# Lessons from Post-processing Climate Data on Modern Flash-based HPC Systems

Adnan Haider[1],  Sheri Mickelson(Advisor)[2], John Dennis(Advisor) [2], Xian-He Sun (Advisor) [1]
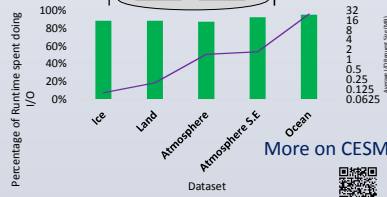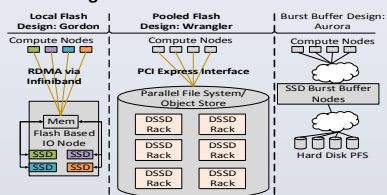
[1]Illinois Institute of Technology, USA; [2]National Center of Atmospheric Research, USA

## Flash-based Systems and Post-processing Software

Flash devices are a plausible solution to accelerate I/O bound applications. However, the tradeoffs associated with different flash architectures is unclear. We quantitatively assess two modern flash architectures using post-processing climate data applications to **facilitate correct matching between I/O workloads and flash storage architectures.**

Local Flash Design: Gordon | Pooled Flash Design: Wrangler | Burst Buffer Design: Aurora

Compute Nodes
RDMA via Infiniband
Mem
Flash Based IO Node
SSD SSD
SSD SSD

Compute Nodes
PCI Express Interface
Parallel File System/ Object Store
DSSD Rack DSSD Rack
DSSD Rack DSSD Rack
DSSD Rack DSSD Rack

Compute Nodes
SSD Burst Buffer Nodes
Hard Disk PFS



Percentage of Runtime spent doing I/O — Dataset: Ice, Land, Atmosphere, Atmosphere S.E, Ocean
Average I/O Request Size (MB): 32, 16, 8, 4, 2, 1, 0.5, 0.25, 0.125, 0.0625

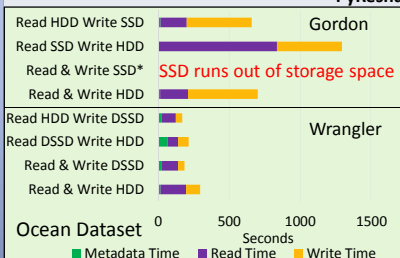**More on CESM**

- % I/O Time — Average I/O Request Size

- PyReshaper and PyAverager
- 90% of execution time is spent waiting for I/O to complete.
- Different datasets have vastly different I/O workloads (i.e. request size).
- IOR used for comparison with other workloads

## Gordon System Results: Local Flash Architecture

- Each compute node has access to a single solid state drive (SSD)
- Remote direct memory access via Infiniband.
- Can cause accesses to become queued
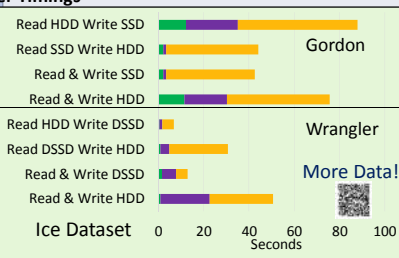1) Single SSD cannot handle rate of parallel accesses and interconnect causes latency.

### PyReshaper Timings

Gordon:
Read HDD Write SSD
Read SSD Write HDD
Read & Write SSD*  — **SSD runs out of storage space**
Read & Write HDD

Wrangler:
Read HDD Write DSSD
Read DSSD Write HDD
Read & Write DSSD
Read & Write HDD

**Ocean Dataset**
Seconds: 0, 500, 1000, 1500
- Metadata Time  - Read Time  - Write Time

### IOR Benchmark

2) Benefits of flash **decreases** at moderate scale and relatively small request sizes.



