August 1, 2022

# Training Deep Learning Models on Habana Gaudi®

Milind S. Pandit, Senior Solutions Architect, mpandit@habana.ai

https://www.habana.ai

**habana®**
An Intel Company

# A little about Habana

- Founded in 2016 to develop purpose-built AI processors

- Launched inference processor in 2018, training processor in 2019

- Acquired by Intel in late-2019

- Fully leveraging Intel's scale, resources and infrastructure

- Accessing Intel ecosystem and customer partnerships

- Delivering aggressive roadmap optimized for AI data center performance and efficiency

intel

&

habana

# Demand for compute for ML training doubles every 3.4 months

- ## Increasing Complexity

  - Businesses need higher precision in their model predictions

  - Results in larger and more complex models

  - Requires frequent retraining of models

- ## Increasing Costs

  - Increasing compute power required for frequent training of larger models drives up cost to train

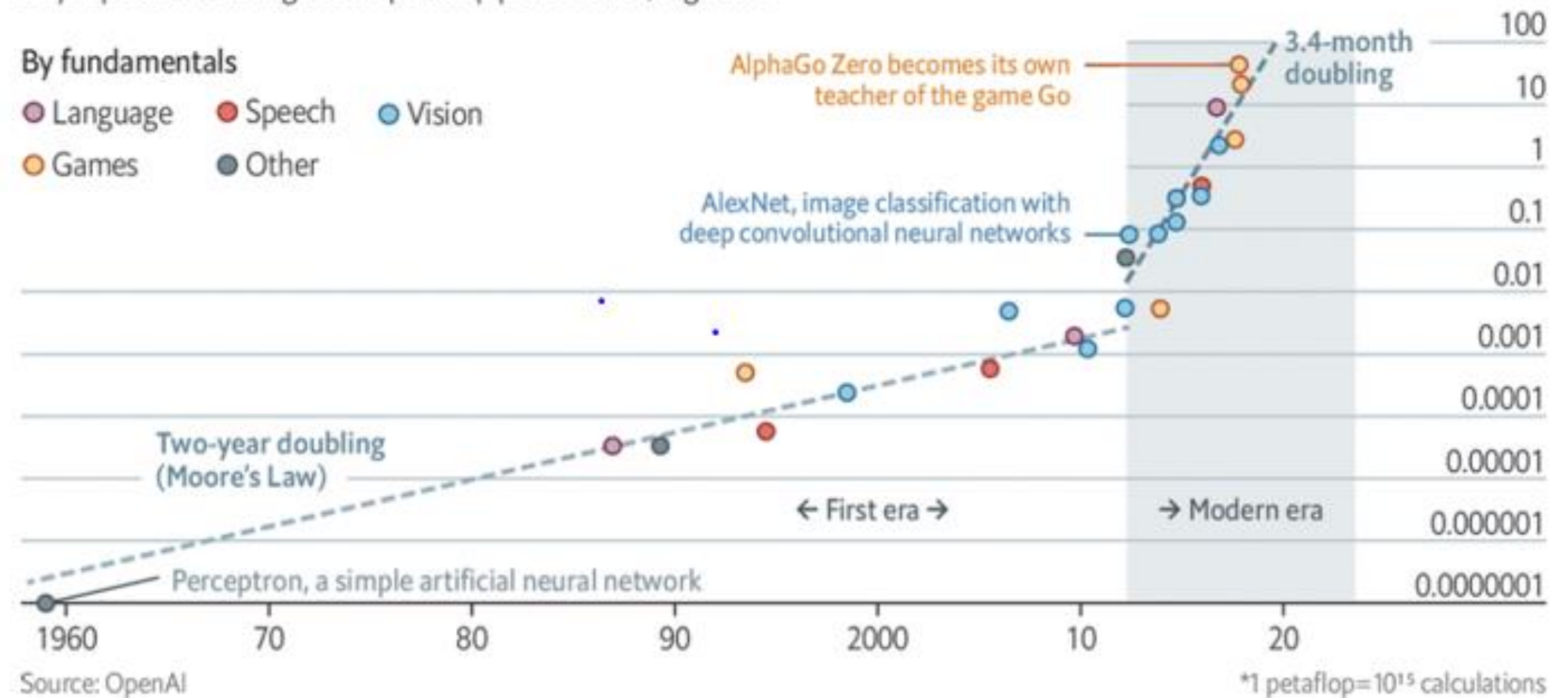  - Becomes a barrier for innovation and growth

**Deep and steep**

Computing power used in training AI systems
Days spent calculating at one petaflop per second*, log scale

By fundamentals
- Language
- Speech
- Vision
- Games
- Other

AlphaGo Zero becomes its own teacher of the game Go

3.4-month doubling

AlexNet, image classification with deep convolutional neural networks

Two-year doubling (Moore's Law)

← First era →   → Modern era

Perceptron, a simple artificial neural network

100
10
1
0.1
0.01
0.001
0.0001
0.00001
0.000001
0.0000001

1960   70   80   90   2000   10   20

Source: OpenAI
*1 petaflop=10$^{15}$ calculations
The Economist

**Need for dedicated AI processors to address the compute, memory and communication challenges**

# Today's Cost to Train: Biggest Barrier to AI Implementation

**"Cost is the most significant challenge to implementing AI/ML solutions."**
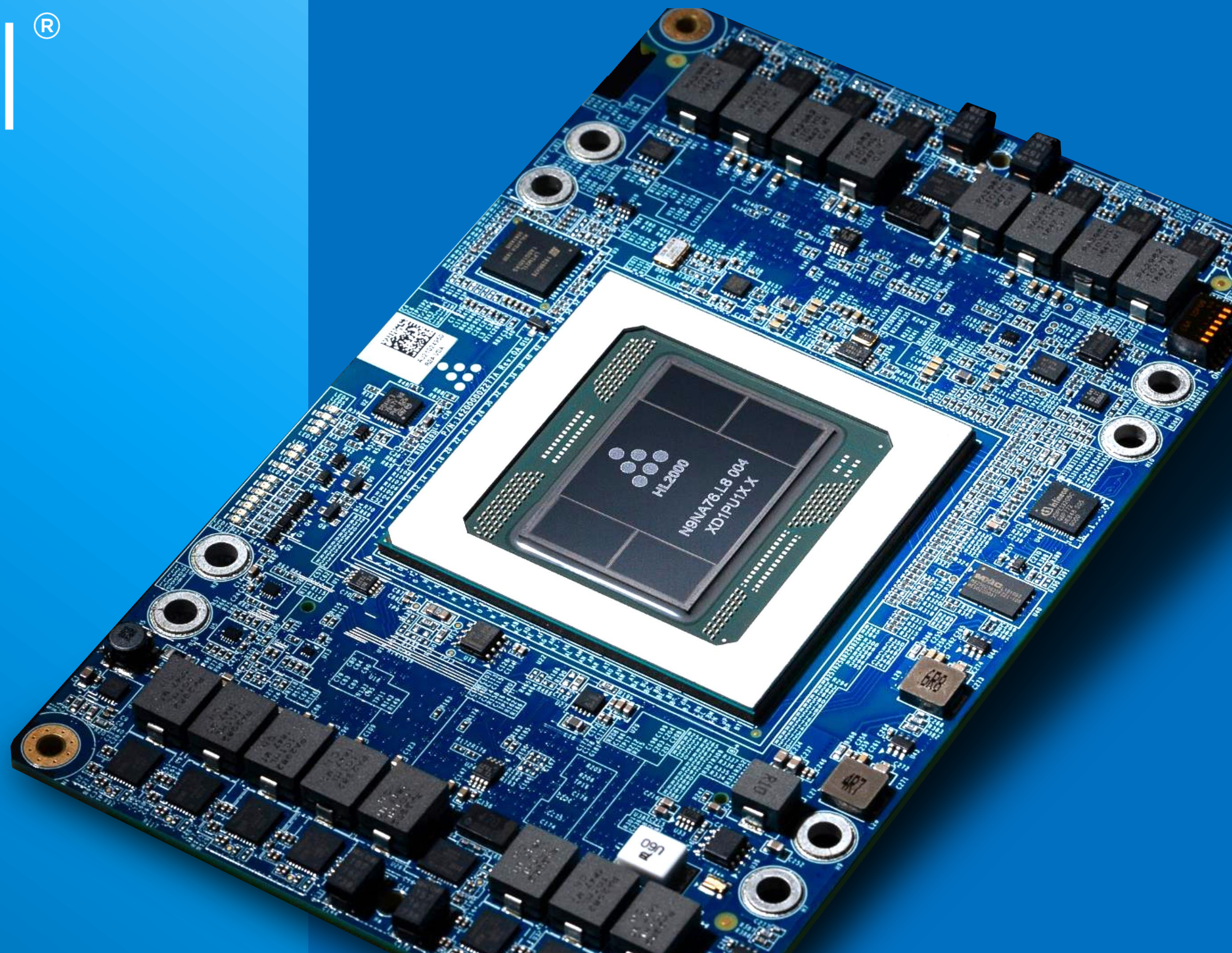
**56%** of AI/ML customers

## Industry Challenge:

How to give more customers access to more AI?

