Al on Intel® Xeon® processors

Partner Enablement Package

Addressing customers' Al business challenges with Intel® Xeon® based solutions



Contents

- > Why Partner with Intel on AI
 - Why Choose Intel[®] Xeon[®] Processors for Developing an AI Solution
- Intel Al Portfolio
 - Scalable Systems for AI
 - Intel AI Software is Enterprise Ready
- > Intel[®] Xeon[®] for AI
 - Intel[®] Xeon[®] the Processor for AI
 - Intel[®] Xeon[®] Processor delivers TCO Value for Mixed General-Purpose and AI Workloads
 - Accelerators
 - 5th Gen Intel[®] Xeon[®] Outperforms Competition Around The Clock
 - AI Case Studies
- Product Availability
- > Introducing Intel[®] Xeon[®] 6 Processor
- Call to Action
- Resources

Why Partner With Intel?

At Intel, our goal is to improve lives and outcomes for everyone and every enterprise on this planet

But we aren't doing this alone!

Together with our partners, we are creating real value for our customers by **bringing Al everywhere** and minimizing the risks in Al solution deployment

<u>Unlock Business Transformation in a Digital-First Economy:</u> <u>Become an Artificial Intelligence Disruptor</u>



When you partner with Intel, you partner with a complete AI ecosystem

Our broad portfolio of Al-enabling technologies and collaboration with hardware, software, and solution ecosystem partners delivers real world solutions and differentiated business outcomes for industries, companies, and communities.

Helping you to grow your business.

Intel Leads the Way in Al

More than

100M

More than **300**

Al-accelerated ISV features throughout 2024¹ processors with Al accelerators through 2025¹ Intel[®] Xeon[®] install base of

provisions Al workloads alongside other workloads²

Join Us On the Journey to Bring AI Everywhere

¹ <u>https://www.intel.com/content/www/us/en/products/docs/processors/core-ultra/ai-pc-acceleration.html</u> ² <u>AI Intel® Xeon® Sales Card</u>

Bringing Al Everywhere

In today's hypercompetitive environment, enterprises that embrace AI are pulling ahead.

Intel infrastructure is engineered for enterprise AI, empowering you to maximize your investments and realize your vision at a lower cost. And, with enterpriseready solutions and open, optimized software, you can go to market fast, even with sensitive and regulated data.

It's time to think differently about enterprise AI.

- Partner Guide: Assessing Today's Enterprise AI Opportunity Landscape
- Infographic: **Optimize Enterprise AI Results with Intel**

LEARN MORE

- Enterprise AI / Generative AI Partner Enablement Package
- Al Partner Enablement Package

Model Creation



4

intel

Why Choose Intel[®] Xeon[®] for Developing an AI Solution



1 Forbes: https://www.forbes.com/sites/gilpress/2019/11/22/top-artificial-intelligence-ai-predictions-for-2020-from-idc-and-forrester/#4fef9821315a

2 Source: VMware: https://www.vmware.com/files/pdf/VMware-Corporate-Brochure-BR-EN.pdf

The Al Hierarchy: Mapping ML, Deep Learning, and GenAl with Intel

Discover how Intel® processors fuel AI workloads across inference, training, and next-gen generative AI applications



6

Intel® products: The Best Solutions for AI Transformation

Use Intel products to avoid high costs and long wait times while adopting AI

5th Gen and Intel[®] Xeon[®] 6 P-core Processors are the Best CPU for AI

Leverage unique <u>Intel® AMX</u> and MRDIMMs

- Reduce TCO up to 74 percent¹
- Consolidate servers up to 10:1²
- Realize up to 1.38x superior LLM inferencing performance vs. AMD³



Use Intel[®] Gaudi[®] accelerators to right size your acceleration deployments

Accelerate the most demanding dedicated inferencing workloads with the best price/performance ratio vs. NVIDIA H100⁴

intel. Gaudi

Add Intel[®] Ethernet products to enhance performance and scale out

Optimize performance for AI and high-performance computing (HPC) applications, reduce TCO and capital expenditures (CapEx) by up to 20% for massive AI clusters⁵

Consume less power vs. NVIDIA Mellanox ConnectX-6⁶ intel. ethernet

Get the best Al performance out of the box with the most common ML and DL tools. Run Al on sensitive data with confidence using <u>Intel® TDX</u>, <u>Intel® SGX</u>, and hardened security ecosystem

¹See [9710] at <u>intel.com/processorclaims</u>; Intel[®] Xeon[®] 6. Results may vary. ²See [7721] at <u>intel.com/processorclaims</u>: Intel[®] Xeon[®] 6. Results may vary. ³See [9A231] at <u>intel.com/processorclaims</u>: Intel[®] Xeon[®] 6 vs. AMD EPYC Turin. Results may vary. ⁴See <u>intel.com/processorclaims</u>: Intel® Gaudi® Al Accelerator. Results may vary. ⁵See [9T3] at <u>intel.com/processorclaims</u>: Intel® Xeon® 6. Results may vary. ⁶See Massed Compute FAQ, "<u>Power Consumption Differences Between NVIDIA Mellanox InfiniBand Adapters and Ethernet Ad for NVIDIA A40 GPUs</u>"



Scalable Systems for Al

From Cloud & Data Center to the Edge, Intel[®] Xeon[®] processors provides optimized performance, scale and efficiency at a cost-effective price



Selling Intel[®] Al Hardware: A Conversation Guide

Intel[®] AI Software is Enterprise Ready



Intel[®] Xeon[®] Processors

ACCESS NOW >

- Sales Play Pitch: Power Your AI Transformation with Intel
- Sales Play Pitch: Modernize Your Data Center to Achieve
 Unprecedented Levels of Performance and Efficiency

Report: CPUs are Key to Enterprise AI

Intel[®] Xeon[®] - The Processor Designed for AI

intelEfficiently run Al inference5th Gen Intel® Xeon® processor	Build and deploy Al everywhere Intel Al software suite of optimized open-source frameworks and tools	Copen Ecosystem Extensive Intel AI products and partnership
The flexibility of Intel® Xeon® with the built-in DL performance of an Al accelerator	Enables out of the box AI performance and E2E productivity	Accelerate end customer time to market
 Up to 29% higher training and up to 42% higher inference performance than our previous generation¹ 	 5x improvement on GPT-J in 10 weeks through software optimizations alone³ 	READ MORE Product Brief ***XeON
 Up to 2.69x higher performance than AMD EPYC 9654 (96C) and 9754 (128C) processors² 	 Optimizing larger models up to 70B parameters to meet customer SLAs Optimized 300+ DL models and 50+ ML and Graph Models 	<section-header><section-header> Stit Genitation Processors</section-header></section-header>

 Based on performance gains of 1.1x to 1.29x for training (ResNet50v1.5, BERT-Large, SSD-ResNet34, RNN-T, MaskRCNN, and DLRM) and 1.19x to 1.42x for inference (ResNet50v1.5, BERT-Large, SSD-ResNet34, RNN-T (BF16 only), Resnext101 32x16d, MaskRCNN (BF16 only), DistilBERT) compared to 4th Gen Intel® Xeon® processor. See A15-A16 at intel.com/processorclaims: 5th Gen Intel® Xeon® Scalable processors. Results may vary.
 Based on performance gains of 1.19x to 2.69x with Intel® Advanced Matrix Extensions (Intel® AMX) for inference on GPT-J, LLaMA-2 13B, DLRM, DistilBERT, BERT-Large, and ResNet50v1.5 compared to AMD EYPC 9654 and 9754. See A201, A202, A208-A211 at intel.com/processorclaims: 5th Gen Intel Xeon

. Based on performance gains of 1.19x to 2.69x with Intel® Advanced Matrix Extensions (Intel® AMX) for inference on GPT-J, LLaMA-2 13B, DLRM, DistilBERT, BERT-Large, and ResNet50v1.5 compared to AMD EYPC 9654 and 9754. See A201, A202, A208-A211 at intel.com/processorclaims: 5th Gen Intel Xeor Scalable processors. Results may vary..

intel

Intel[®] Xeon[®] Processor delivers TCO Value for Mixed General-Purpose and Al Workloads

Case Study: Video Conferencing Service



Accelerate Al Workloads with Intel® Advanced Matrix Extensions (Intel® AMX)

Intel AMX is a **built-in accelerator** that enables 4th and 5th Gen Intel[®] Xeon[®] processors to optimize deep learning (DL) training and inferencing workloads. With Intel[®] AMX, 4th and 5th Gen Intel[®] Xeon[®] processors can quickly pivot between optimizing general computing and AI workloads.



compared to 3rd Gen Intel® Xeon® processors¹

READ MORE

Get Started with Intel[®] AMX

<u>Al framework</u> optimizations Tuning Guide

<u>Quick-start Guide</u>

<u>Al reference</u>

kits

USE CASE EXAMPLES

Recommender systems

Deliver a customized end-user experience, whether recommending movies and books or showing targeted ads. Create a DL-based recommender system that accounts for real-time user behavior signals and context features such as time and location.

