# Accelerating Distributed Deep Learning using HCCL: Habana Collective Communication Library

Karthikeyan Vaidyanathan

Exacomm 2023 Workshop

# Notices and Disclaimers

This document contains information on products, services and/or processes in development.  All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest forecast, schedule, specifications and roadmaps.

Intel technologies may require enabled hardware, software or service activation. Your costs and results may vary.  Cost reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings.  Circumstances will vary.  Intel does not guarantee any costs or cost reduction.  No product or component can be absolutely secure.

Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. For more complete information about performance and benchmark results, visit http://www.intel.com/benchmarks .

Intel does not control or audit third-party benchmark data or the web sites referenced in this document. You should visit the referenced web site and confirm whether referenced data are accurate.

Some results have been estimated or simulated using internal Intel analysis or architecture simulation or modeling and are provided to you for informational purposes.

Intel Advanced Vector Extensions (Intel AVX) provides higher throughput to certain processor operations. Due to varying processor power characteristics, utilizing AVX instructions may cause a) some parts to operate at less than the rated frequency and b) some parts with Intel Turbo Boost Technology 2.0 to not achieve any or maximum turbo frequencies. Performance varies depending on hardware, software, and system configuration and you can learn more at http://www.intel.com/go/turbo.

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.
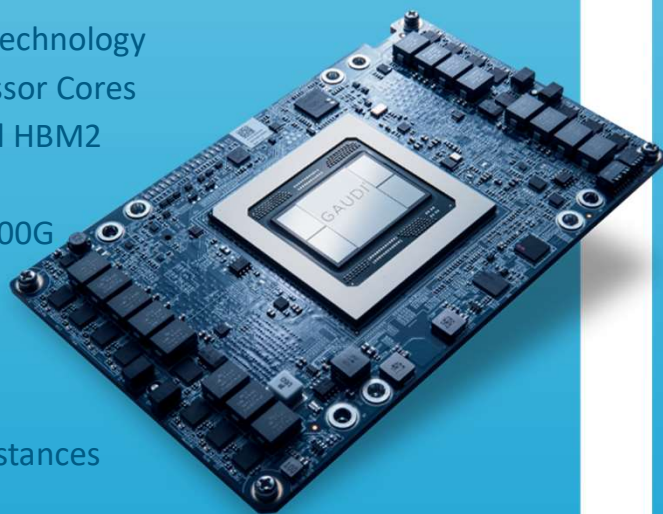
# Habana Deep Learning Solutions

## GAUDI®
### High-performance, high-efficiency (price/performance)

- 16nm process technology
- 8 Tensor Processor Cores
- 32 GB on-board HBM2
- 24 SRAM
- 10 integrated 100G Ethernet ports
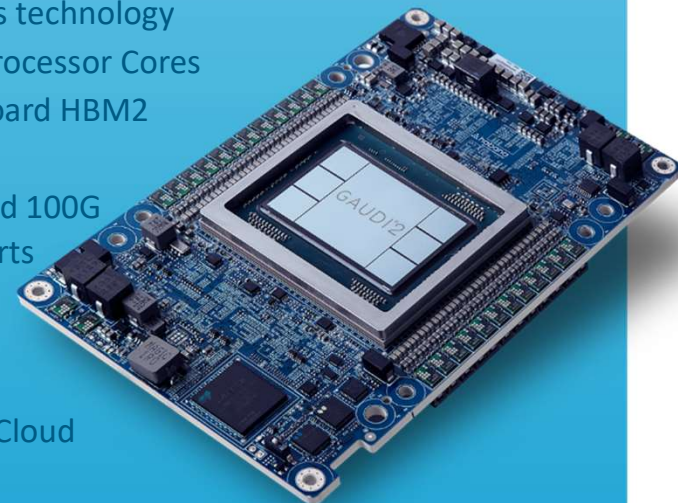
In the cloud:
- Amazon EC2 DL1 Instances

On-premises:
- Supermicro X12 Gaudi Server with 3rd Gen Xeon CPU

## GAUDI®2
### Higher performance, high-efficiency; optimized speed, memory, scalability for large scale models

- 7nm process technology
- 24 Tensor Processor Cores
- 96 GB on-board HBM2
- 48 SRAM
- 24 integrated 100G Ethernet ports

