

Interactive Data Science at Massive Scales using Python and Arkouda

Brad Chamberlain w/ Jade Abraham, Advanced Programming Team, HPE

UW Data Science Seminar
December 2, 2025

A Bit About Me



A Bit About You?



Today's talk as a research question:

“As computer scientists and HPC* experts, what can we do to enable productive data science on massive data sets?”

*HPC = High-Performance Computing

Data Science In Python at scale?

Motivation: Imagine you work with...

- ...Python-based data scientists
- ...HPC-scale data science problems to solve
- ...access to HPC systems



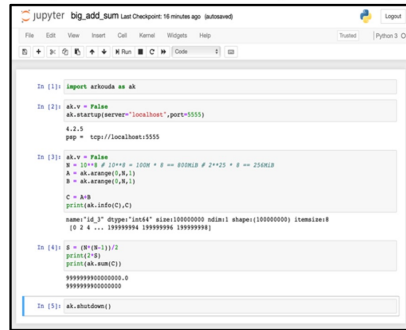
How will you leverage your Python programmers to get your work done?

What is Arkouda?

Q: “What is Arkouda?”



Arkouda Client
(written in Python)

A screenshot of a Jupyter Notebook interface. The title bar says 'jupyter big_add_sum Last Checkpoint: 10 minutes ago (auto-saved)'. The code in the notebook includes:

```
In [1]: import arkouda as ak
In [2]: ak.N = False
      ak.startUpServer("localhost", port=5555)
      4.2.5
      prep = http://localhost:5555
In [3]: ak.N = False
      N = 10**8 # 2**16 = 65536 * 2**12 = 8192 * 8192 * 2**12 = 256K
      A = ak.arange(N, 1)
      B = ak.arange(N, 1)
      C = A*B
      print(ak.info(C, C))
      name: 'A_B' dtype: 'uint64' size: 100000000 ndim: 1 shape: (100000000,) itemsize: 8
      [0 2 4 ... 199999998 199999999]
In [4]: S = ak.zeros(10**7)
      print(S)
      print(ak.sum(C))
      999999900000000.0
      999999900000000.0
In [5]: ak.shutdown()
```



User writes Python code
making familiar NumPy/Pandas calls

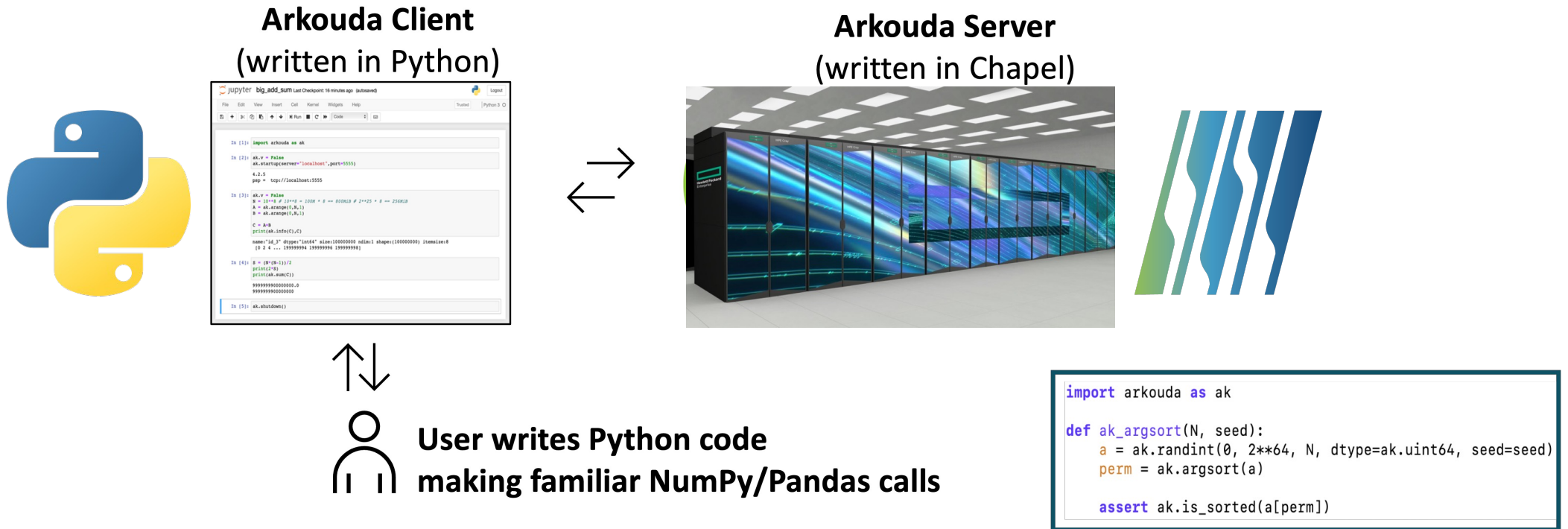
```
import arkouda as ak

def ak_argsort(N, seed):
    a = ak.randint(0, 2**64, N, dtype=ak.uint64, seed=seed)
    perm = ak.argsort(a)

    assert ak.is_sorted(a[perm])
```

What is Arkouda?

Q: “What is Arkouda?”



A1: “A scalable version of NumPy / Pandas routines for data scientists”

A2: “A framework for driving supercomputers interactively from Python”

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)