Natural Language Processing

With a global market projected to reach 80.68 billion USD by 2026,11 NLP applications, including chatbots and sentiment analysis, are critical for businesses to support and scale various functions, including sentiment analysis, chatbots, and machine translation.

Retail e-commerce software solutions

Grow revenue and deliver an exceptional customer experience by minimizing transaction time and effortlessly handling peak demands with DL inference and training, in addition to AI-optimized frameworks like PyTorch and TensorFlow.

¹See linked Solution Brief above for configurations. Results may vary.

Drive Revenue Growth and Improve Customer Experience with Faster, More Effective Al

AMD Benchmarks

Leadership performance with the world's best CPU for AI

READ <u>MORE</u> Leadership Performance

5th Generation Intel[®] Xeon[®] processors with Intel[®] Advanced Matrix Extensions (Intel[®] AMX) Outperform AMD EPYC¹

5th Gen Intel[®] Xeon[®] delivers up to Save up to 883.0 Higher BERT Large Performance than 4th Gen AMD EPYC¹ Reduce TCO Across Your Server Fleet for Al Fewer servers to manage Upto F Lower TCO than 4th Gen AMD EPYC while Servers with 5th Gen Servers with 4th Gen running a BERT Large Intel[®] Xeon[®] Scalable **AMD EPYC** workload¹ processors¹ processors¹ **READ MORE**

SOLVE COMMON PROBLEMS

Better inform business decisions to drive revenue growth

Reduce repetitive tasks, costs, and time for your business

Improve customer retention and acquisition

Faster analysis for large amounts of data

Enable more responsive smart assistants and chatbots

Improve text prediction speed and accuracy

¹See linked AMD Benchmarks paper above for configurations. Results may vary.

Performance Data for Intel[®] AI Data Center Products



5th Gen Intel® Xeon® TCO advantages over AMD

A comparison against 50 4th Gen AMD EPYC 9554 servers

	Web NGINX TLS	Data Services RocksDB	Data Services MySQL	HPC Monte Carlo	AI - NLP DistilBERT
5 th Gen Intel® Xeon® Servers	31 servers	31 servers	30 servers	28 servers	15 servers
Fleet Energy Saved*	489.7 MWh	1218.1 MWh	684.0 MWh	585.8 MWh	1496.5 MWh
Reduced CO2 Emissions*	207,611 kg	516,402 kg	289,967 kg	248,352 kg	634,428 kg
TCO Savings*	\$444K	\$471K	\$509K	\$561K	\$1,300K
TCO Delta	21% savings	22% savings	24% savings	27% savings	62% savings

intel

Al Case Studies on Intel® Xeon® Processors

REAL WORLD RESULTS

Healthcare

Winning Health has introduced the WiNGPT solution based on 5th Gen Intel® Xeon® Scalable processors, through working with Intel, the inference performance has been increased by over 3X compared with the platform based on the 3rd Gen Intel® Xeon® Scalable processors



READ ARTICLE

Media & Entertainment

Gunpowder accelerated rendering times for stunning visual effects while lowering costs with as much as 52% better performance per dollar compared to previous-gen instances with Intel[®] Xeon[®] processors³

GUNPOWDER®

READ THE CASE STUDY

READ THE

ARTICLE

Professional Services

Ropers Majeski increased worker productivity by **18.5%**, saving an average of **75 minutes** per user per day by automating email processing, document filing, and report generation with builtin Al acceleration from Intel[®] Xeon[®] CPUs⁵

MAJESKI

READ THE CASE STUDY

Energy

Storm Reply chose the new Amazon EC2 C7i instances supported by 4th Gen Intel® Xeon® Scalable processors and Intel libraries for LLM modeling. After a HW evaluation process, they matched the priceperformance ratio of GPU-based options by using CPU-based instances. **Netflix** delivered fast and seamless streaming experiences with 2x better Al-enabled video encoding and significant cloud savings by upgrading AWS EC2 instances. Netflix achieved a 3.5x performance improvement per CPU with Intel® Xeon® CPUs and software optimizations, at a lower cost than with GPUs⁴

Retail

Meituan uses vision AI services to **improve a** wide range of customer experiences, and achieved **70% cost savings** by migrating from GPUs to Intel[®] Xeon[®] CPUs and software for AI inference⁶

美团 Meituan

READ THE CASE STUDY

READ <u>ARTICLE</u>

Introducing Intel[®] Xeon[®] 6 Processors



P-core Optimized for performance in compute-intensive and Al workloads **Common** platform foundation and **shared** software stack

E-core

Optimized for **efficiency** in high-density and scale-out workloads

Intel[®] Xeon[®] 6 Processors

The best processors to meet diverse performance and efficiency requirements



Intel[®] Xeon[®] 6 processor with P-cores AI | HPC | IaaS | General Compute

Intel[®] Xeon[®] 6 processors with P-cores

- Industry-leading Performance-cores (P-cores) are architected for compute-intensive workloads which benefit from multiple data elements being processed in parallel
- Choose from a range of SKUs with up to 128 cores and 12 memory channels for higher overall performance
- Maximize data throughput with the latest DDR5 and Multiplexed Combined Rank (MCR) DIMMs
- Scale AI everywhere with Intel Advanced Matrix Extensions (Intel AMX) to accelerate inferencing for INT8, BF16, and newly supported FP16 datatypes



2x

higher Al inference performance vs. 5th Gen Intel® Xeon® processors¹

Up to

higher HPC performance vs. 5th Gen Intel® Xeon® processors¹

SOLUTION BRIEF >

Intel® Xeon® 6 with P-cores for the Cloud

PRODUCT BRIEF >

Intel[®] Xeon[®] 6 Processors with Performance-Cores (P-Cores) Deep Dive Intel[®] Xeon[®] 6700 with P-cores higher average performance for general compute vs. 5th Gen Intel® Xeon® processors¹

Intel[®] Xeon[®] 6 with Performance Cores (P-cores) Server Consolidation



With more cores, double the memory bandwidth, and AI acceleration in every core, Intel[®] Xeon[®] 6 processors with P-cores provide **twice the performance** for the widest range of workloads, including AI and highperformance computing (HPC).¹

Lower your total cost of ownership (TCO) by migrating from 2nd Gen Intel® Xeon® processors (4208) to Intel® Xeon® 6 processors with Pcores (6952P).¹

Recover your cost in 4 months¹

Intel[®] Xeon[®] 6 with Efficient-cores (E-cores)

66 racks

Intel[®] Xeon[®] 6700E

Free up space and power in the data center for new AI projects

200 racks 2nd Gen Intel® Xeon® Processor







80k MWh

Fleet energy saved

34k mt Reduced CO2 emissions

See [7T2] at intel.com/processorclaims: Intel® Xeon® 6. Your costs and results may vary.





