



Lustre* Client I/O Parallelization

Dmitry Eremin

LUG17

* Some names and brands may be claimed as the property of others.

Agenda

- What issues are we solving?
- Current I/O performance
- I/O data flows and Flame Graphs
- Ideas of Improvements
- Results

What issues are we solving?

- Single-threaded applications cannot utilize performance benefits of multiple I/O operations in Lustre* if a single core is slow
- The Network bandwidth cannot be saturated because of
 - Intel® OPA fabric uses a host CPU for packets processing
 - High overhead of single thread execution
 - Slow memory transfer operations
- Blocking the userspace I/O operation
 - To do Read Ahead

Number of cores in client machines is continually increasing, yet single-threaded I/O is still common.

A solution other than adding more compute nodes is needed.

* Some names and brands may be claimed as the property of others.

Hardware

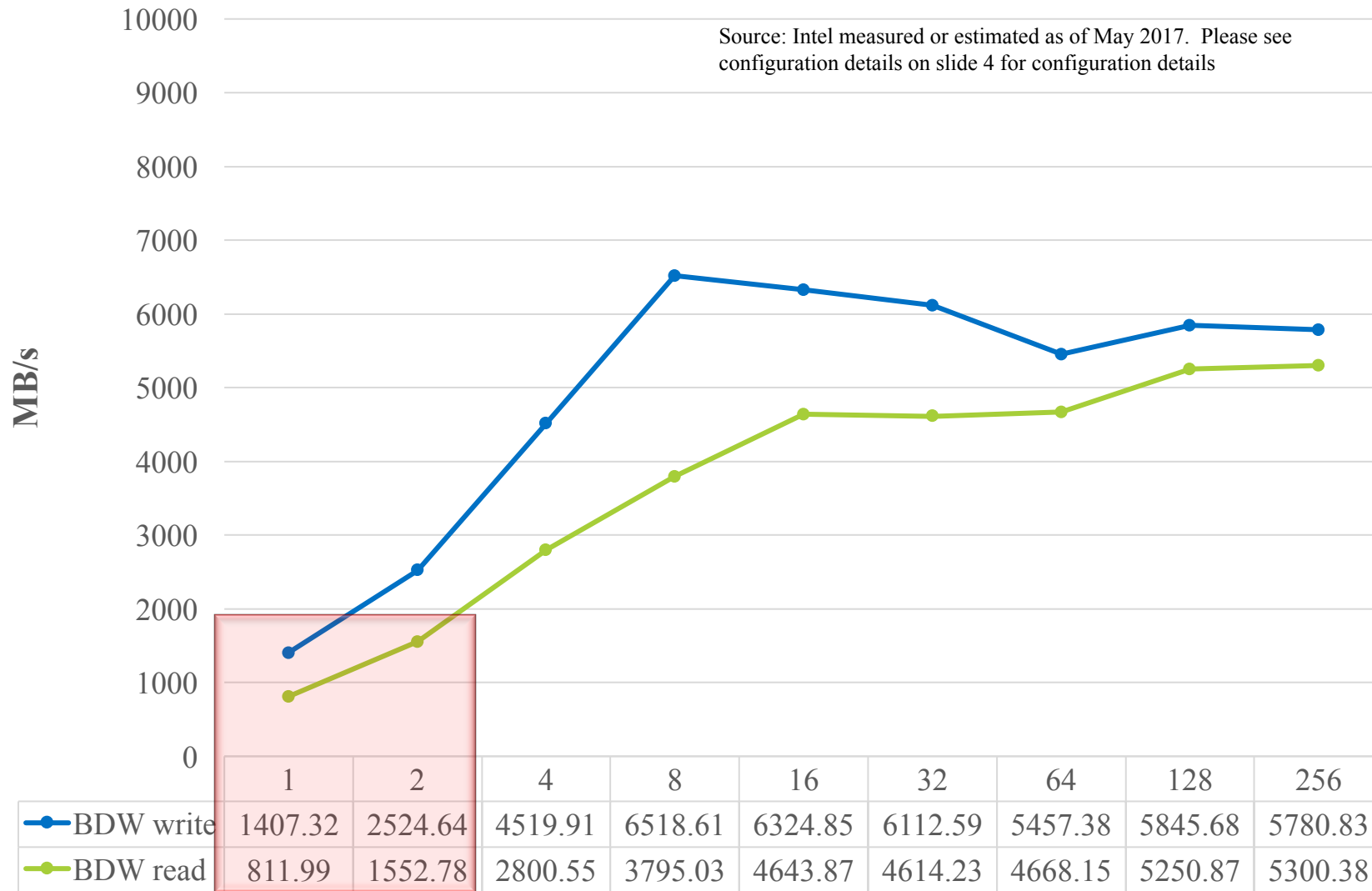
Intel® Xeon® CPU E5-2697 v4 @ 2.30GHz

- L2 cache: 256K, L3 cache: 46080K
- CPUs: 72, Cores: 36, Threads per core: 2
- Intel® OPA HFI Silicon 100 Series [discrete]

