



Optimizing Apache Spark* throughput using Intel® Memory Drive Technology

Ravi Durgavajhala
SSD Solutions Architect
Non-Volatile Memory Solutions Group (NSG), Intel

Ravikanth.Durgavajhala@Intel.com
[@ravykanth](https://twitter.com/ravykanth)

#HWCSAIS1

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

LEGAL NOTICES AND DISCLAIMERS

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. **No computer system can be absolutely secure.** Check with your system manufacturer or retailer or learn more at intel.com

Intel may make changes to specifications and product descriptions at any time, without notice. Designers must not rely on the absence or characteristics of any features or instructions marked "reserved" or "undefined". Intel reserves these for future definition and shall have no responsibility whatsoever for conflicts or incompatibilities arising from future changes to them. All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest Intel product specifications and roadmaps.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors.

The products described in this document may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Contact your local Intel sales office or your distributor to obtain the latest specifications and before placing your product order.

This document contains information on products in the design phase of development.

The benchmark results may need to be revised as additional testing is conducted. The results depend on the specific platform configurations and workloads utilized in the testing, and may not be applicable to any particular user's components, computer system or workloads. The results are not necessarily representative of other benchmarks and other benchmark results may show greater or lesser impact from mitigations.

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. Results have been estimated or simulated based on internal Intel analysis and are provided for informational purposes only. Any difference in system hardware or software design or configuration may affect actual performance.

Cost reduction scenarios described are intended as examples of how a given Intel- based product, in the specified circumstances and configurations, may affect future costs and provide cost savings. Circumstances will vary. Intel does not guarantee any costs or cost reduction.

Intel does not control or audit the design or implementation of third party benchmark data or Web sites referenced in this document. Intel encourages all of its customers to visit the referenced Web sites or others where similar performance benchmark data are reported and confirm whether the referenced benchmark data are accurate and reflect performance of systems available for purchase.

Intel, Intel Xeon, Intel Optane, and 3D XPoint are trademarks of Intel Corporation in the U.S. and/or other countries.

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

Copyright © 2018 Intel Corporation. All rights reserved.

Problem Statement

Optimize the performance of Spark* and get more out of my infrastructure while operating within the budget.

Assumptions

- **Extrapolate overall infrastructure set up.**
- **Match the individual system resources to that of real-world production, as much as possible.**
- **Come up with a representative workload.**
- **Identify a solution along with alternatives.**

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

A quick overview of the K-Means workload

“Definition: K-Means clustering aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean.

Standard Algorithm: “Given an initial set of k means $m_1(1), \dots, m_k(1)$, the algorithm proceeds by alternating between two steps:

1. **Assignment step:** Assign each observation to the cluster whose mean has the least squared **Euclidean distance**, this is intuitively the “nearest” mean.
2. **Update step:** Calculate the new means to be the centroids of the observations in the new clusters.

The algorithm has converged when the assignments no longer change.”

https://en.wikipedia.org/wiki/K-means_clustering

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

Hardware Configuration

Master Node		Data Node (x3)
CPU	Intel® Xeon® Gold 6140 CPU @ 2.30GHz	Intel® Xeon® Gold 6140 CPU @ 2.30GHz
Cores per Socket	18	18
Sockets	2	2
Threads per Core	2	2
Total vcores	72	72
Memory	192GB	192GB
SSD	None	3.7TB Intel® SSD DC P4500 (x2)
		375GB Intel® Optane™ SSD DC P4800X (x2)
Network	10Gbps	

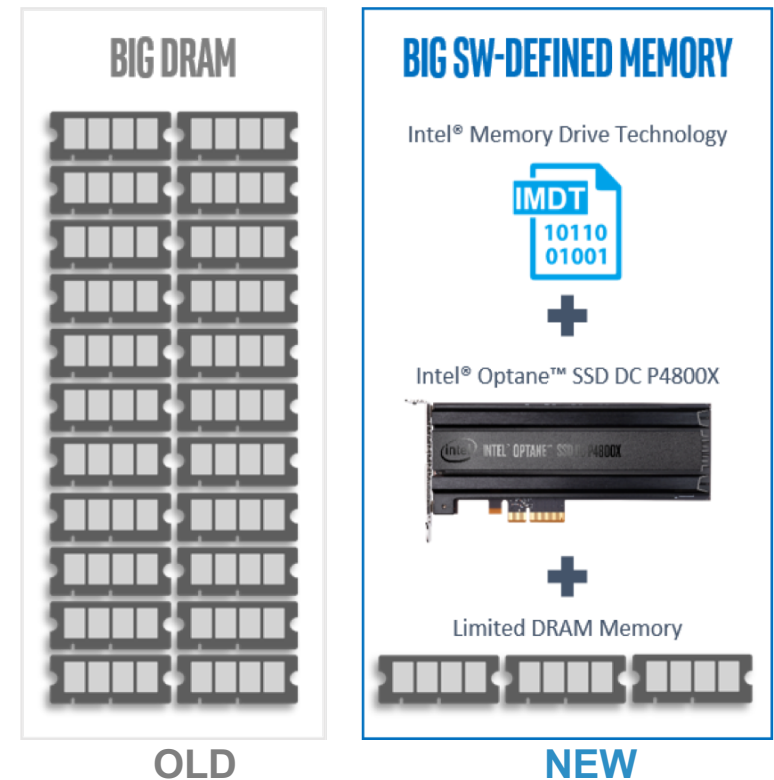
Software Configuration

Stack	Version
Distribution	HDP 2.6.4.0
HDFS*	2.7.3
YARN*	2.7.3
Spark*	2.2.0
OS	CentOS 7.4*
Kernel	4.14.16

