# INSIGHTS. NOW ON DEMAND WITH INTEL® DEEP LEARNING BOOST.

**intel**

Enterprises looking to monetize AI need a powerful hardware infrastructure that delivers timely, precise insights. 2nd Generation Intel® Xeon® Scalable processors with new Intel® Deep Learning Boost (Intel® DL Boost) are enhanced specifically to run performance-hungry AI applications alongside existing cloud and data center workloads.

**Keywords:** *Inference, deep learning, image recognition, object detection, recommendation systems, speech recognition, deep neural network, convolutional neural networks*

## WHEN TO RECOMMEND

Talk about Intel Xeon Scalable processors with Intel DL Boost to customers who want exceptional AI performance with lower memory requirements, enabling their hardware footprint to do more.

## CUSTOMER PAIN POINTS

- Data-center bottlenecks are obstructing **real-time intelligence**
- When computational capacity is an issue, like in **convolutional neural networks (CNNs)** and **deep neural networks (DNNs)**
- Where l**ow-latency, low-power** hardware solutions are needed to drive **inference at scale**

## IDEAL FOR

- Image recognition
- Object detection
- Recommendation systems
- Speech recognition

## WHY UPGRADE

### ACCELERATED INSIGHTS

- Up to 2x faster inference with new int8 instructions vs. previous generations[1]
- Up to 30x improved deep learning performance vs. previous generations[1]

### OPTIMIZED FRAMEWORKS & LIBRARIES

Caffe

PaddlePaddle

mxnet

TensorFlow

INTEL® MKL-DNN

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

## SAY THIS TO YOUR CUSTOMER

"The Intel Xeon Scalable platform introduces a common platform for AI with high throughput for both inference and training, so you can **do both without purchasing a GPU**."[2]

"Intel is partnering with developers to continue optimizing popular frameworks and libraries to further accelerate inference performance."

"Intel DL Boost unlocks insights by optimizing systems for impactful automation. Imagine the efficiency you can deliver to your business by no longer having to purchase dedicated hardware to uncover the data you need."

"Intel DL Boost works by extending the Intel AVX-512 instruction set to do with **one instruction** what took **three instructions** in previous-generation processors. How would your organization benefit from those dramatic increases in efficiency?"

# HARDWARE AND STORAGE INNOVATION

*intel*

## ACCELERATE INFERENCE PERFORMANCE

**2ND GENERATION INTEL XEON PLATINUM PROCESSOR 9200 SERIES**
*with Intel DL Boost*

Up to **30X** better inference performance on image classification
*compared to competing processors[1]*

**2ND GENERATION INTEL XEON PLATINUM PROCESSOR 8200 SERIES**
*with Intel DL Boost*

Up to **14X** better inference throughput
*compared to previous-generation technology[3]*

intel **XEON® PLATINUM** inside™

## ADDED VALUE FOR DEEP LEARNING WITH INTEL® OPTANE™ TECHNOLOGY

Together with the Intel Xeon Scalable processor with Intel DL Boost, Intel® Optane™ technology can enhance both training and inference in deep learning.

### INTEL OPTANE DC PERSISTENT MEMORY
*Lower latency and more memory closer to the CPU enable larger in-memory working datasets and persistence across power cycles.*

### INTEL OPTANE SOLID STATE DRIVES
*Bigger, more affordable datasets and application acceleration help enterprises take advantage of next-level insights.*

## BENEFITS

### FOR TRAINING
Larger datasets and optimized batch training mean AI solutions can get smarter, faster.

### FOR INFERENCE
Larger datasets enable real-time and batch expansion of inference workloads.

Help businesses deliver AI readiness across the data center with Intel Xeon Scalable processors featuring Intel Deep Learning Boost. Contact your Intel Authorized Distributor or visit **ai.intel.com**