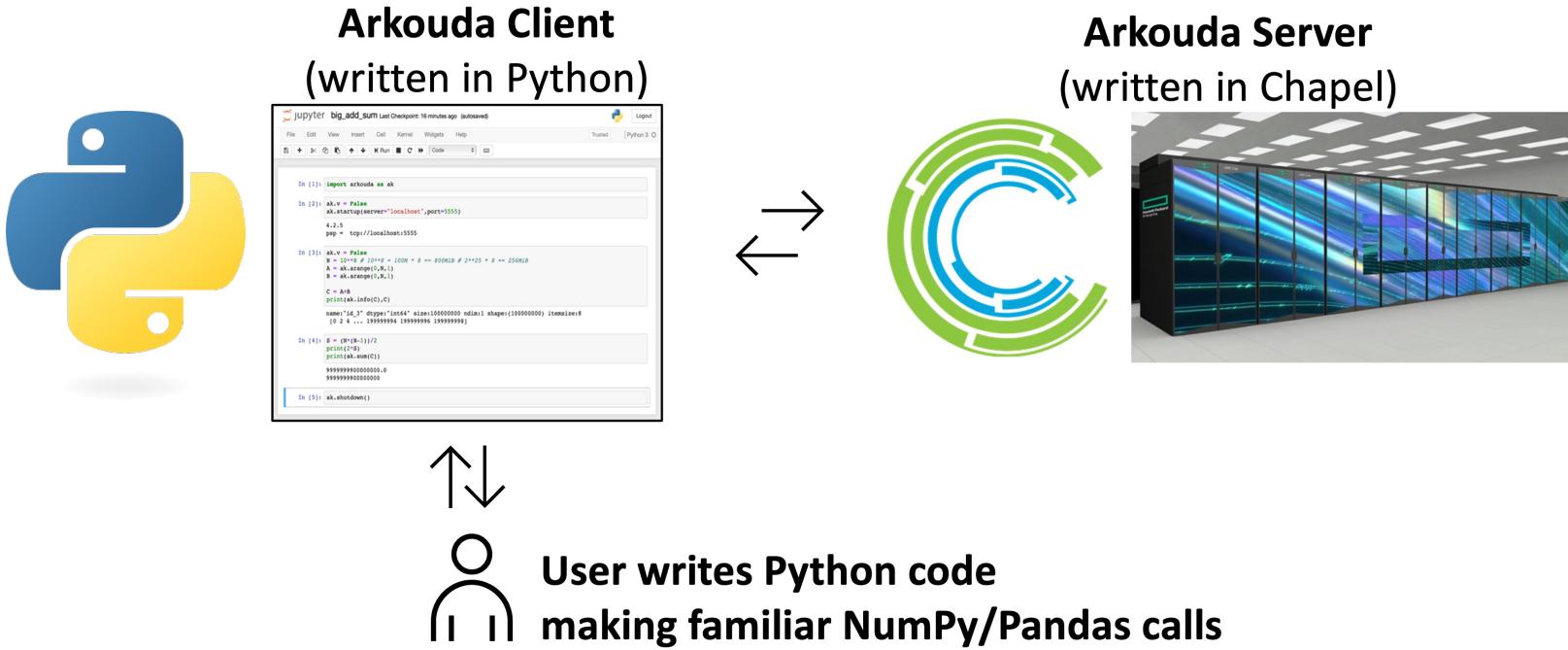# Interactive, HPC-scale Exploratory Data Analysis in Arkouda: Past Successes and Future Challenges

Brad Chamberlain, Advanced Programming Team, HPE

Productive, Performant Software for Large-Scale Scientific Data Analysis, SLAC
October 21, 2025

# What is Arkouda?

**Q:** "What is Arkouda?"

**Arkouda Client**
(written in Python)

**Arkouda Server**
(written in Chapel)



**User writes Python code**
**making familiar NumPy/Pandas calls**

**A1:** "A scalable version of NumPy / Pandas for data scientists"

**A2:** "An extensible framework for using supercomputers interactively from Python"

# Key Properties of Arkouda

— **Columnar:** represents dataframes using a distributed array per column

— **Extensible:** new features can be added to the server and/or client

  • e.g., NJIT's Arachne extension for graph analytics

— **Open-Source:** developed on GitHub, released under the MIT license

— **Portable:** runs on virtually any system (laptop, cluster, cloud instance, supercomputer)

— **Interactive:** operations are designed to complete in seconds to small numbers of minutes

— **Scalable:** has scaled to hundreds of TB, thousands of computes nodes, and over a million processor cores

