



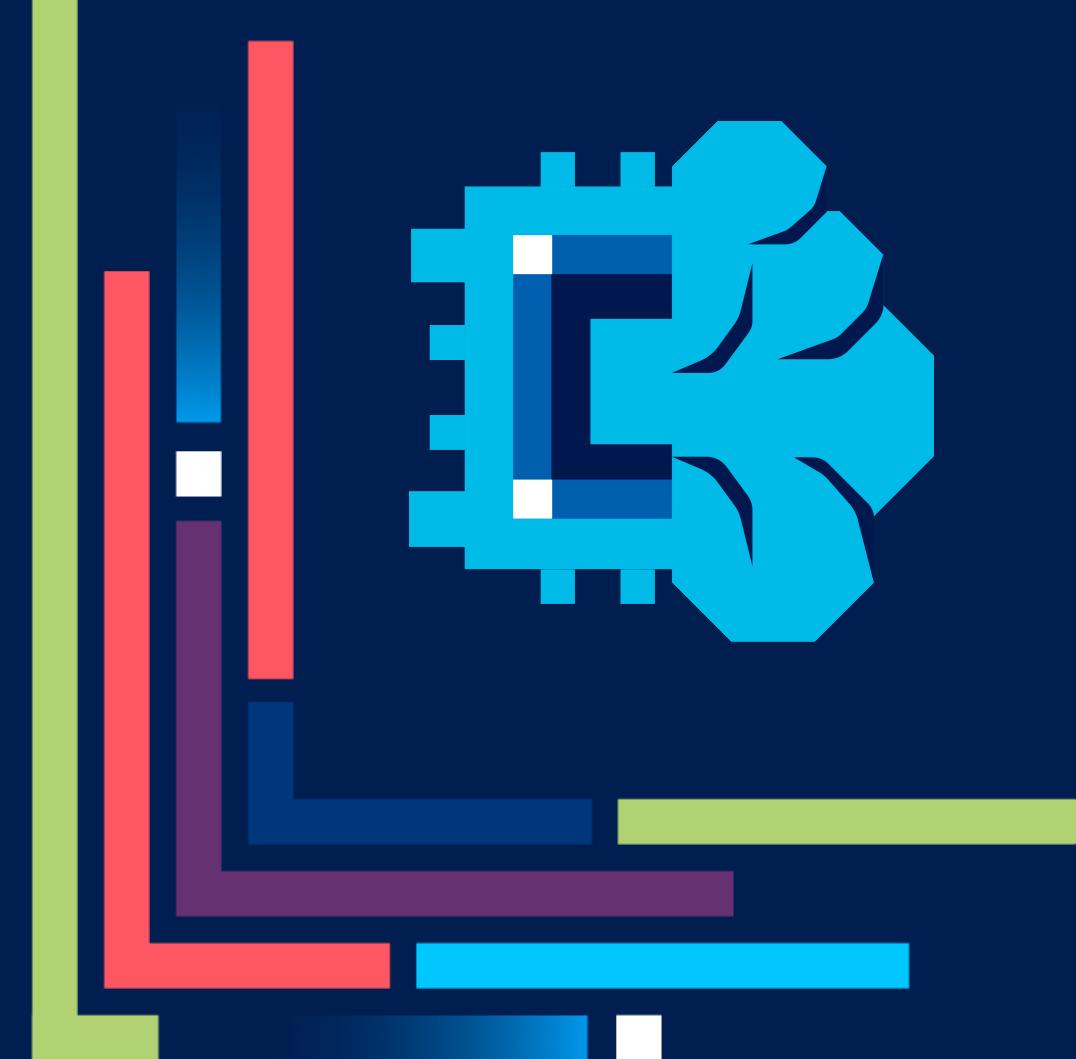
A New Paradigm for AI Productivity & Performance

Seetha . Nookala

Senior Director

Datacenter Graphics, HPC/AI ,APAC & Japan

December 2023



AI Exponentially Increasing In Quality & Size



Yet Inaccessible For Most – Cost of AI Today

Training Cost

GPT-3

\$1.65M

(3,640 petaFLOPS-days) costs if trained on Google TPU v3

GPT-4

\$40M

(450,000 petaFLOPS-days), 7,600 GPUs running for a year

Inferencing Cost

ChatGPT

\$40M

to process prompts per month with 100 million active users

Bing AI Chatbot

\$4B

Bing AI Chatbot to serve responses to all Bing users

Driving AI Ubiquity Will Unlock New Levels of Innovation

Small AI Models

Smartphones & Mobile Devices	Smart Cameras & Surveillance	Industrial IoT	Autonomous Vehicles	Edge Servers & Gateways
Facial Recognition	Object Detection	Predictive Maintenance	Object Detection	Local Decision-Making
Voice Assistants	Anomaly Detection	Quality Control	Driver Monitoring	Data Filtering
Drones & UAVs	Smart Home Devices	Wearables	Robots	Retail and Customer Analysis
Obstacle Avoidance	Thermostats	Health Monitoring	Navigation	Foot Traffic Analysis
Aerial Data Processing	Home Security	Gesture Control	Object Manipulation	Personalized Marketing
Agriculture Sensors & Equipment	Healthcare	Environmental Monitoring	Smart Grids & Energy Management	Sports and Entertainment
Livestock Monitoring	Remote Patient Monitoring	Air Quality Sensing	Power Distribution	Player and Audience Analytics
Precision Farming	Medical Imaging	Earthquake Detection	Grid Security	AR Headsets

