

intel[®] ai
summit
英特爾 AI 科技論壇

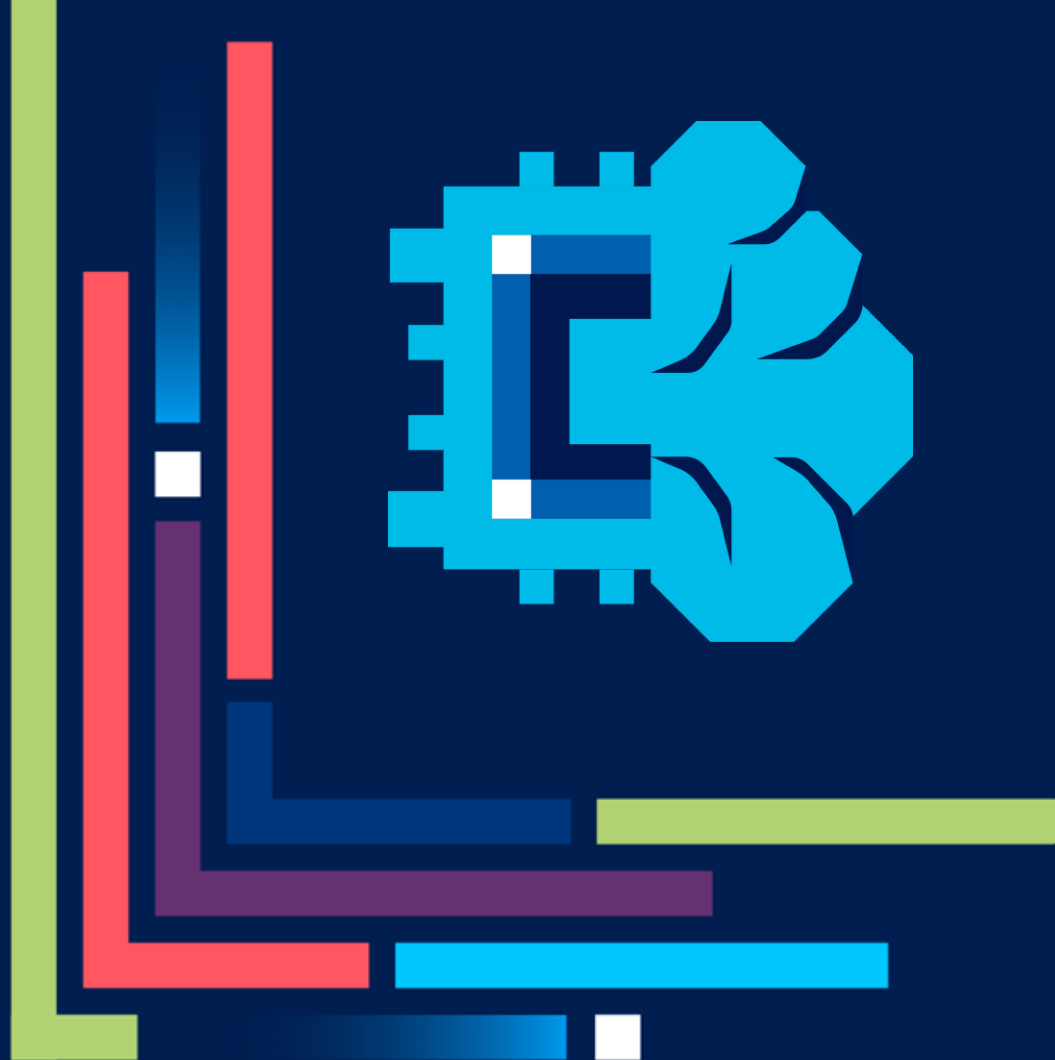
Bringing AI Everywhere

Accelerating GenAI for Enterprise

Sean Kuo

AI Center of Excellence APJ Sales Lead

March 27th, 2024



Parameters in notable artificial intelligence systems

Parameters are variables in an AI system whose values are adjusted during training to establish how input data gets transformed into the desired output; for example, the connection weights in an artificial neural network.

[Linear](#)
[Log](#)
[Select systems](#)

Task domain

- Drawing
- Games
- Language
- Multimodal
- Other
- Recommendation
- Robotics
- Search
- Speech
- Vision

Number of parameters

Publication date

Source: Epoch (2023)
 Note: Parameters are estimated based on published results in the AI literature and some with some uncertainty. The authors express the estimates to be correct within a factor of 10.

OurWorlds-Data.com/artificial-intelligence • CC BY

[CHART](#)
[TABLE](#)
[SOURCES](#)
[DOWNLOAD](#)

Jul 2, 1950
 Jun 20, 2023



Amazing Innovation

But Huge Models are Inaccessible For Most

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

Amazing Innovation

But Huge Models are Inaccessible For Most

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

Specialized AI Models

The answer for the "masses"



Large Foundational Model

Advantages

- + Incredible all-in-one, out-of-the-box versatility: text, programming, continual natural language conversation and plain summarization
- + Surprisingly, compelling outcomes

Challenges

- Big (>100B parameters), expensive- \$4m+ to train, \$3m per month for inferencing
- Hallucinations; lack of explainability, intellectual property issues
- Frozen in time (sampling)



Domain Specific Models

Advantages

- + 10-100x smaller models while maintaining/improving accuracy
- + Economical on general-purpose compute
- + Correctness; Source attribution; Explainability
- + Utilizing private/enterprise data
- + Continuously updated information

Challenges

- Reduced range of tasks
- Requires few-shot fine-tuning and indexing

Specialized AI Models

The answer for the "masses"



Large Foundational Model

Advantages

- + Incredible all-in-one, out-of-the-box versatility: text, programming, continual natural language conversation and plain summarization
- + Surprisingly, compelling outcomes

Challenges

- Big (>100B parameters), expensive- \$4m+ to train, \$3m per month for inferencing
- Hallucinations; lack of explainability, intellectual property issues
- Frozen in time (sampling)



Domain Specific Models

Advantages

- + 10-100x smaller models while maintaining/improving accuracy
- + Economical on general-purpose compute
- + Correctness; Source attribution; Explainability
- + Utilizing private/enterprise data
- + Continuously updated information

Challenges

- Reduced range of tasks
- Requires few-shot fine-tuning and indexing

Specialized Models Enable Scale



Education

Teacher Assistant
Student Study Buddy
Parent Chat Portal



Health

Drug Discovery
Doctor Assistant
Patient Family chatbot



Finance

Algorithmic Trading
Customer Portfolio Assistant
Risk / Credit Assessment



Retail

Product Promotion
Customer Interface and Sentiment Tool
Image Shopping Aid



Government

Gov Services Assistant
Document Search Summarization
Live Language Translation



Energy

Energy Consumption Forecasting
Operational Performance
Energy Trading Assistant



Automotive

Autonomous Car Development
Multi-language in car aid
Supply Chain Optimization



Manufacturing

Factory Automation
Predictive Maintenance
Precision Agriculture



Telco

Personalized Customer Services
Network Automation
Operational Performance

Enterprise options to build specialized GenAI

1.



Large Foundation Models



Fine tuning



Retrieval Augmented Generation

Domain Specific Models



Finance ChatBot



Personalized Marketing Content



Code Generator



Clinical Notes Transcription



Text-to-Speech Generator

2.