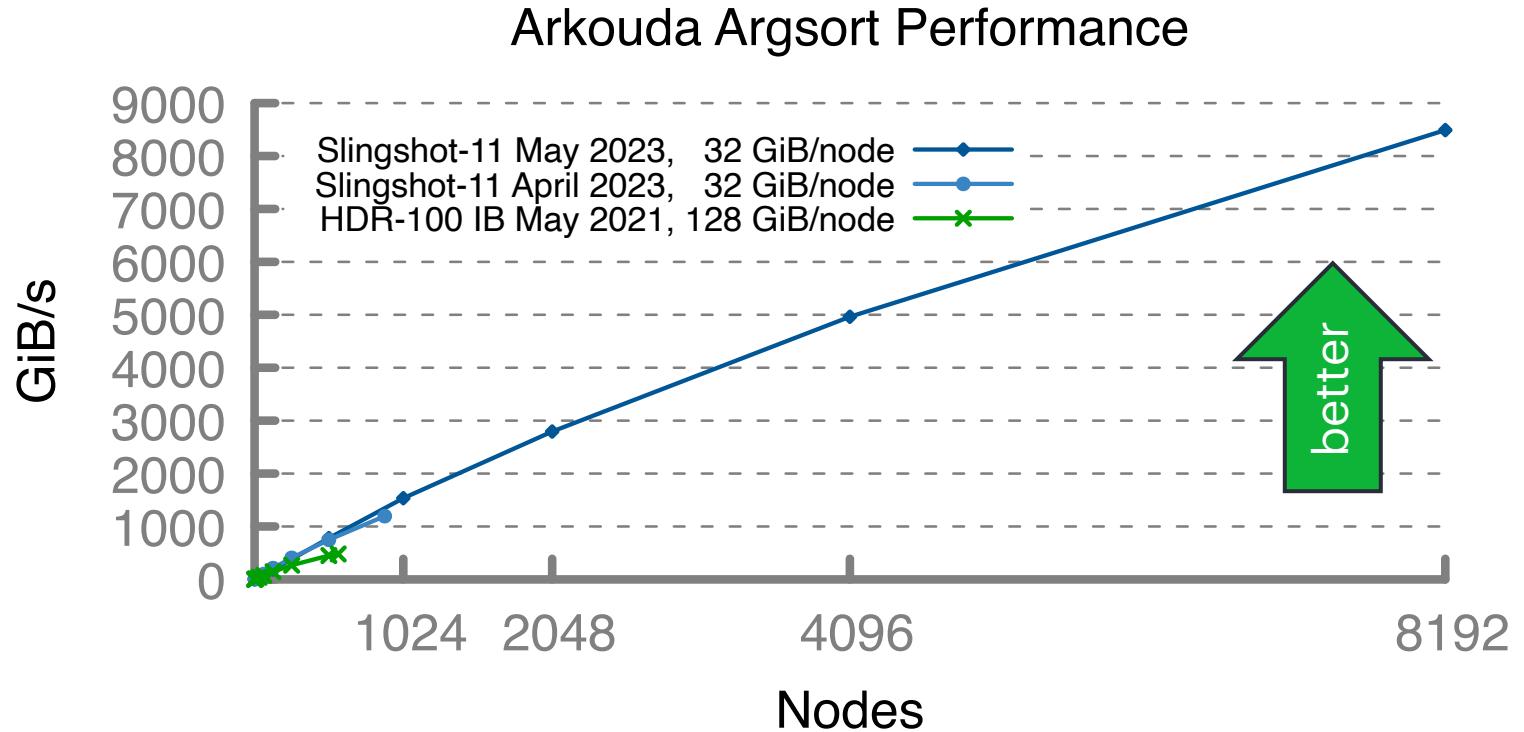# Performance and Productivity: Arkouda Argsort

**HPE Cray EX** ◆——◆

— Slingshot-11 network (200 Gb/s)
— 8192 compute nodes
— 256 TiB of 8-byte values
— ~8500 GiB/s (~31 seconds)

**HPE Cray EX** ●——●

— Slingshot-11 network (200 Gb/s)
— 896 compute nodes
— 28 TiB of 8-byte values
— ~1200 GiB/s (~24 seconds)

