

Optimized Data Analytics
with Intel[®] In-Memory Analytics Accelerator
(Intel[®] IAA) on 4th Gen Xeon[®] Scalable processors

The Intel logo is located in the bottom left corner of the slide. It consists of the word "intel" in a white, lowercase, sans-serif font, with a registered trademark symbol (®) to its upper right. The logo is positioned in front of a dark blue background that features a decorative graphic of several overlapping squares in various shades of blue, arranged in a stepped pattern.

intel[®]

Table of Contents

- Intel IAA Value Proposition
- Workload targets for Intel IAA
- Metrics/Results to help see benefit of Intel IAA
- Software optimizations - Software requirements

Why Intel® In-Memory Analytics Accelerator?

Useful for in-memory databases and beyond.
Accelerates fundamental analytics operations.
Enables gains in CPU efficiency.

Intel IAA Performance Snap Shot

Function

- Integrated accelerator IP accelerating analytics primitives, CRC calculations, compression, and decompression

Business Value

- Increases query throughput for in-memory databases and analytics workloads
- Decreases memory and bandwidth footprint for analytics workloads, freeing up space on CPU

Software Support

- Intel® Query Processing Library, Intel® Data Mover Library

Use Cases

- Commercial in-memory databases, open-source in-memory databases (RocksDB, Redis, Cassandra, MySQL, MongoDB), columnar formats for big data analytics

Performance gains
vs not using these accelerators

Embedded Databases

Up to

2.1x

Higher RocksDB
performance with built-in
Intel® IAA vs.
Zstd software

Performance gains
vs prior generation products

Embedded Databases

Up to

3x

higher RocksDB
performance with

66%

latency reduction with
built-in Intel® IAA vs.
prior generation

See [D1] at <https://edc.intel.com/content/www/us/en/products/performance/benchmarks/4th-generation-intel-xeon-scalable-processors/>

Your results may vary.

These are the kinds of database applications Intel In-Memory Analytics Accelerator can speed-up



RocksDB



ClickHouse



cassandra



MySQL®



mongoDB®

GBASE®



Parquet



Filebeat

Oceanbase

PingCAP

PrestoDB

Velox

Spark-SQL

InnoDB

CockroachDB

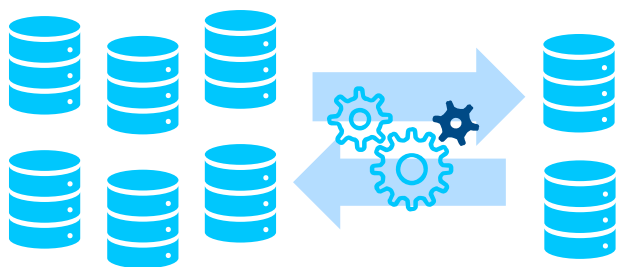
MariaDB

Intel IAA delivers these benefits

Technology Benefits	Bottom Line Benefit
Increases queries per second for analytics workloads; Reduced Latency and bottleneck for CPU access for Data; Reduced memory footprint for data at rest (storage, datalake, columnar databases, etc.) and Data Center cost reduction via memory "Cold Page" Compression.	Faster db performance, and support larger databases.
Increased effective memory capacity to store data in compressed form and decrease storage/memory footprint. Seamless Write/Transfer/Shuffle of compressed data with IAA enabled instances with minor latency gaps	Reduced memory cost, and...
Consumes less data & network bandwidth with deeply compressed data. While executing analytical functions used for database queries "on the fly" optimizes less data move around network	...can do more on the same network;
Offload CPU core-cycles spent on compression to hardware from software, and Faster access and reduced memory consumption for comp/decom data in transit & data at rest	Get more done with fewer CPUs;

Intel IAA key capabilities required to get best performance for data analytics

Compression/Decompression



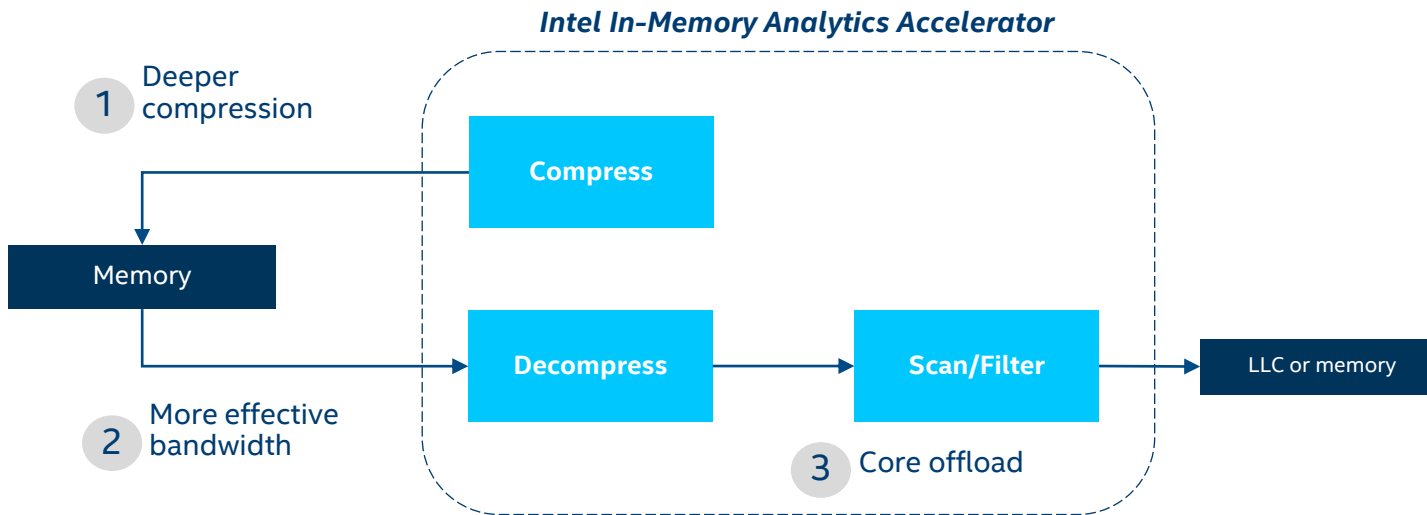
Reduces memory consumption with deeper compression vs. software-only methods

Analytics Primitives e.g. Scan, Filter

Row ID	Fruit	State	Price \$
1	Orange	FL	0.85
2	Apple	NC	0.45
4	Peach	SC	0.60
5	Grape	CA	1.25
18	Lemon	FL	0.25
19	Strawberry	CA	2.45
23	Blueberry	ME	1.5

Identifies relevant data in large data sets, accelerating database query throughput