Large

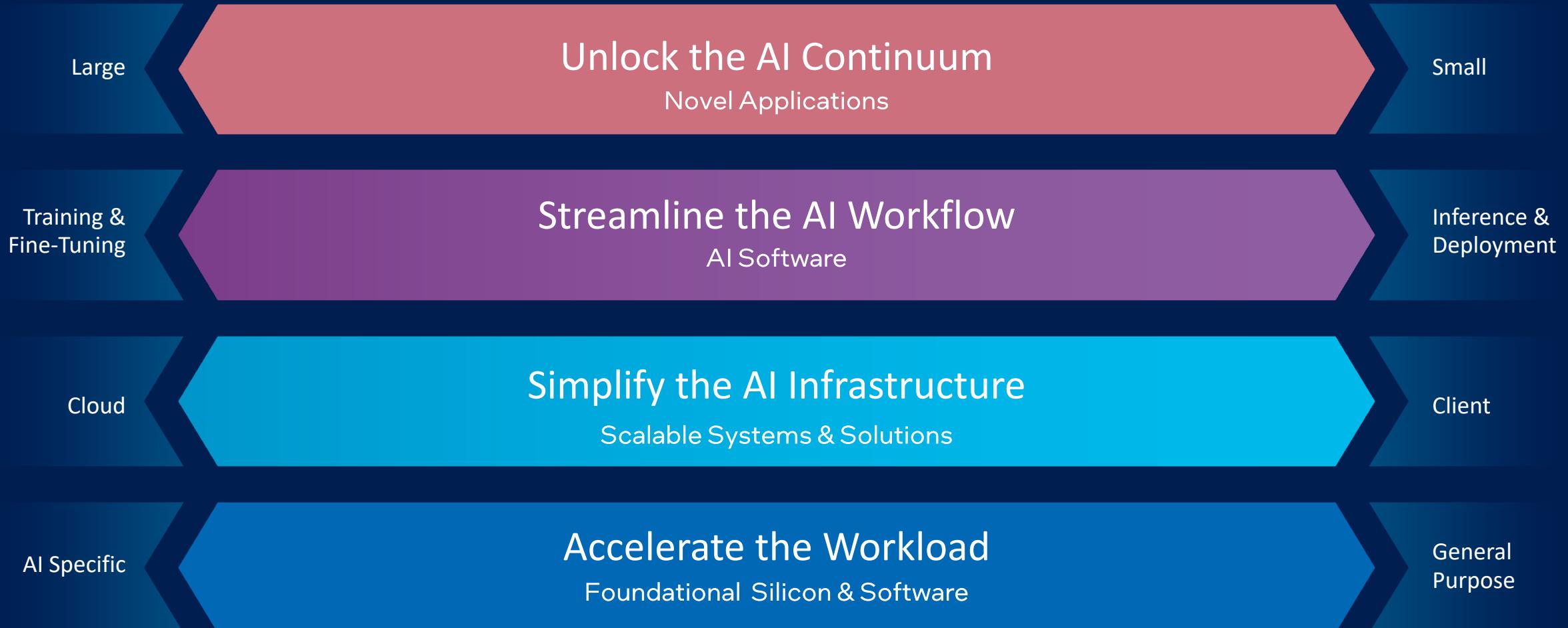
Unlocking the AI Continuum

Small

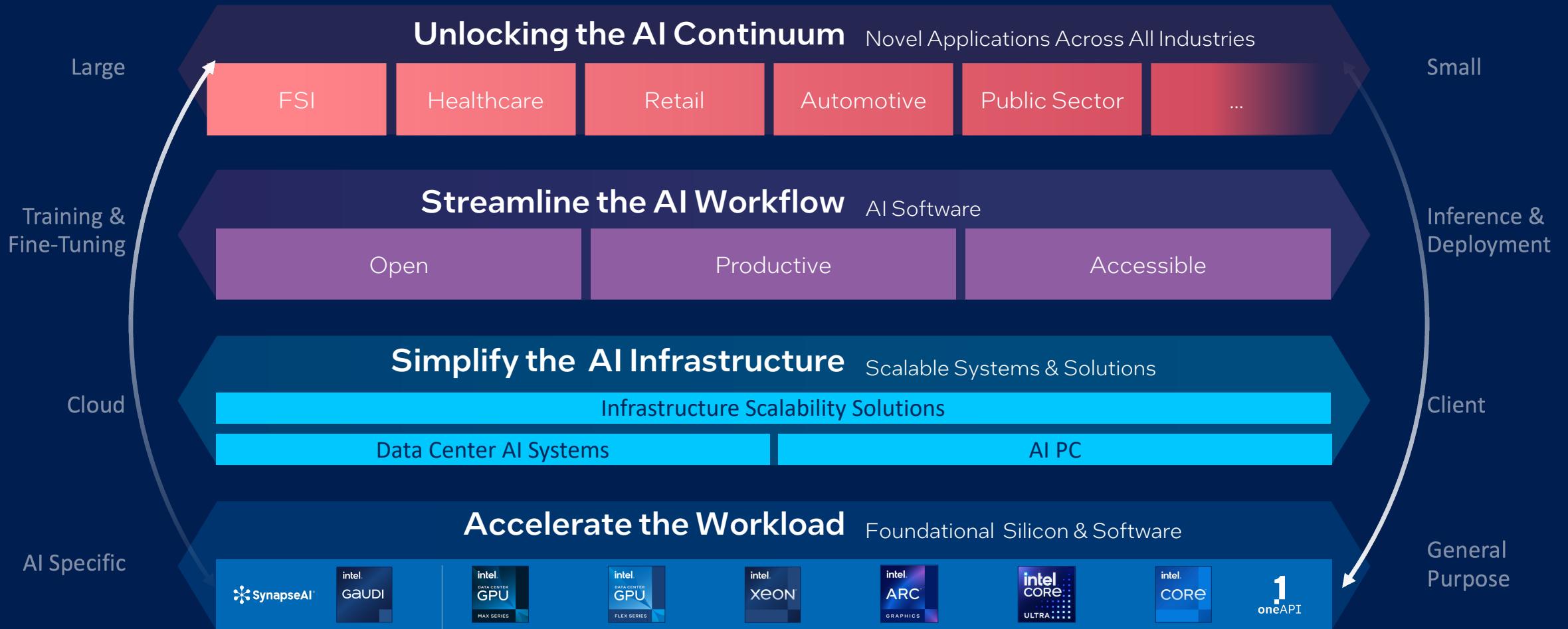


intel innovation

Bringing AI Everywhere



Intel's Approach



Bringing AI Everywhere

Large

Unlock the AI Continuum

Novel Applications

Small

Training &
Fine-Tuning

Streamline the AI Workflow

AI Software

Inference &
Deployment

Cloud

Simplify the AI Infrastructure

Scalable Systems & Solutions

Client

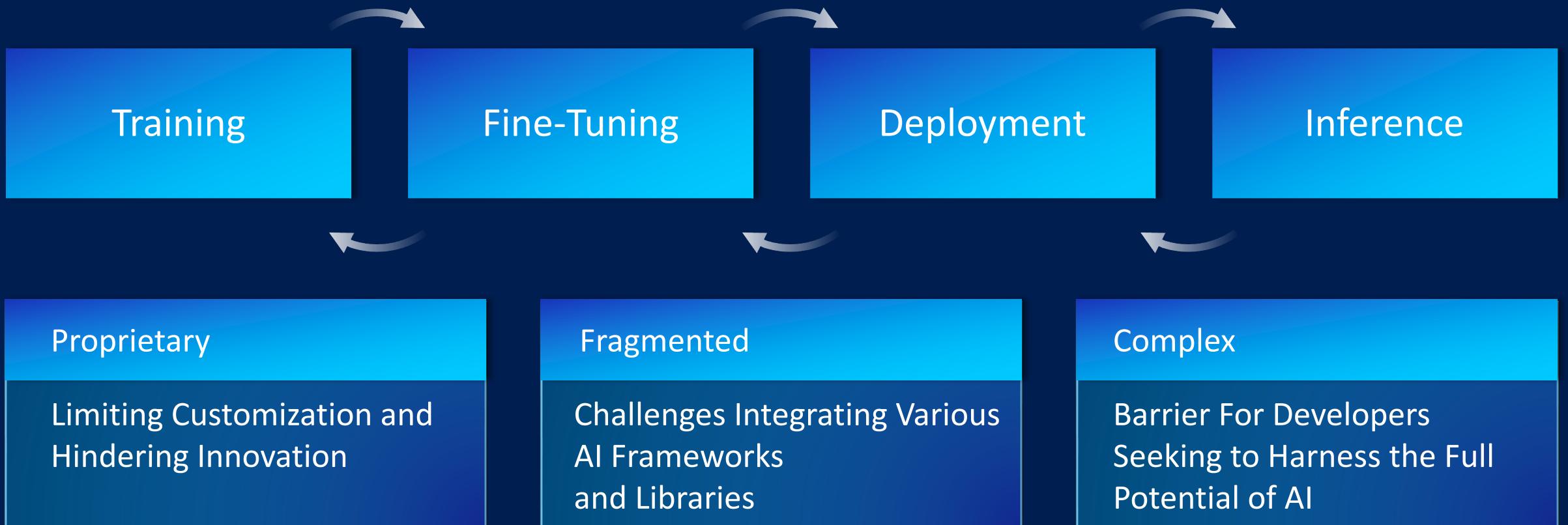
AI Specific

Accelerate the Workload

Foundational Silicon & Software

General
Purpose

AI Workflow Today



Streamline the AI Workflow

Training



SigOpt
AutoML



