# SSD Provisioning for Exascale Storage Systems

## Devesh Tiwari
## Oak Ridge National Laboratory

**Sarp Oral**

**Feiyi Wang**

**Saurabh Gupta**

**Josh Judd**

OAK RIDGE NATIONAL LABORATORY

MANAGED BY UT-BATTELLE FOR THE DEPARTMENT OF ENERGY

WARP MECHANICS®

# SSDs: The Good, the bad and the Ugly

High performance for random workloads

Low power consumption

Shock resistant

Write-endurance

High cost per Byte

# Can we build an exascale storage system out of SSDs?

# Write Endurance

# Is SSD write-endurance a problem for HPC?

Projected SSD storage space: ~5-10 PB

If building blocks are typical 256 GB SSDs

Number of SSDs in the system = 20,000

5 year warranty for max. 40GB write per day*

Allowed write amount: 600TB write per day

*Samsung 840 Pro Data Sheet
  http://www.samsung.com/us/pdf/memory-storage/840PRO_25_SATA_III_Spec.pdf

# Is SSD write-endurance a problem for HPC?

Assuming write amplification factor = 1.3

Allowed user written data = ~460TB per day

Write-endurance becomes a roadblock if an application dumps even 10% of system memory as checkpointing data every hour

System-level checkpointing easy on the programmer, hard on SSD-based storage system

# Is SSD write-endurance a problem for HPC?

Assuming write amplification factor = 1.3

Allowed user written data = ~460TB per day

OLCF: S3D 360 TB per day; GTC 240 TB per day

NERSC: <100 TB per day (Darshan instrumented)

ALCF:  most jobs moving <100 TB

Carns et al., Understanding and improving computational science storage access through continuous characterization, ACM Transactions on Storage, 2011

Computational Requirements of Leadership Computing
http://www.olcf.ornl.gov/wp-content/uploads/2010/03/ORNL_TM-2007_44.pdf

# Is SSD write-endurance a problem for HPC?

Even 1TB SSD as a building block will allow up to 150TB write per day
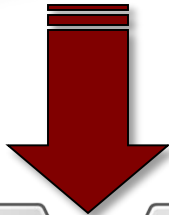
At higher price, up to 3 full writes per day

Write-endurance improving at a fast speed

Reducing the checkpointing size is the key to alleviating SSD write-endurance issue (application-level checkpointing strategies)

# **Where in the System? And, how much?**
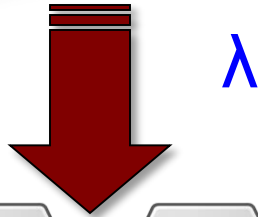
# SSDs as a burst buffer in HPC

λ = Data production rate

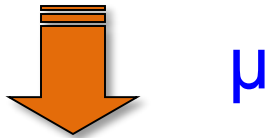SSD Burst Buffer

μ = Bandwidth to Permanent Storage System

Permanent Storage

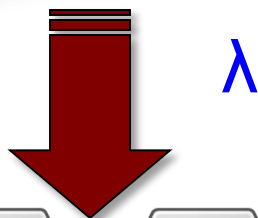# SSDs as a burst buffer in HPC



λ

SSD Burst Buffer

μ

Permanent Storage

If λ > μ :

SSD burst buffer capacity is unbounded

Intuition:
Higher incoming flux than outgoing flux will burst the pipe
Never enough time to drain out

# SSDs as a burst buffer in HPC
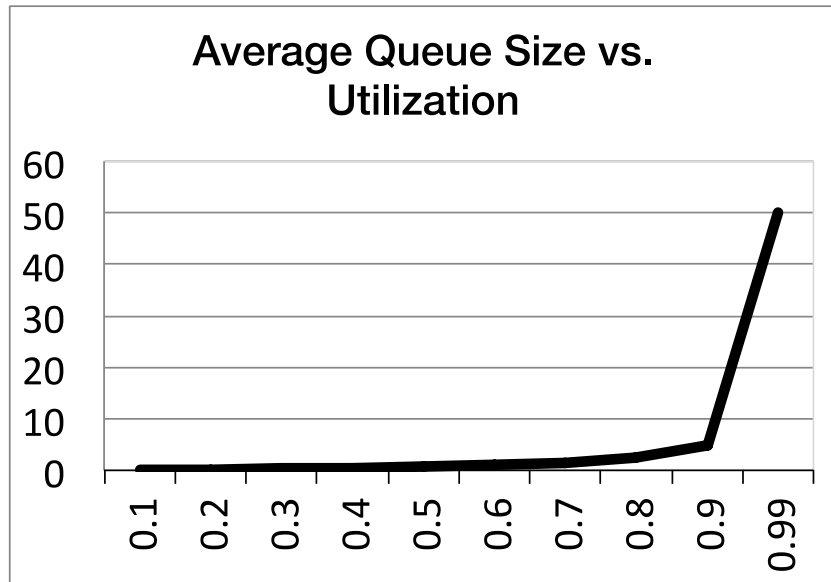
λ

### SSD Burst Buffer

μ

### Permanent Storage

Implications:

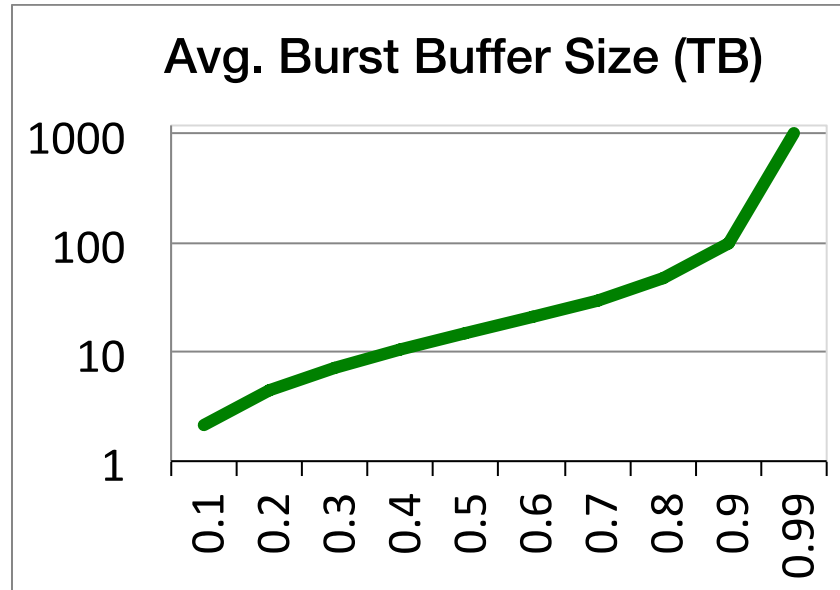Not all the data from burst buffer can be copied just in time