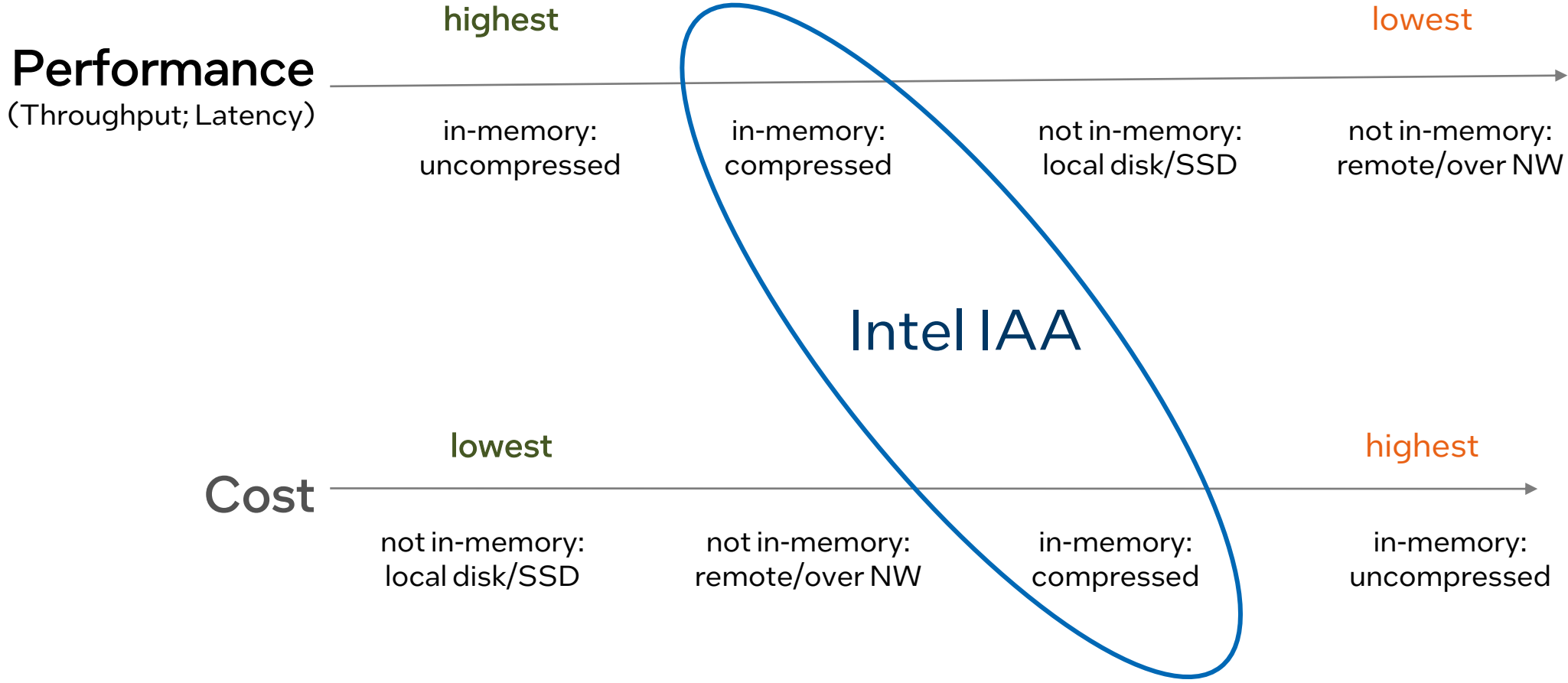
Intel IAA improves analytics workloads



- Intel IAA accelerates fundamental analytics operations
- Intel IAA optimizes memory capacity by reducing memory bandwidth for on-the-fly database queries with no need to un-compress raw data transfer
- Intel IAA enables significant gains in CPU efficiency given these operations are running on an application specific engine (and not on CPU cores).
- Speeds up query processing and provides high throughput compression and decompression (different from Intel QuickAssist Technology*(Intel* (QAT))
- Open-source Intel® Query Processing Library (Intel® QPL) wrapped around Intel IAA Compressor / Decompressor software library abstract Intel IAA hardware for applications access via API

*Intel IAA compression is different from the Deflate, ZStandard & LZ4 which are supported by Intel QAT

Intel IAA helps achieve better perf/\$ versus disk/SSD and improves memory capacity to save cost



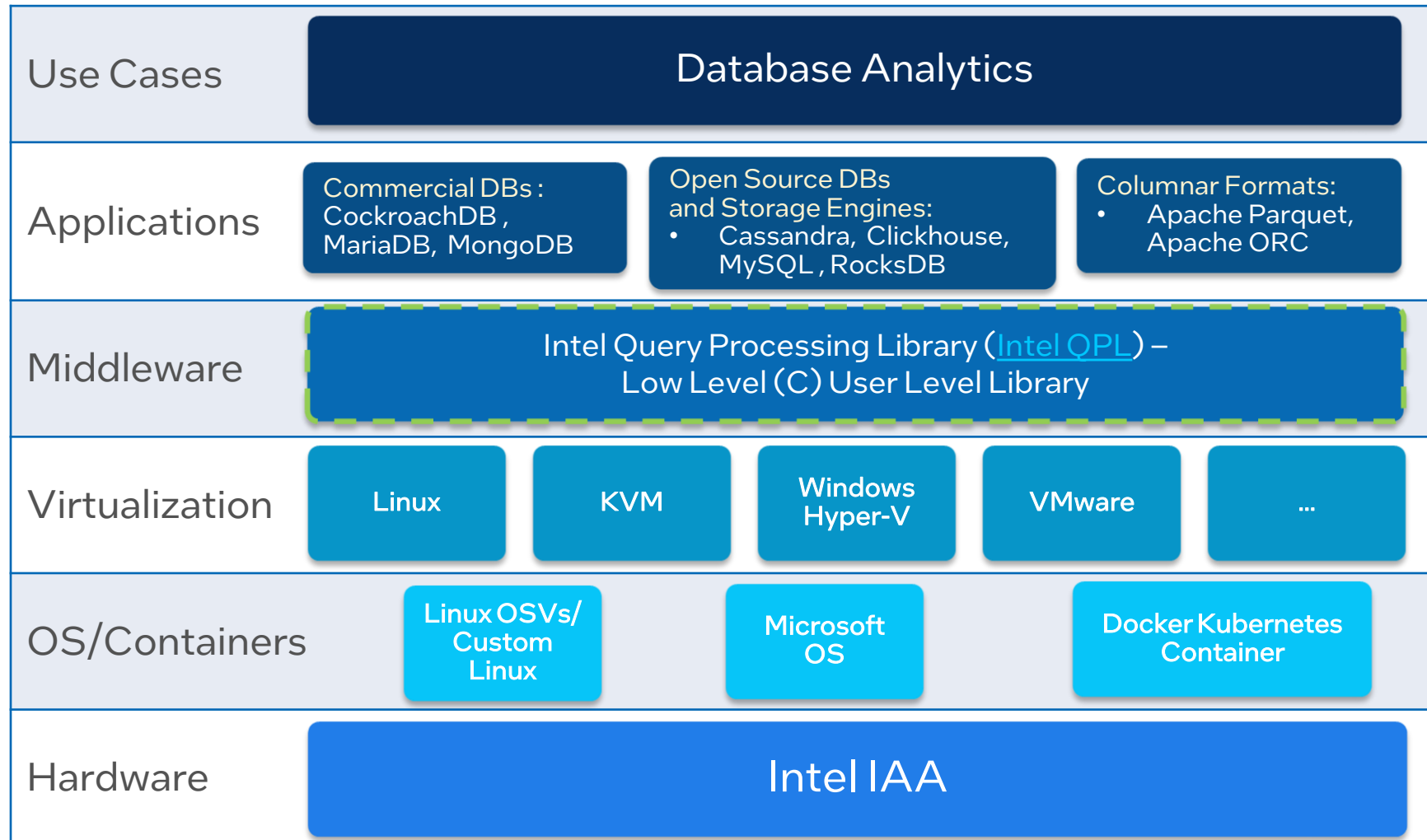
Value Prop Summarized:

With Intel IAA, database applications get:

- Higher throughput and/or lower latency vs. software,
- Improved storage/memory utilization with data compression
- Lower CPU core usage by offloading to accelerator
- Improved efficiency with higher perf/watt

See [E1] at <https://edc.intel.com/content/www/us/en/products/performance/benchmarks/4th-generation-intel-xeon-scalable-processors/>

In the software stack, Intel IAA interfaces with the Intel® Query Processing Library, the highly optimized Intel IAA Driver



- Intel supports implemented open source applications to get the benefits of Intel IAA
- Installing QPL allows the OS to map the Intel IAA accelerators, work queues, VM access, etc.
- The [QPL](#) sits on top of the system drivers. The applications will use the QPL to interface with the Intel IAA and offload analytics operations to the Intel IAA device.

Developer Tools for 4th Gen Intel[®] Xeon[®] Scalable Processors

Energize Performance with Tools to Take Advantage of Built-in Accelerators

Compilers, libraries, & analysis tools

provide support for instruction sets (such as Intel AMX, AVX2 & Intel AVX512) to unleash performance, including faster training and inference for AI workloads.

- **Intel[®] oneAPI Math Kernel Library** for HPC and technical compute
- **Intel[®] oneAPI Deep Neural Network library** for deep learning training + inference
- **Intel[®] Query Processing and Intel[®] Data Mover Library*** for query processing, compression and data movement
- **Intel[®] VTune[™] Profiler** helps locate time-consuming parts of code and identify significant issues affecting application performance.