In the cloud:
- Intel Developer Cloud

On-premises:
- Supermicro Gaudi2 Server with 3rd Gen Xeon CPU

# Gaudi2 Server Delivers Flexible and Efficient Scalability

## GAUDI®2

**HLS-Gaudi2 Reference Server featuring...**

- 8 Gaudi2 mezzanine cards

- 24x 100 GbE ports per card
  - 21 for all-to-all connectivity to other 7 Gaudis within the server
  - 3 to scale out
  - Through 6 QSFP-DD ports
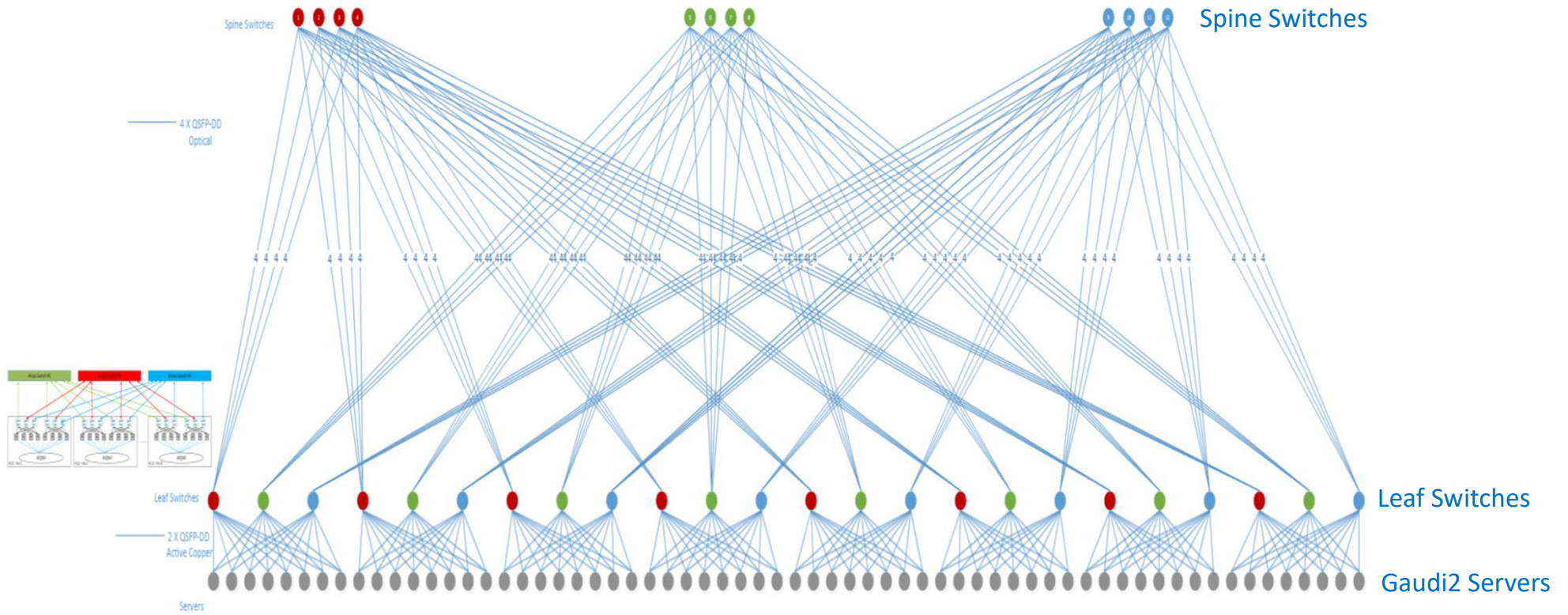
- Dual-socket Host CPU: 3rd Gen Xeon Scalable Processors



