

PCCC HPCオープンソースソフトウェア普及部会ワークショップ

「多様化するアクセラレータ」

# インテルXPU戦略

インテル株式会社

APJ データセンター・グループ・セールス

AI テクニカル・ソリューション・スペシャリスト

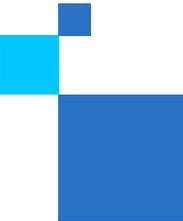
大内山 浩

インテル株式会社

プログラマブル・ソリューションズ営業本部

事業開発マネージャー

高藤 良史



intel®

# 注意事項および免責条項

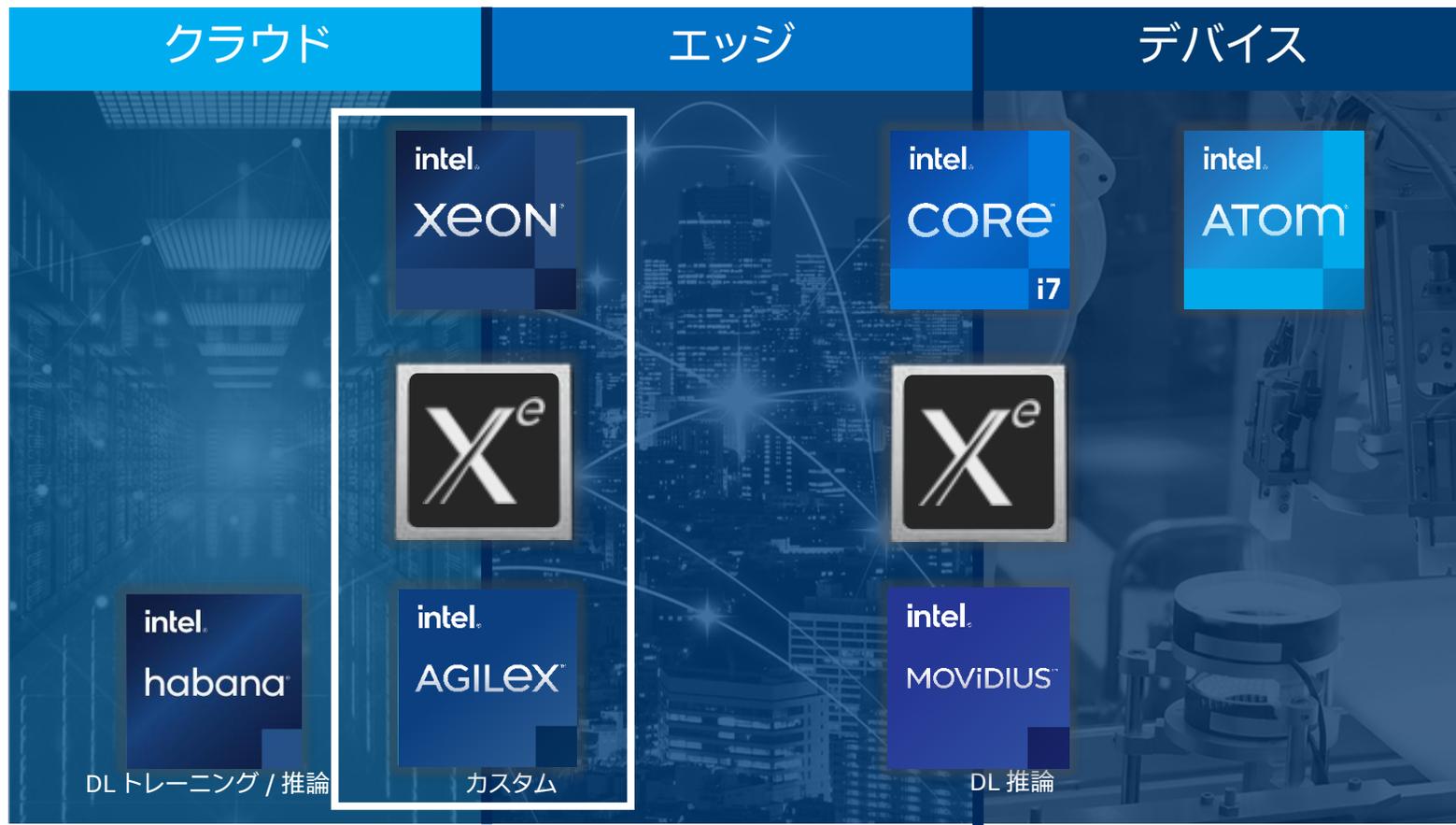
- 性能は、使用状況、構成、その他の要因によって異なります。詳細については、<http://www.Intel.com/PerformanceIndex/> (英語) を参照してください。性能の測定結果は、システム構成に記載された日付時点のテストに基づいています。また、現在公開中のすべての更新プログラムが適用されているとは限りません。構成の詳細については、<http://www.Intel.com/InnovationEventClaims/> (英語) を参照してください。絶対的なセキュリティを提供できる製品やコンポーネントはありません。
- すべての製品計画およびロードマップは、予告なく変更されることがあります。製品版出荷前のシステムやコンポーネントで測定された結果は、インテル・リファレンス・プラットフォーム (新しいシステムの社内サンプルモデル)、インテル社内の分析、アーキテクチャー・シミュレーション、モデリングでの推定 / シミュレーション結果を含め、情報提供のみを目的としています。システム、コンポーネント、仕様、構成に対する今後の変更によって、結果は異なる場合があります。インテルのテクノロジーを使用するには、対応したハードウェア、ソフトウェア、またはサービスの有効化が必要となる場合があります。
- 本資料は、明示されているか否かにかかわらず、また禁反言によるとよらずにかかわらず、いかなる知的財産権のライセンスも許諾するものではありません。開発コード名は、一般向けに発表または出荷されていない製品やテクノロジー、サービスを識別するためにインテルによって使用されているものです。いずれも「商用」の名称ではなく、商標としての機能を前提としたものではありません。
- インテルは、Principled Technologies が統括する BenchmarkXPRT 開発コミュニティを含め、さまざまなベンチマーク制定団体 / 機関へのスポンサーとしての参画、また技術的サポートの提供によって、ベンチマークの開発に貢献しています。
