

世界の難問を紐解く鍵となる

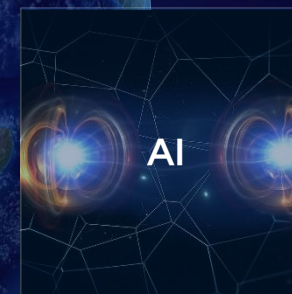
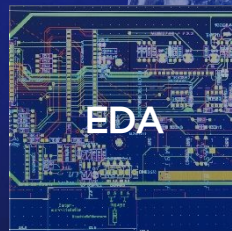
AIで加速するHPC

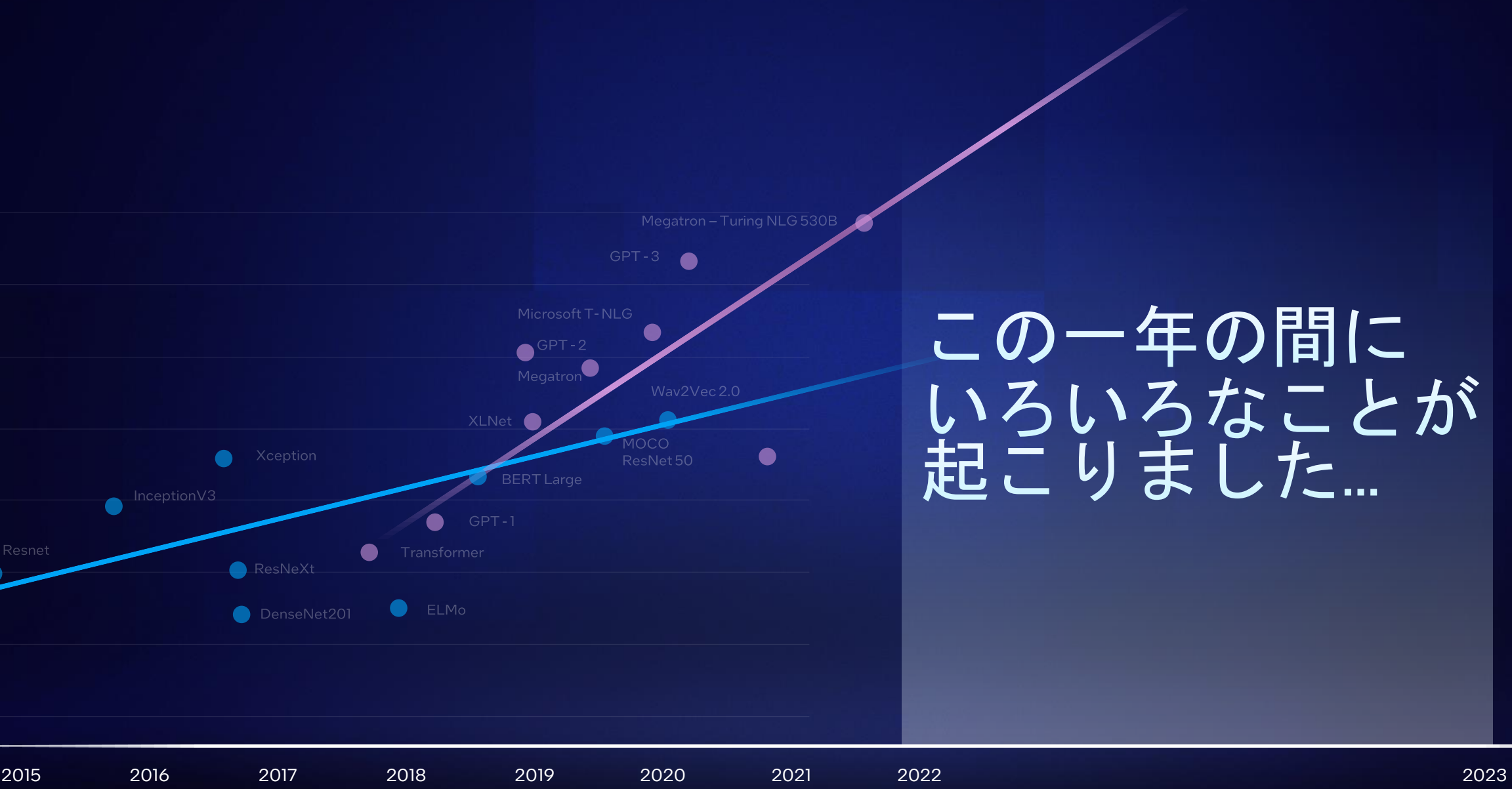
矢澤 克巳

インテル株式会社 インダストリー事業本部

intel[®]

シミュレーションへの 飽くなき要求





この一年の間に
いろいろなことが
起こりました...



Source: Moore, S. (2022), "Nvidia's Next GPU Shows That Transformers Are Transforming AI", IEEE Spectrum.

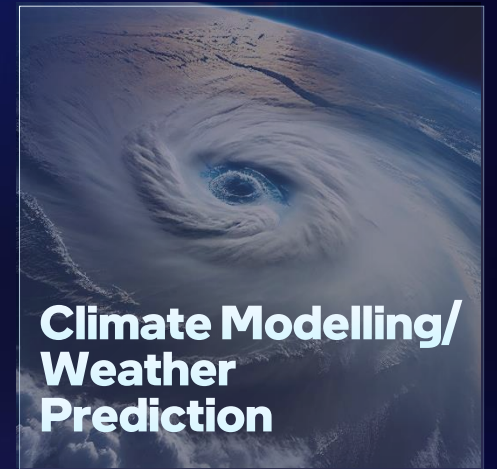
大規模言語モデル(LLM)の爆発的急伸



2022 June July August September October November December January February March April May 2023

Source: Zhao, W. X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., ... & Wen, J. R. (2023). A survey of large language models. arXiv preprint arXiv:2303.18223.

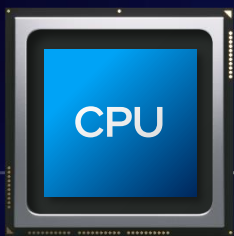
科学分野に AIブームが到来...



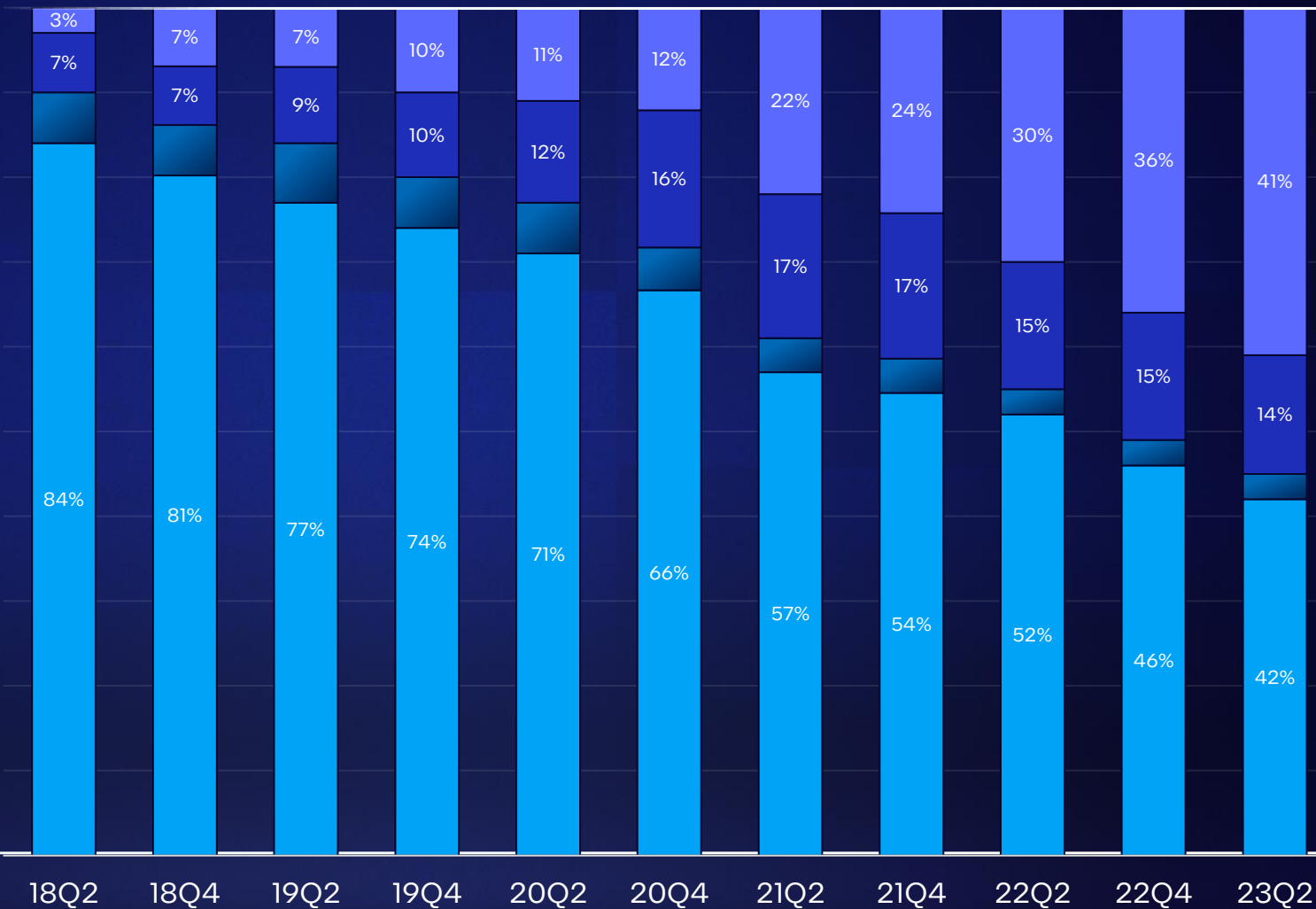
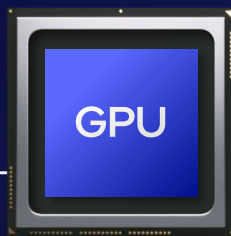
Top100の CPU:GPU比率