PFS Bandwidth still critical

Utilization = production rate/achievable PFS bandwidth = λ/μ

# SSDs as a burst buffer in HPC

### Average Queue Size vs. Utilization



Utilization= λ/μ

### Avg. Burst Buffer Size (TB)


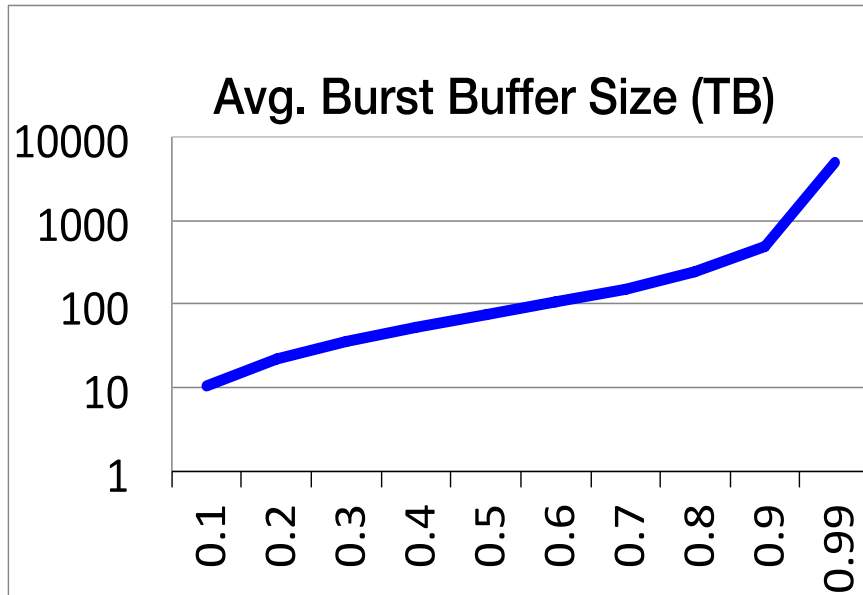
Utilization= λ/μ
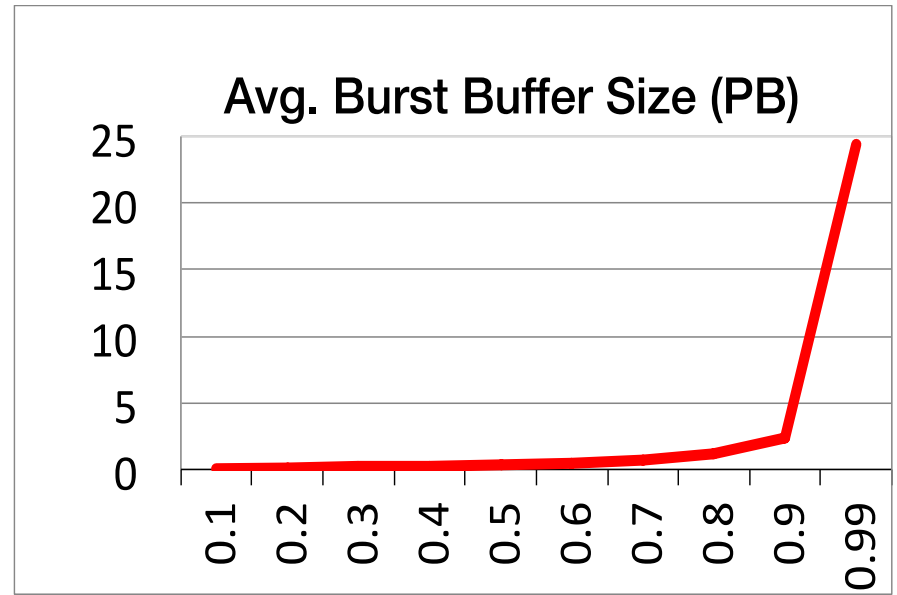
Output data size at once = 20TB

Queuing models suggest that avg. queue size (burst buffer size) increases exponentially with the increase in utilization

# SSDs as a burst buffer in HPC

## Avg. Burst Buffer Size (TB)

y-axis: 10000, 1000, 100, 10, 1

x-axis: 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.99

Utilization= λ/µ
Output data size at once = 100TB

## Avg. Burst Buffer Size (PB)

y-axis: 25, 20, 15, 10, 5, 0
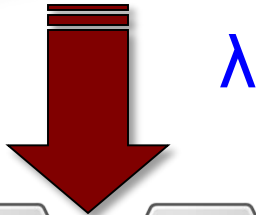
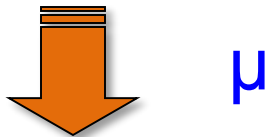x-axis: 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.99

Utilization= λ/µ
Output data size at once = 200TB

**1-5 PB of SSD storage may suffice depending on the amount of data being produced at each step**

# SSDs as a burst buffer in HPC
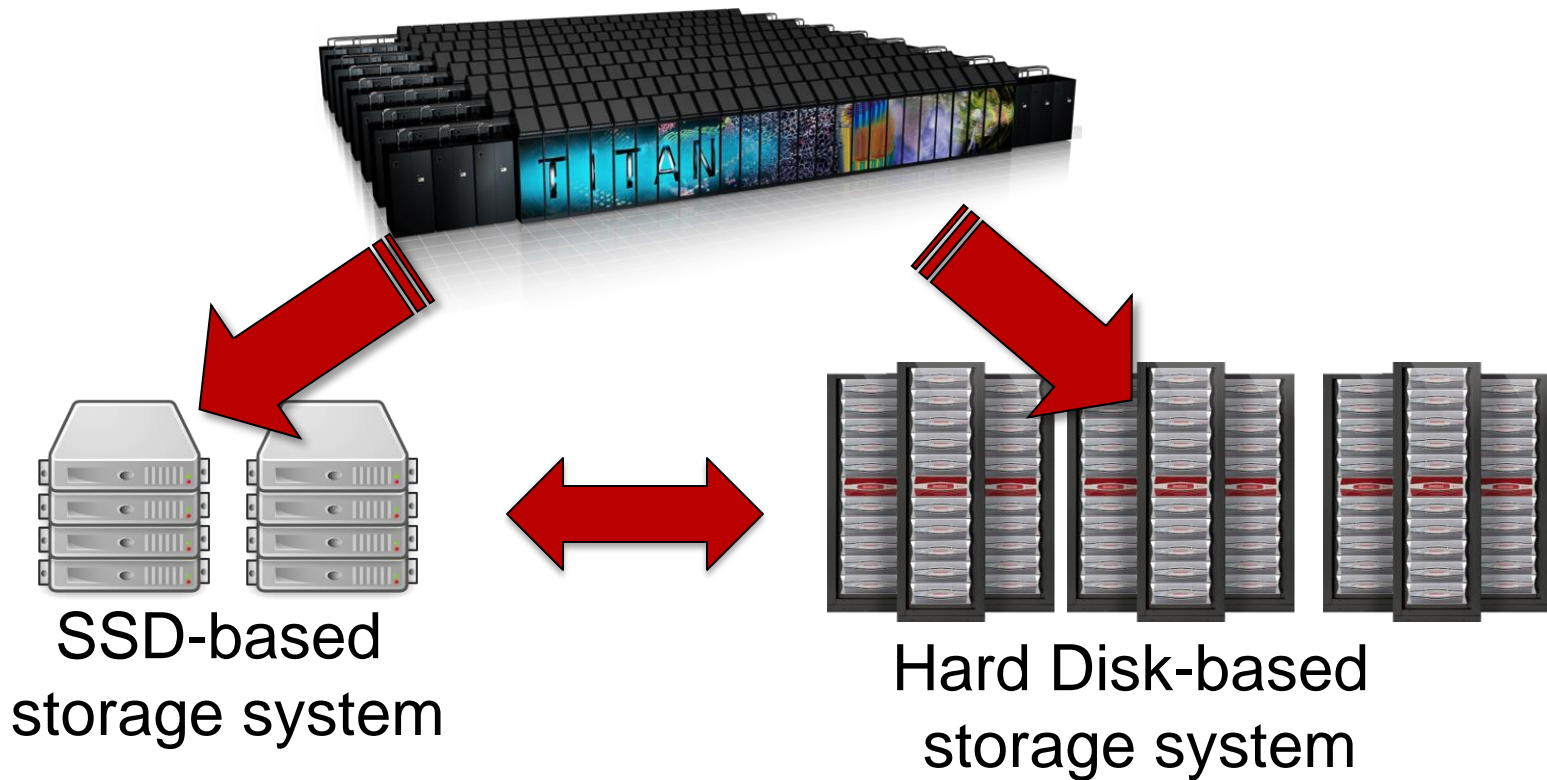


SSD Burst Buffer

$\lambda$

$\mu$

Permanent Storage

## This model causes excessive data-movement

We are using SSDs for draining writes, something they are fundamentally not good at.