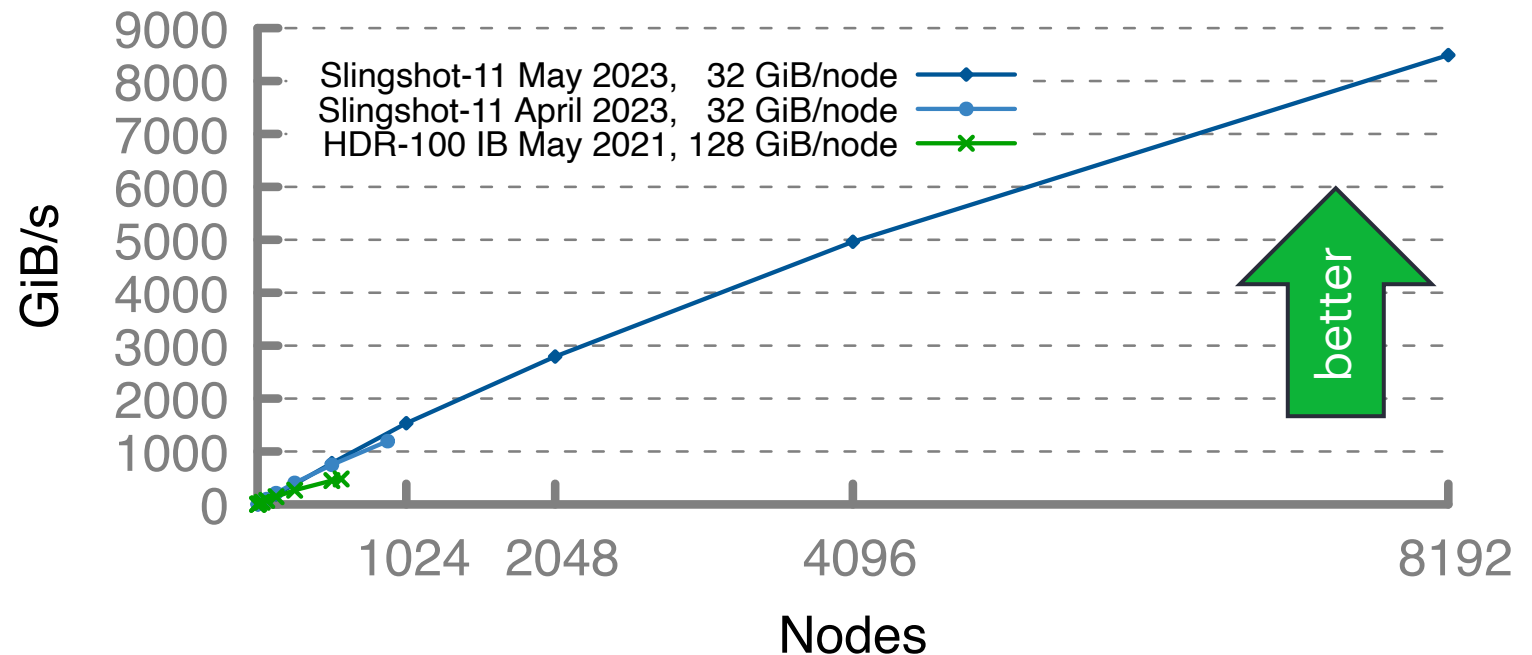
HPE Cray EX

- 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



Implemented using ~100 lines of Chapel code

Key Properties of Arkouda

- **Scalable:** has scaled to hundreds of TB, thousands of computes nodes, and over a million processor cores
- **Interactive:** operations are designed to complete in seconds to small numbers of minutes
- **Portable:** runs on virtually any system (laptop, cluster, cloud instance, supercomputer)
- **Open-Source:** developed on GitHub, released under the MIT license

The image shows a split-screen view. On the left is the GitHub repository page for 'Bears-R-Us / arkouda'. It displays the repository's structure with folders like .configs, .github, LICENSES, arkouda, benchmark_v2, benchmarks, converter, cpp-comparison, dep, docker, examples, pictures, and pydoc. Recent commits and pull requests are listed. On the right is the README file. It features a hand-drawn bear logo with the text 'arkouda' and 'massive scale data science'. Below the logo, it says 'Arkouda (αρκούδα) 🐻' and 'Interactive Data Analytics at Supercomputing Scale'. There are badges for CI, docs, license, and code style. Links for 'Online Documentation' and 'Nightly Arkouda Performance Charts' are provided at the bottom.

GitHub Repository: Bears-R-Us / arkouda

Recent Commits:

- ShreyasKhandekar: Fix CSV Parsing of Quoted Fields, Escaped Quotes, E... (821af11 - 2 hours ago, 4,468 Commits)
- Dec 1, 2025, 1:01 PM PST

Recent Pull Requests:

- #5060: remove isort (#5061) - last year
- #2981: add licenses (#2988) - last year
- #5076: ArkoudaExtensionArray.copy (#5077) - 3 hours ago
- #5060: remove isort (#5061) - last week
- #5060: remove isort (#5061) - last week
- #5060: remove isort (#5061) - last week
- #5060: remove isort (#5061) - last week
- #4209: create pyproject.toml (#4887) - 3 months ago
- #3358: Refactor MessageArgs (#3358) - last year
- #4893: Improve CI to automate builds (Part 3) (#4893) - 2 months ago
- #5060: remove isort (#5061) - last week
- #1143: more variations on the arkouda logo (#1143) - 3 years ago
- #5060: remove isort (#5061) - last week

README

Arkouda (αρκούδα) 🐻
massive scale data science

Interactive Data Analytics at Supercomputing Scale

CI passing docs passing license MIT code style black

Online Documentation

[Arkouda docs at Github Pages](#)

Nightly Arkouda Performance Charts

[Arkouda nightly performance charts](#)

Key Properties of Arkouda

- **Scalable:** has scaled to hundreds of TB, thousands of compute nodes, and over a million processor cores
- **Interactive:** operations are designed to complete in seconds to small numbers of minutes
- **Portable:** runs on virtually any system (laptop, cluster, cloud instance, supercomputer)
- **Open-Source:** developed on GitHub, released under the MIT license
- **Columnar:** represents dataframes using a distributed array per column
 - Current I/O Formats: Parquet, CSV, HDF5
- ...

Outline

Introduction to Arkouda

Performance/Scaling Comparisons

Live Demo

Extensibility

Wrap-Up

Arkouda Performance and Scaling Comparisons

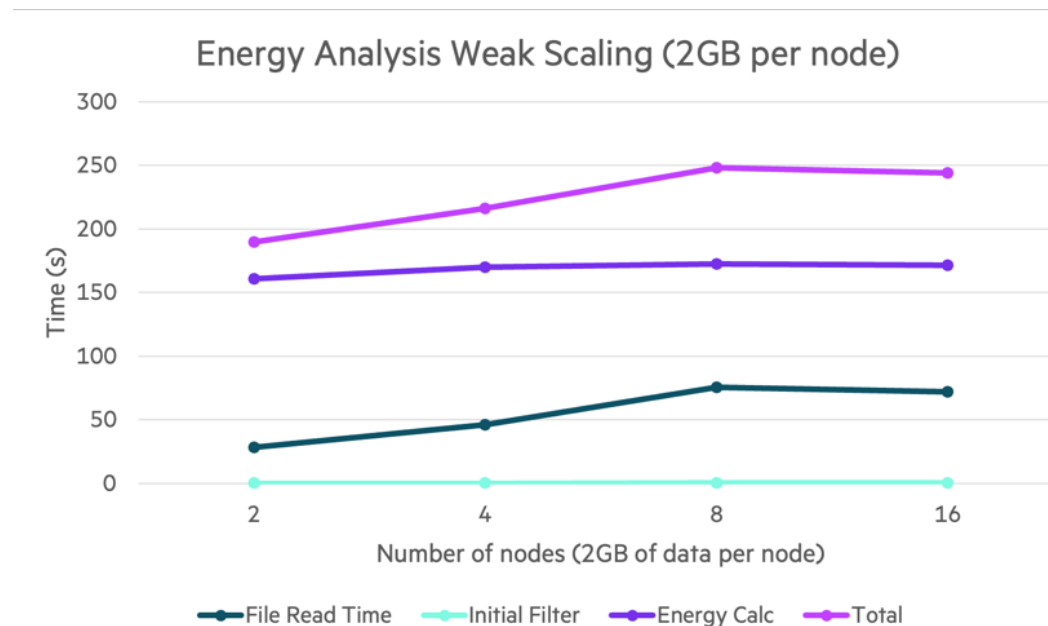
Arkouda/Pandas Comparison

Background: A collaboration with ORNL to analyze telemetry data from their Frontier supercomputer