- 将来的な計画や予測について言及している本プレゼンテーション資料内の記述は、多数のリスクや不確定要素を伴う将来の見通しです。「想定される」、「見込まれる」、「意図する」、「目標とする」、「計画する」、「考えられる」、「求める」、「推定する」、「継続する」、「可能性がある」、「予定である」、「期待する」、「はざである」、「仮定する」などの表現やその変化形および類似表現は、いずれも将来の見通しであることを示しています。推定、予想、予測、不確定な事象、または仮定について言及している、またはこれらに基づいている記述も、今後リリースされる製品やテクノロジー、こうした製品やテクノロジーに期待される利用可能性とメリット、市場機会、インテルの事業および関連市場に見込まれるトレンドに関する記述を含め、将来の見通しであることを示しています。経営陣による現在の予測に基づくものであり、多数のリスクや不確定要素を伴う将来の見通しです。これらの要因によって、実際の結果はこれらの予測的記述に明示的または黙示的に示された結果と著しく異なる可能性があります。実際の業績をインテルの予測と大きく異ならせる重要な要因には、インテルの投資家向け IR サイト (<https://www.intc.com/>) または証券取引委員会 (SEC) のウェブサイト (<https://www.sec.gov/>) で入手可能な Form 10-K および Form 10-Q に関するインテルの最新の報告書を含め、インテルが SEC に提出 / 登録した報告書に記載されています。インテルは、法律で開示が義務付けられている場合を除き、新しい情報、新規開発、その他の結果にかかわらず、本プレゼンテーション資料に記載されたいかなる記述も、更新する義務を一切負わないことを明示的に表明します。
- インテルは、サードパーティーのデータについて管理や監査を行っていません。ほかの情報も参考にしてデータの正確さを評価してください。
- ©2022 Intel Corporation. Intel、インテル、Intel ロゴ、その他のインテルの名称やロゴは、Intel Corporation またはその子会社の商標です。その他の社名、製品名などは、一般に各社の表示、商標または登録商標です。

# インテル® XPU ポートフォリオ

ここから内蔵 AI  
アクセラレーション  
を開始

AI、HPC、  
グラフィックス、  
リアルタイム・  
メディアが中心

ディープ  
ラーニング・  
ワークロードが  
中心



オープンな標準  
ベースの一体型  
ソフトウェア・  
スタック\*



# アクセラレーター

## 内蔵

or

## 外付け

### Integrated Accelerator

- システム構成がコンパクト。TCO最適化。
- CPUとメモリーを共有する構成
- CPUと同じパッケージの半導体に入れるので、設計サイズや使える電力に制限あり