**HPE Apollo** ✕——✕

— HDR-100 InfiniBand network (100 Gb/s)
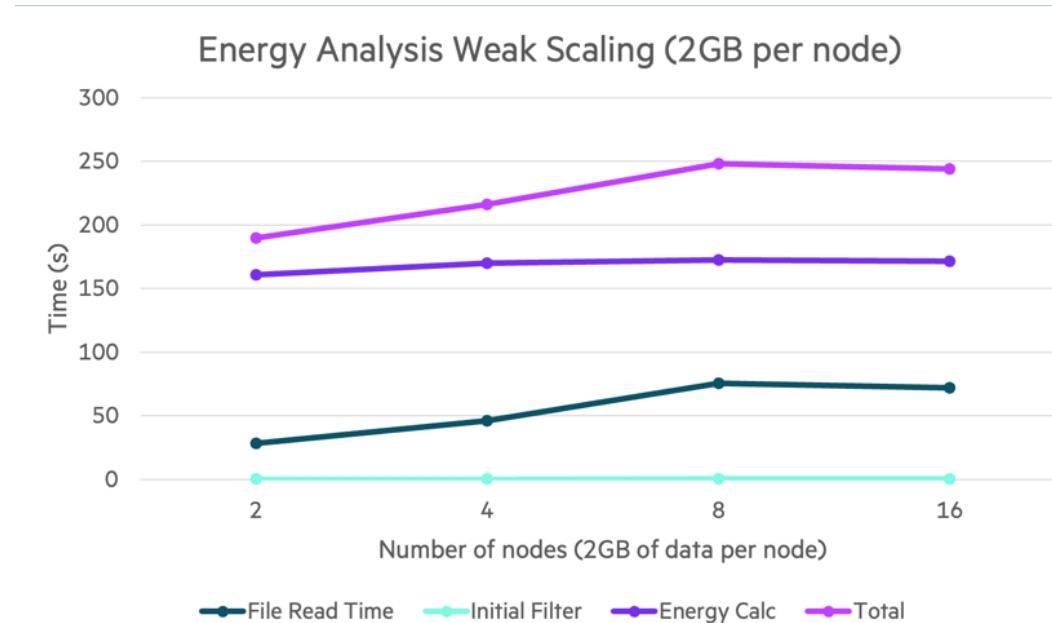— 576 compute nodes
— 72 TiB of 8-byte values
— ~480 GiB/s (~150 seconds)

### Arkouda Argsort Performance

Legend:
- Slingshot-11 May 2023,   32 GiB/node ◆——◆
- Slingshot-11 April 2023,   32 GiB/node ●——●
- HDR-100 IB May 2021, 128 GiB/node ✕——✕

