

Hewlett Packard Enterprise



Exploring Data at Scale with Arkouda: A Practical Introduction to Scalable Data Science

Ben McDonald, Software Engineer Supercomputing 2024

Challenge: Data scientists need to work interactively with massive-scale data sets

Introducing Arkouda: Python package designed for interactive, massive-scale data analysis

Demo

Recent advances: Rethinking Arkouda as a highly-extensible HPC framework

Data Science

Working with big data is much like exploring uncharted territory



Workflows

- Data scientists are taught to work with data interactively, learning with their data as they go
 - In-memory, interactive analysis is preferred by many and how data science is taught (NumPy/Pandas)
 - Batch jobs are useful when workflow already developed, but doesn't provide the same intimacy
- To address the difficulty of working with big data, a typical workflow could be:
 - 1. Cut down to a subset of the data to fit in memory (**downsampling**)
 - 2. Work interactively on a single node to learn about the subset data (NumPy/Pandas)
 - 3. Take what has been learned and send that off in a batch job on full data (Spark, etc.)

The Streetlight Effect

Faced with the unknown (downsampling), data scientists can suffer from the **"lighthouse effect"**







Workflows

Without Arkouda

- To address the difficulty of working with big data, a typical workflow could be:
 - 1. Cut down to a subset of the data to fit in memory (**downsampling**)
 - 2. Work interactively on a single node to learn about the data (NumPy, Pandas, etc.)
 - 3. Take what has been learned and send that off in a batch job (Spark, etc.)
- But what about the outliers? What about the streetlight effect?

With Arkouda

- 1. Work **interactively** with full data set on as many nodes as needed (Arkouda)
- 2. Pass subset of data to NumPy/Pandas to work with as usual (upside-downsampling)
- 3. Take what has been learned and send that off in a batch job (Spark, etc.)

Data Science Beyond the Laptop

Data sets today

Python must scale **beyond the laptop**, without sacrificing **interactivity**



L

Data Science Beyond the Laptop

Arkouda

Interactivity

Arkouda Client (written in Python)

B +	9< (3 🚯 🛧 🔸 M Run 🔳 C 🗰 Code 💠 🔤	
In	[1]:	import arkouda as ak	
In	[2]:	ak.v = False ak.startup(server="localhost".port=5555)	
		4.2.5 psp = tcp://localhost:5555	
In	[3]:	<pre>ak.v = False N = 10**8 ≠ 10**8 = 1000 * 8 == 8000tiB # 2**25 * 8 == 256xtiB A = ak.arange(0,R,1) B = ak.arange(0,R,1) C = A*B print(ak.info(0),C)</pre>	
		name:"id_3" dtype:"int64" size:100000000 ndim:1 shape:[100000000] itemsize:8 [0 2 4 199999994 199999996 199999998]	
In	[4]:	<pre>S = {N*(N-1))/2 print(2*5) print(ak.sun(C))</pre>	
		99999990000000 .0 999999900000000	
In	[5]:	ak.shutdown()	



Scalability

Arkouda Server

(written in Chapel)



Demo and Following along...

Installation

• To install Arkouda, see https://bears-r-us.github.io/arkouda/setup/install_menu.html

Docker Containers

• <u>https://github.com/Bears-R-Us/arkouda-contrib/tree/main/arkouda-docker</u>

Tutorial and Codespace

• <u>https://github.com/bmcdonald3/arkouda-codespace</u>

View the video that was embedded here at <u>https://www.youtube.com/watch?v=eSuyl9ogDfl</u>

Recent Development

- Interest in Arkouda outside of data science community has led to rearchitecting
- Arkouda has traditionally been thought of as "NumPy for HPC"
 - Rethinking Arkouda as a general framework for rapidly developing HPC-ready Python packages



- What previously required 97 lines in Arkouda can now be written as 7
 - Any Chapel function can be called from Python by adding a 'registerCommand' annotation

Questions?



© 2025 Hewlett Packard Enterprise Development LP