Dynamic Load Balancer (Intel[®] DLB)

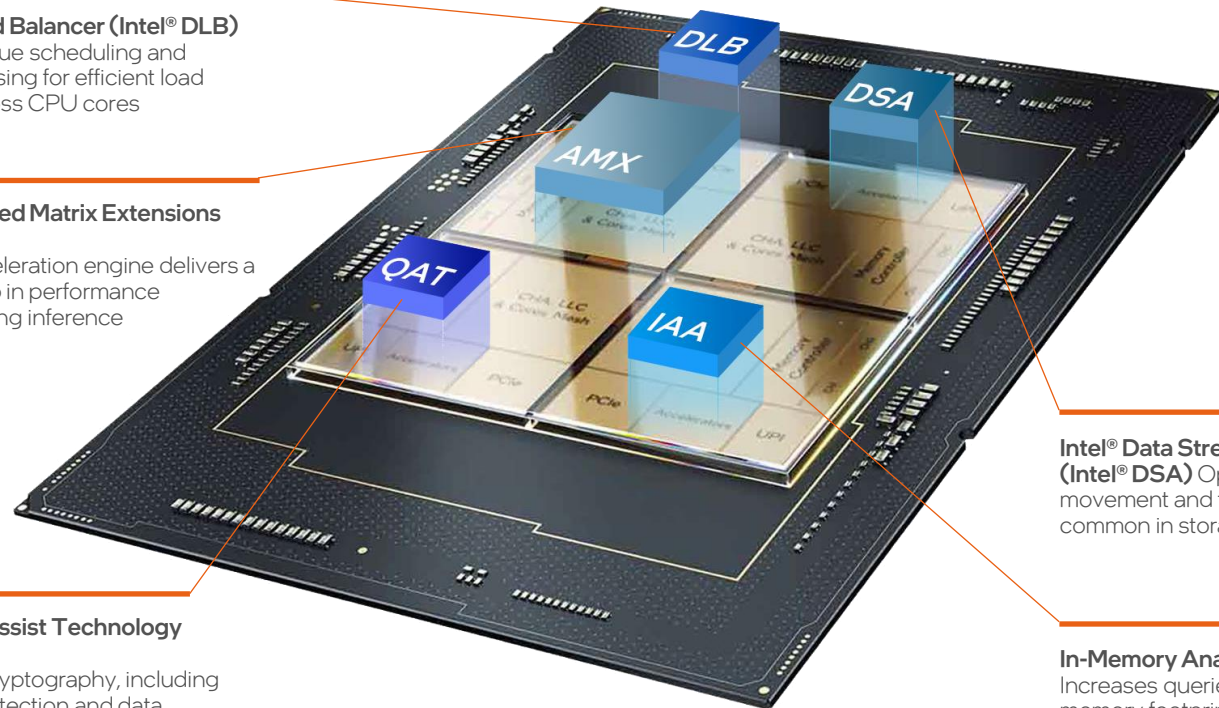
Optimizes queue scheduling and packet processing for efficient load balancing across CPU cores

Intel[®] Advanced Matrix Extensions (Intel[®] AMX)

Built-in AI acceleration engine delivers a significant leap in performance for deep learning inference and training.

Intel[®] Quick Assist Technology (Intel[®] QAT)

Accelerates cryptography, including private key protection and data de/compression.



Intel[®] Data Streaming Accelerator (Intel[®] DSA) Optimizes streaming data movement and transformation operations common in storage, networking, and analytics.

In-Memory Analytics Accelerator (Intel[®] IAA) Increases queries per second and reduces memory footprint for analytics workloads

*Open-source [Intel[®] OPL](#) v1.0.0 now available;

*Open-source [Intel[®] DML](#) in beta, v1.0.0 coming shortly

Metrics for Intel IAA

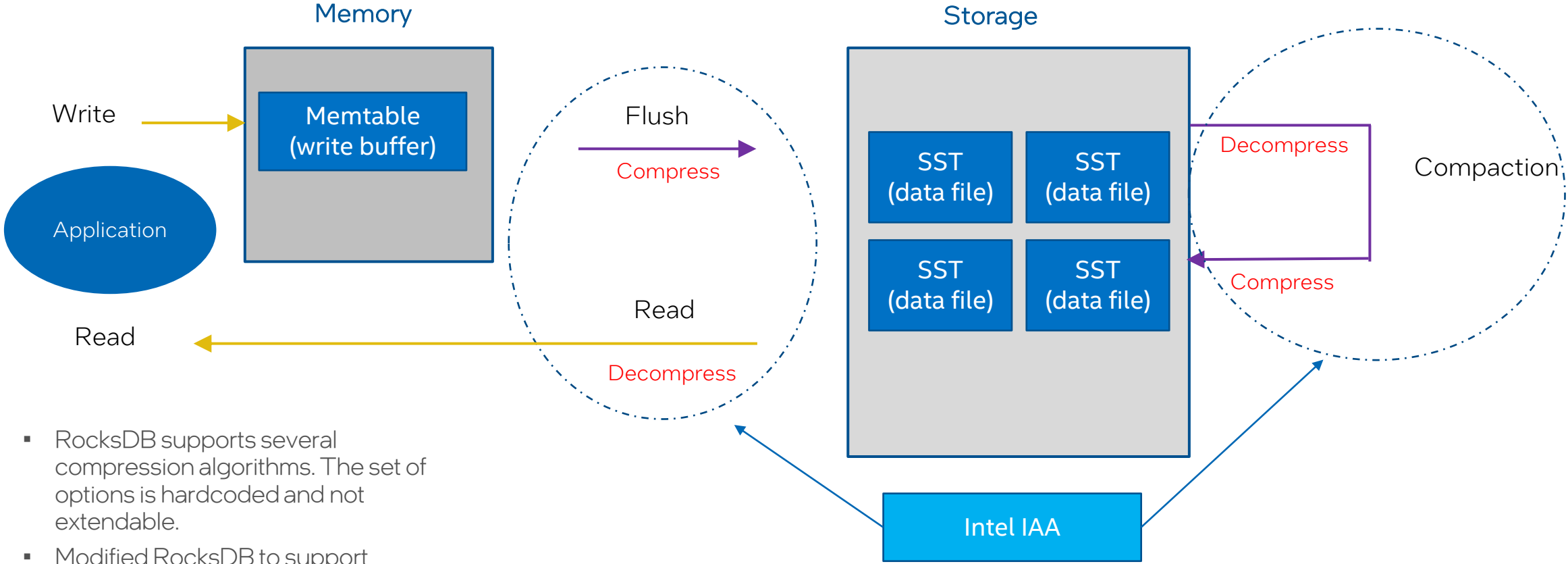
Why is RocksDB an especially good representative Workload?

Because many companies use it, and it showcases all aspects of Intel IAA value.

- RocksDB is an opensource database building block developed and maintained by Facebook DB Engineering Team
- RocksDB is the core building block for a fast key-value server, especially suited for storing data on flash drives.
- Very suitable for storing multiple terabytes of data in a single database.

- Description: [RocksDB](#) is an embedded persistent key-value store
 - Embedded: library, not a standalone database
 - Used as storage engine in many popular databases (MySQL, MariaDB, MongoDB, Redis...)
- Who uses it?
 - RocksDB is used in production systems at various [webscale](#) enterprises

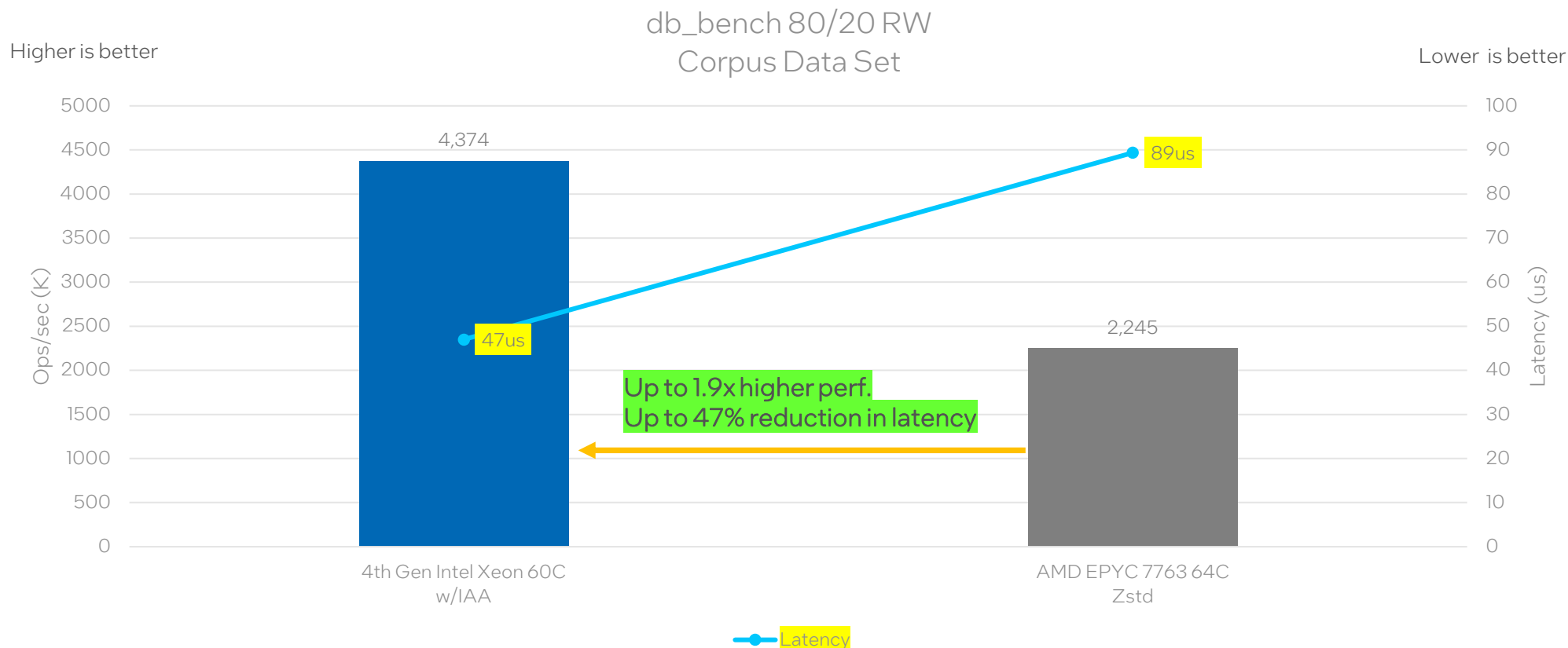
This is how RocksDB helps save memory capacity and bandwidth while speeding up throughput and saving valuable CPU cores



- RocksDB supports several compression algorithms. The set of options is hardcoded and not extendable.
- Modified RocksDB to support compressors as plugins (code to be upstreamed to RocksDB project).