DEEP DIVE >

Modernize for Datacenter Transformation

Make room to expand your Al infrastructure. Consolidate your data center with Intel® Xeon® processors to save space, reduce costs, and take on new workloads.

	From 2 nd Gen Intel Xeon to Intel Xeon 6700 with P-cores			From 2 nd Gen Intel Xeon to Intel Xeon 6900 with P-cores					
	8260>6760P	6230>6787P	5218R>6760P		8260>6979P	6240>6980P	5220>6972P	4208>6952P	
Reduce servers	50 to 14 (72%)	50 to 10 (80%)	50 to 12 (76%)		50 to 8 (84%)	50 to 6 (88%)	50 to 6 (88%)	50 to 3 (94%)	
Reduce energy and CO_2	48%	51%	44%		56%	69%	59%	77%	
Reduce TCO*	51%	60%	57%		62%	64%	72%	87%	
Recover costs (months)	15	11	13		14	11	9	4	
	From 2 nd Gen Intel Xeon to Intel Xeon 6700 with E-cores				From 2 nd Gen Intel Xeon to 5th Gen Intel Xeon				
	6252>6710E	6252>6740E	6252>6780E		8260>8558	6230>8558	5218>6538Y+	4214>5520+	
Reduce servers	50 to 21 (58%)	50 to 15 (70%)	50 to 11 (78%)		50 to 18 (64%)	50 to 15 (70%)	50 to 16 (68%)	50 to 15 (70%)	
Reduce energy and CO_2	45%	53%	56%		33%	30%	47%	38%	
Reduce TCO*	43%	53%	55%		45%	53%	55%	61%	
Recover costs (months)	14	13	16		15	12	9	6	

*Based on US energy costs, estimated over 4 years

Performance varies by use, configuration and other factors. | Please reference Intel Node TCO & Power Calculator

Modernize the Datacenter to Run Al Workloads More Efficiently





*Based on US energy costs, estimated over 4 years Intel Xeon 6 leveraging MRDIMMs See [T8, T7, 7T21] intel.com/processorclaims: 5th gen and Intel Xeon 6. Results may vary

Reduce Your Energy Consumption and Costs on New Server Purchases for AI

Intel® Xeon® processors deliver a significant performance advantage and lower TCO compared to AMD EPYC servers on AI workloads.

	Intel Xeon 6700P vs. AMD EPYC 9535 (Turin)		Intel Xeor vs. AMD EPYC 9755/9	6900P 9654 (Turin/Genoa)	5th Gen Intel Xeon vs. AMD EPYC 9554 (Genoa)	
	Vision Transformation INT8 (Batched)	Stable Diffusion INT8 (Realtime)	Recommendation systems DLRM	Image classification ResNet-50	Natural language processing BERT-Large	Natural language processing DistilBERT
Performance advantage (per server)	2.09x	1.53x	1.87x	4.91x	2.22x	3.49x
Server consolidation	289 to 140	204 to 130	170 to 90	496 to 100	50 to 23	50 to 15
Fleet energy saved	5,788 MWh	2,011 MWh	5,497 MWh	16,717 MWh	1,205 MWh	1,496 MWh
Reduced CO ₂ emissions	2,454 metric tons	853 metric tons	2,330 metric tons	7,087 metric tons	5,109 metric tons	6,344 metric tons
TCO savings (estimated over 4 years)	\$6.3M	\$2.6M	\$4.0M	\$18.8M	\$883K	\$1.3M
% TCO savings*	Up to 52%	Up to 31%	Up to 46%	Up to 74%	Up to 41%	Up to 62%

*Based on US energy costs, estimated over 4 years Intel Xeon 6 leveraging MRDIMMs See [7T221, 7T222, 9T222, 9T10, T205, T206] intel.com/processorclaims: 5th gen and Intel Xeon 6. Results may vary

Access the Intel® Xeon® Processor Advisor Suite to calculate the best route to lower TCO and ROI

Ecosystem Enablement for AI on 5th Gen Intel[®] Xeon[®] and Intel[®] Xeon[®] 6 systems Solutions, OEMs, and foundational software

	5 th Gen Intel® Xec	utions			
OEM	ululu cisco		Lenovo		Red Hat
Deep Learning	Inferencing on Cisco UCS X-Series Inference Operations on UCS X-Series Cisco UCS X-Series M7 Blade Servers Creativity with GenAl	<u>Dell PowerEdge R760,</u> <u>Cloudera</u>		<u>NLF</u> <u>Rec</u> <u>[</u>	⁹ applications on I Hat OpenShift I Hat OpenShift Data Science
GenAl	<u>GenAl Advancements: Dell PowerEdge R760</u> <u>UCS X-Series on Red Hat OpenShift Al</u> <u>Mainstreaming GenAl Inference Operations</u> <u>GenAl Inferencing with UCS X-Series M7</u> <u>Creativity with GenAl</u>	<u>Dell PowerEdge R760</u>	<u>RH OpenShift L</u> ThinkSystem SR	<u>enovo</u> <u>650 V3</u>	
RAG		Evaluating RAG with OpenVINO			
Intel re	ecommended SKUs for Al	CPU	Good	Better	Best
Al Worl	cloads:	5th Gen Intel® Xeon®	8562Y+	8568Y+	8592+

	Linux Kernel
For specific versions check with the ven	dor as they become publicly annot
Available Al systems	
5 th Gen Intel Xeon	Intel Xeon 6
Cisco M7 c220, M7 c240, M7 x210c	Cisco M8 c220, M8 c240, M8 x210c
Dell R660, R660XS, R760, R760XA, R760XD2, R760XS, XR5610, XR7620, Xeon as Host CPU: XE9680	Dell R670, R770, R470, R570, Xeon as CPU: XE9780, XE7740
HPE DL380G11, DL380aG11, SY480G11, ML350G11	HPE DL580G12, DL380G12, DL380aG DL360G12, DL340G12, DL320G1 SY480G12, ML350G12
Lenovo ThinkSystem SR680a V3, HG660x V3, WA5480 G3 (China), WA7780 G3 (China)	Lenovo ThinkSystem SR680a V4, HG660x V4, WA548 (China), SR650a V4
Supermicro SYS-82IGE-TNHR, SYS-A2IGE- NBRT, SYS-82IGE-TNMR2, SYS- 42IGU-TNXR, SYS-42IGE-NBRT- LLC, SYS-42IGE-TNHR2-LCC, SYS- 22IGE-TNHT-LLC	Supermicro SYS-A22GA-NBRT, SYS-422GA NBRT-LCC, SYS-822GA-NBRT, S 822GA-NGR3, SYS-522GA- NRT, SYS-422GL-NR

Foundational Software Enablement

Operating System

SUSE Linux Enterprise Server

nced.

MS Windows Server Red Hat Enterprise Linux

Ubuntu

Hypervisor

Microsoft Hyper-V WS

VMware ESXi

Linux

KVM

intel. 26/26

YS

Small and Medium LLM	2	Computer Vision
Classical ML (Non-DL)		NLP

DLRM

Classical ML (Non-DL)RAG

5th Gen Intel® Xeon®	8562Y+ 8548Y+	8568Y+ 8558P	8592+
Intel® Xeon® 6 with P-Cores	6960P 6737P 6730P	6740P 6747P	6787P

Call to Action



- Understand your partner / customers' usage model and performance needs FIRST
- When AI is just another workload in a mixed general purpose and AI environment, lead with the Intel[®] Xeon[®] processors that are already running your customers' business
- For dedicated AI deployments, Intel[®] Xeon[®] processors paired with Intel[®] Gaudi[®] AI Accelerators will deliver the optimal TCO

Access the <u>Intel[®] Xeon[®] Processor Advisor Suite</u> to calculate the best route to lower TCO and path to ROI for your partners

Additional Resources

AssetType	Title and Link
Pitch Cards	5th Gen Intel® Xeon® Pitch Cards PUBLIC SET
Business Brief	Performance. Efficiency. Security. All Across Your Clouds

Al Enablement Zones

<u>Access a comprehensive resource hub</u> designed to help grow your business and solve your customers' most pressing business challenges. Find exclusive, value-added technical and sales enablement resources to help you build and sell solutions with Intel technology.