GiB/s (y-axis): 0, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000

Nodes (x-axis): 1024, 2048, 4096, 8192

(better ↑)
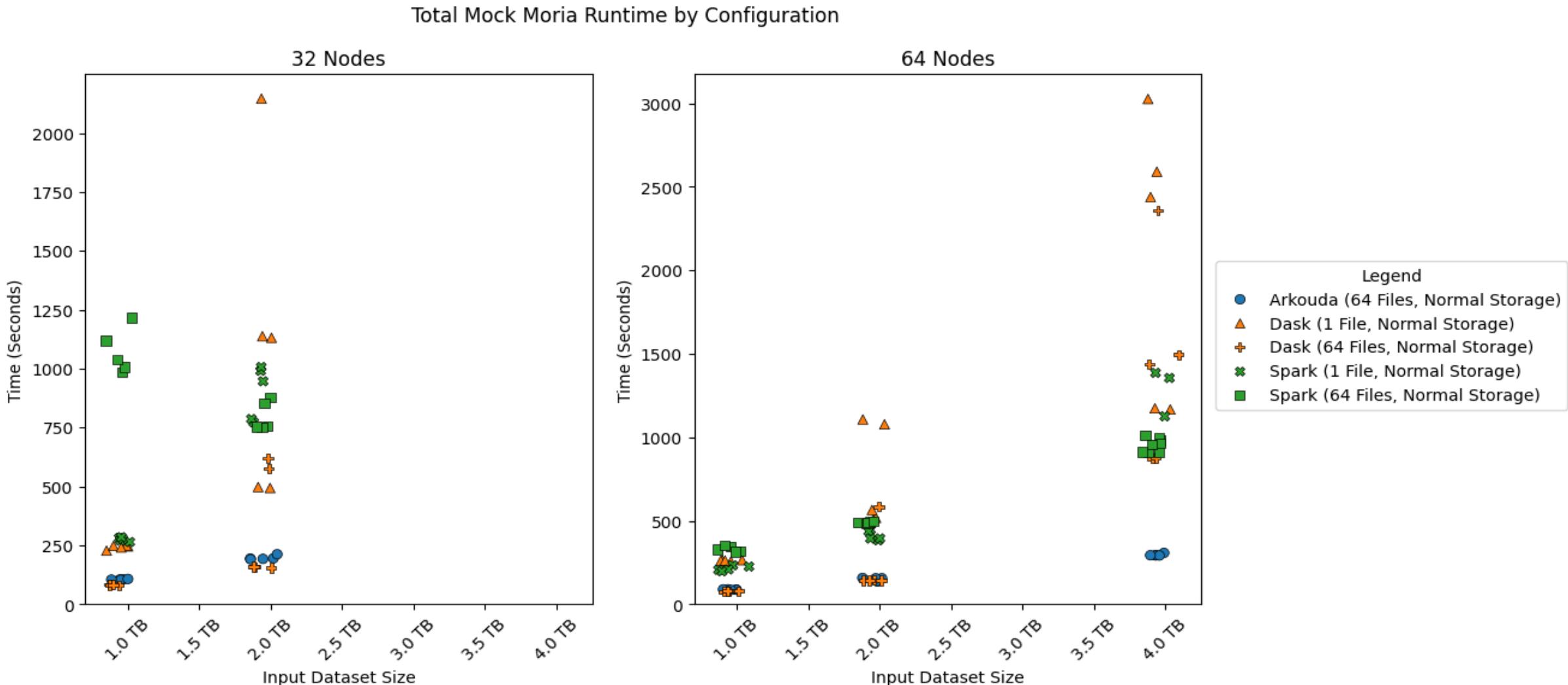
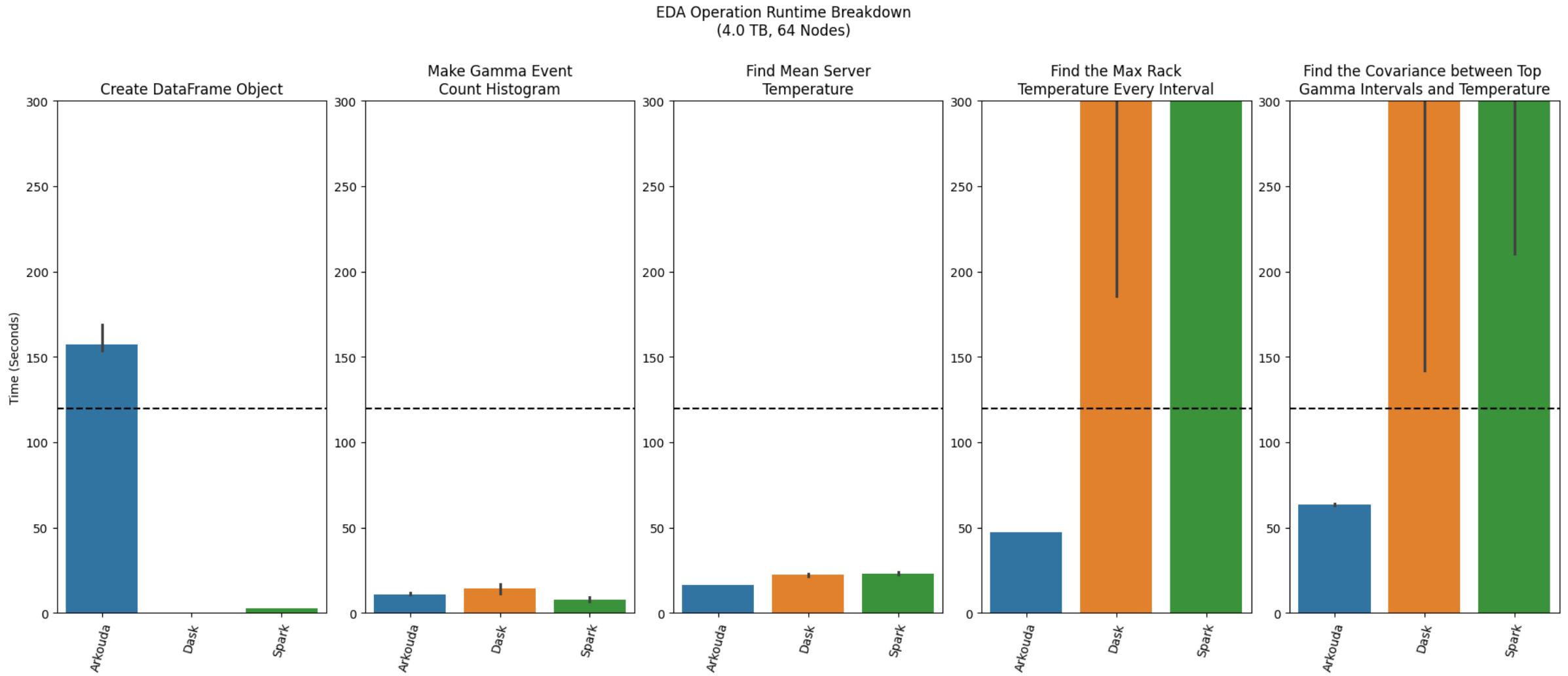## Implemented using ~100 lines of Chapel

