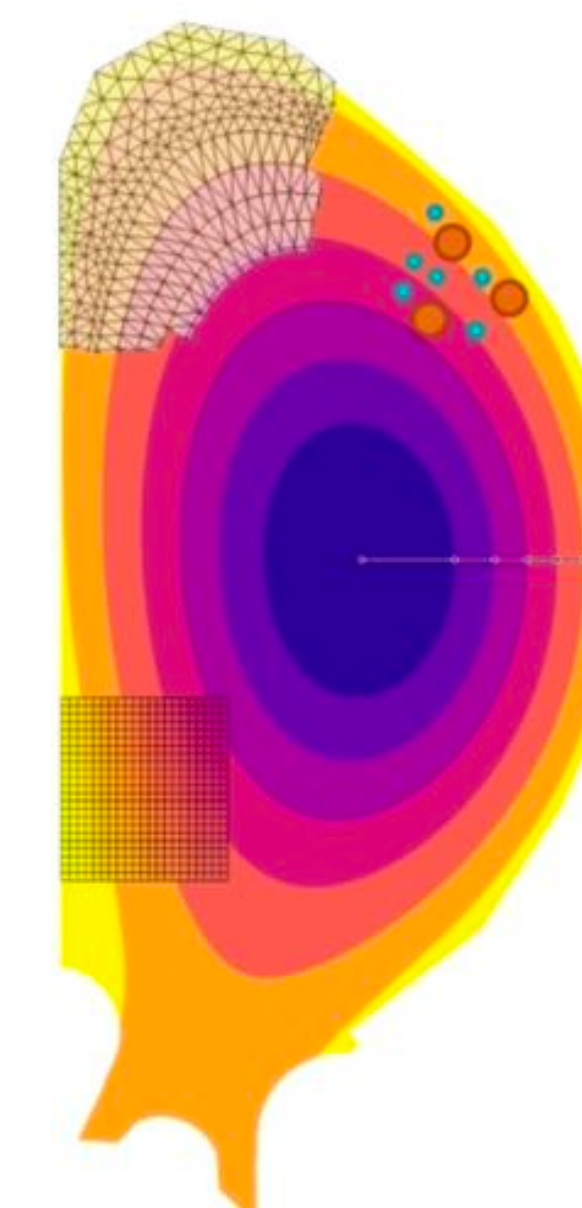
**Status of implementation:**
- Based on original PUMI structures – new GPU focused structures will be integrated when complete
- Core mesh/particle interaction operations in place
- Mesh solve in place
- Ion and electron push (including subcycling) implemented
- Initial $\delta f$ simulations executed
- Performance evaluation and improvement underway
- Initial push results show 25% improvement on many core system
- Other steps slower due to need to modify mesh copies (underway)

Snapshot of electrostatic potential fluctuation (a) at toroidal angle $\zeta=0,\pi/2,\pi,3\pi/2$ from left to right and (b) in local domain of each group at $\zeta=0$

