# Testing methodology for large-scale file systems

**Sarp Oral, PhD**

**Technology Integration Group**

*Lustre User Group Meeting, April 12, 2011*

# Large-scale FS testing

- Benchmarking and testing large-scale file and storage systems is not straight forward
  - Scale
  - Complexity
  - Limited toolset

- A black art
  - Tricks/tips/knowledge needs to passed down to new testers

- Developing a complete and quantitative test methodology
  - Work in progress

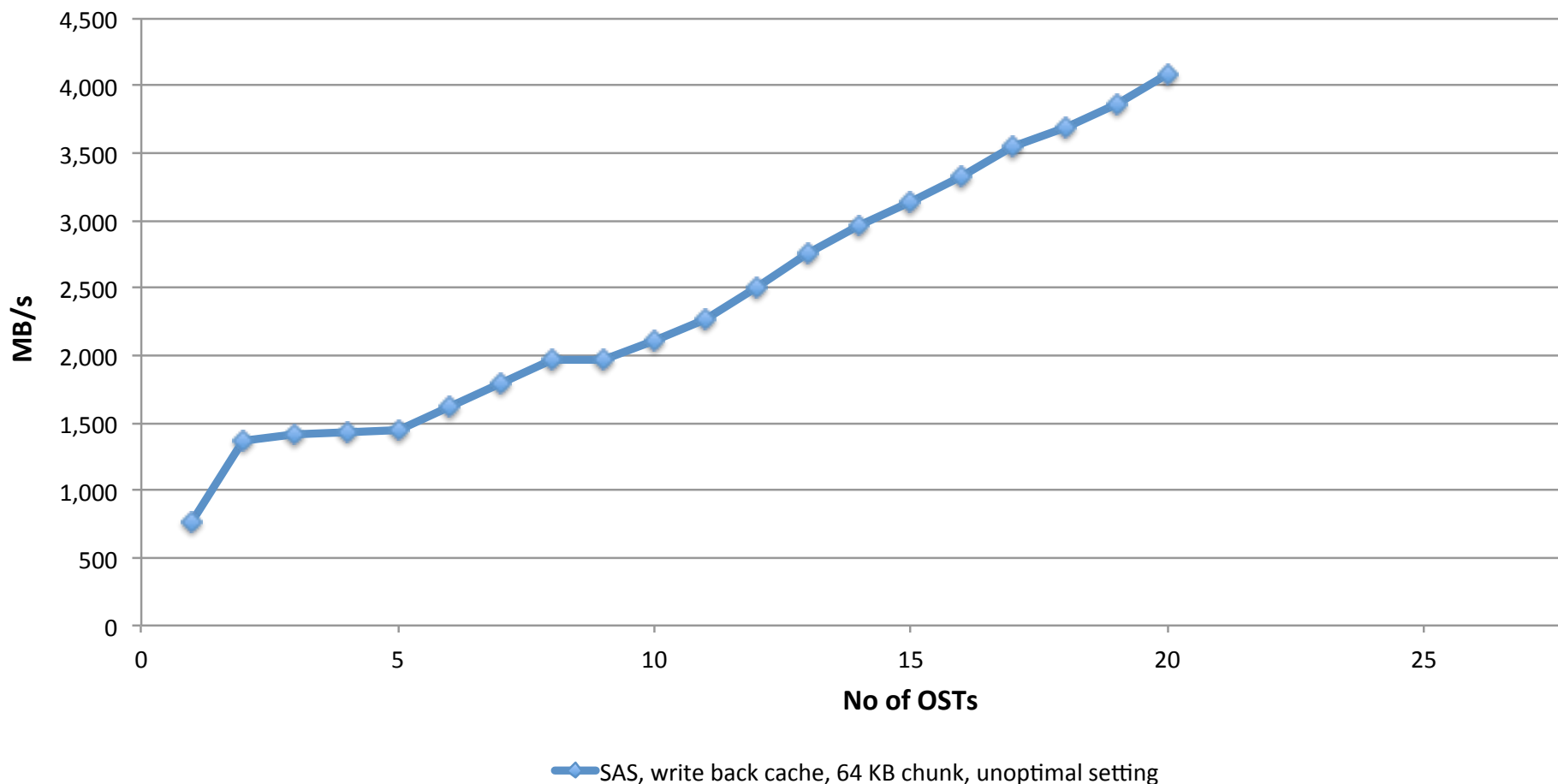ロLCF ● ● ● ●

OAK RIDGE
National Laboratory

# Rules of thumb, #1

- Know thy hardware!
  - Identify all components of the disk backend system and how they are setup and configured
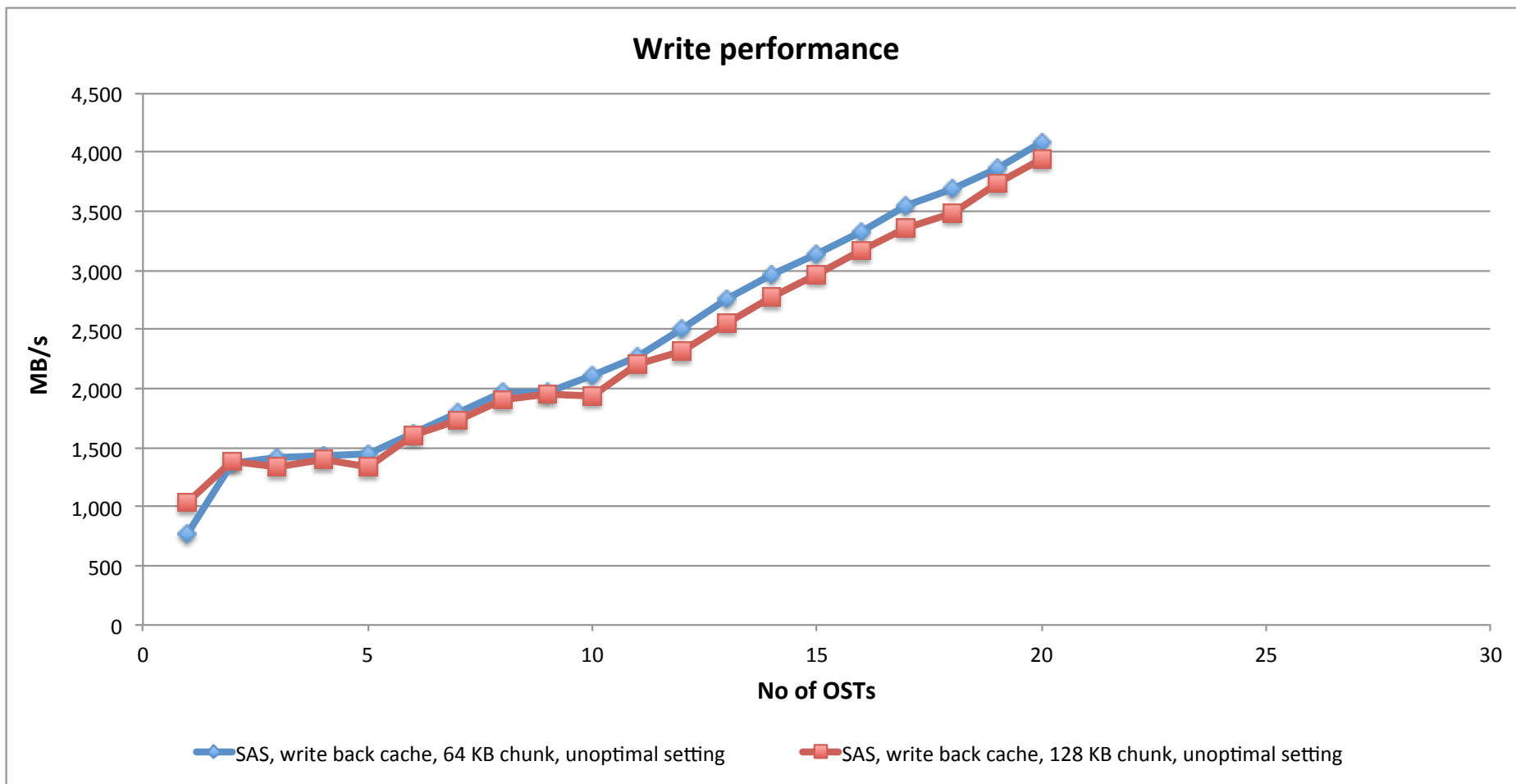    - Setup/configuration directly impacts the expected/observed performance

