

Table of contents

Which Intel built-in accelerators are right for your business? AI — Intel® Deep Learning Boost HPC — Intel® Advanced Vector Extensions 512 Security — Intel® Software Guard Extensions	What is built-in acceleration and why should you use it?	3
right for your business? AI — Intel® Deep Learning Boost HPC — Intel® Advanced Vector Extensions 512 Security — Intel® Software Guard Extensions		4
HPC — Intel® Advanced Vector Extensions 512 Security — Intel® Software Guard Extensions		5
Security — Intel® Software Guard Extensions	AI — Intel® Deep Learning Boost	6
	HPC — Intel® Advanced Vector Extensions 512	8
The next generation of Intel built-in accelerators	Security — Intel® Software Guard Extensions	10
	The next generation of Intel built-in accelerators	12
Conclusion 1	Conclusion	12

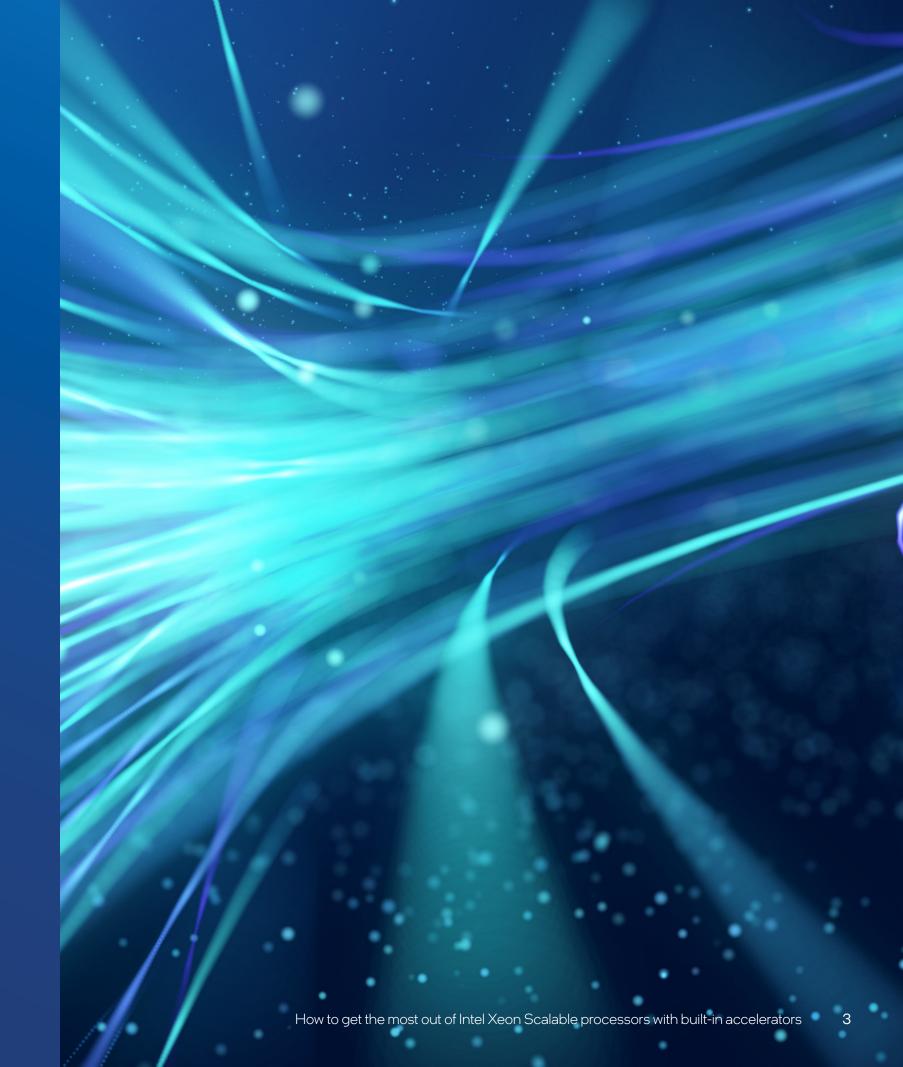


What is built-in acceleration and why should you use it?