Technical Enablement

Sales & Marketing Enablement

Technical Enablement

Sales & Marketing Enablement

Technical Enablement

Sales & Marketing Enablement

Legal Notices and Disclaimers

Notices and Disclaimers.

© 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.





Al Use Case Benchmarks on 5th Gen Intel[®] Xeon[®]

Generative AI

GPT-J

- Up to 2.3x performance speedup and 1.81x higher performance/watt with 5th Gen Intel Xeon processor vs 3rd Gen Intel® Xeon on GPT-J first token latency (int8)¹
- Up to 1.64x performance speedup and 1.30x higher performance/watt with 5th Gen Intel® Xeon processor vs 3rd Gen Intel® Xeon on GPT-J second token latency (int8)¹

LLaMA

- Up to 2.1x performance speedup and 1.58x higher performance/watt with 5th Gen Intel Xeon processor vs 3rd Gen Intel® Xeon® on Llama 2 13B first token latency (int8)²
- Up to 1.48x performance speedup and 1.11x higher performance/watt with 5th Gen Intel® Xeon® processor vs 3rd Gen Intel® Xeon® on Llama 2 13B second token latency (int8)²

Recommender Systems

DIEN

Up to 1.12x higher end-to-end recommendation Speedup with 5th Gen Intel® Xeon® processor compared with prior generation on Deep Interest Evolution Network (DIEN) Recommendation System⁵

DLRM

Up to 8.7x higher batch Recommendation System inference performance (DLRM) and 6.2x higher performance/watt on 5th Gen Intel® Xeon® Platinum 8592+ with AMX BF16 vs. 3rd Gen Intel® Xeon® processor with FP32⁶

DLRM (Competitive)

- Up to 2.34x higher batched Recommendation System inference performance (DLRM) and 2.26x higher performance/watt on 5th Gen Intel® Xeon® Platinum 8592+ with AMX INT8 vs. AMD EPYC 9654 (Genoa)⁷
- Up to 1.90x higher batched Recommendation System inference performance (DLRM) and 1.72x higher performance/watt on 5th Gen Intel® Xeon® Platinum 8592+ with AMX INT8 vs. AMD EPYC 9754 (Bergamo)⁷

Natural Language Processing

BERT-Large

- Up to 9.9x higher real-time Natural Language Processing inference (BERT-large) performance and 7.7x higher performance/watt on 5th Gen Intel® Xeon® Platinum 8592+ with AMX BF16 vs. 3rd gen Intel® Xeon® processor with FP32³
- Up to 10x higher batch Natural Language Processing inference (BERT-large) performance and 7.1x higher performance/watt on 5th Gen Intel® Xeon® Platinum 8592+ with AMX BF16 vs. 3rd gen Intel® Xeon® processor with FP32³

DistilBERT

- Up to 7x higher real-time Natural Language Processing inference (DistilBERT) performance and 5.6x higher performance/watt on 5th Gen Intel[®] Xeon[®] Platinum 8592+ with AMX BF16 vs. 3rd gen Intel[®] Xeon[®] processor with FP32⁴
- Up to 10x higher batch Natural Language Processing inference (DistilBERT) performance and 8x higher performance/watt on 5th Gen Intel® Xeon® Platinum 8592+ with AMX BF16 vs. 3rd gen Intel® Xeon® processor with FP32⁴

Computer Vision

ResNet-50

- Up to 8x higher performance and 6.4x higher performance/watt on real-time image classification(ResNet50) inference (BF16) with 5th Gen Intel® Xeon® 8592+ (64c) vs. 3rd Gen Intel® Xeon® 8380 (fp32)⁸
- Up to 8.8x higher performance and 7.4x higher performance/watt on batch image classification(ResNet50) inference (BF16) with 5th Gen Intel® Xeon® 8592+ (64c) vs. 3rd Gen Intel® Xeon® 8380 (fp32) ⁸

ResNet-34

- Up to 14x higher real time object detection inference performance (SSD-ResNet34) and 9.6x higher performance/watt on 5th Gen Intel® Xeon® Platinum 8592+ with AMX BF16 vs. 3rd Gen Intel® Xeon® processor with FP32⁹
- Up to 13.8x higher batch object detection inference performance (SSD-ResNet34) and 10x higher performance/watt on 5th Gen Intel® Xeon® Platinum 8592+ with AMX BF16 vs. 3rd Gen Intel® Xeon® processor with FP32⁹

Machine Learning

Up to 1.39x faster end-to-end Census workload performance using 5th Gen Intel® Xeon® processor compared to prior generation(FP32)¹⁰

^{1,2} See [A1, A2] at intel.com/processorclaims: 5th Gen Intel® Xeon® Scalable processors. Results may vary ^{3,4} See [A19, A24] at intel.com/processorclaims: 5th Gen Intel® Xeon® Scalable processors. Results may vary ^{5,6,7} See [A7, A20, A208] at intel.com/processorclaims: 5th Gen Intel® Xeon® Scalable processors. Results may vary ^{8,9,} See [A26, A21] at intel.com/processorclaims: 5th Gen Intel® Xeon® Scalable processors. Results may vary

¹⁰.See [A8] at intel.com/processorclaims: 5th Gen Intel® Xeon® Scalable processors. Results may vary

intel

Configuration details: Al performance

Performance varies by use, configuration and other factors. See backup for configuration details. Results may vary.

ResNeXT101_32x16d Inference: BS1, BSx

8592+: 1-node, 2x Intel(R) Xeon(R) Platinum 8592+, 64 cores, HT On, Turbo On, NUMA 4, Total Memory 1024GB (16x64GB DDR5 5600 MT/s]), BIOS EGSDCRBI.SYS.0105.D74.2308261933, microcode 0x21000161, 1x Ethernet Controller 1225-LM, 1x 3.6T INTEL SSDPE2KX040T8, 1x 931.5G INTEL, CentOS Stream 9, 6.2.0-emr.bkc.6.2.3.6.31.x86_64. PyTorch: torch:2.1.0.dev20230825+cpu, oneDNN:v3.2.1 IPEX:2.1.0+git31b5ee1, Test by INTEL as of 09/07/2023.

8480+: 1-node, 2x Intel(R) Xeon(R) Platinum 8480+, 56 cores, HT On, Turbo On, NUMA 2, Total Memory 1024GB (16x64GB DDR5 4800 MT/s], BIOS EGSDCRB1.SYS.0102.D37.2305081420, microcode 0x2b0004b1, 1x Ethernet Controller 1225-LM, 1x 7.2G Cruzer Blade, 2x 3.7T INTEL SSDPE2KX040T8, 1x 894.3G INTEL SSDSC2KG96, CentOS Stream 8, 5.15.0-spr.bkc.pc.16.4.24.x86_64. PyTorch: torch:2.1.0.dev20230825+cpu oneDNN:v3.2.1 IPEX:2.1.0+git31b5ee1. Test by INTEL as of 09/05/2023.

ResNeXT101_32x16d, BS=1: 4cores/instance, BS=x: 1 instance/numa node; Resnext101: ImageNet

RNN-T Inference : BS1, BSx

8592+: 1-node, 2x Intel(R) Xeon(R) Platinum 8592+, 64 cores, HT On, Turbo On, NUMA 4, Total Memory 1024GB (16x64GB DDR5 5600 MT/s [5600 MT/s]), BIOS EGSDCRBI.SYS.0105.D74.2308261933, microcode 0x21000161, 1x Ethernet Controller 1225-LM, 1x 3.6T INTEL SSDPE2KX040T8, 1x 931.5G INTEL, CentOS Stream 9, 6.2.0-emr.bkc.6.2.3.6.31.x86_64. PyTorch: torch:2.1.0.dev20230825+cpu, oneDNN:v3.2.1 IPEX:2.1.0+git31b5ee1, Test by INTEL as of 09/07/2023.