General Open Models



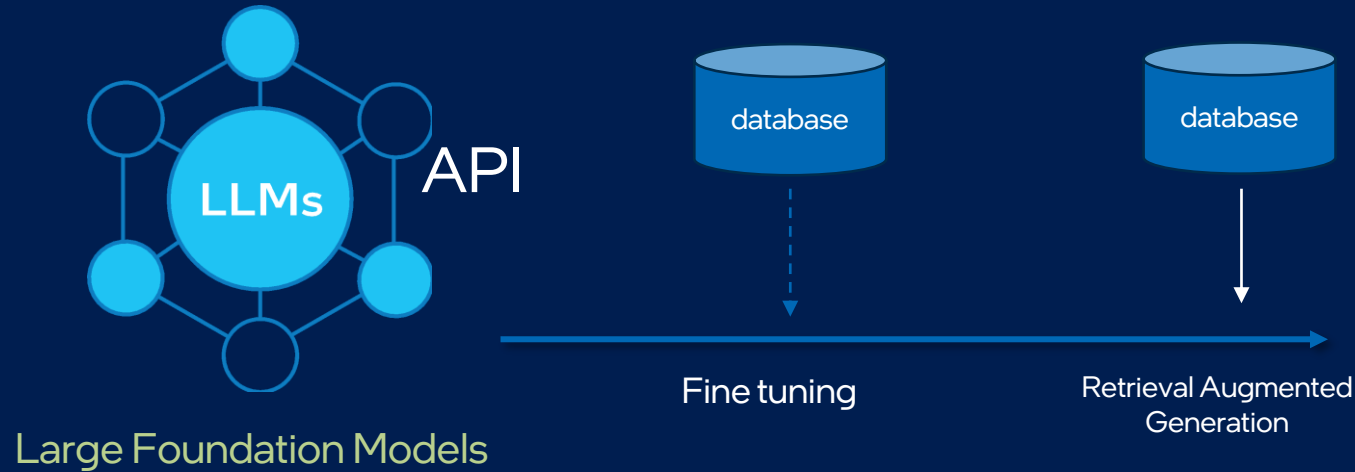
Fine tuning



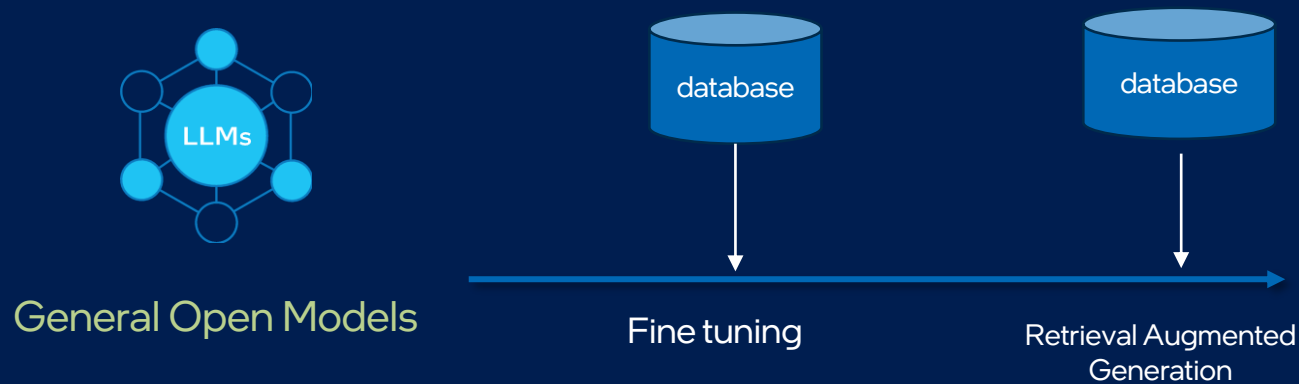
Retrieval Augmented Generation

Enterprise options to build specialized GenAI

1.



2.



Domain Specific Models



Finance ChatBot



Personalized Marketing Content



Code Generator



Clinical Notes Transcription



Text-to-Speech Generator

Intel has the software tools developers use to scale AI Everywhere

Upstream

Integrated acceleration to popular open-source software

PyTorch, TensorFlow, ONNX RT, more ...

Intel Extension

Easily pluggable extensions to open-source software

Intel Extension for PyTorch,
Intel Extension for TensorFlow, more ...

Intel Tools

Tools / Kits which improve productivity and performance on Intel hardware

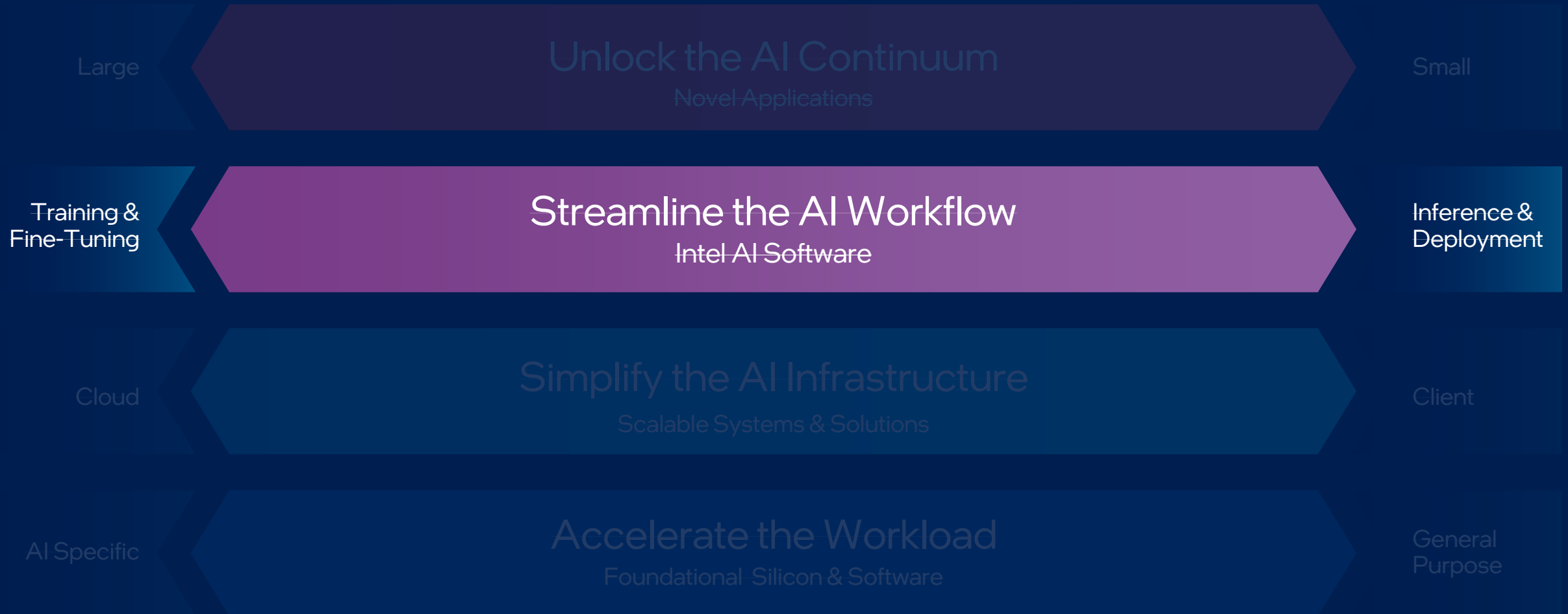
OpenVINO™, oneContainer Portal,
Intel Neural Compressor,
SynapseAI®, Developer Cloud

Across major software channels (PyPI, Anaconda, Intel, Apt, Yum, Docker) and ecosystems, (Optimum Intel through Hugging Face)

Bringing AI Everywhere



Bringing AI Everywhere



Streamline the AI Workflow

Training

Fine-Tuning

Deployment

Inference



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

Streamline the AI Workflow



Open

Logos for open-source and cross-platform technologies:

- julia
- SYCL
- Java
- OpenCL
- OpenXLA
- OpenMP
- C, C++
- oneAPI

Productive

Solutions	Pre-configured containers	AI tool Selector
Tooling	Optimized Extensions	OpenVINO
	Modin	CNVRG
References	AI Reference Kits	Hugging Face Collaboration