# Arkouda/Pandas Comparison

— A collaboration with ORNL to analyze server telemetry data

  • Goal: to understand the impact of energy capping on application performance

— Translated ORNL Pandas script into Arkouda

  • Using the same data on a single node, Arkouda **outperformed Pandas by ~3.5x**
  • Moreover, the same script shows **promising weak scaling** enabling **much larger data** to be analyzed
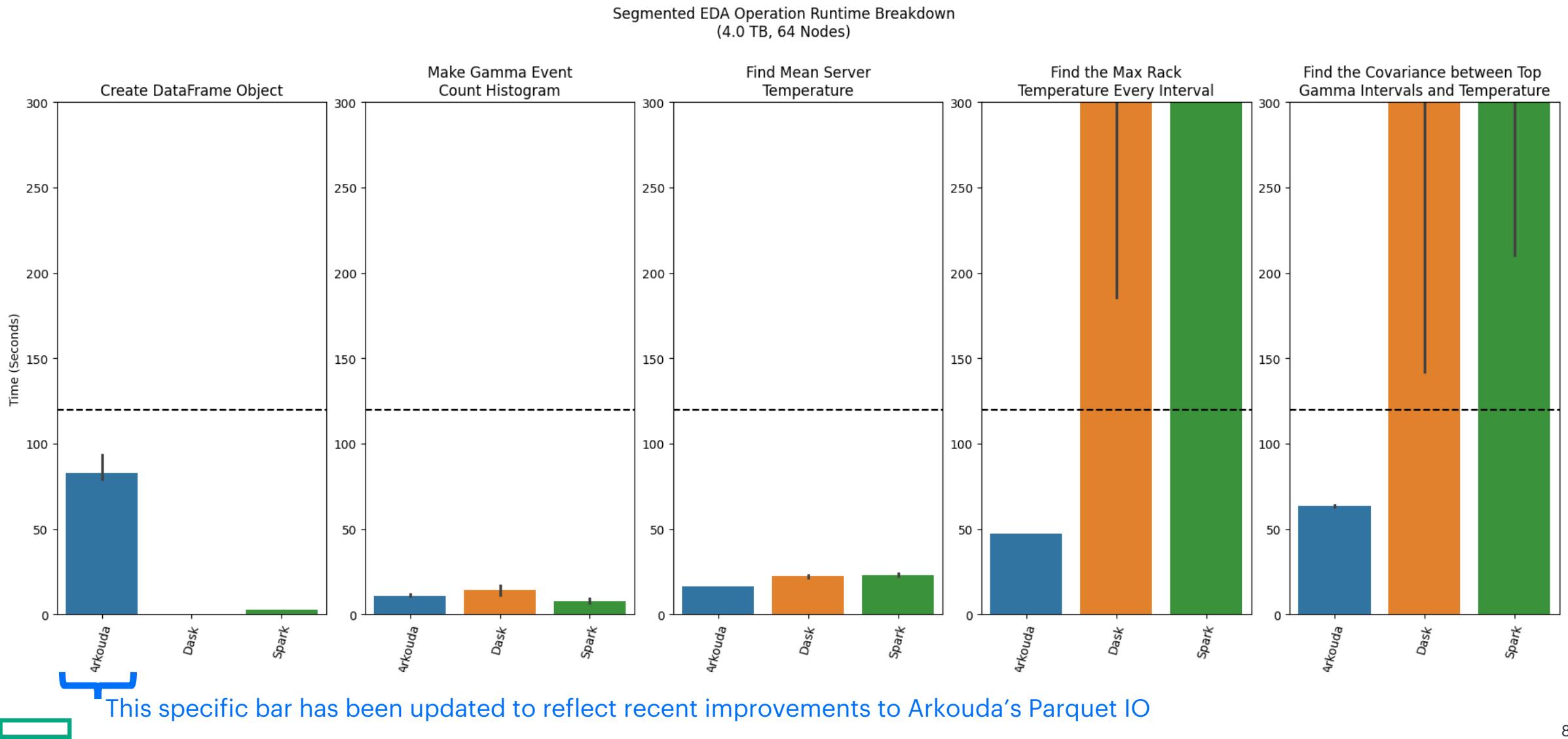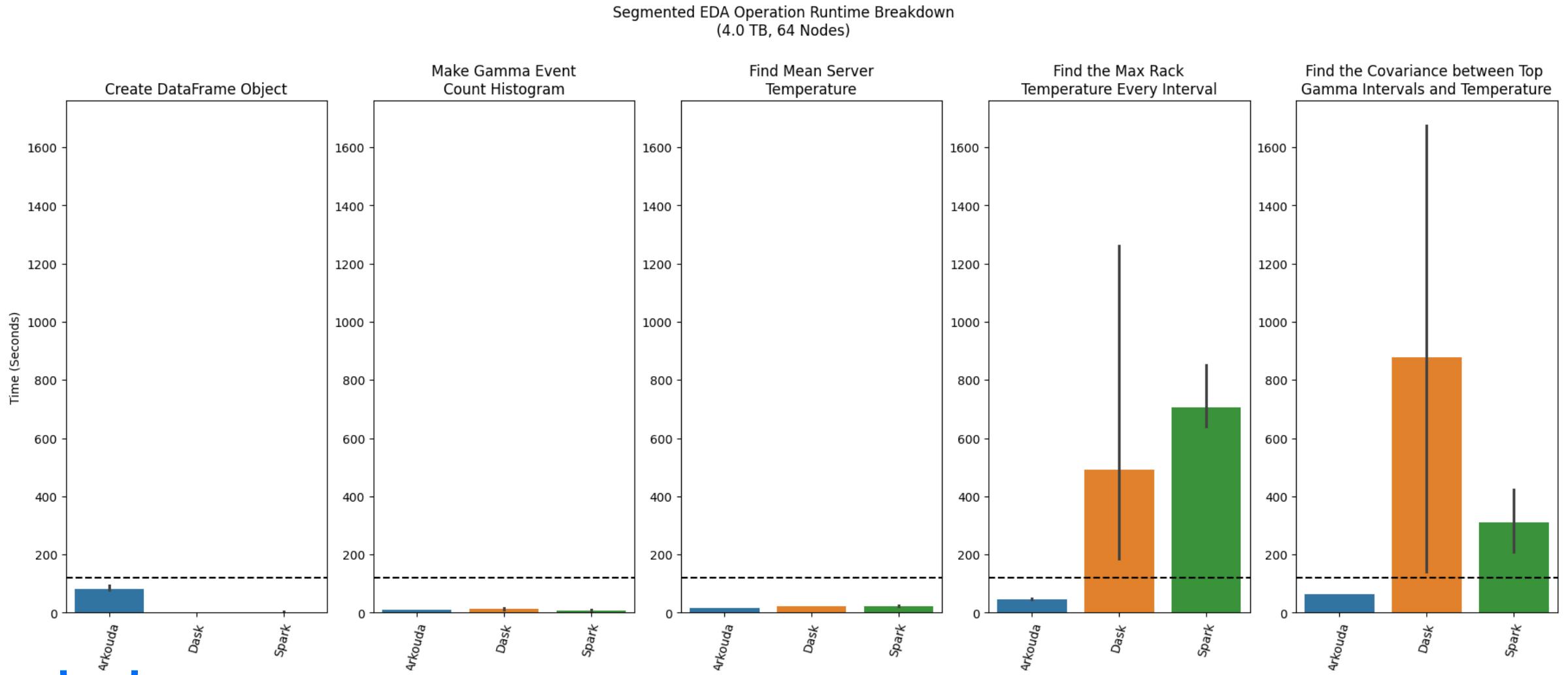