SST: Sorted Sequence Table

For RocksDB, Intel IAA offers increased throughput and latency reduction



Intel IAA compressed data size within 9% of Zstd (level 3), but close to twice the throughput

Results using pre-production 4th Gen Intel® Xeon® Scalable Processor, systems, and software.
Performance varies by part, use, configuration and other factors. Learn more at www.intel.com/PerformanceIndex.
See backup for workloads and configurations. Results may vary.
All product plans, roadmaps, and performance are subject to change without notice.

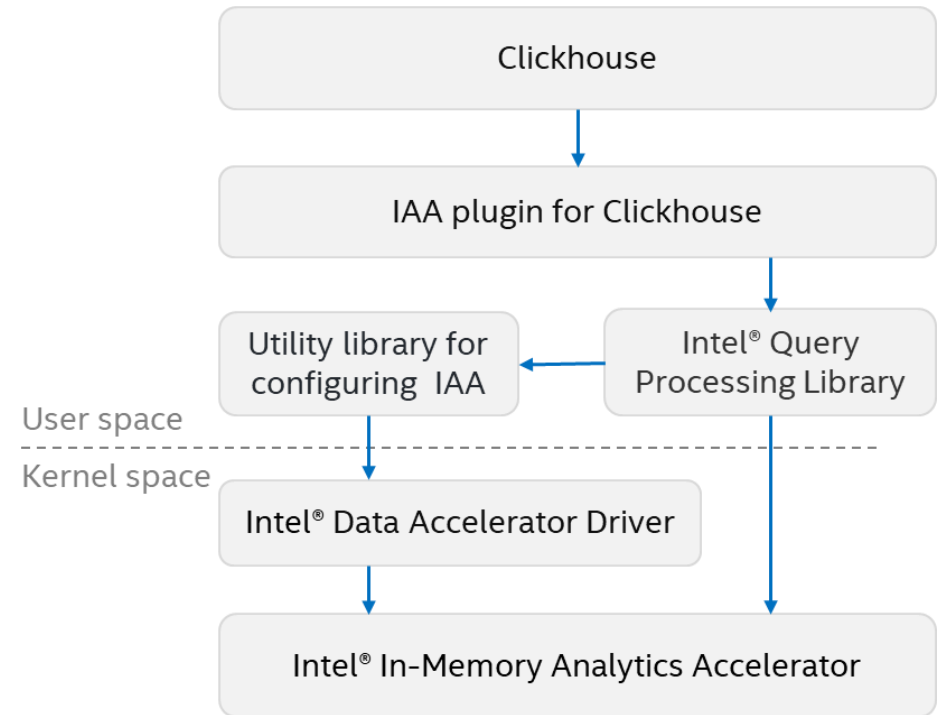
ClickHouse database demonstrates throughput, CPU savings, and bandwidth savings

▪ ClickHouse

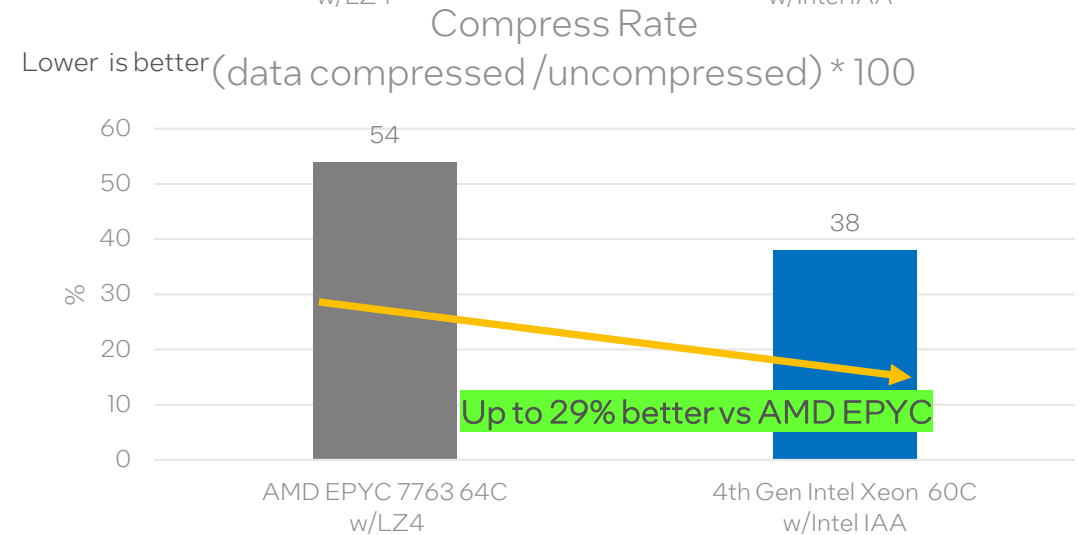
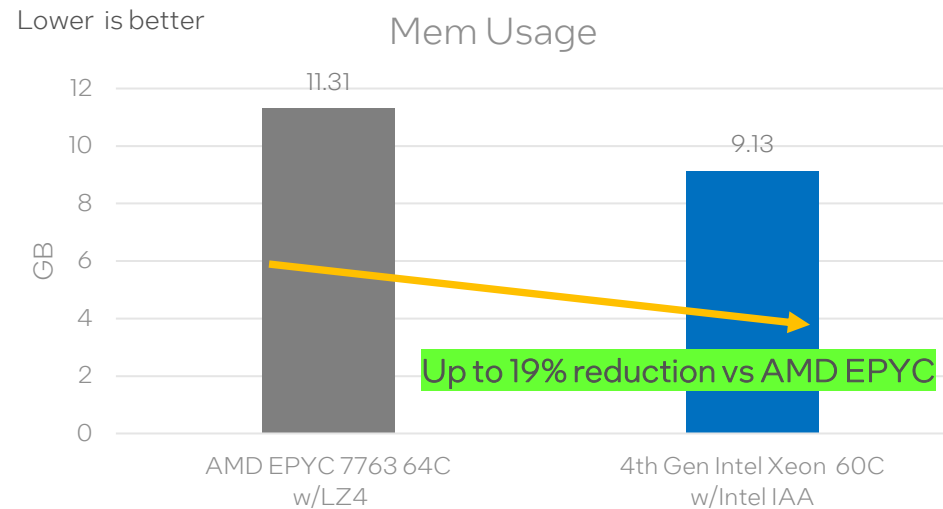
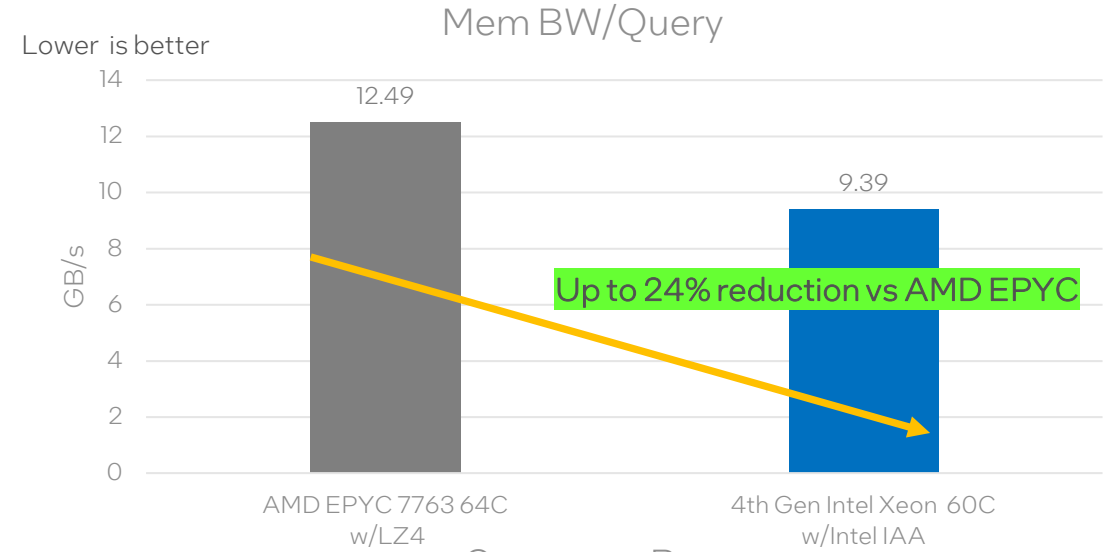
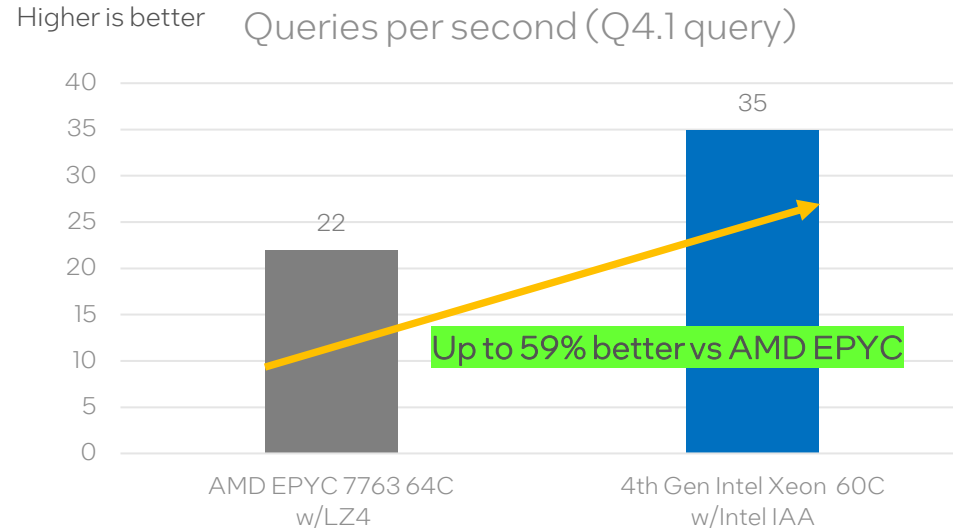
- Open-source column-oriented DBMS for online analytical processing
- Can run everywhere – bare metal or cloud, or Kubernetes
- Linearly scalable and can be scaled up to store and process trillions of rows and petabytes of data.