Intel® Xeon Phi™ CPU 7250 @ 1.40GHz

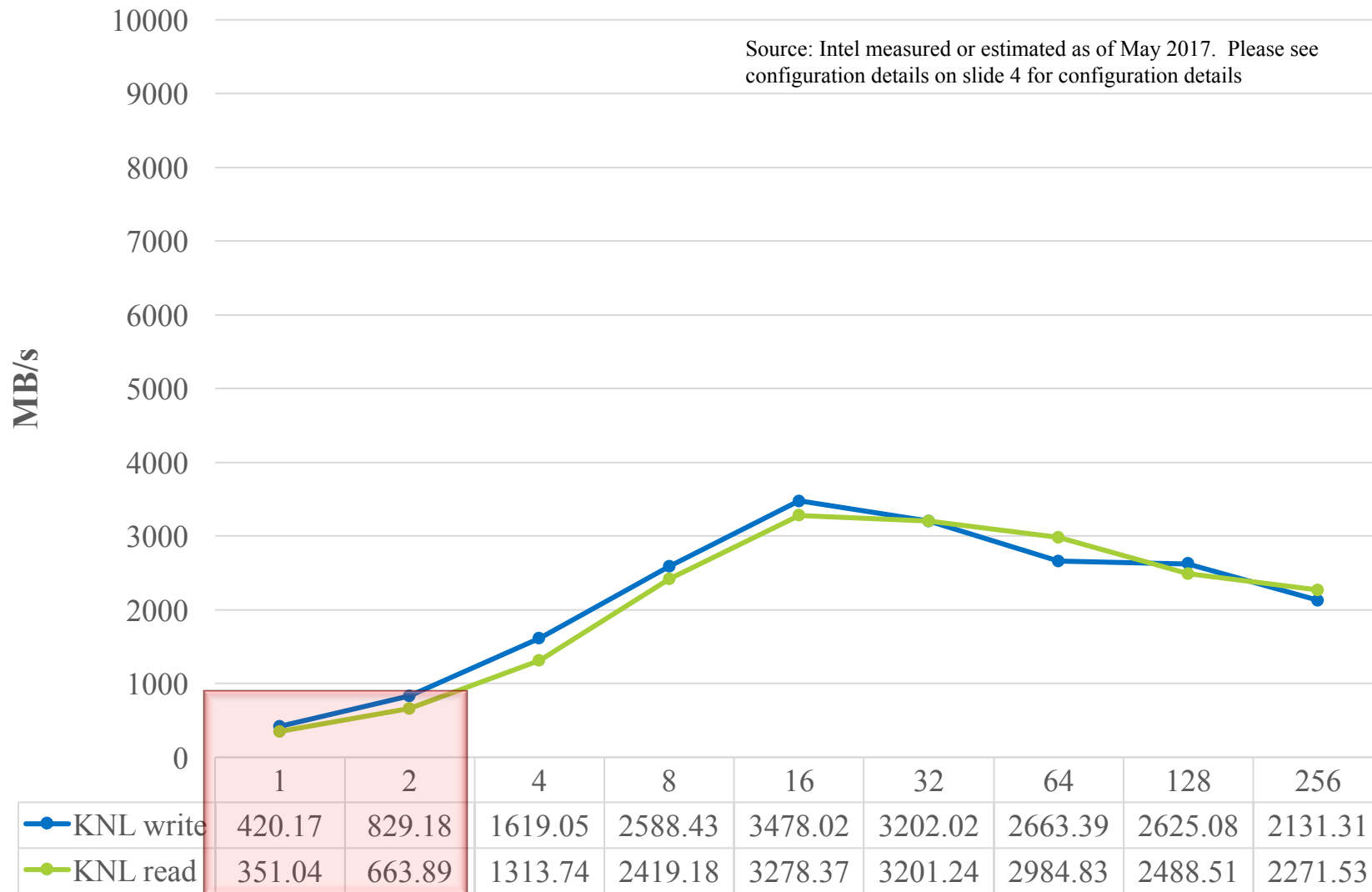
- L2 cache: 1024K
- CPUs: 272, Cores: 68, Threads per core: 4
- Intel® OPA HFI Silicon 100 Series [discrete]

Intel® Xeon® E5-2697 IOR results current version of Lustre*



* Some names and brands may be claimed as the property of others.

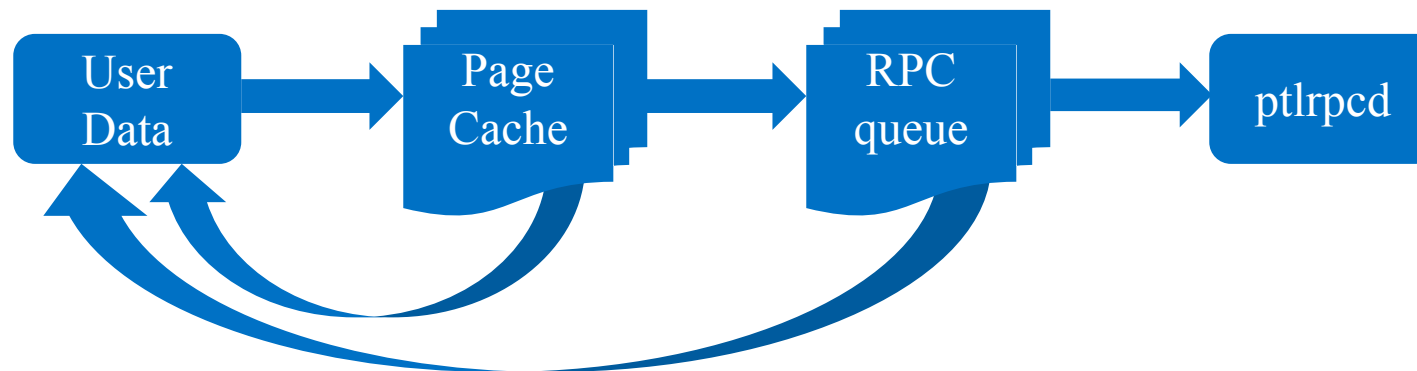
Intel® Xeon Phi™ 7250 IOR results current version of Lustre*



* Some names and brands may be claimed as the property of others.