**Dual-socket Intel® Xeon® 3rd Gen Scalable Host CPU**
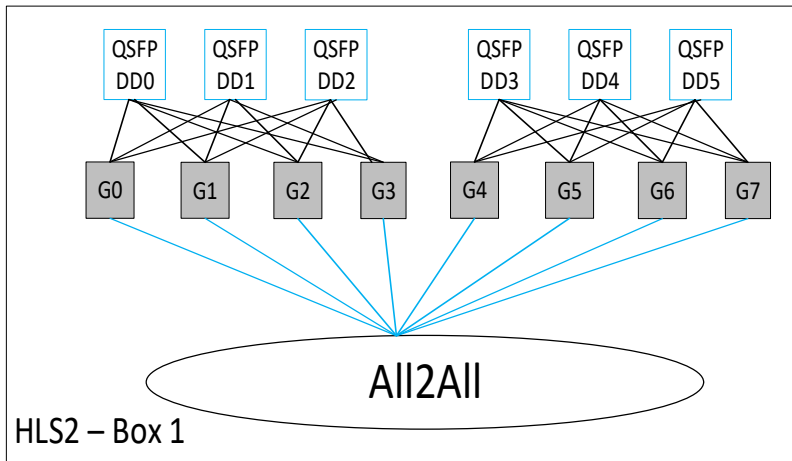
# 16 servers (128 Gaudi2s) datacenter

# Gaudi Network for 512 Gaudi2s (64 Servers)

# HCCL Hierarchical Collectives



**Step 1: Intra-HLS box communication**

QSFP DD0  QSFP DD1  QSFP DD2  QSFP DD3  QSFP DD4  QSFP DD5

G0  G1  G2  G3  G4  G5  G6  G7

All2All

HLS2 – Box 1

Scale-up algorithm

**Step 2: Inter-HLS box communication**

All steps exchange W/128 message

Gaudi #0  Gaudi #8  Gaudi #16  Gaudi #24  ......  Gaudi #96  Gaudi #104  Gaudi #112  Gaudi #120

Scale-out, pair-wise exchange algorithm

# High-level flow of hcclAllreduce in a box

| Original | | | | | | | | | Reduce-Scatter | | | | | | | | | All-gather | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| P0 | P1 | P2 | P3 | P4 | P5 | P6 | P7 | | P0 | P1 | P2 | P3 | P4 | P5 | P6 | P7 | | P0 | P1 | P2 | P3 | P4 | P5 | P6 | P7 |
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | | 8 | | | | | | | | | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 |
| 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | | | 16 | | | | | | | | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 |
| 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | | | | 24 | | | | | | | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 24 |
| 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | | | | | 32 | | | | | | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | | | | | | 40 | | | | | 40 | 40 | 40 | 40 | 40 | 40 | 40 | 40 |
| 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | | | | | | | 48 | | | | 48 | 48 | 48 | 48 | 48 | 48 | 48 | 48 |
| 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | | | | | | | | 56 | | | 56 | 56 | 56 | 56 | 56 | 56 | 56 | 56 |
| 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | | | | | | | | | 64 | | 64 | 64 | 64 | 64 | 64 | 64 | 64 | 64 |

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

- Address exchange the buffers

- Reduce-scatter + allgather + remainder, slice large allreduce into small ones; on each slice:

- Reduce-scatter
  - (Send) RDMA write 1/8th of the tensor to every Gaudi and Recv 1/8th from every Gaudi
  - SRAM reduction: read from HBM, reduce in SRAM and write result back to HBM
  - Ordering guarantees, reproducible reductions

- Allgather
  - (Send) RDMA write the tensor to every other Gaudi

# HCCL Collective Performance (Gaudi1 and Gaudi2) in a box



We achieve full utilization in message sizes that are much smaller than other RDMA solutions

# 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();

# Near-linear scaling with Gaudi2

– High scaling efficiency across multiple workloads drives performance and TCO advantages

# Stable Diffusion Training on Gaudi2

Linear scaling efficiency > 99% up to 64x cards



Model source: https://github.com/HabanaAI/Model-References/tree/master/PyTorch/generative_models/stable-diffusion-training, Dataset laion2B-en.
Training with BF16, batch size = 16, global batch size = 1024, for 1K iterations. Image size 256x256. Measurements using SynapseAI 1.9.0

# Performance of Reference Models



https://developer.habana.ai/resources/habana-models-performance/

# Q&A

# Legal Notices and Disclaimers

Performance varies by use, configuration and other factors. See details on performance and other related claims at the QR code below or 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. No product or component can be absolutely secure.

Your costs and results may vary.

© Habana, Gaudi and SynapseAI marks and logos are trademarks of Habana Labs. Other names and brands may be claimed as the property of others.

# Legal Notices and Disclaimers

– For notices, disclaimers, and details about performance claims, visit https://habana.ai/habana-claims-validation

or scan the QR code:



© 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.

# Thank You.