- **Goal:** to understand the application performance impact of energy-capping GPUs

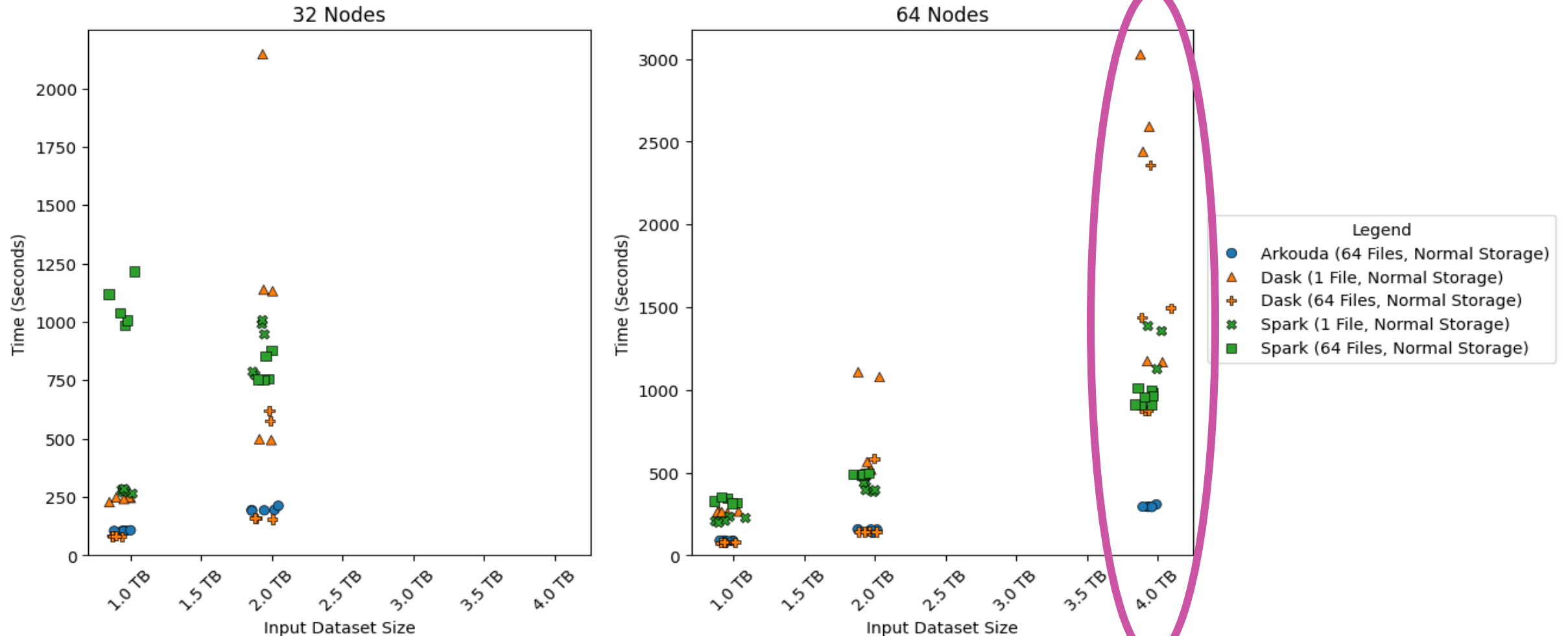
Experience: Translated ORNL Pandas script into Arkouda