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

Introducing Intel® Memory Drive Technology (IMDT)

- Intel® Optane™ Technology - Write in place, Bit addressable, Low latency.
- Use Intel® Optane™ SSD DC P4800X transparently as memory.
- Grow beyond system DRAM capacity, or replace high-capacity DIMMs for lower-cost alternative, with similar performance.
- Leverage storage-class memory today!
 - No change to software stack: unmodified Linux* OS, applications, and programming.
 - No change to hardware: runs bare-metal, loaded before OS from BIOS or UEFI.



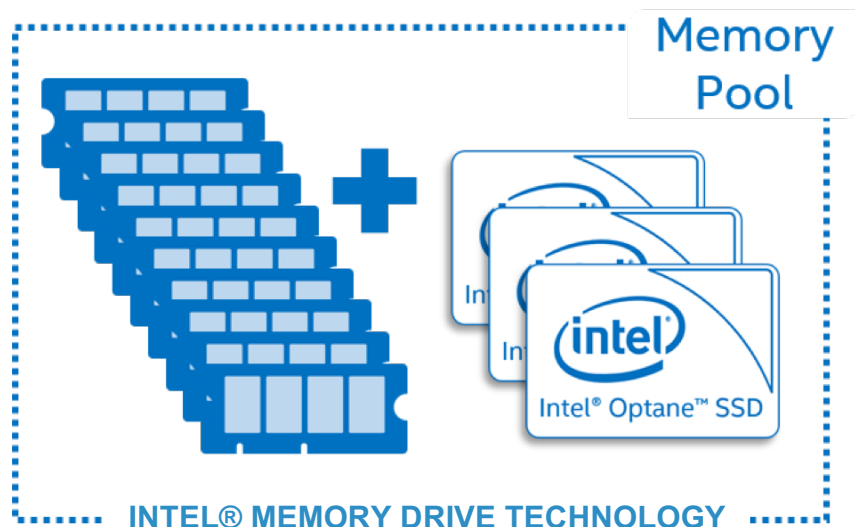
*Other names and brands may be claimed as the property of others

Intel® Memory Drive Technology delivers big, affordable memory

use case

1

EXPAND beyond limited DRAM capacity

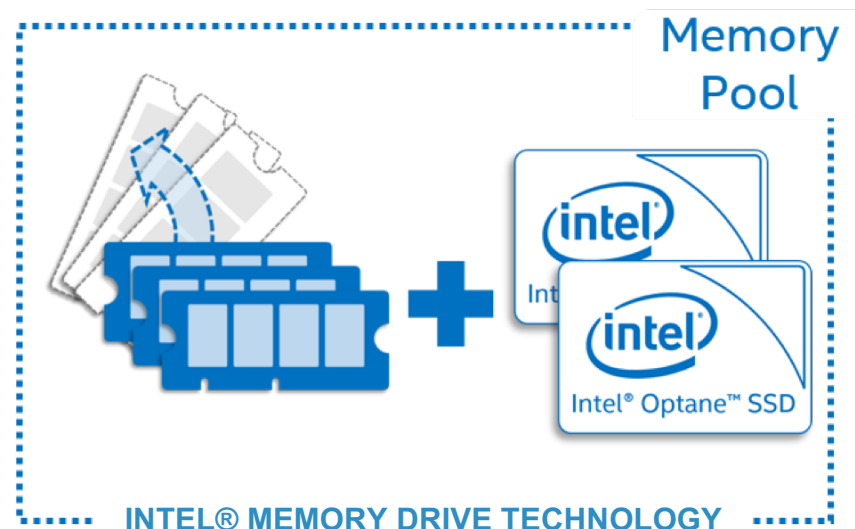


**Expand Insights with
Massive Data Pools**

use case

2

DISPLACE DRAM with affordable SSDs



**Reduce High-capacity DRAM
CAPEX Expenditures**

Note: Intel® Memory Drive Technology supports Linux* x86_64 (64-bit), kernels 2.6.32 or newer.

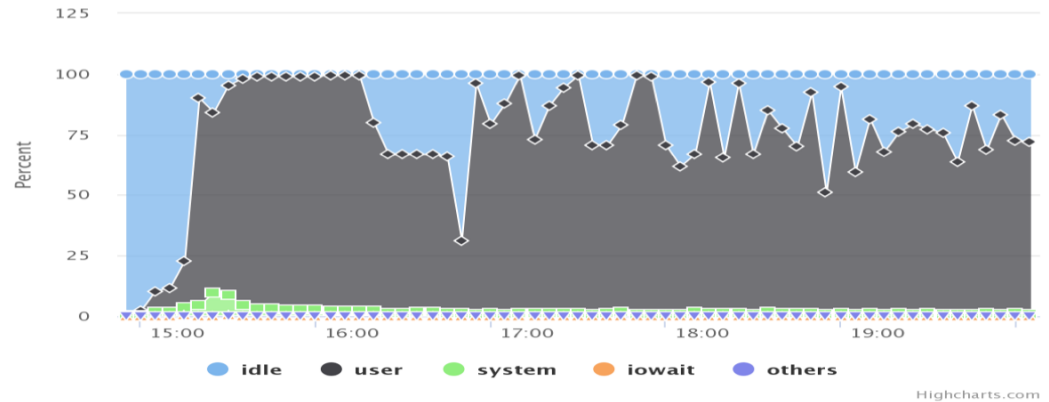
*Other names and brands may be claimed as the property of others