# GAUDI®

**Purpose-built for
AI training efficiency,
usability and scale**

# DL1 Model Training Cost Savings

**ResNet50 $/image**
(lower is better)

**BERT-Large $/seq**
(lower is better)

Pre-training Phase 1

Pre-training Phase 2

77%

49%
45%

64%

26%
22%

75%

52%
55%

■ Gaudi ■ A100-80G ■ A100-40G ■ V100-32G

■ Gaudi ■ A100-80G ■ A100-40G ■ V100-32G

■ Gaudi ■ A100-80G ■ A100-40G ■ V100-32G

# The Habana® Gaudi® AI Training Processor

GAUDI®

## Designed to optimize AI performance, delivering higher AI efficiency than traditional CPUs and GPUs

**Heterogeneous compute architecture enables high-efficiency on large AI workloads**
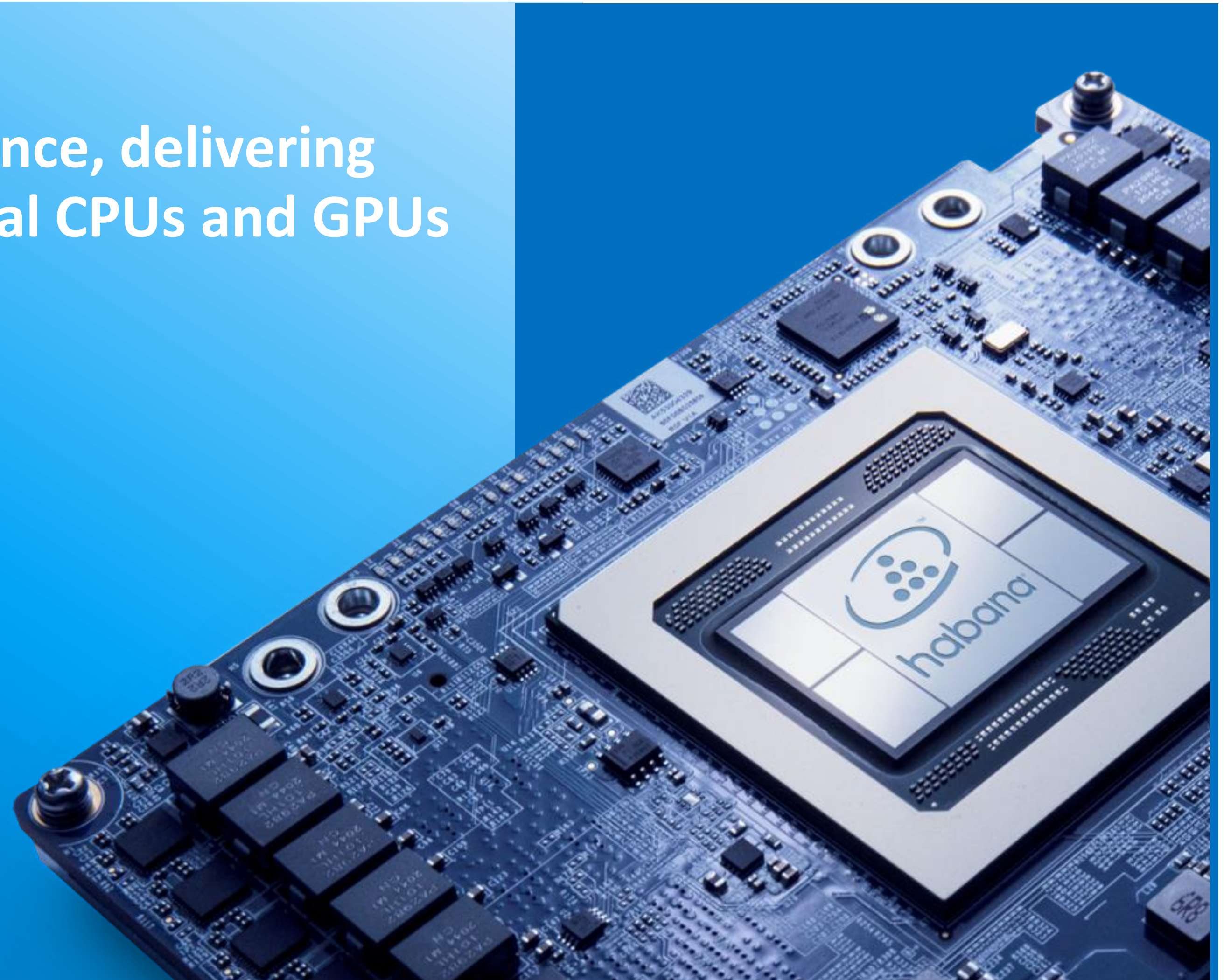
- GEMM engine (MME) excels at matrix multiplication
- While TPC runs non-linear and element wise ops

**Software-managed memory architecture**
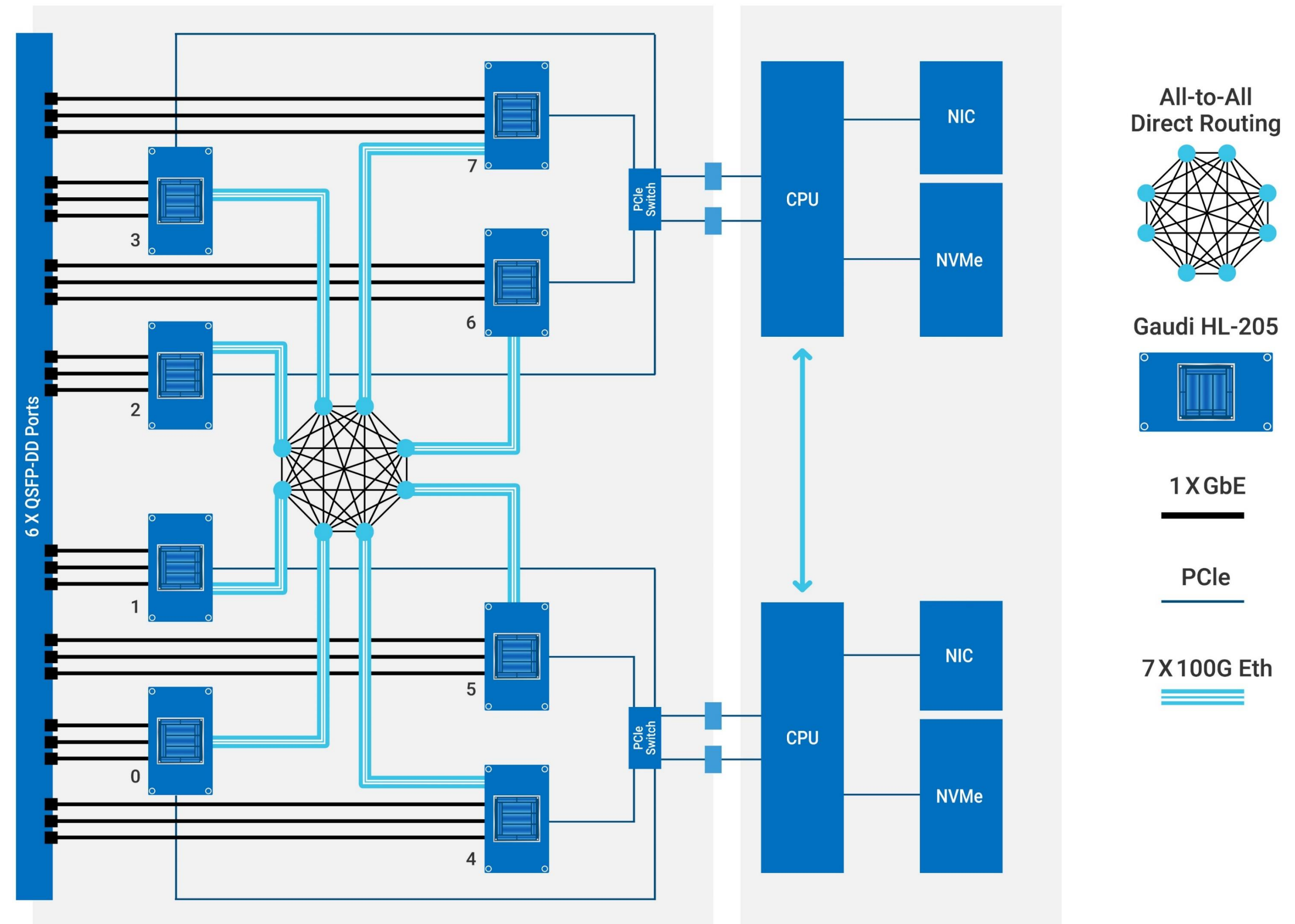
- 32 GB of HBM2 memory

**Integrates ten 100Gb Ethernet RoCE ports**

- Scaling capacity
- Flexibility based on industry standard
- Cost-efficiency with integrated NIC

