

Intel® AI Engines
5th Gen Intel® Xeon® Scalable Processors

Boost Entire AI Pipeline Performance with 5th Gen Intel® Xeon® Scalable Processors and Intel® AI Engines

65%

of data center AI inferencing runs on Intel® Xeon® processors¹

Up to

14x higher

real-time object detection inference performance (SSD-ResNet34) on 5th Gen Intel Xeon processors with AMX BF16 vs. 3rd Gen Intel Xeon processors²

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 processors with AMX BF16 vs. 3rd gen Intel Xeon processors³

Up to

8.7x higher

batch Recommendation System inference performance (DLRM) and 6.2x higher performance/watt on 5th Gen Intel Xeon vs. 3rd Gen Intel Xeon processors⁴

AI spans a wide range of workloads and use cases, from data preprocessing and classical machine learning (ML) to deep-learning usages, such as natural-language processing and image recognition. Intel Xeon Scalable processors deliver powerful computing performance for the entire AI pipeline. They include built-in accelerators optimized for specific AI workloads across ML, data analysis and deep learning.

Built-in power for AI across the enterprise

AI is pervasive and stretches across diverse and critical workloads. Classic ML and deep learning are becoming basic building blocks for how business gets done, from core enterprise applications to automated voice attendants. Putting AI to work at scale depends on a lengthy development pipeline that flows from data preprocessing to training to deployment. Each step has its own development toolchains, frameworks and workloads — all of which create unique bottlenecks and place distinct demands on computing resources. Intel Xeon Scalable processors feature built-in accelerators that can be used to run the entire pipeline right out of the box and to increase AI performance across the board.

Intel® Accelerator Engines are purpose-built integrated accelerators that support the most demanding emerging workloads

The 5th Gen Intel Xeon Scalable processors excel in general-purpose computing and will continue to serve as the foundation supporting many of today's critical AI workloads. These processors feature Intel® Advanced Matrix Extensions (Intel® AMX), a built-in AI accelerator designed to speed up deep-learning inference and training on the CPU. In many cases, this can eliminate the additional cost and complexity of a discrete accelerator. The latest generation of Intel Xeon processors is positioned exceptionally well for large language models (LLM) with less than 20 billion (20B) parameters — typically meeting the SLAs of customers.⁵ Intel AMX also excels in transfer learning and fine-tuning, so you can train models in as little as four minutes — not hours or days — without the need for additional hardware. With 65% of data center inferencing running on Intel Xeon processors, customers will benefit from their existing architecture for general-purpose AI, rather than navigating the complexities of a move to a GPU infrastructure.

Future innovation is here with 5th Gen Intel Xeon Scalable processors and Intel® Accelerator Engines

Whether you're using Intel Xeon processors for your workloads on premises, in the cloud or at the edge, Intel Xeon processors with built-in Intel Accelerator Engines can help your business reach new heights. They provide a range of benefits, including stronger data protection and better infrastructure utilization.



**Customer success:
Real-world acceleration
on Intel Xeon Scalable
processors**

Tencent Cloud delivers real-time speech synthesis with Intel Xeon Scalable processors.

[Get the details >](#)

Gunpowder runs Google Cloud C3 instances with 4th Gen Intel Xeon CPUs to accelerate rendering performance.

[Read the story >](#)

Intel Accelerator Engines can also help increase virtual and physical CPU utilization and reduce per-core solution licensing. Above all, these built-in accelerators provide increased application performance and reduced costs and improved platform-level efficiency.

Accelerate deep learning with Intel Advanced Matrix Extensions

Intel AMX is Intel's latest advancement for deep-learning training and inference on 5th Gen Intel Xeon Scalable processors. Ideal for workloads like natural-language processing, recommendation systems and image recognition, Intel AMX helps customers achieve up to 7.2x higher real-time object classification inference performance and 5.3x higher performance/watt on 5th Gen Intel Xeon with AMX BF16 versus 3rd Gen Intel Xeon processors.⁶

Intel AMX also provides a workload boost for AI models and enables more customers to meet their SLAs on the platforms that they are already running. 5th Gen Intel Xeon Scalable processors can offer improved turbo frequencies for workloads with affinity to vector and matrix operations, including high-performance computing and AI, adding five levels of turbo ratios.

Intel AMX improves matrix multiply operations performance with higher throughput (Ops/Cycle) compared to Intel® Advanced Vector Extensions 512 (Intel® AVX-512) on CPU cores.⁷ This facilitates faster completion of deep-learning training workloads and enables more customers to meet their SLAs on the platforms that are already running their businesses.

Supporting natural-language processing and generative AI

5th Gen Intel Xeon Scalable processors with Intel AMX offer a big performance boost for natural-language processing — and without additional hardware. Intel libraries are optimized for and integrated with TensorFlow and PyTorch, giving developers the benefits of built-in AI acceleration out of the box. Developers can also more easily migrate code from different hardware environments — a process that can be both lengthy and costly.

By accelerating deep-learning inferencing and training, the 5th Gen Intel Xeon Scalable processor featuring Intel AMX helps you meet your SLAs while balancing the total cost of ownership (TCO). It does this with a deep-learning-based recommender system that factors in real-time user behavior signals and additional context features like time and location.

The 5th Gen processor also runs generative AI models that mimic human-centric content, supporting large language models and text-to-image generation. For more intensive generative AI tasks, the purpose-built Intel® Gaudi® AI accelerator, the Intel® Data Center GPU and other hardware components can be used to expand the CPU's capabilities.

Intel AVX-512 for faster ML

Intel Xeon processors can hash SSL encryption for websites, crunch massive databases and run simulations for pharmaceutical research, chip design or Formula 1 engines.