Workload that fits entirely into DRAM

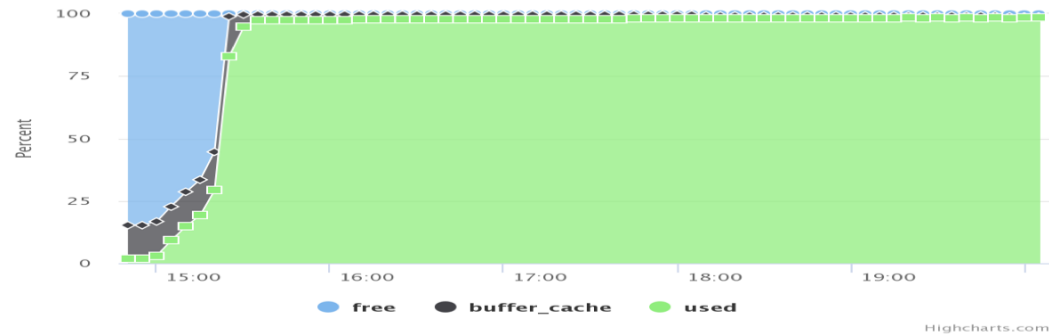
Spark* Workload Configuration	
# of Executors across all Nodes	42
# of Cores per Executor	5
Memory per Executor	12 GiB
Memory Overhead per Executor	3 GiB
Driver Memory	1 GiB
Driver Memory Overhead	1 GiB
K-Means workload Scale Factor	1.2 Billion samples
Time taken to run the workload is 5.3 min ¹	

- Spark* configuration is based on generally understood guidelines.
- Data set fits entirely into memory, without any spill.
- The objective is to utilize maximum available resources on the system to get best possible run-time.

Summarized CPU usage



Summarized Memory usage



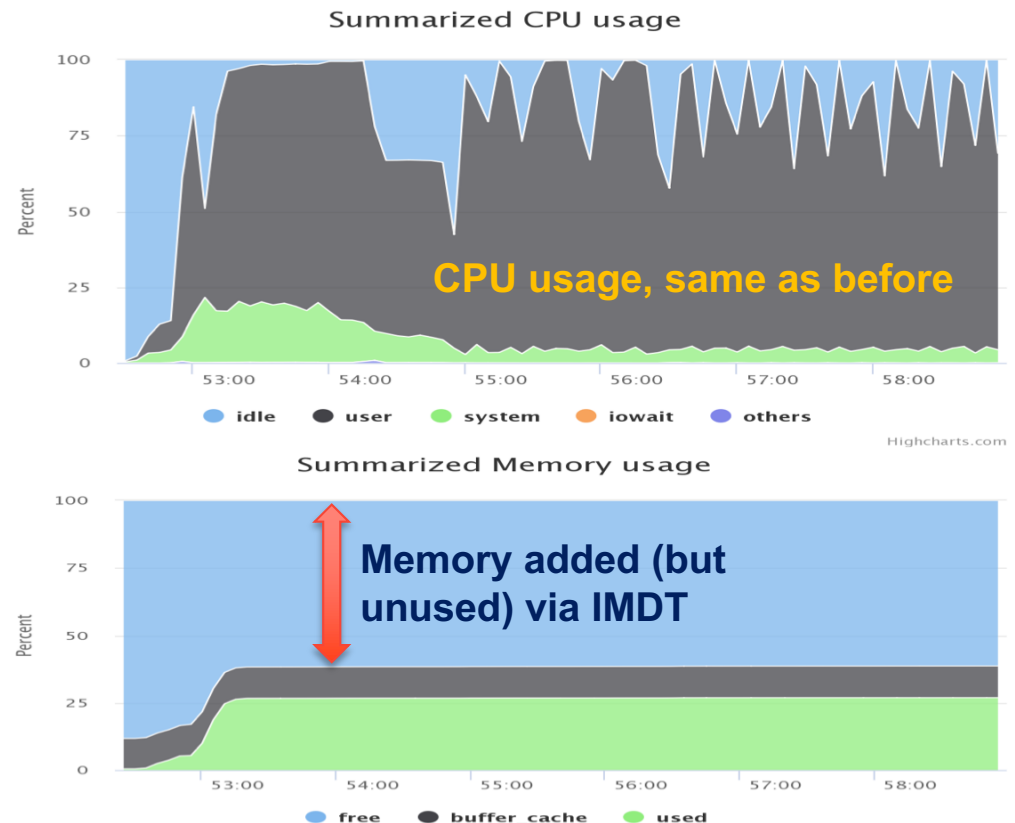
¹For system configuration details, please refer to Slide #5. Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

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

Workload that fits entirely into DRAM (+IMDT)

Spark* Workload Configuration	
# of Executors across all Nodes	42
# of Cores per Executor	5
Memory per Executor	12 GiB
Memory Overhead per Executor	3 GiB
Driver Memory	1 GiB
Driver Memory Overhead	1 GiB
K-Means workload Scale Factor	1.2 Billion samples
Time taken to run the workload is 5.3 min ¹	

- Objective is to ensure performance did not get impacted when running the same workload using same resource configuration, except for memory expansion using IMDT.