Accessible

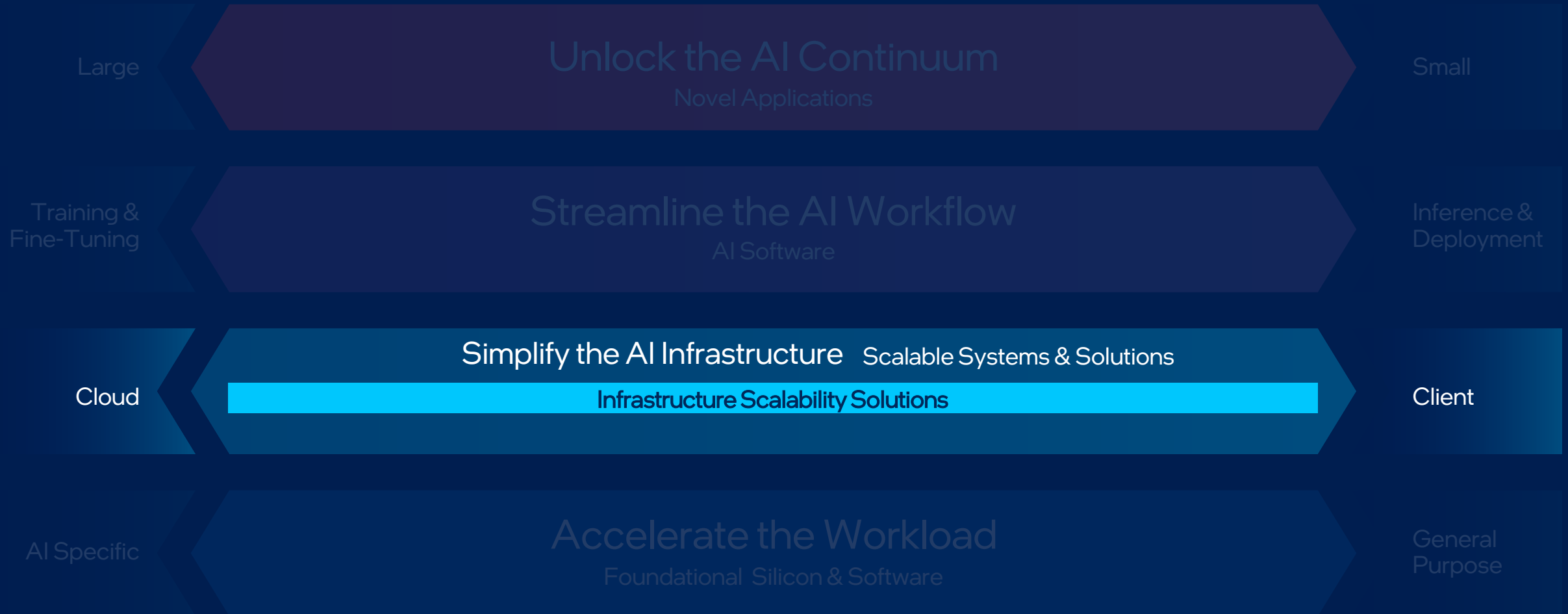
Ecosystem Engagement

- Industry & Academia
- Solutions Marketplace
- High-Touch Support

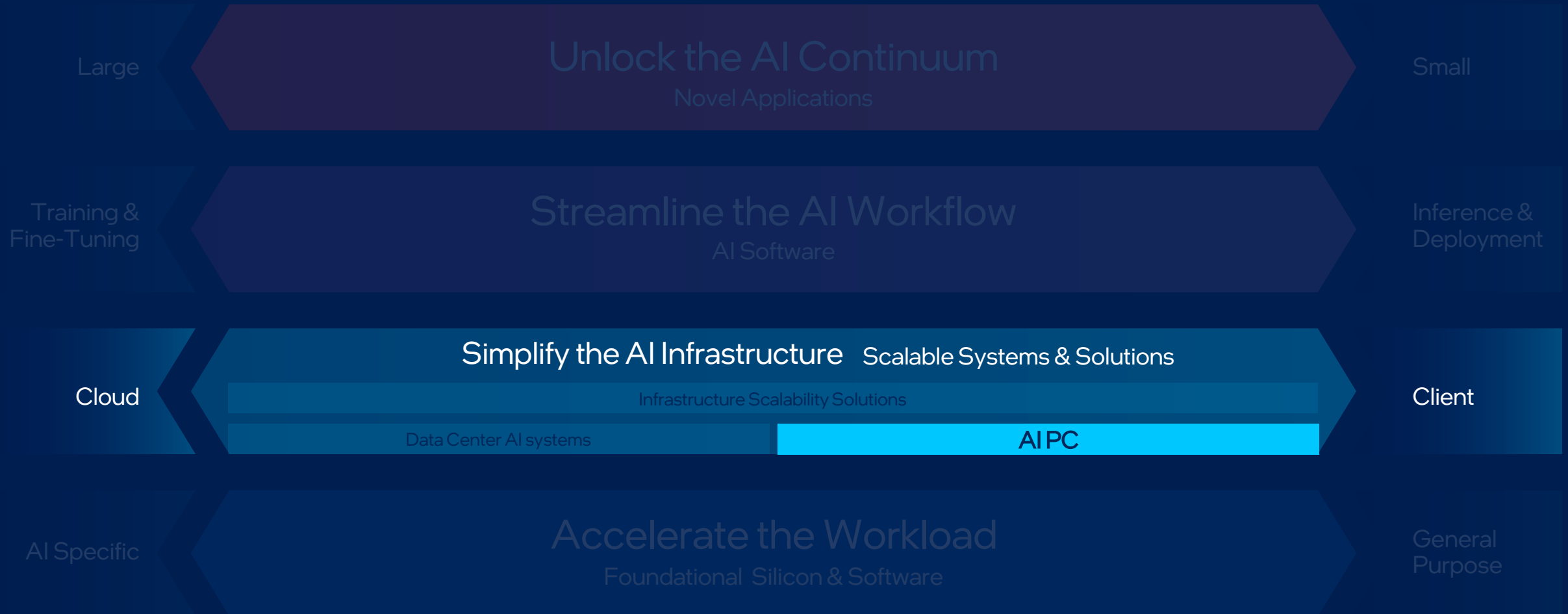
Developer Training

- MLOPS training
- Centers of Excellence
- Documentation & Tutorials
- Training Videos
- Summits & Hackathons
- Liftoff Program