CPU : GPU Ratio

- 1:0
- 2:1, 2:2, 1:2
- 2:4, 2:6
- 2:8, 1:4, 2:16



Or/
And



HPC利用でのCPU:アクセラレータ比率*



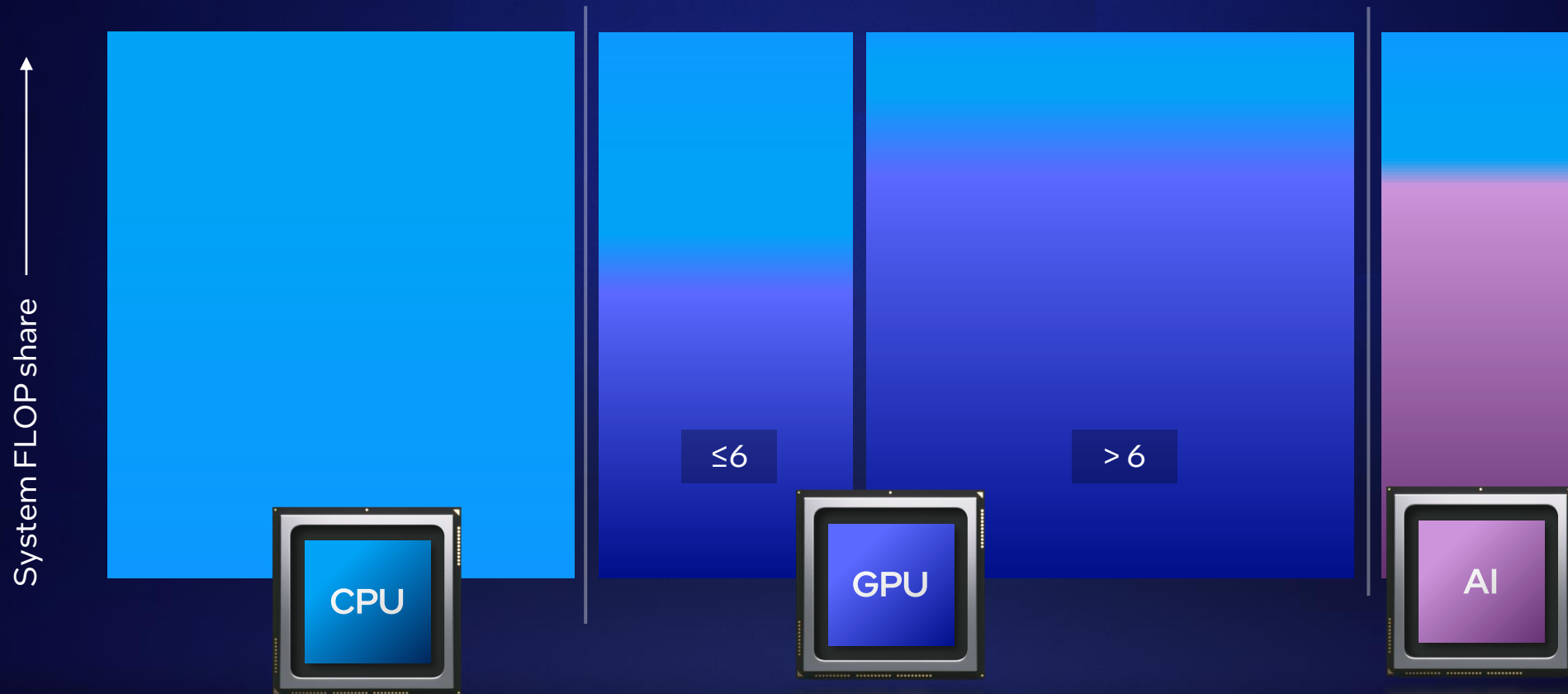
* 抽象化した図

大規模AI利用でのCPU:アクセラレータ比率**



**市場動向から見たコンセプト

HPCとAIの融合をどう実現させるか？

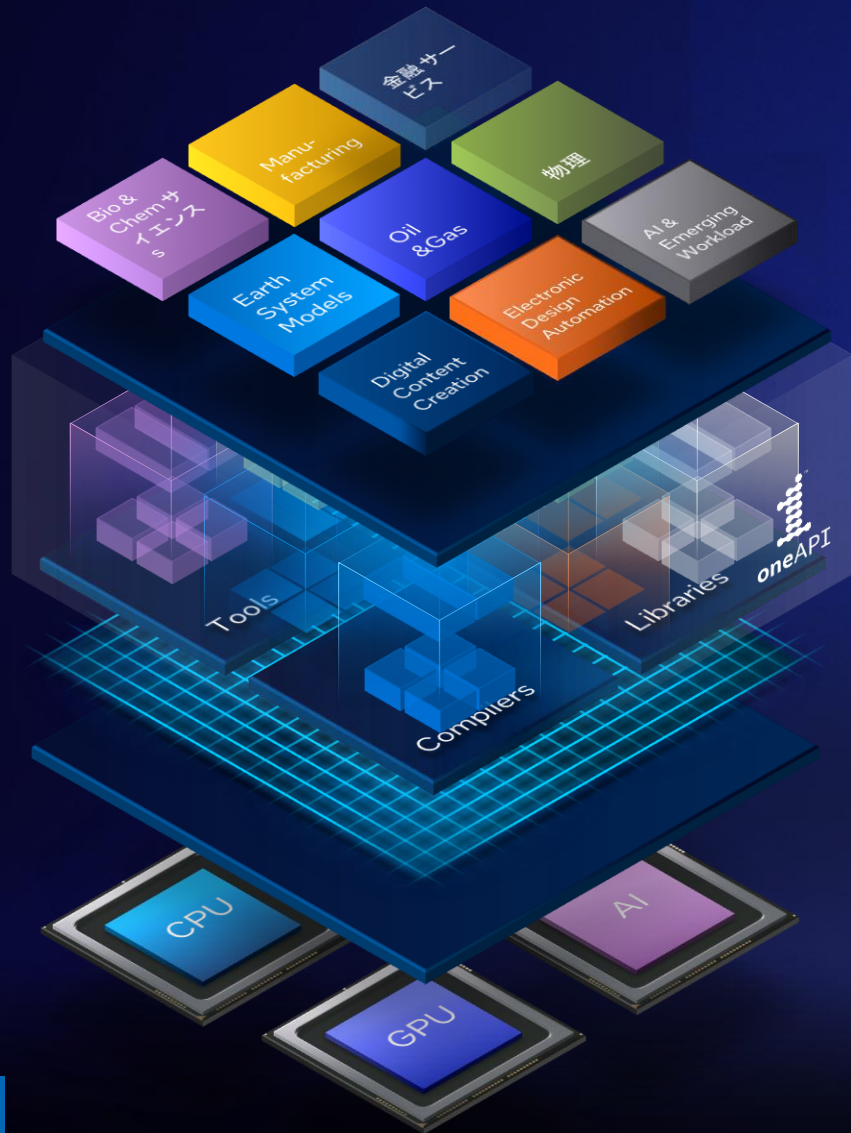


世界でもっとも困難な 技術的課題の解決に向け

業界と協力して提供する
サステナブルな HPC

アーキテクチャやベンダーを問わず
オープンかつ統一されたソフトウェア

AIで加速するHPCの未来を支える
計算エンジン



2023年1月10日（日本時間11日）正式発表

第4世代 インテル® Xeon® スケーラブル・プロセッサ

インテル® Xeon® CPU マックス・シリーズ

インテル® データセンター GPU マック ス・シリーズ



インテル データセンター ロードマップ

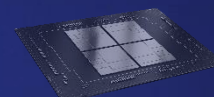
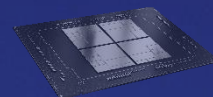
E-Core
CPU

Intel® Xeon® Processor
コードネーム Sierra Forest

Intel® Xeon® Processor
コードネーム Clearwater Forest

P-Core
CPU

第四世代 Intel® Xeon®
Scalable processor



Intel® Xeon®
CPU Max Series

第五世代 Intel® Xeon®
Processors
コードネーム Emerald Rapids

Intel® Xeon®
Processors
コードネーム Granite Rapids

AI
専用アクセラレータ

Habana®
Gaudi® 2

Habana®
Gaudi® 3

次世代 GPU
コードネーム Falcon Shores

HPC & AI
GPU

Intel® Data Center GPU Max Series

Visual Cloud
GPU

Intel® Data Center GPU Flex Series

Intel® Data Center GPU Flex Series
コードネーム Melville Sound

FPGA

Intel STRATIX 10 eASIC AGILEX 15+ new FPGAs on
schedule to PRQ in 2023

Intel AGILEX eASIC Next Gen FPGAs

2023

2025+

第四世代 Intel® Xeon® Processors

400+

Design Wins

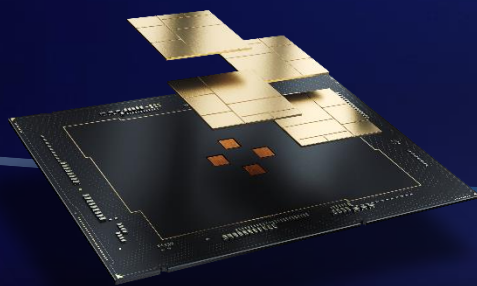
The most ever for any Xeon family

Top 10

Global CSPs*
deploying now
and throughout 2023

* Cloud service providers

提供中



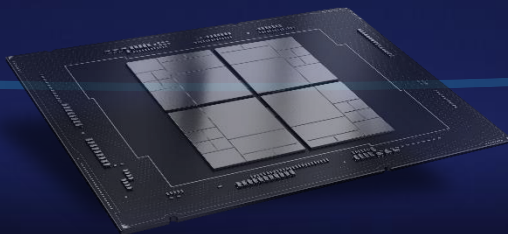
第五世代 Intel® Xeon® Processors

コードネーム Emerald Rapids

より高い消費電力あたりの性能

Same platform as
第四世代 Xeon

Q4 2023 予定



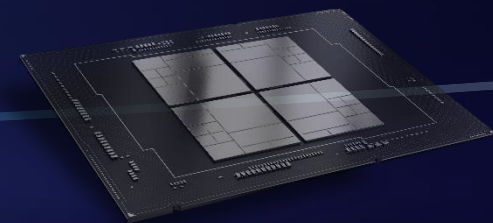
Intel® Xeon® Processors

コードネーム Granite Rapids

Intel 3 で作る初の P-Core Xeon

Increased core density,
memory & I/O innovations

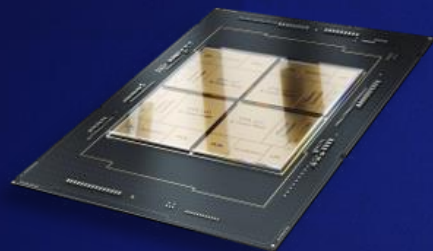
2024 予定



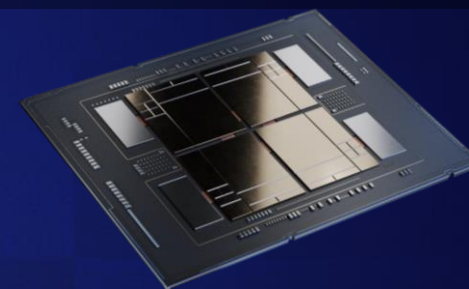
インテル® Xeon® プロセッサ

HPC とAI アクセラレーションに最適化された特徴的な新機能

第4世代インテル® Xeon® スケーラブル・プロセッサ



インテル® Xeon® CPU マックス



ブレイクスルー・テクノロジー

DDR5

強化されたメモリバンド幅

PCIe 5

高いスループット

CXL 1.1

次世代 IO

内蔵AI アクセラレーション

インテル® Advanced Matrix Extensions (AMX)

ディープラーニングの推論および学習処理性能を向上

広帯域メモリ

HBM2e

広帯域を必要とするワークロードの性能を飛躍的に向上

インテル® Advanced Matrix Extensions (AMX)