intel Neural Compressor



DirectML



cnvrg.io

Open



Productive

Solutions

Pre-configured
containers

AI tool
Selector

Tooling

Optimized
Extensions

OpenVINO

Modin

CNVRG

References

AI Reference
Kits

Hugging Face
Collaboration

Accessible

Ecosystem Engagement

Industry &
Academia

Solutions
Marketplace

High-Touch
Support

Developer Training

MLOPS training

Centers of
Excellence

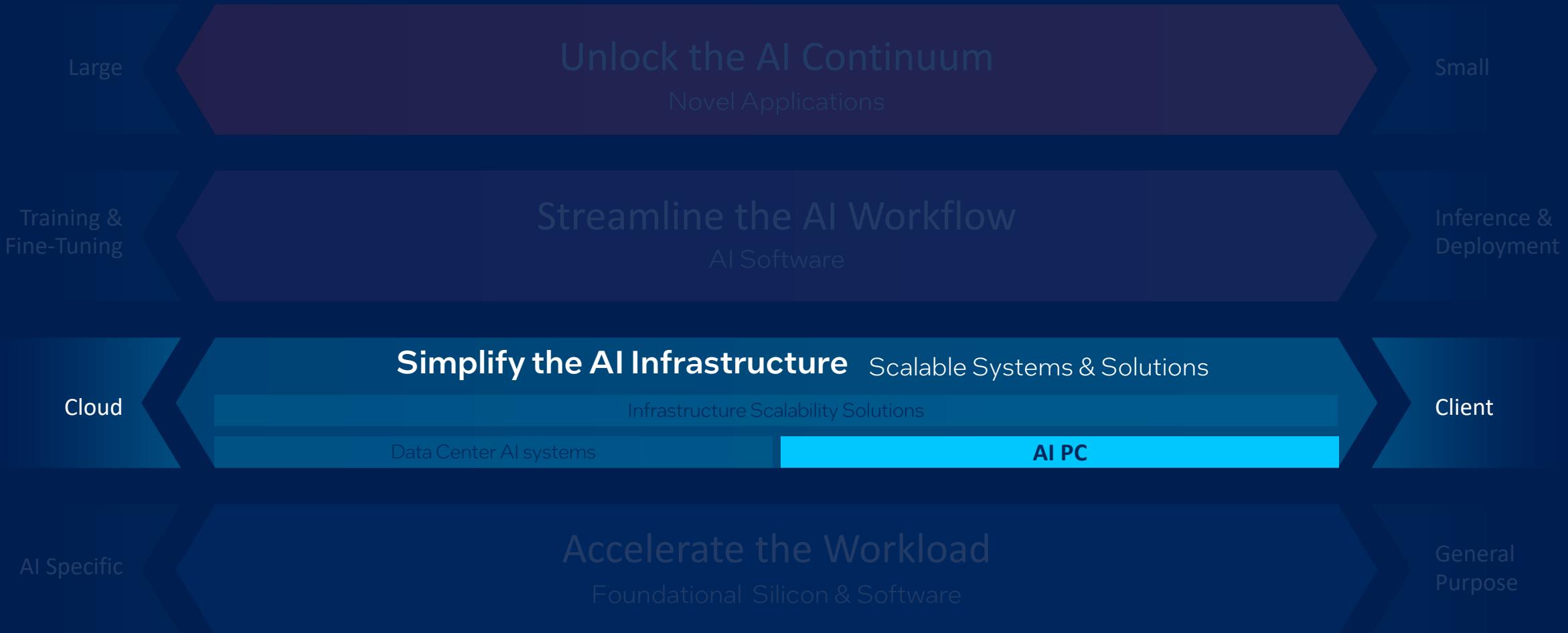
Documentation
& Tutorials

Training Videos

Summits &
Hackathons

Liftoff Program

Bringing AI Everywhere



Age of the AI PC

Enabling a World-class, AI-ready PC
with High Performance and Broad Compatibility



