# Profiling and Visualizing I/O Access Patterns of HPC Applications

Abby Kapocius,  UNC Charlotte
Md Hasanur Rashid, College of Computing and Informatics

**UNIVERSITY OF NORTH CAROLINA CHARLOTTE**

## Introduction

I/O performance remains a critical point of investigation in high performance computing (HPC). The exponential increases in computing power in HPC systems have made I/O performance the bottleneck of HPC facilities. Moreover, the accelerated innovations in the hardware design, platform architecture, and software ecosystem of HPC facilities, as well as the diversity of HPC applications, keep the trends of I/O performance evolving continuously. As leading HPC facilities continue scaling, we must understand the existing trends of I/O performance to improve resource allocation strategies and alleviate possible I/O contentions observed in a parallel file system. The imbalance of I/O operations across shared storage in a parallel file system can severely impact the performance of scientific applications.

## Motivations

The dynamic nature of I/O behaviors, along with the necessity of understanding I/O performance, motivates us to investigate the current I/O access pattern trends with the help of recently developed advanced tools. Our work aims to characterize and profile I/O access behaviors of traditional HPC applications, I/O benchmarks, and machine learning applications. We will monitor, collect, and analyze I/O behaviors with state-of-the-art tools like Darshan. With the help of profiling information and tools like DXT-Explorer and Drishti, we will create visuals and ultimately seek out I/O patterns that lead us to strategies that will help improve the parallel file system performance. Our findings will contribute to a better understanding of I/O performance and achieve more efficient and accelerated application run times.
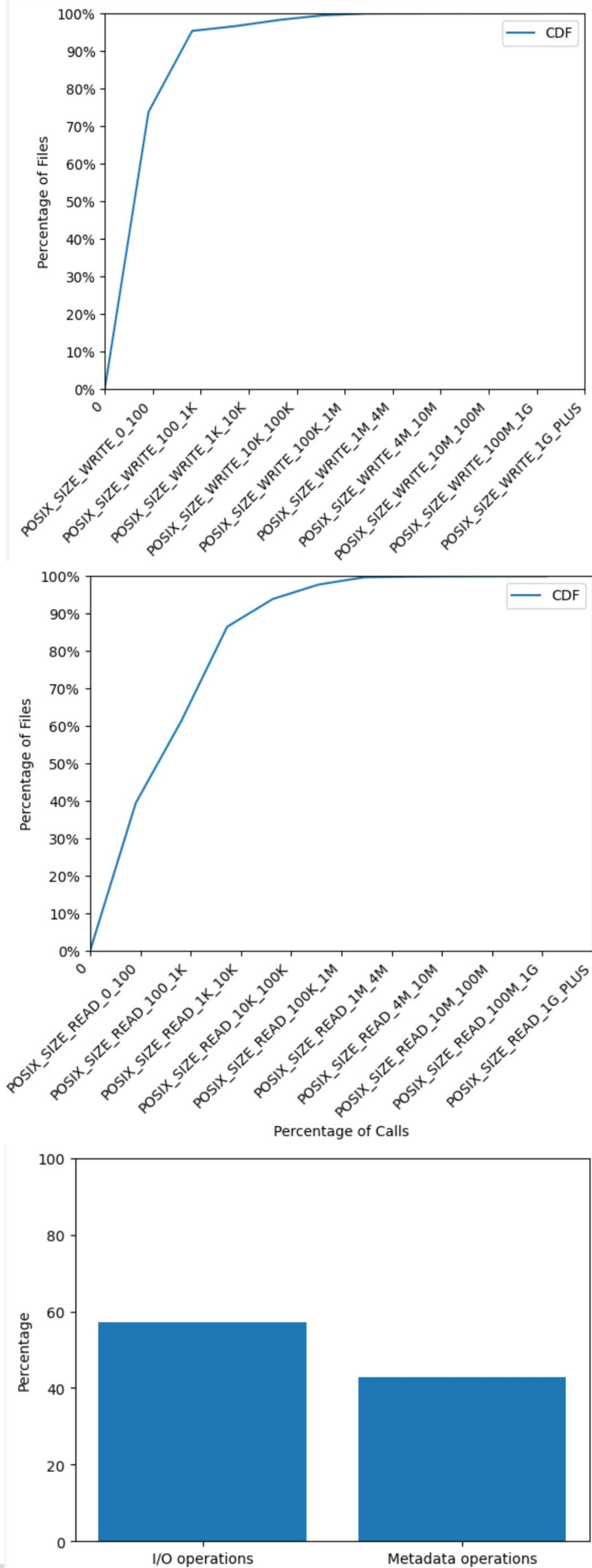
## Method