ディープラーニングの推論および学習処理性能を向上

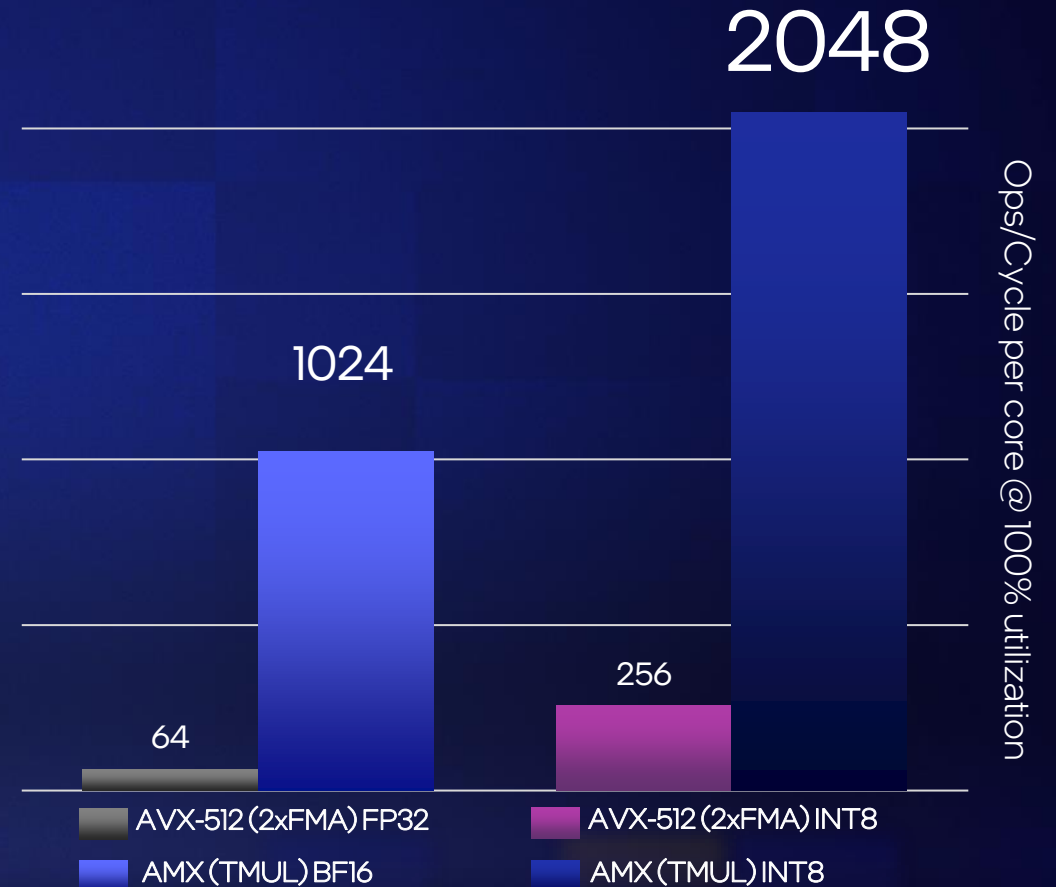
ディープラーニング
データタイプ

- int8 with int32 accumulation
- Bfloat16 with IEEE SP accumulation

ISAレベルでの加速

- 完全なインテルアーキテクチャの
プログラマビリティ
- 低遅延

業界関連のフレームワーク
およびライブラリが利用可能

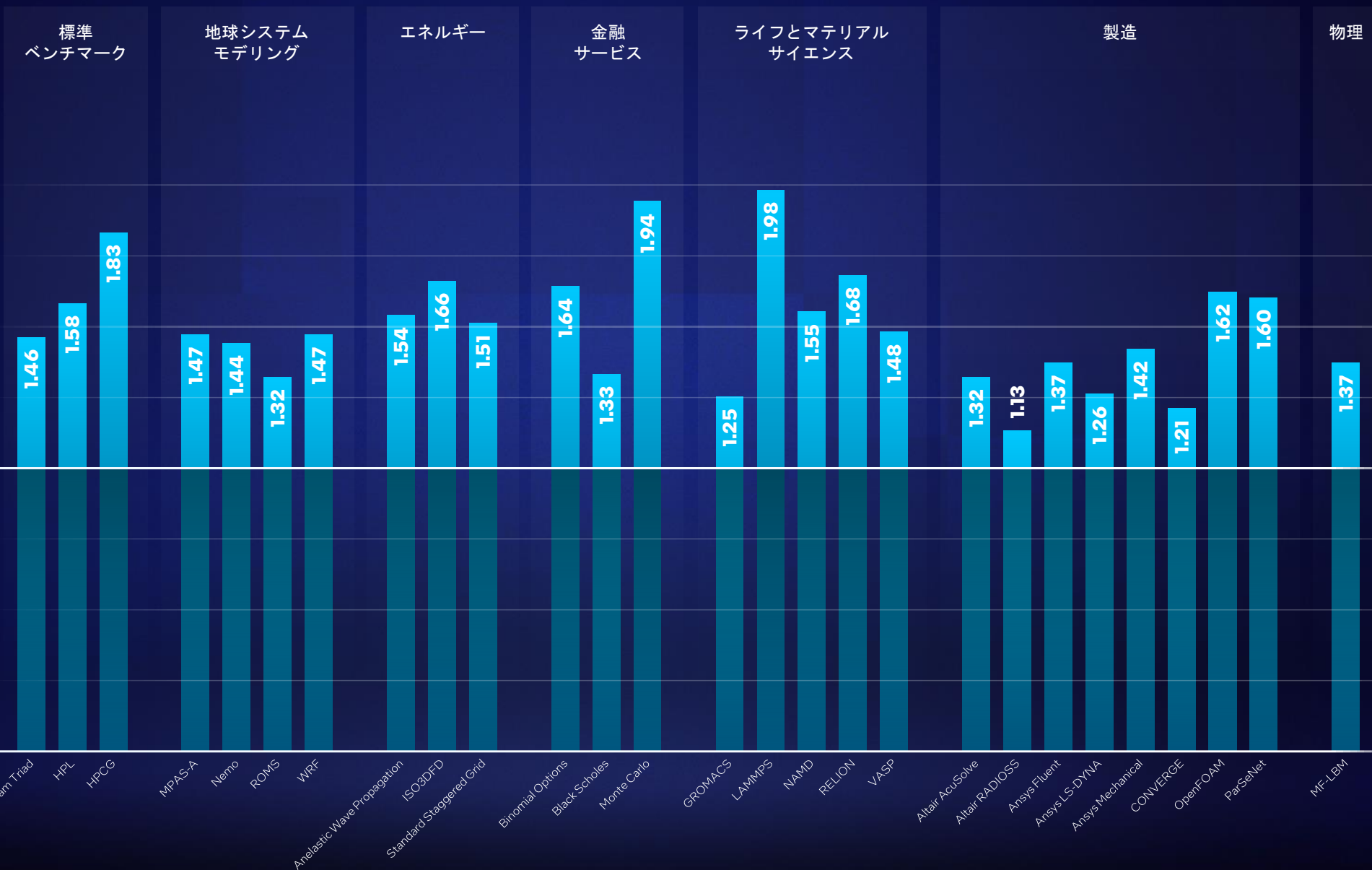


Results have been simulated. For workloads and configurations visit www.intel.com/ArchDay21claims. Results may vary



最大 1.5倍の 性能

2S Intel® Xeon® 8480+
vs. 2S AMD EPYC 7763



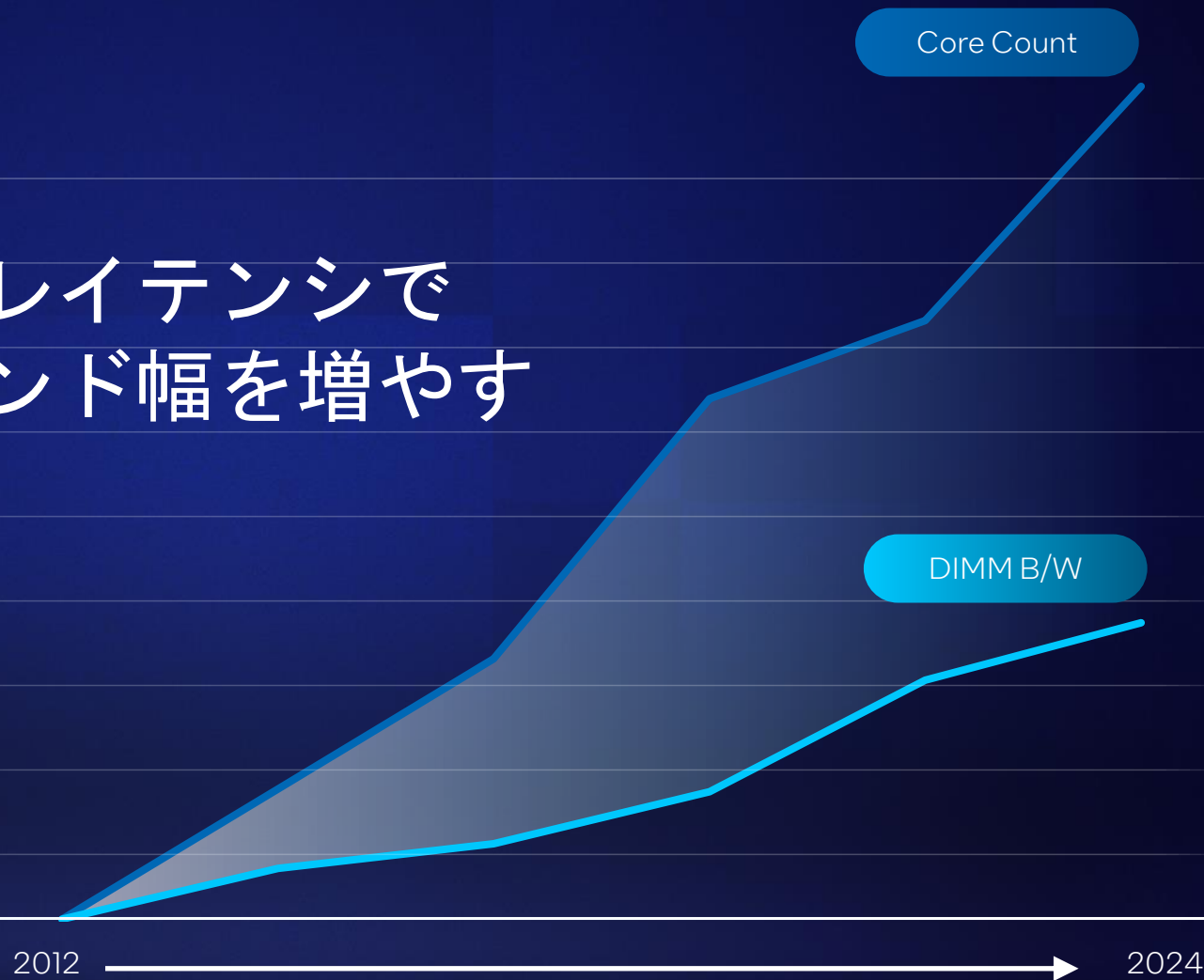
Relative performance (Higher is better)

Visit www.intel.com/performanceindex for workloads and configurations. Results may vary

This offering is not approved or endorsed by OpenCFD Limited, producer and distributor of the OpenFOAM software via www.openfoam.com, and owner of the OPENFOAM® and OpenCFD® trademark