Energy Analysis Weak Scaling (2GB per node)

# Arkouda/Dask/Spark Comparison



Total Mock Moria Runtime by Configuration

# Arkouda/Dask/Spark Comparison: 64 nodes w/ 4 TB



EDA Operation Runtime Breakdown
(4.0 TB, 64 Nodes)

# Arkouda/Dask/Spark Comparison: w/ Parquet Improvements



Segmented EDA Operation Runtime Breakdown
(4.0 TB, 64 Nodes)

This specific bar has been updated to reflect recent improvements to Arkouda's Parquet IO

# Arkouda/Dask/Spark Comparison: Zoomed out



Segmented EDA Operation Runtime Breakdown
(4.0 TB, 64 Nodes)

This specific bar has been updated to reflect recent improvements to Arkouda's Parquet IO

9

# For More Information on Arkouda

## Arkouda website:



## Interview with founding co-developer, Bill Reus:



> "I was on the verge of resigning myself to learning MPI when I first encountered Chapel. After writing my first Chapel program, I knew I had found something much more appealing."
>
> ...
>
> "Chapel's separation of concerns immediately felt like the most natural way to think about large-scale computing. I would highly encourage anyone wanting to get into HPC programming to start with Chapel."

# What is Chapel?

**Chapel:** A modern parallel programming language

- Portable & scalable
- Open-source & collaborative
- An HPSF / Linux Foundation project



**Goals:**

- Support general parallel programming
- Make parallel programming at scale far more productive

# HPCC Stream Triad / RA: C+MPI+OpenMP vs. Chapel

**STREAM TRIAD: C + MPI + OPENMP**

```c
#include <hpcc.h>
#ifdef _OPENMP
#include <omp.h>
#endif

static int VectorSize;
static double *a, *b, *c;

int HPCC_StarStream(HPCC_Params *params) {
  int myRank, commSize;
  int rv, errCount;
  MPI_Comm comm = MPI_COMM_WORLD;

  MPI_Comm_size( comm, &commSize );
  MPI_Comm_rank( comm, &myRank );

  rv = HPCC_Stream( params, 0 == myRank);
  MPI_Reduce( &rv, &errCount, 1, MPI_INT, MPI_SUM, 0, comm );

  return errCount;
}

int HPCC_Stream(HPCC_Params *params, int doIO) {
  register int j;
  double  scalar;

  VectorSize = HPCC_LocalVectorSize( params, 3, sizeof(double), 0 );

  a = HPCC_XMALLOC( double, VectorSize );
  b = HPCC_XMALLOC( double, VectorSize );
  c = HPCC_XMALLOC( double, VectorSize );
```

```chapel
use BlockDist;

config const n = 1_000_000,
             alpha = 0.01;
const Dom = blockDist.createDomain({1..n});
var A, B, C: [Dom] real;


B = 2.0;
C = 1.0;


A = B + alpha * C;
```
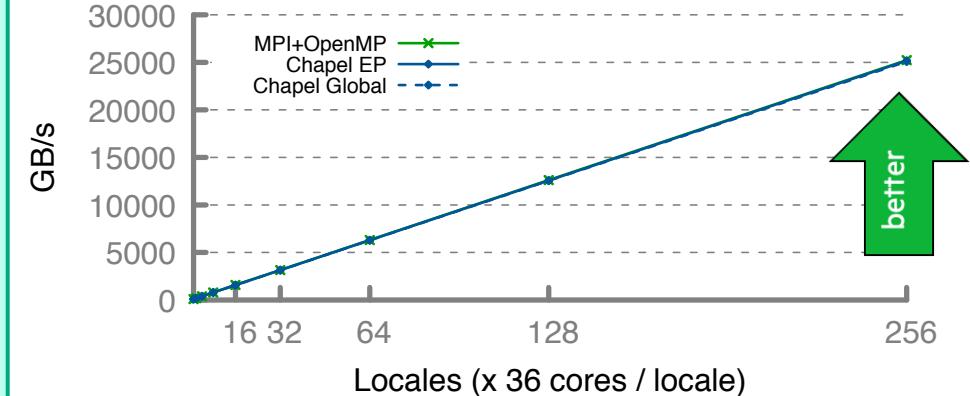
**HPCC RA: MPI KERNEL**

```chapel
...
forall (_, r) in zip(Updates, RandVals()) do
  T[r & indexMask].xor(r);
...
```