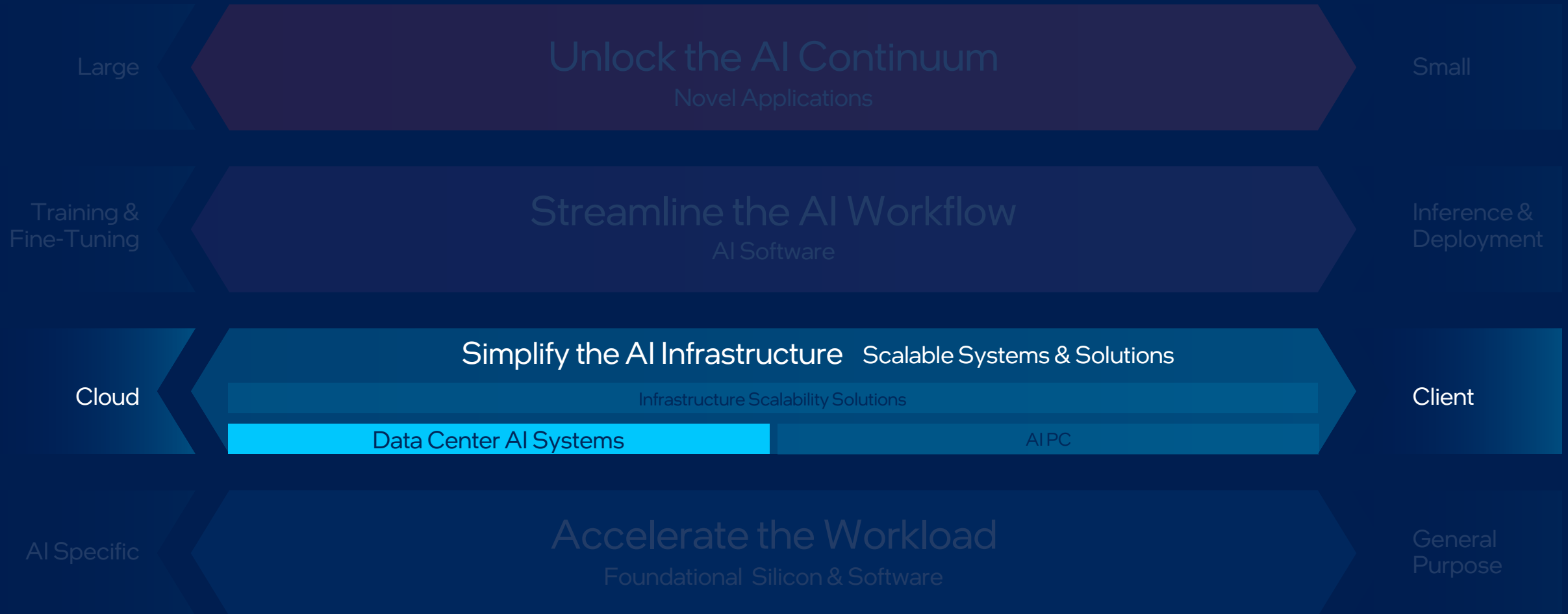
Bringing AI Everywhere



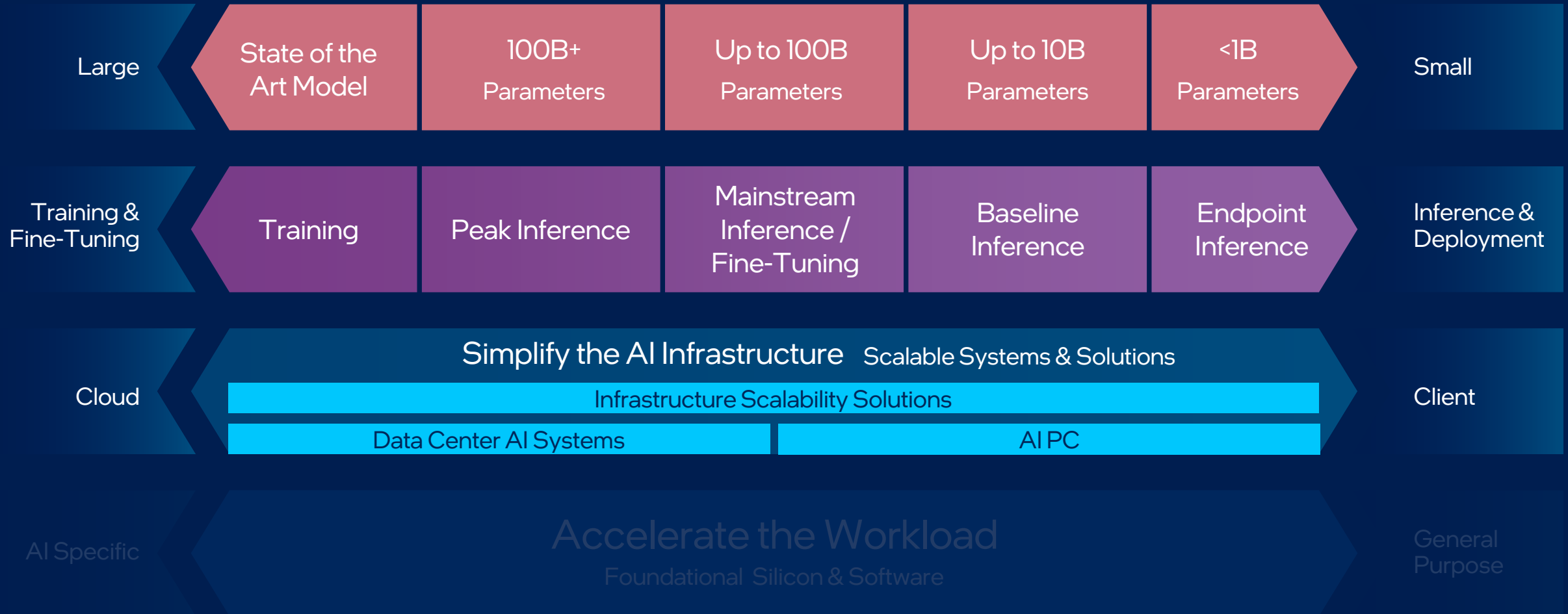
Bringing AI Everywhere



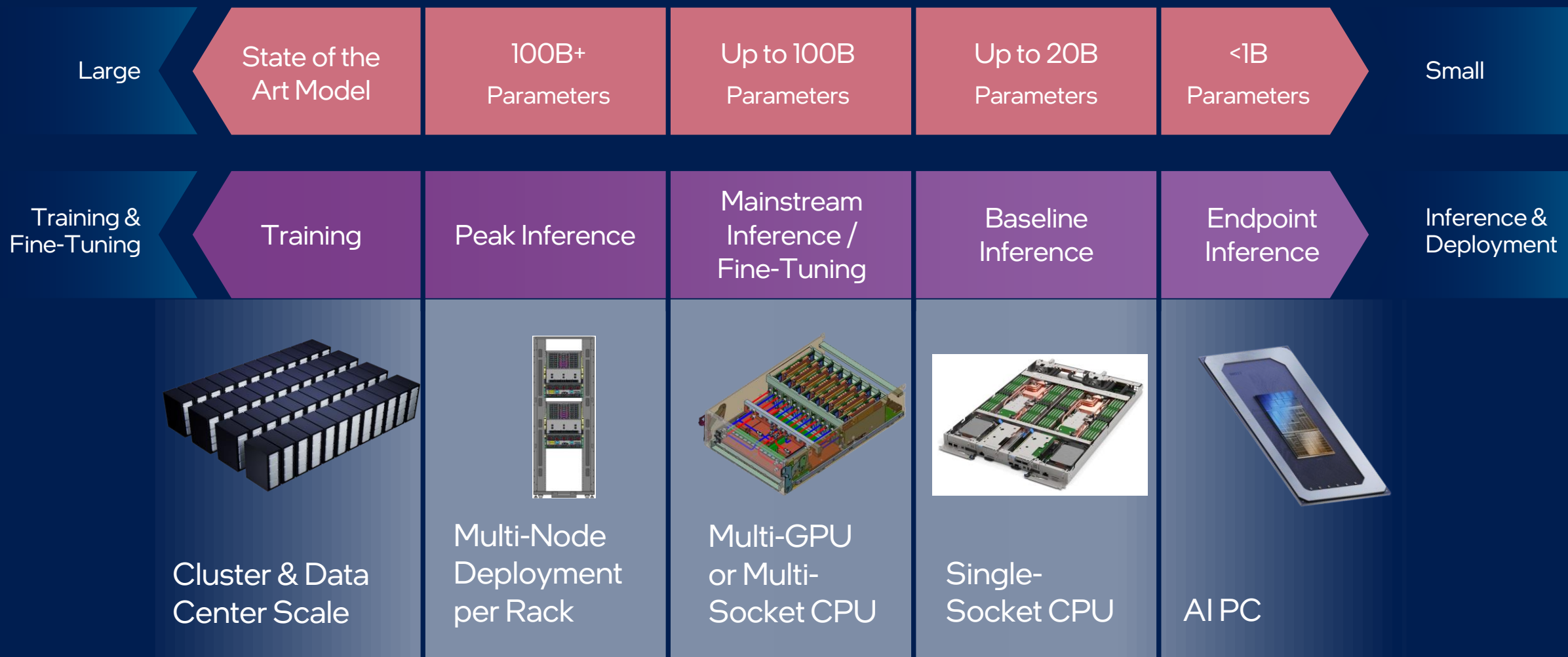
Bringing AI Everywhere



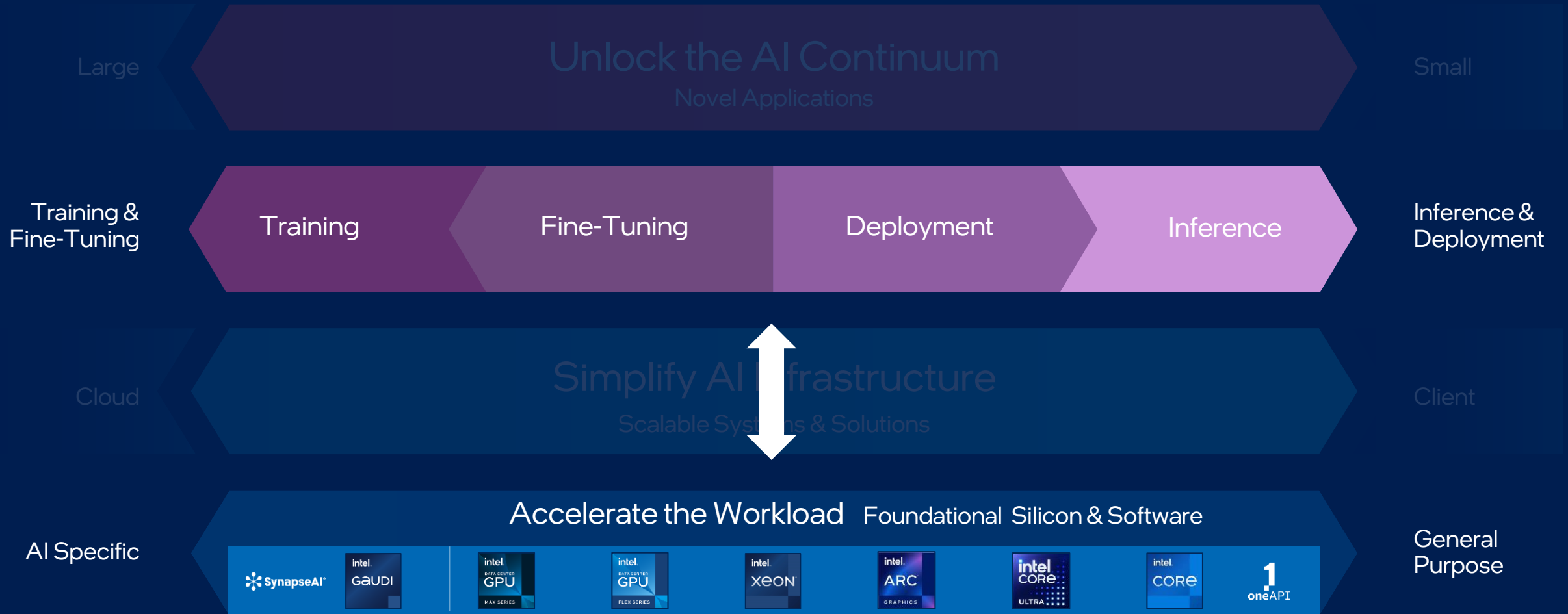
Scalable AI Workloads From Cloud to Client