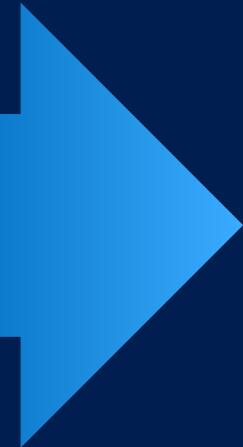
Enhanced
Audio Effects

Elevated
Video
Collaboration
& Streaming

Creator &
Gaming
Effects



Today – Enhancements



AI Assistants Know
Your Daily Context



More Creative,
Productive, &
Collaborative

Across Everything
You Do

Tomorrow – Everything

Bringing AI Everywhere

Large

Unlock the AI Continuum

Novel Applications

Small

Training &
Fine-Tuning

Streamline the AI Workflow

AI Software

Inference &
Deployment

Cloud

Simplify the AI Infrastructure

 Scalable Systems & Solutions

Infrastructure Scalability Solutions

Data Center AI Systems

AI PC

Client

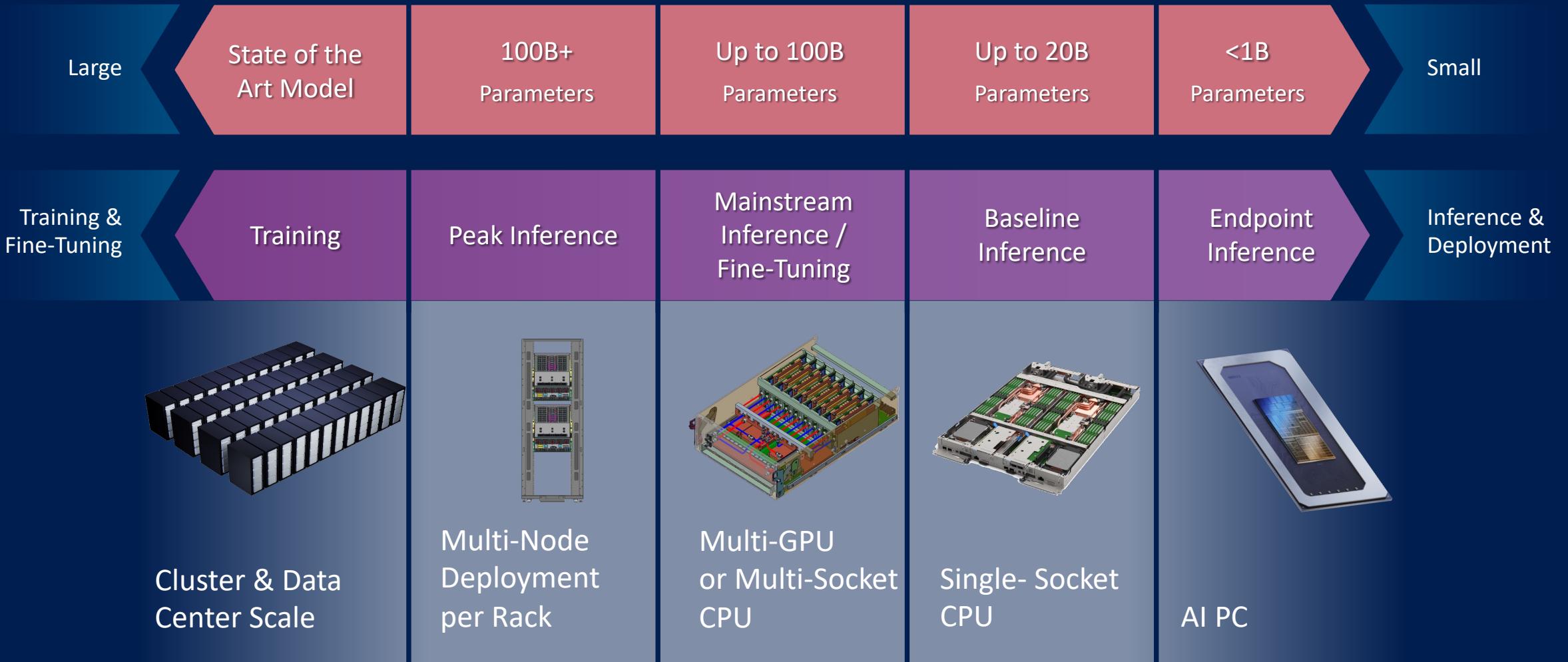
AI Specific

Accelerate the Workload

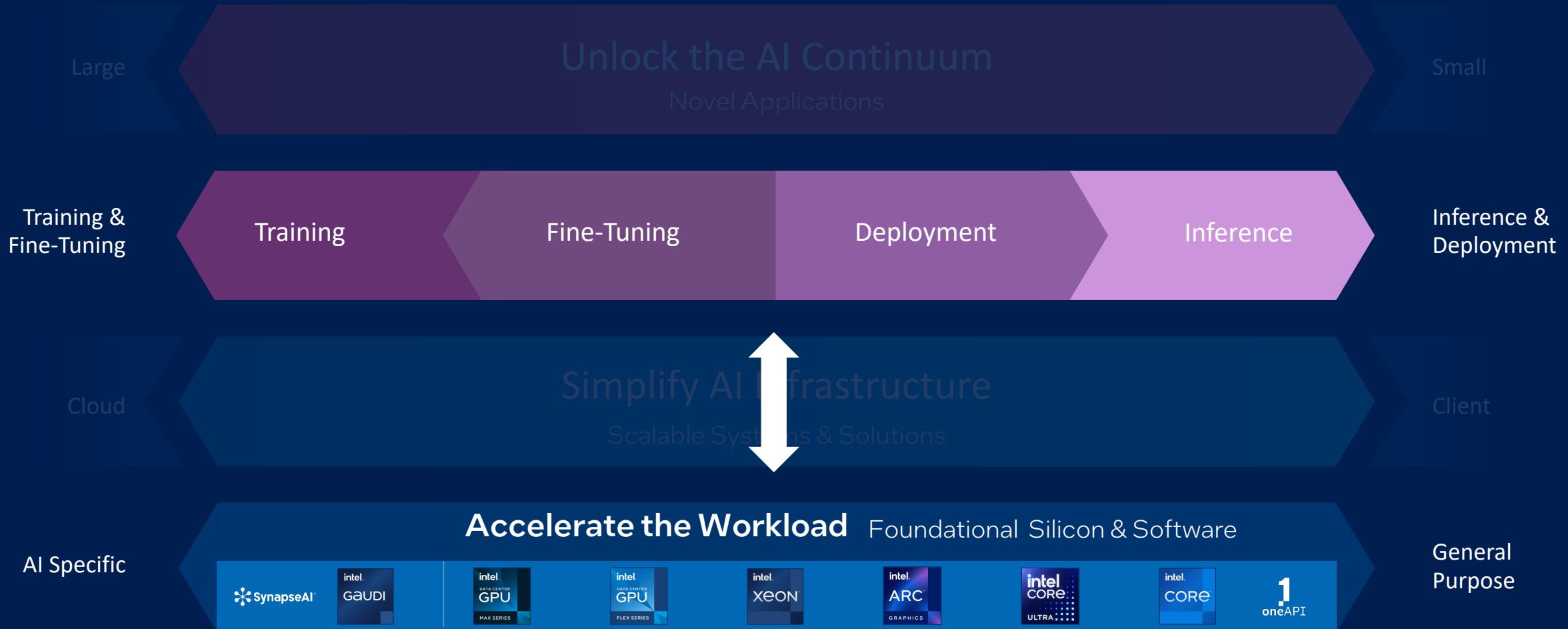
Foundational Silicon & Software

General
Purpose

Scalable Systems for Simple AI Infrastructures



Bringing AI Everywhere



Intel® Gaudi® 2 AI Accelerator



7nm
Process
Technology

24
Tensor
Processor Cores

96 GB
On-Board
HBM2

48 MB
SRAM

24
Integrated
Ethernet ports

Proven Performance

- The ONLY alternative to H100 for training LLMs based on MLPerf
- Trained GPT-3* model in 311 minutes on 384 Intel Gaudi2 accelerators

Price Performance

- Intel Gaudi2 accelerators with FP8 estimated to deliver price-performance >H100
- ~2x price-performance to A100

Scalability

- 95% linear scaling on MLPerf GPT-3 training benchmark
- Access large Intel Gaudi2 cluster on the Intel Developer Cloud

Ease of Use

- Software optimized for deep learning training and inference
- PyTorch, Hugging Face, Optimum Library optimizations

Seamless Code Transitioning

Performant AI Code with Minimal Changes



deepspeed



...