# SSD in Exascale System Architecture

SSD-based
storage system

Hard Disk-based
storage system

Users will be charged differently for different kind of storage system they use
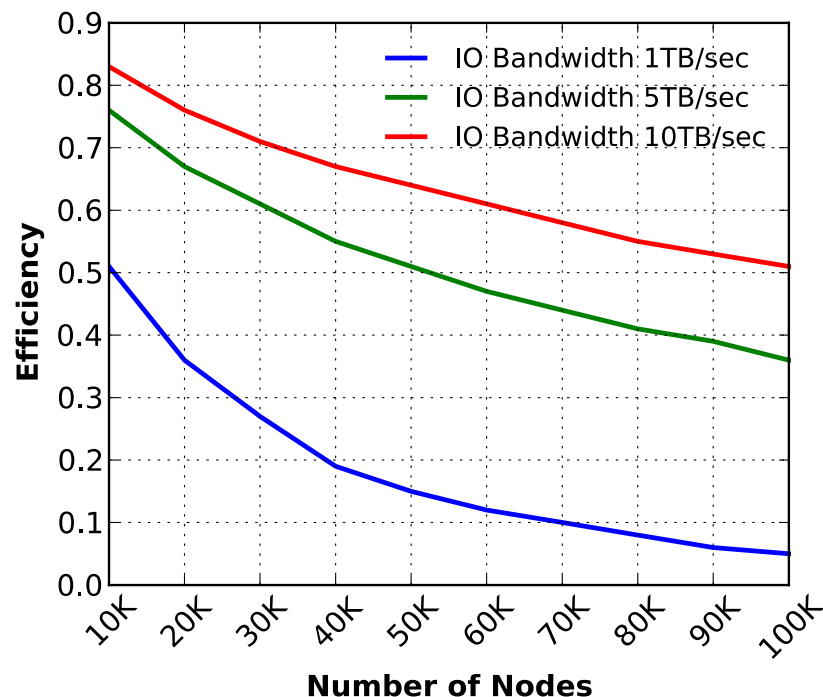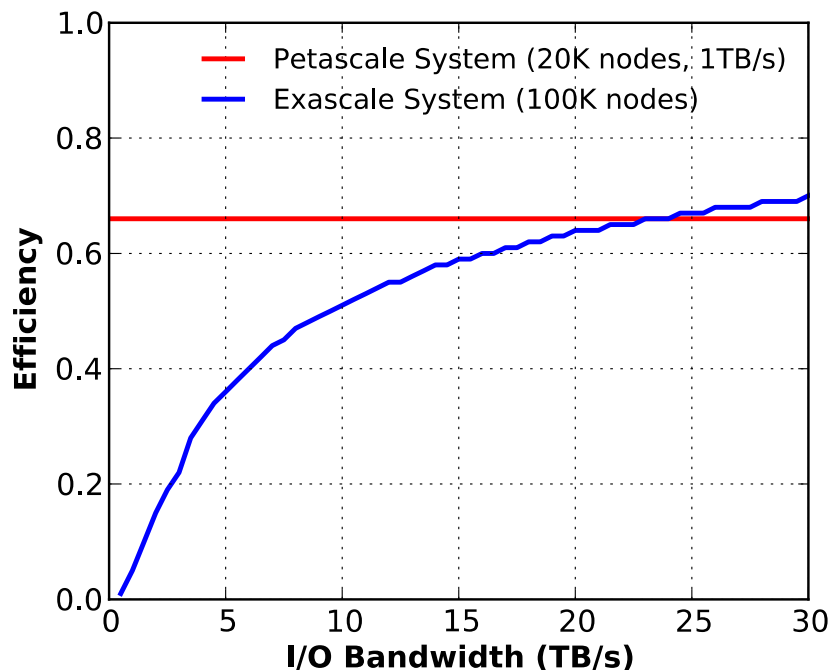
Disk system will become archival storage

Amount of waste work increases at larger scale system if we don't have a fast enough storage system to quickly take a checkpoint