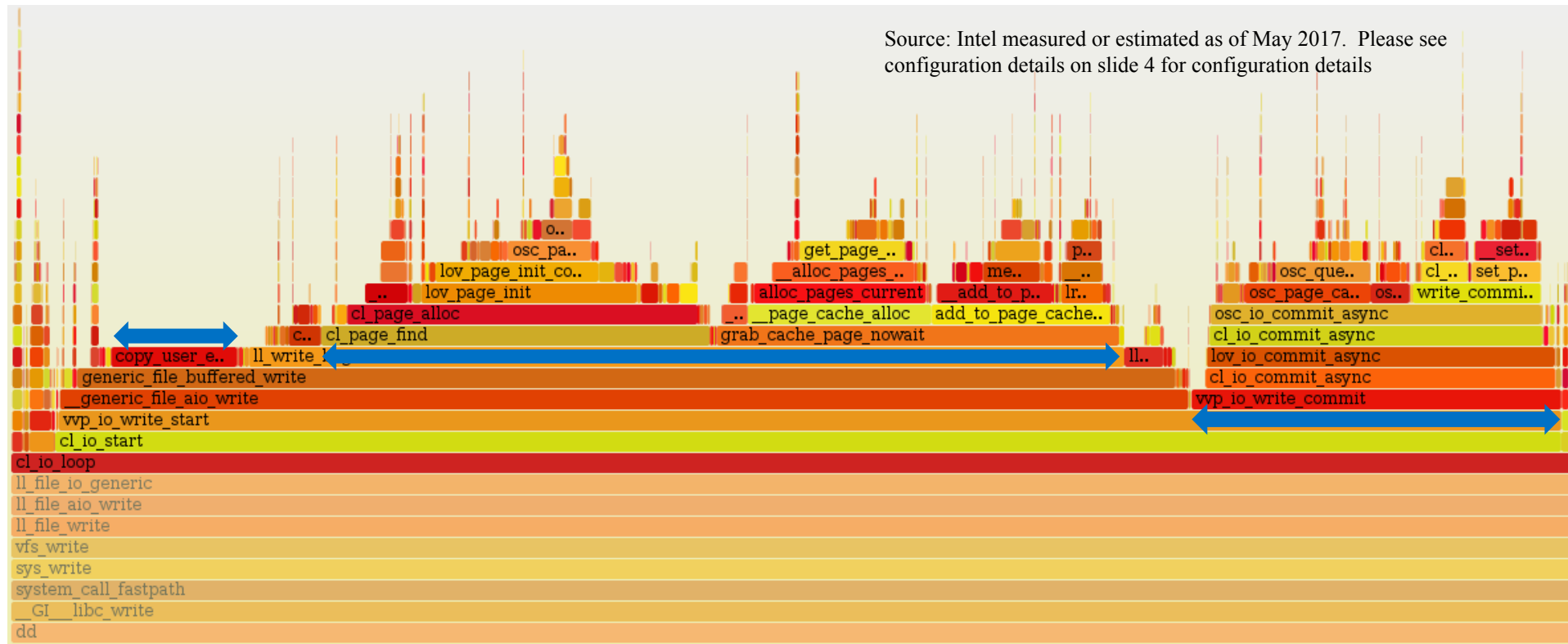
Write operation data flow

- For each PAGE_SIZE of user data:
 - Allocate new page in page cache
(`cl_page_alloc()`, `grab_cache_page()`)
 - Copy data from user space to page cache
(`copy_user_enhanced_fast_string()`)
- Submit pages to RPC queue (`vvp_io_write_commit()`)
- Wake up `ptlrpcd` thread and wait for completion



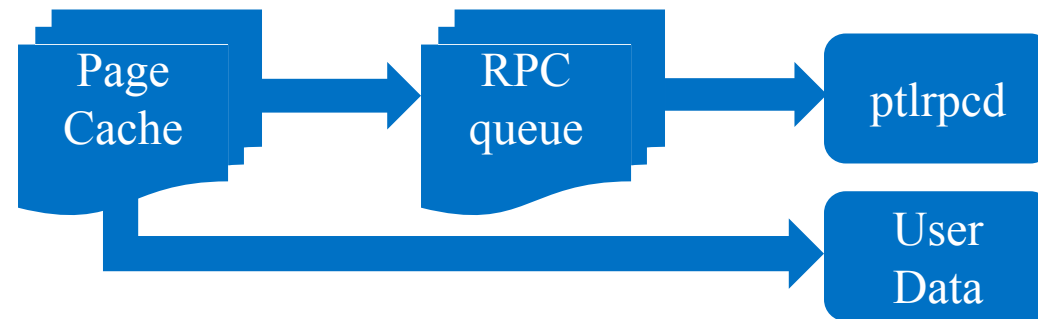
Flame Graph of Write operation

- Page cache allocation is $\sim 1/2$ of all write time
- Submit pages to RPC queue
- Copy user data



Read operation data flow

- Read Ahead (`ll_readahead()`) for each page:
 - Allocate new page in page cache (`cl_page_alloc()`, `grab_cache_page()`)
 - Submit page to RPC queue for read (`cl_io_submit_rw()`)
- Wake up `ptlrpcd` thread and wait for completion
- ...
- Copy data from page cache to user space (`copy_user_enhanced_fast_string()`)



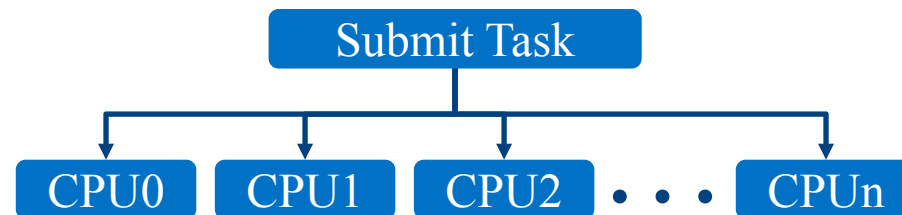
Ideas of Improvements

- Split time consuming operations into several threads
 - Page cache allocation in parallel
 - Submit pages into several RPC queues in parallel
 - Copy user data in parallel**
- Move Read Ahead into asynchronous threads
 - Resume main read thread as soon as possible
 - Avoid cost of unused reads (false read ahead)
- Utilize more ptlrpcd threads
 - Process network packets in parallel**
- Multiple network connection
 - Use multiple IB/OPA endpoints between nodes

** *Very important for less performant cores.*

Introduce new Parallel Tasks API

- Farm tasks out to be done in parallel on multiple CPUs
- Use system thread pool (don't create new thread pool)
- Provide ordered or disordered task completion
- Submit for execution in Round Robin fashion