OLCF ● ● ● ●

OAK RIDGE
National Laboratory

# Understanding observed data

- 4 Hosts, QDR IB, 200 SAS disks, R6 (8+2), a pair of HW RAID controllers, obdfilter-survey

**Write performance**



SAS, write back cache, 64 KB chunk, unoptimal setting

OLCF ● ● ● ●

OAK RIDGE
National Laboratory

# Understanding observed data

- 4 Hosts, QDR IB, 200 SAS disks, R6 (8+2), a pair of HW RAID controllers, obdfilter-survey

**Write performance**



Chart: Y-axis "MB/s" ranging 0 to 4,500; X-axis "No of OSTs" ranging 0 to 30.

Legend:
- SAS, write back cache, 64 KB chunk, unoptimal setting
- SAS, write back cache, 128 KB chunk, unoptimal setting

OLCF

OAK RIDGE
National Laboratory

# Understanding observed data

- 4 Hosts, QDR IB, 200 SAS disks, R6 (8+2), a pair of HW RAID controllers, obdfilter-survey



**Write performance**
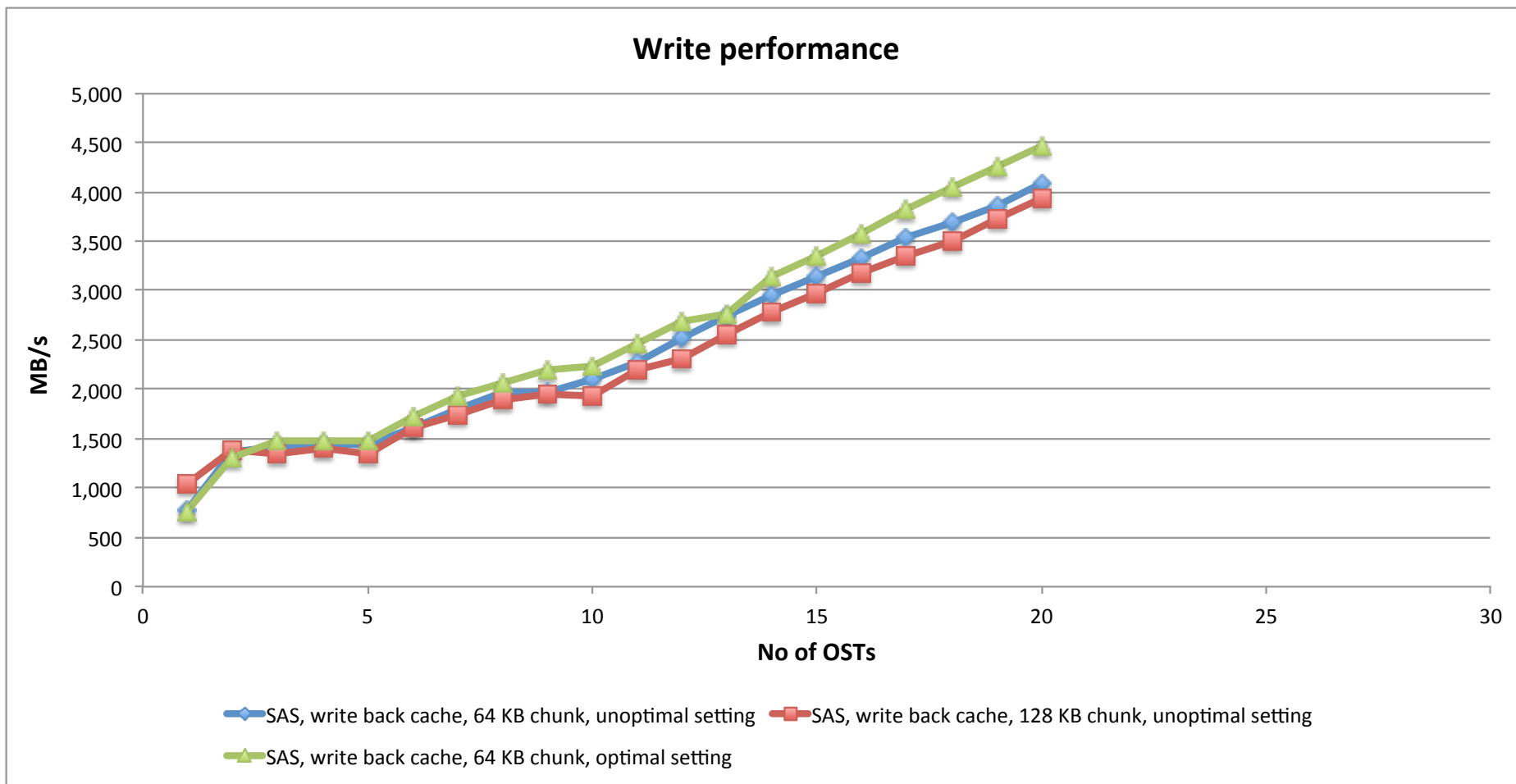
Legend:
- SAS, write back cache, 64 KB chunk, unoptimal setting
- SAS, write back cache, 128 KB chunk, unoptimal setting
- SAS, write back cache, 64 KB chunk, optimal setting