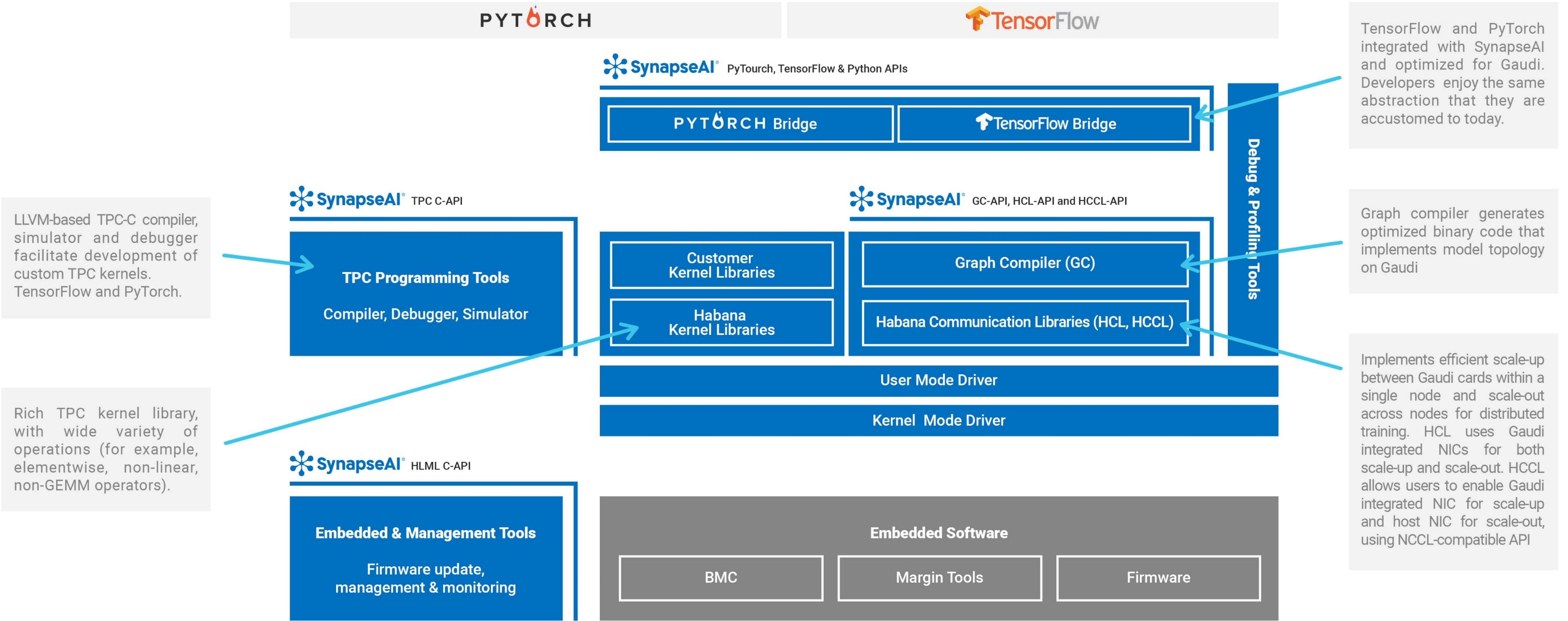# Scaling within a Gaudi Server

- 8 Gaudi OCP OAM cards

- 24 x 100GbE RDMA RoCE for scale-out

- Non-blocking, all-2-all internal interconnect across Gaudi AI processors

- Separate PCIe ports for external Host CPU traffic



Example of Integrated Server with eight Gaudi AI processors, two Xeon CPU and multiple Ethernet Interfaces

# Software Suite Detail



**PYTORCH**

**TensorFlow**

**SynapseAI** PyTourch, TensorFlow & Python APIs

**PYTORCH** Bridge

**TensorFlow** Bridge

**SynapseAI** TPC C-API

**SynapseAI** GC-API, HCL-API and HCCL-API

**Debug & Profiling Tools**

**TPC Programming Tools**

Compiler, Debugger, Simulator

Customer Kernel Libraries

Habana Kernel Libraries

Graph Compiler (GC)

Habana Communication Libraries (HCL, HCCL)

User Mode Driver

Kernel Mode Driver

**SynapseAI** HLML C-API

**Embedded & Management Tools**

Firmware update, management & monitoring

**Embedded Software**

BMC

Margin Tools

Firmware

TensorFlow and PyTorch integrated with SynapseAI and optimized for Gaudi. Developers enjoy the same abstraction that they are accustomed to today.

LLVM-based TPC-C compiler, simulator and debugger facilitate development of custom TPC kernels. TensorFlow and PyTorch.

Rich TPC kernel library, with wide variety of operations (for example, elementwise, non-linear, non-GEMM operators).

Graph compiler generates optimized binary code that implements model topology on Gaudi

Implements efficient scale-up between Gaudi cards within a single node and scale-out across nodes for distributed training. HCL uses Gaudi integrated NICs for both scale-up and scale-out. HCCL allows users to enable Gaudi integrated NIC for scale-up and host NIC for scale-out, using NCCL-compatible API

# Mobileye

Custom object detection
(2D and 3D) models trained on Gaudi

*"Multiple teams across Mobileye have chosen to use Gaudi-accelerated training machines, either on Amazon EC2 DL1 instances or on-prem; Those teams consistently see significant cost-savings relative to existing GPU-based instances across model types, enabling them to achieve much better Time-To-Market for existing models or training much larger and complex models...We're excited to see Gaudi2's leap in performance"*

Gaby Hayon, EVP R&D, Mobileye

intel  **Intel**
Jul 19 · 5 min read · ▶ Listen

TECHNOLOGY

# Mobileye journey towards scaling Amazon EKS to thousands of nodes leveraging Intel® Xeon® Scalable Processors and Habana's Gaudi AI accelerators

Authors: Diego Bailon Humpert, AWS EMEA and Global Automotive GTM Lead & David Peer, Mobileye AI Engineering DevOps specialist & team leader.

Mobileye is a company that develops autonomous driving technologies and advanced driver-assistance systems (ADAS) including cameras, computer chips, and software.

©Habana 2022

11

# Accelerating Medical Imaging Applications

## Objective

Demonstrate Gaudi DL1 AI processor cost-efficiency (price-performance ratio) for training deep learning models to detect novel coronavirus pneumonia in frontal chest X-ray images.

## Models

- Pretraining: CheXNet, to detect and localize multiple kinds of diseases from chest X-ray images.
- Finetuning: COVID-CXNet, to detect novel coronavirus pneumonia in frontal chest X-ray images
  - Transfer learning of CheXNet with a focus on Grad-CAM visualizations.

## Datasets

- 3200 normal images from NIH CXR dataset excluding age < 18 images based on paper
- 845 COVID-19 images from dataset used in the paper excluding age < 18 and early stage images.