- 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

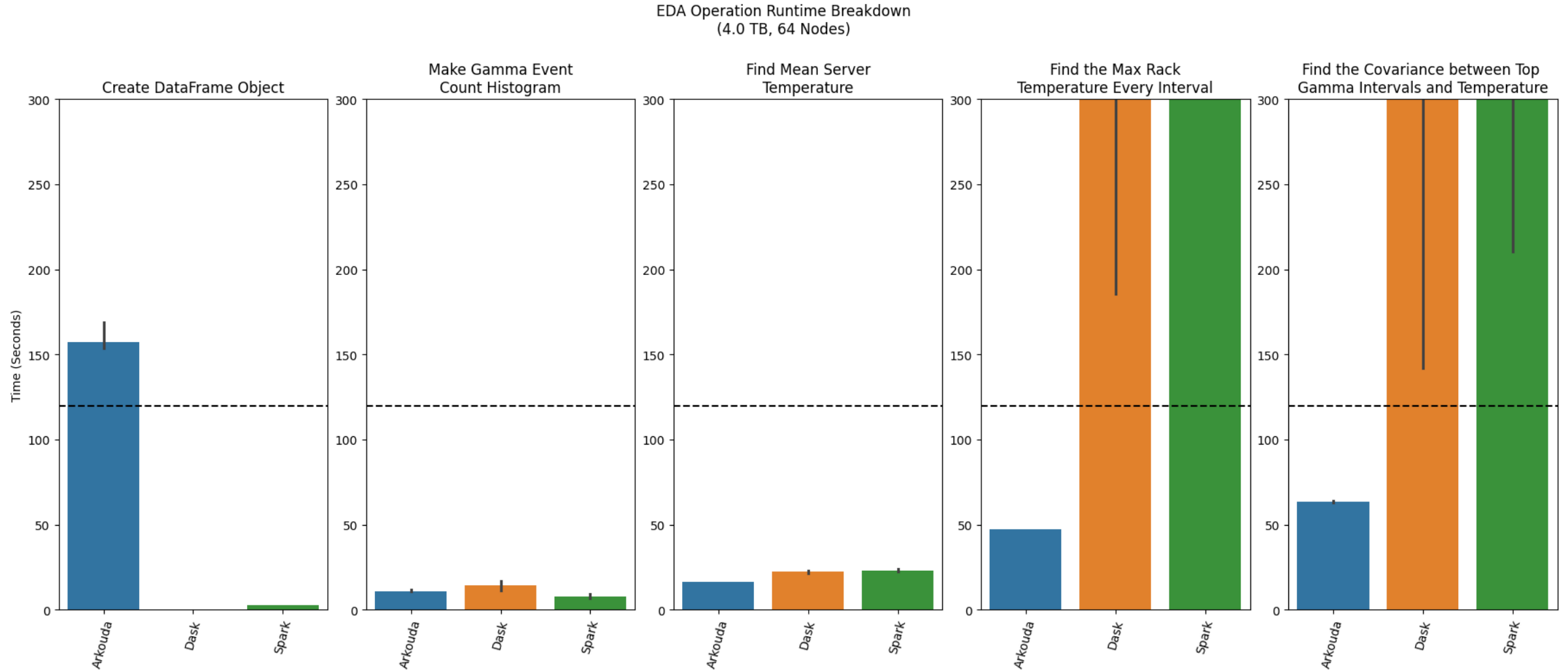


Arkouda/Dask/Spark Comparison

Total Mock Moria Runtime by Configuration

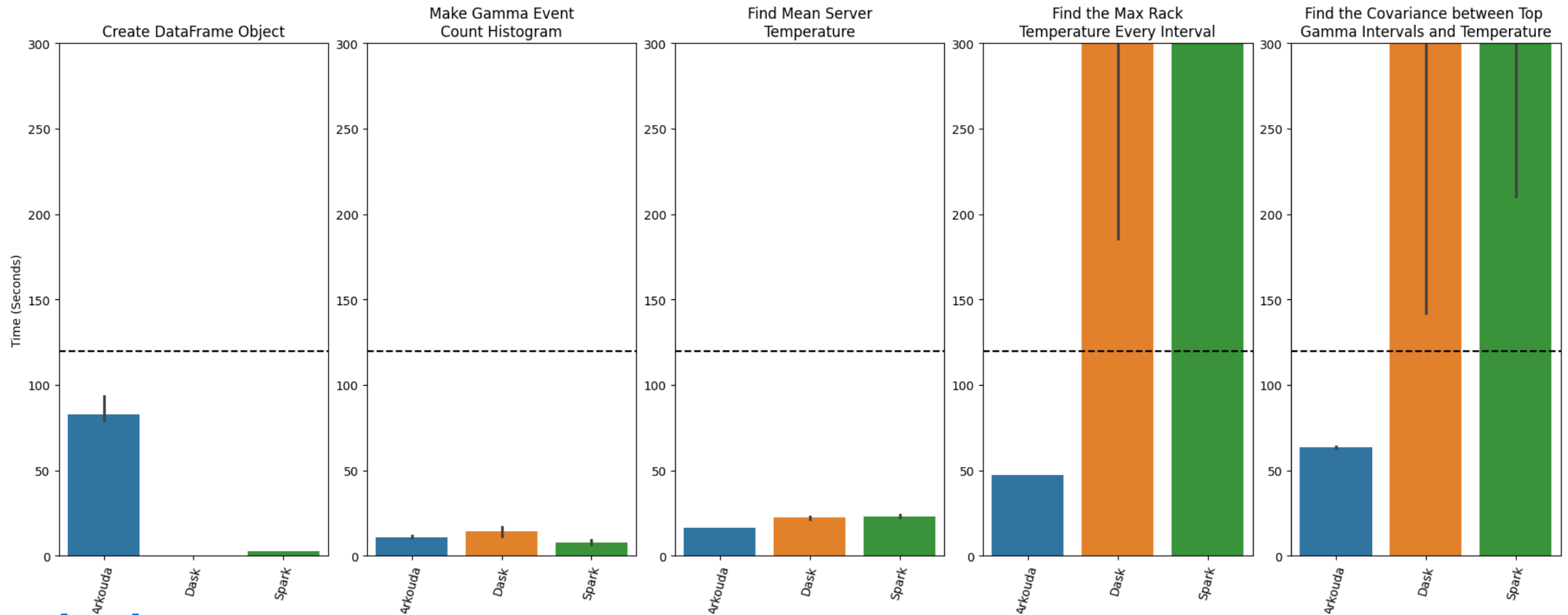


Arkouda/Dask/Spark Comparison: 64 nodes w/ 4 TB



Arkouda/Dask/Spark Comparison: w/ Parquet Improvements

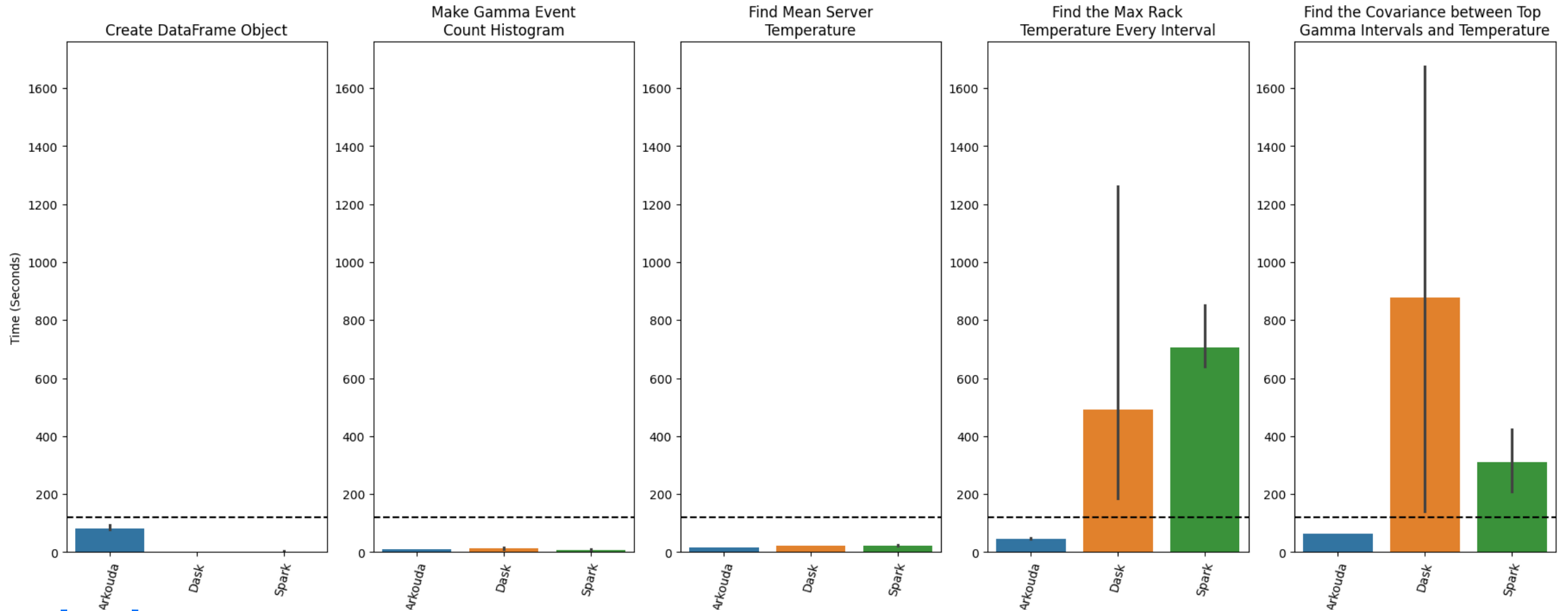
Segmented EDA Operation Runtime Breakdown
(4.0 TB, 64 Nodes)