あらたなHPC & AIニーズ

適正なコスト・電力・レイテンシで
容量あたりのメモリバンド幅を増やす

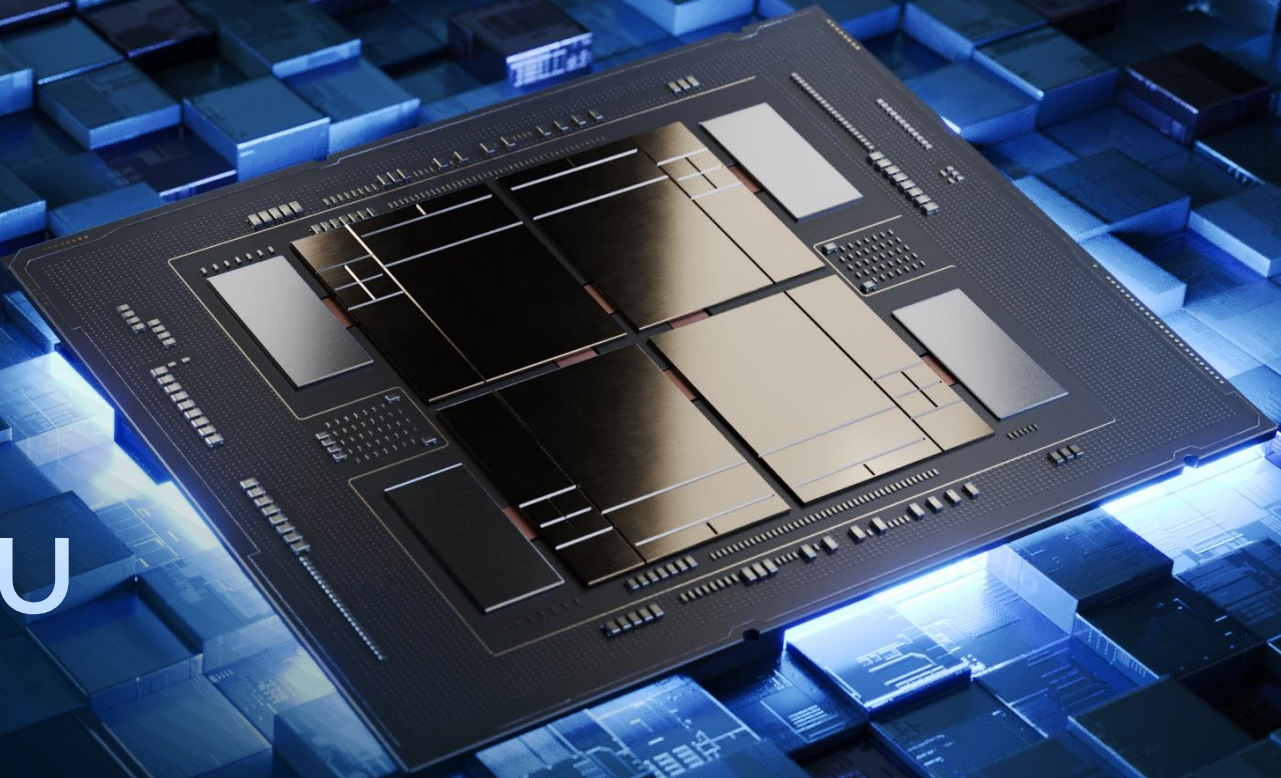


Max Core Count vs. DIMM BW Cumulative Growth



X86初 HBMを持つ唯一のCPU

最適なメモリ選択肢となりうる



Memory Modes

64GB
HBM2e
4 stacks of 16GB

最大
220GF/s
HPCG

最大
2GB
HBM per Core

HBM Only
Bootable from HBM
No code change

HBM ~~DDR~~

HBM Flat
2 Memory Regions
SW Optimization Needed

HBM DDR

HBM Caching
HBM as cache for DDR
No code change

HBM — DDR



See backup for workloads and configurations. Results may vary.



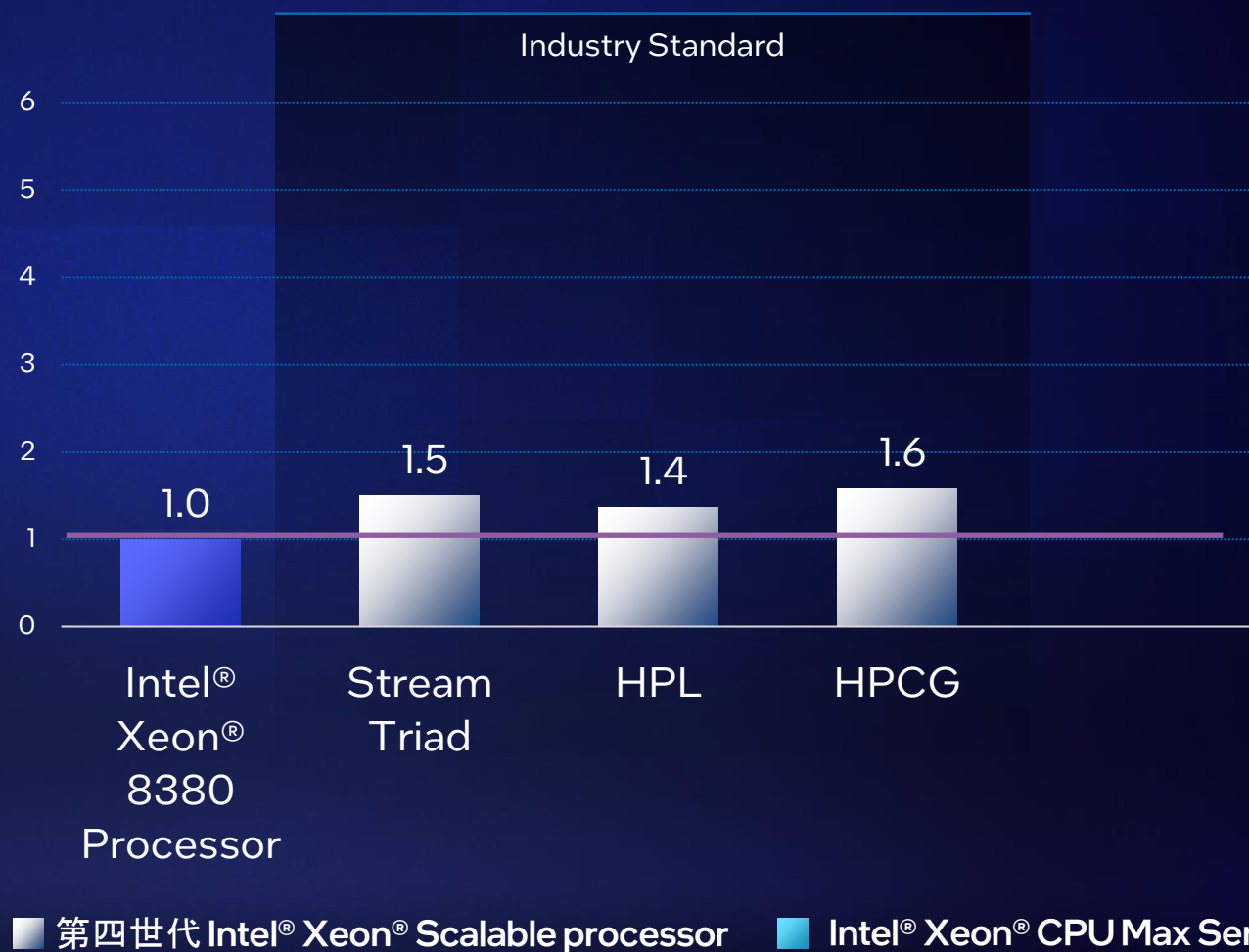
従来のXeonでは最大

1.6倍

Performance Industry Standard Benchmarks

2S 4th Gen Intel® Xeon® processor vs.
2S 第三世代 Intel® Xeon® 8380 processor

Relative Perf. Higher is better



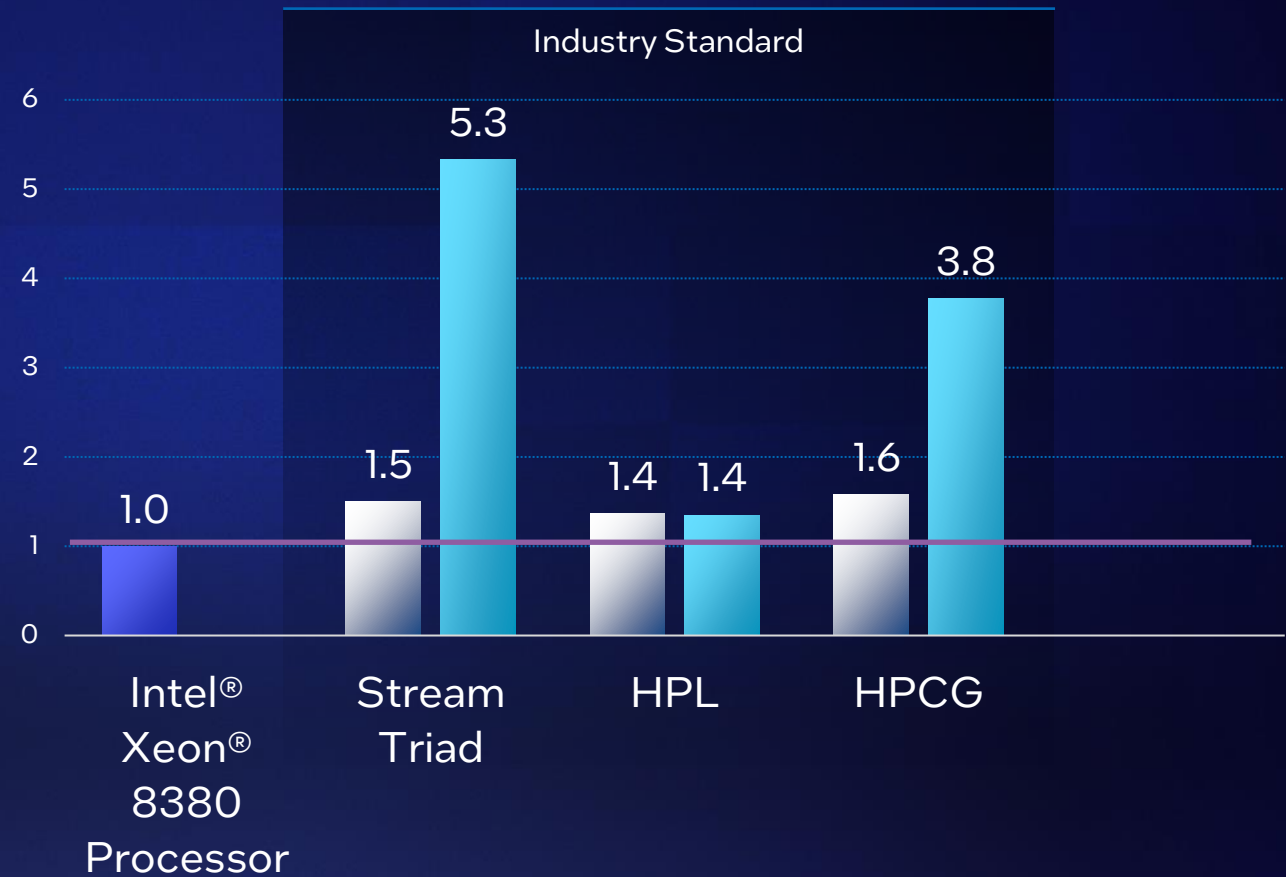
See backup for workloads and configurations. Results may vary.