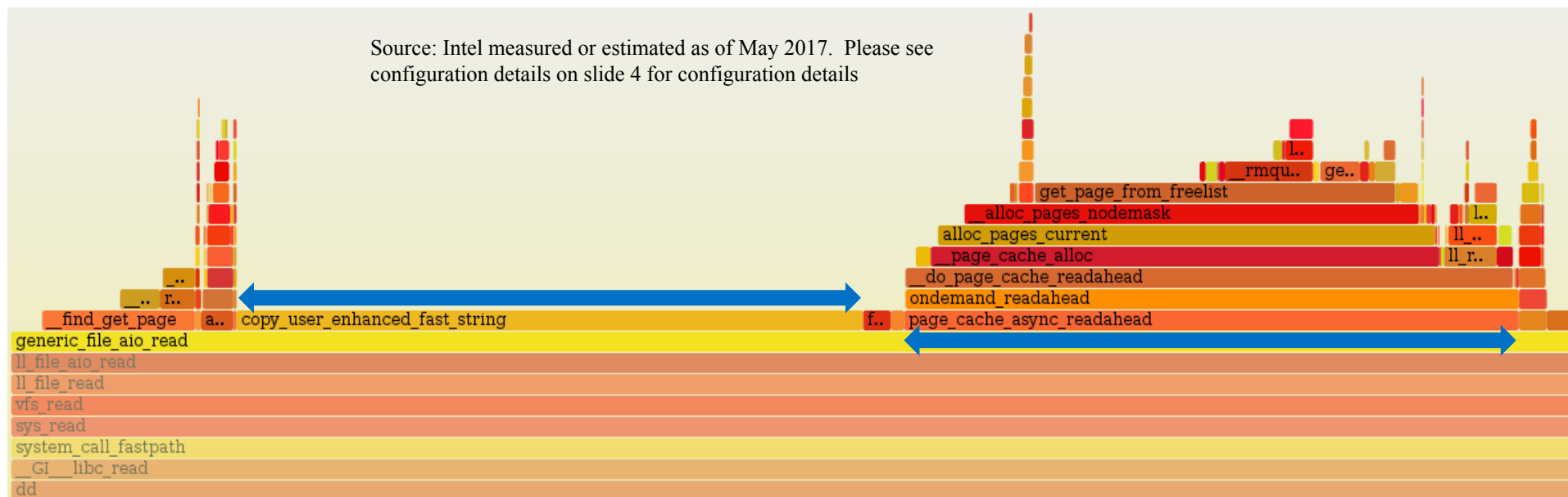


Multiple IB/OPA Endpoints

- Lustre* use a single network channel between nodes
 - Intel® OPA fabric use a host CPU for packets processing, so on machine with slow cores it cannot utilize full network bandwidth
- Optimization of IB/OPA LND driver
 - create multiple endpoints and balance the traffic over them
- <https://jira.hpdd.intel.com/browse/LU-8943>

Flame Graph of Read with async RA

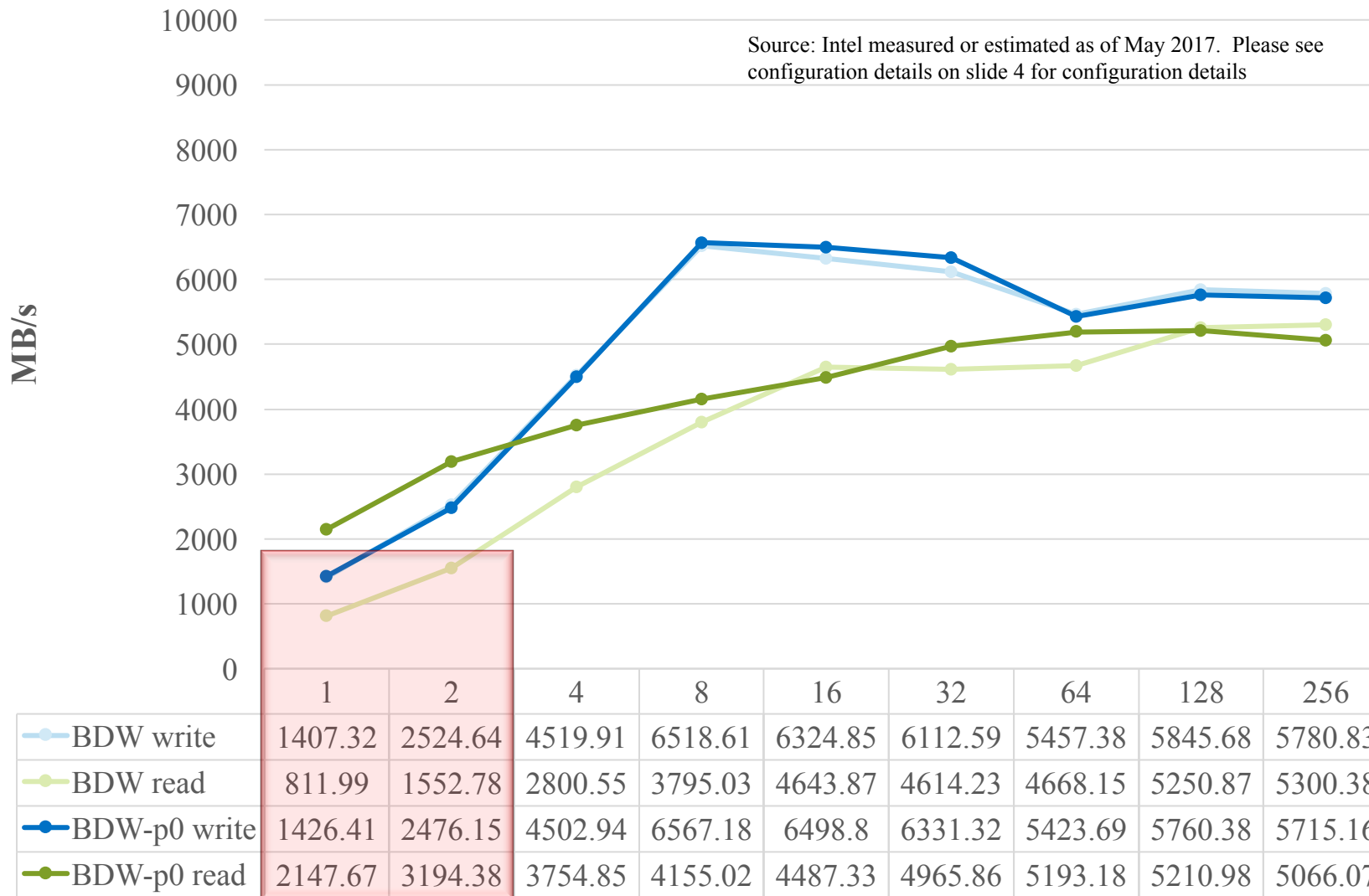
- Copy user data is $\sim 1/2$ of all read time
- Initial Read Ahead