## XGC on Summit

- XGC is part of Early Science Programs on Summit, Aurora and Perlmutter
- XGC is an ECP code
- XGC uses an unstructured grid in poloidal plane, each MPI rank gets particles from a section of poloidal plane
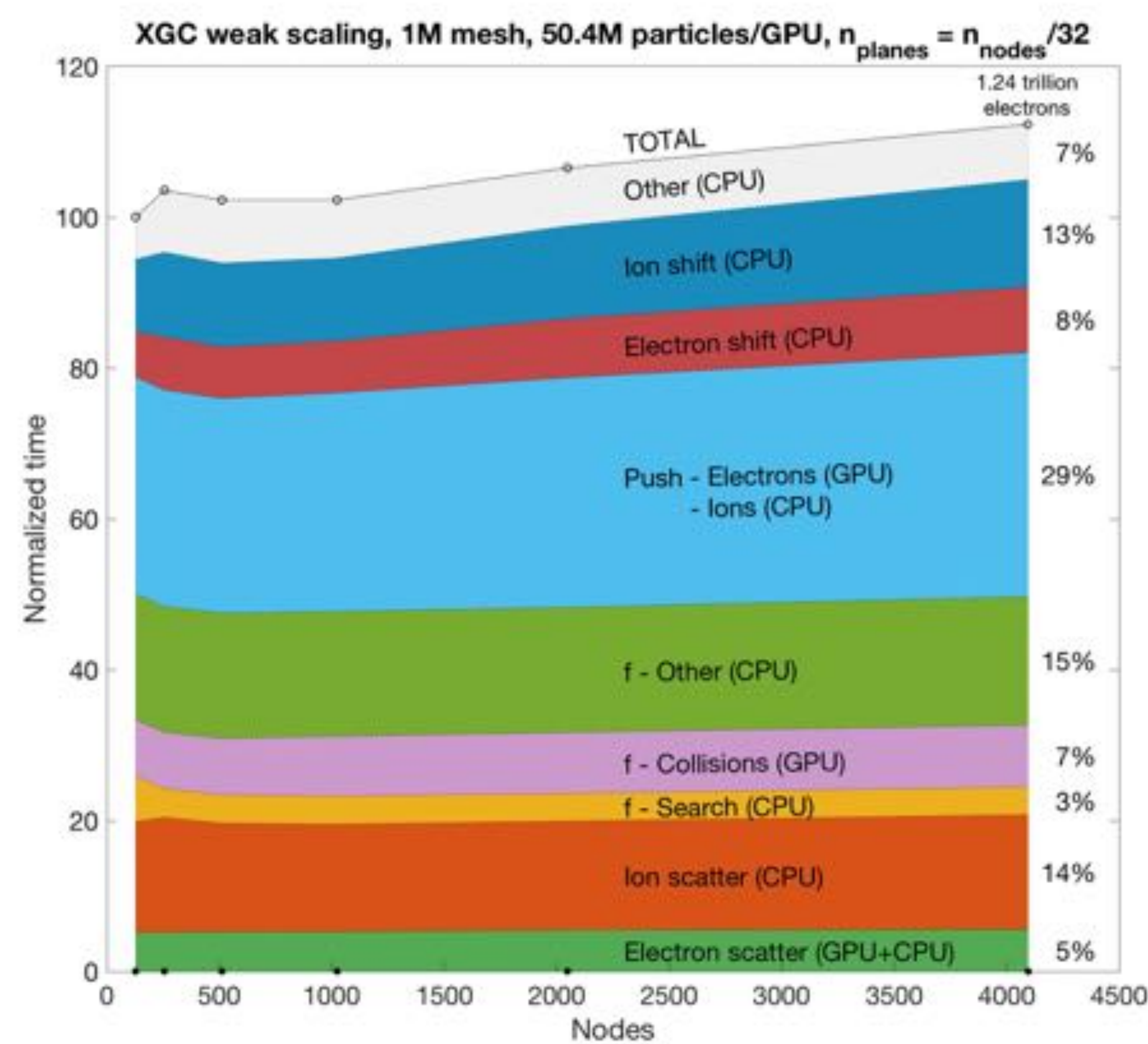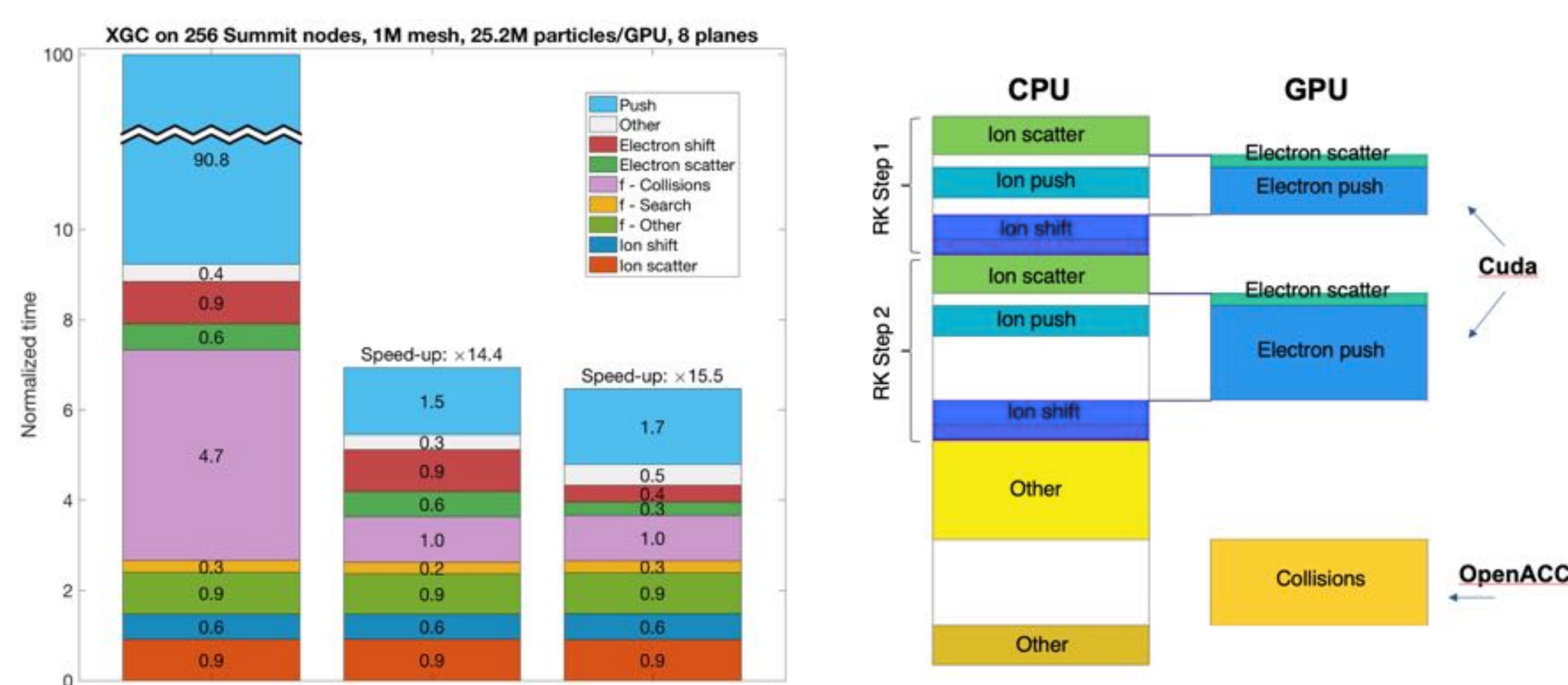- Main computational kernel is electron push
- Utilizes Kokkos

```
XGC_core/pushe.F90:

subroutine pushe
  call sort_particles      ! Sort particles by grid cell
  do iptl=1, n_particles   ! Loop over particles
    do ic=1, n_cycles      ! Subcycle electrons
      do irk=1, n_runge_kutta  ! RK4 loop
        call search            ! Determine which grid cell particle inhabits
        call gather_field    ! Interpolate field at particle location
        call calculate_dx    ! Solve physics: dx/dt = f(E,…)
        call advance_particles  ! Update particle position and velocity
      end do
    end do
  end do
end subroutine pushe
```

