

IT@Intel: Training Machine-Learning Models on Intel® Gaudi® Accelerators

In our quest for more efficient and cost-effective AI in Intel’s factories, the Manufacturing Automation team is migrating GPU-based deep-learning workloads to affordable, readily available Intel Gaudi accelerators

Authors

Ali Ahmadi

Director of AI Solutions and Equipment Productivity, Manufacturing IT

Vishal Kumar Pandey

Machine Learning Engineer, Manufacturing IT

Table of Contents

Executive Summary	1
Background.....	2
Solution.....	3
Benchmarking Details	3
Migration Process.....	3
Intel® Xeon® Scalable Processors for DL Inference.....	3
Results.....	3
Solution Architecture.....	4
Conclusion.....	4
Related Content.....	4

Executive Summary

Intel’s smart manufacturing environment relies on AI to improve yield, perform root-cause analysis, and accelerate anomaly detection. As Intel transforms to an internal foundry model, it is crucial that the Manufacturing Automation team helps improve factory efficiency by enabling faster, more cost-efficient deep-learning (DL) and machine-learning (ML) workloads.¹

We recently benchmarked a manufacturing object detection DL model training workload. We compared the performance of our legacy GPU-based system, a vendor upgrade to the GPUs, and Intel® Gaudi® accelerators. Our AI-based solution using Intel Gaudi accelerators produced the following benefits:

- A 20% improvement in DL model training time
- Reduced manual processes for employees and faster defect detection
- Streamlined workflows and enhanced manufacturing efficiency
- Improved cost efficiency and better supply chain availability for reduced operational expenses

Intel’s AI hardware and software ecosystem is a competitive and efficient solution for high-volume manufacturing (HVM) use cases. The manufacturing efficiency offered by Intel Gaudi accelerators will help us improve production yields and increase Intel factory revenue.

¹ Based on internal Intel IT testing and observations.

Contributors

Prasenjit Bose, Software Research Engineer/
Scientist, Manufacturing IT

Gobind Bisht, Software Application Engineering
Manager, Manufacturing IT

Robert Vaughn, Industry Engagement Manager,
IT@Intel

Acronyms

DL	deep learning
EOL	end of line
GPU	graphics processing unit
HVM	high-volume manufacturing
ML	machine learning
WISTA	Wafer Image Scan and Tool Analytics

Background

Intel's semiconductor high-volume manufacturing (HVM) factories process a few hundred thousand wafers every day using lithography exposure. Due to the sheer volume of wafers being processed, only a small subset of them undergoes defect metrology (Defmet). However, at every lithography exposure step, a scanner measures the dense topography of the wafer to aid in precise exposure control for accurate feature printing.

Due to the limited inline Defmet data, manufacturing engineers can easily miss infrequent but catastrophic excursions. This often results in large excursions escaping detection for extended periods of time (several days to months), sometimes to only be detected at end of line (EOL), leading to huge lost revenue.

To improve defect detection, Intel's Manufacturing Automation team has been on a multiyear journey of continuously improving our use of machine learning (ML) and deep learning (DL). Use cases include yield improvement, root-cause analysis, and anomaly detection. Detecting anomalies using inline data (as opposed to EOL data) is a crucial step toward excursion detection because it enables manufacturing engineers to detect and contain problems more quickly. DL models are very successful in such detection and localization.

AI has transformed Intel's manufacturing environment to help increase yield, reduce costs, and improve profitability. As Intel transitions to an internal foundry model, it's more important than ever that the Manufacturing Automation team supports IDM 2.0² and maintains Intel's manufacturing health by using proactive inline monitoring of every wafer and die. However, it is impossible for fab workers to manually perform excursion monitoring on 100% of wafers.

We developed the Wafer Image Scan and Tool Analytics (WISTA) system to replace manual excursion monitoring. Our unique system performs extraction, parsing, fitting, and image representation of scanner topography sensor data. The result is a two-dimensional rectangular-grid heat map image (see Figure 1 for an example image). WISTA processes approximately a half-million of these images around the globe daily.

Next, the system uses classification and quantitative extraction of defect morphologies using a [fine-tuned model](#) with YOLO-based architecture³ and various computer vision (CV) techniques. This is followed by an estimation of how the defective dies will impact wafers based on defect pixels captured in the previous step. The impacted dies across various wafer layers are then compared to EOL yield analysis data to obtain the die yield probability.

WISTA can detect abnormal patterns early in the wafer processing workflow (the entire wafer manufacturing process takes a couple of months), improving yield and enabling engineers to fix the manufacturing line quickly. WISTA enables us to analyze metrology measurements for 100% of the processed wafers and provide real-time information on defects.

We estimate that WISTA reduces the cost of a manufacturing excursion from millions to a few thousand dollars.⁴ However, even as WISTA saves Intel Manufacturing millions of dollars, we knew we could further improve the efficiency and cost-effectiveness of our DL training workloads in manufacturing.