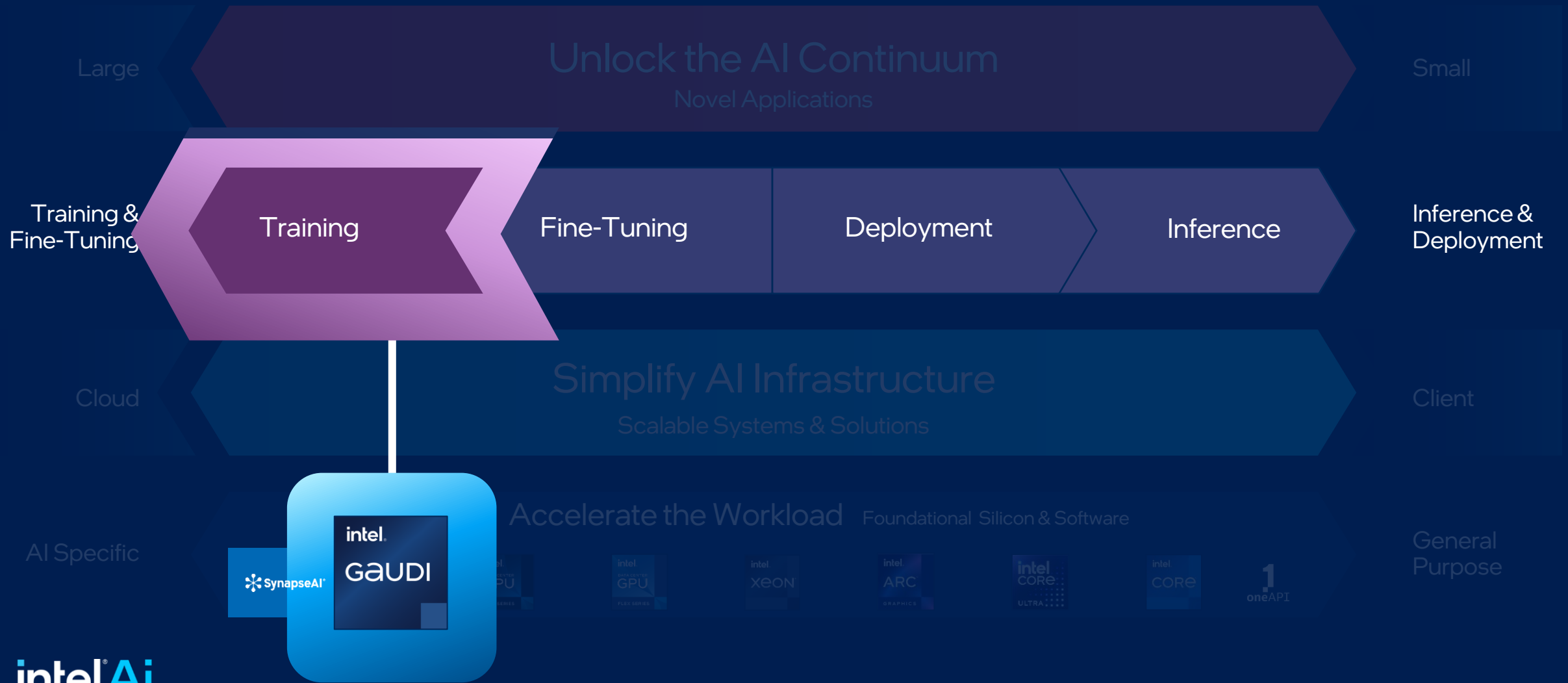
Scalable Systems for Simple AI Infrastructures



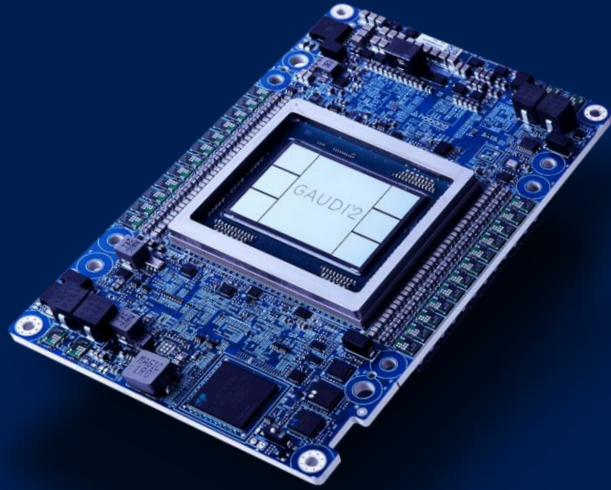
Bringing AI Everywhere



Bringing AI Everywhere



Intel® Gaudi® 2 AI Accelerator



Proven Performance

- The ONLY alternative to H100 for training LLMs based on MLPerf
- Trained GPT-3* model TTT doubled in 2023 from 311 minutes in Jun to 153 min in Nov'23 on 384 Intel Gaudi2 accelerators

Price Performance

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

7nm

Process Technology

24

Tensor Processor Cores

96 GB

On-Board HBM2

Scalability

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

48 MB

SRAM

24

Integrated Ethernet ports

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

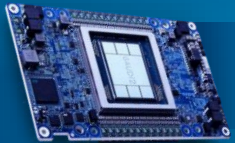
 TensorFlow

 PyTorch

 deepspeed



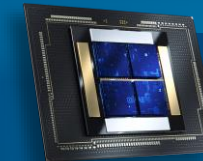
Across Generations & Architectures



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

1
oneAPI

Unified Programming Model

Stability.AI

Performance Evaluation of Gaudi

- Multimodal Diffusion Transformer model (MMDiT) on Stable Diffusion 3 (in preview now)
- MMDiT: 50% faster training than H100-80GB & 3x faster than A100-80GB
- Stable Beluga 2.5, 70B language model-- 28% faster inference speed vs. the A100.
- Better Performance, TCO and took less than one day to port the code base

stability.ai



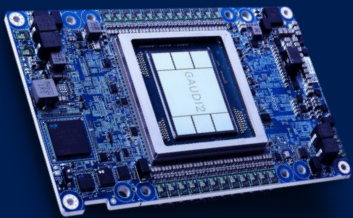
intel

intel[®]
Ai
summit

See blog for workloads and configurations. Results may vary
<https://stability.ai/news/putting-the-ai-supercomputer-to-work>



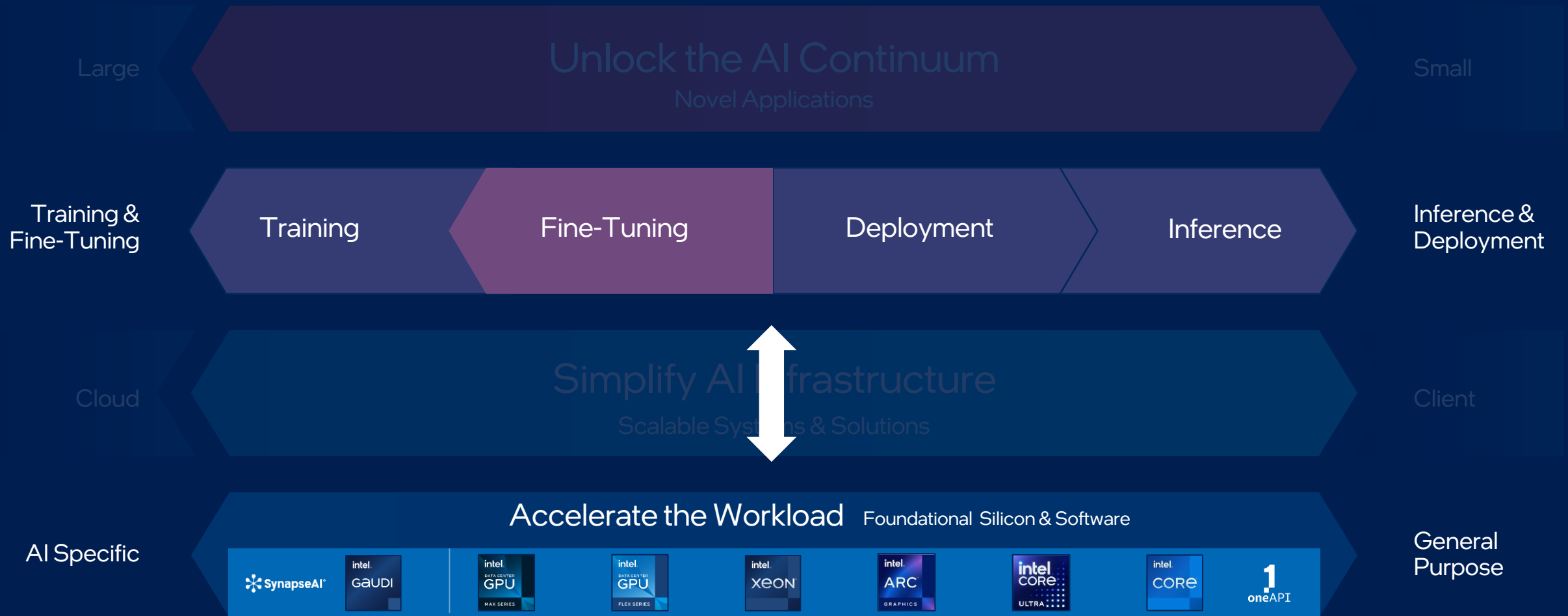
Intel® Gaudi® 2 AI Accelerator Up To 55% Faster Than NVIDIA H100 In Stable Diffusion, 3X Faster Than A100 In AI Benchmark Showdown