CPU MAXでは最大 5倍 Better Performance for Memory Bandwidth

2S Intel® Xeon® CPU Max Series vs.
2S 第三世代 Intel® Xeon® 8380 processor

Relative Perf. Higher is better



■ 第四世代 Intel® Xeon® Scalable processor ■ Intel® Xeon® CPU Max Series

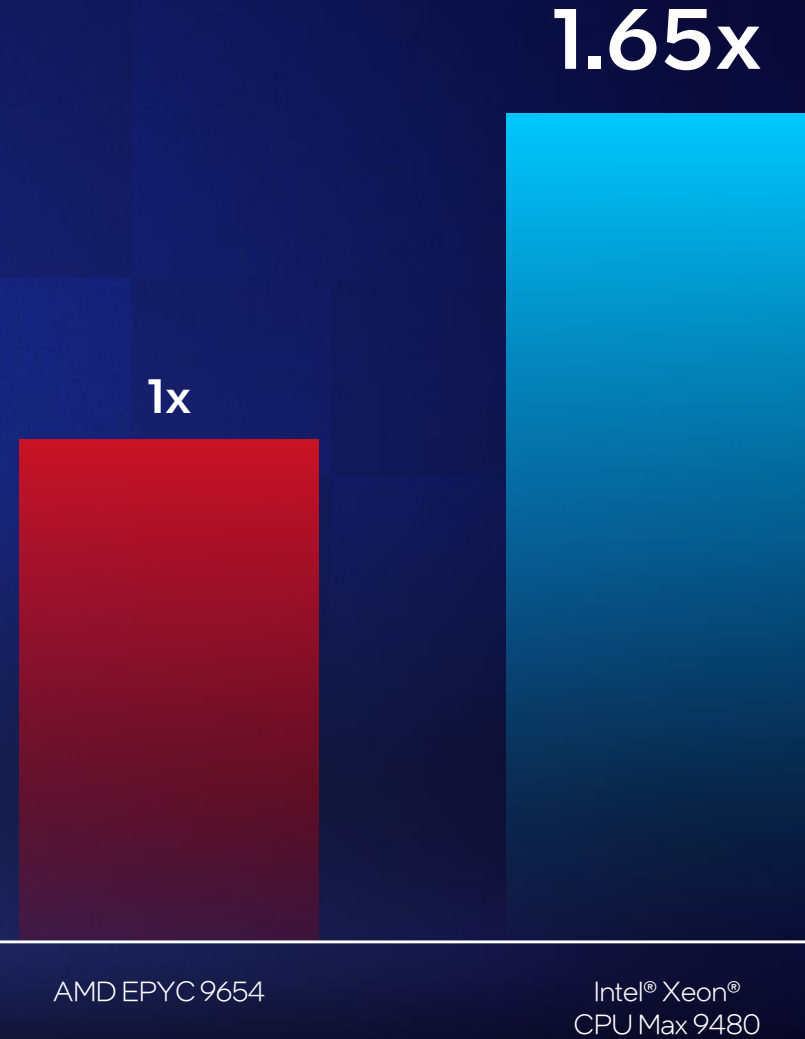
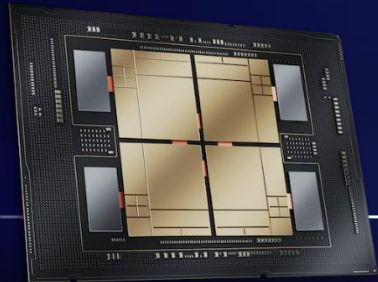
See backup for workloads and configurations. Results may vary.





コア数差を超えた性能

HPCG Performance



Relative performance (Higher is better)

AMD EPYC 9654

Intel® Xeon®
CPU Max 9480

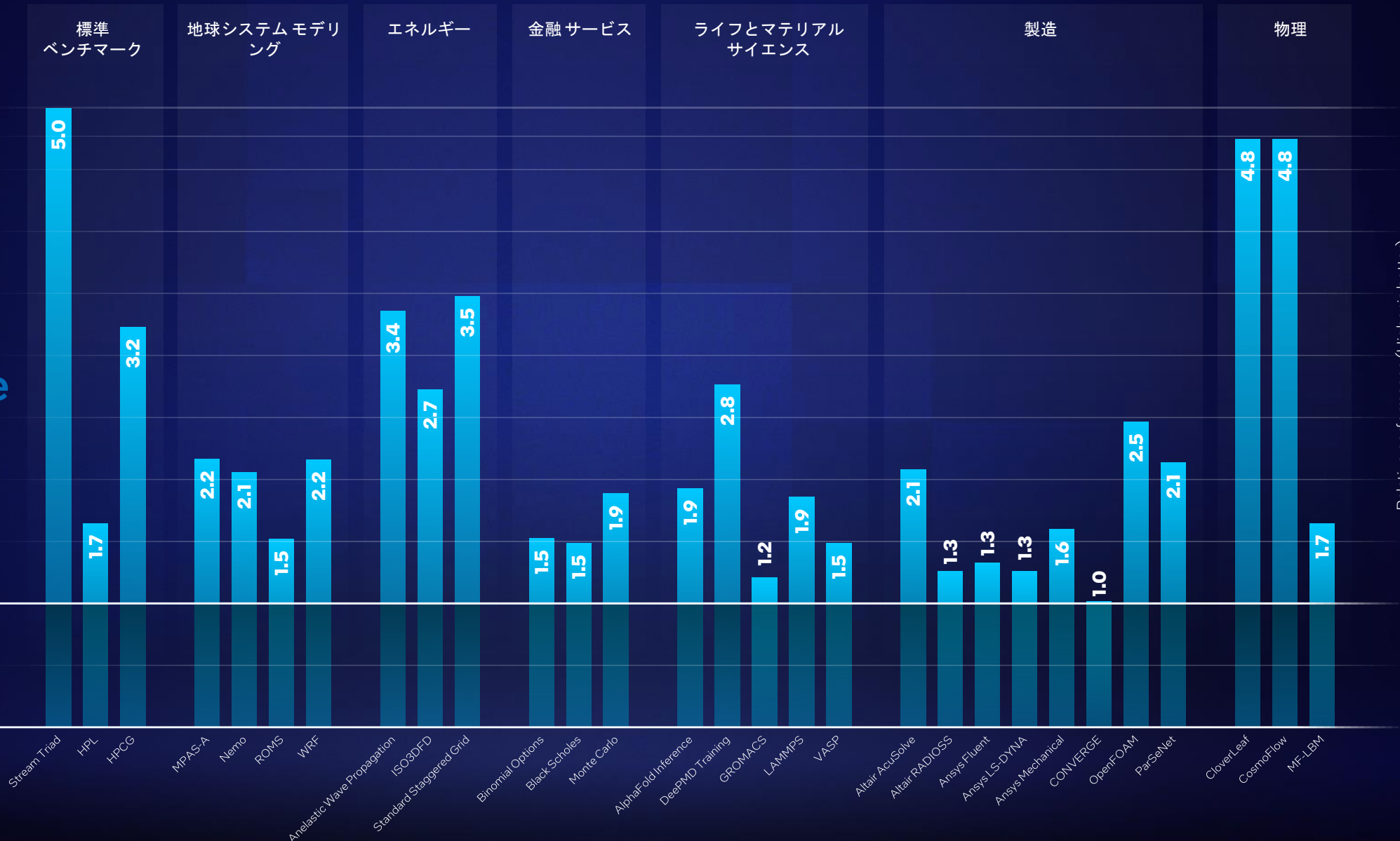
Visit www.intel.com/performanceindex for workloads and configurations. Results may vary





CPU MAXでは平均
2倍
performance

2S Intel® Xeon® CPU Max 9480
vs. 2S AMD EPYC 7773X



Relative performance (Higher is better)

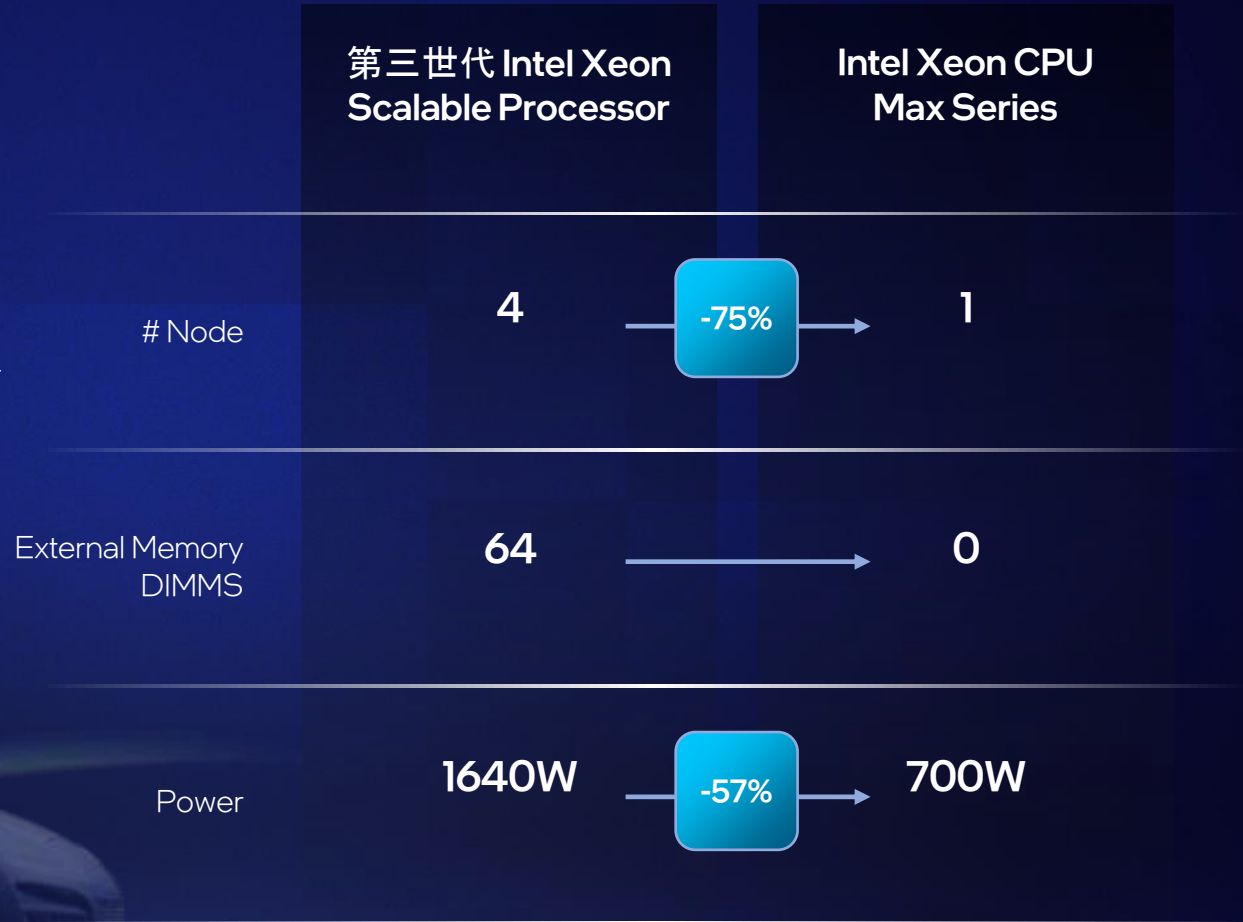
See backup for workloads and configurations. Results may vary.
 This offering is not approved or endorsed by OpenCFD Limited, producer and distributor of the OpenFOAM software via www.openfoam.com, and owner of the OPENFOAM® and OpenCFD® trademark.
 MLPerf™ HPC-AI v0.7 Training ベンチマーク Performance. Result not verified by MLCommons Association. Unverified results have not been through an MLPerf™ review and may use measurement methodologies and/or workload implementations that are inconsistent with the MLPerf™ specification for verified results. The MLPerf™ name and logo are trademarks of MLCommons Association in the United States and other countries. All rights reserved. Unauthorized use strictly prohibited. See www.mlcommons.org for more information.