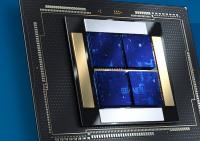
Across Generations & Architectures



Intel® Gaudi® 2
AI Accelerator



Intel® Gaudi® 3 AI
Accelerator



Next Gen GPU (Codename
Falcon Shores)

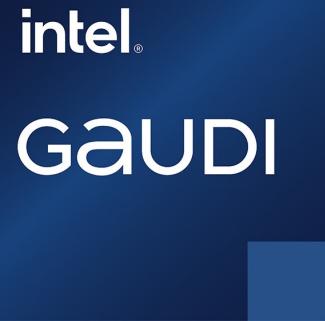


Gaudi Software Suite

Powered by Transitioning Into a
Single Software Environment



Unified Programming Model



intel.
GAUDI



Fine Tuning



intel.
xeON®

Fine-tune with
Intel® Gaudi® 2 Processor
When Optimal Speed is Desired

Fine-tune On Intel® Xeon®,
Exploiting Its Industry-leading
Ubiquity In the Data Center

Intel Provides Solution Options for
Fine-tuning Gen AI and LLMs
to Fit Workload Needs

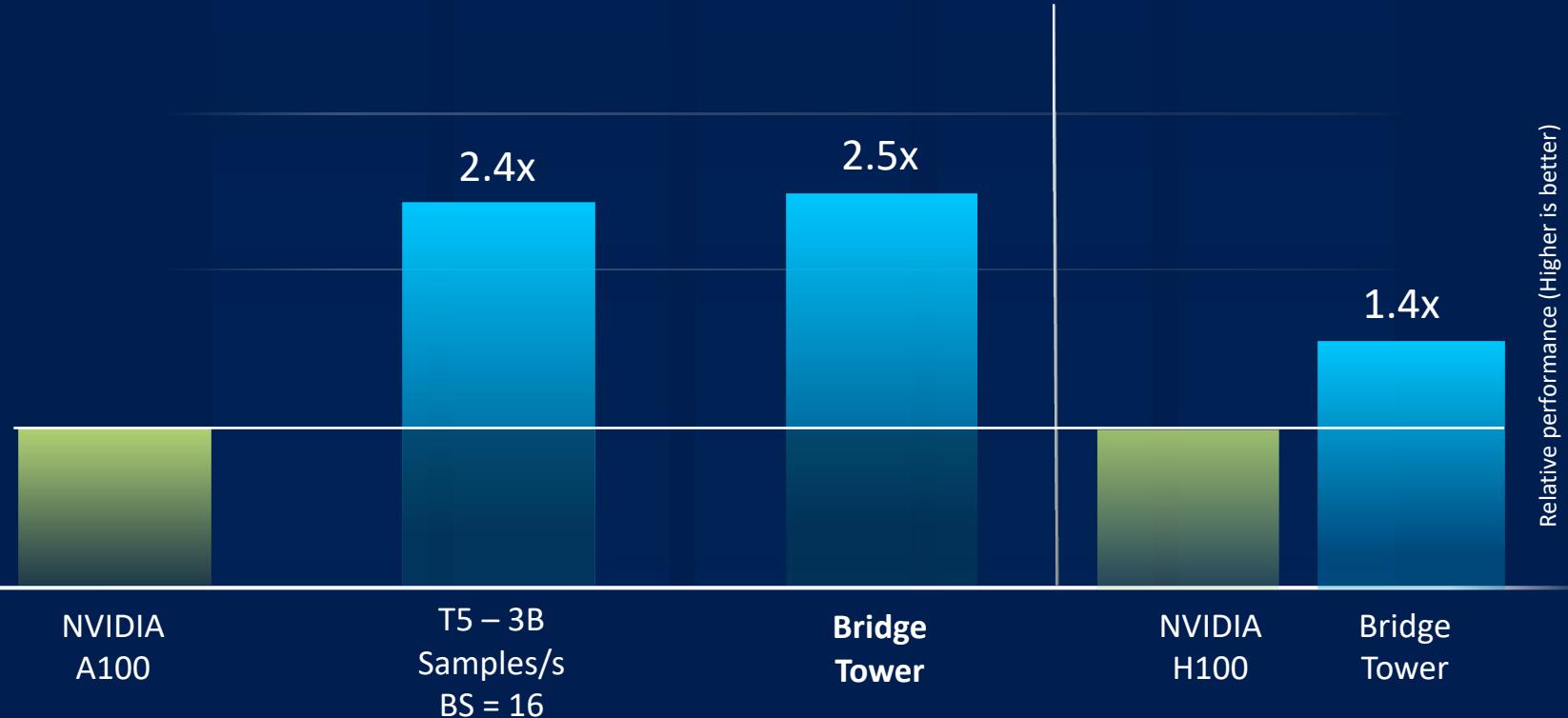
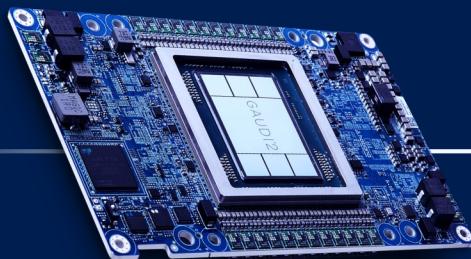
intel.

GAUDI



Hugging Face Evaluations Substantiate Intel® Gaudi® 2 Accelerator LLM Performance vs. Nvidia A100 and H100

Fine-tuning Across Numerous LLMs



intel.innovation

Visit <https://habana.ai/habana-claims-validation> for workloads and configurations. Results may vary.
<https://huggingface.co/blog/habana-gaudi-2-benchmark>
<https://huggingface.co/blog/bridgetower>