▪ Benchmarks

- Star Schema Benchmark: Focused on Query Q4.1 which shows the highest CPU utilization for this benchmark.
- Demonstrated query performance
- Config for all compress level is default for LZ4 (1)/ZSTD (1)/IAA-Deflate (qpl_default_level)



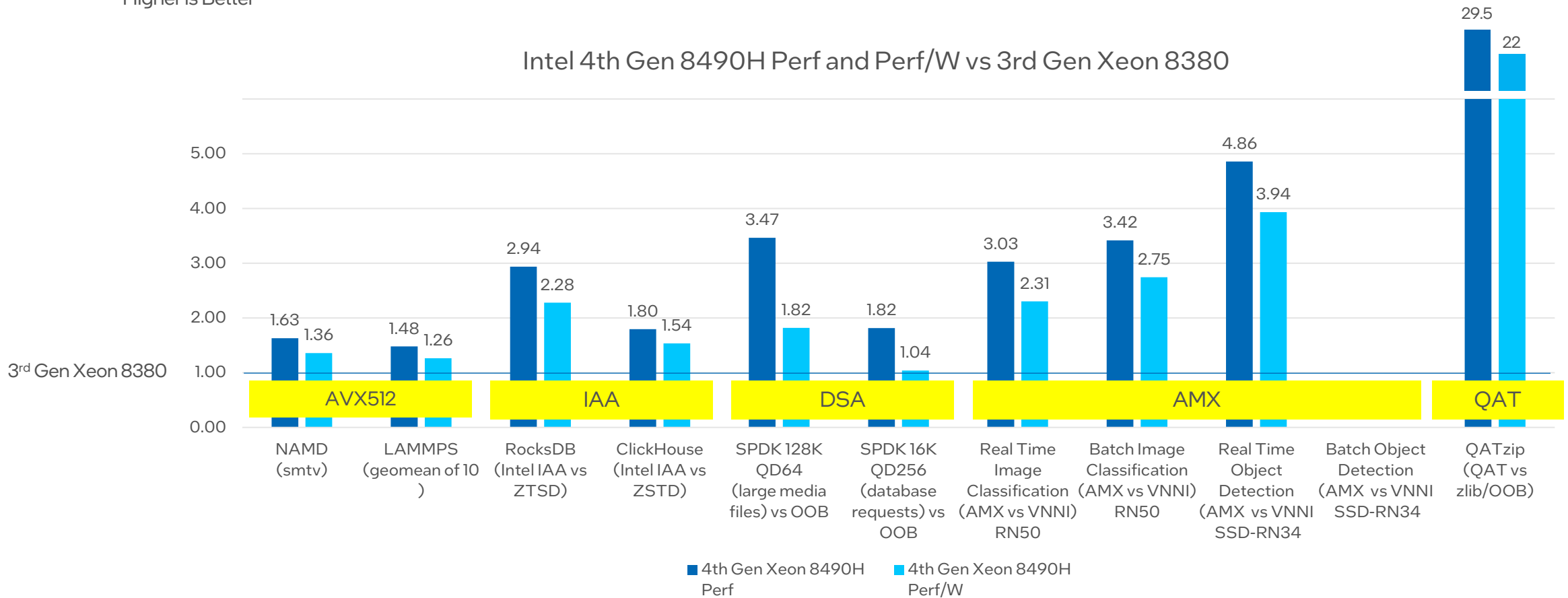
More performance and CPU resource savings with Intel IAA on ClickHouse DB



4th Gen Generational Accelerator Performance and Efficiency

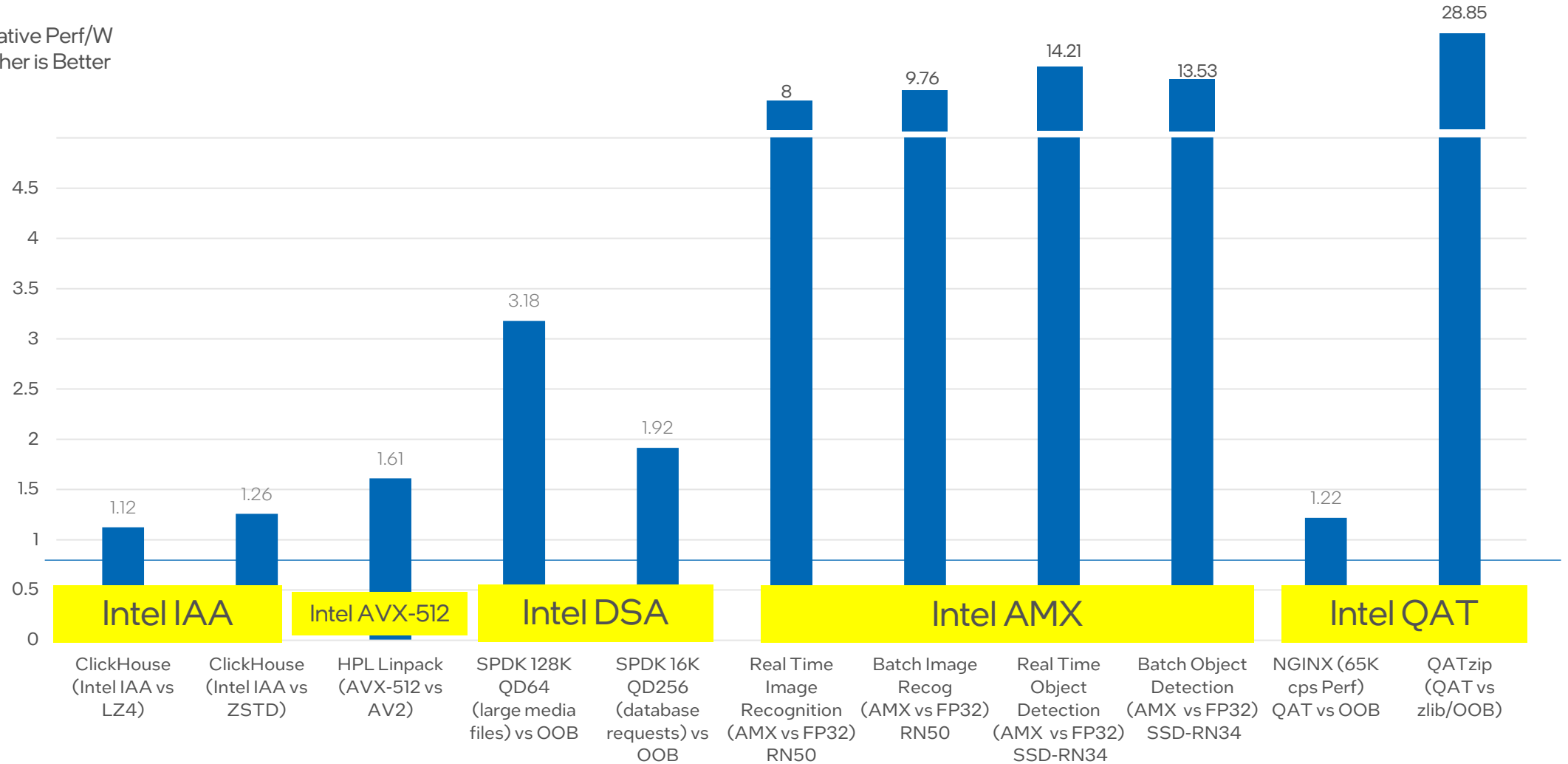
Relative Perf and Perf/W
Higher is Better