Improved over multiple generations, Intel AVX-512 allows Intel Xeon Scalable processors to pack more operations into each clock cycle and improve performance for parallel processing applications. The Intel AVX-512 instruction set architecture (ISA) includes extensions built to augment performance of diverse workloads across AI, HPC, networking and storage.

Turbo performance increases from four to five levels of turbo ratios in the new generation of processors. This improves turbo frequencies for certain HPC and AI workloads leveraging Intel AMX and Intel AVX-512.

Fewer steps mean faster processing

Math can be very smart — and very elegant. Intel AVX-512 on 5th Gen Intel Xeon Scalable processors uses a lot of smart, beautiful math to condense, combine and fuse common computing operations into fewer steps. Here's a primitive example: You could instruct a CPU to calculate $3 \times 3 \times 3 \times 3 \times 3$, which would take five clock cycles. Or you could create an instruction for 3^5 that the CPU can do in one cycle. Intel AVX-512 takes that logic and applies it to hundreds of workload-specific operations, including some of the toughest operations in AI.

Counting by eight is a lot faster than counting by one

The "512" in Intel AVX-512 refers to the second way that these instructions increase the number of bits at the CPU's disposal with every clock cycle. Forty years ago, a 16-bit PC was pretty impressive. Soon, 32-bit machines took over. Today, your smartphone runs at 64 bits. Bit count refers to the number of registers — the memory slots where the CPU holds data — that the CPU can address per clock cycle. As the name suggests, Intel AVX-512 expands the number of registers to 512 bits. When an application takes advantage of Intel AVX-512, it runs up to eight times faster than the CPU's base 64-bit speed simply by expanding the number of registers. It's like counting to 96 by 1, 2, 3... versus 8, 16, 24.

Engines that require less power to run more powerful AI

Because Intel Xeon Scalable processors featuring Intel AI Engines require fewer hardware resources, they offer a more powerful and energy-efficient solution for running AI workloads.

Intel Xeon Scalable processors with built-in accelerator engines can also help provide improved workload results, like lowering TCO and delivering a better return on investment (ROI) for today's demanding AI workloads.

Faster AI is practically automatic with Intel Xeon processors

AI acceleration on Intel Xeon Scalable processors is built into the CPU's instruction set architecture (ISA). This means it's ready and available for any piece of software that can take advantage of it. Intel software engineers are constantly optimizing open-source AI toolchains and passing those optimizations back to the community. For example, TensorFlow 2.9 ships with Intel® oneAPI Deep Neural Network Library (Intel® oneDNN) optimizations by default. Download the latest edition, and TensorFlow will automatically take advantage of Intel optimizations.

For other applications in the AI pipeline, data scientists and developers can download free open-source Intel distributions, libraries and development environments that take advantage of every built-in accelerator in our ISA for Intel Xeon Scalable processors. Why should data scientists and AI developers recode their tools and recompile them for Intel AVX-512 when it can be done for them?

Organizations today need to get more workload performance out of their infrastructure — and do so with more power efficiency and at lower costs. The purpose-built Intel AI Engines integrated into Intel Xeon Scalable processors will help you get the most out of the AI workloads that matter most to your business.

Learn more about what Intel Xeon Scalable processors with built-in Intel Accelerator Engines can accomplish for the AI workloads that matter most to your business.

Learn more

[AI and deep learning on Intel Xeon Scalable Processors](#) ›

[Intel AVX-512](#) ›

[Intel® AI Analytics Toolkit](#) ›

[Developing on Intel® Hardware and Software](#) ›

Start accelerating AI workloads now — in the cloud or on your own infrastructure — with Intel optimizations for AI and ML.

[Learn more](#) ›



1. Based on Intel market modeling of the worldwide installed base of data center servers running AI inference workloads as of December 2022.
2. See [A21] at [intel.com/processorclaims](https://www.intel.com/processorclaims): 5th Gen Intel Xeon Scalable processors. Results may vary.
3. See [A19] at [intel.com/processorclaims](https://www.intel.com/processorclaims): 5th Gen Intel Xeon Scalable processors. Results may vary.
4. See [A20] at [intel.com/processorclaims](https://www.intel.com/processorclaims): 5th Gen Intel Xeon Scalable processors. Results may vary.
5. Based on Intel internal modeling as of December 2023.
6. See [A22] at [intel.com/processorclaims](https://www.intel.com/processorclaims): 5th Gen Intel Xeon Scalable processors. Results may vary.
7. <https://edc.intel.com/content/www/us/en/products/performance/benchmarks/vision-2022/>, Session Benchmark #41 and #42. Results may vary.

Notices and 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.

For workloads and configurations, visit 5th Gen Xeon Scalable processors at www.intel.com/processorclaims. Results may vary.

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

© 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.

Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy.

Availability of accelerators varies depending on SKU. Visit the [Intel Product Specifications page](#) for additional product details.

Intel® Advanced Vector Extensions (Intel® AVX) provides higher throughput to certain processor operations. Due to varying processor power characteristics, utilizing AVX instructions may cause, a) some parts to operate at less than the rated frequency and, b) some parts with Intel® Turbo Boost Technology 2.0 to not achieve any or maximum turbo frequencies. Performance varies depending on hardware, software, and system configuration, and you can learn more at [intel.com/content/www/us/en/architecture-and-technology/turbo-boost/intel-turbo-boost-technology.html](https://www.intel.com/content/www/us/en/architecture-and-technology/turbo-boost/intel-turbo-boost-technology.html).

Intel is committed to respecting human rights and avoiding complicity in human rights abuses. See Intel's [Global Human Rights Principles](#). Intel® products and software are intended only to be used in applications that do not cause or contribute to a violation of an internationally recognized human right.

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