x-axis: No of OSTs
y-axis: MB/s

# Understanding observed data

- 4 Hosts, QDR IB, 200 SAS disks, R6 (8+2), a pair of HW RAID controllers, obdfilter-survey



**Write performance**
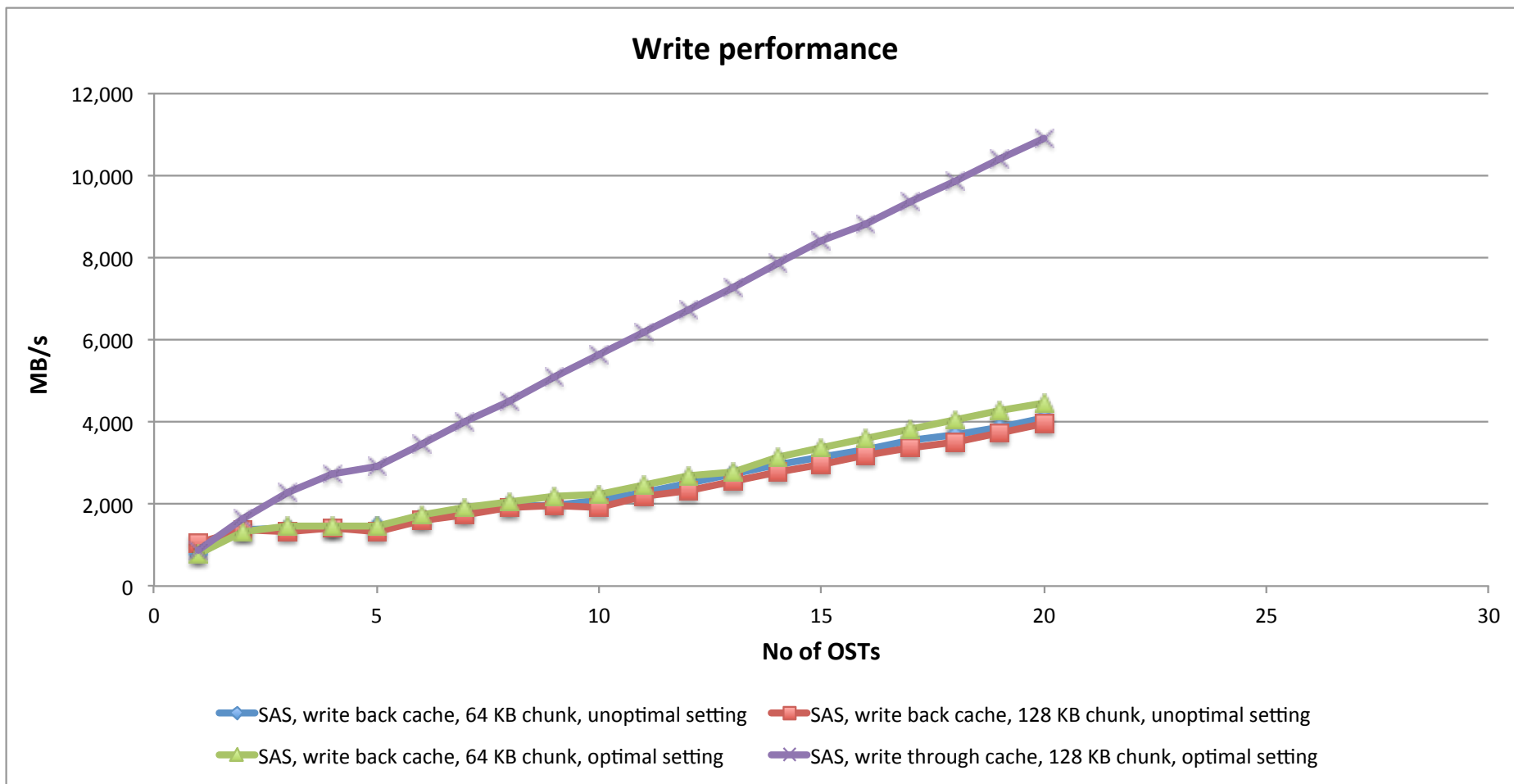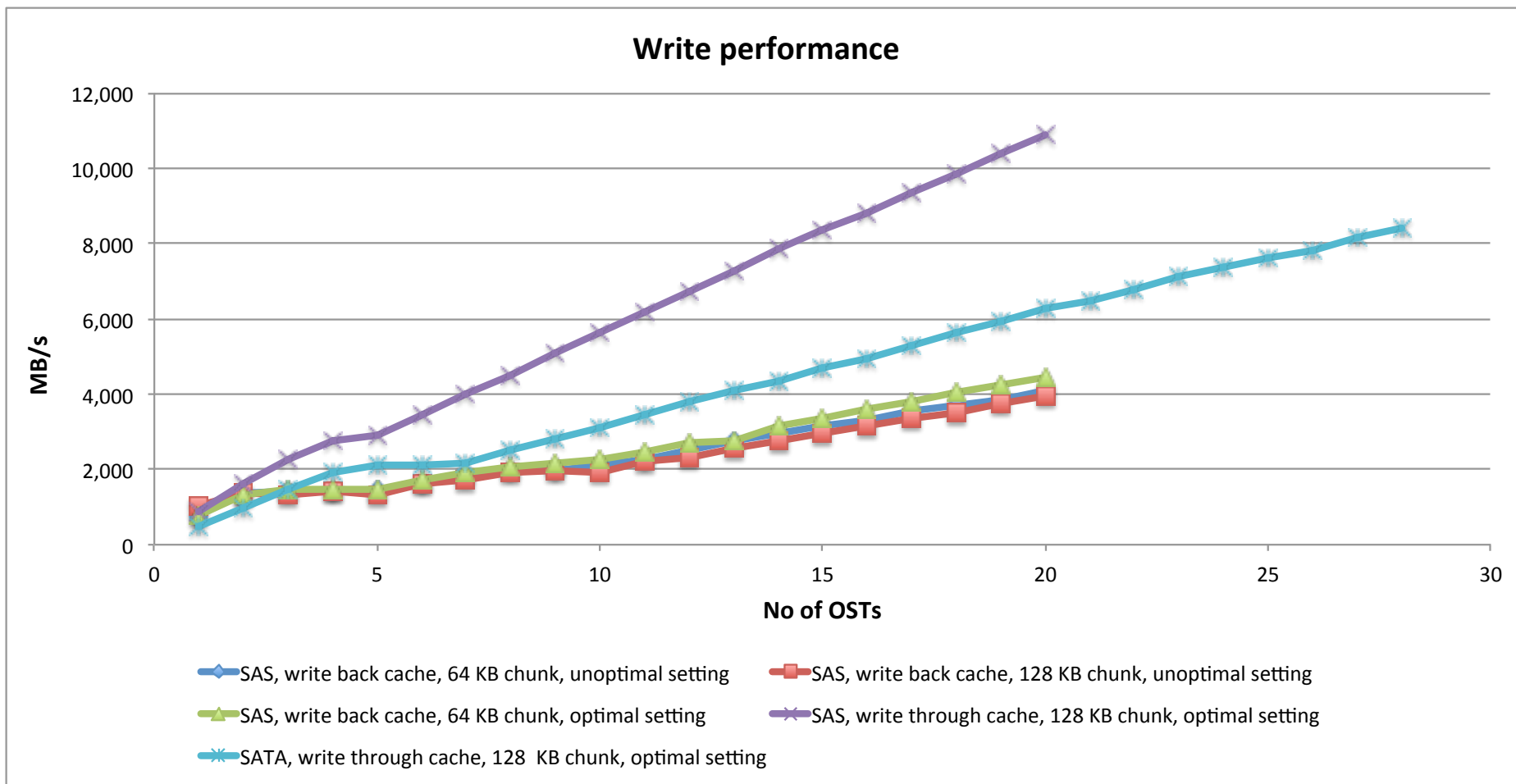
Chart axes: Y-axis "MB/s" (0 to 12,000), X-axis "No of OSTs" (0 to 30)