- Downloaded Darshan Logs from a published 2022 paper to retrieve the access patterns and performance behaviors of Cori Supercomputer in California
- Used Dask Dataframes to preprocess the abundance of data, as well as numpy and matplotlib to access the tools needed to make visualizations
- Concatenated all the csv files and removed all missing data
- Singled out specific columns in the csv files to perform calculations needed for the creation of various CDF and bar charts.
- Referred back to the original published paper to assure that the visualizations are an accurate and similar representation of data distribution
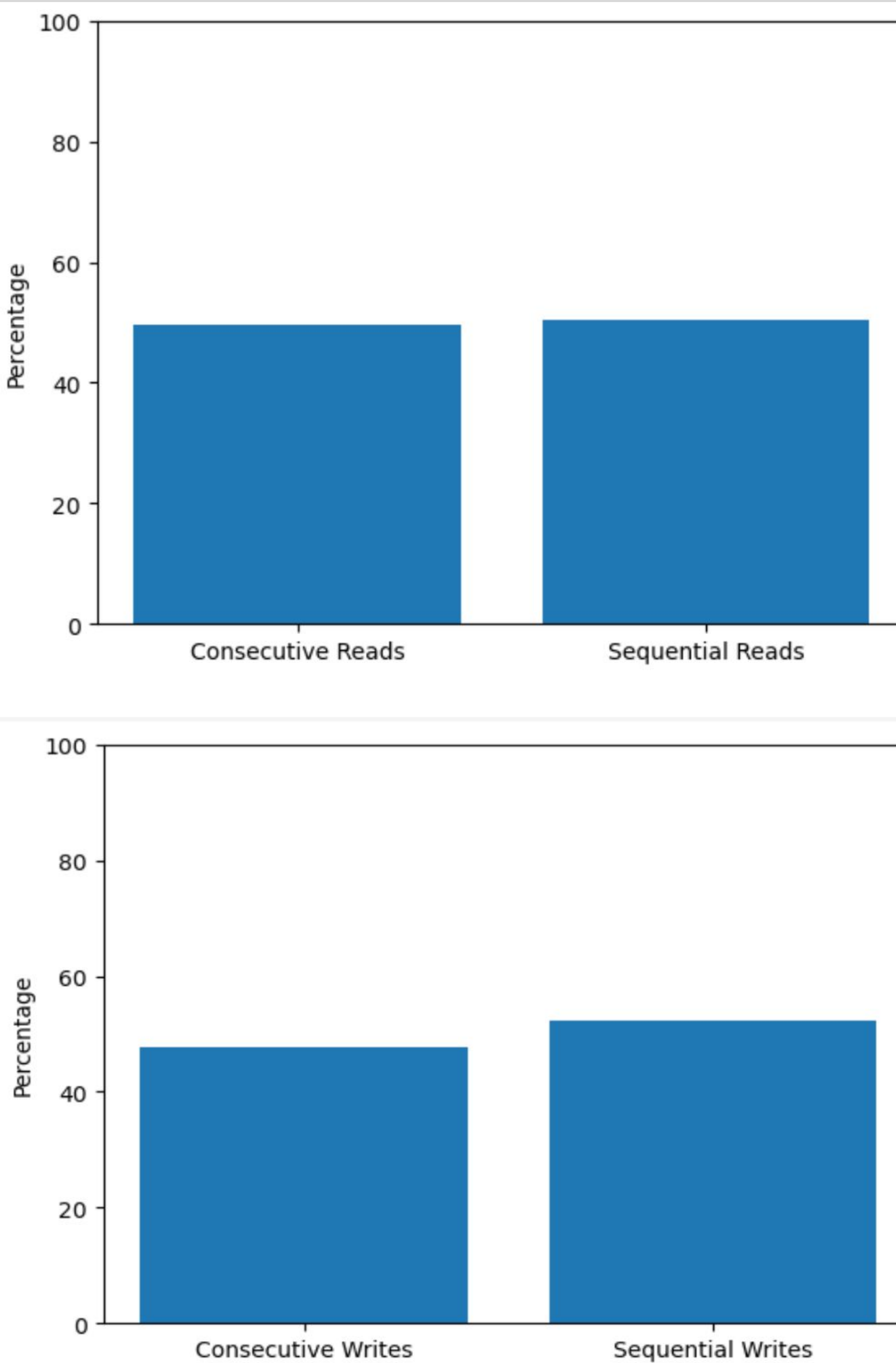
## Collected Data

- Data collected was primarily regarding I/O operations and meta operations for POSIX:
  - POSIX_SIZE_READ (10 columns)
  - POSIX_SIZE_WRITE (10 columns)
  - POSIX_CONSEC_READS
  - POSIX_CONSEC_WRITES
  - POXIS_SEQ_READS
  - POSIX_SEQ WRITES
  - POSIX_OPENS
  - POSIX_SEEKS
  - POSIX_STATS
- The following graphs represent:
  - Prevalence of different sized read and write operations
  - Consequential vs. sequential read and write operations
  - I/O operations vs. the amount of metadata operations

## Results







## Results





## Conclusions

Processing the data found in the Darshan logs of Cori Supercomputer allows us to look for consistencies in access patterns, which will allow us to continue working towards alleviating contention caused by I/O and metadata operations. Identifying patterns and visualizing data is the first step in understanding and recognizing problems that cause bottlenecks at HPC facilities, and with the pre-processed data at hand, we can lean on our peer's discoveries to work towards improving parallel file systems with the goal of increasing efficiency.