<https://jira.hpdd.intel.com/browse/LU-8964>

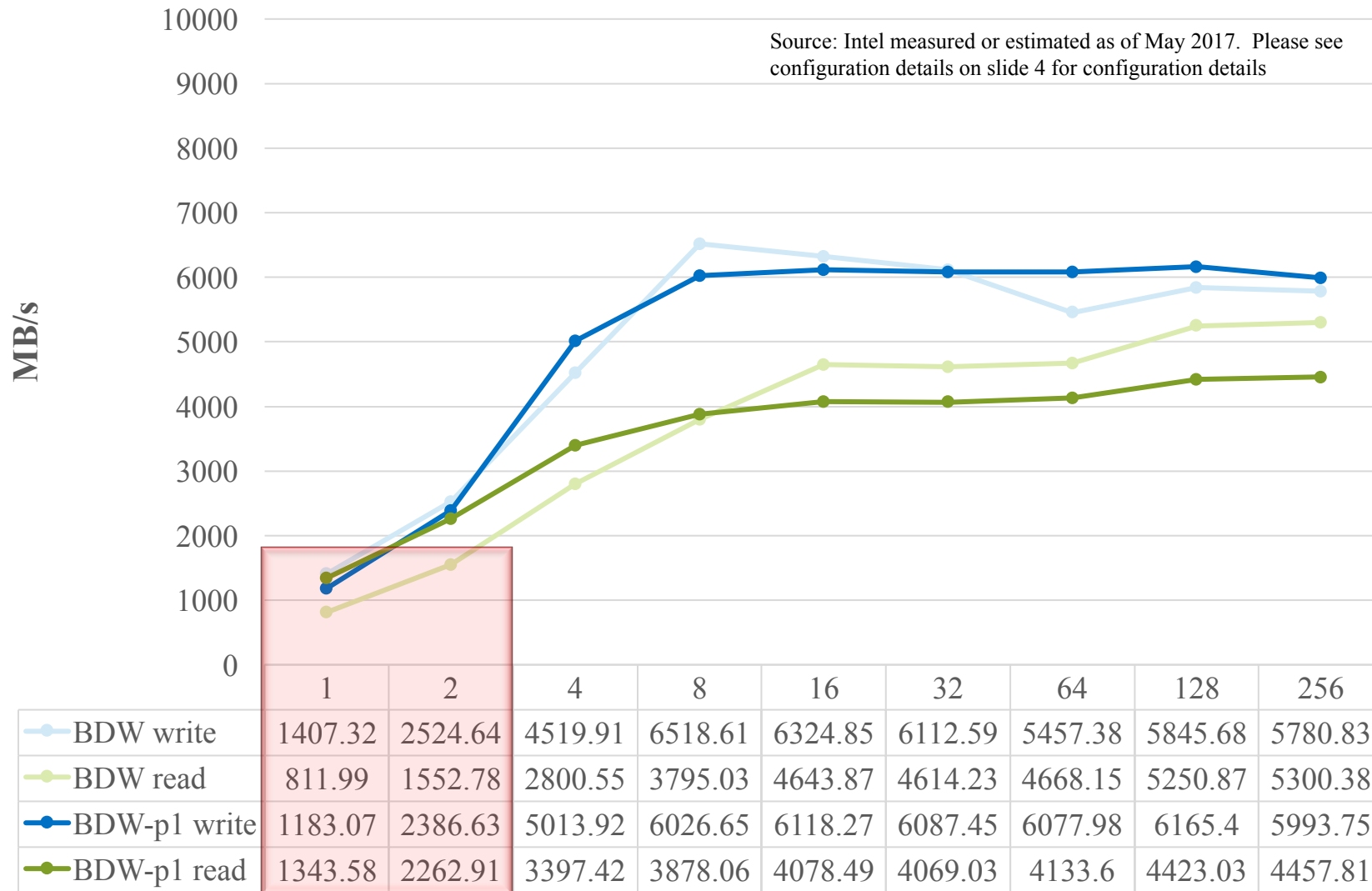
Intel® Xeon® E5-2697 IOR results

Parallel I/O off, async Read Ahead



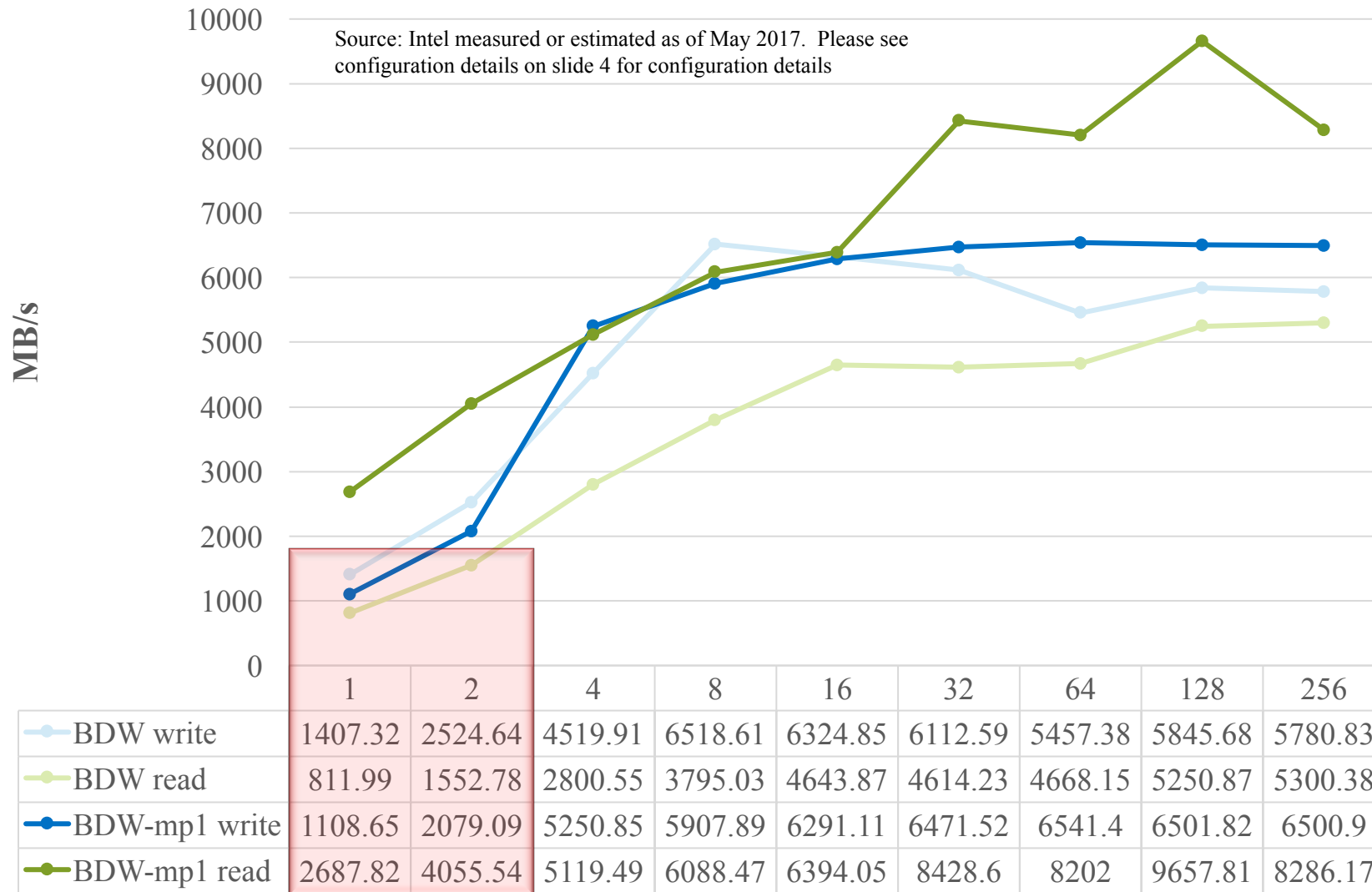
Intel® Xeon® E5-2697 IOR results

Parallel I/O on, async Read Ahead



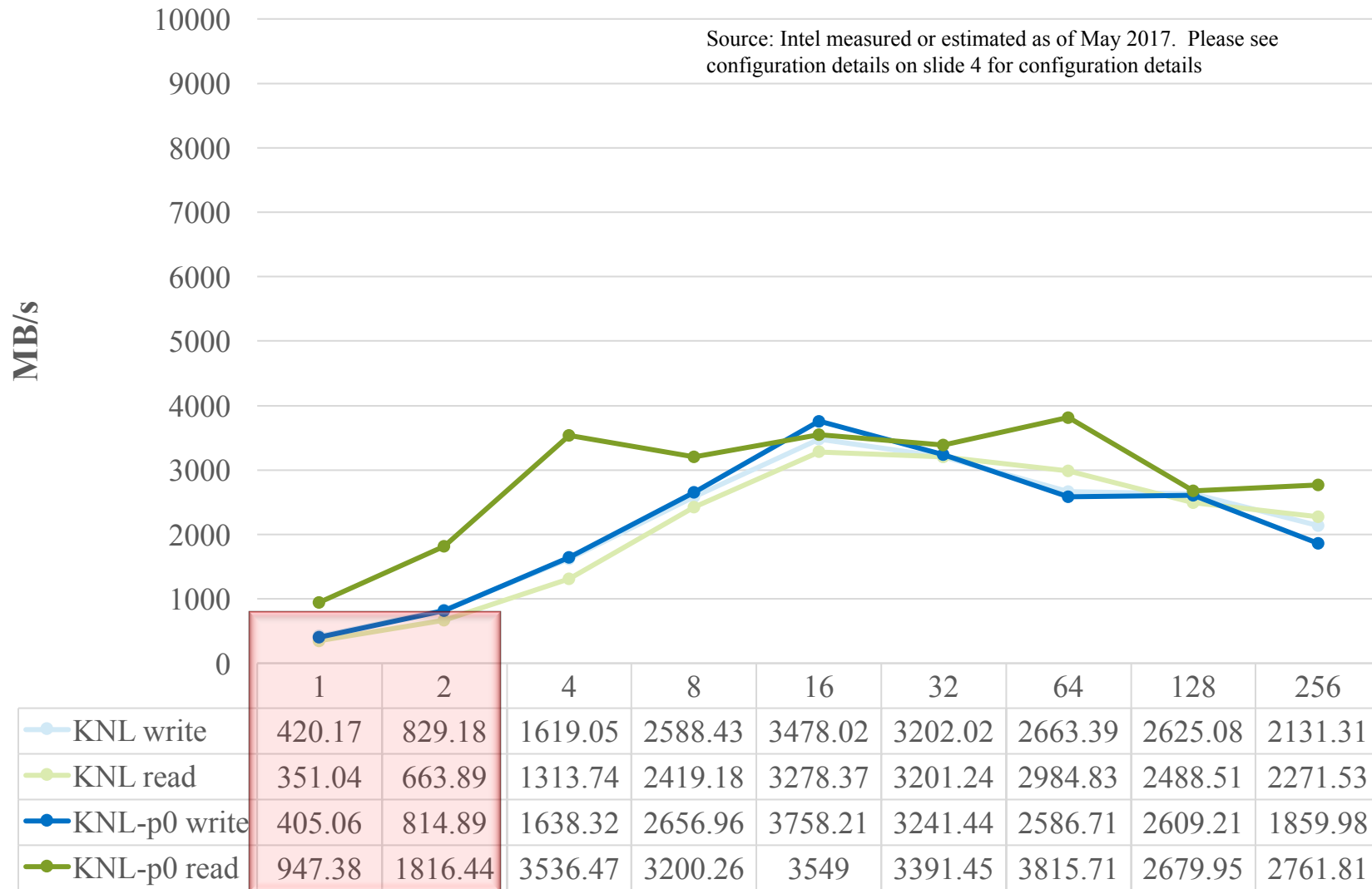
Intel® Xeon® E5-2697 IOR results

Parallel I/O on, async Read Ahead, MQP



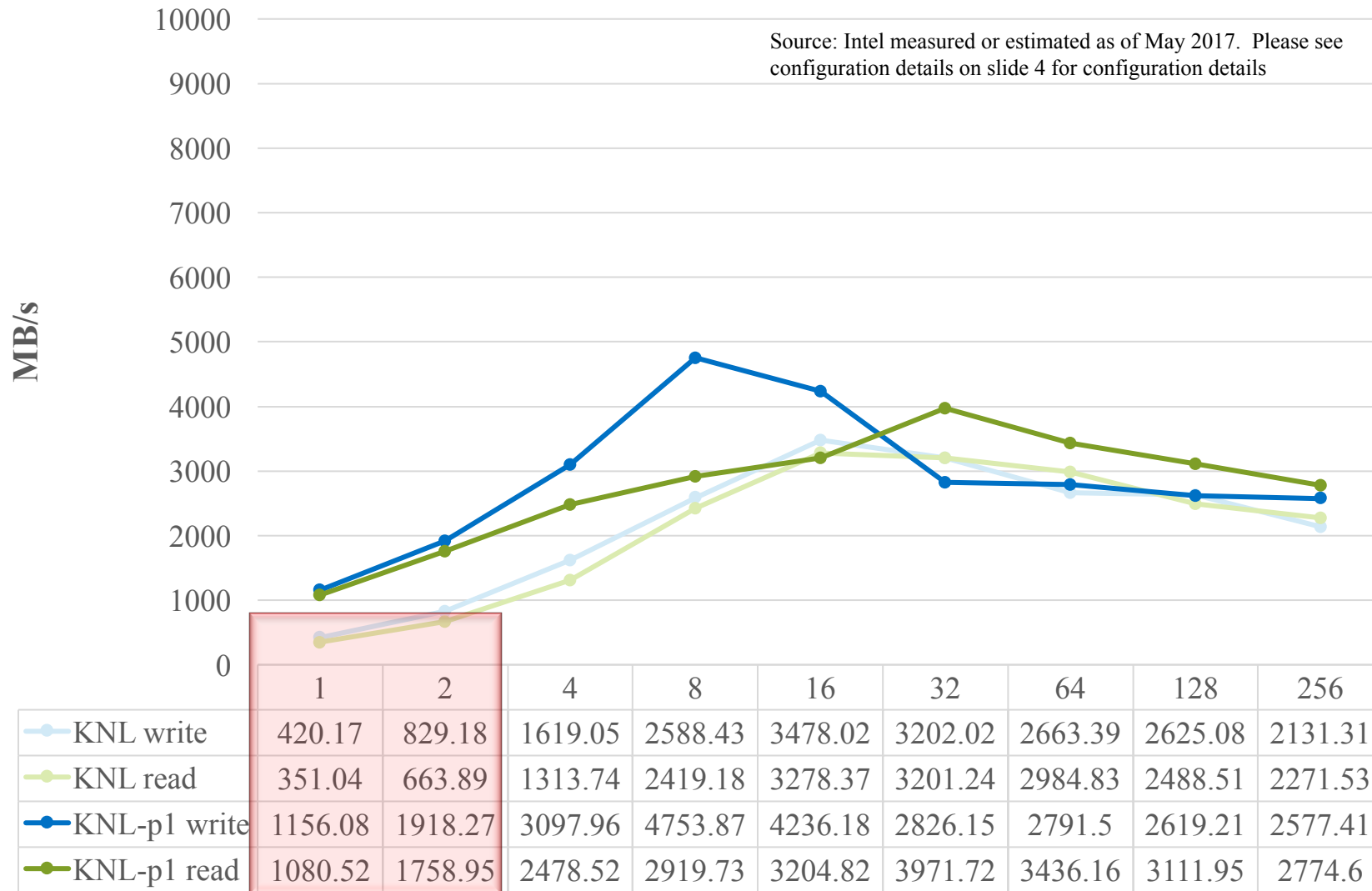
Intel® Xeon Phi™ 7250 IOR results

Parallel I/O off, async Read Ahead



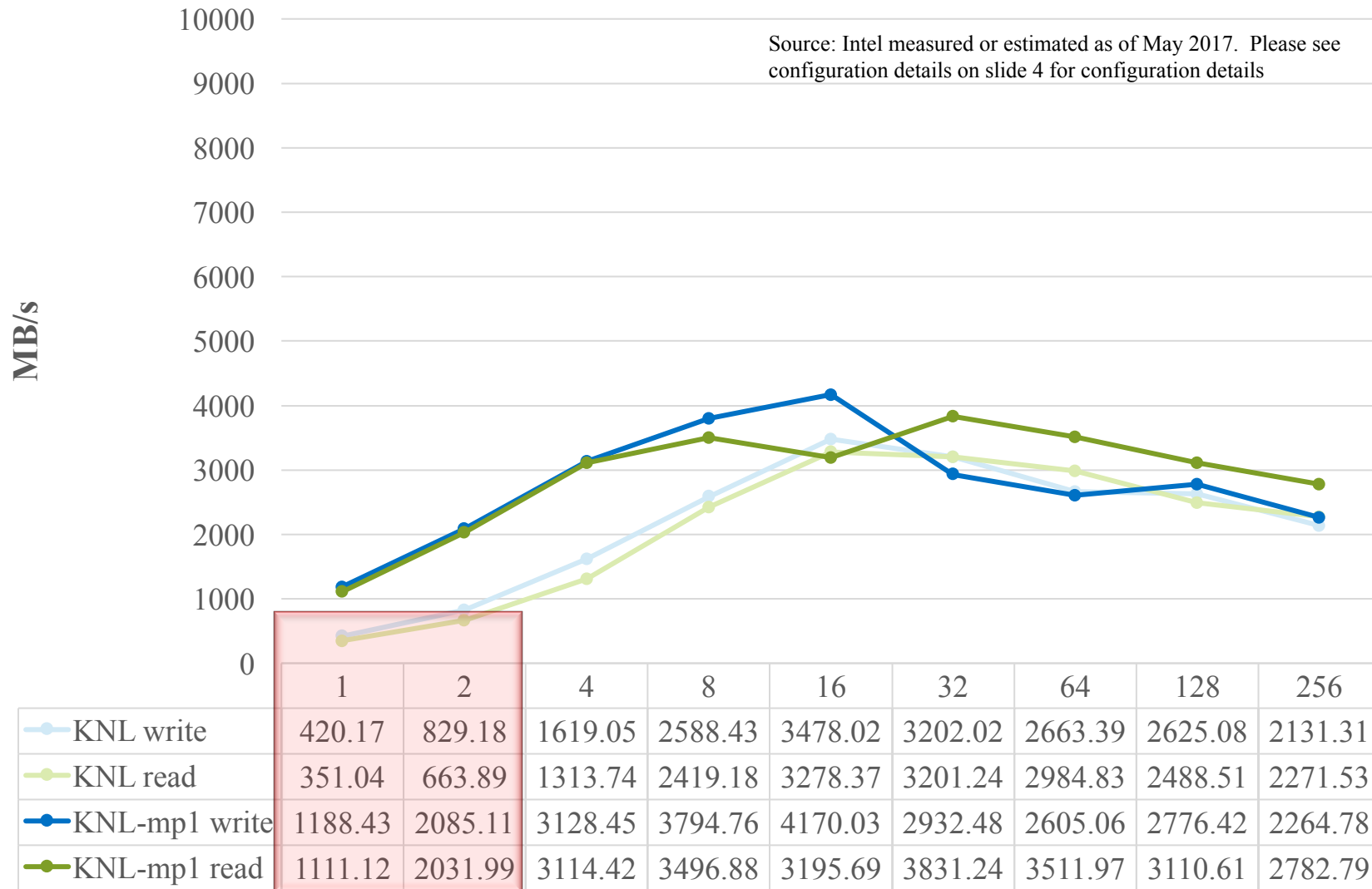
Intel® Xeon Phi™ 7250 IOR results

Parallel I/O on, async Read Ahead



Intel® Xeon Phi™ 7250 IOR results

Parallel I/O on, async Read Ahead, MQP



Hybrid Coordinate Ocean Model (HYCOM) application

Source: Intel measured or estimated as of May 2017. Please see configuration details on slide 4 for configuration details