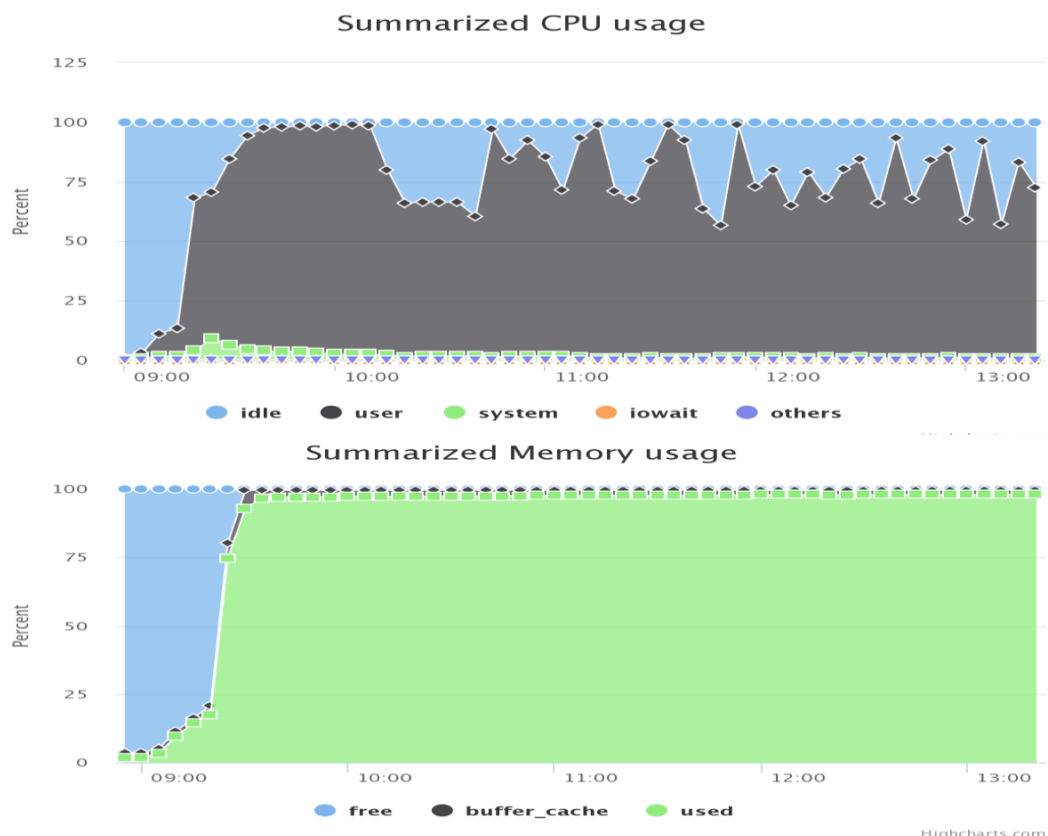
¹For system configuration details, please refer to Slide #5. Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

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

Workload that fits entirely into DRAM – fine tuned

Spark* Workload Configuration	
# of Executors across all Nodes	30
# of Cores per Executor	7
Memory per Executor	17 GiB
Memory Overhead per Executor	3 GiB
Driver Memory	1 GiB
Driver Memory Overhead	1 GiB
K-Means workload Scale Factor	1.2 Billion samples
Time taken to run the workload is 4.5 min ¹	

- Spark* configuration is fine tuned based on Memory and CPU utilization.
- Not all workloads are alike, so each workload needs to be custom-adjusted for better resource utilization.



