

How to make the right hardware choice for DEEP LEARNING TRAINING AND INFERENCE APPLICATIONS

There are many factors to bear in mind when addressing deep learning training and inference. Here's a quick guide to help you identify the optimal tool to meet your specific needs.

BEFORE YOU GET STARTED, CONSIDER THESE QUESTIONS:

- 1** How often do you need to train your model, both initially and after it has been deployed?
- 2** What kind of data are you using as your training set?
- 3** Can you train during off-peak hours?
- 4** What are your cost constraints?
- 5** What type of hardware architecture do you already have in your workstations or server rack?
- 6** Can you deal with the software complexities of multiple architectures?

TRAINING

Though only 10 to 15% of a typical workflow, training a model is a key step in harnessing the power of artificial intelligence. In this phase, the algorithm seeks to learn features and patterns from the data you feed it, later applying that knowledge to unseen data. Think of showing millions of examples of bone images to an algorithm that will be trained to identify bone density for radiologists. Training can take several hours to several days.

CPU

For many applications—such as high-definition-, 3D-, and non-image-based deep learning on language, text, and time-series data—CPUs shine. This is especially true for memory-intensive data, including massive amounts of unstructured data as well as sparse data. Infrequent training (fewer than 10 times per year) may also be a factor in staying on CPUs.

Pros

- AI applications can be run side by side with other applications, maximizing hardware utilization
- Larger, memory-intensive models can be trained with greater ease
- Data can be accessed from the same infrastructure on which you train, saving the time it would take to port data from one architecture to another
- Utilization can be maximized, which can contribute to improved total cost of ownership
- Public cloud instances of CPUs are typically far less expensive than GPUs, especially when training for extended periods on large models
- If your training workload is not time sensitive, CPUs may be a viable alternative to GPUs

Cons

- If you are training models frequently, the lower speed can potentially cost you valuable time
- The same training will typically require more cycles with a CPU than when using a GPU

GPU

For deep learning training with several neural network layers or on massive sets of certain data, like 2D images, GPUs or other accelerators can be your best bet. GPUs also tend to be the better choice for fast deep learning as the simple matrix math calculations greatly benefit when computations can be done in parallel.

Pros

- Fast training on certain types of data
- Memory limitations require you to break down images or your model, resulting in more work and/or subpar results
- Data must be ported from one architecture to another when transferred from GPUs to CPUs in the data center
- Can add cost, complexity, and operational expenses

Cons

- Hardware may not be utilized to full capacity
- Memory limitations require you to break down images or your model, resulting in more work and/or subpar results
- Data must be ported from one architecture to another when transferred from GPUs to CPUs in the data center
- Can add cost, complexity, and operational expenses

INFERENCE

After your model is trained, you'll put it to work with inference, or the inferring of something about data it has never seen before. Inference cycles already surpass those of training by anywhere from 10x to 1000x, depending on the application. Today, inference primarily runs on CPUs, and a typical deployment will use a mix of neural networks and other types of compute power. Facebook uses CPUs for 100% of its current inference applications.

As deep learning has been adopted more broadly, there has been a clear shift in the ratio between cycles of inference and training. Our conservative internal estimate predicts a move from 1:1 in the early days of deep learning to potentially more than 10:1 by 2020. With inference taking a large majority of the workflow, it is critical to use hardware architectures well suited to those needs, meaning low latency and often low power.

CPU

Recent hardware changes have enhanced the naturally good inference performance of CPUs. Coupled with all-new software optimizations and tools, this means the CPU has never been more performant for AI inference. There is also unprecedented diversity in acceleration hardware for robust and sustained inference applications. As inference is not as resource heavy as training, CPUs are economical as well.

Pros

- Real-world applications have increasingly stricter latency, a CPU strength
- The newest CPUs can support much more of the system memory required for complex models
- CPUs can extend further to the edge, including devices, unlike GPUs, allowing a similar architecture from training through deployment
- CPUs are well suited to new options like using pretrained models or transfer learning
- CPUs already power laptops, workstations, robotics, some phones and smart speakers, and vehicles, etc.

Cons

Given today's inference demands and the capabilities of current technology there are no significant downsides to using a CPU for inference.

GPU

GPUs excel at performing matrix operations that relied on by graphics, AI, and many scientific algorithms. Due to their parallel computing capabilities, GPUs can be useful for inference. But as AI models continue to grow in complexity, inference for real-world deployments increasingly favors CPUs.