Intel 4th Gen 8490H Perf and Perf/W vs 3rd Gen Xeon 8380



4th Gen Xeon Accelerators Efficiency

Relative Perf/W
Higher is Better



Baseline
4th Gen Xeon with
No Acceleration

Utilizing Intel IAA for database applications delivers Cost Savings (data shown versus 3rd Gen Intel Xeon Processors)

Comparisons to deploying 50 servers with 3 rd Gen Intel Xeon processor	Database (Rocks DB w/Intel [®] IAA)
Number of Intel Xeon processor-based servers	18 servers <small>with 4th Gen Intel[®] Xeon processors</small>
Lower Fleet Power (kilowatts)	15.4 kW
Reduced CO2 emissions (kg)*	366,000 kg
TCO savings (\$)*	\$1.2M
	55% Lower TCO

See [E7] at <https://edc.intel.com/content/www/us/en/products/performance/benchmarks/4th-generation-intel-xeon-scalable-processors/>

Notices & Disclaimers

Performance varies by use, configuration and other factors. Learn more on the [Performance Index site](#).

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See backup for configuration details. No product or component can be absolutely secure.

Your costs and results may vary.

Intel technologies may require enabled hardware, software or service activation.

Availability of accelerators varies by SKU. Visit

<https://ark.intel.com/content/www/us/en/ark/products/series/228622/4th-generation-intel-xeon-scalable-processors.html>

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

The Intel logo is centered on a solid blue background. It consists of the word "intel" in a white, lowercase, sans-serif font. A small blue square is positioned above the letter 'i'. To the right of the word "intel" is a registered trademark symbol (®).

intel®

Configurations

Configuration for RocksDB

1-node, 2x pre-production 4th Gen Intel Xeon Scalable Processor (60 cores) with integrated Intel In-Memory Analytics Accelerator (Intel IAA), on pre-production Intel platform and software, HT On, Turbo On, Total Memory 1024GB (16x64GB DDR5 4800), microcode 0xf000380, 1x 1.92TB INTEL SSDSC2KG01, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, QPL v0.1.21, accel-config-v3.4.6.4, ZSTD v1.5.2, RocksDB v6.4.6 (db_bench), tested by Intel September 2022.

1-node, 2x AMD* EPYC 7763 64 core Processor on GIGABYTE R282-Z92 platform, SMT On, Boost On, NPS=1, Total Memory 1024GB (16x64GB DDR4 3200), microcode 0xa001144, 1x 1.92TB INTEL SSDSC2KG01, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, ZSTD v1.5.2, RocksDB v6.4.6 (db_bench), tested by Intel September 2022.

Configurations: ClickHouse (Intel vs AMD)

- **INTEL 4th Gen Xeon Scalable Processor**

1-node, 2x pre-production 4th Gen Intel Xeon Scalable processor (60 cores) with integrated Intel In-Memory Analytics Accelerator (Intel IAA), on pre-production Intel platform and software, HT On, Turbo On, SNC off, Total Memory 1024GB (16x64GB DDR5 4800), microcode 0xf000380, 1x 1.92TB INTEL SSDSC2KG01, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, QPL v0.1.21, accel-config-v3.4.6.4, gcc 11.2, Clickhouse 21.12, Star Schema Benchmark, tested by Intel September 2022.

- **AMD EPYC 7763**

1-node, 2x AMD EPYC 7763 64 core Processor on GIGABYTE R282-Z92 platform, SMT On, Boost On, NPS=1, Total Memory 1024GB (16x64GB DDR4 3200), microcode 0xa001144, 1x 1.92TB INTEL SSDSC2KG01, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, gcc 11.2, Clickhouse 21.12, Star Schema Benchmark, tested by Intel September 2022.

Resources and Configurations

A More Energy Efficient Server Architecture

Up to 1.12x and 1.26x higher performance/W using 4th Gen Xeon Scalable w/Intel Analytics Accelerator vs LZ4 and Zstd on ClickHouse

1-node, 2x pre-production 4th Gen Intel Xeon Scalable processor (60 cores) with integrated Intel In-Memory Analytics Accelerator (Intel IAA), Number of IAA device utilized=8(2 sockets active), on pre-production Intel platform and software, HT On, Turbo On, SNC off, Total Memory 1024GB (16x64GB DDR5 4800), microcode 0x2b0000a1, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, QPL v0.1.21, accel-config-v3.4.6.4, gcc 11.2, Clickhouse 21.12, Star Schema Benchmark, tested by Intel November 2022.

Up to 2.01x higher performance/W using 4th Gen Xeon Scalable w/Intel Analytics Accelerator vs Zstd on RocksDB

1-node, 2x pre-production 4th Gen Intel Xeon Scalable Processor (60 cores) with integrated Intel In-Memory Analytics Accelerator (Intel IAA), on pre-production Intel platform and software, HT On, Turbo On, Total Memory 1024GB (16x64GB DDR5 4800), microcode 0x2b0000a1, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, QPL v0.2.1, accel-config-v3.4.6.4, ZSTD v1.5.2, RocksDB v6.4.6 (db_bench), tested by Intel November 2022.

Up to 1.61 higher performance/W using 4th Gen Xeon Scalable w/AVX-512 vs AVX2 on Linpack

1-node, 2x pre-production 4th Gen Intel® Xeon® Scalable processor (60 core), on pre-production Supermicro SYS-221H-TNR and software with 1024GB DDR5 memory (16x64 GB), microcode 0x2b0000c0, HT On, Turbo On, SNC 4, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, One API BaseKit 2022.2.0.262, One API HPC 2022.2.0.191, Linpack ver 2.3, tested by Intel November 2022.

Up to 3.18x and 1.92x higher performance/W using 4th Gen Xeon Scalable w/Data Streaming Accelerator vs out-of-box OS software on SPDK NVMe TCP

1-node, 2x pre-production 4th Gen Intel Xeon Scalable processor (60 core) with integrated Intel Data Streaming Accelerator (Intel DSA), DSA device utilized=1(1 active socket), on pre-production Intel platform and software with 1024GB DDR5 memory (16x64 GB), microcode 0x2b0000a1, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x 1.92TB Intel® SSDSC2KG01, 4x 1.92TB Samsung PMI733, 1x Intel® Ethernet Network Adapter E810-2CQDA2, 2x100GbE, FIO v3.30, SPDK 22.05, tested by Intel November 2022.

Up to 8x and 9.76x higher performance/W using 4th Gen Xeon Scalable w/Advanced Matrix Extensions using AMX vs VNNI instructions on ResNet50 Image Processing

1-node, 2x pre-production 4th Gen Intel® Xeon® Scalable processor (60 core) with Intel® Advanced Matrix Extensions (Intel AMX), on pre-production Supermicro SYS-221H-TNR with 1024GB DDR5 memory (16x64 GB), microcode 0x2b0000c0, HT On, Turbo On, SNC Off, CentOS Stream 8, 5.19.16-301.fc37.x86_64, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Intel TF 2.10, AI Model=Resnet 50 v1_5, best scores achieved: BS1 FP32 8 cores/instance (max. 15ms SLA), BS1 INT8 2 cores/instance (max. 15ms SLA), BS1 AMX 1 core/instance (max. 15ms SLA), BS16 FP32 5 cores/instance, BS16 INT8 5 cores/instance, BS16 AMX 5 cores/instance, using physical cores, tested by Intel November 2022.