Out-of-the Box Intel® Xeon® Fine Tuning

Optimized Models & Spaces

Dolly

LLAMA2

MPT

LDM3D

Whisper

100k's
Mode

Intel Optimized Hugging Face Libraries & Tools

Transformers

Fine Tuning

Diffusers

Use Cases

Accelerate

Fine Tuning at Scale

PEFT

Efficient Fine
Tuning

Optimum

Performance
Optimization

Foundational Stack

PyTorch + IPEX

deepspeed + IDEX

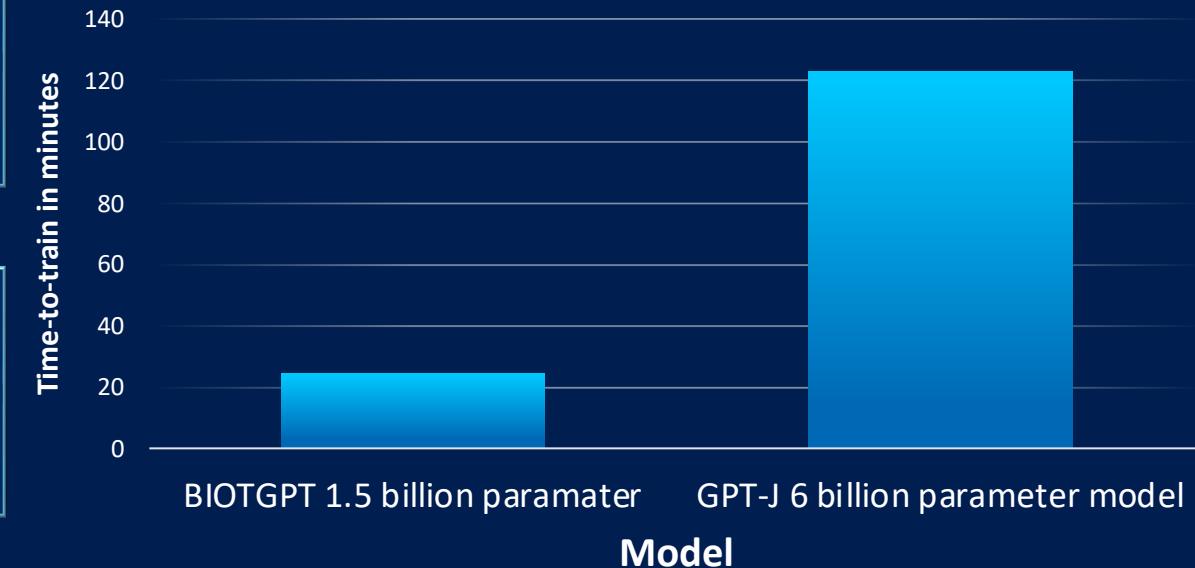
ITREX

ONNX
RUNTIME

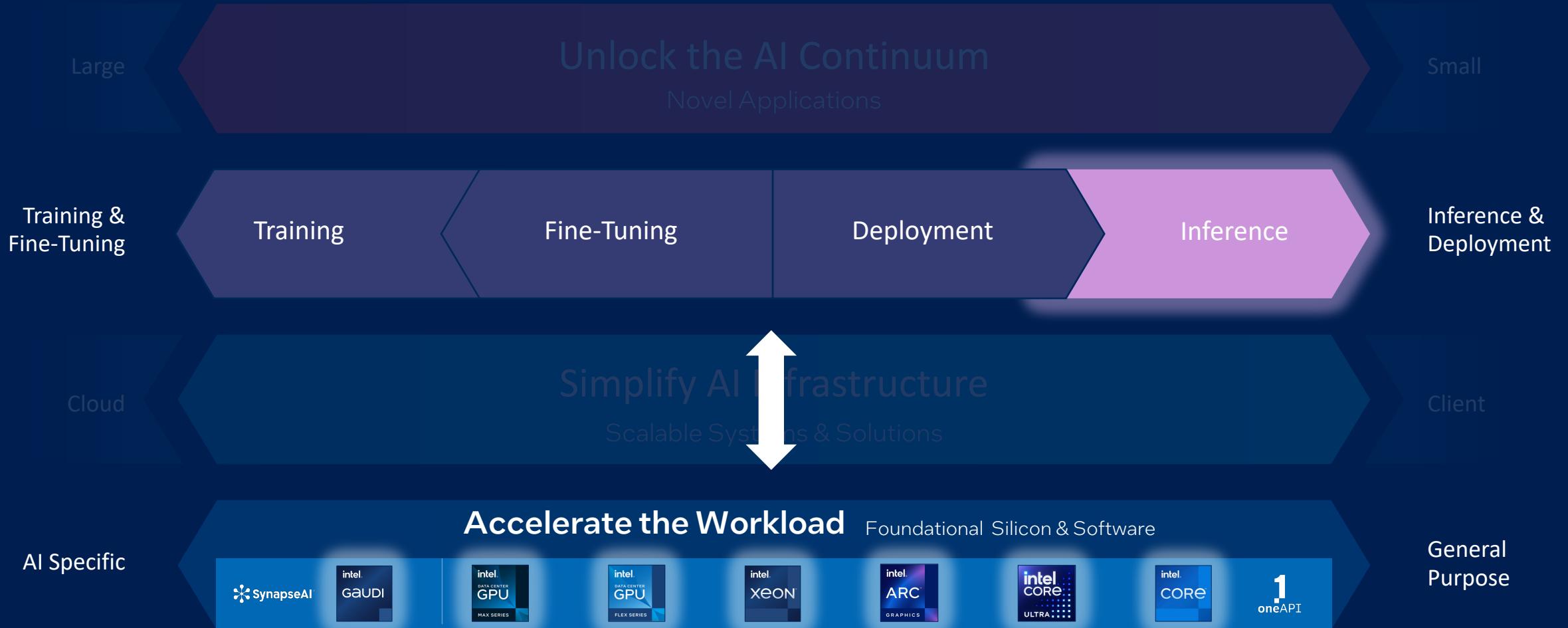
OpenVINO®

Multi-node Fine Tuning Open-source
Commercial Large Foundational Models In
Minutes To Hours

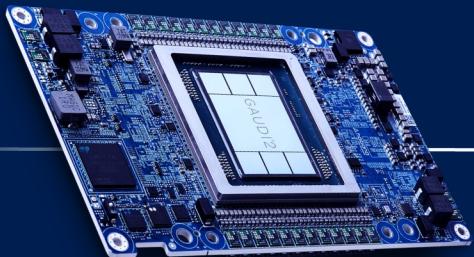
BIOGPT 1.5 Billion Parameter and
GPT-J 6 Billion Parameter Model



Inference



Inference Advantage Across Multiple LLM Performance Metrics



NVIDIA
A100

2.84x

Stable Diffusion
Latency
BS =8; fp32

1.42x

BLOOMz176B
Inference
BS=1, BF16

2.89x

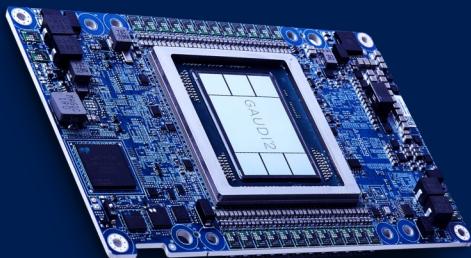
BLOOM 7B
BS=1, BF16

Relative performance (Higher is better)

Energy Efficiency

Throughput-per-Watt on BLOOMZ 176B Inference is 1.79x
better than H100; 1.61x better than A100

Intel® Gaudi® 2 AI Accelerator: Solving LLM Challenges



Inference on GPT-J

Intel Gaudi 2 Accelerator with FP8

- Near-parity* on GPT-J with H100
- Outperformed A100 by 2.4x (Server) and 2x (Offline)
- Achieved 99.9% accuracy with FP8