# DL1 Cost Savings

## CheXNet-Keras

Dataset: ChestXray-NIHCC

Batch size: 32

Precision: FP32

Device count: 8

| Instance | On-Demand hourly rate of EC2 instance [$/Hour] | Time per epoch [Seconds] | Cost per epoch [$] | DL1 Cost Savings to EC2 Customers [%] |
|---|---|---|---|---|
| 8x V100-32 GB* (p3dn.24xlarge) | $31.21 | 4.6 | $143.57 | 59% |
| 8x Gaudi DL1.24xlarge** | $13.11 | 4.47 | $58.56 | |

## COVID-CXNet

Dataset: COVID-CXNet

Batch size: 16

Precision: BF16

Device Count: 1

| Instance | On-Demand hourly rate of EC2 instance [$/Hour] | Time per epoch [Seconds] | Cost per epoch [$] | DL1 Cost Savings to EC2 Customers [%] |
|---|---|---|---|---|
| 8x V100-32 GB* (p3dn.24xlarge) | $31.21 | 718 | $6.22 | 67% |
| 8x Gaudi DL1.24xlarge** | $13.11 | 565 | $2.06 | |

Source: Leidos

# Summary

- Using Amazon EC2 DL1 instances for Chest X-Ray COVID Detection model pretraining and finetuning resulted in 60%+ savings in cost of training

- Successfully trained deep learning models on EC2 DL1 platform with minimal code changes

- Excellent support and documentation available on Habana Developer Site https://developer.habana.ai and GitHub with reference models

# Accelerating Medical Benefit Application Processing

Leidos customer using NLP-based deep learning solution to facilitate medical benefit application processing

**Objective:** Demonstrate price/performance of Gaudi based EC2 DL1.24xlarge instance versus GPU based G4DN.12xlarge EC2 instance used by customer

- TensorFlow DistilBERT Model finetuned for a multi-labeling classification task

- Trained with 737k labeled examples and tested against 184k test examples

| Attribute | G4DN.12xlarge (GPU based) | DL1.24xlarge (Gaudi Based) |
|---|---|---|
| Memory | 192 GB | 768 GB |
| Accelerator | 4 x Tesla T4 | 8 x Gaudi |
| Accelerator Type | GPU | HPU |
| On Demand Cost (per hour) | $3.91 | $13.11 |

# Results

## Cost Performance

- Compared to GPU, training took only 45% of the time with 1x Gaudi processor and only 10% of the time when all 8 Gaudi processors were used

- Although the DL1 instance costs more per hour, due to shorter training time, the total cost for model training ended up being only 1/3 of the baseline, i.e., 66% cost savings.

| Device | Training Time | Training Cost | DL1 cost savings |
|---|---|---|---|
| GPU (g4dn) | 10 hrs | $39.12 | |
| Gaudi x1 (dl1) | 4.5 hrs | $59.00 | |
| Gaudi x8 (dl1) | 1 hr | $13.11 | 66% |



Time per epoch



Cost per epoch

## User Experience

- Performing single card training is quite intuitive and simple

- Distributed training is simple when using Horovod and OpenMPI

Source: Leidos

# Summary

- Gaudi performs significantly better than NVIDIA Tesla T4 when looking at the time and cost metrics

- DL1 is worthy of strong consideration as part of cloud-forward strategy, especially when an organization anticipates using deep learning models.

- For Leidos' customers who are interested in continuing training, low cost is a very attractive feature

Areas for Future Work:

- With lower cost to train with Gaudi, one can potentially train/update more complex and accurate models

- Pre-training or continue to pre-train domain specific model, which is a more computing expensive task. This is critical for domain adaptation, which impacts a broad range of NLP tasks and is relevant to lot of our clients.

# > 60% cost savings with DL1 vs. GPU instances

**leidos**

*"Given Leidos and its customers' need for quick, easy, and cost-effective training for deep learning models, we are excited to have begun this journey with Intel and AWS to use Amazon EC2 DL1 instances based on Habana Gaudi AI processors."*

Chetan Paul, CTO Health and Human Services at Leidos

# Habana AI Powers SDSC's Voyager Research Program

## 336 Gaudi Training accelerators with native RoCE scaling and 16 Goya Inference processors

- In service since Fall of 2021

- Funded by $5M grant from National Science Foundation

  - Matching funds targeting community support and operation

- AI research conducted across range of science and engineering domains

  - Astronomy, climate sciences, chemistry, particle physics,

- Announced by SDSC in July 2020, more information here.



VOYAGER EXPLORING AI PROCESSORS in SCIENCE and ENGINEERING

**3-YEAR TESTBED PHASE**
Focused Select Projects
Workshops, Industry Interaction

**2-YEAR ALLOCATIONS PHASE**
NSF Allocations to the Broader Community
User Workshops

**INNOVATIVE AI RESOURCE**
Specialized Training Processors
Specialized Inference Processors
High-Performance Interconnect
X86 Standard Compute nodes
Rich Storage Hierarchy

**IMPACT & ENGAGEMENT**
Large-Scale Models
AI Architecture Advancement
Improved Performance of AI Applications
External Advisory Board of AI & HPC Experts
Wide Science & Engineering Community
Advanced Project Support & Training
Accelerating Scientific Discovery
Industrial Engagement

**OPTIMIZED AI SOFTWARE**
Community Frameworks
Custom user-developed AI Applications
PyTorch, Tensorflow

# Combining Fire Science With AI For Wildfire Mitigation

DL algorithms of satellite images determine land covers across geographies in the context of wildfire management

*"With innovative solutions optimized for deep learning operations and AI workloads, Habana accelerators are excellent choices to power Voyager's forthcoming AI research"*

Amit Majumdar, Director of Data Enabled Scientific Computing Division, SDSC

# First-generation Gaudi

## IN THE CLOUD

- AWS EC2 DL1 Instances
- Leading AWS AI training efficiency

## ON PREMISES

- Supermicro X12 Gaudi Server
- DDN AI 400X2 storage solution

MLOps SOFTWARE

cnvrg.io

# Getting Started with TensorFlow on Gaudi®

```python
import tensorflow as tf

from TensorFlow.common.library_loader import load_habana_module
load_habana_module()

(x_train, y_train), (x_test, y_test) =
tf.keras.datasets.mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

model = tf.keras.models.Sequential([
                    tf.keras.layers.Flatten(input_shape=(28, 28)),
                    tf.keras.layers.Dense(10),
])
loss = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
optimizer = tf.keras.optimizers.SGD(learning_rate=0.01)

model.compile(optimizer=optimizer, loss=loss, metrics=['accuracy'])

model.fit(x_train, y_train, epochs=5, batch_size=128)
model.evaluate(x_test, y_test)
```

Load the Habana® libraries needed to use Gaudi aka **HPU** device

Once loaded, the **HPU** device is registered in TensorFlow

When an Op is available for both CPU and HPU, Op is assigned to the HPU

When an Op is not supported on HPU, it runs on the CPU

# Getting Started With PyTorch Lightning On Gaudi

```python
import pytorch_lightning as pl
from pytorch_lightning.plugins import HPUPrecisionPlugin

# mixed precision distributed training with 8 Gaudis
trainer = pl.Trainer(accelerator="hpu", devices=<n>, precision=16)
```

All you need is to provide **accelerator="hpu"** parameter to the Trainer class

Select the number of Gaudi devices, **n**=1..8

For mixed precision training, import **HPUPrecisionPlugin** and set "**precision=16**"

Lightning 1.6 now supports HPU with SynapseAI 1.4: https://pytorch-lightning.readthedocs.io/en/stable/accelerators/hpu.html

# Getting Started With Huggingface On Gaudi

```python
from optimum.habana import GaudiConfig, GaudiTrainer, GaudiTrainingArguments
from transformers import BertTokenizer, BertModel
…
tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
model = BertModel.from_pretrained("bert-base-uncased")
gaudi_config = GaudiConfig.from_pretrained("Habana/bert-base-uncased")
args = GaudiTrainingArguments(
    output_dir="/tmp/output_dir",
    use_habana=True,
    use_lazy_mode=True,
)
trainer = GaudiTrainer(
    model=model,
    gaudi_config=gaudi_config,
    args=args,
    tokenizer=tokenizer,
)

trainer.train()
```

- Uses Optimum Habana library

- Model instantiated the same way as in the Transformers library

- Only difference is to load Gaudi configuration and provide to the Gaudi trainer

# Habana Developer Platform---developer.habana.ai

# Habana Developer Documentation---docs.habana.ai

**🏠 Gaudi Documentation**
latest

Search docs

**GETTING STARTED**

Gaudi Architecture and Software Overview

Support Matrix

Release Notes

Installation

**GUIDES**

TensorFlow

PyTorch

PyTorch Lightning

Profiling

Management and Monitoring

Orchestration

AWS Quick Start Guides

APIs

TPC Programming

**LEGAL NOTICE**

Legal Notice and Disclaimer

## Welcome to Habana® Gaudi® v1.5 Documentation

Find detailed documentation to learn how to use the Habana Gaudi solutions - first-generation Gaudi and Gaudi2. This will cover the details on how to migrate models to Habana, code samples, diagrams, best practices for debug and optimization, API references, and more.