Up to 14.21x and 13.53x higher performance/W using 4th Gen Intel Xeon Scalable w/Advanced Matrix Extensions using AMX vs VNNI instructions on SSD-ResNet34 on Object Detection

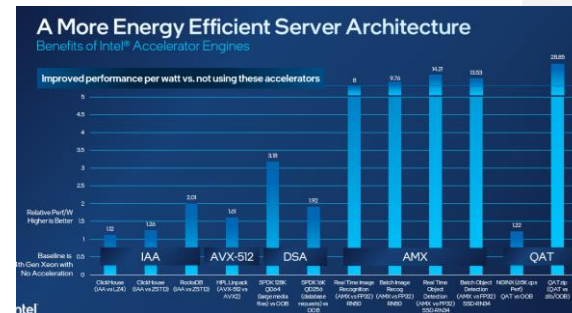
1-node, 2x pre-production 4th Gen Intel® Xeon® Scalable processor (60 core) with Intel® Advanced Matrix Extensions (Intel AMX), Intel platform with 1024GB DDR5 memory (16x64 GB), microcode 0x2b0000a1, HT On, Turbo On, SNC Off, CentOS Stream 8, 5.19.16-301.fc37.x86_64, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Intel TF 2.10, AI Model=SSD-ResNet34, best scores achieved: BS1 FP32 60 cores/instance (max. 100ms SLA), BS1 INT8 4 cores/instance (max. 100ms SLA), BS1 AMX 4 core/instance (max. 100ms SLA), BS8 FP32 8 cores/instance, BS2 INT8 1 cores/instance, BS2 AMX 1 cores/instance, using physical cores, tested by Intel November 2022.

Up to 1.22x higher performance/W using 4th Gen Intel Xeon Scalable w/QuickAssist Accelerator vs out-of-box software on NGINX TLS Handshake.

QAT Accelerator: 1-node, 2x pre-production 4th Gen Intel Xeon Scalable Processor (60 cores) with integrated Intel QuickAssist Accelerator (Intel QAT), Number of QAT device utilized=4(1 socket active), on pre-production Intel platform and software with DDR5 memory total 1024GB (16x64 GB), microcode 0x2b0000a1, HT On, Turbo Off, SNC Off, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, 1x Intel® Ethernet Network Adapter E810-2CQDA2, 1x100GbE, QAT engine v0.6.14, QAT v20.1.0.9.1, NGINX 1.20.1, OpenSSL 1.1.1l, IPP crypto v2021_5, IPsec v1.1, TLS 1.3 AES_128_GCM_SHA256, ECDHE-X25519-RSA2K, 65K CPS target SLA, tested by Intel November 2022. Out of box configuration: 1-node, 2x pre-production 4th Gen Intel Xeon Scalable Processor (60 cores) with integrated Intel QuickAssist Accelerator (Intel QAT), Number of QAT device utilized=0, on pre-production Intel platform and software with DDR5 memory total 1024GB (16x64 GB), microcode 0x2b0000a1, HT On, Turbo Off, SNC Off, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, 1x Intel® Ethernet Network Adapter E810-2CQDA2, 1x100GbE, NGINX 1.20.1, OpenSSL 1.1.1l, TLS 1.3 AES_128_GCM_SHA256, ECDHE-X25519-RSA2K, 65K CPS target SLA, tested by Intel November 2022.

Up to 28.85x higher performance/W using 4th Gen Intel Xeon Scalable w/QuickAssist Accelerator vs out-of-box zlib on QATzip compression

1-node, 2x pre-production 4th Gen Intel® Xeon Scalable Processor (60 core) with integrated Intel QuickAssist Accelerator (Intel QAT), QAT device utilized=8(2 sockets active), on pre-production Intel platform and software with DDR5 memory Total 1024GB (16x64 GB), microcode 0x2b0000a1, HT On, Turbo Off, SNC Off, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, QAT v20.1.0.9.1, QATzip v1.0.9, tested by Intel November 2022.



Resources and Configurations

Significant Performance and Performance /Watt Gains:

NAMD

New Configuration: 1-node, 2x pre-production 4th Gen Intel® Xeon® Scalable processor (60 core), on pre-production Supermicro SYS-221H-TNR and software with 1024GB DDR5 memory (16x64 GB), microcode 0x2b0000c0, HT On, Turbo On, SNC 4, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, One API BaseKit 2022.2.0.262, One API HPC 2022.2.0.191, NAMD release-2-15-alpha-1, charm v6.10.2, tcl core-8-5-branch, benchmark from NAMD v2.13, tested by Intel November 2022. Baseline: 1-node, 2x production 3rd Gen Intel Xeon Scalable Processors(40 cores) on Supermicro SYS-220U-TNR, DDR4 memory total 1024GB (16x64 GB), microcode 0xd000375, HT On, Turbo On, SNC On, Turbo On, SNC On, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, One API BaseKit 2022.2.0.262, One API HPC 2022.2.0.191, NAMD release-2-15-alpha-1, charm v6.10.2, tcl core-8-5-branch, benchmark from NAMD v2.13, tested by Intel November 2022.

LAMMPS

New Configuration: 1-node, 2x pre-production 4th Gen Intel® Xeon® Scalable processor (60 core), on pre-production Supermicro SYS-221H-TNR and software with 1024GB DDR5 memory (16x64 GB), microcode 0x2b0000c0, HT On, Turbo On, SNC 4, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, One API BaseKit 2022.2.0.262, One API HPC 2022.2.0.191, LAMMPS update 2 for Stable release 23 June 2022, tested by Intel November 2022. Baseline: 1-node, 2x production 3rd Gen Intel Xeon Scalable Processors(40 cores) on Supermicro SYS-220U-TNR, DDR4 memory total 1024GB (16x64 GB), microcode 0xd000375, HT On, Turbo On, SNC On, HT On, Turbo On, SNC On, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, One API BaseKit 2022.2.0.262, One API HPC 2022.2.0.191, LAMMPS update 2 for Stable release 23 June 2022, tested by Intel November 2022.

RocksDB

New Configuration: 1-node, 2x pre-production 4th Gen Intel Xeon Scalable Processor (60 cores) with integrated Intel In-Memory Analytics Accelerator (Intel IAA), on pre-production Intel platform and software, HT On, Turbo On, Total Memory 1024GB (16x64GB DDR5 4800), microcode 0x2b0000a1, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, OPL v0.2.1, accel-config-v3.4.6.4, ZSTD v1.5.2, RocksDB v6.4.6 (db_bench), tested by Intel November 2022. Baseline: 1-node, 2x production 3rd Gen Intel Xeon Scalable Processors(40 cores) on Supermicro SYS-220U-TNR, HT On, Turbo On, SNC Off, Total Memory 1024GB (16x64GB DDR4 3200), microcode 0xd000375, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, ZSTD v1.5.2, RocksDB v6.4.6 (db_bench), tested by Intel November 2022.

ClickHouse

New Configuration: 1-node, 2x pre-production 4th Gen Intel Xeon Scalable processor (60 cores) with integrated Intel In-Memory Analytics Accelerator (Intel IAA), Number of IAA device utilized=8(2 sockets active), on pre-production Intel platform and software, HT On, Turbo On, SNC off, Total Memory 1024GB (16x64GB DDR5 4800), microcode 0x2b0000a1, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, OPL v0.1.21, accel-config-v3.4.6.4, gcc 11.2, Clickhouse 21.12, Star Schema Benchmark, tested by Intel November 2022. Baseline: 1-node, 2x production 3rd Gen Intel Xeon Scalable Processors(40 cores) on Supermicro SYS-220U-TNR, HT On, Turbo On, SNC Off, Total Memory 1024GB (16x64GB DDR4 3200), microcode 0xd000375, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, gcc 11.2, Clickhouse 21.12, Star Schema Benchmark, tested by Intel November 2022.

SPDK

New Configuration: 1-node, 2x pre-production 4th Gen Intel Xeon Scalable processor (60 core) with integrated Intel Data Streaming Accelerator (Intel DSA), DSA device utilized=1(1 active socket), on pre-production Intel platform and software with 1024GB DDR5 memory (16x64 GB), microcode 0x2b0000a1, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x 1.92TB Intel® SSDSC2K01, 4x 1.92TB Samsung PM1733, 1x Intel® Ethernet Network Adapter E810-2CQDA2, 2x 100GbE, FIO v3.30, SPDK 22.05, tested by Intel November 2022. Baseline: 1-node, 2x production 3rd Gen Intel Xeon Scalable Processors(40 cores) on Supermicro SYS-220U-TNR, DDR4 memory total 1024GB (16x64 GB), HT On, Turbo On, SNC Off, microcode 0xd000375, 10GbE x540-AT2, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x 1.92TB Intel SSDSC2K01, 4x 1.92TB Samsung PM1733, 1x Intel Ethernet Network Adapter E810-2CQDA2, 2x 100GbE, FIO v3.30, SPDK 22.05, tested by Intel November 2022.