What if, instead of buying new equipment every time you need to establish new capabilities, you could rely on the technology already built into your CPU? With Intel® Xeon® Scalable processors, you can. These CPUs contain features known as built-in accelerators to deliver enhanced benefits to your workloads.

Every Intel Xeon Scalable processor features the broadest and most unique set of integrated accelerators to help boost performance and efficiency, reducing the need for additional specialized hardware. In both cloud and on-prem environments, these purpose-built features support today's most common and demanding workloads spanning AI, security, HPC, analytics, storage and networking.

In this guide, we'll focus on AI, HPC and security.



The real-world benefits of Intel built-in accelerators

Whether you're using Intel Xeon Scalable processors for your workloads on prem, in the cloud or at the edge, our integrated accelerators can help your business reach new heights. They provide a range of benefits including faster security processing, stronger data protection and better infrastructure utilization. Above all, these built-in accelerators provide increased application performance, reduced costs and improved energy efficiency:

Performance



As a purpose-built component, Intel builtin accelerators can often deliver higher performance for a targeted workload.¹

Cost savings



Intel built-in accelerators allow you to improve performance without having to purchase additional specialized hardware.

Energy savings



Intel built-in accelerators help users improve energy efficiency by sparing users from having to include additional cores to the server rack.



Which Intel built-in accelerators are right for your business?

While Intel Xeon Scalable processors have a full array of accelerators built in, certain accelerators are better suited for specific tasks/workloads. To help you decide which Intel technologies can best support your business, let's take a deep dive into three of our key offerings in the AI, HPC and security segments.

Al

Intel® Deep Learning Boost (Intel® DL Boost) brings significantly accelerated performance to inferencing and training across common AI and HPC workloads.²

HPC

Intel® Advanced Vector Extensions 512 (Intel® AVX-512) is purpose-built to accelerate performance for the most demanding computational workloads in science, business and beyond.

Security

Intel® Software Guard Extensions (Intel® SGX) helps protect data in use via unique application-isolation technology.



AI — Intel Deep Learning Boost

What is Intel Deep Learning Boost?

Intel DL Boost is an accelerator designed to enhance performance and provide greater efficiency for AI and deep learning-related tasks and workloads.

Intel Xeon Scalable processors introduced purpose-built AI acceleration in 2019 with Intel Vector Neural Network Instructions (VNNI), now Intel DL Boost.

Based on the Intel AVX-512 accelerator, the VNNI component of Intel DL Boost can combine three instruction sets into one, drastically reducing the amount of time it takes to complete a task.

What are the most common use cases for Intel DL Boost?

Intel DL Boost accelerates various AI inferencing tasks such as image classification, language translation and object detection.

How are companies using Intel DL Boost in the real world?



Software company <u>rinf.tech</u> used Intel DL Boost to provide quicker and more accurate image analysis to support better real-time decision-making in retail, automotive, video surveillance and business intelligence use cases. Inference performance was up to 7.4 times faster than the baseline.³



Huiyi Huiying Medical Technology (HYHY) used Intel DL Boost to help optimize the performance of full-cycle AI medical imaging solutions. Thanks to the advantages brought by synergistic software-hardware acceleration, the company saw significant improvements in inference speed for image analysis scenarios like COVID-19 screening and breast cancer detection.⁴



Konfoong Biotech International Co., Ltd. (KFBIO), which specializes in digital pathology system development and production, used Intel DL Boost to complete scanning and diagnosis of M. tuberculosis specimens up to 11.4 times faster than the baseline.⁵



What performance advantages does Intel DL Boost offer?

Customers who use Intel® optimization for TensorFlow and Intel DL Boost will **gain over**

11 x more Al inference performance

on 3rd Gen Intel Xeon Scalable processors compared to 2nd Gen Intel Xeon Scalable processors.⁶

Intel DL Boost can also deliver more performance per watt for AI workloads, allowing organizations to reduce costs and energy consumption.

HPC — Intel Advanced Vector Extensions 512

What is Intel Advanced Vector Extensions 512?