Pros

- When inference speed is a bottleneck, GPUs can provide financial and time gains
- By design, GPUs can work well for tasks like image recognition
- GPUs may be a cheaper option today

Cons

- GPUs have inherent memory constraints
- GPUs can be power hungry, which conflicts with the needs of edge devices
- GPUs are less common in infrastructure that may already need to be utilized and upgraded to support AI applications

Though we focus here on CPUs and GPUs, there are certain environments that are not especially well suited to CPUs and GPUs, such as very small devices. Even for data center or cloud-based inference, sometimes the workloads are simply too sustained or of too high a volume for CPUs or GPUs to make economical sense. In these cases, other tools may be preferred, including vision processing units (VPUs), application-specific integrated circuits (ASICs), and field-programmable gate arrays (FPGAs).

VPUs – VPUs are ultralow-power computer vision engines and offer capabilities for inference for low-power devices such as smart cameras and network video recorders.

ASICs – Technically, a GPU is an ASIC used for processing graphics algorithms. A custom ASIC is used to performing fixed operations extremely fast since the entire chip's logic area can be focused on a narrow set of functions, making them suitable for a high degree of parallelism.

FPGAs – With particular value for high-throughput, low-latency inference applications, FPGAs offer a unique blend of flexibility, performance, and extensibility unmatched by custom ASICs or GPUs.

Intel® Xeon® Scalable processors are optimized specifically to run high-performance deep learning inference. Most deep learning, including computer vision and inference, already runs on Intel® Xeon® processors as they are the foundation of many of the world's data centers. The hardware performance of AI applications can further benefit from software optimizations.

Our new 2nd Generation Intel Xeon Scalable processors have AI built in, offering significant leaps in inference performance, memory, and bandwidth that accelerate complex AI applications. The processors have been enhanced with substantial improvements in software optimizations and hardware instructions, giving you the flexibility you need for both AI and the vast range of data-centric applications.

Intel® Deep Learning Boost technology – This cross-platform tool features a model optimizer and inference engine to streamline and simplify model deployment. Access a new set of embedded accelerators (vector neural network instructions, or VNNIs) to speed up the dense computations of convolutional neural networks (CNNs) and deep neural networks (DNNs). The low-precision integer operations deliver up to 30x improvement in inference performance.²

Intel® Optane™ DC persistent memory – Achieve up to triple the maximum storage per node and enable more memory closer to the CPU so data can be sustained even throughout power cycles or system maintenance.

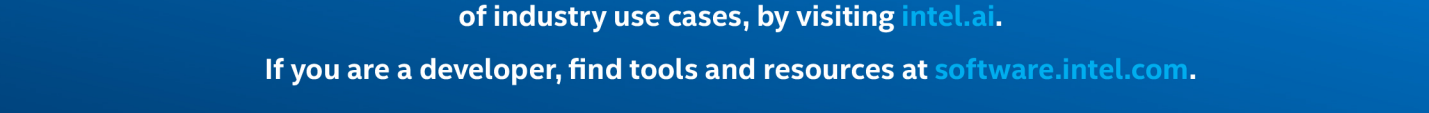
This is AI on Intel.

UP TO 277x¹ INFERENCE THROUGHPUT
achieved using an Intel® Xeon® Platinum 8180 processor running Intel-optimized Caffe®/GBoLeNet™ v1 with Intel® Math Kernel Library for Deep Neural Networks (Intel® MKL-DNN), vs. an Intel® Xeon® processor E5-2699 v3 with BVLC-Caffe®

UP TO 30x² INFERENCE THROUGHPUT IMPROVEMENT
on Intel® Xeon® Platinum 9282 processor with Intel® MKL-DNN, vs. an Intel Xeon processor E5-2699 v3 with BVLC-Caffe

UP TO 241x¹ TRAINING THROUGHPUT
achieved using an Intel-optimized Caffe AlexNet™ with Intel MKL-DNN, vs. an Intel Xeon processor E5-2699 v3 with BVLC-Caffe

To boost performance, we've also optimized the software tools and frameworks widely used today for Intel Xeon Scalable processor-based platforms:



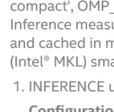
The best results are achieved when the right tool is used for the job. Today's CPUs and GPUs each boast distinct benefits, but neither is perfect for every environment or goal. To unlock the greatest possible impact for your deep learning application, ensure that you are using the optimal solution throughout training and inference.

Learn more about deep learning, and how Intel is powering AI across an exciting set of industry use cases, by visiting intel.ai.

If you are a developer, find tools and resources at software.intel.com.

For more information on Intel Enterprise Solutions for AI, please contact us today.

Enterprise Technology International



Software and workloads used in performance tests may have been optimized for performance only on Intel® microprocessors. Performance tests, such as SPECint_rate_base2000 and MobileMark, are measured using specific computer systems, components, software, and configurations. Intel, the Intel logo, Intel Inside, the Intel Inside logo, Intel Optane, and Xeon are trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries. *Other names and brands may be claimed as the property of others. © Intel Corporation 02/19/2018/ETI/PDF