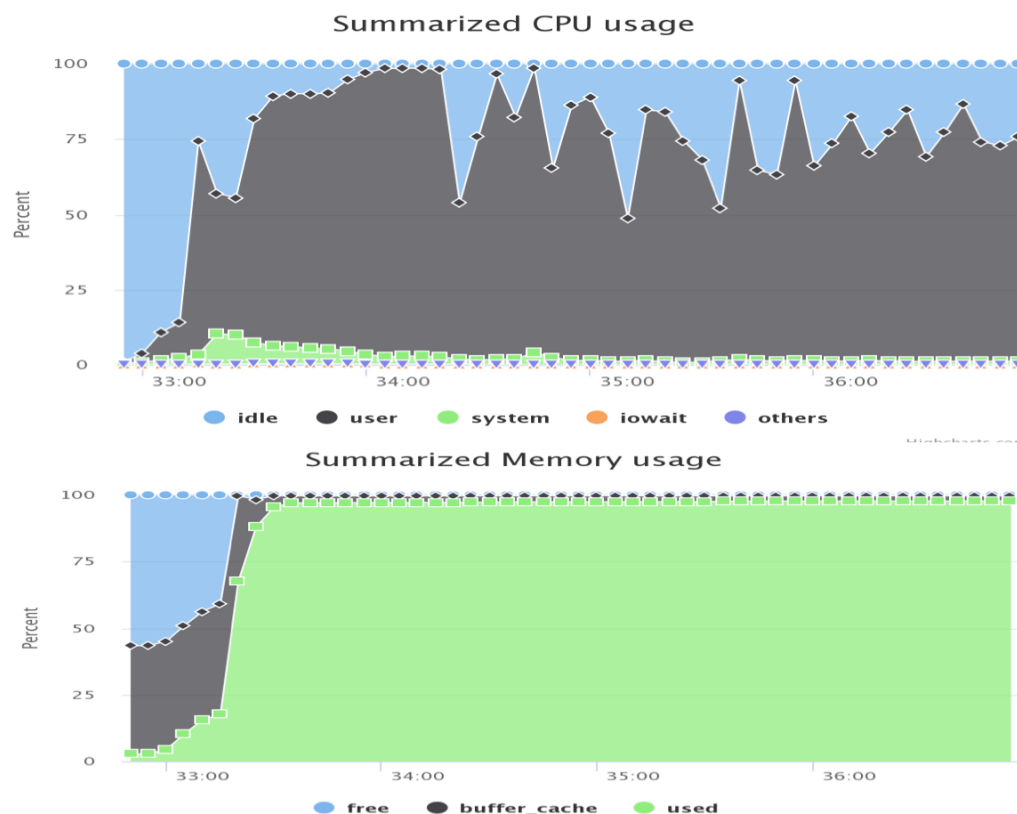
¹For system configuration details, please refer to Slide #5. Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

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

Workload that fits entirely into DRAM – fine tuned

Spark* Workload Configuration	
# of Executors across all Nodes	30
# of Cores per Executor	5
Memory per Executor	17 GiB
Memory Overhead per Executor	3 GiB
Driver Memory	1 GiB
Driver Memory Overhead	1 GiB
K-Means workload Scale Factor	1.2 Billion samples
Time taken to run the workload is 4.1 min ¹	

- Utilizing max of resources available does not always yield best possible performance.
- Performance varies based on memory and other resource utilization within the application code.



¹For system configuration details, please refer to Slide #5. Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

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

Bigger Workload using DRAM

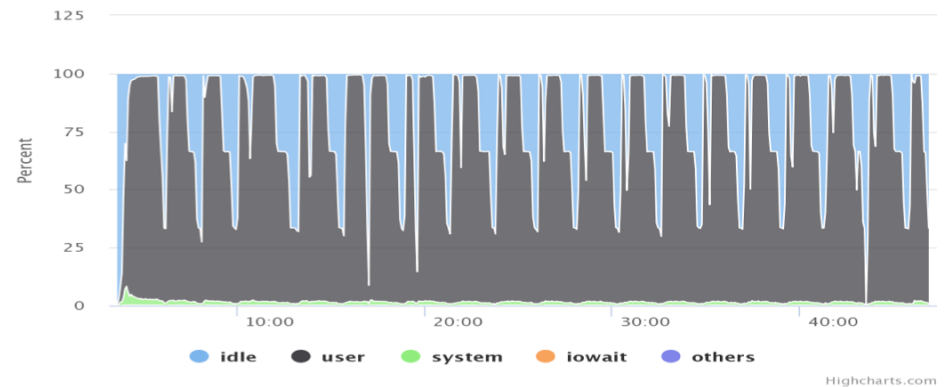
Spark* Workload Configuration	
# of Executors across all Nodes	30
# of Cores per Executor	7
Memory per Executor	17 GiB
Memory Overhead per Executor	3 GiB
Driver Memory	1 GiB
Driver Memory Overhead	1 GiB
K-Means workload Scale Factor	2 Billion samples
Time taken to run the workload is 43 min ¹	

- Spark* shuffles the data between memory and storage when dataset does not fit entirely in memory.
- If the workload is large enough that it cannot fit with fully populated memory channel, the next logical move is to scale out and add more nodes.
- Storage: 2x Intel® Optane® SSD DC P4800X (375GB)

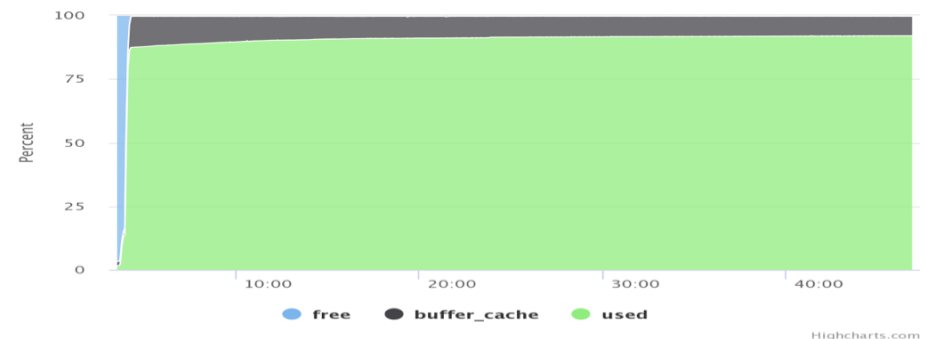
¹For system configuration details, please refer to Slide #5. Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