Device	Attention	# Nodes	# Accelerators (total)	Batch Size per Accelerator	Total Batch Size	Images / sec (100-MA)
Gaudi2	FusedSDPA	2	16	32	512	1,254
Gaudi2	FusedSDPA	2	16	16	256	927
H100-80GB	xFormers	2	16	16	256	595
A100-80GB	xFormers	2	16	16	256	381

Device	# Nodes	# Accelerators (total)	Batch Size per Accelerator	Total Batch Size	Images / sec	Images / sec / device
Gaudi2	32	256	16	4096	12,654	49.4
A100-80GB	32	256	16	4096	3,992	15.6

Fine Tuning





intel.
GAUDI

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

Fine Tuning

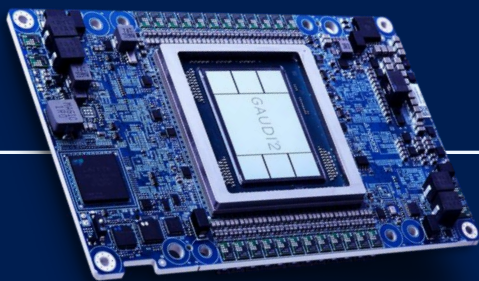
intel.
XEON®

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

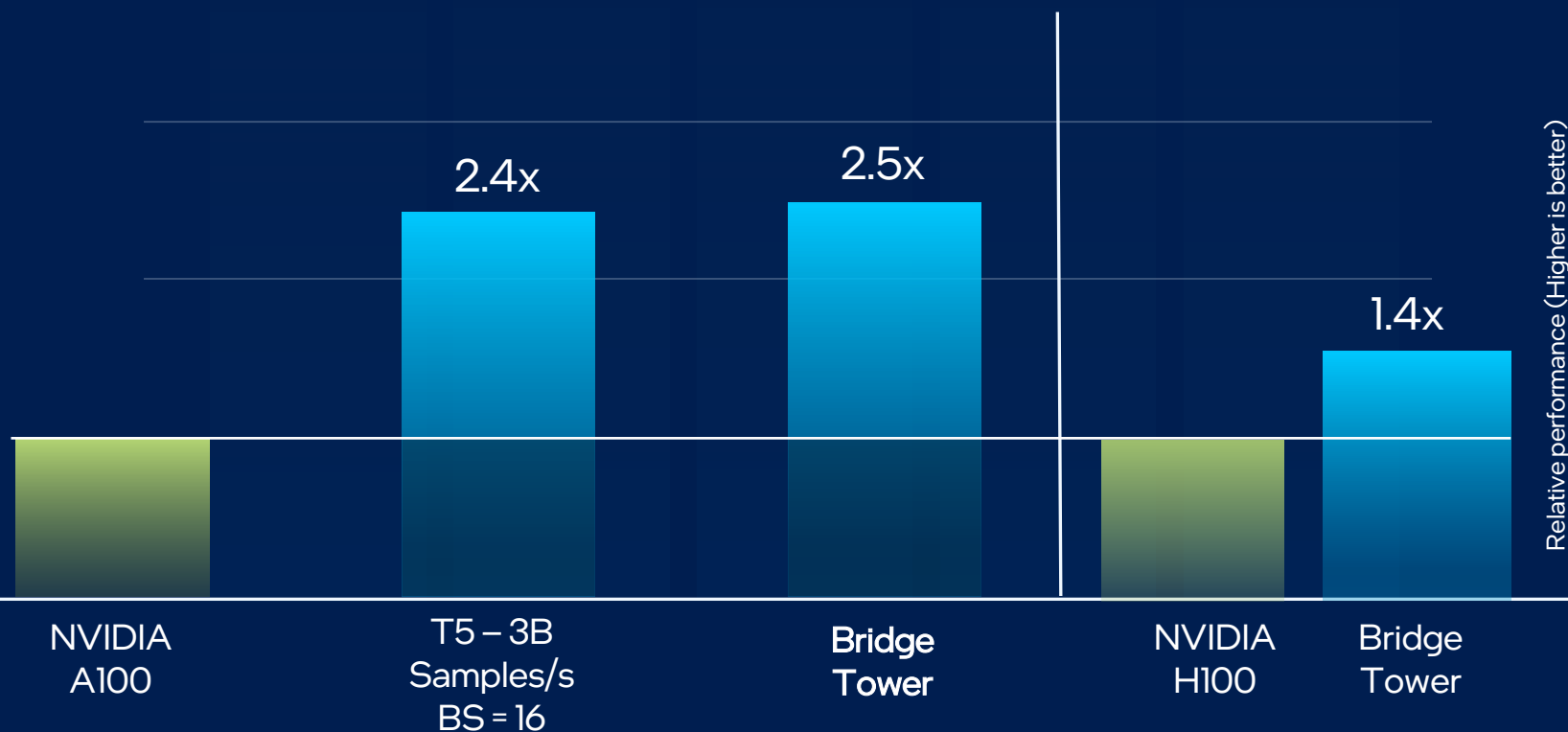
Intel Provides Solution Options for
Fine-tuning Gen AI and LLMs
to Fit Workload Needs



Fine-tuning Across Numerous LLMs



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





5th Gen Intel® Xeon® Fine Tuning

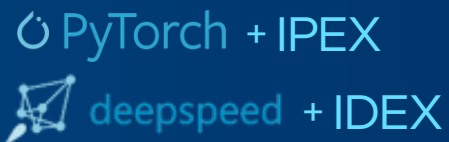
Optimized Models & Spaces

Dolly	LLAMA 2	MPT	LDM3D	Whisper	100k's Mode
-------	---------	-----	-------	---------	-------------

Intel Optimized Hugging Face Libraries & Tools

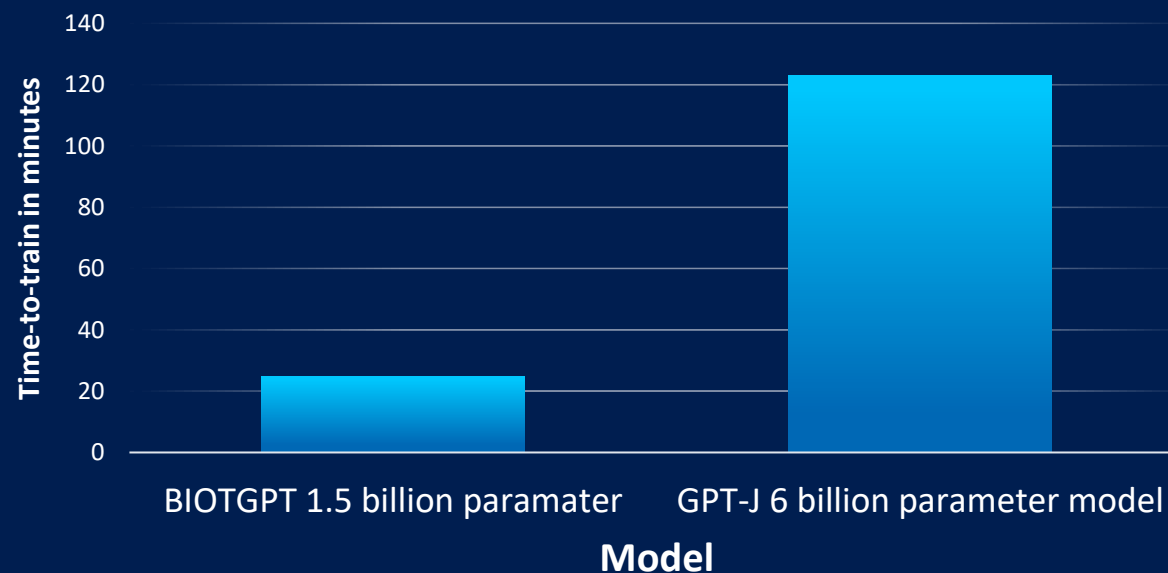
Transformers	Diffusers	Accelerate	PEFT	Optimum
Fine Tuning	Use Cases	Fine Tuning at Scale	Efficient Fine Tuning	Performance Optimization