少ないノード数と電力で 同等の性能を実現

 ALTAIR AcuSolve™



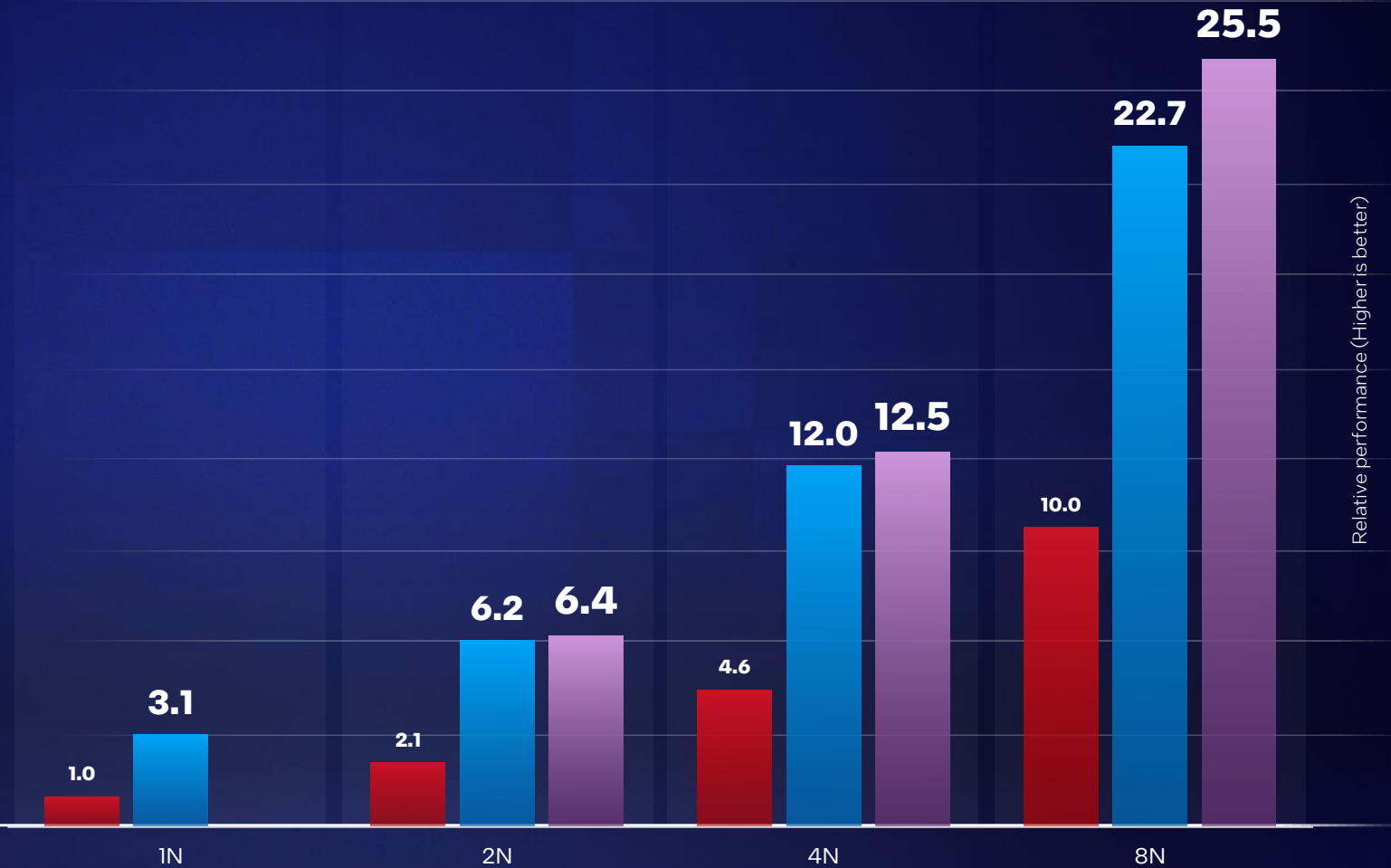
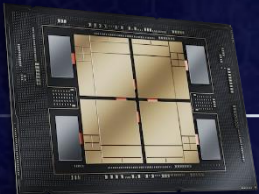


地球システム モデリング

MPAS-A – dyncore 30-km

91% Scaling Efficiency in cache mode

- AMD EPYC 7763
- Intel® Xeon® CPU Max 9480 (Cache Mode)
- Intel® Xeon® CPU Max 9480 (HBM Only)

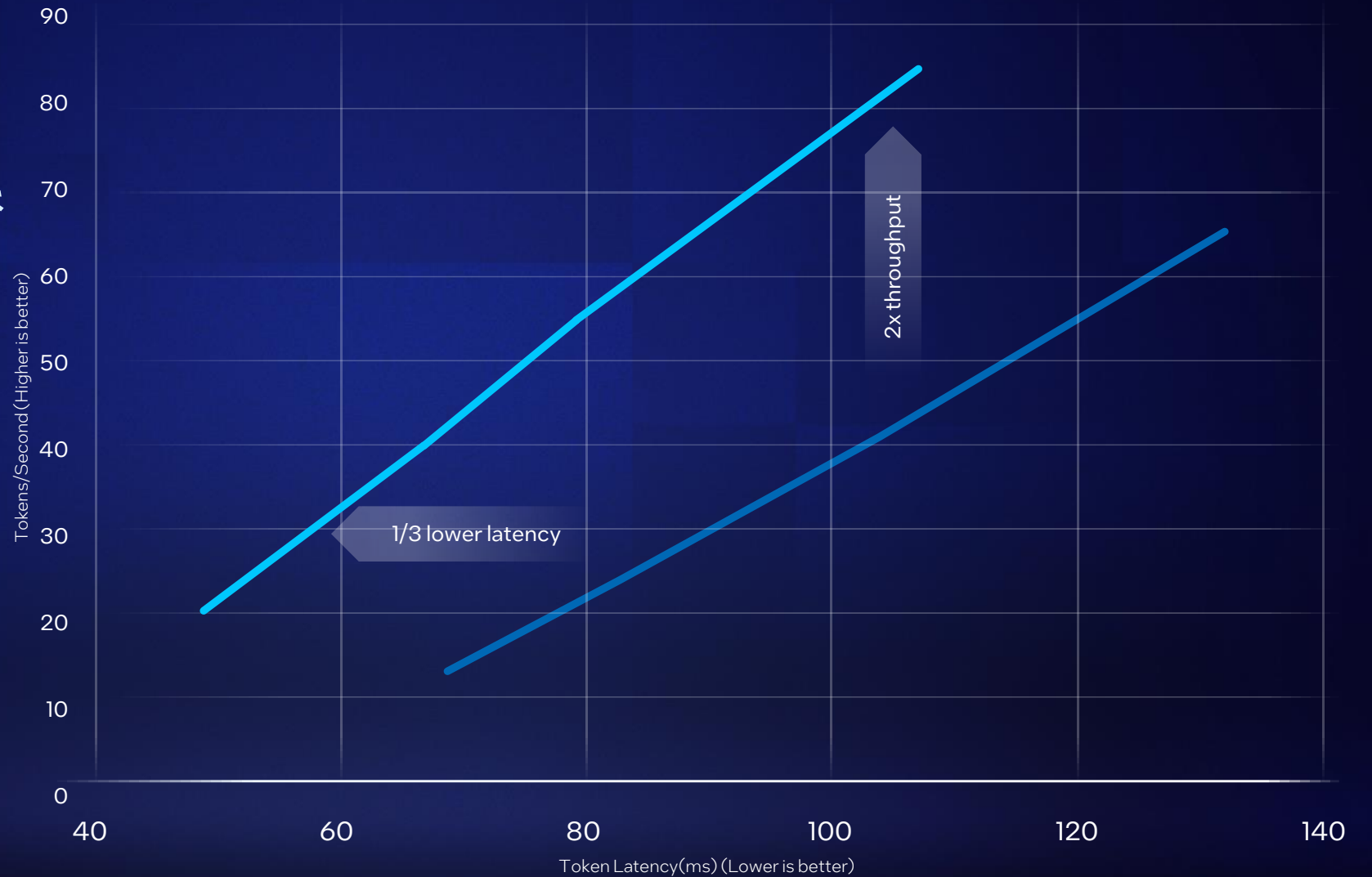
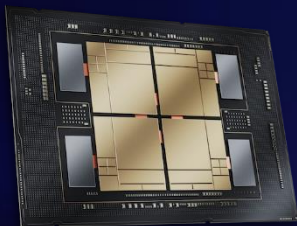




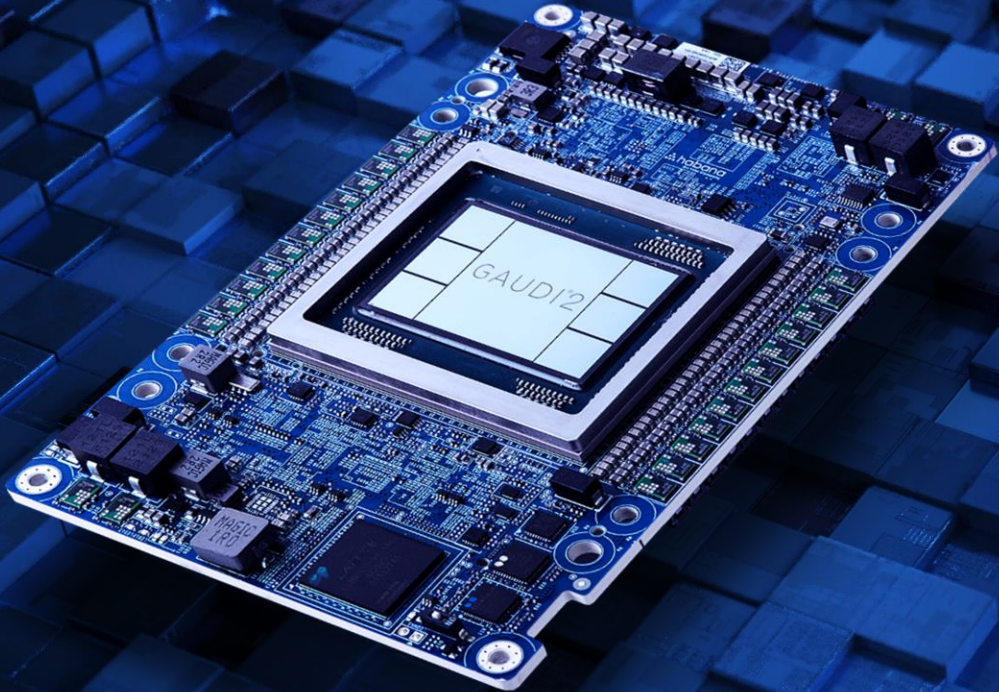
1/3のレイテンシで 2倍の推論 スループット

GPT-J 6B BF16 Inference Perf
(Batch Sizes 1,2,4,8)