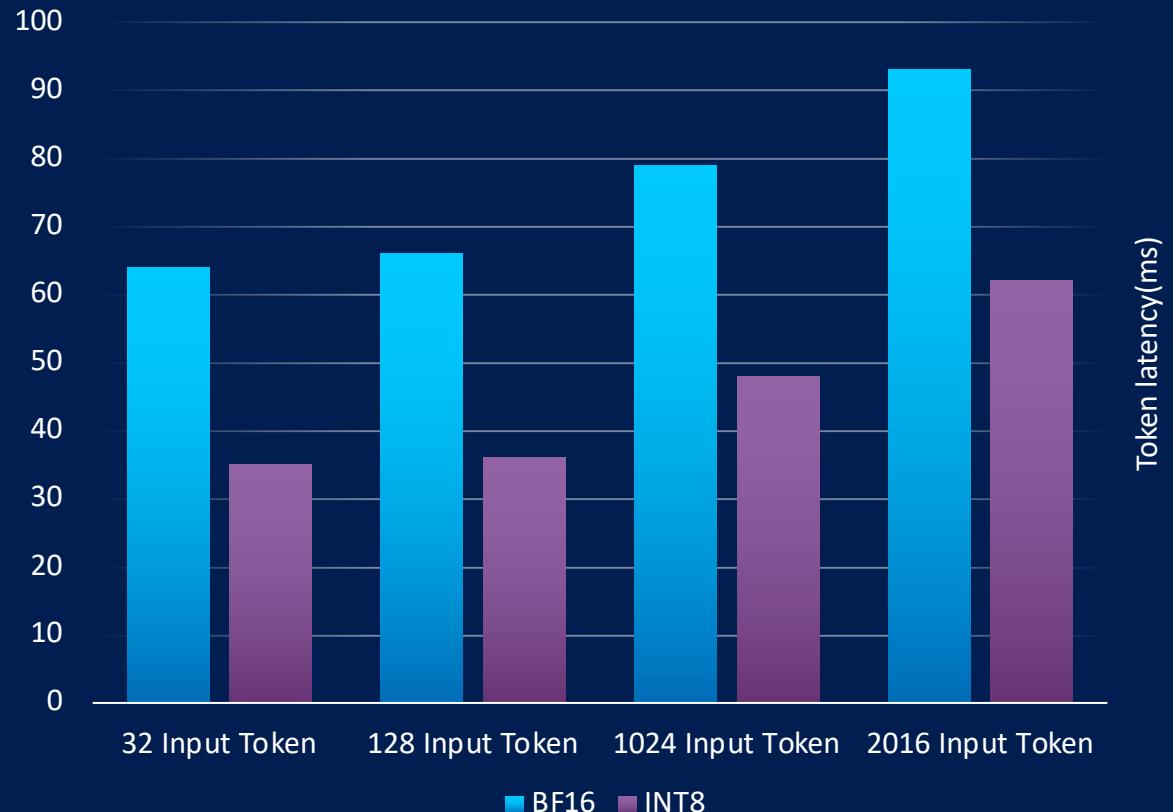
GPT-J On MLPerf Inference Benchmark



LLaMA2 (7B) Inference with 4th Gen Intel® Xeon® Processors

- Use any popular industry standard AI libraries
- Intel AI Platform validated with over 300 inference models
- One socket of 4th Gen Intel® Xeon® processors can run LLaMa2 chatbots in under 100ms 2nd token latency

LLaMA2 7B : Intel Xeon 4th Gen 8480 1S P90 Latency
Batch Size 1, Beam Width 4, PyTorch + IPEX
(Lower is better)



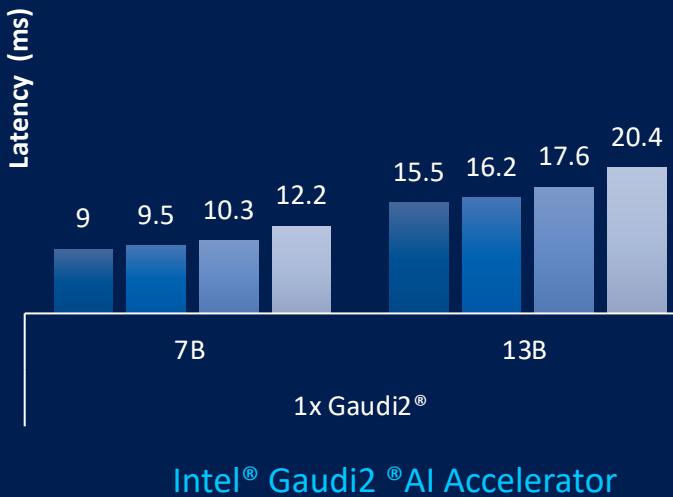
Results shown for bare metal.

Inference Across Multiple Products

LLaMA2 7B & 13B Inference

On 1 Socket Intel® Xeon®
Scalable Processor

Greedy Mode, Mixed Precision
(bfloating16), BS = 1, 256 Output Tokens



intel innovation

Llama 2 Next Token Latency (Lower is Better)

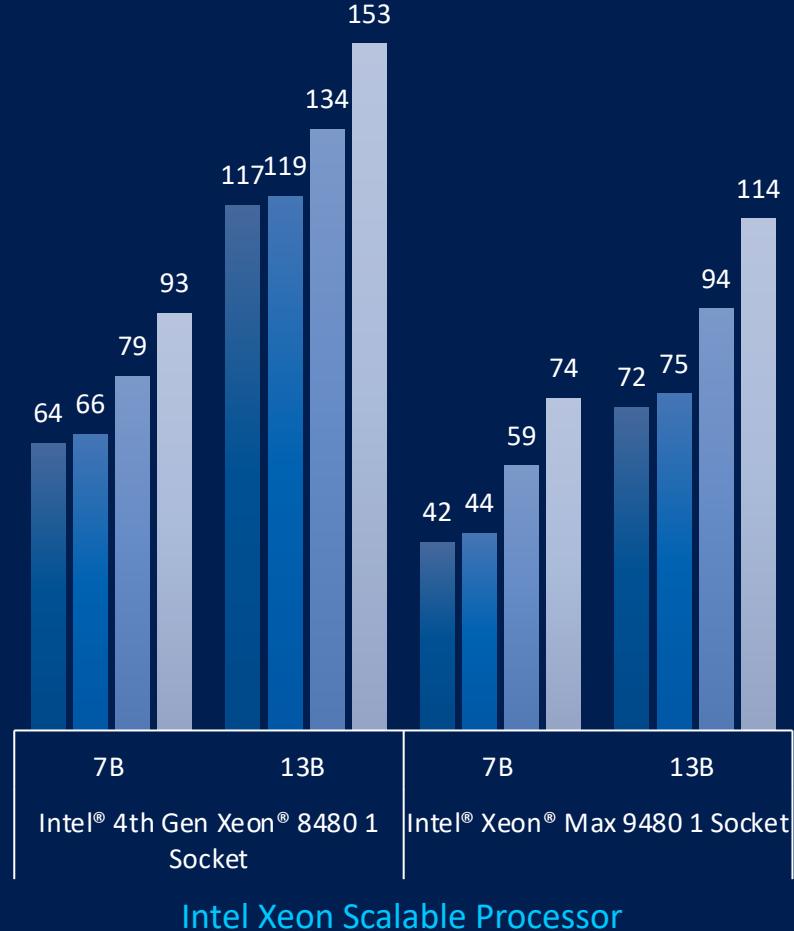
On 1 Tile (out of 2 tiles per card) Intel®
Data Center GPU Max 1550



7B
13B
Intel® Data Center GPU Max 1550 single tile (1 card
has 2 tiles)