## STREAM Performance (GB/s)

- MPI+OpenMP
- Chapel EP
- Chapel Global

GB/s axis: 0, 5000, 10000, 15000, 20000, 25000, 30000

Locales (x 36 cores / locale): 16 32 64 128 256

better

## RA Performance (GUPS)

- Chapel
- MPI

GUPS axis: 0, 2, 4, 6, 8, 10, 12, 14

Locales (x 36 cores / locale): 16 32 64 128 256

better

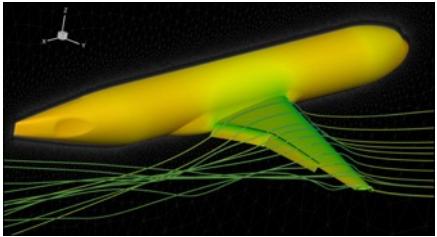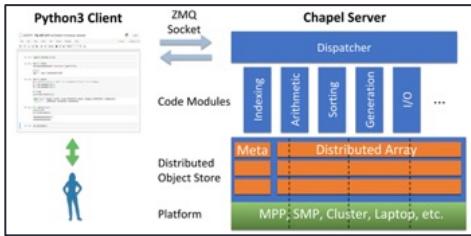# Why was Arkouda written in Chapel?

— **productivity**, readability, writability
  • Pythonic syntax is attractive to Python users who want to add features
— **parallelism** and **distributed arrays** as first-class features
— **performance:** competitive with conventional approaches
— **portability**: developed on laptop, deployed on supercomputer
— **interoperability:** can call to existing libraries
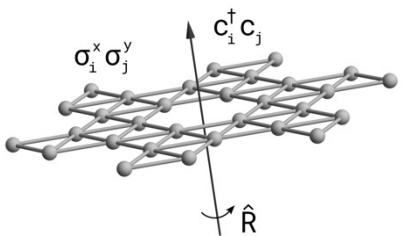
# Applications of Chapel



**CHAMPS: 3D Unstructured CFD**
Laurendeau, Bourgault-Côté, Parenteau, Plante, et al.
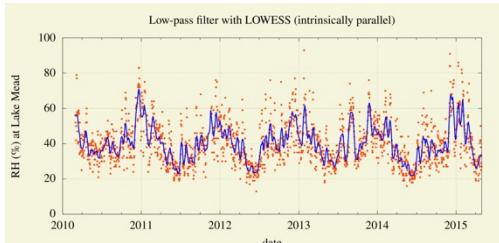*École Polytechnique Montréal*



**Arkouda: Interactive Data Science at Massive Scale**
Mike Merrill, Bill Reus, et al.
*U.S. DoD*



**ChOp: Chapel-based Optimization**
T. Carneiro, G. Helbecque, N. Melab, et al.
*INRIA, IMEC, et al.*



**ChplUltra: Simulating Ultralight Dark Matter**
Nikhil Padmanabhan, J. Luna Zagorac, et al.
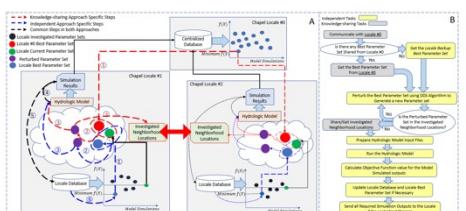*Yale University et al.*



**Lattice-Symmetries: a Quantum Many-Body Toolbox**
Tom Westerhout
*Radboud University*
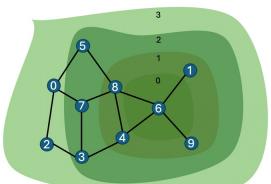