We retrain our DL and CV models every one to two weeks. Retraining is necessary if a new abnormal pattern is detected or if the model is drifting. Our reliance on traditional graphics processing units (GPUs) for model training workloads led to higher operational costs and potential supply chain constraints, impacting manufacturing efficiency. We needed a hardware solution that could improve supply chain resilience, cut costs, and demonstrate comparable or improved performance without disrupting existing processes.

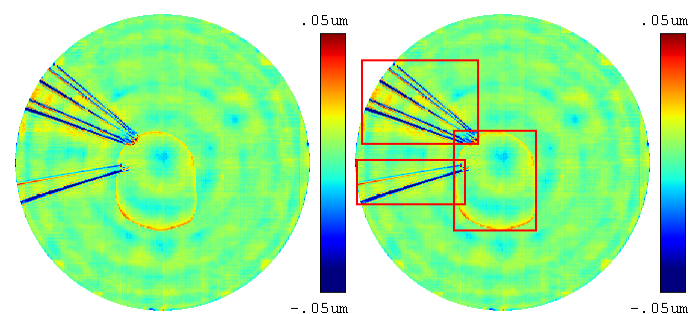


Figure 1. A raw image (left side) is fed to the model, which then overlays one or more bounding boxes around the identified defect (right side).

² "IDM 2.0" is a major evolution of Intel's integrated device manufacturing (IDM) model. For more information, read the [press release](#).

³ YOLO stands for "You Only Look Once." Find out more at <https://www.datacamp.com/blog/yolo-object-detection-explained>.

⁴ Costs based on internal assessments, assumes a manufacturing cost of \$10/die.

Solution

We evaluated three hardware options for model training:

- Continuing with the legacy GPU platform
- Exploring other externally sourced AI accelerators
- Migrating to [Intel® Gaudi® AI accelerators](#)

After carefully benchmarking these three options for eight weeks, we chose Intel Gaudi accelerators for their cost-effectiveness, performance improvements, and ease of integration. Migrating DL workloads to Intel Gaudi accelerators was straightforward and took only one week. Our success with Intel Gaudi accelerators validates that Intel's AI hardware and software ecosystem is a competitive and efficient solution for HVM use cases.

Benchmarking Details

We began benchmarking by evaluating the relative performance and price/performance of industry-standard generic workloads running on our legacy GPU hardware, a vendor upgrade to our legacy hardware, and Intel Gaudi 2 accelerators. These workloads were as follows:

- Object detection (YOLOv5)
- Sequence modeling (VGG16)
- Generative image synthesis (PIX2PIXHD)

The results of the generic testing were promising, so we moved on to testing with a real object detection workload from our manufacturing environment. Our testing revealed that the cost delta between a single legacy GPU card and a single Intel Gaudi accelerator is substantial, which means we can obtain more hardware for the same cost using Intel Gaudi accelerators. In addition to the cost benefits of Intel Gaudi accelerators, we also improved DL model performance by working with Intel Gaudi accelerator engineers to quickly customize our DL model training to take advantage of Intel Gaudi accelerators. This optimized model training code is publicly available on [GitHub](#).

Migration Process

When testing was complete and we were confident that migrating our DL workload to Intel Gaudi accelerators could deliver business value—both from a performance and cost perspective—we were ready to proceed with the full migration. To deploy the pre-trained model that was tuned on a custom dataset in our production environment, we removed dependencies on the legacy hardware and switched to the [Intel Gaudi Software Suite and libraries](#). We completed the migration in less than 10 hours.

Now that the infrastructure setup is complete, we can expand our efforts by tuning this model and others for other manufacturing use cases. In addition, other groups outside of Intel Manufacturing can also take advantage of our learnings and Intel Gaudi platform to cut AI costs and improve model performance.

Intel® Xeon® Scalable Processors for DL Inference

While we have found using Intel® Gaudi® AI accelerators useful for streamlining deep-learning (DL) model training, we run our manufacturing DL model inference workloads on Intel Xeon Scalable processors. These processors, already installed in our manufacturing data centers, enable us to take advantage of existing investments by running new workloads on existing hardware.

The latest generation of Intel Xeon processors has built-in AI acceleration through Intel® Advanced Matrix Extensions (Intel® AMX). This integrated feature improves the performance of DL inference on the CPU and is ideal for workloads like natural-language processing, recommendation systems, and image recognition. In addition to accelerating AI inference, Intel Xeon Scalable processors offer higher performance per watt than other systems.⁵

Near-zero effort is required to improve performance with Intel AMX. This is because default frameworks are optimized with Intel® oneAPI Deep Neural Network Library (oneDNN). Windows and Linux operating systems, kernel-based virtual machines (KVMs), and popular hypervisors expose the Intel AMX instruction set. INT8 and BF16 operations are automatically optimized in open-source frameworks like TensorFlow and PyTorch. The Intel® Distribution of OpenVINO™ toolkit allows developers to automate, optimize, tune, and run AI inferencing with little or no coding knowledge.