Capital investment in more expensive storage system promises higher pay-off during operation
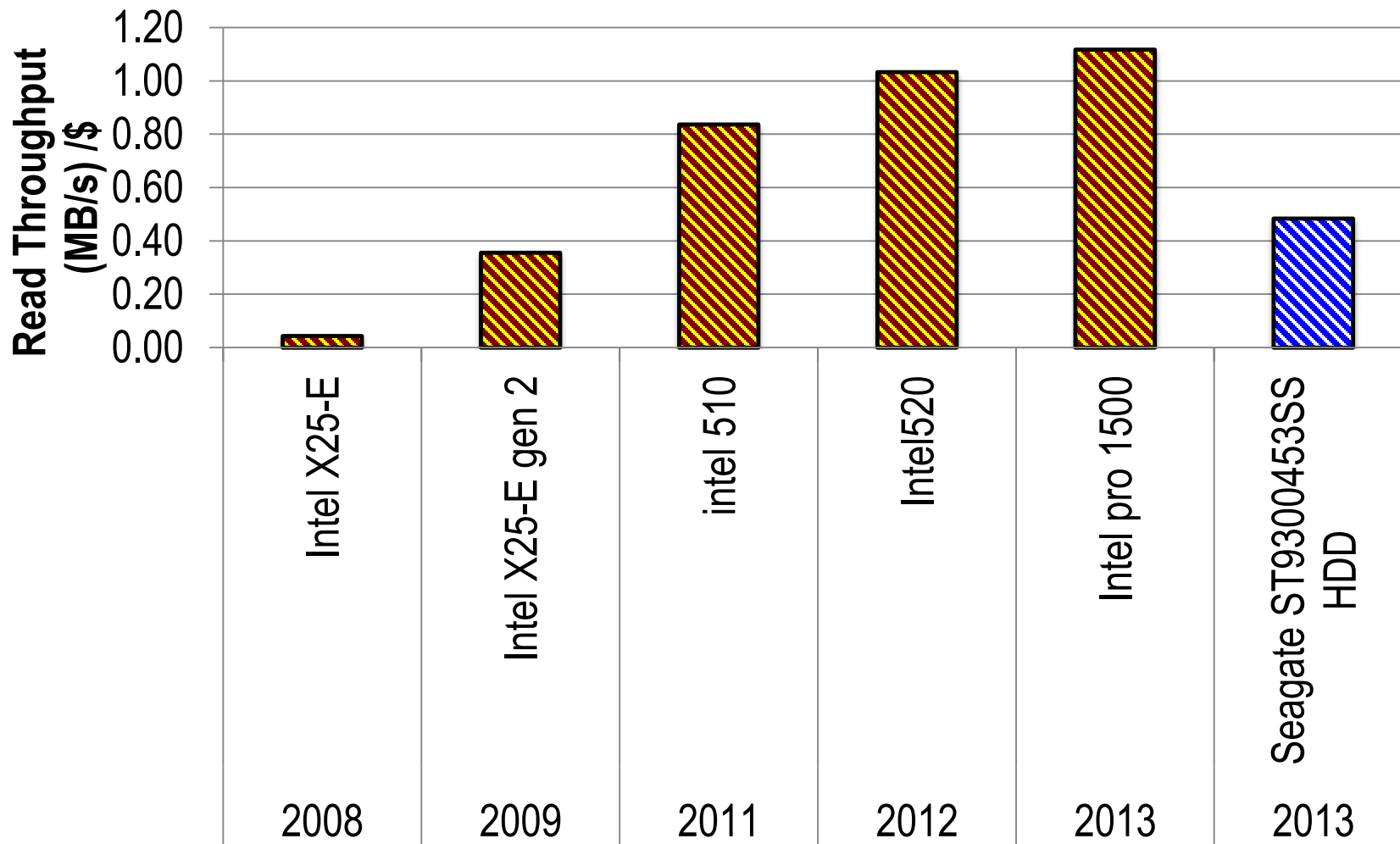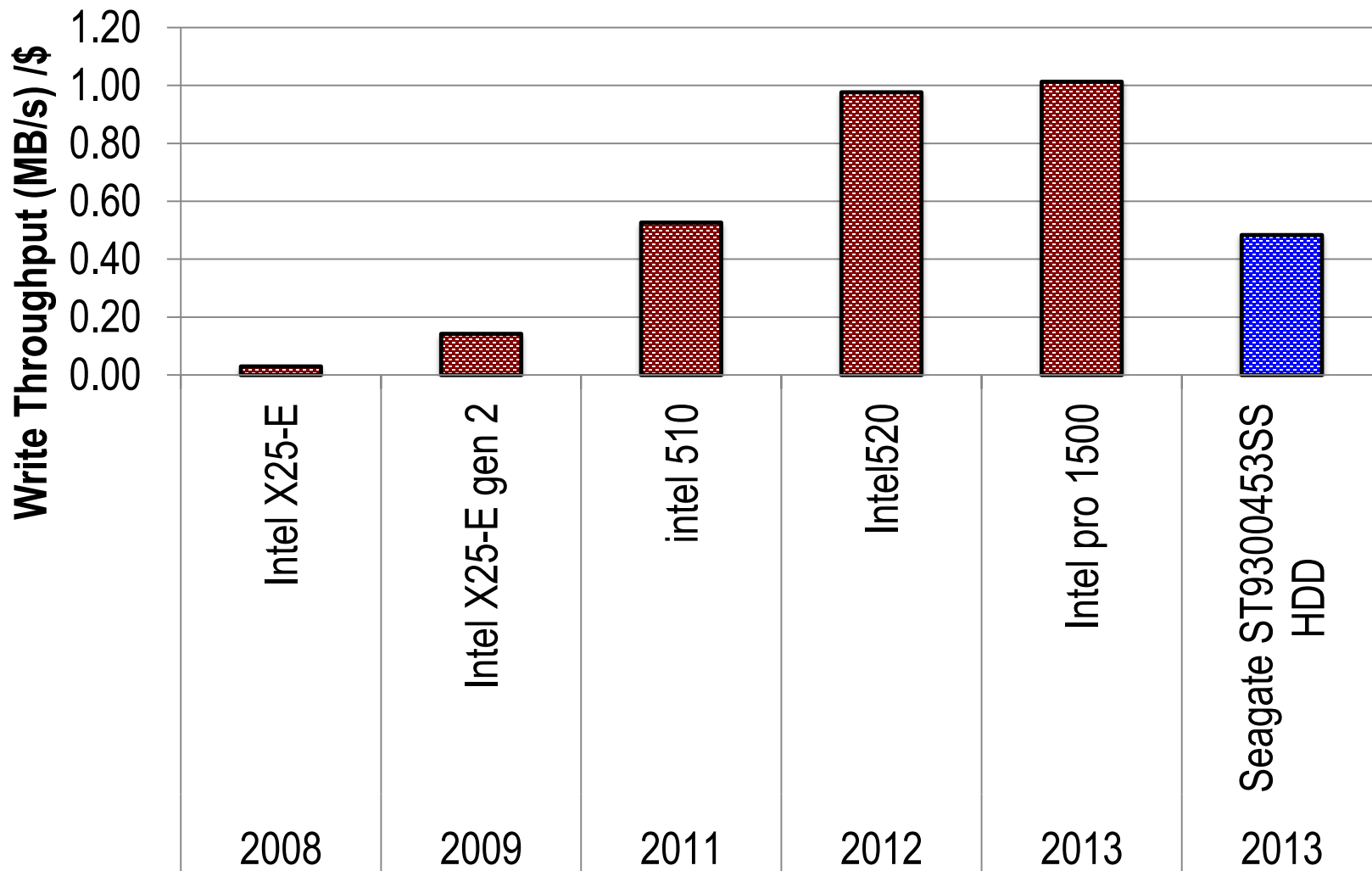
# Storage System Bandwidth



Storage system bandwidth determines the overall efficiency (i.e., amount of lost work due to failures)

Based on Daly's optimal checkpointing frequency (MTBF and time to checkpoint).
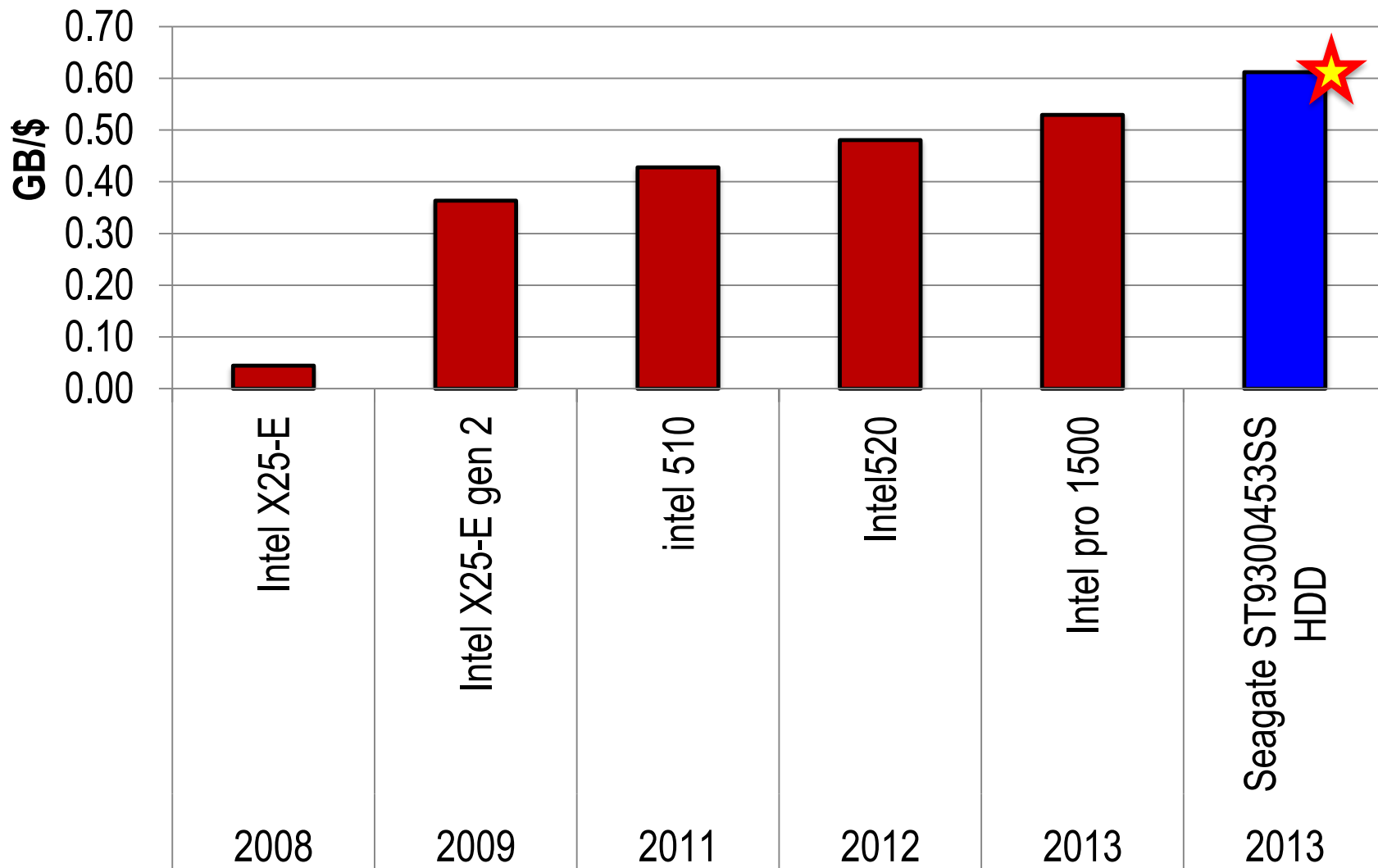
# SSD vs HDD