Legend:
- SAS, write back cache, 64 KB chunk, unoptimal setting
- SAS, write back cache, 128 KB chunk, unoptimal setting
- SAS, write back cache, 64 KB chunk, optimal setting
- SAS, write through cache, 128 KB chunk, optimal setting
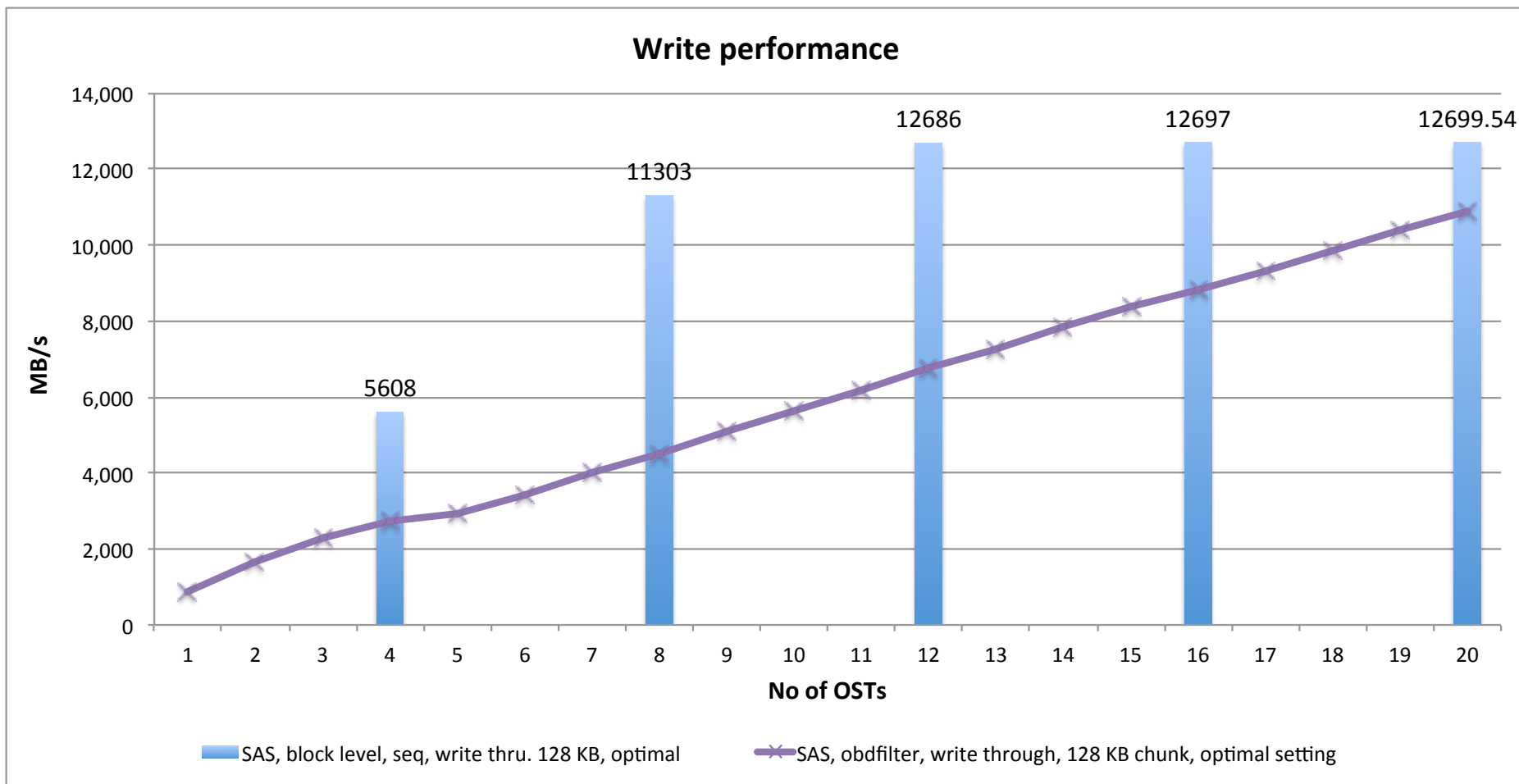
OLCF

