# Analyzing Lustre File System Performance With Splunk
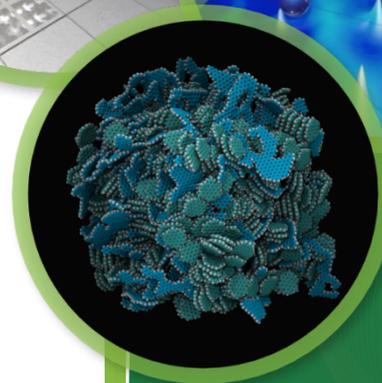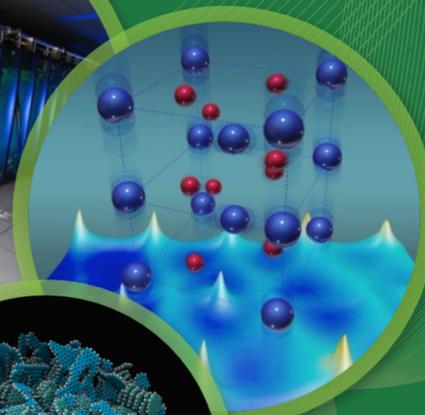
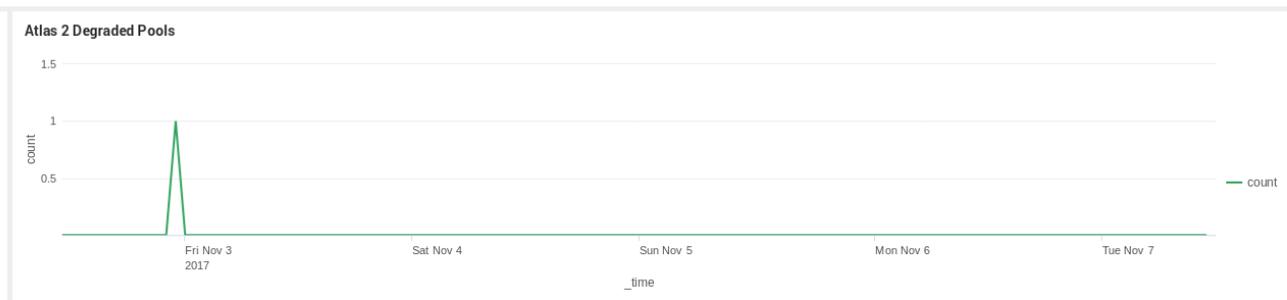OAK RIDGE | LEADERSHIP COMPUTING FACILITY
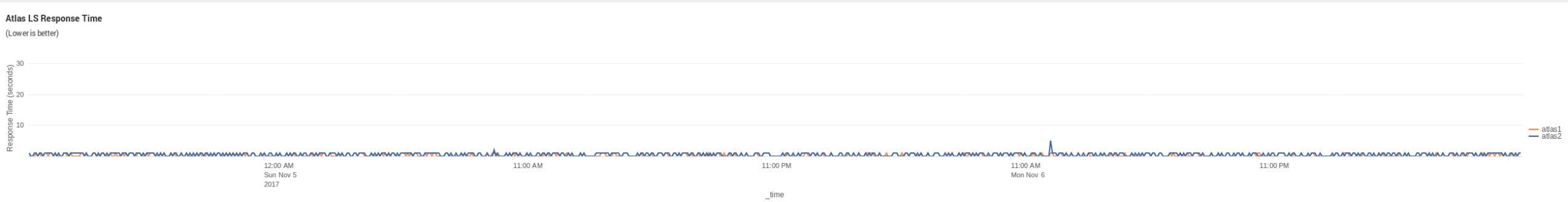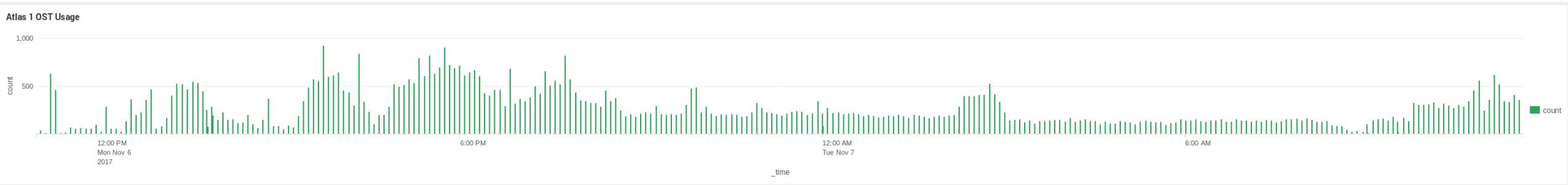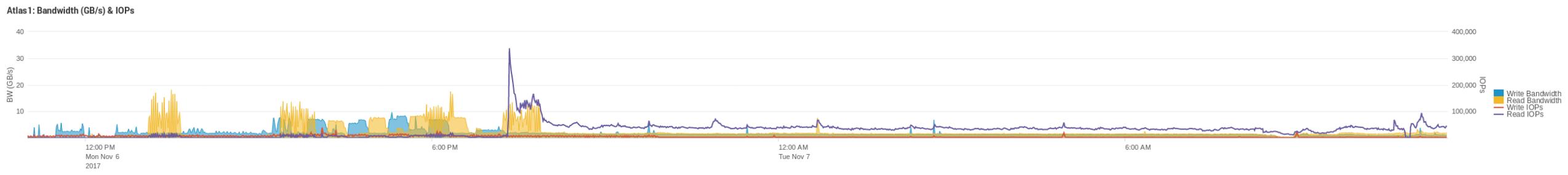National Laboratory