Chart: Read Throughput (MB/s) / $ by device and year.

| Device | Year | Read Throughput (MB/s) /$ |
|---|---|---|
| Intel X25-E | 2008 | ~0.03 |
| Intel X25-E gen 2 | 2009 | ~0.35 |
| intel 510 | 2011 | ~0.83 |
| Intel520 | 2012 | ~1.03 |
| Intel pro 1500 | 2013 | ~1.12 |
| Seagate ST9300453SS HDD | 2013 | ~0.48 |

# SSD vs HDD

# SSD vs HDD

# Building a HPC Storage System

Building blocks (cost factors):

   SSD/Hard Disk

   I/O Controller

   JBOD I/O module

   Enclosure

   Power and cooling

   RAID and firmware cost

Our preliminary study focuses on SSD/Hard disk and controllers costs.

# Building a HPC Storage System

- HDD 1TB $200, 4TB $500, performance* 200MB/s

- SSD  256GB $250, 1TB $1000, performance* 500MB/s

- A pair of mid-scale controllers: throughput 8GB/s, $8K


- Each shelf can hold up to 60 drives

- 40 hard-disks saturate one pair of controllers

- 16 SSDs saturate one pair of controllers

  *depends on the workload, lost on the route etc.

How to build a cost-effective high performance/high capacity HPC storage system?

# Building a HPC Storage System

If planning to build a high capacity system:

GB/$ for SSDs is too low

Controller and other components do not matter

Hard-disks tend to become cheaper as density increases, but not true for SSDs

# Building a HPC Storage System

Building 10TB/s system*

1280 pairs of controllers = ~$10million

Hard-disk based system 1280*40 (1TB) hard disks

disks cost = 1280*40*$200 = ~$10million

capacity = 50 PB, total cost = ~$20million

SSD based system  1280*40 (1TB) hard disks

disks cost = 1280*16*$250 = ~$5million

capacity = 5 PB, total cost = ~$15million

*numbers here only represent the trend, and are not (un)official quotes from WarpMechanics

# Building a HPC Storage System

Controller cost very important factor in deciding the total capital cost of the system

This is is fundamentally different from enterprise computing that focuses on per drive IOPS
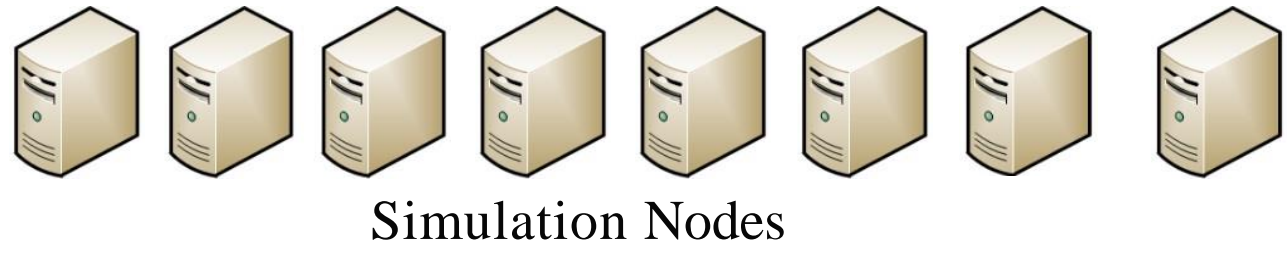
Need to account firmware, and other hardware cost

On-going work: workload characterization to explore   a hybrid storage system, random vs. sequential performance, Lustre and RAID overhead etc.

# What else can SSDs do for us?

Tiwari et al., Active Flash: Towards Energy-Efficient, In-Situ Data Analytics on Extreme-Scale Machines, USENIX FAST 2013.

# Traditional Scientific Data Analysis Approach



Simulation Nodes

Parallel File System
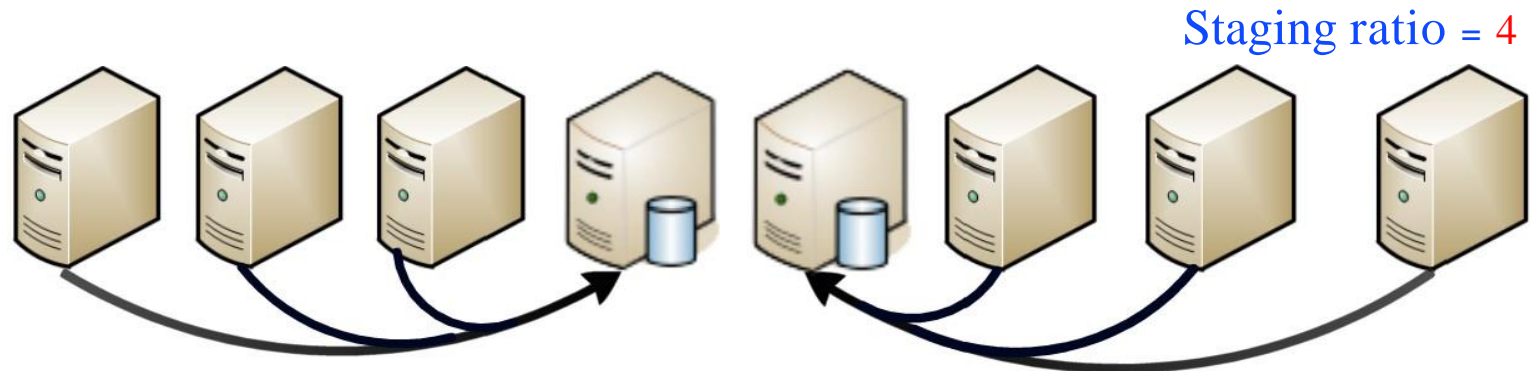
Offline Data Analysis Cluster
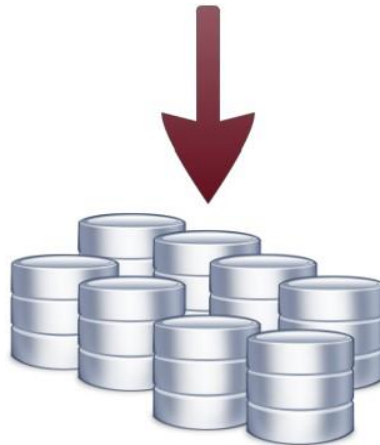
Regex matching, statistics collection, clustering, compression, etc.