### Getting Started

Start using Habana Gaudi Processors

Click here to get started

### Tutorials

Tutorials to show basic examples of how to run on TensorFlow and PyTorch

Click Here

### Model Catalog

Start with TensorFlow and PyTorch models already running on Gaudi

Click Here

### User Forum

Post questions and get help in the User Forum

Click Here

Next ➡

Feedback

© Copyright 2022, Habana Labs Revision 78a40cbe.

developer.habana.ai | Legal Notice and Disclaimer | Outbound Software License Agreement | Send Feedback

# Habana Developer Software---vault.habana.ai

# Habana GitHub Repositories---github.com/HabanaAI

# Gaudi Reference Models---August 2022

## TensorFlow

| | |
|---|---|
| ResNet50 Keras | BERT |
| ResNeXt101 | DistilBERT |
| SSD | ALBERT |
| Mask R-CNN | Transformer |
| DenseNet | T5 Base |
| UNet 2D | Electra |
| UNet 3D | |
| UNet Industrial | |
| CycleGAN | |
| EfficientDet | |
| RetinaNet | |
| SegNet | |
| Vision Transformer | |
| MobileNet V2 | |

## PyTorch

| | |
|---|---|
| ResNet50, ResNeXt101, ResNet152 | BERT Pretraining |
| MobileNet V2 | BERT Finetuning |
| UNet 2D, Unet 3D | DeepSpeed BERT-1.5B |
| SSD | RoBERTa |
| GoogLeNet | ALBERT |
| Vision Transformer | DistilBERT |
| Swin Transformer | Electra |
| DINO | Transformer |
| | BART |
| | GPT2 |

# Habana Developer Forum---forum.habana.ai



30

# GAUDI®2

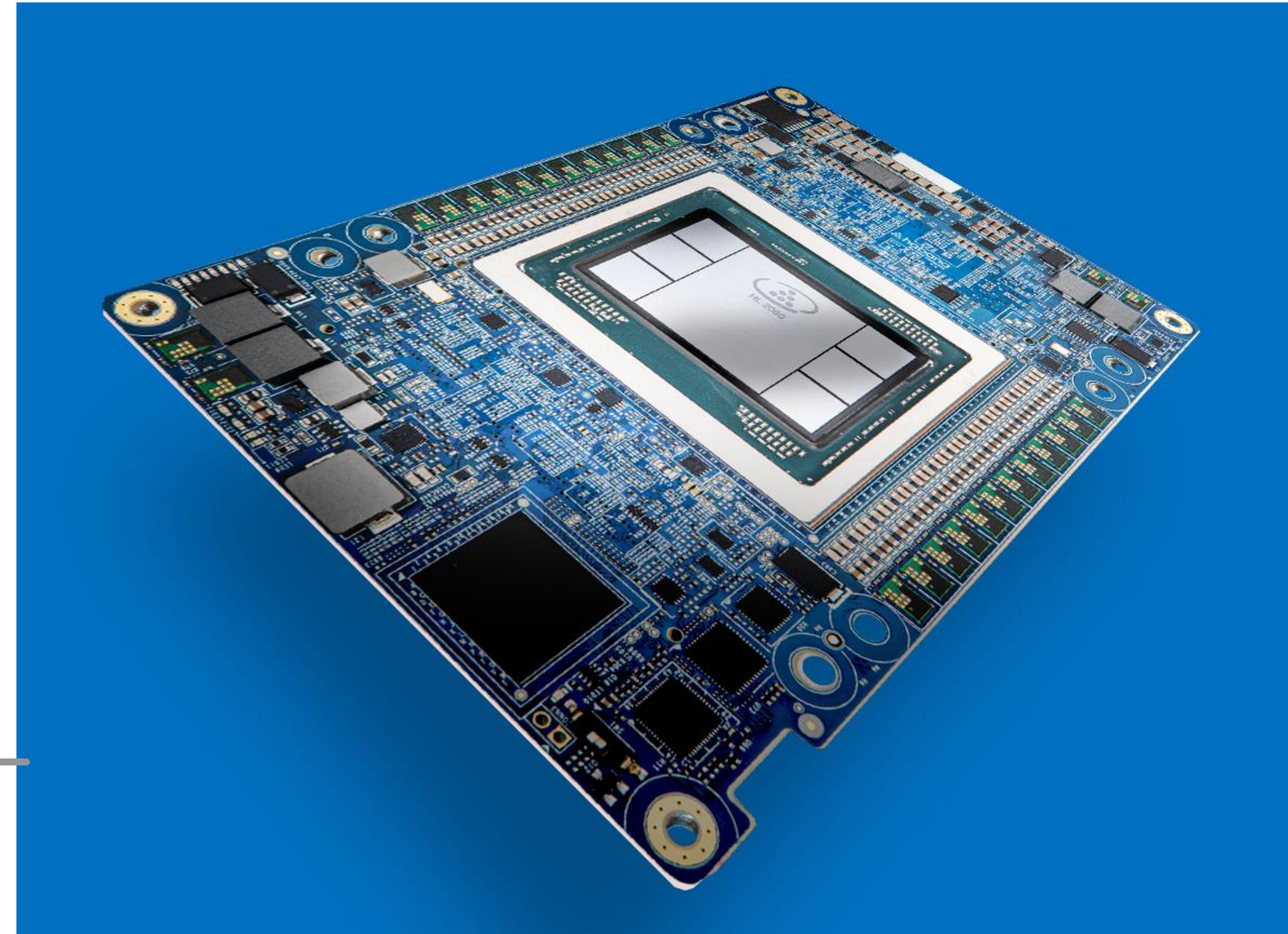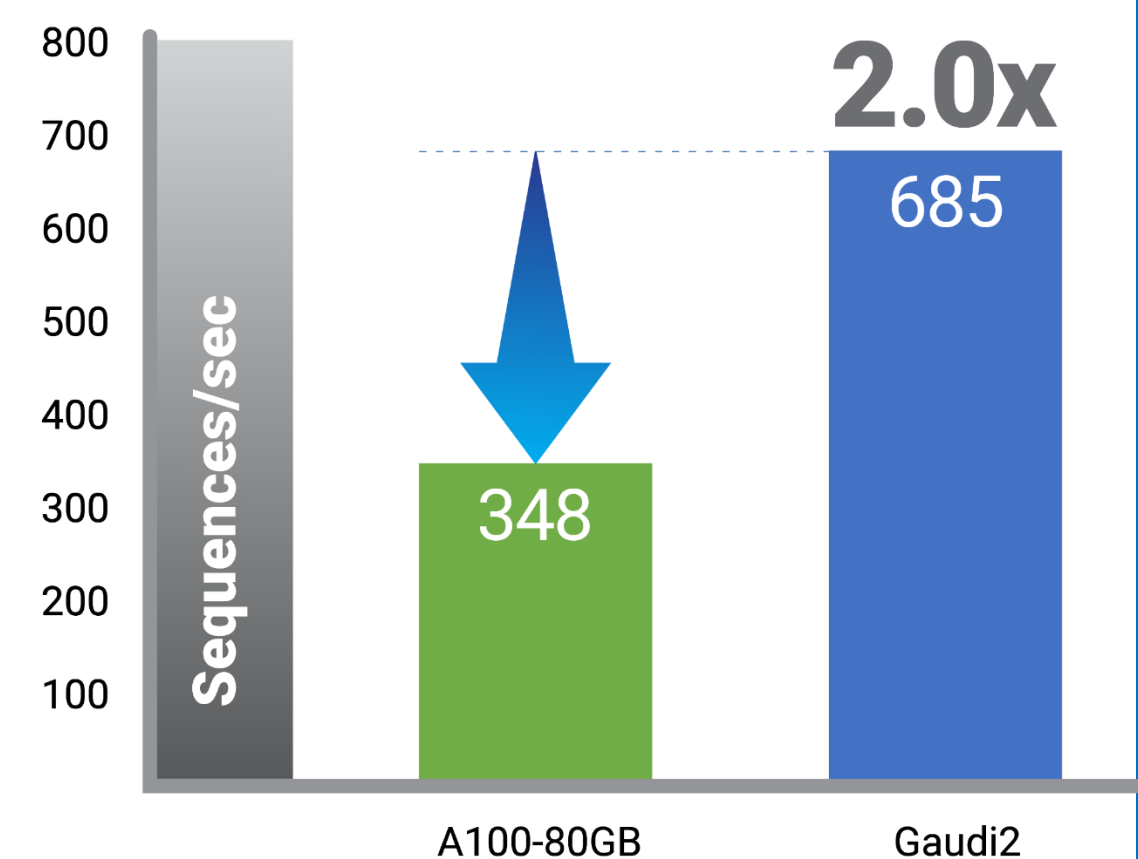## Leadership Performance

