# Intel® Gaudi®2 AI Accelerator for Deep Learning Training and Inference

Karthikeyan Vaidyanathan

November 2023

intel.

# Notices and Disclaimers

Performance varies by use, configuration and other factors. Learn more at https://habana.ai/habana-claims-validation/

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.

# Agenda

- Gaudi2 programming model and recent MLPERF results
- Experience at scale
  - Use case1: SWIFT congestion control
  - Use case2: Packet/message spraying

# Intel Xeon and Gaudi2 Processors for Models E2E

**Train and deploy large scale GenAI and LLMs**

**Fine tune and run thousands of domain specialized models with targeted curated data sets from the data center and the factory floor to devices**
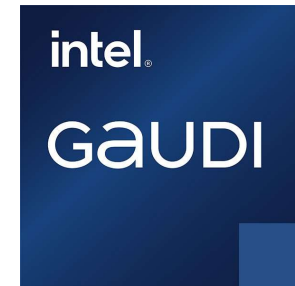
Gaudi2 Clusters and systems for models from billions to trillions of parameters

Intel Xeon for fine tuning and inferencing models up to tens of billions of parameters

# Intel® Gaudi® Accelerator Roadmap

**Available Now**

## GAUDI®

Native RoCE
Scaleup & out

**Available via:**
- HLS-1 Server (x8)
- SMC Server (x8)
- SDSC
- Public Cloud AWS: EC2

**Available Now**

## GAUDI®2
(7nm)

Native RoCE
Scaleup & out

**Available via:**
- HLS-Gaudi2 Server (x8)
- SMC Server (x8)
- Aivres/IEI Server (x8)
- Intel Dev Cloud

**In Development**

## GAUDI®3
(5nm)

Native RoCE
Scaleup & out

**2024**

**In Development**

## Next Generation AI Accelerator: Falcon Shores 1

Native RoCE
Scaleup & out

**2025**

# Intel delivers increasingly competitive Training Performance

- One of only three accelerators submitting GPT-3 results: Intel, Nvidia, Google

- Xeon continues to be the only CPU to submit training results on the MLPerf Benchmark.

# Intel® Gaudi®2 Accelerator Performance Doubled with FP8

- Intel Gaudi team projected to customers +90% performance gain with FP8
- Delivered more than promised: 103% on GPT-3 industry benchmark

**MLPerf Training 3.1GPT-3 Benchmark**
*- Lower is better -*

Time-to-Train in Minutes
384 Accelerators

311.94

153.58

Jun '23

Nov '23

See backup for workloads and configurations. Results may vary.

intel.

# Intel® Gaudi®2 performance advances strengthen competitive price-performance vs. H100

- Gaudi2 performance on ResNet near that of H100.

- H100 with FP8 outperformed Gaudi2 with BF16 on BERT.

- Vs.TPU, Gaudi2 delivered 3x performance on GPT-3.

- Given its significantly lower server cost vs. H100 server cost, Intel Gaudi2 delivers price-performance advantage vs. H100 across models.

**Time-to-Train Relative Performance**

**Lower is better**

Legend:
- Nvidia H100 (Nov)
- Intel Gaudi2 (Nov)
- Intel Gaudi2 (May)
- Google TPU

Bar values:
- GPT-3: 0.5x, 1, 2.03x, 3.09x
- Stable Diffusion: 0.5x, 1
- ResNet: 0.85x, 1
- BERT: 0.41x, 1

GPT-3
512x GPUs
384x G2s
4096 TPUs
FP8

Stable Diffusion
64x GPUs/ G2s
H100/FP16
G2/BF16

ResNet
8x GPUs/ G2s
H100/FP16
G2/BF16

BERT
8x GPUs/ G2s
H100/FP8
G2/BF16

For complete results information and configurations, see MLCommons publication: https://mlcommons.org/en/inference-datacenter-31/

intel.

# Outstanding Intel®Gaudi®2 AI Accelerator performance on MLPerf v3.1 Inference Benchmark

## Intel Gaudi2 Accelerator with FP8: near-parity performance on GPT-J (Server) with H100

- Gaudi 2 inference performance on GPT-J: -9% (Server) and -28% (Offline) vs H100

- Gaudi 2 outperformed A100 by 2.4x (Server) and 2x (Offline)

- Gaudi 2 employed FP8 and reached 99.9% accuracy

**GPTJ-99 Server Throughput Samples/sec, 8 accelerators (Higher is better)**

| | Gaudi2 | H100 | GH200-96G | A100-SXM-8 |
|---|---|---|---|---|
| | 1X | 1.09X | 1.12X | 0.42X |