Foundational Stack

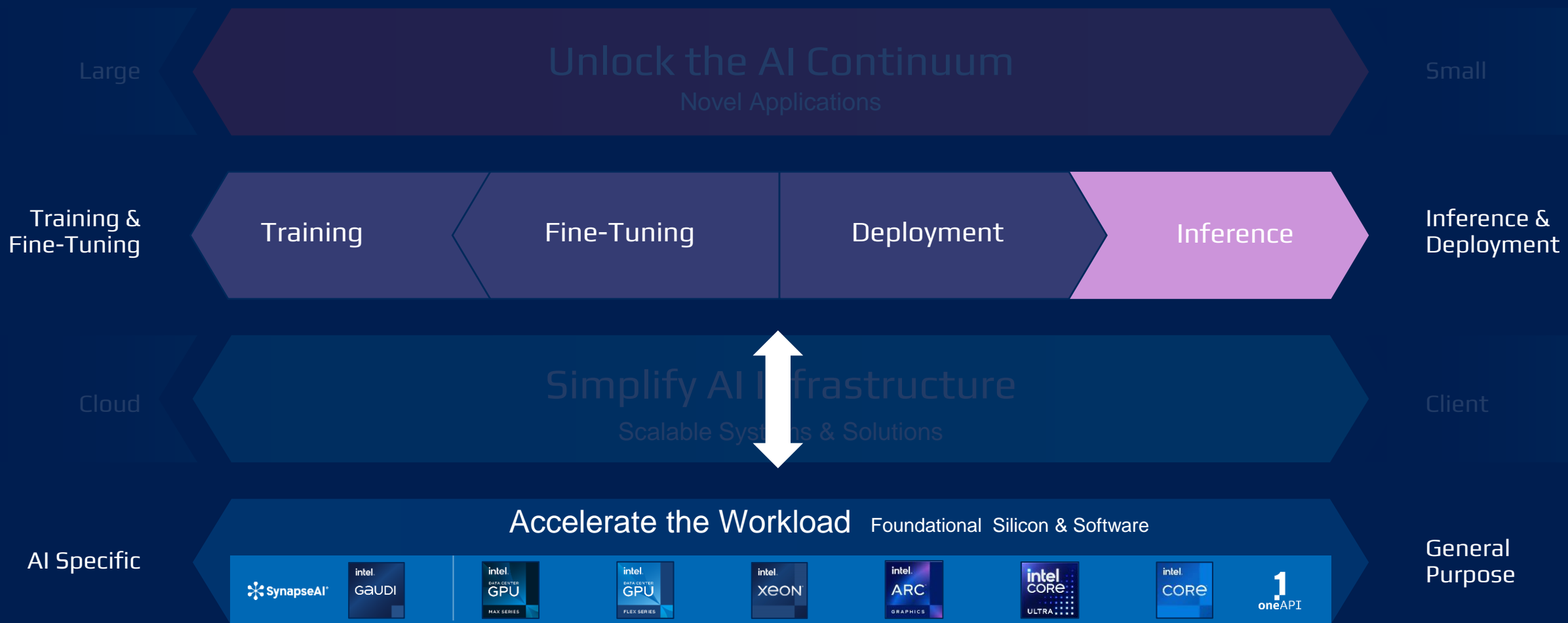


Multi-node Fine Tuning Open-source Commercial Large Foundational Models In Minutes To Hours