"Energy-cost for data movement at Exascale is likely to be of the same order of computation cost, if not more!"

-- Exascale Computing Study, 2008
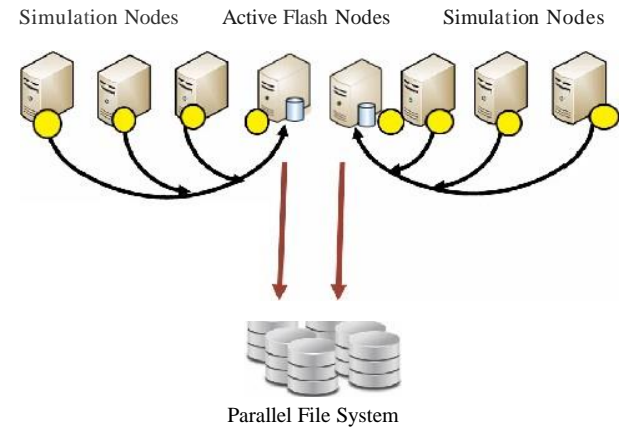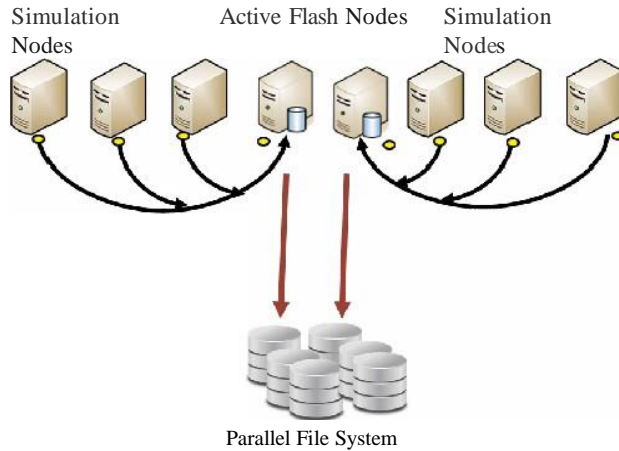Principle Investigator: Peter Kogge

# Active Computation on SSDs

Staging ratio = 4



Scientific data analysis performed    on SSD controllers
in-parallel with simulation without affecting it
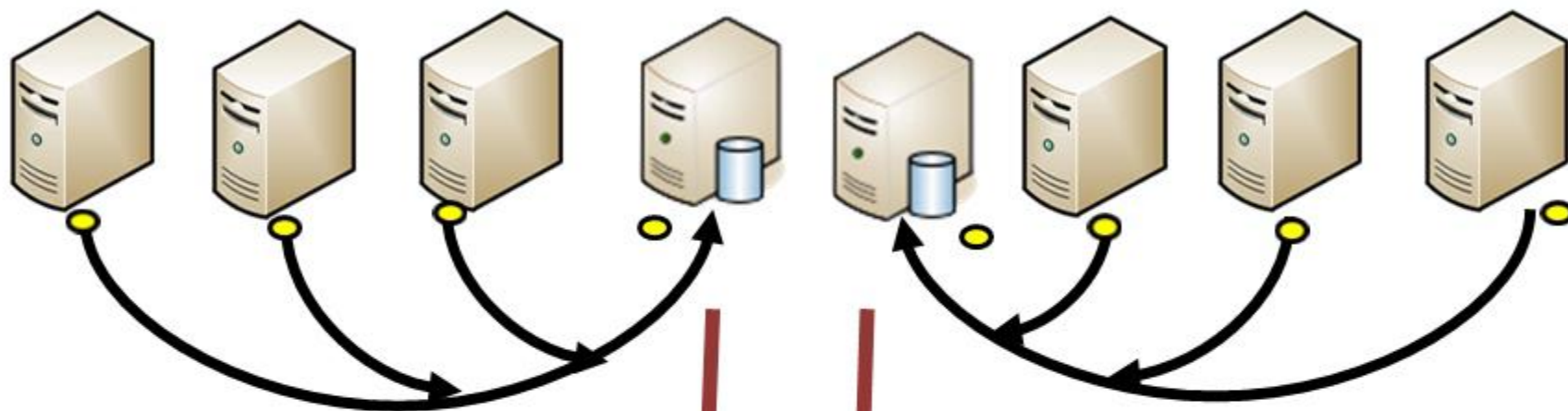
Parallel File System

**An analysis kernel needs to meet a "threshold compute throughput" to be placed on SSD controllers**

$$T_{SSD\_k} > \frac{\lambda_a \cdot R_{SSD}}{1 - \lambda_a \cdot R_{SSD} \cdot \left(\frac{1}{BW_{fm2c}} + \frac{1}{BW_{c2m}}\right) - \frac{N \cdot (\alpha \cdot \lambda_a + \lambda_c)}{BW_{PFS}} - \frac{t_i}{t_{iter}}}$$

Simulation Nodes    Active Flash Nodes    Simulation Nodes

○ Simulation Output (per unit time)

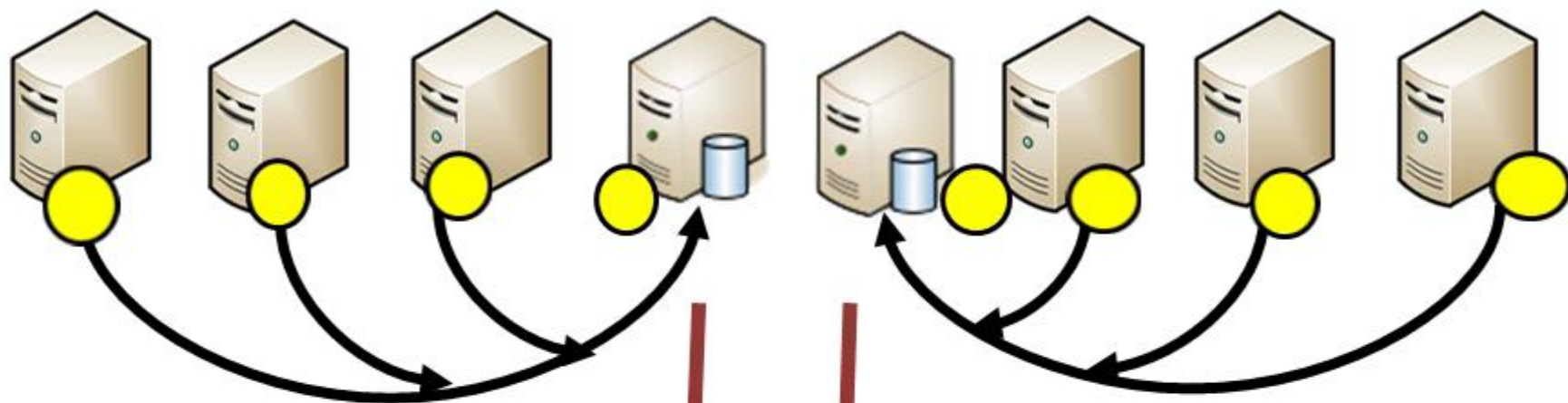✳ Data Analysis Kernel (high compute intensive) Low compute throughput