### Good Weak Scaling to Full Summit

- On 256 nodes of Summit, GPU version has **15X speedup over CPU only**
- Good weak scaling up to full Summit using **1.24 trillion electrons on GPU and 1.24 trillion ions on CPU**

XGC on 256 Summit nodes, 1M mesh, 25.2M particles/GPU, 8 planes

Speed-up: ×14.4
Speed-up: ×15.5

CPU / GPU
Cuda
OpenACC

XGC weak scaling, 1M mesh, 50.4M particles/GPU, $n_{planes} = n_{nodes}/32$

1.24 trillion electrons

| | |
|---|---|
| Other (CPU) | 7% |
| Ion shift (CPU) | 13% |
| Electron shift (CPU) | 8% |
| Push - Electrons (GPU) - Ions (CPU) | 29% |
| f - Other (CPU) | 15% |
| f - Collisions (GPU) | 7% |
| f - Search (CPU) | 3% |
| Ion scatter (CPU) | 14% |
| Electron scatter (GPU+CPU) | 5% |

## Details on XGC-Kokkos

- XGC in Fortran, Kokkos in C++
- Fortran interface (Cabana) enables easy porting of new kernels
- Single code for CPU and GPU
- Electron push kernel in CUDA Fortran (C++ version under development)

Must cast Cabana array into predefined Fortran type for use in Fortran kernels using ISO_C_BINDING

## Performance on KNL

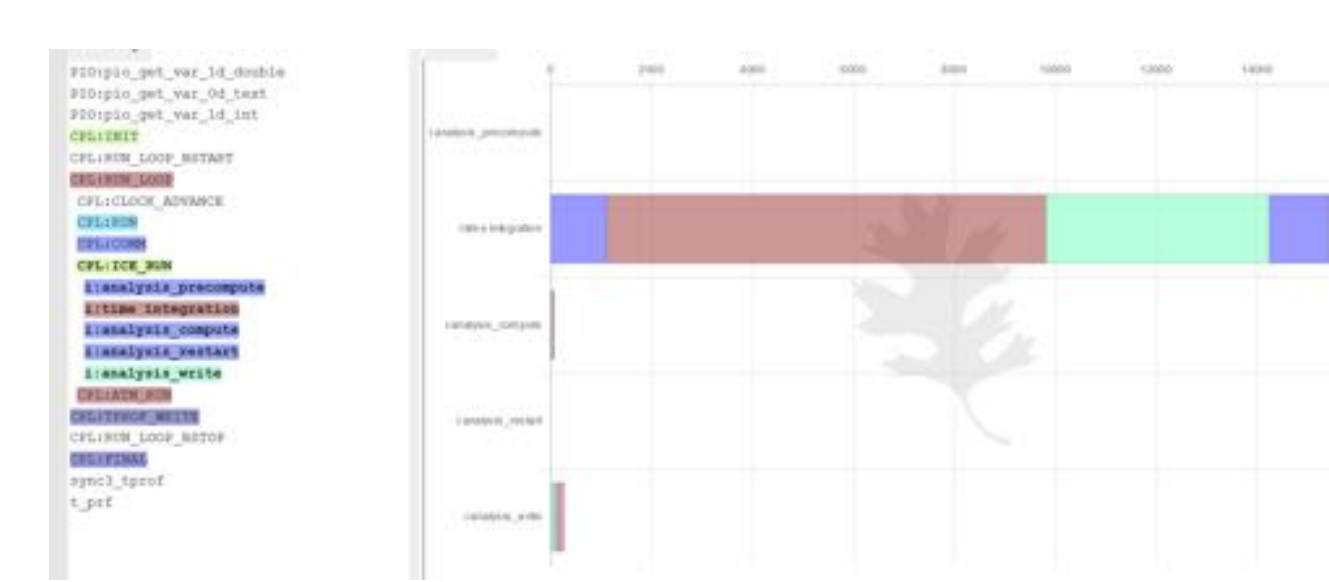- Kokkos version of XGC has been ported to Cori KNL
- Roofline analysis of vectorized version of XGC shows in-lining and re-factoring useful in optimizing use of wide-vector registers. However, vector dependences and data type conversions limiting peak performance

XGC versions on Cori KNL, 370k mesh, 12M particles/node, 2 planes, 512 KNL nodes

Speed-up: ×1.5
Speed-up: ×1.5

## Performance Analytics for Computational Experiments for XGC

- Central hub of performance data, already used in Climate application
- Interactively deep-dive and track performance benchmark
- Facilitate performance analysis:
  - Load balancing
  - Identification of bottlenecks
  - Inform targeted optimization efforts

https://pace.ornl.gov

Tree and Flame Graphs