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

Arkouda Demo

Jade Abraham



Extending Arkouda

Extending Arkouda

After the initial NumPy/Pandas features, Arkouda added support for extending its feature set:

- new operations in the client and/or server
- new data types

Using this capability, new features have been added for:

- Multidimensional arrays
- The Python Array API / XArrays
- Zarr I/O
- SciPy operations
- Sparse matrices and sparse matrix-matrix multiplication
- Graph analytics

Graph Analytics in Arkouda (by NJIT)

Motivation: Interactive, massive-scale graph analytics are useful for:

- ...finding communities in co-authorship & citation graphs
- ...categorizing NetFlow packets as malicious vs. benign
- ...finding subgraphs within neuroscience connectomes

Arachne: an Arkouda extension supporting graph computations

- performance competes with or beats leading approaches
- **representations:** vertex- and edge-centric property graphs, ...
- **algorithms:** BFS, triangle counting, connected components, ...
 - subgraph isomorphism via a novel parallel algorithm and visualizer

For more information:

- [On the Design of a Framework for Large-Scale Exploratory Graph Analytics](#), Oliver Andres Alvarado Rodriguez, NJIT Ph.D. dissertation, May 2025
- [HiPerMotif: Novel Parallel Subgraph Isomorphism in Large-Scale Property Graphs](#), Mohammad Dindoost et al., ChapelCon '25, October 2025

Motivation 2: We *Need* Large-Scale Graph Analytics

12 March 2025 Oliver Alvarado Rodriguez 9

NJIT New Jersey Institute of Technology

MOMO: Use Case

Our Collaboration with Harvard

Subgraphs	Arachne (s)	NetworkX (s)
	2.48	336.45
	3.62	173.75
	2.88	5,980.54
	339.46	16,436.85
	1.56	435.07
	78.77	810.23
	4.10	1,018.23
	38.06	>12,000

Dataset: 13,000 neurons with over 500,000 synaptic connections

using Arachne.HiPerMotif vs NetworkX VF2: Up to 2,000 X faster!

M. Shewarega, J. Trodi, et al. / MOMO

12

NJIT New Jersey Institute of Technology

Page 10

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



```

STREAM TRIAD: C + MPI + OPENMP

#include <hpc.h>
#ifdef OPENMP
#include <omp.h>
#endif

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

int HPCC_StartStream(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);

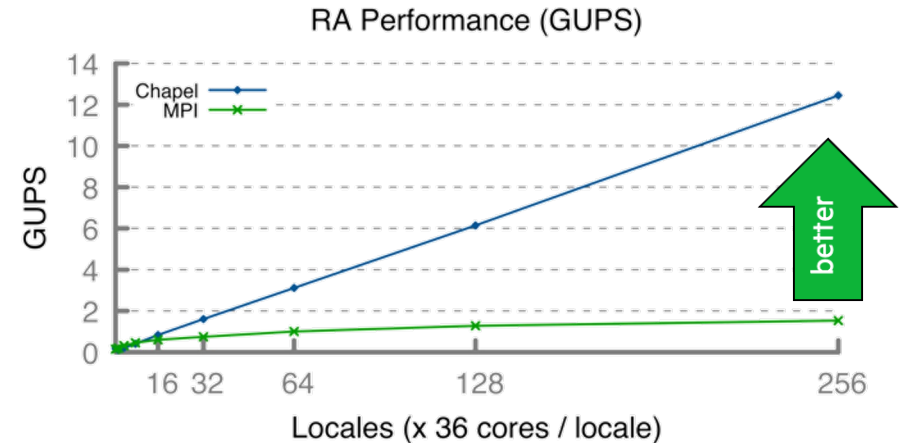
    return errCount;
}

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);

```

[illegible]

Bale IG in Chapel vs. SHMEM on HPE Cray EX (Slingshot-11)

Chapel

```
forall (d, i) in zip(Dst, Inds) do
    d = Src[i];
```

SHMEM (Exstack version)

```
i=0;
while( exstack_proceed(ex, (i==l_num_req)) ) {
    i0 = i;
    while(i < l_num_req) {
        l_indx = pckindx[i] >> 16;
        pe = pckindx[i] & 0xffff;
        if(!exstack_push(ex, &l_indx, pe))
            break;
        i++;
    }

    exstack_exchange(ex);

    while(exstack_pop(ex, &idx, &fromth)) {
        idx = ltable[idx];
        exstack_push(ex, &idx, fromth);
    }
    lgp_barrier();
    exstack_exchange(ex);

    for(j=i0; j<i; j++) {
        fromth = pckindx[j] & 0xffff;
        exstack_pop_thread(ex, &idx, (uint64_t)fromth);
        tgt[j] = idx;
    }
    lgp_barrier();
}
```