Parallel File System

Simulation Nodes   Active Flash Nodes   Simulation Nodes
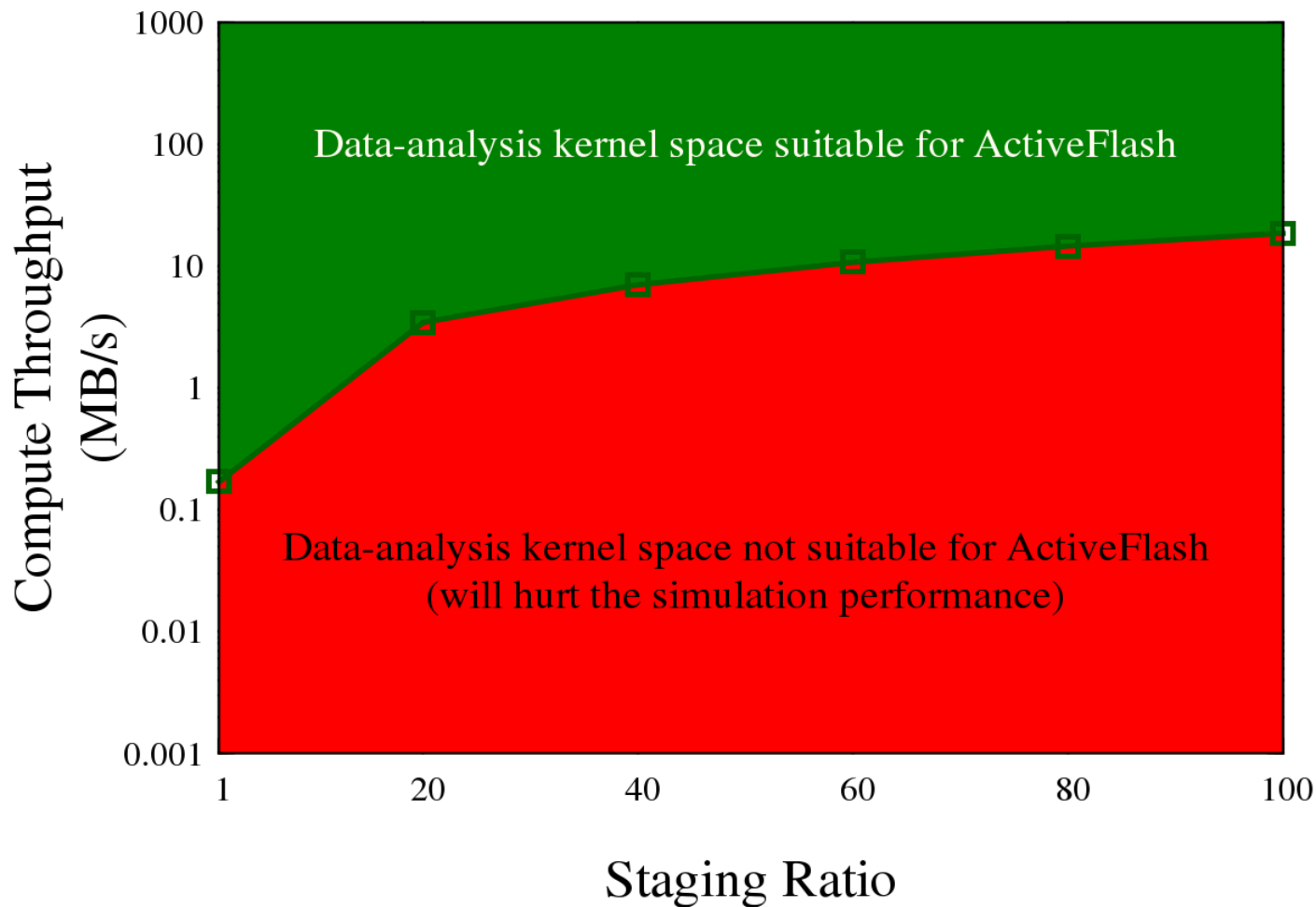
○ Simulation Output (per unit time)