ResNet-50

New Configuration: 1-node, 2x pre-production 4th Gen Intel® Xeon® Scalable processor (60 core) with Intel® Advanced Matrix Extensions (Intel AMX), on pre-production Supermicro SYS-221H-TNR with 1024GB DDR5 memory (16x64 GB), microcode 0x2b0000c0, HT On, Turbo On, SNC Off, CentOS Stream 8, 5.19.16-301.fc37.x86_64, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Intel TF 2.10, AI Model=Resnet 50 v1.5, best scores achieved: BS1 FP32 8 cores/instance (max. 15ms SLA), BS1 INT8 2 cores/instance (max. 15ms SLA), BS1 AMX 1 core/instance (max. 15ms SLA), BS16 FP32 5 cores/instance, BS16 INT8 5 cores/instance, BS16 AMX 5 cores/instance, using physical cores, tested by Intel November 2022. Baseline: 1-node, 2x production 3rd Gen Intel Xeon Scalable Processors(40 cores) on Supermicro SYS-220U-TNR, DDR4 memory total 1024GB (16x64 GB), microcode 0xd000375, HT On, Turbo On, SNC Off, CentOS Stream 8, 5.19.16-301.fc37.x86_64, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Intel TF 2.10, AI Model=Resnet 50 v1.5, best scores achieved: BS1 FP32 8 cores/instance (max. 15ms SLA), BS1 INT8 2 cores/instance (max. 15ms SLA), BS16 FP32 5 cores/instance, BS16 INT8 5 cores/instance, using physical cores, tested by Intel November 2022.

SSD-ResNet-34

New Configuration: 1-node, 2x pre-production 4th Gen Intel® Xeon® Scalable processor (60 core) with Intel® Advanced Matrix Extensions (Intel AMX), Intel platform with 1024GB DDR5 memory (16x64 GB), microcode 0x2b0000a1, HT On, Turbo On, SNC Off, CentOS Stream 8, 5.19.16-301.fc37.x86_64, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Intel TF 2.10, AI Model=SSD-ResNet34, best scores achieved: BS1 FP32 60 cores/instance (max. 100ms SLA), BS1 INT8 4 cores/instance (max. 100ms SLA), BS1 AMX 4 core/instance (max. 100ms SLA), BS8 FP32 8 cores/instance, BS2 INT8 1 cores/instance, BS2 AMX 1 cores/instance, using physical cores, tested by Intel November 2022. Baseline: 1-node, 2x production 3rd Gen Intel Xeon Scalable Processors(40 cores) on Supermicro SYS-220U-TNR, DDR4 memory total 1024GB (16x64 GB), microcode 0xd000375, HT On, Turbo On, SNC Off, CentOS Stream 8, 5.19.16-301.fc37.x86_64, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Intel TF 2.10, AI Model=SSD-ResNet34, best scores achieved: BS1 FP32 40 cores/instance (max. 100ms SLA), BS1 INT8 10 cores/instance (max. 100ms SLA), BS16 FP32 4 cores/instance, BS16 INT8 1 cores/instance, using physical cores, tested by Intel November 2022.

QAT.zip

New Configuration: 1-node, 2x pre-production 4th Gen Intel® Xeon Scalable Processor (60 core) with integrated Intel QuickAssist Accelerator (Intel OAT), OAT device utilized=8(2 sockets active), on pre-production Intel platform and software with DDR5 memory Total 1024GB (16x64 GB), microcode 0x2b0000a1, HT On, Turbo Off, SNC Off, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, QAT v2.0.10.9.1, QATzip v1.0.9, tested by Intel November 2022. Baseline: 1-node, 2x production 3rd Gen Intel Xeon Scalable Processors(40 cores) on Supermicro SYS-220U-TNR, DDR4 memory total 1024GB (16x64 GB), microcode 0xd000375, HT On, Turbo Off, SNC Off, Ubuntu 22.04.1 LTS, 5.15.0-52-generic, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, QAT v2.0.10.9.1, QATzip v1.0.9, tested by Intel November 2022.



Resources and Configurations

A More Cost-Efficient Server Architecture

ResNet50 Image Classification

New Configuration: 1-node, 2x pre-production 4th Gen Intel® Xeon® Scalable 8490H processor (60 core) with Intel® Advanced Matrix Extensions (Intel AMX), on pre-production SuperMicro SYS-221H-TNR with 1024GB DDR5 memory (16x64 GB), microcode 0x2b0000c0, HT On, Turbo On, SNC Off, CentOS Stream 8, 5.19.16-301.fc37.x86_64, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Intel TF 2.10, AI Model=Resnet 50 v1_5, best scores achieved: BS1 AMX 1 core/instance (max. 15ms SLA), using physical cores, tested by Intel November 2022. Baseline: 1-node, 2x production 3rd Gen Intel Xeon Scalable 8380 Processor (40 cores) on SuperMicro SYS-220U-TNR, DDR4 memory total 1024GB (16x64 GB), microcode 0xd000375, HT On, Turbo On, SNC Off, CentOS Stream 8, 5.19.16-301.fc37.x86_64, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Intel TF 2.10, AI Model=Resnet 50 v1_5, best scores achieved: BS1 INT8 2 cores/instance (max. 15ms SLA), using physical cores, tested by Intel November 2022.

For a 50 server fleet of 3rd Gen Xeon 8380 (RN50 w/DLBoost), estimated as of November 2022:

CapEx costs: \$1.64M

OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$739.9K

Energy use in kWh (4 year, per server): 44627, PUE 1.6

Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394

For a 17 server fleet of 4th Gen Xeon 8490H (RN50 w/AMX), estimated as of November 2022:

CapEx costs: \$799.4K

OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$275.3K

Energy use in kWh (4 year, per server): 58581, PUE 1.6

Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394

RocksDB

New Configuration: 1-node, 2x pre-production 4th Gen Intel Xeon Scalable 8490H Processor (60 cores) with integrated Intel In-Memory Analytics Accelerator (Intel IAA), on pre-production Intel platform and software, HT On, Turbo On, Total Memory 1024GB (16x64GB DDR5 4800), microcode 0x2b0000a1, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, QPL v0.2.1, accel-config-v3.4.6.4, ZSTD v1.5.2, RocksDB v6.4.6 (db_bench), tested by Intel November 2022. Baseline: 1-node, 2x production 3rd Gen Intel Xeon Scalable 8380 Processor (40 cores) on SuperMicro SYS-220U-TNR, HT On, Turbo On, SNC Off, Total Memory 1024GB (16x64GB DDR4 3200), microcode 0xd000375, 1x3.84TB P5510 NVMe, 10GbE x540-AT2, Ubuntu 22.04.1 LTS, 5.18.12-051812-generic, ZSTD v1.5.2, RocksDB v6.4.6 (db_bench), tested by Intel November 2022.

For a 50 server fleet of 3rd Gen Xeon 8380 (RocksDB), estimated as of November 2022:

CapEx costs: \$1.64M

OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$677.7K

Energy use in kWh (4 year, per server): 32181, PUE 1.6

Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394

For a 18 server fleet of 4th Gen Xeon 8490H (RockDB w/IAA), estimated as of November 2022:

CapEx costs: \$846.4K

OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$260.6K

Energy use in kWh (4 year, per server): 41444, PUE 1.6

Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394

OpenFOAM

New Configuration: 1-node, 2x pre-production 4th Gen Intel Xeon CPU Max Series (56 cores) on pre-production Intel platform and software, HT On, Turbo On, SNC4 mode, Total Memory 128 GB (8x16GB HBM2 3200MT/s), microcode 0x2c000020, 1x3.5TB INTEL SSDPF2KX038TZ NVMe, CentOS Stream 8, 5.19.0-rc6.0712.intel_next.1.x86_64+server, OpenFOAM 8, Motorbike 20M @ 250 iterations, Motorbike 42M @ 250 iterations, Tools: ifort:2021.6.0, icc:2021.6.0, impi:2021.6.0, tested by Intel December 2022. Baseline: 1-node, 2x production 3rd Gen Intel Xeon Scalable 8380 Processor (40 cores) on SuperMicro SYS-220U-TNR, HT On, Turbo On, 512GB (16x32GB DDR4 3200 MT/s), microcode 0xd000375, 1x2.9TB INTEL SSDPE2KE032T8 NVMe, CentOS Stream 8, 4.18.0-408.el8.x86_64, OpenFOAM 8, Motorbike 20M @ 250 iterations, Motorbike 42M @ 250 iterations, Tools: ifort:2021.6.0, icc:2021.6.0, impi:2021.6.0, tested by Intel December 2022

For a 50 server fleet of 3rd Gen Xeon 8380 (OpenFOAM), estimated as of December 2022:

CapEx costs: \$1.50M

OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$780.3K

Energy use in kWh (4 year, per server): 52700, PUE 1.6

Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394

For a 16 server fleet of Intel Xeon CPU Max Series 56 core, estimated as of December 2022:

CapEx costs: \$507.2K

OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$274.9K

Energy use in kWh (4 year, per server): 74621, PUE 1.6

Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394