8480+: 1-node, 2x Intel(R) Xeon(R) Platinum 8480+, 56 cores, HT On, Turbo On, NUMA 2, Total Memory 1024GB (16x64GB DDR5 4800 MT/s] (4800 MT/s]), BIOS EGSDCRB1.SYS.0102.D37.2305081420, microcode 0x2b0004b], 1x Ethernet Controller 1225-LM, 1x 7.2G Cruzer Blade, 2x 3.7T INTEL SSDPE2KX040T8, 1x 894.3G INTEL SSDSC2KG96, CentOS Stream 8, 5.15.0-spr.bkc.pc.16.4.24.x86_64. PyTorch: torch:2.1.0.dev20230825+cpu oneDNN:v3.2.1 IPEX:2.1.0+git31b5ee]. Test by INTEL as of 09/05/2023.

RNN-T, BS=1: 4cores/instance, BS=x: 1 instance/numa node; RNNT: LibriSpeech

DistilBERT Inference : BS1, BSx

8592+: 1-node, 2x INTEL(R) XEON(R) PLATINUM 8592+, 64 cores, HT On, Turbo On, NUMA 2, Total Memory 1024GB (16x64GB DDR5 5600 MT/s]), BIOS 2.0, microcode 0x21000161, 2x Ethernet Controller X710 for 10GBASE-T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.2 LTS, 5.15.0-78-generic, Test by Intel as of 10/10/23.

8480+: 1-node, 2x Intel(R) Xeon(R) Platinum 8480+, 56 cores, HT On, Turbo On, NUMA 2, Total Memory 1024GB (16x64GB DDR5 4800 MT/s [4800 MT/s]), BIOS 2.0, microcode 0x2b0004d0, 2x Ethernet Controller X710 for 10GBASE-T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.2 LTS, 5.15.0-78-generic, Test by Intel as of 10/25/23.

Software configuration: DistilBERT, Intel Model Zoo:https://github.com/IntelAl/models, gcc=12.3, OneDNN3.2, Python 3.9, PyTorch 2.0, IPEX 2.0, Transformer version 4.18.0, physical cores only.

MaskRCNN Inference: : BS1, BSx

8592+: 1-node, 2x Intel(R) Xeon(R) Platinum 8592+, 64 cores, HT On, Turbo On, NUMA 4, Total Memory 1024GB (16x64GB DDR5 5600 MT/s]), BIOS EGSDCRBI.SYS.0105.D74.2308261933, microcode 0x21000161, 1x Ethernet Controller 1225-LM, 1x 3.6T INTEL SSDPE2KX040T8, 1x 931.5G INTEL, CentOS Stream 9, 6.2.0-emr.bkc.6.2.3.6.31.x86_64. PyTorch: torch:2.1.0.dev20230825+cpu, oneDNN:v3.2.1 IPEX:2.1.0+git31b5ee1, Test by INTEL as of 09/07/2023.

8480+: 1-node, 2x Intel(R) Xeon(R) Platinum 8480+, 56 cores, HT On, Turbo On, NUMA 2, Total Memory 1024GB (16x64GB DDR5 4800 MT/s], BIOS EGSDCRB1.SYS.0102.D37.2305081420, microcode 0x2b0004b1, 1x Ethernet Controller 1225-LM, 1x 7.2G Cruzer Blade, 2x 3.7T INTEL SSDPE2KX040T8, 1x 894.3G INTEL SSDSC2KG96, CentOS Stream 8, 5.15.0-spr.bkc.pc.16.4.24.x86_64. PyTorch: torch:2.10.dev20230825+cpu oneDNN:v3.2.1 IPEX:2.1.0+git3lb5ee1. Test by INTEL as of 09/05/2023.

MaskRCNN, BS=1: 4cores/instance, BS=x: 1 instance/ numa node; Mask RCNN: COCO 2017



Resources and Configurations



• 5th Gen Intel[®] Xeon[®] Outperforms Competition Around The Clock

- ResNet50v1.5
- Intel® Xeon® 8592+: 1-node, 2x Intel(R) Xeon(R) Platinum 8592+, 64 cores, HT On, Turbo On, NUMA 2, Total Memory 1024GB (16x64GB DDR55600 MT/s [5600 MT/s]), BIOS 2.0, microcode 0x21000161, 2x Ethernet Controller X710 for 10GBASE-T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.2 LTS, 5.15.0-78-generic. TensorFlow= Intel TF 2.13, OneDNN=3.2, Python 3.8, AI Model=ResNet50v1.5 Large(https://github.com/IntelAI/models/), Batched Results: best scores achieved using BFloat16, INT8-AMX (BS >1),, Test by INTEL as of 10/10/2023.
- AMD EPYC 9654: 1-node, 2x AMD EPYC 9654 96-Core Processor, 96 cores, SMT On, Turbo On, NUMA 2, Total Memory 1536GB (24x64GB DDR5 4800 MT/s [4800 MT/s]), BIOS 1.5, microcode 0xa10113e, 2x Ethernet Controller 10G X550T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.3 LTS, 5.15.125-0515125-generic, ZenDNN 4.1 TensorFlow 2.12.1, Python 3.8, AI Model=ResNet50v1.5 Large(https://github.com/IntelAI/models/), Batched Results: best scores achieved using (BS >1), Test by INTEL as of 09/11/23.

NGINXTLS

- Intel® Xeon® 8592+: 1-node, 2x 5th Gen Intel® Xeon® Scalable processor (64 core) with integrated Intel Quick Assist Technology (Intel QAT), QAT device utilized=4 (1 active socket), HT On, Turbo Off, SNC On, with 1024GB DDR5 memory (16x64 GB 5600), microcode 0x21000161, Ubuntu 22.04.3 LTS, 5.15.0-78-generic, 1x 1.7T SAMSUNG MZWLJ1T9HBJR-00007, 1x Intel® Ethernet Network Adapter E810-2CQDA2, 1x100GbE, NGINX Async v0.5.1, OpenSSL 3.1.3, IPP Crypto 2021.8, IPsec MB v 1.4, QAT _Engine v 1.4.0, QAT Driver 20.1.1.1.20-00030, TLS 1.3 Webserver: ECDHE-X25519-RSA2K, tested by Intel October 2023.
- AMD EPYC 9554: 1-node, AMD platform with 2x 4th Gen AMD EPYC processor (64 cores), SMT On, Core Performance Boost Off, NPS1, Total Memory 1536GB (24x64GB DDR5 4800), microcode 0xa10113e, Ubuntu 22.04.3 LTS, 5.15.0-78-generic, 1x 1.7T SAMSUNG MZWLJ1T9HBJR-00007, 1x Intel® Ethernet Network Adapter E810-2CQDA2, 1x100GbE, NGINX Async v0.5.1, OpenSSL 3.1.3, TLS 1.3 Webserver: ECDHE-X25519-RSA2K,tested by Intel October 2023.