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

Summarized CPU usage



Summarized Memory usage



Bigger Workload using IMDT

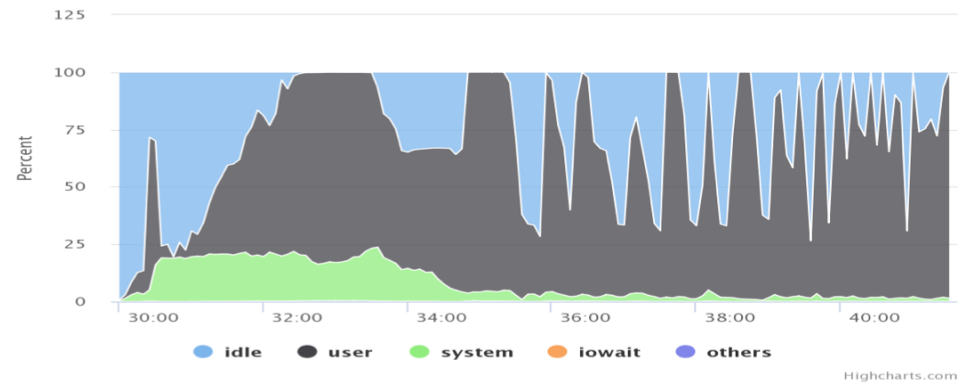
Spark* Workload Configuration	
# of Executors across all Nodes	42
# of Cores per Executor	10
Memory per Executor	40 GiB
Memory Overhead per Executor	3 GiB
Driver Memory	1 GiB
Driver Memory Overhead	1 GiB
K-Means workload Scale Factor	2 Billion samples
Time taken to run the workload is 12 min ¹	

- IMDT helps to bring more memory resources without having to scale out.
- IMDT can expand memory capacity to grow x8 beyond system spec.
- That directly translates to more Spark* executors that can run in parallel.

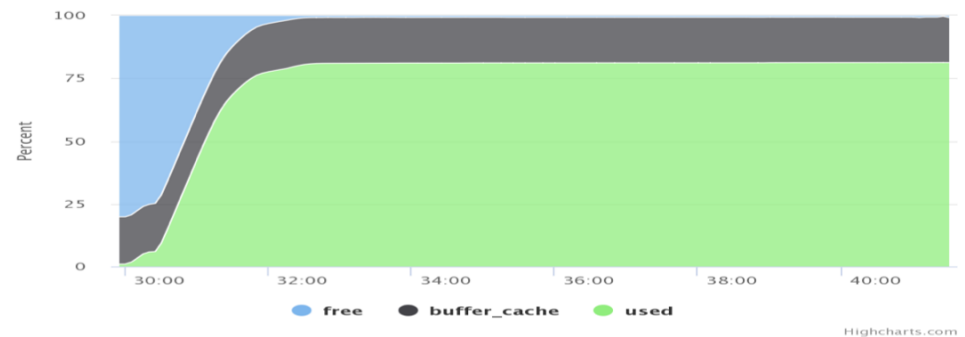
¹For system configuration details, please refer to Slide #5. Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

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

Summarized CPU usage

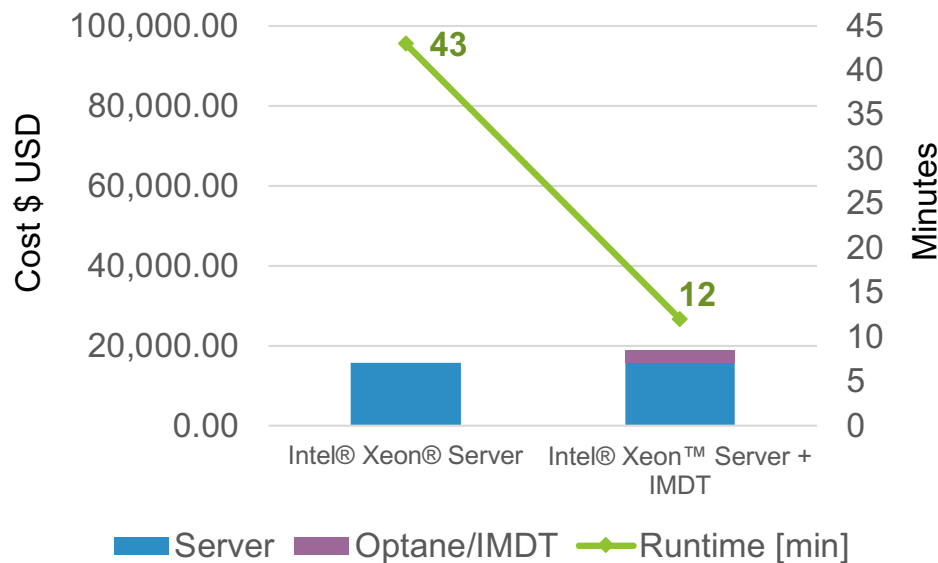


Summarized Memory usage



Solution Economics

Per-node Configuration Cost Comparison



Master Node		Data Node (x3)
CPU	Intel® Xeon® Gold 6140 CPU @ 2.30GHz	
Cores/Socket	18	
Sockets	2	
Threads per Core	2	
Total vcores	72	
Memory	192GB	
SSD	None	3.7TB Intel® SSD DC P4500 (x2)
		375GB Intel® Optane™ SSD DC P4800X (x2)
Network	10Gbps	