**GPTJ-99 Offline Throughput Samples/sec, 8 accelerators (Higher is better)**

| | Gaudi2 | H100 | GH200-96G | A100-SXM-8 |
|---|---|---|---|---|
| | 1X | 1.28X | 1.27X | 0.50X |

For complete results information and configurations, see MLCommons publication: https://mlcommons.org/en/inference-datacenter-31/

# SynapseAI Software:
## Optimized for Intel® Gaudi® Performance and Ease of Use



- Shared software suite for <u>training and inference</u>
- Start running on Intel Gaudi accelerators with minimal code changes
- Integrated with PyTorch and TensorFlow
- Rich library of performance-optimized kernels
- Advanced users can write their custom kernels
- <u>Docker container images</u> and Kubernetes orchestration
- <u>Habana Developer Site</u> & <u>HabanaAI GitHub</u>
- <u>Habana Developer Forum</u>

# Distributed Pytorch

# HCCL API

// Communicator creation

hcclGetUniqueId(hcclUniqueId* uniqueId);

hcclCommInitRank(hcclComm_t* comm, int nranks, hcclUniqueId commId, int rank);

// Communicator destruction

hcclCommDestroy(hcclComm_t comm);

// Collectives communication

hcclReduceScatter(void* sbuff, void* rbuff, size_t recvcount, hcclDataType_t datatype, hcclRedOp_t op, hcclComm_t comm, synStreamHandle stream_handle);

hcclAllReduce(void* sbuff, void* rbuff, size_t count, hcclDataType_t datatype, hcclRedOp_t op, hcclComm_t comm, synStreamHandle stream_handle);

hcclBroadcast(void* sbuff, void* rbuff, size_t count, hcclDataType_t datatype, int root, hcclComm_t comm, synStreamHandle stream_handle);

hcclAllGather(void* sbuff, void* rbuff, size_t sendcount, hcclDataType_t datatype, hcclComm_t comm, synStreamHandle stream_handle);

hcclReduce(void* sbuff, void* rbuff, size_t count, hcclDataType_t datatype, hcclRedOp_t op, int root, hcclComm_t comm, synStreamHandle stream_handle);

hcclAlltoAll(...);

// Point-to-point communication

hcclSend(void* sbuff, size_t count, hcclDataType_t datatype, int peer, hcclComm_t comm, synStreamHandle stream);

hcclRecv(void* rbuff, size_t count, hcclDataType_t datatype, int peer, hcclComm_t comm, synStreamHandle stream);

// Aggregation/Composition

hcclGroupStart();

hcclGroupEnd();

# Incast Congestion



- PFC is great but does not work for multi-tenant and multi-level switches
- When packet drops occur, utilization is poor

# SWIFT congestion control for Habana Gaudi2

| 7:1 congestion, No PFC | #packets dropped | Bandwidth utilization |
|---|---|---|
| Default | 276787 | 10-50% |
| SWIFT<br>for 4KB/8KB MTU (targetdelay=20usecs, ai=2, beta=0.5,min_cwnd=2, max_cwnd=32) | ~1 | ~98% |



7:1 incast congestion



SWIFT:
https://dl.acm.org/doi/pdf/10.1145/3387514.3406591

# Packet collision at large-scale



Collision due to mapping to same output port → leads to performance degradation

ECMP hashing
- src ip
- dst ip
- src port
- dst port
- protocol

Spine Switches

Leaf Switches

Gaudi2 Servers

G0    G4        G16    G20

Only certain output ports have traffic and rest are idle and unutilized

# Solution: Packet spraying



Spine Switches

ECMP hashing
- src ip
- dst ip
- src port
- dst port
- protocol

Leaf Switches

Gaudi2 Servers

G0

G16

# Packet spraying solution

| HCCL collectives BW | With collisions | Packet spraying | Expected |
|---|---|---|---|
| All2All | 22 GB/s | 64 GB/s | 65 GB/s |
| Allgather | 183 GB/s | 272 GB/s | 272 GB/s |

Almost all ports are utilized

# Developer Resources

**Gaudi Developer Site: developer.Habana.ai**

**Habana GitHub**

**Habana Optimum Library on Hugging Face Hub**

# Summary

- Intel Gaudi2 continues to be the only viable alternative to NVIDIA's H100 for GenAI/LLM compute, with a significant price-performance advantage.

- 4th Gen Intel Xeon processors help customers train small- to mid-sized deep learning models, as well as fine tuning and transfer learning.

- Intel is well positioned to address every phase of the AI continuum across AI workloads, from large to small models—giving customers choice.

# Thank you