- 1S Intel® Xeon® CPU Max 9480
- 1S Intel® Xeon® CPU 8480



GAUDI[®]2



7nm

Process
Technology

24

Tensor
Processor Cores

96 GB

On-Board
HBM2

48 MB

SRAM

24

Integrated
Ethernet ports

Availability

Cloud

Intel Developer Cloud

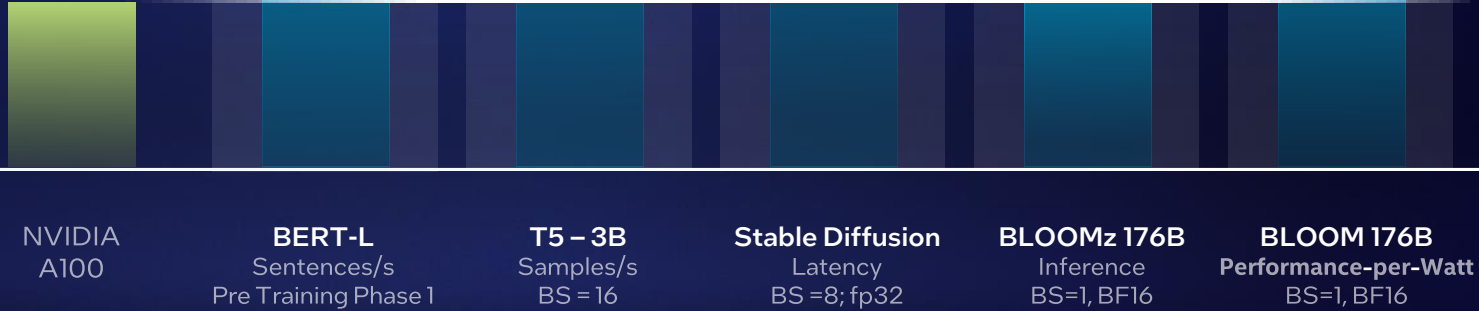
On-Prem

Supermicro Gaudi2 Server



Training ▶ Fine-tuning ▶ Inference ▶ Inference ▶ Power/Perf

数々の性能指標で 優位性を示す



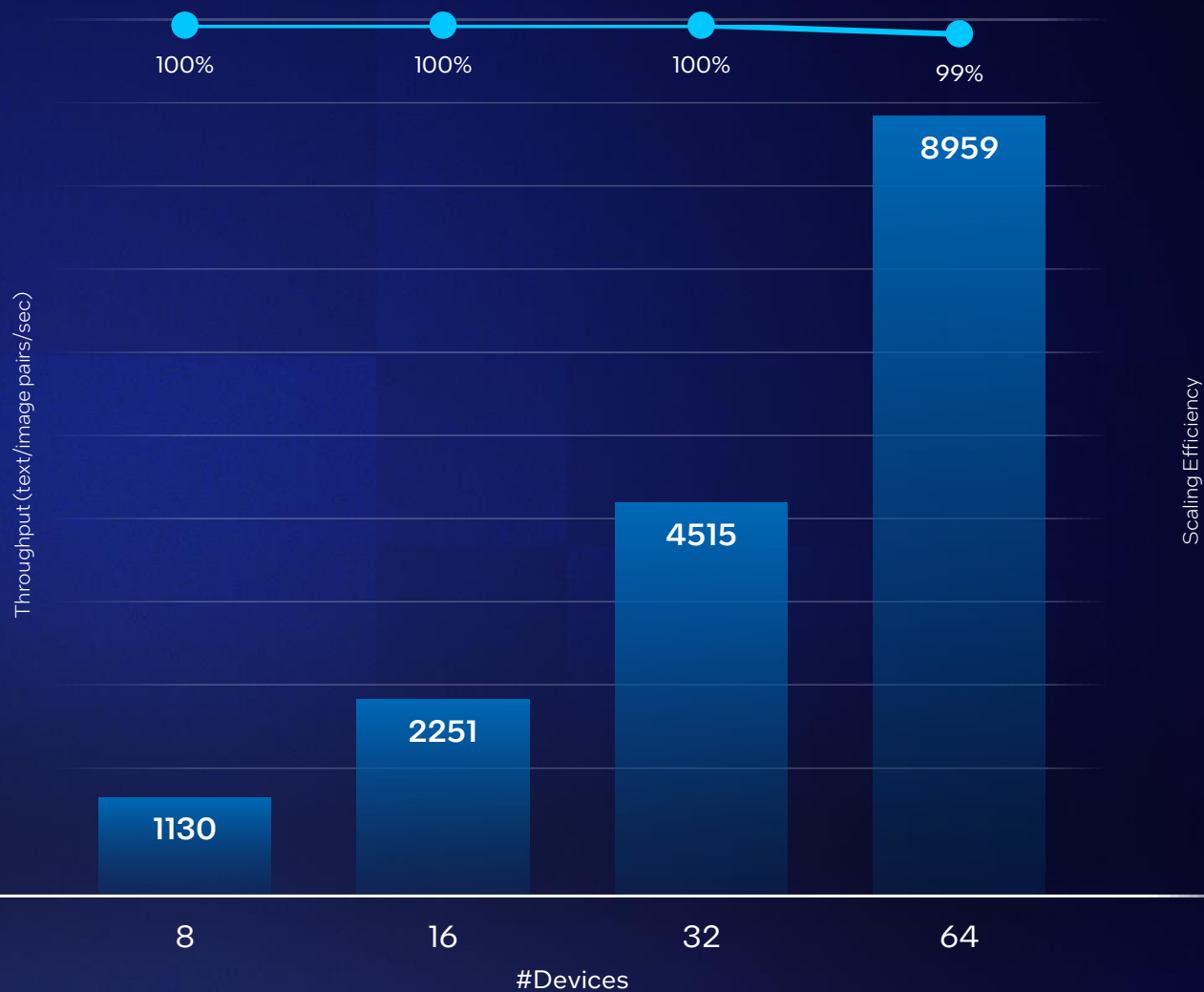
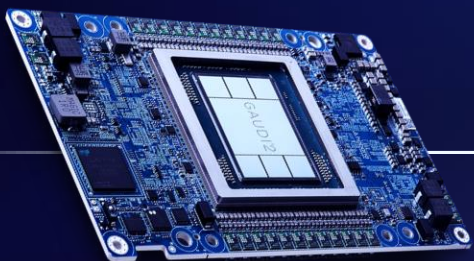
Relative performance (Higher is better)

Visit www.intel.com/performanceindex for workloads and configurations. Results may vary
<https://huggingface.co/blog/habana-gaudi-2-ベンチマーク>, <https://habana.ai/habana-claims-validation>



大規模で柔軟な統合された ネットワークによる ニアなスケーリング

最大64枚のカードまでニアなスケーリング効率
> 99%



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

intel
habana®

生成AIに最適

Large-scale

Training on Intel AI supercomputing cluster

Fast inference Speed

Habana Gaudi2 on 176B parameter BLOOM LLM

Real-time multimodal semantic search

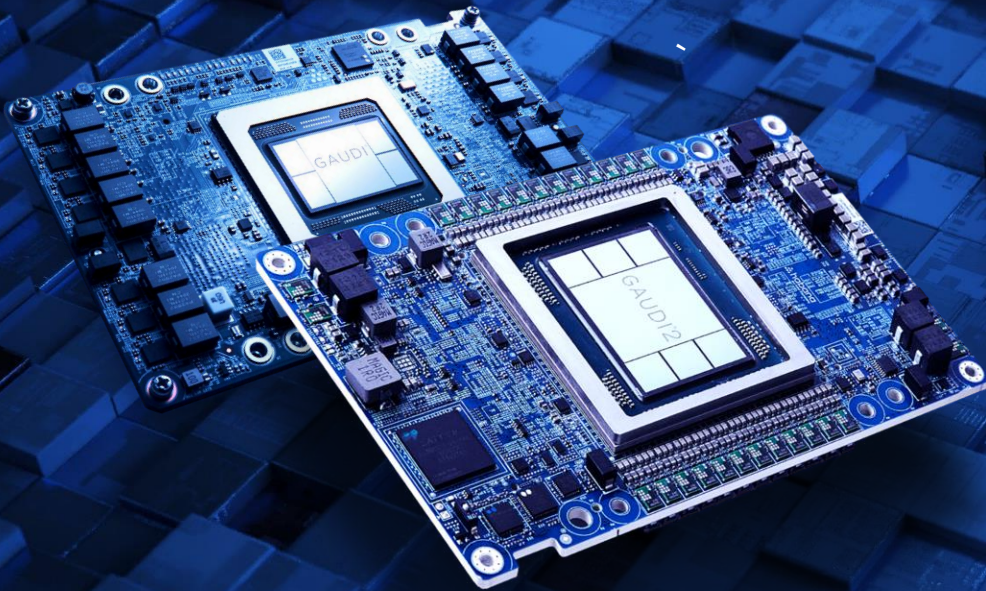
Running in the Intel Developer Cloud on Gaudi 2

12B Parameter Dolly 2.0 with GPT-3 like performance

Running on Habana Gaudi 2

New Model transforms text into immersive 3D

Trained on Intel AI Supercomputer





Intel® Data Center GPU Max Series

最大
128
Xe HPC
Cores

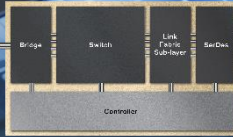
A diagram of the Xe-core architecture showing a grid of vector lanes. Each lane is labeled "Vector Lane" and contains "XMX" units. A "Load/Store" unit is also shown at the bottom of the grid.

52TF
Peak FP64
Throughput

839TF
Peak BF16
Throughput

最大
128 MB
HBM2e
Memory

16
Xe Links

A diagram of the Xe Link architecture showing a "Bridge" connected to a "Switch", which is connected to a "Link Exp. Sub-layer" and a "Controller".

976 GB/s
GPU-to-GPU comms
via Xe Links



100+ HPC Apps Running



標準ベンチ	D,S,H,I,BF-GEMM	DAOS	Graph500	HPCG	HPL	IO500	MLPerf 2.0
	MLPerf.HPC	OSU	SPECchpc	Stream Triad	MLBench	SPEC ACCEL	SpMV
	RINF	ELPA	Ginkgo	HeFFTe	HYPRE	MFIX (AMReX)	
物理	HACC	DPEcho	GENE	NEKBONE	nekRS	OpenMC	
	XGC	CloverLeaf	Deep Galaxy	Gadget	GRID QCD	MILC	
	QUDA	Chroma	HotQCD	BQCD	CERN 3D GAN		
ライフサイエンス	Autodock-GPU	miniBUDE	AMBER	GROMACS	LAMMPS	NAMD	OpenMM
	Relion	Quantum Espresso	BerkeleyGW	CP2K	NWChem	QMCPACK	DeePMD
製造	ANSYS CoMLSim	Jacobi Solver	Commercial EDA ISV	Commercial CFD ISV			
	ANSYS ParSeNet (SplineNet)	Commercial Multi 物理 ISV	Commercial CFD ISV	Proprietary CFD code			
金融	Binomial Options	Black-Scholes	European Monte Carlo	American Monte Carlo	Riskfuel Risk Calculations	STAC-A2	
地球システム	SPECFEM3D GLOBE	ES3M/MMF	SeisSol				
エネルギー	RTM Stencil Kernel	ISO3DFD					

AI モデルサポート



Computer Vision Image Classification	ResNet-50 v15	ResNeXt-101	ResNet-101	EfficientNet-B7	SE-ResNeXt50	TSM	
	Adorym	CosmoFlow	RegNetY-32Y	ResNeXt-101	Candle Uno	Swin Transformer	
画像セグメンテーション	Cosmic Tagger	Mask R-CNN	DenseNet169	FFN	3D-Unet		
	PointNet	DeepCAM	DRN-D-54	ResNeXt3D-101			
物体検知	SSD-ResNet34	SSD-ResNet50	EfficientDet	ShuffleNet	YOLO-v3	YOLO-v4	RetinaNet-ResNet50
	Deep Fusion	CascadeRCNN-	MobileNet v3	SSD-MobileNet	MMA	ResNet101-FPN	
NLP 言語モデリング	BERT-Large	Stable Diffusion	ALBERT	FastFormers	Transformer-LT	Big Bird	Faster Transformer
	BERT-base	GP-J	BLOOM	DistilBERT	RoBERTa	XLNet	
音声認識	RNN-T	LAS - Listen Attend & Smell	Wave2Vec	QuartzNet			
音声合成	FastSpeech2	Tacotron-2 with LPCNet					
リコメンデーション	DLRM	DSSM	ESSM	Wide & Deep	DeepFM		
	DIN	AttRec	DIEN	MMOE			



平均

1.7倍の性能

Intel DC GPU Max Series 1550 vs. Nvidia A100 80G PCIe

1.0

Nvidia A100 80G PCIe

BabelStream_Triad

RINF

ISO3DFD

SpecFEM3D_Globe

American Monte Carlo

Binomial Options

Black-Scholes

European Monte Carlo

RiskFuel Training

Autodock

LAMMPS

miniBUDE

NAMD

CoML Sim Inference

CoML Sim Training

Jacobi Solver

3D-GANS for Particle...