~2x better throughput vs A100
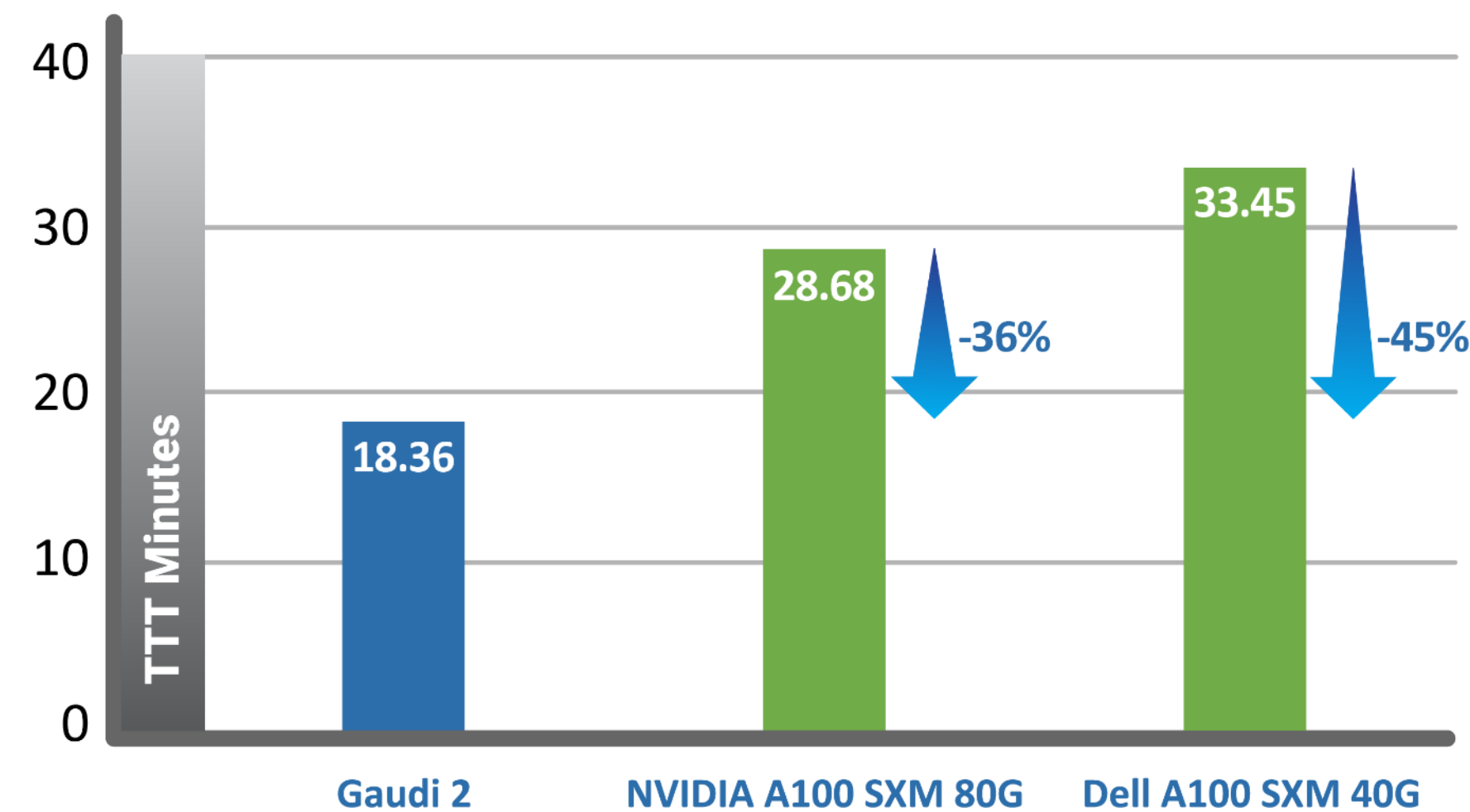for popular vision and language models

**ResNet50 Training Throughput**

**1.9x**

| | A100-80GB | Gaudi2 |
|---|---|---|
| Images/sec | 2,930 | 5,425 |

**BERT Training Throughput**

**2.0x**

| | A100-80GB | Gaudi2 |
|---|---|---|
| Sequences/sec | 348 | 685 |

# Gaudi2: Second-generation Training & Inference

## GAUDI® / GAUDI® 2

### Node
**TPC** [X] MME

1 16nm → **7nm**

### Compute

1 8 TPCs New → **24** TPCs
**Media decode & processing**
**FP8**

### Memory

1 1 TB/s, 32GB HBM2, 24 MB → **2.45** TB/s, **96** GB HBM2e, **48** MB SRAM

### Networking

1 10 x100 → **24 x 100 GbE**

### TDP

1 350w → **600w**

# Gaudi2 MLPerf June '22 Training Benchmark Results

Gaudi2 outperformed Nvidia A100 MLPerf submissions on both ResNet and BERT

**MLPERF ResNet-50 Training Time** [lower is better]
( 8 accelerator server )

TTT Minutes

| | | |
|---|---|---|
| Gaudi 2: 18.36 | NVIDIA A100 SXM 80G: 28.68 (-36%) | Dell A100 SXM 40G: 33.45 (-45%) |

**MLPERF BERT Training Time** [lower is better]
( 8 accelerator server )

TTT Minutes

| | | |
|---|---|---|
| Gaudi 2: 17.2 | NVIDIA A100 SXM 80G: 18.44 (-7%) | Dell A100 SXM 40G: 26.52 (-35%) |

...and First-gen Gaudi achieved near-ideal linear scale on 128- and 256-accelerators

## *Gaudi2 Time-to-Train (TTT) improved by 3 to 4.7x compared to First-gen Gaudi*

# MLPerf Press Coverage

# Greco: Second-Generation Inference for Deep Learning

GOYA™
GRECO™

**Node**

16nm → **7nm**

**Memory**

| 40GB/s DDR4 | **204GB/s LPDDR5** |
| 16GB | **16GB** |
| 50MB | **128MB on chip SRAM** |

**Compute**

BF16, FP16, INT4

Media decode and processing

**Form Factor**

Dual-slot PCIe >
Single-slot HHHL

**TDP**

200w → **75w**

# TensorFlow integration with SynapseAI



SynapseAI receives a computational graph of the model from the framework

It identifies subgraphs (blue nodes) that can be accelerated by Gaudi

The rest of the graph runs on CPU (yellow node)

The original graph is modified to replace the Gaudi subgraphs with encapsulated nodes (blue)

The framework runtime executes the modified graph

For each encapsulated node, SynapseAI generates optimized binary code that runs on Gaudi

# Software Installation and Deployment

[Setup and Install](#) repository on Habana GitHub provides instructions on how to setup your environment with the SynapseAI software stack