Resources and Configurations



• 5th Gen Intel[®] Xeon[®] Outperforms Competition Around The Clock

- HammerDB Microsoft SQL Server + Backup
- Intel® Xeon® 8592+: 1-node, 2x 5th Gen Intel® Xeon® Scalable processor 8592+ (64 cores) with integrated Intel Quick Assist Technology (Intel QAT), Number of IAA device utilized=8(2 sockets active), HT On, Turbo On, SNC off, Total Memory 1024GB (16x64GB DDR5 5600), microcode 0x21000161, 2x Ethernet Controller 10-Gigabit X540-AT2, 7x 3.5T INTEL SSDPE2KE032T807, QATZip 2.0.W.1.9.0-0008, Microsoft Windows Server Datacenter 2022, Microsoft SQL Server 2022, SQL Server Management Studio 19.0.1, HammerDB 4.5, tested by Intel October 2023.
- AMD EPYC 9554: 1-node, AMD platform with 2x 4th Gen AMD EPYC processor (64 cores), SMT On, Core Performance Boost On, NPS1, Total Memory 1536GB (24x64GB DDR5 4800), microcode 0xa10113e, 2x Ethernet Controller 10G X550T, 7x 3.5T INTEL SSDPE2KE032T807, Microsoft Windows Server Datacenter 2022, Microsoft SQL Server 2022, SQL Server Management Studio 19.0.1, HammerDB 4.5, tested by Intel October 2023.
- RocksDB
- Intel® Xeon® 8592+: 1-node, 2x 5th Gen Intel® Xeon® Scalable processor 8592+ (64 cores) with integrated Intel In-Memory Analytics Accelerator (Intel IAA), Number of IAA device utilized=8(2 sockets active), HT On, Turbo On, SNC off, Total Memory 1024GB (16x64GB DDR5 5600), microcode 0x21000161, 2x Ethernet Controller 10-Gigabit X540-AT2, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.3 LTS, 6.5.0-060500-generic, QPL v1.2.0, accel-config-v4.0, iaa_compressor plugin v0.3.0, ZSTD v1.5.5, gcc 10.4.0, RocksDB v8.3.0 trunk (commit 62fc15f) (db_bench), 4 threads per instance, 64 RocksDB instances, tested by Intel October 2023.
- AMD EPYC 9554: 1-node, AMD platform with 2x 4th Gen AMD EPYC processor (64 cores), SMT On, Core Performance Boost On, NPS1, Total Memory 1024GB (16x64GB DDR5 4800), microcode 0xa10113e, 2x Ethernet Controller 10G X550T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.3 LTS, 6.5.0-060500-generic, ZSTD v1.5.5, gcc 10.4.0, RocksDB v8.3.0 trunk (commit 62fc15f) (db_bench), 4 threads per instance, 28 RocksDB instances, tested by Intel October 2023.
- Monte Carlo
- Intel® Xeon® 8592+: 1-node 2x Intel® Xeon® 8592+, HT On, Turbo On, SNC2, 1024 GB DDR5-5600, ucode 0x21000161, Red Hat Enterprise Linux 8.7, 4.18.0-425.10.1.el8_7.x86_64, Monte Carlo v1.2, cmkl:2023.2.0 icc:2023.2.0 tbb:2021.10.0. Test by Intel as of October 2023.
- AMD EPYC 9554: 1-node, 2x AMD EPYC 9554, SMT On, Turbo On, CTDP=360W, NPS=4, 1536GB DDR5-4800, ucode= 0xa101111, Red Hat Enterprise Linux 8.7, Kernel 4.18, Monte Carlo v1.2, cmkl:2023.2.0 icc:2023.2.0 tbb:2021.10.0. Test by Intel as of March 2023

Resources and Configurations



• 5th Gen Intel[®] Xeon[®] Outperforms Competition Around The Clock

- DLRM
- Intel® Xeon® 8592+: 1-node, 2x Intel(R) Xeon(R) Platinum 8592+, 64 cores, HT On, Turbo On, NUMA 2, Total Memory 1024GB (16x64GB DDR55600 MT/s [5600 MT/s]), BIOS2.0, microcode 0x21000161, 2x Ethernet Controller X710 for 10GBASE-T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.2 LTS, 5.15.0-78-generic. Framework=Pytorch 2.1, IPEX=2.1, Python 3.8, AI Model= DLRM(https://github.com/IntelAI/models/), Batched Results: best scores achieved using BS>1, Precision=INT8-AMX, Test by INTEL as of 10/10/2023.
- AMD EPYC 9654: 1-node, 2x AMD EPYC 9654 96-Core Processor, 96 cores, SMT On, Turbo On, NUMA 2, Total Memory 1536GB (24x64GB DDR5 4800 MT/s [4800 MT/s]), BIOS 1.5, microcode 0xa10113e, 2x Ethernet Controller 10G X550T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.3 LTS, 5.15.125-0515125-generic, Framework=Pytorch 2.1, IPEX=2.1, Python 3.8, AI Model= DLRM(https://github.com/IntelAI/models/), Batched Results: best scores achieved using BS>1, Precision=INT8. Test by INTEL as of 09/11/23.
- HammerDBMySQL
- Intel® Xeon® 8592+: 1-node, 2x INTEL(R) XEON(R) PLATINUM 8592+, 64 cores, HT On, Turbo On, NUMA 2, Integrated Accelerators Available [used]: DLB 8 [0], DSA 8 [0], IAX 8 [0], QAT 8 [0], Total Memory 1024GB (16x64GB DDR5 5600 MT/s [5600 MT/s]), BIOS 2.0, microcode 0x21000161, 2x Ethernet Controller X710 for 10GBASE-T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, 2x 1.7T SAMSUNG MZWLJ1T9HBJR-00007, Ubuntu 22.04.3 LTS, 5.15.0-84-generic, HammerDB Mv4.4, MySQL 8.0.33. Test by Intel as of 10/04/23.
- AMD EPYC 9554: 1-node, 2x AMD EPYC 9554 64-Core Processor, 64 cores, HT On, Turbo On, NUMA 2, Integrated Accelerators Available [used]: DLB 0 [0], DSA 0 [0], IAX 0 [0], QAT 0 [0], Total Memory 1536GB (24x64GB DDR5 4800 MT/s [4800 MT/s]), BIOS 1.5, microcode 0xa10113e, 2x Ethernet Controller X710 for 10GBASE-T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, 2x 1.7T SAMSUNG MZWLJ1T9HBJR-00007, Ubuntu 22.04.3 LTS, 5.15.125-0515125-generic, HammerDB v4.4, MySQL 8.0.33. Test by Intel as of 10/05/23.

Configurations: 5th Gen Intel[®] Xeon[®] TCO Advantages (1 of 6)

Claim: 5th Gen Intel® Xeon® delivers up to 24% lower TCO than the 4th Gen AMD Epyc while running a HammerDB MySQL OLTP database workload.

- Based on 5th Gen Intel[®] Xeon[®] delivers up to a 1.70x faster than the 4th Gen AMD Epyc while running a HammerDB MySQL OLTP workload. This performance drives a fleet reduction from 50 to 30
- servers which, over 4 years, saves: 684.0 MWH of energy, 289,967 kgCo2 emissions, and \$508.9k of cost.
- 1-node, 2x INTEL(R) XEON(R) PLATINUM 8592+, 64 cores, HT On, Turbo On, NUMA 2, Integrated Accelerators Available [used]: DLB 8 [0], DSA 8 [0], IAX 8 [0], QAT 8 [0], Total Memory 1024GB (16x64GB DDR5 5600 MT/s [5600 MT/s]), BIOS 2.0, microcode 0x21000161, 2x Ethernet Controller X710 for 10GBASE-T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, 2x 1.7T SAMSUNG MZWLJIT9HBJR-00007, Ubuntu 22.04.3 LTS, 5.15.0-84-generic, HammerDB Mv4.4, MySQL 8.0.33. Test by Intel as of 10/04/23.
- I-node, 2x AMD EPYC 9554 64-Core Processor, 64 cores, HT On, Turbo On, NUMA 2, Integrated Accelerators Available [used]: DLB 0 [0], DSA 0 [0], IAX 0 [0], QAT 0 [0], Total Memory 1536GB (24x64GB DDR5 4800 MT/s [4800 MT/s]), BIOS 1.5, microcode 0xa10113e, 2x Ethernet Controller X710 for 10GBASE-T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, 2x 1.7T SAMSUNG MZWLJ1T9HBJR-00007, Ubuntu 22.04.3 LTS, 5.15.125-0515125-generic, HammerDB v4.4, MySQL 8.0.33. Test by Intel as of 10/05/23.
- For a 50 server fleet of AMD EPYC 9554, estimated as of October 2023:
- CapEx costs: \$1.40M
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$755.1K
- Energy use in kWh (4 year, per server): 47654, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394
- For a 30 server fleet of 5th Gen Intel[®] Xeon[®] 8592+ as of October 2023
- CapEx costs: \$1.17M
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$480.0K
- Energy use in kWh (4 year, per server): 58625, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394
- Costs based on Intel estimates and information from thinkmate.com as of October 2023.
 Global Partners and Support
 Intel (

Intel Confidential

Configurations: 5th Gen Intel[®] Xeon[®] TCO Advantages (2 of 6)

Claim: 5th Gen Intel® Xeon® delivers up to 22% lower TCO than the 4th Gen AMD Epyc while running a RocksDB database workload.