Testing a real world application - see <https://hycom.org>

Current version of Lustre*

zaiio**	calls =	6	time =	0.39432	time/call =	0.06572032
zaiord	calls =	99	time =	2.24527	time/call =	0.02267950
zaiowr	calls =	726	time =	81.19858	time/call =	0.11184378
total	calls =	1	time =	5203.09187	time/call =	5203.09187198

Parallel I/O on, async Read Ahead, MQP

zaiio**	calls =	6	time =	0.31822	time/call =	0.05303649
zaiord	calls =	99	time =	3.48087	time/call =	0.03516028
zaiowr	calls =	726	time =	65.78485	time/call =	0.09061274
total	calls =	1	time =	5149.86002	time/call =	5149.86001992

* Some names and brands may be claimed as the property of others.

Summary

- Using parallel I/O and multiple IB/OPA endpoints
Lustre* now can utilize performance benefits of
 - Multi-cores systems even for single-threaded applications
 - Multiple I/O operations in Lustre even if a single core is less performant
 - Intel® OPA fabric which uses a host CPU for packets processing
- Using asynchronous Read Ahead Lustre now
 - Doesn't block the userspace I/O operation to do Read Ahead
- For some workloads up to 2x speed-up in reads/writes has been observed
- Write work expected to be in 2.10.0; read work in 2.10.x maintenance release

Notices and Disclaimers

No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.

Intel disclaims all express and implied warranties, including without limitation, the implied warranties of merchantability, fitness for a particular purpose, and non-infringement, as well as any warranty arising from course of performance, course of dealing, or usage in trade.

This document contains information on products, services and/or processes in development. All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest forecast, schedule, specifications and roadmaps.

The products and services described may contain defects or errors known as errata which may cause deviations from published specifications. Current characterized errata are available on request.

Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. Consult other sources of information to evaluate performance as you consider your purchase. For more complete information about performance and benchmark results, visit www.intel.com/benchmarks

© Intel Corporation. Intel, Intel Inside, the Intel logo, Xeon, Intel Xeon Phi, Intel Xeon Phi logos and Xeon logos are trademarks of Intel Corporation or its subsidiaries in the United States and/or other countries.

Questions?