**SynapseAI Orchestration**
(Kubernetes Gaudi plugin, Kubeflow mpi-operator)

**SynapseAI TensorFlow Container Image**
(TensorFlow frontend, horovod, open-mpi)

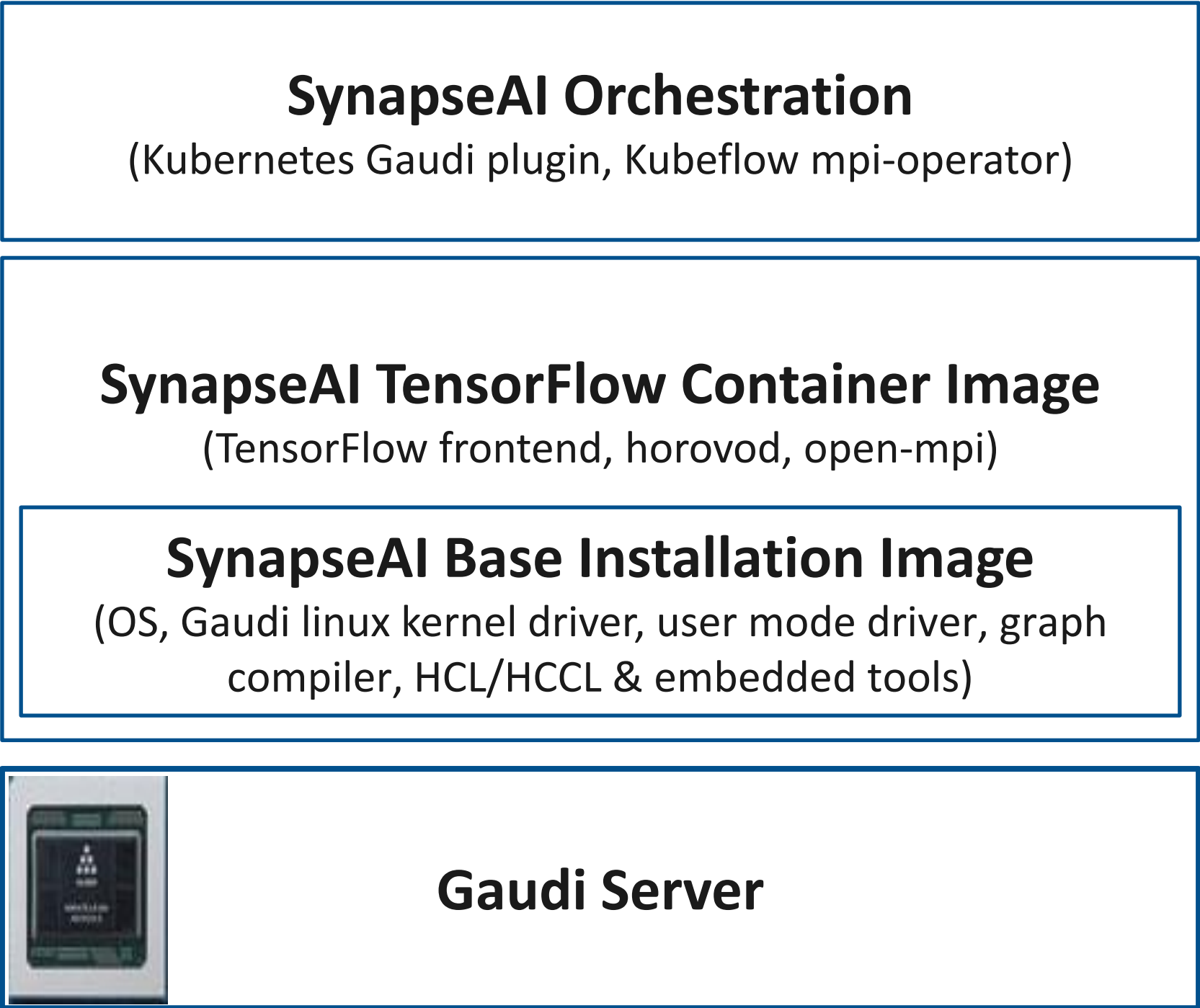**SynapseAI Base Installation Image**
(OS, Gaudi linux kernel driver, user mode driver, graph compiler, HCL/HCCL & embedded tools)

**Gaudi Server**

Gaudi-optimized Docker container images with all necessary dependences*

Official releases publicly available on Habana Vault

| | |
|---|---|
| Orchestration | Kubernetes (1.19) |
| Frameworks | TensorFlow2 and PyTorch |
| Operating Systems | Ubuntu 18.04 and 20.04 |
| Container Runtimes | Docker (Docker CE version 18.09) |
| Distributed Training Schemes | TensorFlow with Horovod and tf.distribute PyTorch distributed (native) |

*Habana GitHub will have repository with Dockerfiles to "build your own" Docker images*

# DL1 Vision Model Training Performance



**ResNet50 Throughput (images/second)**

| | | |
|---|---|---|
| 8x | 8 x A100 | 17400 |
| | 8 x Gaudi | 12880 |
| | 8 x V100 | 9510 |
| 1x | | 2375 |
| | | 1695 |
| | | 1270 |