20% added cost¹ → reduce runtime by factor of x3.5²

¹ Cost estimates based on quote from Colfax International as of May 27, 2018

² For system configuration details, please refer to Slide #5. Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

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

Bigger Workload using IMDT and fewer nodes

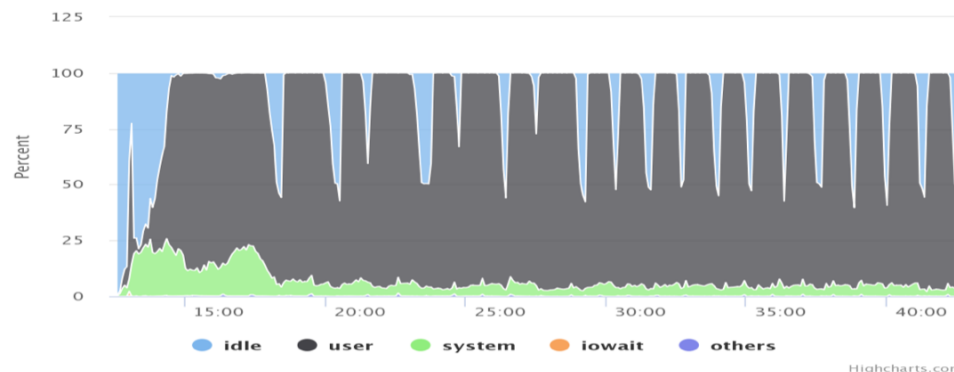
Spark* Workload Configuration (2 Data Nodes only)	
# of Executors across all Nodes	28
# of Cores per Executor	10
Memory per Executor	40 GiB
Memory Overhead per Executor	3 GiB
Driver Memory	1 GiB
Driver Memory Overhead	1 GiB
K-Means workload Scale Factor	2 Billion samples
Time taken to run the workload is 30 min ¹	

- For workloads that are not fully utilizing CPU resources in a given infrastructure, IMDT can help increase CPU utilization.
- Increasing CPU utilization allows for savings on data center footprint by reducing node-count, with larger memory per node.
- Savings can be put back into improved networks, higher-core-count CPUs, etc.

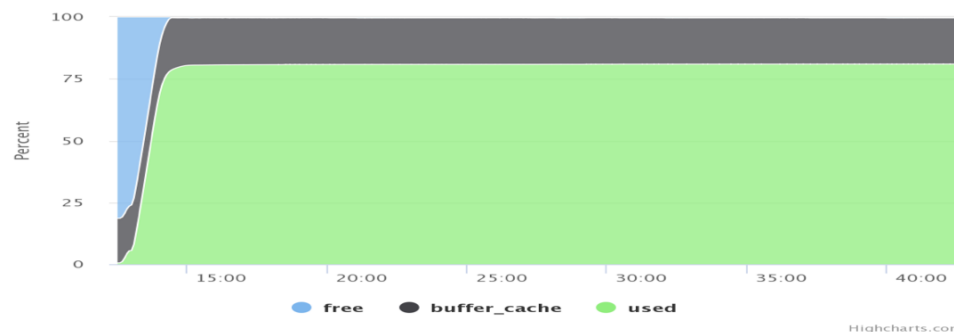
¹For system configuration details, please refer to Slide #5. Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

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

Summarized CPU usage

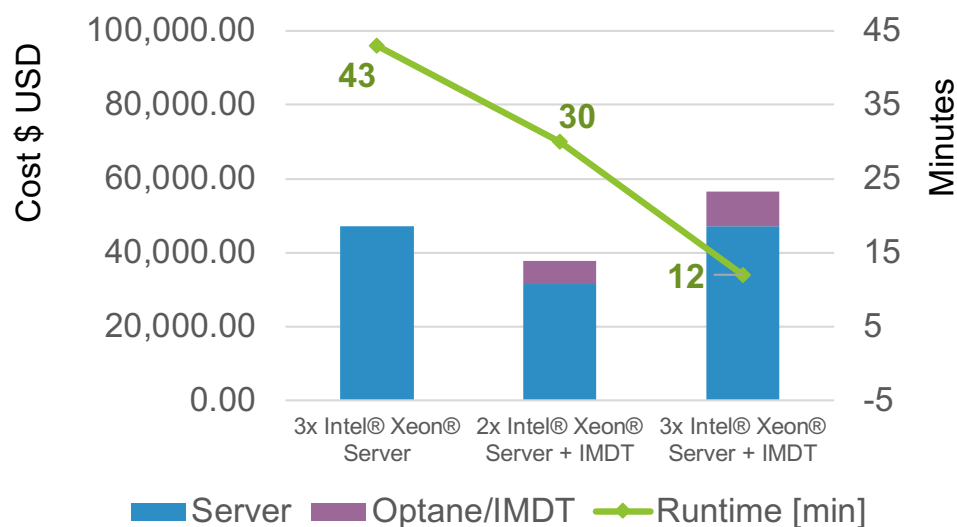


Summarized Memory usage



Solution Economics

Cluster (workers) Configuration Cost Comparison



	Master Node	Data Node (x2)
CPU	Intel® Xeon® Gold 6140 CPU @ 2.30GHz	
Cores/Socket	18	
Sockets	2	
Threads per Core	2	
Total vcores	72	
Memory	192GB	
SSD	None	3.7TB Intel® SSD DC P4500 (x2)
		375GB Intel® Optane™ SSD DC P4800X (x2)
Network	10Gbps	