OAK RIDGE
National Laboratory

# Understanding observed data

- 4 Hosts, QDR IB, a pair of HW RAID controllers, obdfilter-survey
    - 200 SAS disks, R6 (8+2)
    - 280 SATA disks, R6 (8+2)



**Write performance**

Legend:
- SAS, write back cache, 64 KB chunk, unoptimal setting
- SAS, write back cache, 128 KB chunk, unoptimal setting
- SAS, write back cache, 64 KB chunk, optimal setting
- SAS, write through cache, 128 KB chunk, optimal setting
- SATA, write through cache, 128 KB chunk, optimal setting

OAK RIDGE National Laboratory

# Understanding observed data

- 4 Hosts, QDR IB, a pair of HW RAID controllers, 200 SAS disks, R6 (8+2)
  - Obdfilter-survey
  - In-house coded benchmark



**Write performance**

SAS, block level, seq, write thru. 128 KB, optimal

SAS, obdfilter, write through, 128 KB chunk, optimal setting

# Rules of thumb, #2

- Establish a clear understanding of the critical data path

    – From the disk backend host port to disk

    – Theoretical performance of all components on the critical data path

    – The system is a combination of serial and parallel (and combination of these two in some cases) connected set of devices

    – Lowest performing component will lower the overall performance

OAK
RIDGE
National Laboratory

# Rules of thumb, #3

- Understand how each component interacts with each other
  - A component's performance response may be different if it exercised differently
    - A RAID set might perform differently when there is background disk scrubbing going on
      - Contention on disks
    - A group of RAID sets may perform differently when only one of them is exercised compared to when all in the group exercised concurrently
      - Contention on RAID controller

OAK RIDGE
National Laboratory

# Rules of thumb, #4

- Repeat the tests, verify the variability
  - An obtained performance may vary with time
    - Even under the same test conditions

# Rules of thumb, #5

- Learn disk backend system's reporting and statistics mechanisms

  – Benchmarks will report the observed performance for a given test

  – Disk backend systems **_should_** provide some level of internal statistics and performance figures

  – Identify and familiarize with these and compare them with benchmark reported figures

OAK RIDGE
National Laboratory

# Rules of thumb, #6

- Test from bottom-to-top
  - Start testing with the very basics and from the bottom of the disk backend system

  - Compare observed results with theoretical/expected result

  - Identify performance gaps, if any and explain them

  - Keep testing by adding one more component to the test setup and move towards the top

  - Bottom, naturally is the disks and top is the host ports for a given disk backend system

OAK RIDGE
National Laboratory

# Rules of thumb, #7

- Have an expected result (hypothesis) before actually running any test
  - This will help analyzing the results and tester's understanding of the underlying hardware

OAK
RIDGE
National Laboratory

# Rules of thumb, #8

- Use the right tool for the right job
  - Do not use a cannon to kill a fly!
    - If you want to exercise the disks at the very basic level (e.g. block level), do not use a user-level MPI-based benchmark (e.g. IOR)

OAK
RIDGE
National Laboratory

# Rules of thumb, #9

- Be a good recorder; document everything!
  - You will forget
    - *What you did*
    - *How you did it*
    - *What were the results*
  - Down the road you will need to revisit the results or the tests

OAK RIDGE
National Laboratory

# Rules of thumb, #9 (continued)

- Document everything after a given test
- Keep notes of
  - Benchmark command lines
  - Configurations settings
  - Client configurations
  - Storage and file system settings
- If possible, write a few sentences about the test and the results
- Archive and time stamp

OAK RIDGE
National Laboratory

# Putting all together

- Know thy hardware

- Understand the critical data path

- Understand how component interacts with each other

- Repeat the tests, verify the variability

- Gather stats and performance data from disk backend

- Test bottom-to-top

- Have a hypothesis before running the test

- Use right tools

- Record and document everyhting

OAK RIDGE
National Laboratory

# Thank You

- Questions?

> Sarp Oral
>
> Research and Development Staff Member
>
> 865-574-2173
>
> oralhs@ornl.gov