BigDFT

CloverLeaf

DeepGalaxy

DPEcho

GENIE

GRIDQCD

標準
ベンチマーク

4.9

エネルギー

1.85

地球
システム
モデル

1.85

金融サービス

2.03

2.69

2.67

1.98

2.42

ライフとマテリアル
サイエンス

1.48

1.16

1.99

0.85

製造

1.78

2.25

1.72

物理

2.3

1.33

1.12

2.27

1.51

1.56

1.31

Relative performance (Higher is better)

See backup for workloads and configurations. Results may vary.



平均
1.3倍の性能

Nvidia H100 PCIe vs.
Intel Data Center GPU Max 1550

Nvidia H100 PCIe

BabelStream Triad

RINF

ISO3DFD

SpecFEM3D_Globe

American Monte Carlo

Binomial Options

Black-Scholes

European Monte Carlo

RiskFuel Training

Autodock

LAMMPS

miniBUE

NAMD

CoMLSim Inference

CoMLSim Training

Jacobi Solver

3D-GANS for Particle...

BigDFT

CloverLeaf

DeepGalaxy

DPEcho

GENE

GRIDQCD

標準
ベンチマーク

エネル
ギー

地球
システム
モデル

金融 サービス

ライフとマテリアル
サイエンス

製造

物理

3.8

1.12

1.55

0.98

1.89

1.16

1.58

1.17

1.51

1.45

0.92

1.31

0.82

1.86

1.83

1.45

1.54

1.13

0.91

1.77

1.2

1.35

0.92

Relative performance (Higher is better)

最小限の努力で 最適なソリューションを選択可能にする



Open

Ecosystem



Choice

Compatibility



Trust

Workloads



Scale

Delivery & Deployment



オープン、
マルチアーキテクチャ、
マルチベンダープログラミング

オープンな業界仕様

ハードウェア選択の自由度

性能、生産性、ポータビリティ

コミュニティ主導の標準化



ツールセットとしてイン
テルが実装

インテルハードウェアに最適化

確かな機能と性能

Fortran, Python, OpenMP, MPI...
サポート

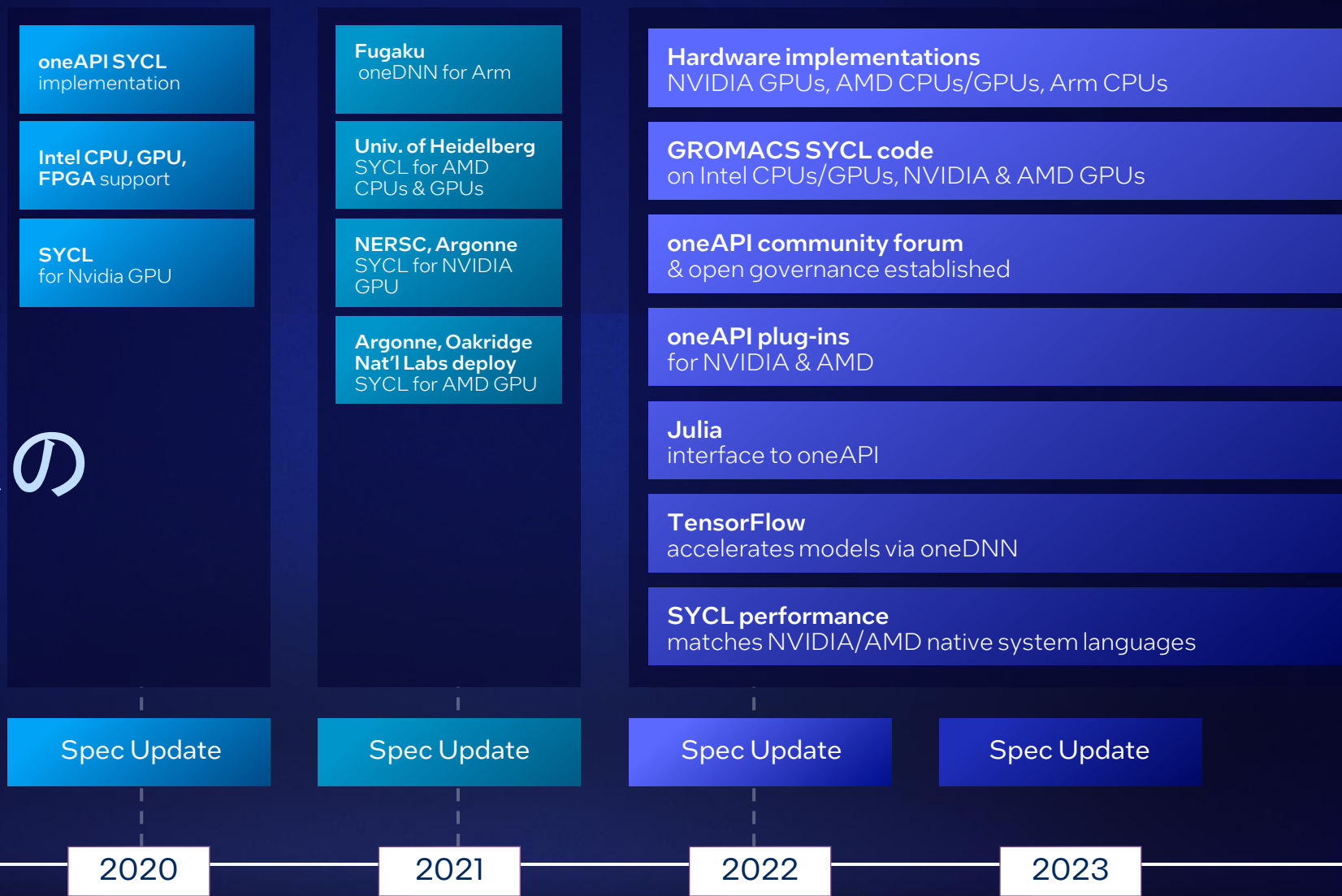
ダウンロードは無償、
商用サポートあり





オープン エコシステムの 進展

採用が堅調に進む





Intel oneAPI および AI ツール

インテルハードウェアの価値を最大化

新機能

Accelerated performance & AI
by enabling 第四世代 Intel® Xeon® CPUs & Max Series CPUs & GPUs

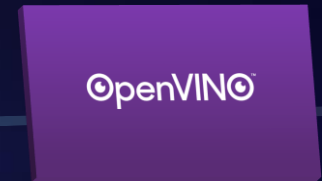
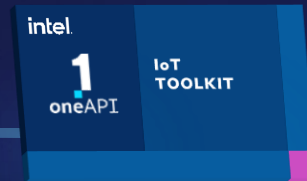
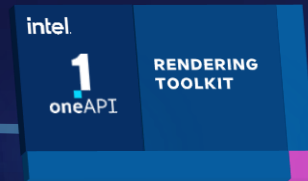
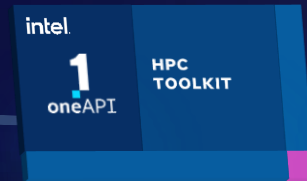
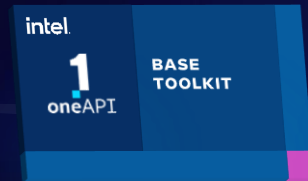
Speed up HPC applications
with OpenMP GPU offload & extended support for OpenMP & Fortran
Base & HPC toolkit

Real-time ray tracing
with hardware acceleration on Intel® GPUs & AI-based denoising in milliseconds*

Easier CUDA*-to-SYCL* migration
with improved tools
Base Kit

oneAPI

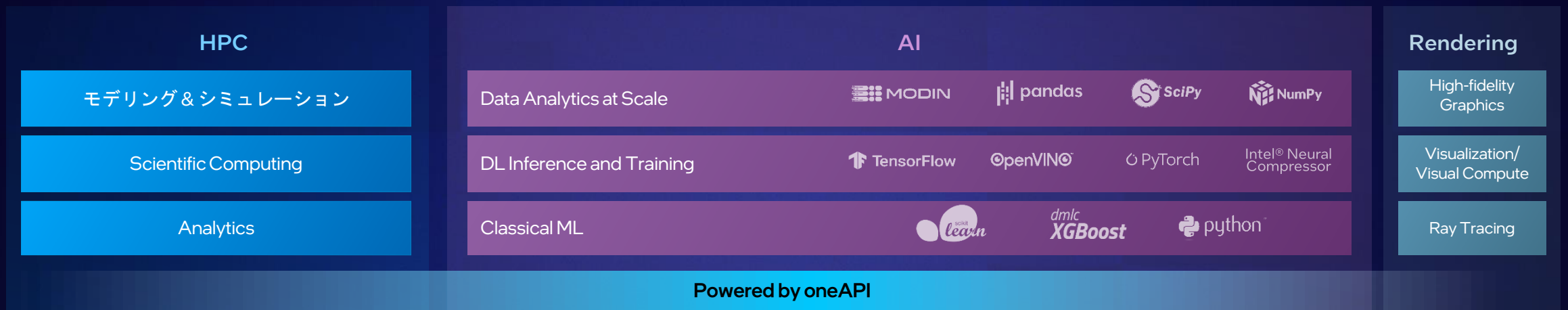
Powered by oneAPI



*avail. in open source Intel® Embree & Intel® Open Image Denoise, they will be in the Render Kit's next release

柔軟かつ包括的なオープンソフトウェアスタック

アプリケーション性能と生産性の高い開発環境によりインテルハードウェアの価値を最大化する



この一年の間に いろいろなことが起こりました...



AIとHPCアプリケーション性能をリード

強化された統一ソフトウェア層

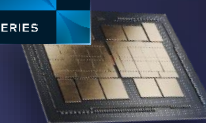
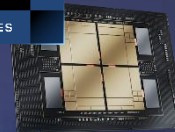
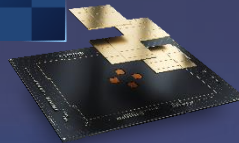
Scale

Open

Trusted

Choice

すべてのHPCとAIニーズに対応する製品群



The Intel logo is centered on a dark blue background. It features the word "intel" in a white, lowercase, sans-serif font. A small, light blue square is positioned above the letter "i". To the right of the word "intel" is a registered trademark symbol (®).

intel®

Notices and Disclaimers

Statements in this document that refer to future plans or expectations are forward-looking statements. These statements are based on current expectations and involve many risks and uncertainties that could cause actual results to differ materially from those expressed or implied in such statements. For more information on the factors that could cause actual results to differ materially, see our most recent earnings release and SEC filings at www.intc.com.

All product plans and roadmaps are subject to change without notice.

Performance varies by use, configuration and other factors. Learn more on the [Performance Index site](#). Intel technologies may require enabled hardware, software or service activation.

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 does not control or audit third-party data. You should consult other sources to evaluate accuracy.

Code names are used by Intel to identify products, technologies, or サービス that are in development and not publicly available. These are not "commercial" names and not intended to function as trademarks.

© 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.