✹ Data Analysis Kernel (less compute intensive) High compute throughput

Parallel File System
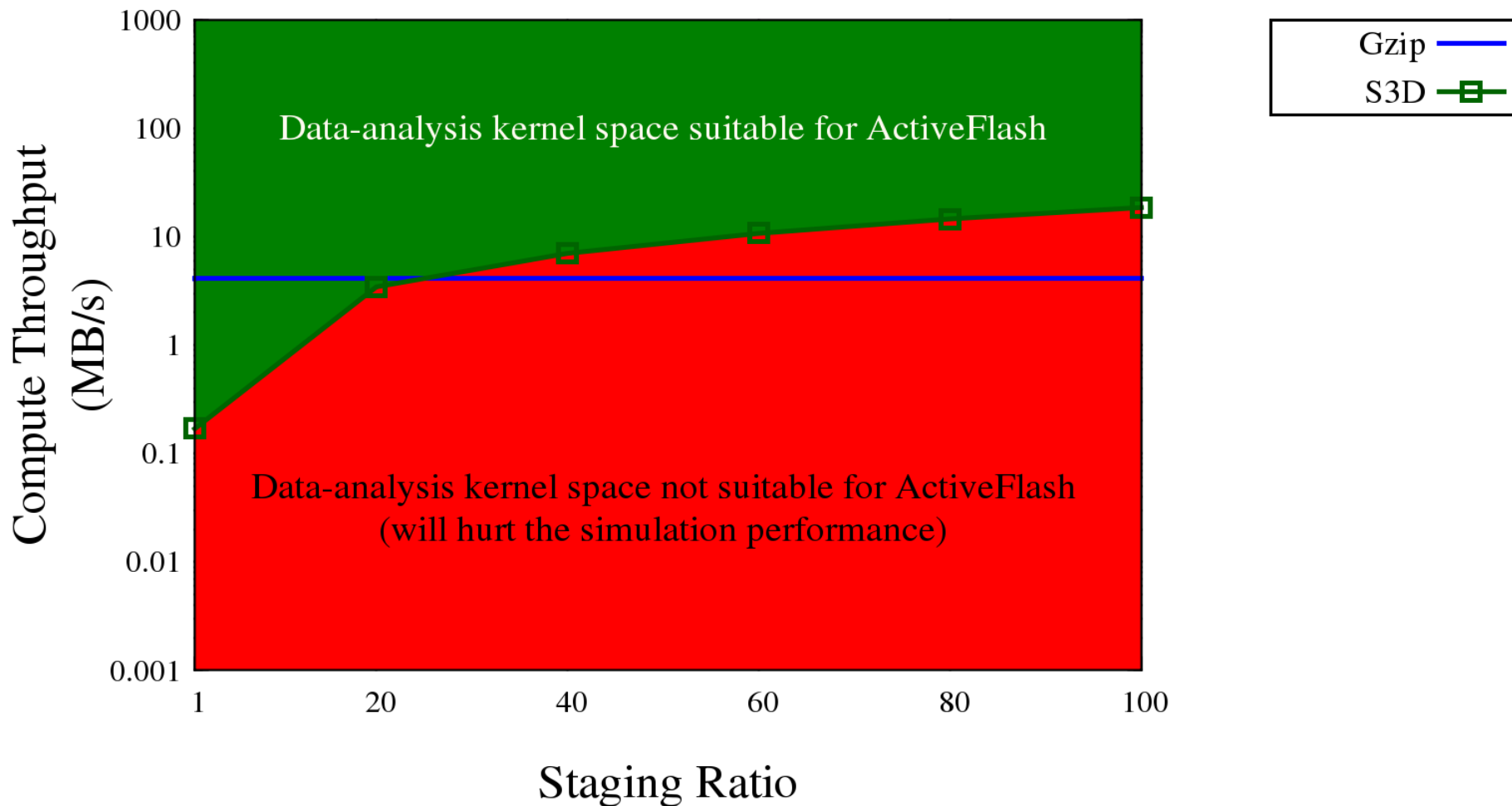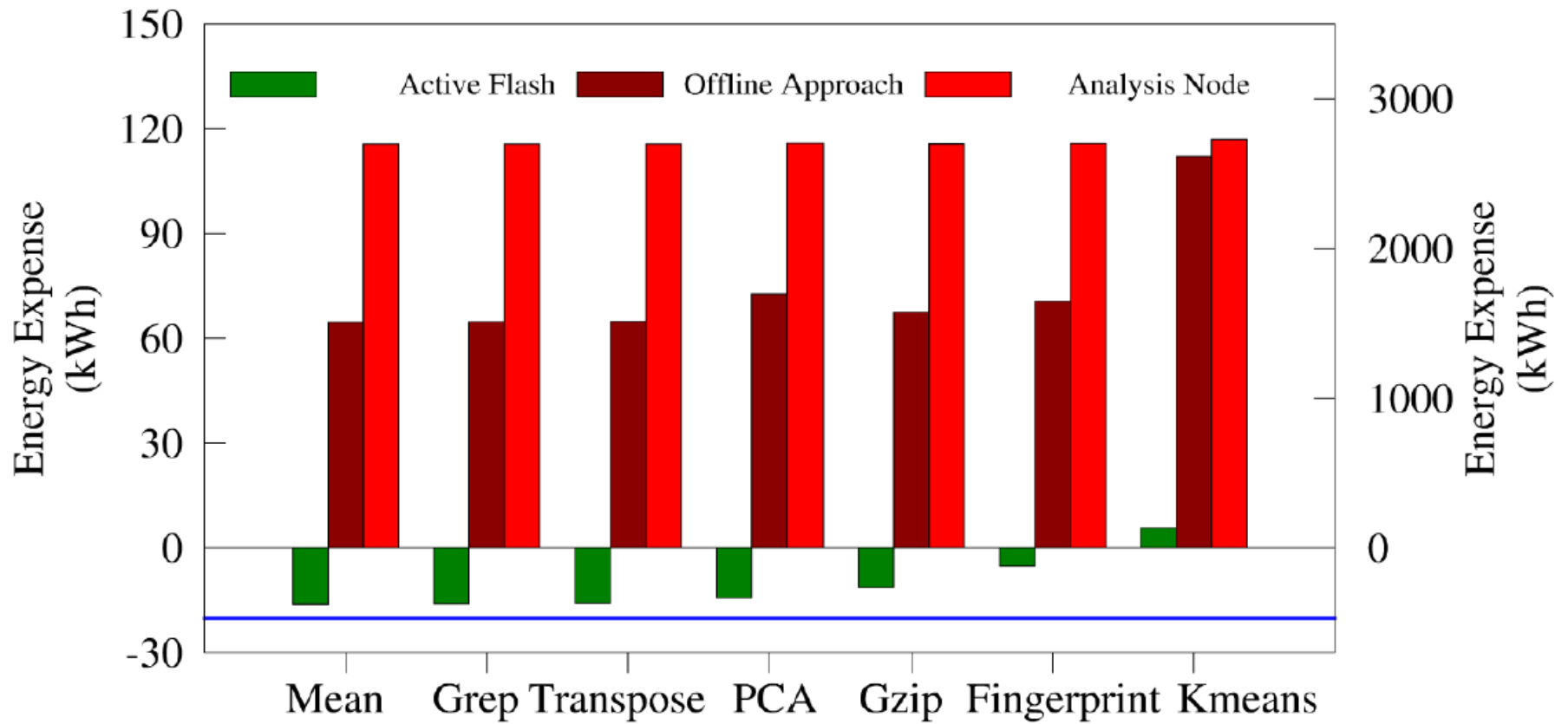
# Feasibility of Active Flash Approach



Compute Throughput (MB/s) vs. Staging Ratio.

Data-analysis kernel space suitable for ActiveFlash

Data-analysis kernel space not suitable for ActiveFlash
(will hurt the simulation performance)

S3D

# Feasibility of Active Flash Approach



Legend:
- Gzip
- S3D

Y-axis: Compute Throughput (MB/s), values: 1000, 100, 10, 1, 0.1, 0.01, 0.001

X-axis: Staging Ratio, values: 1, 20, 40, 60, 80, 100

Data-analysis kernel space suitable for ActiveFlash

Data-analysis kernel space not suitable for ActiveFlash
(will hurt the simulation performance)

Application: POP

**Finding:** Most data analysis kernels can be placed on SSD controllers without degrading simulation performance

Tiwari et al., Active Flash: Towards Energy-Efficient, In-Situ Data Analytics on Extreme-Scale Machines, USENIX FAST 2013.

**Finding:** Additional SSDs are not required for supporting in-situ data analysis on SSDs, beyond what is needed for sustaining the I/O requirements of scientific applications

Tiwari et al., Active Flash: Towards Energy-Efficient, In-Situ Data Analytics on Extreme-Scale Machines, USENIX FAST 2013.

# ActiveFlash Prototype based on OpenSSD Platform

Prototype demonstrates the viability of our approach

Changes only in the FTL, no hardware changes

Preemption based scheduling

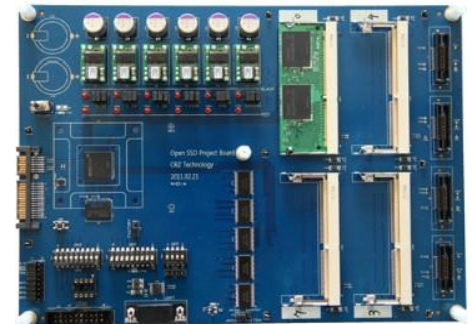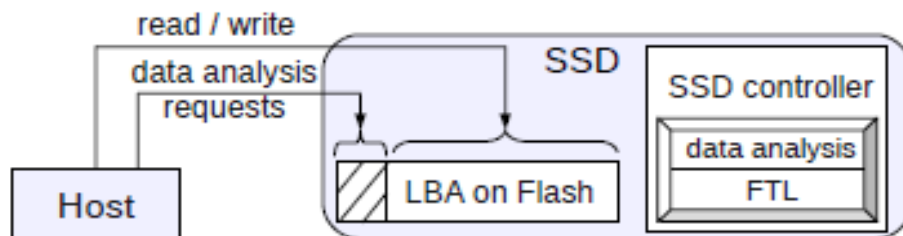See USENIX paper for the details and evaluation results



Figure courtesy: open-ssd project

# Contact
## Devesh Tiwari
## tiwari@ornl.gov