Intel Data Center GPU Max 1550



7B
13B
Intel® 4th Gen Xeon® 8480 1
Socket



7B
13B
Intel Xeon Scalable Processor

Bringing AI Everywhere

Large Scale
AI Models

Unlock the AI Continuum
Novel Applications

Small AI
Models

Training &
Fine-Tuning

Training

Fine-Tuning

Deployment

Inference

Inference &
Deployment

Cloud

Simplify AI Infrastructure

Scalable Systems & Solutions

Client

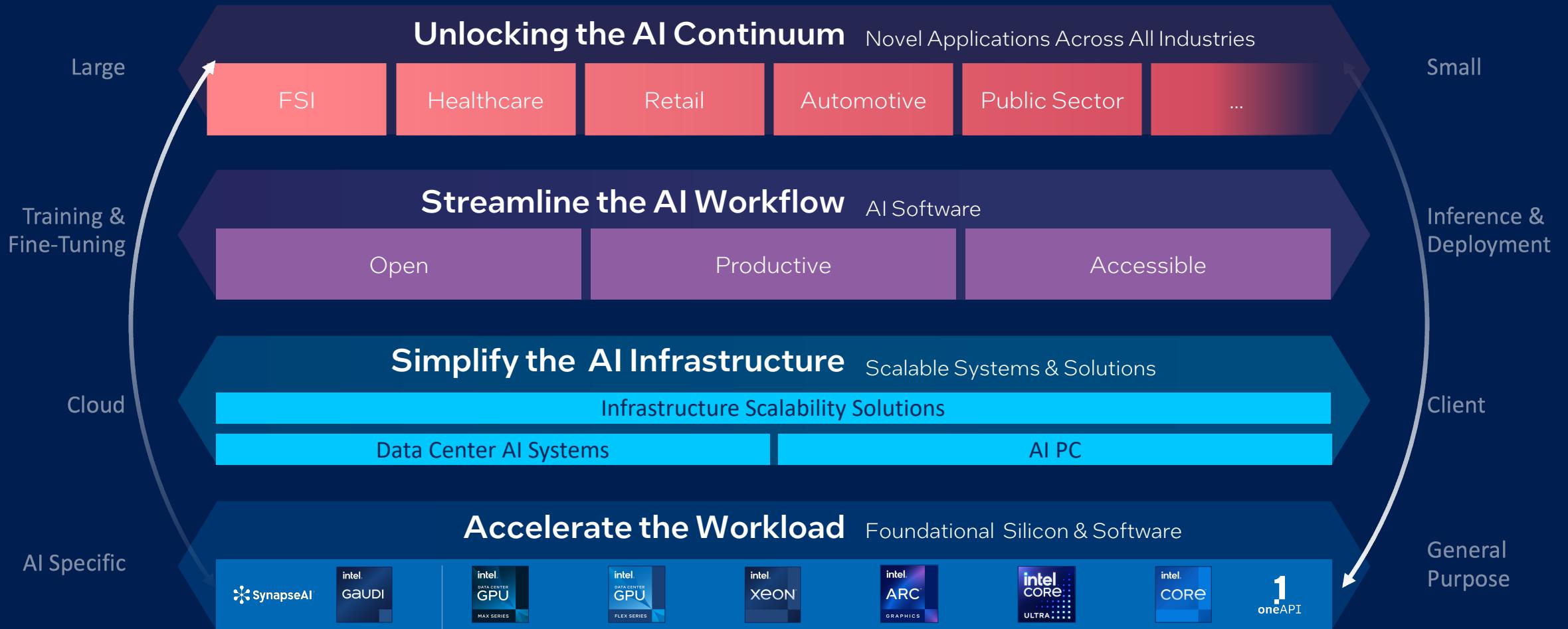
AI Specific

Accelerate the Workload Foundational Silicon



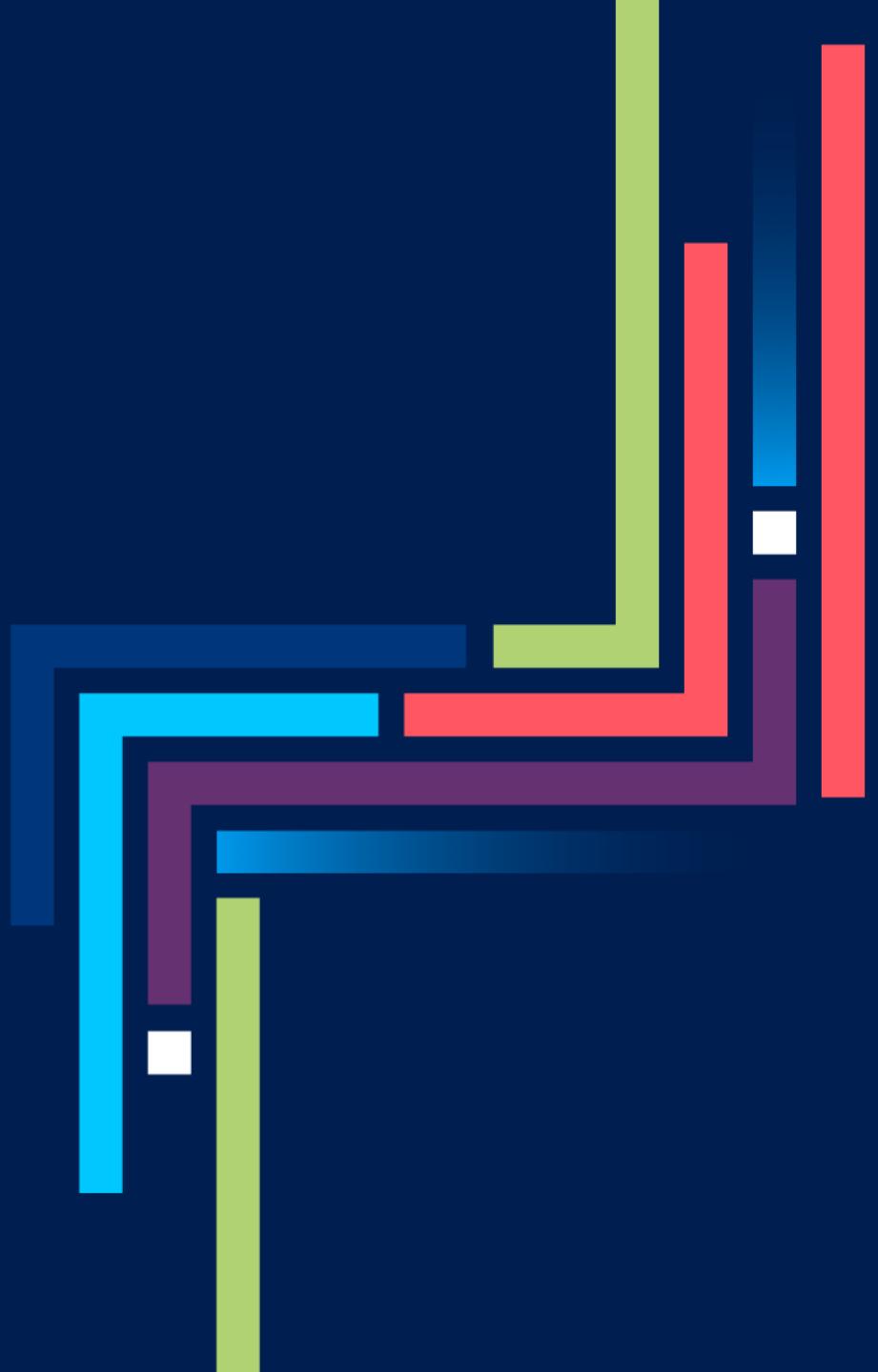
General
Purpose

Intel's Approach



intel.[®]
Innovation

Thank You!



Notices and Disclaimers

For notices, disclaimers, and details about performance claims, visit www.intel.com/PerformanceIndex 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.