- Based on 5th Gen Intel[®] Xeon[®] delivers up to a 1.62x faster than the 4th Gen AMD Epyc while running a RocksDB database workload. This performance drives a fleet
 a reduction from 50 to 31 servers which, over 4 years, saves: 1,218 MWH of energy, 516,402 kgCo2 emissions, and \$471.8k of cost.
- 1-node, 2x 5th Gen Intel[®] Xeon[®] Scalable processor 8592+ (64 cores) with integrated Intel In-Memory Analytics Accelerator (Intel IAA), Number of IAA device utilized=8(2 sockets active), HT On, Turbo On, SNC off, Total Memory 1024GB (16x64GB DDR5 5600), microcode 0x21000161, 2x Ethernet Controller 10-Gigabit X540-AT2, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.3 LTS, 6.5.0-060500-generic, QPL v1.2.0, accel-config-v4.0, iaa_compressor plugin v0.3.0, ZSTD v1.5.5, gcc 10.4.0, RocksDB v8.3.0 trunk (commit 62fc15f) (db_bench), 4 threads per instance, 64 RocksDB instances, tested by Intel October 2023.
- I-node, AMD platform with 2x 4th Gen AMD EPYC processor (64 cores), SMT On, Core Performance Boost On, NPSI, Total Memory 1024GB (16x64GB DDR5 4800), microcode 0xa10113e, 2x Ethernet Controller 10G X550T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.3 LTS, 6.5.0-060500-generic, ZSTD v1.5.5, gcc 10.4.0, RocksDB v8.3.0 trunk (commit 62fc15f) (db_bench), 4 threads per instance, 28 RocksDB instances, tested by Intel October 2023.
- For a 50 server fleet of AMD EPYC 9554, estimated as of October 2023:
- CapEx costs: \$1.36M
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$809.5K
- Energy use in kWh (4 year, per server): 58531, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394
- For a 31 server fleet of 5th Gen Intel® Xeon® 8592+ as of October 2023
- CapEx costs: \$1.21M
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$491.3K
- Energy use in kWh (4 year, per server): 55111, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394

Clobappartners and Btoppstimates and information from thinkmate.com as of Optoperrigentia

Configurations: 5th Gen Intel[®] Xeon[®] TCO Advantages (3 of 6)

Claim: 5th Gen Intel® Xeon® delivers up to 21% lower TCO than the 4th Gen AMD Epyc while running a NGINX TLS Handshake workload.

- Based on 5th Gen Intel[®] Xeon[®] delivers up to a 1.66x faster than the 4th Gen AMD Epyc while running a NGINX TLS Handshake workload. This performance drives a
 fleet a reduction from 50 to 31 servers which, over 4 years, saves: 489.7 MWH of energy, 207,611 kgCo2 emissions, and \$443.5k of cost.
- I-node, 2x 5th Gen Intel® Xeon® Scalable processor (64 core) with integrated Intel Quick Assist Technology (Intel QAT), Integrated Accelerators Available [used]: DLB 2 [0], DSA 2 [0], IAA 2 [0], QAT 2 [0], HT On, Turbo Off, SNC On, with 1024GB DDR5 memory (16x64 GB 5600), microcode 0x21000161, Ubuntu 22.04.3 LTS, 5.15.0-78-generic, 1x 1.7T SAMSUNG MZWLJ1T9HBJR-00007, 1x Intel® Ethernet Network Adapter E810-2CQDA2, 1x100GbE, NGINX Async v0.5.1, OpenSSL 3.1.3, IPP Crypto 2021.8, IPsec MB v 1.4, QAT _Engine v 1.4.0, QAT Driver 20.1.1.1.20-00030, TLS 1.3 Webserver: ECDHE-X25519-RSA2K, tested by Intel October 2023.
- I-node, 9554: I-node, AMD platform with 2x 4th Gen AMD EPYC processor (64 cores), SMT On, Core Performance Boost Off, NPS1, Total Memory I536GB (24x64GB DDR5 4800), microcode 0xa10113e, Ubuntu 22.04.3 LTS, 5.15.0-78-generic, 1x 1.7T SAMSUNG MZWLJ1T9HBJR-00007, 1x Intel® Ethernet Network Adapter E810-2CQDA2, 1x100GbE, NGINX Async v0.5.1, OpenSSL 3.1.3, TLS 1.3 Webserver: ECDHE-X25519-RSA2K,tested by Intel October 2023.
- For a 50 server fleet of AMD EPYC 9554, estimated as of October 2023:
- CapEx costs: \$1.41M
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$698.7K
- Energy use in kWh (4 year, per server): 36386, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394
- For a 31 server fleet of 5th Gen Intel® Xeon® 8592+ as of October 2023
- CapEx costs: \$1.21M
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$453.4K
- Energy use in kWh (4 year, per server): 42889, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394

Globas Partners and the petimates and information from thinkmate.com as of October 2002 Antial

Configurations: 5th Gen Intel[®] Xeon[®] TCO Advantages (4 of 6)

Claim: 5th Gen Intel® Xeon® delivers up to 27% lower TCO than the 4th Gen AMD Epyc while running Monte Carlo workload.

- Based on 5th Gen Intel[®] Xeon[®] delivers up to a 1.83x faster than the 4th Gen AMD Epyc while running a Monte Carlo workload. This performance drives a fleet a reduction from 50 to 28 servers which, over 4 years, saves: 585.8 MWH of energy, 248,352 kgCo2 emissions, and \$561.0k of cost.
- 1-node 2x Intel[®] Xeon[®] 8592+, HT On, Turbo On, SNC2, 1024 GB DDR5-5600, ucode 0x21000161, Red Hat Enterprise Linux 8.7, 4.18.0-425.10.1.el8_7.x86_64, Monte Carlo v1.2, cmkl:2023.2.0 icc:2023.2.0 tbb:2021.10.0. Test by Intel as of October 2023.
- 1-node, 2x AMD EPYC 9554, SMT On, Turbo On, CTDP=360W, NPS=4, 1536GB DDR5-4800, ucode= 0xa101111, Red Hat Enterprise Linux 8.7, Kernel 4.18, Monte Carlo v1.2, cmkl:2023.2.0 icc:2023.2.0 tbb:2021.10.0. Test by Intel as of March 2023
- For a 50 server fleet of AMD EPYC 9554, estimated as of October 2023:
- CapEx costs: \$1.36M
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$739.3K
- Energy use in kWh (4 year, per server): 44505, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394
- For a 28 server fleet of 5th Gen Intel® Xeon® 8592+ as of October 2023
- CapEx costs: \$1.09M
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$453.3K
- Energy use in kWh (4 year, per server): 58550, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394
- Costs based on Intel estimates and information from thinkmate.com as of October 2023.

Configurations: 5th Gen Intel[®] Xeon[®] TCO Advantages (5 of 6)

Claim: 5th Gen Intel® Xeon® delivers up to 46% lower TCO than the 4th Gen AMD Epyc while running a real-time Natural Language Processing inference (BERT-Large) workload.