SHMEM (Conveyors version)

```
i = 0;
while (more = convey_advance(requests, (i == l_num_req)),
        more | convey_advance(replies, !more)) {

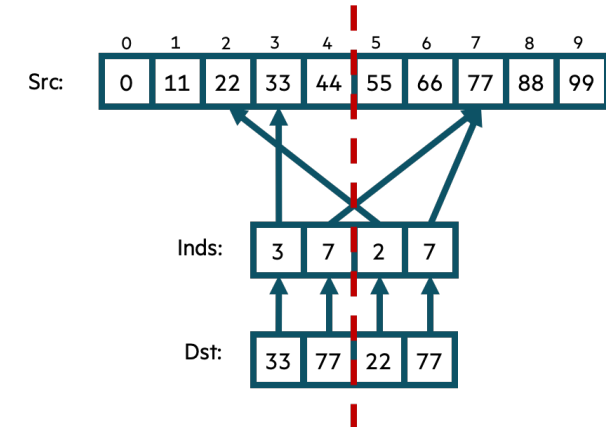
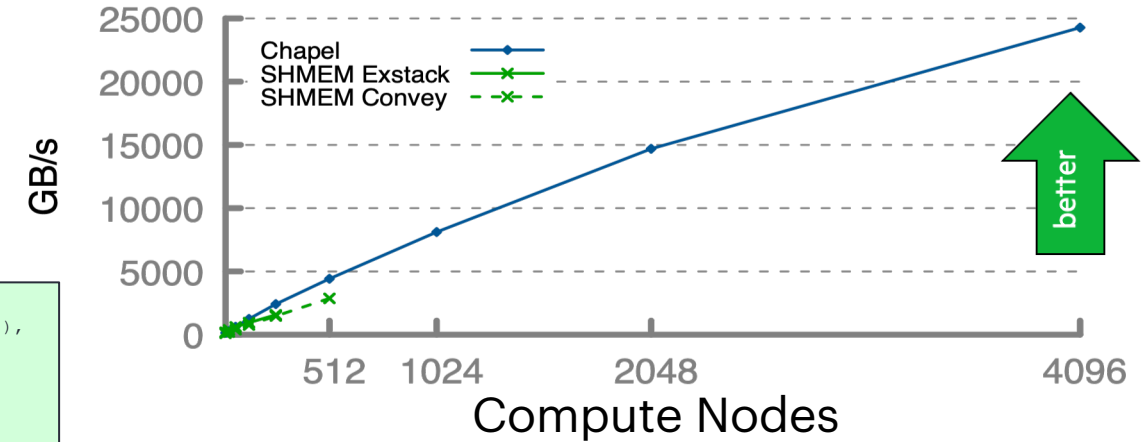
    for (; i < l_num_req; i++) {
        pkg.idx = i;
        pkg.val = pckindx[i] >> 16;
        pe = pckindx[i] & 0xffff;
        if (!convey_push(requests, &pkg, pe))
            break;
    }

    while (convey_pull(requests, ptr, &from) == convey_OK) {
        pkg.idx = ptr->idx;
        pkg.val = ltable[ptr->val];
        if (!convey_push(replies, &pkg, from)) {
            convey_unpull(requests);
            break;
        }
    }

    while (convey_pull(replies, ptr, NULL) == convey_OK)
        tgt[ptr->idx] = ptr->val;
}
```

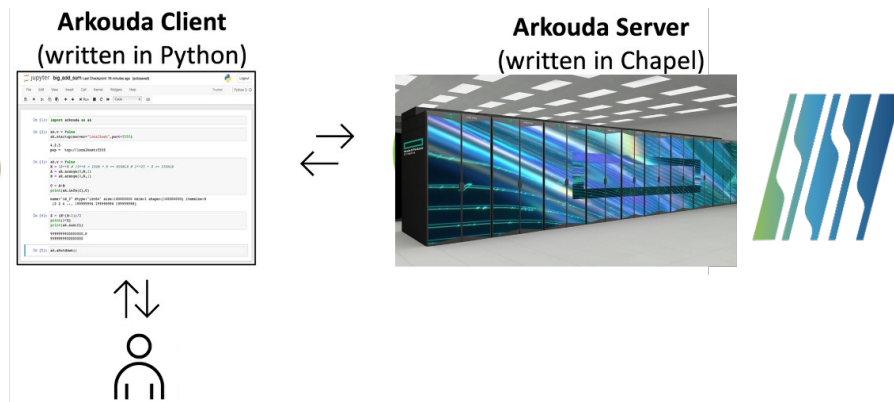
Bale Indexgather Performance

HPE Cray EX (Slingshot-11)



Why was Arkouda written in Chapel?

- **productivity**, readability, writability
 - Python-level 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 C, C++, Fortran libraries



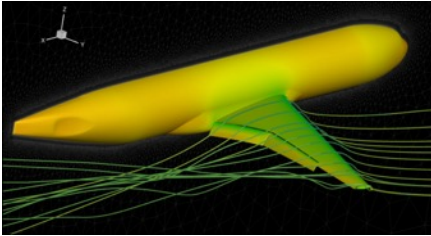
The screenshot shows a blog post from the Chapel Language Blog. The title is '7 Questions for Bill Reus: Interactive Supercomputing with Chapel for Cybersecurity'. The post is dated February 12, 2025, and is tagged with 'User Experiences', 'Interviews', 'Data Analysis', and 'Arkouda'. The author is listed as 'By: Engin Kayraklioglu, Brad Chamberlain'. A 'Table of Contents' is provided, listing four questions. The first question, '1. Who are you?', is partially visible, showing the beginning of Bill Reus's introduction.

"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."

...

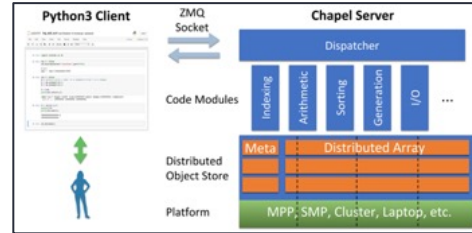
"Chapel's separation of concerns immediately felt like the most natural way to think about large-scale computing. I would highly encourage anyone wanting to get into HPC programming to start with Chapel."

Applications of Chapel



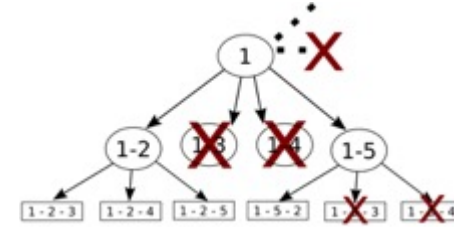
CHAMPS: 3D Unstructured CFD

Laurendeau, Bourgault-Côté, Parenteau, Plante, et al.
École Polytechnique Montréal



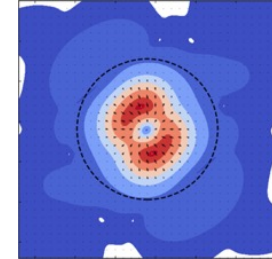
Arkouda: Interactive Data Science at Massive Scale

Mike Merrill, Bill Reus, et al.
U.S. DoD



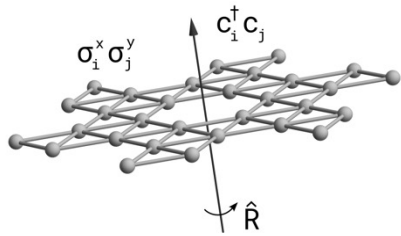
ChOp: Chapel-based Optimization

T. Carneiro, G. Helbecque, N. Melab, et al.
INRIA, IMEC, et al.



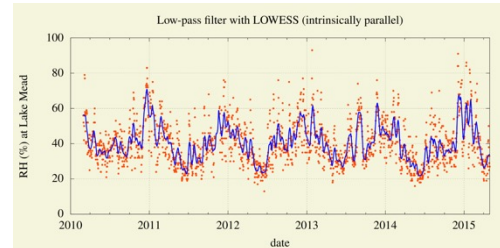
ChplUltra: Simulating Ultralight Dark Matter

Nikhil Padmanabhan, J. Luna Zagorac, et al.
Yale University et al.