⁵ For details, read the [Intel® AMX product brief](#).

Results

With the WISTA solution that we developed—coupled with our efficient migration to Intel Gaudi accelerators—we achieved several business benefits:

- We improved model training time by 20% and the system can process half a million images daily.
- Employees experienced reduced manual workload and faster defect detection.
- We streamlined workflows and enhanced manufacturing efficiency.
- The cost-effectiveness of Gaudi and better supply chain availability reduced operational expenses.

Overall, migrating our DL training workload has improved efficiency, which can lead to higher production yields and revenue.

Solution Architecture

Our ML/DL training solution architecture is shown in Figure 2. The architecture is based on open-source tools to curb costs and enable fast innovation. We use Docker containers, which contain a description of the container along with all required build elements. The model is executed using the PyTorch ML framework, and the Intel Gaudi PyTorch Bridge enables the execution of PyTorch models on Intel Gaudi accelerators. The bridge includes features such as the following:

- Memory management for memory allocation and release.
- Pinned memory, which allows tensors to be created with pinned memory to reduce direct memory access time.
- Operations placement to distribute operations between the CPU and the Intel Gaudi accelerator for optimal performance.

Once the model is trained, we convert it to an intermediate representation using the Intel® Distribution of OpenVINO™ toolkit. We then run batched inferencing jobs on Intel® Xeon® Scalable processors using a software stack that is similar to the training stack.

We are currently using Intel Gaudi 2 accelerators. Every accelerator is interconnected to every other accelerator within a node in an all-to-all configuration. Each accelerator integrates 24x100 GbE RDMA over Converged Ethernet (RoCE) ports, of which 21 are used for scale-up connectivity within an eight-card universal baseboard, and three ports are used for scale-out connectivity. To scale outside a node, each of six QSFP-DD ports aggregates four Gaudi 100 GbE connections. 400 Gb switches are used for the fabric that interconnects all the accelerators in a cluster.

We are in the process of acquiring Intel Gaudi 3 accelerators, which offer additional performance and scalability advantages. Intel Gaudi 3 accelerators differ from the previous generation by the following improvements:

- 200 GbE RDMA instead of 100 GbE RDMA
- Six OSFP ports that aggregate 200 GbE connections instead of QSFP-DD ports and 100 GbE connections
- 800 Gb switches instead of 400 Gb switches

Solution Stack for DL/ML Model Training

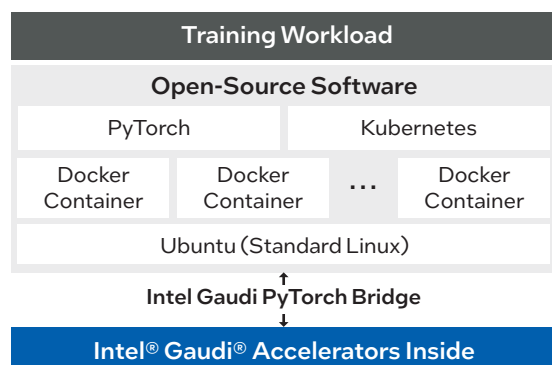


Figure 2. An open-source-based software stack helps keep costs down.

Conclusion

As Intel’s smart manufacturing increasingly relies on DL and ML, Intel Gaudi AI accelerators can provide the high performance that is necessary for rapid training at the scale and cost efficiency we need. Our production-grade Intel Gaudi accelerator deployment has proven to be a viable alternative to the competition in DL compute capability, pricing, and market availability. We can also take advantage of these accelerators’ excellent scalability.

Using open-source software and industry-standard Ethernet connections, these accelerators can help drive down the cost of AI solutions. We will continue to invest in Intel Gaudi accelerators and expand our use of them for model training across Intel’s manufacturing environment.

Related Content

If you liked this paper, you may also be interested in these related stories:

- Faster, More Accurate Defect Classification Using Machine Learning
- Transforming Manufacturing Yield Analysis with AI
- Democratizing the Use and Development of Generative AI Across Intel
- Push-Button Productization of AI Models
- Autonomous Quality in AI Model Productization: A Journey
- Transforming Intel with AI
- Smart Manufacturing Using Computer Vision and AI for Inline Inspection
- Revolutionizing Product Validation Using AI
- Machine Learning: The Next Step in Advanced Analytics

For more information on Intel IT best practices, visit intel.com/IT.

IT@Intel

We connect IT professionals with their IT peers inside Intel. Our IT department solves some of today’s most demanding and complex technology issues, and we want to share these lessons directly with our fellow IT professionals in an open peer-to-peer forum.

Our goal is simple: improve efficiency throughout the organization and enhance the business value of IT investments.

Follow us and join the conversation on [X](#) or [LinkedIn](#). Visit us today at intel.com/IT if you would like to learn more.



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

No product or component can be absolutely secure.

Your costs and results may vary.

© 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.

0425/WWES/KC/PD