- Based on 5th Gen Intel[®] Xeon[®] delivers up to a 2.44x faster than the 4th Gen AMD Epyc while running a real-time Natural Language Processing inference (BERT-Large) workload. This performance drives a fleet a reduction from 50 to 21 servers which, over 4 years, saves: 1231.2 MWH of energy, 521,941 kgCo2 emissions, and \$982.9k of cost.
- Intel® Xeon® 8592+: 1-node, 2x Intel(R) Xeon(R) Platinum 8592+, 64 cores, HT On, Turbo On, NUMA 2, Total Memory 1024GB (16x64GBDDR55600 MT/s [5600 MT/s]), BIOS 2.0, microcode 0x21000161, 2x Ethernet Controller X710 for 10GBASE-T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.2 LTS, 5.15.0-78-generic. Framework=Pytorch 2.0, IPEX=2.0, Python 3.8, AI Model=BERT Large(https://github.com/IntelAI/models/), INT8-AMX, Real Time (BS=1) results while maintaining 130ms latency SLA, Test by INTEL as of 10/10/2023.
- 9554: 1-node, 2x AMD EPYC 9554 64-Core Processor, 64 cores, SMT On, Turbo On, NUMA 2, Total Memory 1536GB (24x64GB DDR5 4800 MT/s [4800 MT/s]), BIOS 1.5, microcode 0xa10113e, 2x Ethernet Controller 10G X550T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.3 LTS, 5.15.125-0515125-generic, Framework=Pytorch 2.0, IPEX=2.0, Python 3.8, AI Model=BERT Large(https://github.com/IntelAI/models/), INT8, Real Time (BS=1) results while maintaining 130ms latency SLA. Test by INTEL as of 09/11/23.
- For a 50 server fleet of AMD EPYC 9554, estimated as of October 2023:
- CapEx costs: \$1.36
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$766.9K
- Energy use in kWh (4 year, per server): 50021, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394
- For a 21 server fleet of 5th Gen Intel® Xeon® 8592+ as of October 2023
- CapEx costs: \$801K
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$344.0K
- Energy use in kWh (4 year, per server): 60472, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394
- Costs based on Intel estimates and information from thinkmate.com as of October 2023.

Global Partners and Support

Intel Confidential

Configurations: 5th Gen Intel[®] Xeon[®] TCO Advantages (6 of 6)

Claim: 5th Gen Intel® Xeon® delivers up to 62% lower TCO than the 4th Gen AMD Epyc while running real-time Natural Language Processing inference (DistilBERT) workload.

- Based on 5th Gen Intel[®] Xeon[®] delivers up to a 3.49x faster than the 4th Gen AMD Epyc while running a real-time Natural Language Processing inference (DistilBERT) workload. This performance drives a fleet a reduction from 50 to 15 servers which, over 4 years, saves: 1496.5 MWH of energy, 634,428 kgCo2 emissions, and \$1,300k of cost.
- 1-node, 2x Intel(R) Xeon(R) Platinum 8592+, 64 cores, HT On, Turbo On, NUMA 2, Total Memory 1024GB (16x64GB DDR5 5600 MT/s [5600 MT/s]), BIOS 2.0, microcode
 0x21000161, 2x Ethernet Controller X710 for 10GBASE-T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.2 LTS, 5.15.0-78-generic. Framework=Pytorch 2.0, IPEX=2.0, Python 3.8, AI Model=DistilBERT (https://github.com/IntelAI/models/), INT8-AMX, Real Time (BS=1) results while maintaining 5ms latency SLA, Test by INTEL as of 10/10/2023.
- I-node, 2x AMD EPYC 9554 64-Core Processor, 64 cores, SMT On, Turbo On, NUMA 2, Total Memory 1536GB (24x64GB DDR5 4800 MT/s [4800 MT/s]), BIOS 1.5, microcode 0xa10113e, 2x Ethernet Controller 10G X550T, 1x 1.7T SAMSUNG MZQL21T9HCJR-00A07, Ubuntu 22.04.3 LTS, 5.15.125-0515125-generic, Framework=Pytorch 2.0, IPEX=2.0, Python 3.8, AI Model=DistilBERT Large(https://github.com/IntelAl/models/), INT8, Real Time (BS=1) results while maintaining 5ms latency SLA. Test by INTEL as of 09/11/23.
- For a 50 server fleet of AMD EPYC 9554, estimated as of October 2023:
- CapEx costs: \$1.36
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$749.7K
- Energy use in kWh (4 year, per server): 46573, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394
- For a 15 server fleet of 5th Gen Intel® Xeon® 8592+ as of October 2023
- CapEx costs: \$572K
- OpEx costs (4 year, includes power and cooling utility costs, infrastructure and hardware maintenance costs): \$238.3K
- Energy use in kWh (4 year, per server): 55475, PUE 1.6
- Other assumptions: utility cost \$0.1/kWh, kWh to kg CO2 factor 0.42394
- Costs based on Intel estimates and information from thinkmate.com as of October 2023.

Confidential computing and supporting technologies

<u>Confidential computing</u> with 5th Gen Intel[®] Xeon[®] Scalable processors isolates workloads inside trusted execution environments (TEEs). It helps strengthen data privacy, regulatory compliance, data access control and sovereignty, and application and data security.

- Intel® Software Guard Extensions (Intel® SGX) offers hardware-based security that helps protect data in use via unique application-isolation technology. Intel SGX is the most researched and most trusted execution environment for the data center, and it provides the smallest attack surface within the system.
- Intel[®] Trust Domain Extensions (Intel[®] TDX) provides virtual machine (VM) isolation from cloud providers and other tenants, and it more readily supports existing applications.¹
- Intel® Trust Authority takes Confidential Computing to the next level with a Zero Trust attestation SaaS that verifies the trustworthiness of compute assets at the network, edge, and in the cloud. Intel Trust Authority attests to the validity of Intel Confidential Computing environments.

Why it matters

Cloud: Confidently take advantage of the cloud while remaining compliant and in control of your data. Even confidential or regulated data can be protected when in use in the public cloud.

Collaboration: Confidential computing enables you to engage in multi-party analysis in a way that keeps each party's data private. Realize the benefits of shared analysis without losing privacy.

Regulatory compliance and data sovereignty: Confidential computing adds technological controls that help ensure data will be handled in compliance with proper procedures and regulatory frameworks.

Workloads: Confidential computing "armors-up" your workloads, helping protect sensitive data, content, and software intellectual property (IP) from advanced attack, tampering, and theft.

Learn more

- intel.com/confidentialcomputing
- intel.com/security-engines
- Accelerated AI inference with confidential computing
- <u>Confidential computing is combating modern slavery</u>
- Intel® Trust Authority | Intel Software (video)
- Securing Your Trust Boundary with Intel SGX and Intel TDX (video)

Business benefits

- Migrate workloads to the cloud with confidence
- Activate sensitive data in new services while maintaining privacy and regulatory compliance
- Enable beneficial multi-party data collaborations with full confidentiality
- Harden application security and strengthen compliance or data-sovereignty programs
- Customers can see up to 11% higher VM performance on 5th Gen Intel[®] Xeon[®] processors with Intel TDX vs. 4th Gen Intel[®] Xeon[®] processors without Intel TDX on integer, floating point and BERT-large²

Value differentiators

- Isolation: Separation of the TEE from the underlying software, admins, and other cloud tenants
- Encryption and control: Workload owner holds key to decrypt data, retaining control and preventing access by cloud provider or other entities
- Verification: Cryptographic confirmation that TEE is genuine and correctly configured, and that software is exactly as expected

Use cases

Protect

Protect data in use

• Data is routinely encrypted at rest and in transit. Confidential computing helps protect it while in use in the CPU and memory.

Sectors (often heavily

- regulated)
- Healthcare
- Financial
- Ad tech
- Government
- Retail
- Industrial and edge

- Confidential artificial intelligence (AI)
- Collaborative analytics
- Privacy-preserving ad tech
- Privacy-preserving blockchains
- Data and software IP control
- ¹ Intel TDX will be available on select 4th Gen Intel® Xeon® Scalable instances through four leading cloud providers. Previews have begun with select providers. Check with your provider for availability. Intel TDX becomes generally available with 5th Gen Intel® Xeon® Scalable processors.

² See [S1] at <u>intel.com/processorclaims</u>: 5th Gen Intel® Xeon® processors. Results may vary.