**Desk dot chpl: Utilities for Environmental Eng.**
Nelson Luis Dias
*The Federal University of Paraná, Brazil*



**RapidQ: Mapping Coral Biodiversity**
Rebecca Green, Helen Fox, Scott Bachman, et al.
*The Coral Reef Alliance*



**ChapQG: Layered Quasigeostrophic CFD**
Ian Grooms and Scott Bachman
*University of Colorado, Boulder et al.*



**Chapel-based Hydrological Model Calibration**
Marjan Asgari et al.
*University of Guelph*



**Arachne Graph Analytics**
Bader, Du, Rodriguez, et al.
*New Jersey Institute of Technology*



**Modeling Ocean Carbon Dioxide Removal**
Scott Bachman Brandon Neth, et al.
*[C]Worthy*



**CrayAI HyperParameter Optimization (HPO)**
Ben Albrecht et al.
*Cray Inc. / HPE*

[images provided by their respective teams and used with permission]

# "7 Questions with Chapel Users" Interviews

Read about users' Chapel experiences in the "*7 Questions with Chapel Users*" series on our blog



Chapel Language Blog

**7 Questions for Éric Laurendeau: Computing Aircraft Aerodynamics in Chapel**

Posted on September 17, 2024.

Tags:  Computational Fluid Dynamics   User Experiences   Interviews

By: Engin Kayraklioglu, Brad Chamberlain

**7 Questions for David Bader: Graph Analytics at Scale with Arkouda and Chapel**

Posted on November 6, 2024.

Tags:  User Experiences   Interviews   Graph Analytics   Arkouda

By: Engin Kayraklioglu, Brad Chamberlain

**7 Questions for Marjan Asgari: Optimizing Hydrological Models with Chapel**

Posted on September 15, 2025.

Tags:  User Experiences   Interviews   Earth Sciences

By: Engin Kayraklioglu, Brad Chamberlain

**7 Questions for Scott Bachman: Analyzing Coral Reefs with Chapel**

Posted on October 1, 2024.

Tags:  Earth Sciences   Image Analysis   GPU Programming   User Experiences   Interviews

By: Brad Chamberlain, Engin Kayraklioglu

**7 Questions for Bill Reus: Interactive Supercomputing with Chapel for Cybersecurity**

Posted on February 12, 2025.

Tags:  User Experiences   Interviews   Data Analysis   Arkouda

By: Engin Kayraklioglu, Brad Chamberlain

**7 Questions for Nelson Luís Dias: Atmospheric Turbulence in Chapel**

Posted on October 15, 2024.

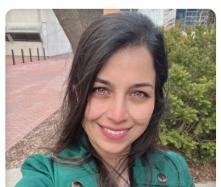Tags:  User Experiences   Interviews   Data Analysis   Earth Sciences   Computational Fluid Dynamics

By: Engin Kayraklioglu, Brad Chamberlain

**7 Questions for Tiago Carneiro and Guillaume Helbecque: Combinatorial Optimization in Chapel**

Posted on July 30, 2025.

Tags:  User Experiences   Interviews

By: Engin Kayraklioglu, Brad Chamberlain

# Ways to engage with the Chapel Community

## Synchronous Community Events

- Project Meetings, weekly
- Deep Dive / Demo Sessions, weekly timeslot
- ChapelCon (formerly CHIUW), annually

## Asynchrounous Communications

- Chapel Blog, typically ~2 articles per month
- Community Newsletter, quarterly
- Announcement Emails, around big events

## Social Media

**FOLLOW US**

- BlueSky
- Facebook
- LinkedIn
- Mastodon
- Reddit
- X (Twitter)
- YouTube

## Discussion Forums

**GET IN TOUCH**

- Discord
- Discourse
- Email
- GitHub Issues
- Gitter
- Stack Overflow

## Ways to Use Chapel

**GET STARTED**

- Attempt This Online
- Docker
- E4S
- GitHub Releases
- Homebrew
- Spack

(from the footer of chapel-lang.org)

# Next Steps: SUF Characterizations / Speed-Dating?

**Big Q:** With current capabilities, can Arkouda support Scientific User Facility (SUF) workloads?

- correct file formats?

- required operations?

- performance and scalability?

If not, what is lacking, and what would be required to address them?

# Next Steps: Research Questions and Challenges

**GPUs:**
- Would scientific data analysis (SDA) operations benefit from GPU acceleration?  Or are other things a bottleneck?
- Would such use cases require new features from Arkouda/Chapel?  (e.g., GPU-initiated communication?)

**IO subsystems and file formats:**
- What new IO systems or file formats might be beneficial, or do we have what we need?
- What changes to system-level software would be necessary to (better) leverage such IO capabilities?

**Custom Hardware Accelerators:**
- What role might exotic new chips play in the SDA space?
- Will these be generally programmable, or more like library operations in silicon?

**Extensibility:**
- How can Arkouda's extensibility be streamlined to add new capabilities for rapidly changing requirements?
- What is required to dynamically add new Arkouda modules to a running server?

**Community:**
- How can HPC break its cycle of failing to broadly adopt new productive software systems?
- How should innovative HPC software be fostered and sustained over time?

# Thank You

@ChapelLanguage