Lattice-Symmetries: a Quantum Many-Body Toolbox

Tom Westerhout
Radboud University



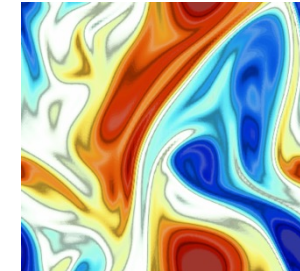
Desk dot chpl: Utilities for Environmental Eng.

Nelson Luis Dias
The Federal University of Paraná, Brazil



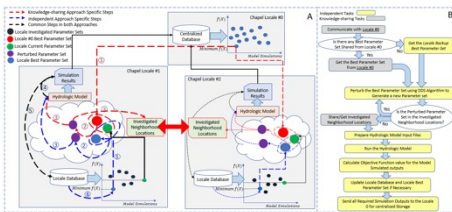
RapidQ: Mapping Coral Biodiversity

Rebecca Green, Helen Fox, Scott Bachman, et al.
The Coral Reef Alliance



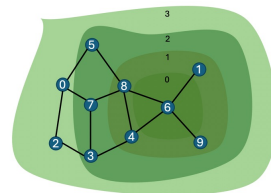
ChapQG: Layered Quasigeostrophic CFD

Ian Grooms and Scott Bachman
University of Colorado, Boulder et al.



Chapel-based Hydrological Model Calibration

Marjan Asgari et al.
University of Guelph



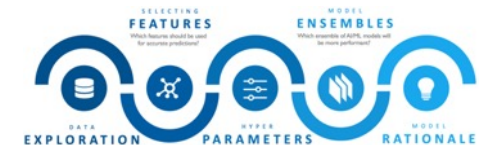
Arachne Graph Analytics

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 user experiences in the “[7 Questions with Chapel Users](#)” interview series on our blog

 Chapel Language Blog

AboutChapel WebsiteFeaturedSeriesTagsAuthorsAll Posts



7 Questions for Éric Laurendeau: Computing Aircraft Aerodynamics in Chapel

Posted on September 17, 2024.

Tags: Computational Fluid DynamicsUser ExperiencesInterviews

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 ExperiencesInterviewsGraph AnalyticsArkouda

By: [Engin Kayraklioglu](#), [Brad Chamberlain](#)




7 Questions for Marjan Asgari: Optimizing Hydrological Models with Chapel

Posted on September 15, 2025.

Tags: User ExperiencesInterviewsEarth Sciences

By: [Engin Kayraklioglu](#), [Brad Chamberlain](#)



7 Questions for Scott Bachman: Analyzing Coral Reefs with Chapel

Posted on October 1, 2024.

Tags: Earth SciencesImage AnalysisGPU Programming

User ExperiencesInterviews

By: [Brad Chamberlain](#), [Engin Kayraklioglu](#)




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

Posted on February 12, 2025.

Tags: User ExperiencesInterviewsData AnalysisArkouda

By: [Engin Kayraklioglu](#), [Brad Chamberlain](#)




7 Questions for Nelson Luís Dias: Atmospheric Turbulence in Chapel

Posted on October 15, 2024.

Tags: User ExperiencesInterviewsData Analysis

Earth SciencesComputational 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 ExperiencesInterviews

By: [Engin Kayraklioglu](#), [Brad Chamberlain](#)

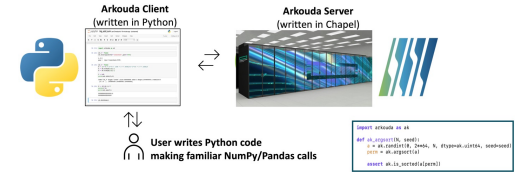


Wrap-up

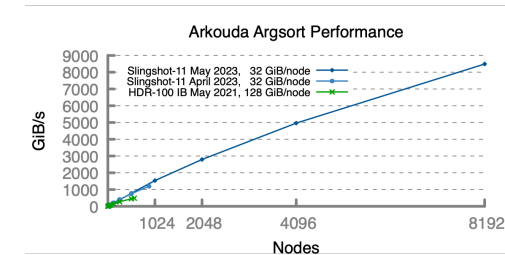


Summary

Arkouda is an open-source Python Library for driving HPC systems from Python

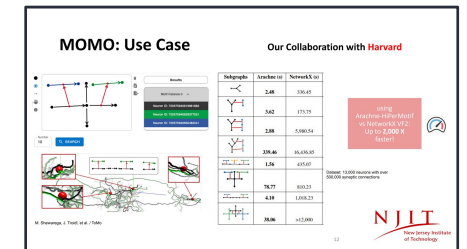


Arkouda operations have demonstrated scalability and interactivity



Arkouda's framework is extensible

- initial features focused on NumPy/Pandas APIs
- more recent work has focused on graph analytics, sparse linear algebra, SciPy, ...



Arkouda is written in Chapel for productivity, portability, and performance



If you or your colleagues would like an interactive intro to Arkouda or Chapel, let us know!




What's Next?

- Continue expanding Arkouda's supported operations
 - NumPy feature completeness in support of Pandas and SciPy
 - Sparse linear algebra benchmarking and optimization
 - ...
- GPUs: Chapel supports GPU computing but, to date, Arkouda operations haven't used them
 - **Big Q:** what motivating workflows and operations would benefit?
- Improve Arkouda's flexibility
 - E.g., add the ability to extend the Arkouda server's capabilities dynamically
- Project Honeycomb: Next-generation version of Arkouda
 - **Goals:** improved modularity and extensibility, language-agnostic, agentic-capable UI
- Grow Arkouda's open-source community of users and developers



For More Information on Arkouda

Arkouda website:

Arkouda

[github](#) [documentation](#) [gitter](#)

Massive-scale data science, from the comfort of your laptop

Arkouda

Ready for supercomputers

Numpy

Industry standard

```
# Launch an Arkouda server: ./arkouda_server -n1 <number-of-locates>

import arkouda as ak

# connect to the server
ak.connect('localhost', 5555)

# Generate two large arrays
a = ak.random.randint(0,2**32,2**38) # ----> Won't fit on a single machine!
b = ak.random.randint(0,2**32,2**38) # 178 of random integers.

# add them
c = a + b

# Sort the array and print first 10 elements
c = ak.sort(c)
print(c[0:10])
```

Try it Out

Tutorial Video

Chat on Gitter

Arkouda v2024.12.06 released!

The new release includes a refactored server making it easier to add new features, more Sparse Matrix functionality, new pdarray manipulation functions, and bug fixes.

[Read the release notes -->](#)

Arkouda is...