Legend: ■ A100 / P4d  ■ Habana Gaudi / DL1  ■ V100 / P3

Habana ResNet50 Model: https://github.com/HabanaAI/Model-References/tree/master/TensorFlow/computer_vision/Resnets/resnet_keras
Habana SynapseAI Container: https://vault.habana.ai/ui/repos/tree/General/gaudi-docker/1.2.0/ubuntu20.04/habanalabs/tensorflow-installer-tf-cpu-2.7.0
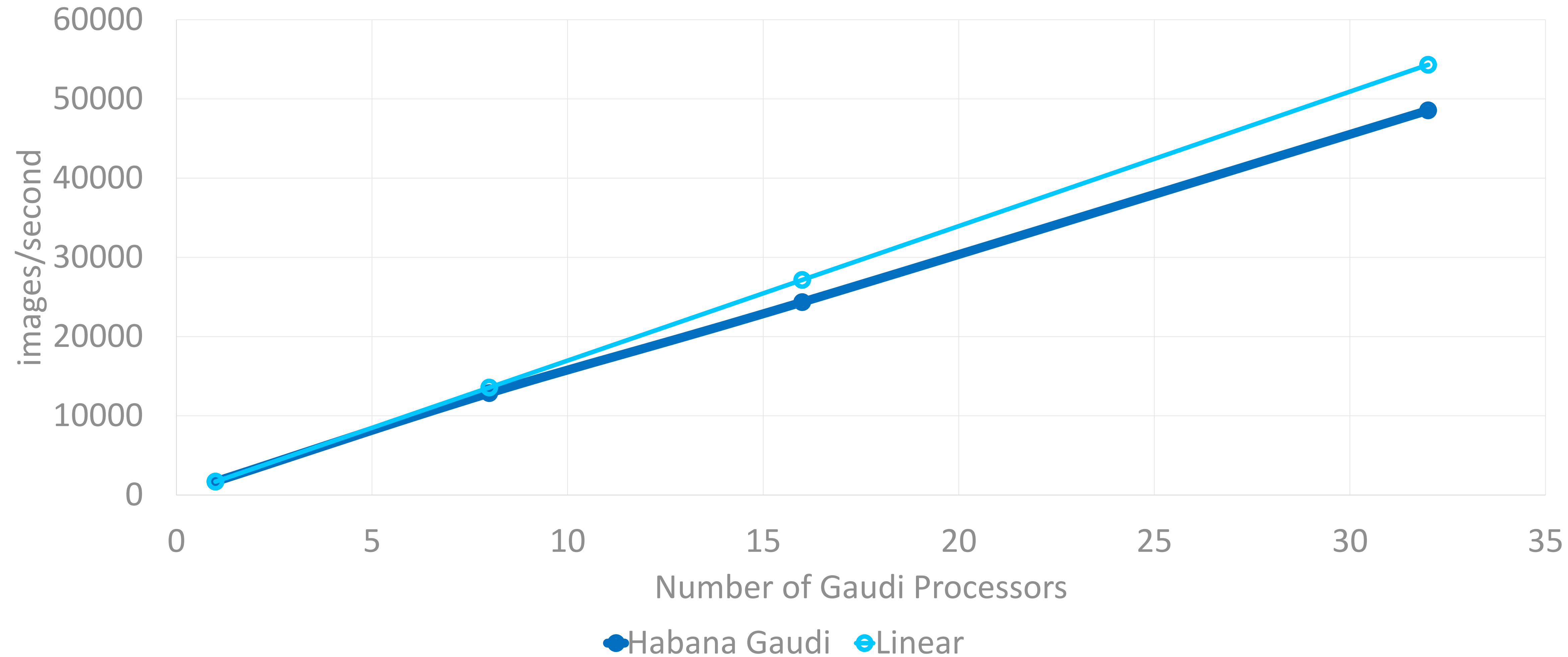Habana Gaudi Performance: https://developer.habana.ai/resources/habana-training-models/
A100 / V100 Performance Source: https://ngc.nvidia.com/catalog/resources/nvidia:resnet_50_v1_5_for_tensorflow/performance, results published for DGX A100-40G and DGX V100-32G
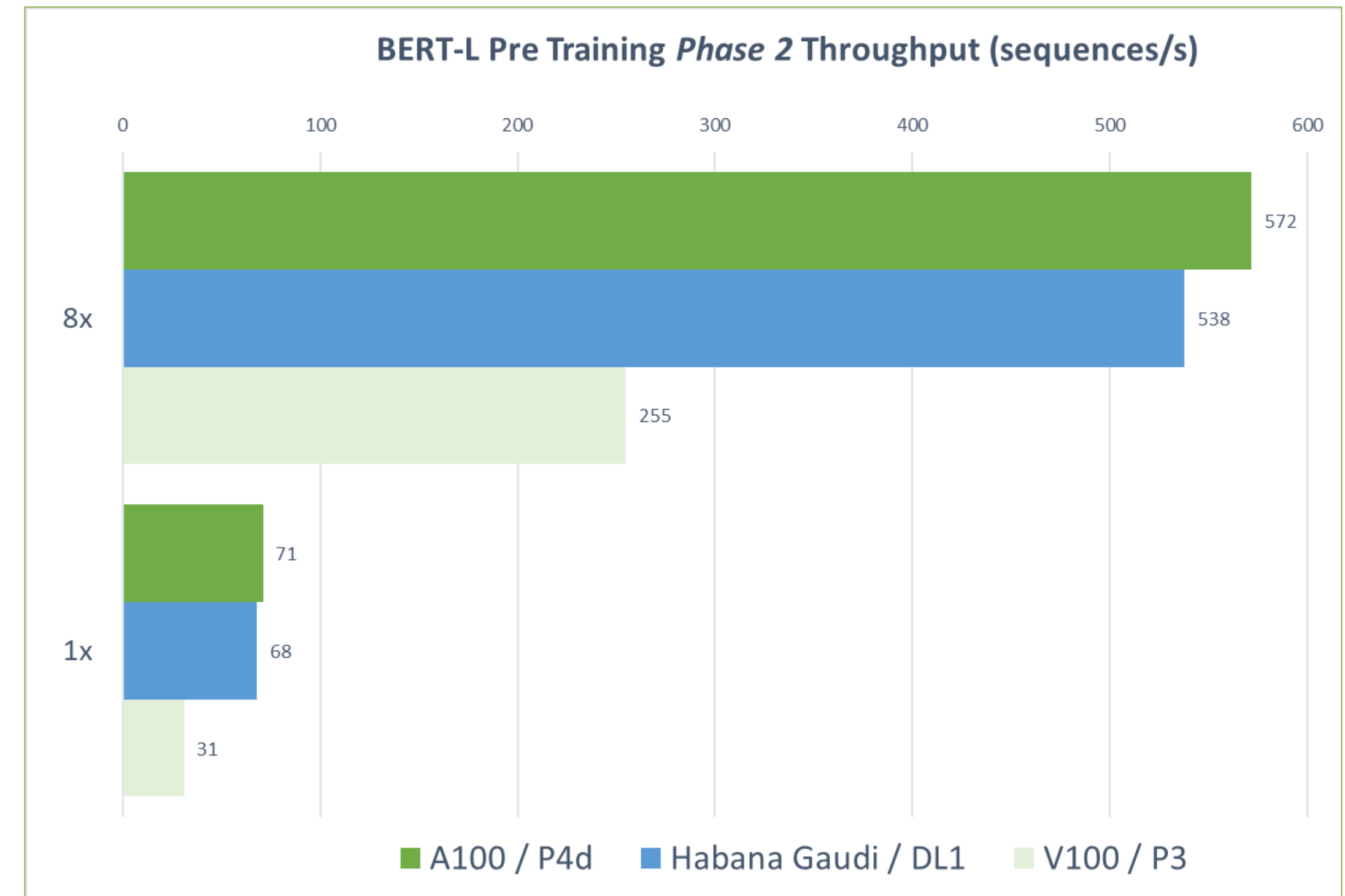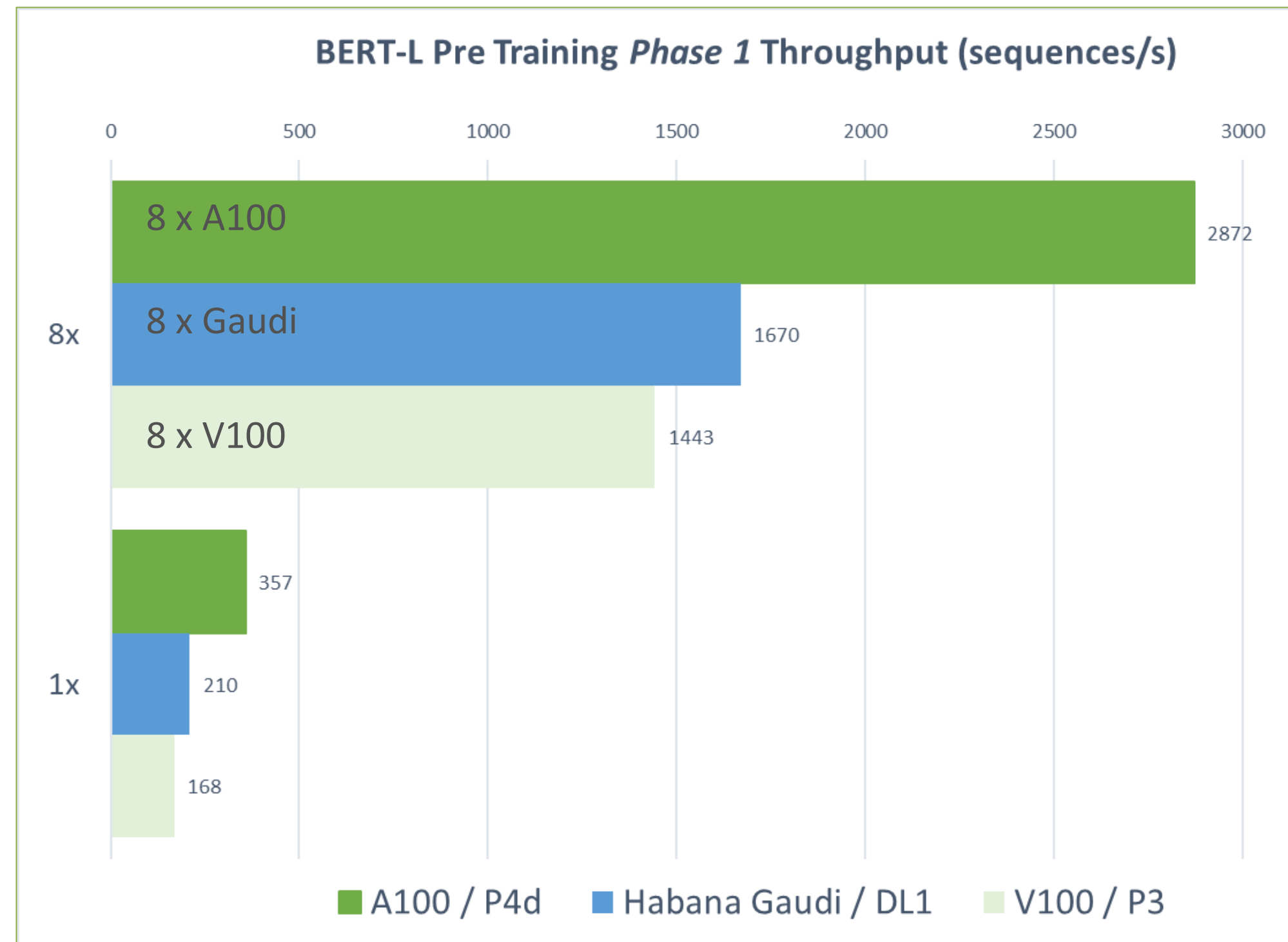
# Vision Model Training Scalability



ResNet 50 Training Throughput

# DL1 NLP Model Training Performance



**BERT-L Pre Training *Phase 1* Throughput (sequences/s)**

| | |
|---|---|
| 8 x A100 | 2872 |
| 8 x Gaudi | 1670 |
| 8 x V100 | 1443 |
| 1x (A100) | 357 |
| 1x (Gaudi) | 210 |
| 1x (V100) | 168 |

Legend: ■ A100 / P4d  ■ Habana Gaudi / DL1  ■ V100 / P3

**BERT-L Pre Training *Phase 2* Throughput (sequences/s)**

| | |
|---|---|
| 8x (A100) | 572 |
| 8x (Gaudi) | 538 |
| 8x (V100) | 255 |
| 1x (A100) | 71 |
| 1x (Gaudi) | 68 |
| 1x (V100) | 31 |

Legend: ■ A100 / P4d  ■ Habana Gaudi / DL1  ■ V100 / P3

# NLP Model Training Scalability



BERT Large Pretraining Throughput
(Phase 1 & Phase 2)

BERT Large Finetuning Throughput
(Phase 1)

BERT Large Pretraining Throughput
(Phase 2)