Intel AVX-512 is a general-purpose performance-enhancing accelerator with a wide range of uses. With ultrawide 512-bit vector operations capabilities, Intel AVX-512 is especially suited to handle the most demanding computational tasks commonly encountered in the HPC segment.

What are the most common use cases for Intel AVX-512?

Intel AVX-512 is used by organizations across educational, municipal, financial, enterprise, engineering and medical industries for a wide range of complex tasks. These include AI, scientific simulations, 3D modeling and analysis, financial analytics, audio and video processing, cryptography, and data compression.



Intel AVX-512 enables real-time analytics for financial service workloads to improve customer experience, compliance and data security.



Intel AVX-512 enables you to run complex workloads on existing hardware, accelerating performance for tasks like 3D modeling and simulation.

How are companies using Intel AVX-512 in the real world?



The University at Buffalo's Center for Computational Research uses Intel AVX-512 to offer rich computing resources to businesses in Western New York like Marion Surgical, a company that uses virtual reality (VR) and augmented reality (AR) to teach surgeons complex procedures.⁷



The Broad Institute of MIT and Harvard uses Intel AVX-512 to improve processing speed and reduce costs for genomics workloads on Google Cloud N1 and N2 instances.⁸



Researchers studying fundamental particles at **CERN**, the European Organization for Nuclear Research, used Intel AVX-512 to accelerate simulation workloads via quantization. It offered **1.8 times higher performance** and slightly improved accuracy.⁹

What performance advantages does Intel AVX-512 offer?

Because HPC workloads involve massive amounts of data, CPU performance is crucial to ensuring accurate results in a timely manner.

Intel AVX-512 helps improve the vector processing power of Xeon CPUs compared to previous-gen solutions, allowing organizations to tackle intense workloads with greater speed. Applications can pack 32 double-precision and 64 single-precision floating point operations per clock cycle within the 512-bit vectors, as well as eight 64-bit and 16 32-bit integers with up to two 512-bit fused-multiply add (FMA) units, thus doubling the width of data registers, number of registers and width of FMA units compared to Intel® Advanced Vector Extensions 2 (Intel® AVX2).

Compared to AMD solutions, 3rd Gen Intel Xeon Scalable processors working with Intel AVX-512 offered superior performance across 12 HPC-related benchmarks and real-world applications.¹⁰



When running the biomolecular simulation algorithm NAMD, 3rd Gen Intel Xeon Scalable processors offered

1.27x higher performance

compared to **AMD Milan.**¹¹

When running the RELION algorithm used for 3D imaging in structural biology, 3rd Gen Intel Xeon Scalable processors offered

1.32 x higher performance

compared to **AMD Milan.**¹²

When running the financial investment simulation algorithm Monte Carlo FSI, 3rd Gen Intel Xeon Scalable processors offered

1.50x higher performance

compared to AMD Milan.¹³

When running the modeling and simulation algorithm LINPACK, 3rd Gen Intel Xeon Scalable processors offered

compared to **AMD Milan.**¹⁴

Security — Intel Software Guard Extensions

What is Intel Software Guard Extensions?

Intel SGX offers a hardware-based security solution that helps protect data in use via application-isolation technology. By protecting selected code and data from inspection or modification, developers can run sensitive data operations inside enclaves to help increase application security or protect data confidentiality. This adds another layer of defense by helping to reduce the attack surface of the system.

What are the most common use cases for Intel SGX?

Intel SGX enables confidential computing solutions to better protect data on prem, at the edge and in the cloud. This accelerator supports organizations looking to protect sensitive data and code to help ensure compliance with regulations related to data privacy, sovereignty and confidentiality.

How are companies using Intel SGX in the real world?



British financial institution <u>Nationwide Building Society</u> used Intel SGX to build "Know Your Customer," a system that allows multiple confidential data sets to be more safely processed inside of an enclave, in compliance with data regulations.¹⁵



Swiss Re Group, one of the world's largest reinsurance providers, is successfully exploring additional **protection of confidential data** from multiple parties with a confidential computing proof of concept using Intel SGX.¹⁶



The <u>University of California San Francisco</u> used Intel SGX to develop privacy-preserving analytics that accelerate the development and validation of clinical algorithms. The platform will provide a "zero trust" environment to help **protect both the intellectual property of an algorithm and the privacy of health care data.**¹⁷

What security advantages does Intel SGX offer?