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

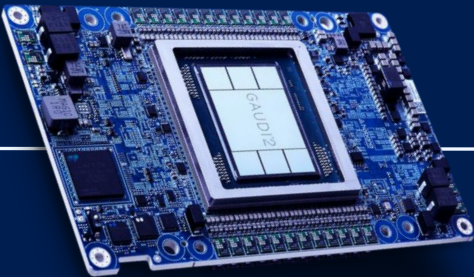


Inference





Inference Advantage Across Multiple LLM Performance Metrics

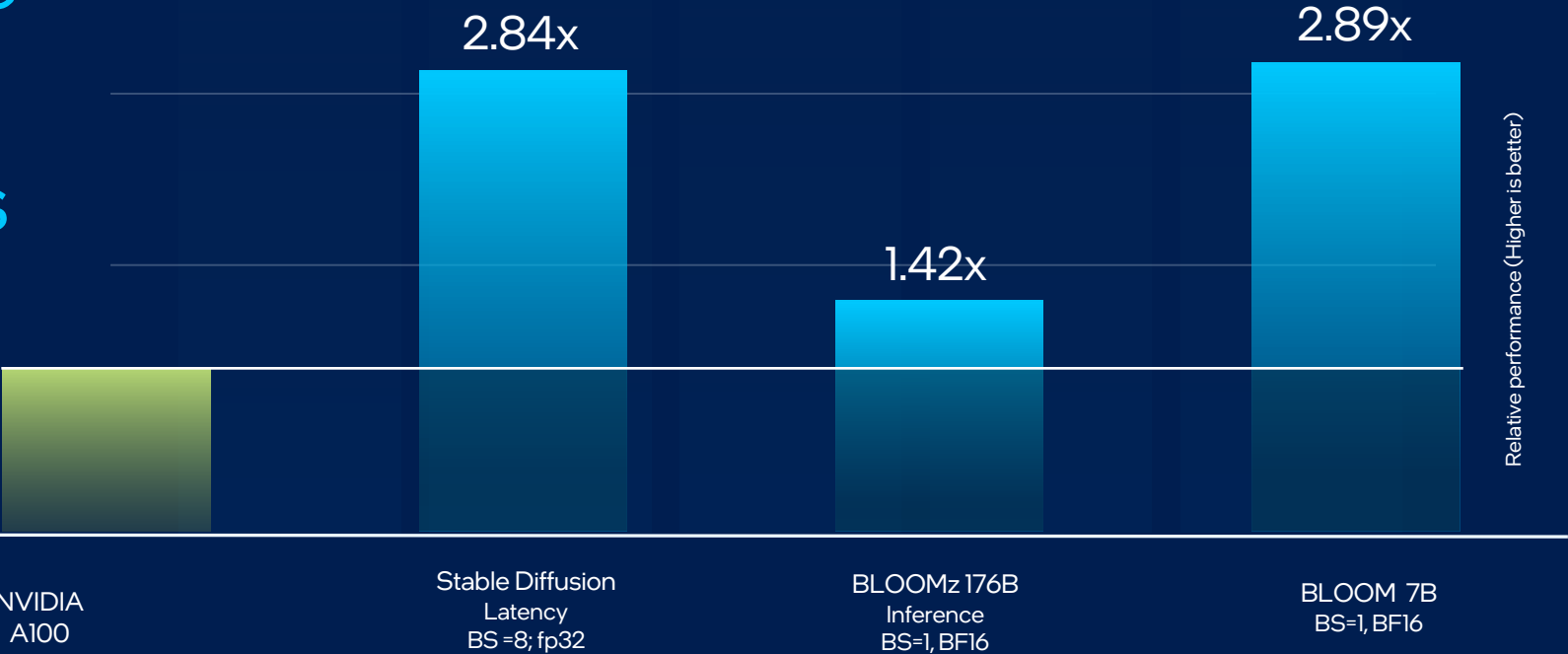


NVIDIA
A100

Stable Diffusion
Latency
BS=8; fp32

BLOOMz 176B
Inference
BS=1, BF16

BLOOM 7B
BS=1, BF16

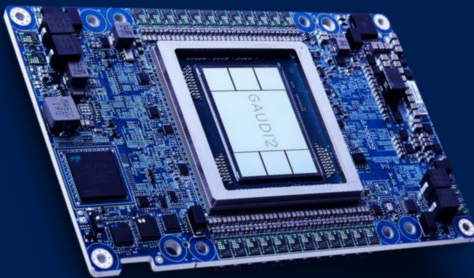


Energy Efficiency

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

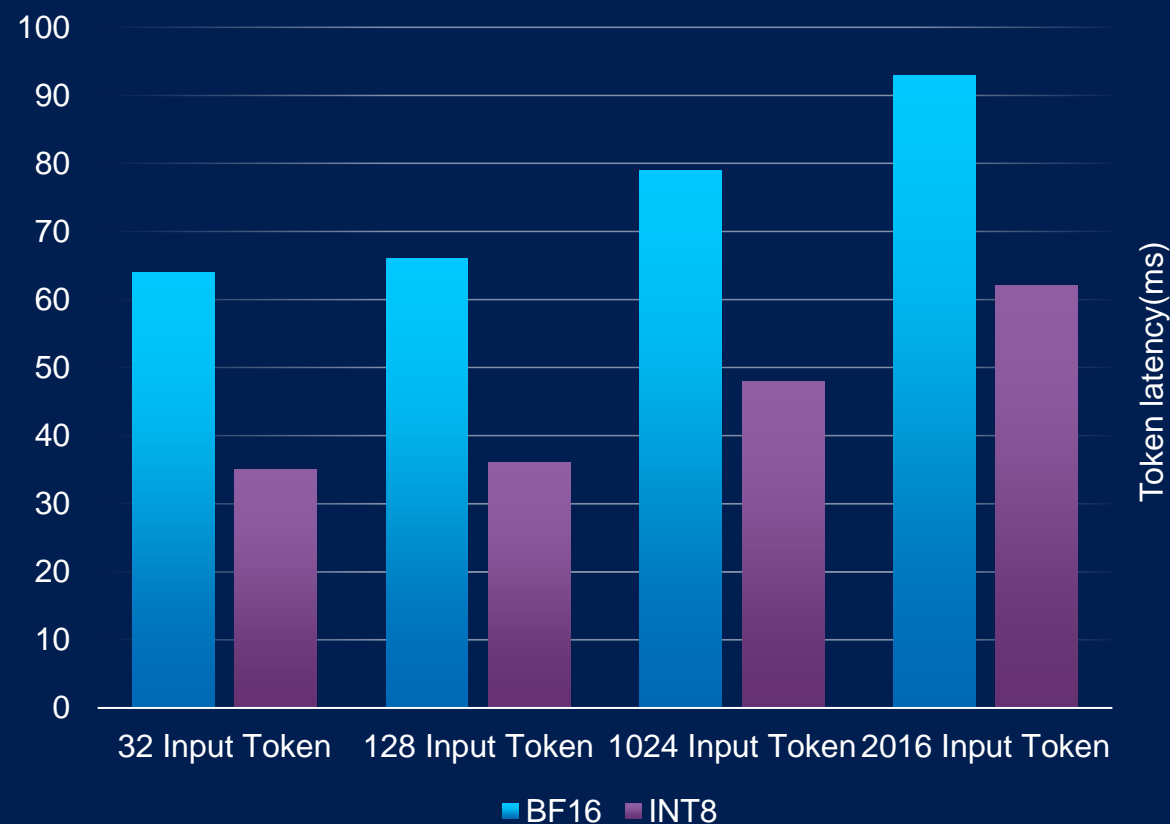
GPT-J On MLPerf Inference Benchmark



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

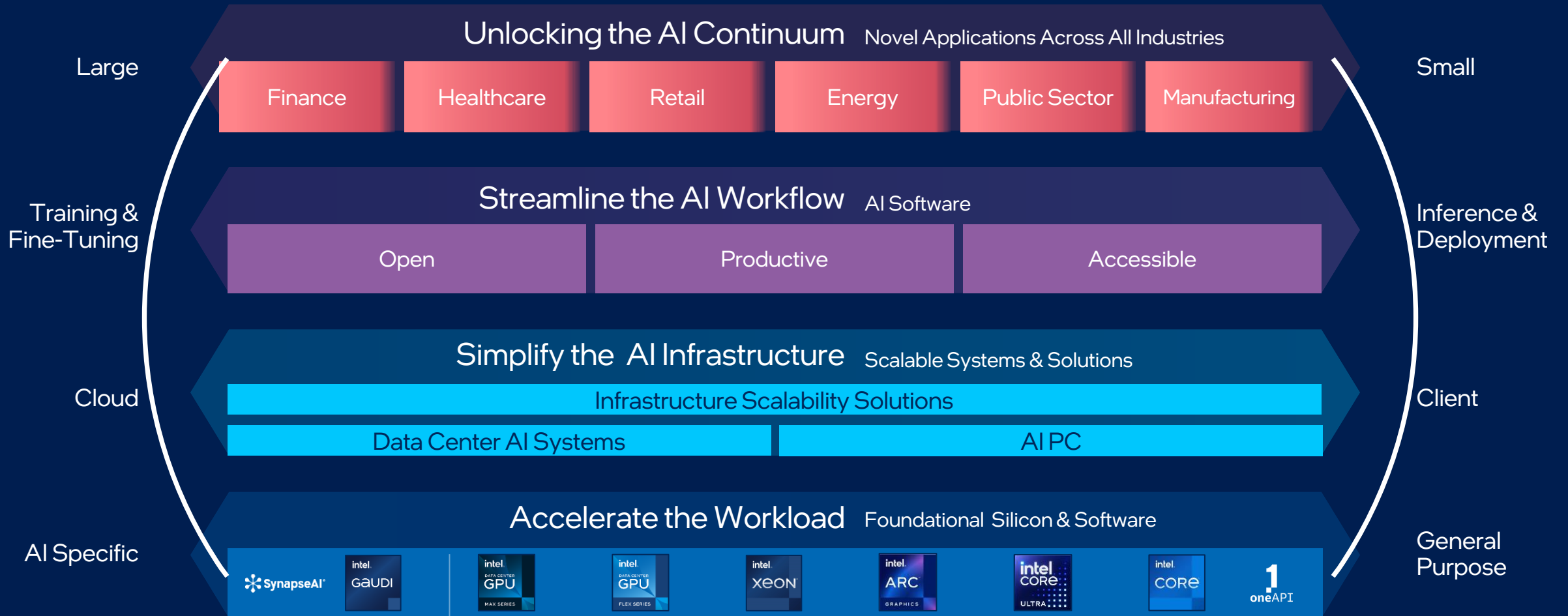
- Use any popular industry standard AI libraries
- Intel AI Platform validated with over 300 inference models
- One socket of 4th Gen Intel® Xeon® processors can run LLaMa2 chatbots in under 100ms 2nd token latency

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

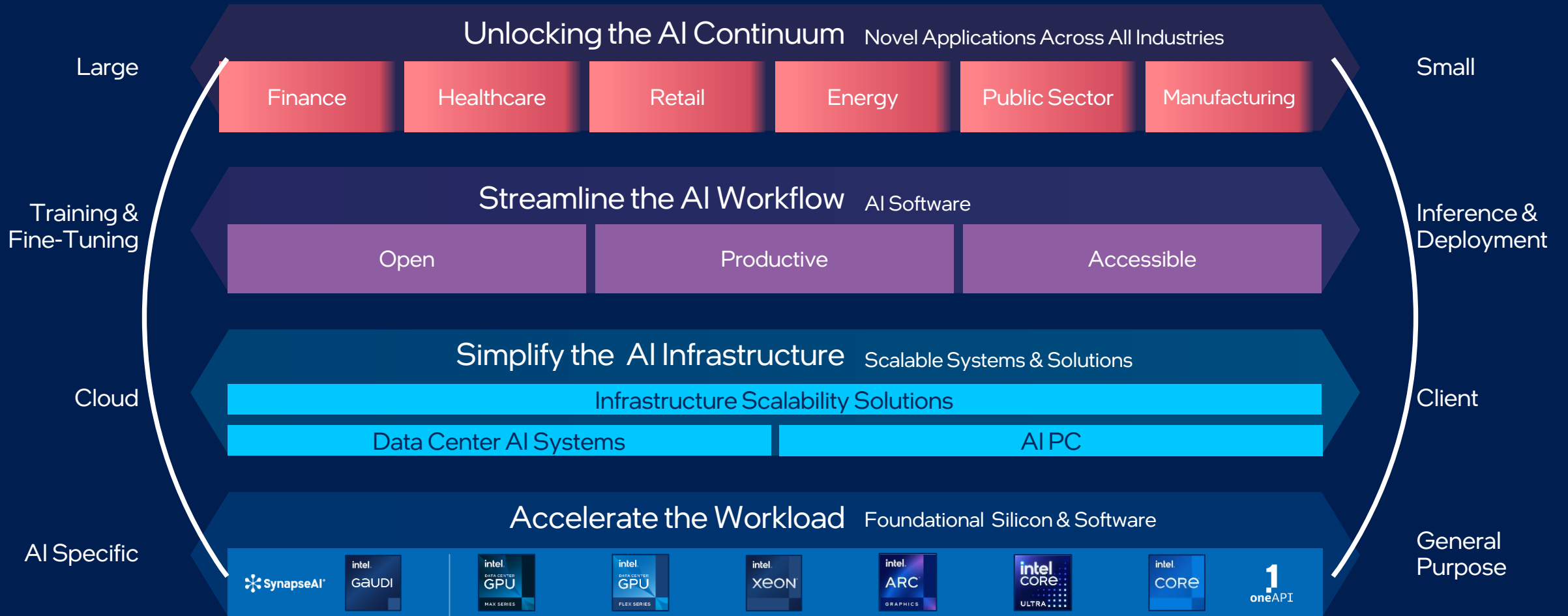


Results shown for bare metal.

Intel's Approach



Intel's Approach



您是否有興趣申請 Intel® Developer Cloud 優惠券？

掃描右側 QR Code
將能獲得一組使用優惠代碼
搶先體驗 Intel® 的硬體和軟體雲端 AI 服務
價值美金 \$250，有效期限至 2024 年 7 月 30 日



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.

intel[®] Ai
summit

Thank You!