# Introduction

- OLCF is a Lustre shop (for now…)
  - Actually have several filesystems (2 production + various testbeds)

- Main filesystem is Atlas
  - 30PB, 20K disks, 288 OSS's, 72 DDN controllers
  - Actually split into two (Atlas1 & Atlas2) for metadata performance
  - Center-wide: filesystem is mounted on multiple compute resources

- Also have a few other filesystems (NOAA, testbeds)

- We've developed custom tools
  - We tried some other projects (like Robinhood), but they just couldn't handle the scale

**OAK RIDGE** National Laboratory | LEADERSHIP COMPUTING FACILITY

# Monitoring Tools – Capturing The Raw Data

- Block-level data from the DDN controllers
  - Read & write bandwidth and IOPs for each drive
  - Number of OST's (LUN's) in use inferred from bandwidth data
  - Gathered via a Python API

- Filesystem responsiveness tests
  - Run 'ls' from several different servers and record the times
  - Sounds simple (and it is), but it's also useful

- File Size Distributions
  - Scan the filesystem from the client side and record stats
  - Tool is distributed and scalable – and thus capable of overloading the MDS

OAK RIDGE
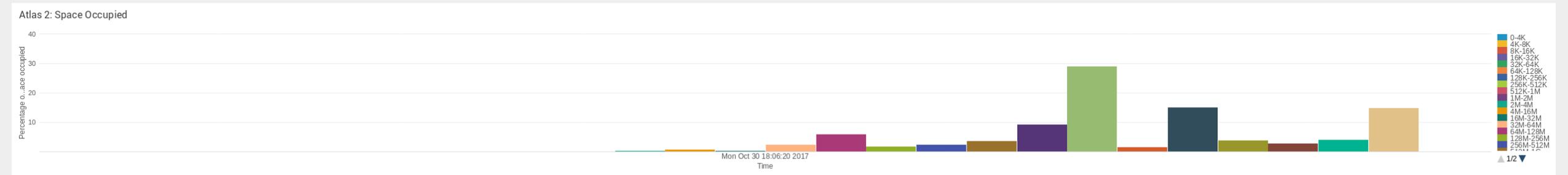National Laboratory
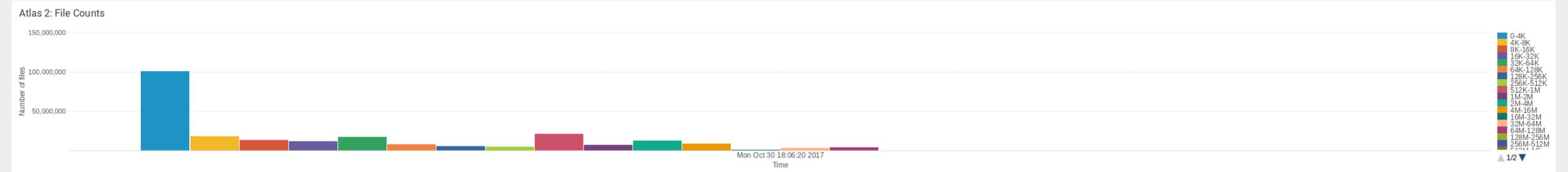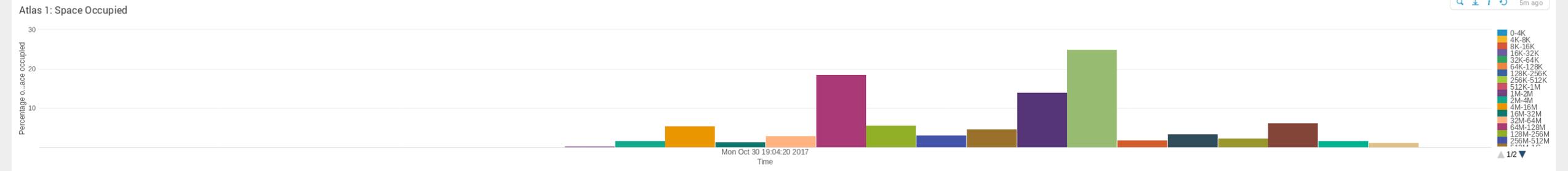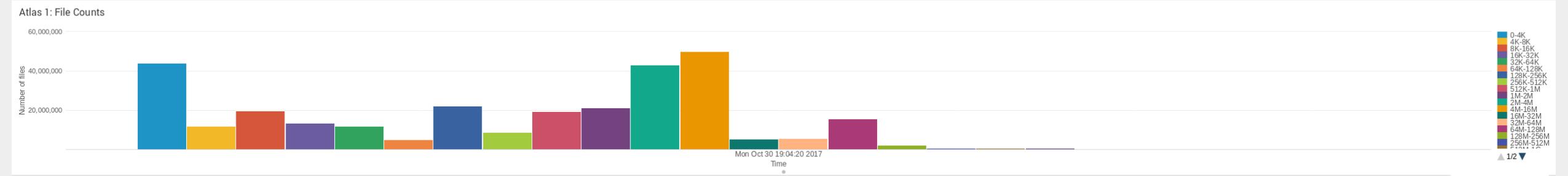LEADERSHIP COMPUTING FACILITY

# Splunk – Indexing & Searching The Data

- Splunk is an enterprise data-mining package

- Can ingest, index and store data from multiple sources
  - Very useful for tying all the different data sources together
  - Also provides mundane but necessary capabilities like user authorization

- Individual monitoring tools all feed their results into Splunk

- Allows us to do complex queries across multiple data sources
  - For filesystem monitoring, we actually don't need to

- Splunk license is based on ingest rate (GB / day)
  - This has implications for what data we collect and how frequently we sample

OAK RIDGE | LEADERSHIP COMPUTING FACILITY
National Laboratory

**Atlas1: Bandwidth (GB/s) & IOPs**



**Atlas 1 OST Usage**



**Atlas LS Response Time**

(Lower is better)



**Atlas 1 Degraded Pools**



**Atlas 2 Degraded Pools**

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

# Questions?

**OAK RIDGE** | LEADERSHIP
National Laboratory | COMPUTING
| FACILITY