20% cost reduction¹ → reduce runtime by 30%²

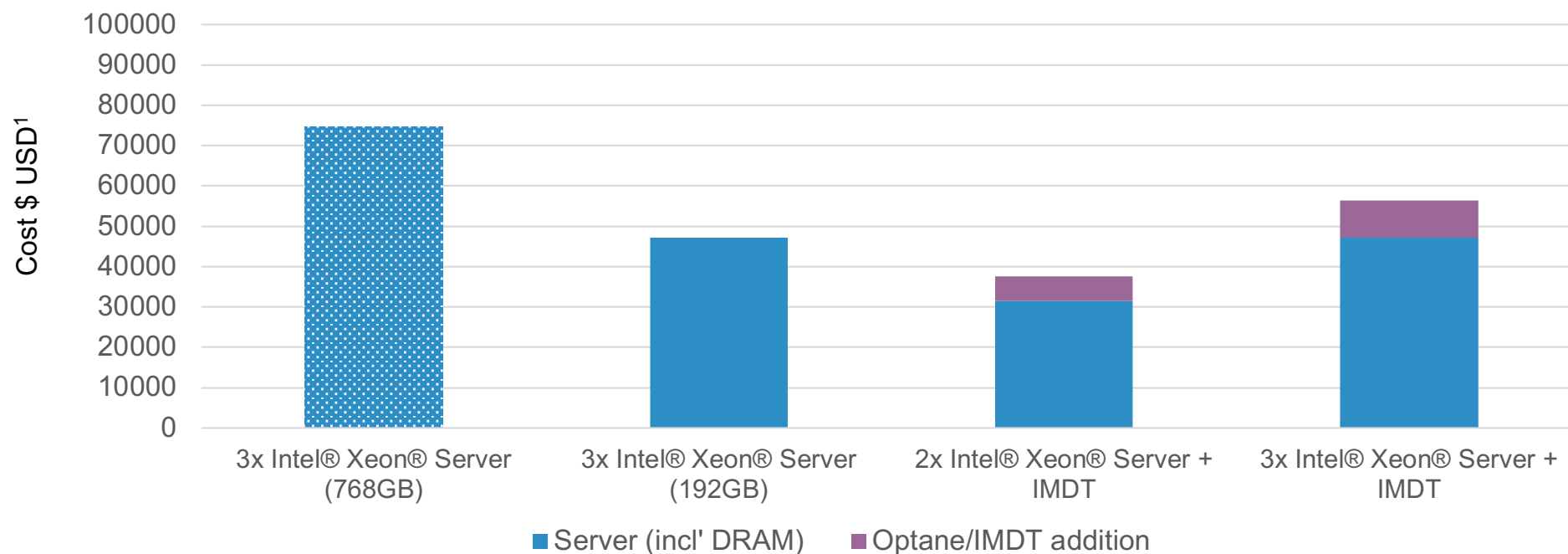
¹ Cost estimates based on quote from Colfax International as of May 27, 2018

² For system configuration details, please refer to Slide #5. Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

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

Solution Alternatives

Cluster (workers) Configuration Cost Comparison – adding the expanded all-DRAM option²



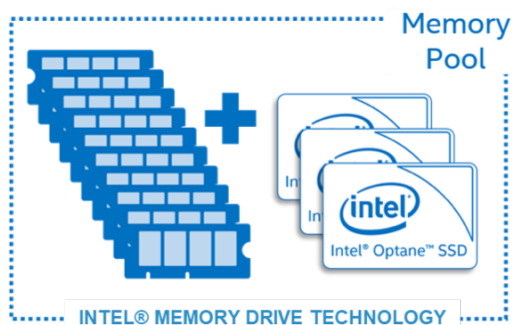
¹ Cost estimates based on quote from Colfax International as of May 27, 2018

² For system configuration details, please refer to Slide #5. Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

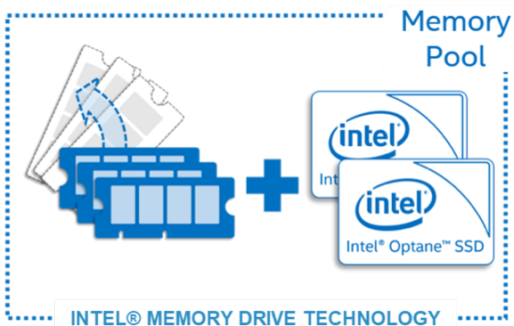
*Other names and brands may be claimed as the property of others

Summary - Optane/IMDT Benefits for Spark*

use case 1 EXPAND beyond limited DRAM capacity



use case 2 DISPLACE DRAM with affordable SSDs



- Reduce manual optimization work by having more memory available
- For workloads with underutilized CPUs:
 - Significantly reduce runtime
 - Increase CPU utilization
 - Reduce cluster node-count. Reinvest free budget in higher-core-count processors

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

Resources

- www.intel.com/optane
- www.intel.com/imdt
- <https://www.intel.com/content/www/us/en/software/apache-spark-optimization-technology-brief.html>

QUESTIONS?