Intel SGX provides a crucial building block for confidential computing by restricting access to sensitive data and code while in use, helping keep it safe from inspection or corruption by other software. Only Intel SGX has the flexibility to support virtualized, bare-metal and cloud-native container deployments.



The next generation of Intel built-in accelerators

Intel built-in accelerators are currently featured within Intel Xeon Scalable processors — and the next generation of built-in accelerators is on the horizon.

The 4th Generation Intel Xeon Scalable processor will include Intel® Advanced Matrix Extensions (Intel® AMX), Intel® QuickAssist Technology (Intel® QAT) and Intel® Data Streaming Accelerator (Intel® DSA), among other built-in features. Here's a quick glimpse at what you can expect.



Our next-generation DL Boost to advance deeplearning performance, Intel AMX features a set of matrix multiplication instructions that will significantly advance Al inference and training, with up to 4.5 times INT8 image inference per second compared to the prior generation.¹⁸



Now a part of Intel Xeon Scalable processors, it will offer users even faster data encryption and more efficient data compression for applications from networking to enterprise, cloud to storage and content delivery to database.



Intel DSA is a high-performance accelerator designed to optimize streaming data movement and transformation operations common in networking, data processing-intensive applications and high-performance storage.

Conclusion

Intel's long history of innovation and integration has uniquely positioned us to develop this new category of integrated accelerators with Intel Xeon Scalable processors. There are various workloads for which purpose-built integrated circuitry will deliver greater business value to customers.

Whether you're looking to increase performance, support sustainability initiatives or help ensure the protection of your most sensitive data, Intel's family of built-in accelerators offers a wide range of solutions without the need for additional hardware. As our own internal testing along with real-world use cases have shown, these accelerators deliver unmatched value compared to other CPU options on the market.

Learn more about Intel Xeon Scalable processors by visiting https://www.intel.com/xeonscalable.

```
See [123] at https://www.intel.com/3gen-xeon-config
```

stories/hypy-customer-story.html

⁵See https://www.intel.com/content/www/us/en/customer-spotlight/ stories/kfbio-ai-customer-story.html

⁶See [118] at https://www.intel.com/content/www/us/en/customer-spotlight/stories/university-at-buffalo-customer-story.html">https://www.intel.com/content/www/us/en/customer-spotlight/stories/university-at-buffalo-customer-story.html

⁸See https://www.intel.com/content/www/us/en/newsroom/news/ broad-institute-intel-google-advance-biomedical-research.html

⁹See https://www.intel.com/content/www/us/en/customer-spotlight/ stories/cern-inference-customer-story.html ¹⁰See [104] at https://www.intel.com/3gen-xeon-config

¹¹See [36] at https://www.intel.com/3gen-xeon-config

¹²See [38] at https://www.intel.com/3gen-xeon-config

¹³See [37] at https://www.intel.com/3gen-xeon-config

¹⁴See [39] at https://www.intel.com/3gen-xeon-config

¹⁵See https://www.intel.co.uk/content/www/uk/en/customer-

spotlight/stories/nationwide-building-society-customer-story.html

Gee https://www.intel.com/content/www/us/en/customer-spotlight/ stories/swiss-re-customer-story.html

17See https://www.intel.com/content/www/us/en/newsroom/news/ ucsf-propel-medical-device-innovations.html

18See Session Benchmark #41 and #42 at https://edc.intel.com/content/www/us/en/products/performance/benchmarks/vision-2022/. Results may vary.

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.

© 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.

²See [121] at https://www.intel.com/3gen-xeon-config

³See page 12 at https://www.intel.com/content/dam/www/public/us/en/documents/product-overviews/dl-boost-product-overview.pdf ⁴See https://www.intel.com/content/www/us/en/customer-spotlight/