Fast

Arkouda is powered by Chapel, a programming language built from the ground up to support parallelism and distributed computing. Make the most out of every core and every node in your system.

Interactive

By distributing your data across multiple nodes, Arkouda allows you to rapidly transform and wrangle datasets in real time that are simply intractable for a laptop or desktop.


Extensible

One can expand on Arkouda's capabilities, thus enabling arbitrary scalable computations to be performed from Python.

Powered by Chapel

Arkouda's backend is implemented in Chapel, an open-source parallel programming language. Chapel is unique among mainstream languages as it puts parallelism and locality in the forefront, while not sacrificing productivity or portability. Chapel enables Arkouda to perform well and scale on many different architectures, from multicore laptops to cloud systems to world's fastest supercomputers.

To learn more about Chapel, check out its blog, presentations, tutorials and demos, and the [How Can I Learn Chapel?](#) page.



Arkouda users are saying...


“ ...solving problems in a matter of seconds, as opposed to days...”

— Toss Hayes, Bytola


“ [I'm] working with more data than I ever thought possible as a data scientist!”

— Jake Trookman, Erias

Interview with founding co-developer, Bill Reus:

Chapel Language Blog

[About](#) [Chapel Website](#) [Featured](#) [Series](#) [Tags](#) [Authors](#) [All Posts](#)



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](#)

Table of Contents

1. Who are you?
2. What do you do? What problems are you trying to solve?
3. How does Chapel help you with these problems?
4. What initially drew you to Chapel?
5. What are your biggest successes that Chapel has helped achieve?
6. If you could improve Chapel with a finger snap, what would you do?
7. Anything else you'd like people to know?

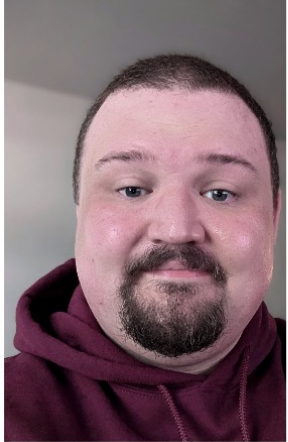
1. Who are you?

My name is Bill Reus, and I live near Annapolis, MD and the beautiful Chesapeake Bay. I am currently a data scientist doing statistical modeling and simulation for the United States government, but I began my career as an experimental chemist. In graduate school, I measured electron transport through thin films of organic molecules using an apparatus that our group invented to collect large volumes of noisy data quickly and with low cost. This approach contrasted with the typical means of studying molecular electronics, which was to spend weeks or months collecting a small number of exquisite measurements in ultra-high vacuum and at ultra-low temperature.

Recent Talk: [Arkouda Bulletin: A Year of Progress in Exploratory Data Analytics at Scale](#), Amanda Potts and Engin Kayraklioglu, ChapelCon '25, October 2025



The Advanced Programming Team at HPE



Ways to engage with the Chapel (and Arkouda) Communities

Synchronous Community Events

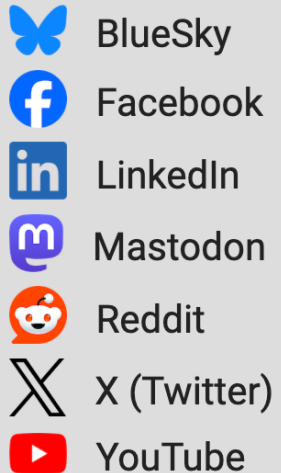
- [Project Meetings](#), weekly
- [Deep Dive / Demo Sessions](#), weekly timeslot
- [Chapel Teaching Meet-up](#), monthly
- [ChapelCon](#) (formerly CHI UW), annually

Asynchronous Communications

- [Chapel Blog](#), typically ~2 articles per month
- [Community Newsletter](#), quarterly
- [Announcement Emails](#), around big events

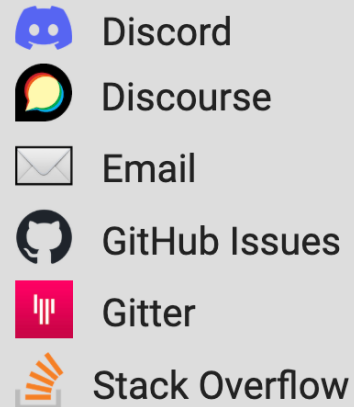
Social Media

FOLLOW US



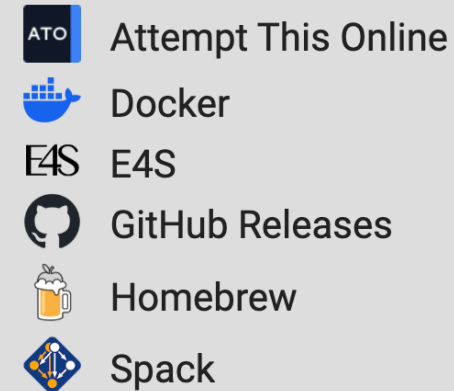
Discussion Forums

GET IN TOUCH



Ways to Use Chapel

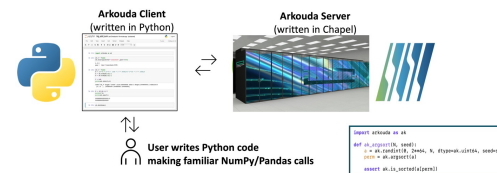
GET STARTED



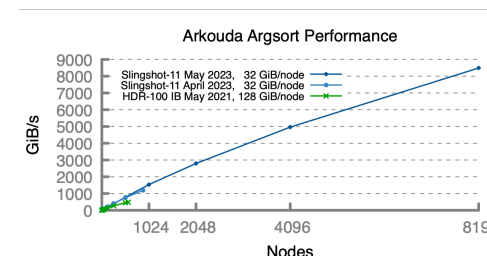
(from the footer of chapel-lang.org)

Summary

Arkouda is an open-source Python Library for driving HPC systems from Python

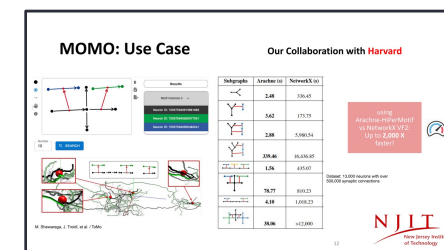


Arkouda operations have demonstrated scalability and interactivity



Arkouda's framework is extensible

- initial features focused on NumPy/Pandas APIs
- more recent work has focused on graph analytics, sparse linear algebra, SciPy, ...



Arkouda is written in Chapel for productivity, portability, and performance



If you or your colleagues would like an interactive intro to Arkouda or Chapel, let us know!



Thank You

@ChapelLanguage