Throughput of SSD / Throughput of HDD
# of Processes / Amount of Data Written (GB)
I/O Request Size (K8): 64, 32, 16, 8, 4, 2
- 0-2  - 2-4  - 4-6  - 6-8  - 8-10  - 10-12

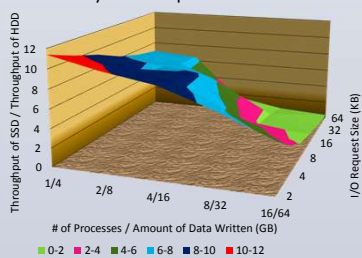## Wrangler System Results: Pooled Flash Architecture

- Uses DSSD devices which are faster than SSD.
- Each compute node has access to all DSSD devices (Pooled) via PCI Express
- Deploys parallel file system
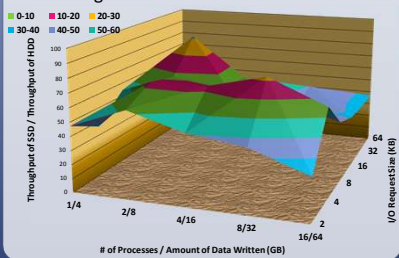1) Multiple DSSD and high throughput interconnect provide 2x to 6x improvements.

Gordon:
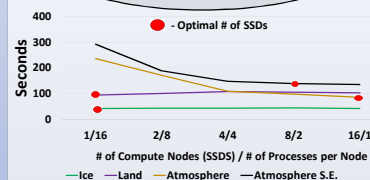Read HDD Write SSD
Read SSD Write HDD
Read & Write SSD
Read & Write HDD

Wrangler:
Read HDD Write DSSD
Read DSSD Write HDD
Read & Write DSSD
Read & Write HDD

**More Data!**

**Ice Dataset**
Seconds: 0, 20, 40, 60, 80, 100

2) **Consistent** benefits for all configurations when using flash.



Throughput of SSD / Throughput of HDD
# of Processes / Amount of Data Written (GB)
I/O RequestSize (K8): 64, 32, 16, 8, 4, 2
- 0-10  - 10-20  - 20-30  - 30-40  - 40-50  - 50-60

## Comparison of I/O Architectures



● - Optimal # of SSDs
Seconds: 400, 300, 200, 100, 0
# of Compute Nodes (SSDS) / # of Processes per Node: 1/16, 2/8, 4/4, 8/2, 16/1
- Ice  - Land  - Atmosphere  - Atmosphere S.E.
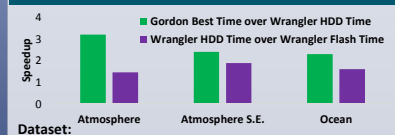
Multiple flash devices per compute node are needed to accommodate rate of parallel accesses issued by post-processing applications.



Speedup Provided by Flash: 6, 4, 2, 0
Dataset: Ice, Land, ATM, ATM S.E.
- Gordon  - Wrangler

A local architecture provides similar speedups as a pooled architecture if using multiple flash devices per compute node.



Speedup: 4, 3, 2, 1, 0
- Gordon Best Time over Wrangler HDD Time
- Wrangler HDD Time over Wrangler Flash Time
Dataset: Atmosphere, Atmosphere S.E., Ocean

Using a three-year newer system while not using flash (green bar) provides **more** speedup than using flash while keeping other hardware constant (purple bar)

## Lessons Learned

- An incorrect matching between storage architecture and I/O workload can hide the benefits of flash devices by increasing runtime by 2x.
- Hybrid I/O decreases flash storage consumption by half while decreasing runtime by 6x.   **Video Presentation→**
- Local flash could be a cheaper alternative to a pooled architecture if scalability and interconnect bottlenecks are alleviated.
- Three main criteria which determine performance on flash systems. 1) Number of flash devices in job. 2) Interconnect 3) Data availability of data stored on flash.
- Three years of more advanced hardware without flash devices provides more speedup than flash devices for some datasets, lessening the need for flash.

2012 Gordon — Flash devices on remote node - *local*
Flash on each compute node - *local* — Catalyst 2013
2015 Wrangler — All to All Connection - *pooled*
750 TB of flash and 750 GB/s bandwidth - *burst buffer* — NERSC Cori 2016

## Acknowledgements