### Discrete Accelerator

- システム構成は大きい電力消費も上がる
- CPUとはPCIeバスを介して接続する構成
- GPUが専用で使えるメモリーを搭載
- 独立した半導体パッケージのため、設計サイズや電力制限大幅に緩和

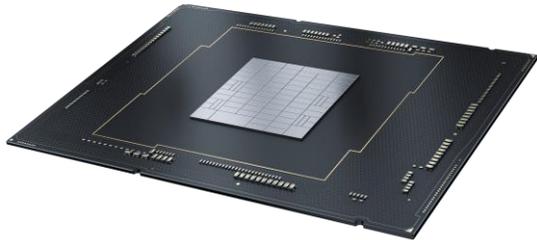


CPU

## 第3世代 Xeon スケーラブル・ プロセッサ

コードネーム: Ice Lake

AVX-512、DL Boost、SGX 等々..

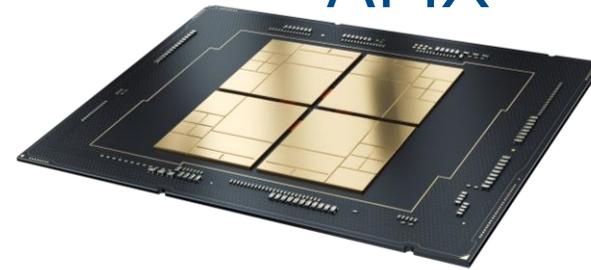


2021

## 次世代 Xeon スケーラブル・ プロセッサ

コードネーム: Sapphire Rapids

AMX



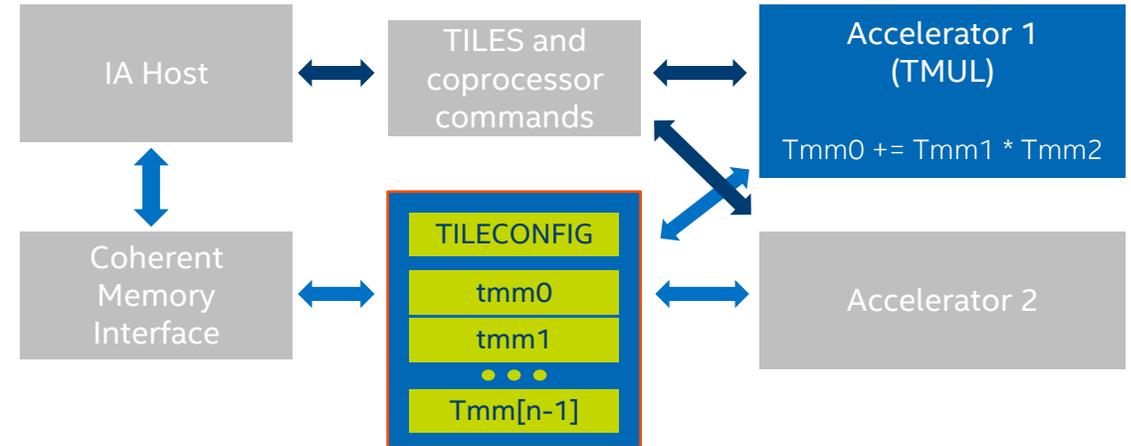
2022

# Advanced Matrix Extensions (AMX)

## Sapphire Rapidsから搭載される AI アクセラレーション・エンジン

- **AMXはディープラーニングの推論と学習性能を飛躍的に向上**
- **サポートされるデータ型**
  - INT8
  - BF16
- **ソフトウェア・サポート**
  - インテル® oneAPI DNNL
  - インテルが最適化したDL フレームワーク
    - TensorFlow\*、Pytorch\*、OpenVINO™ ツールキット、MXNet\* など

### AMX 概要イメージ



- New State to be managed by OS
- Commands and status delivered synchronously via TILE/accelerator instructions
- Dataflow – accelerators communicate to host through memory

### レジスター

- 複数の TILE と呼ばれる 2 次元レジスターファイルにより構成
- 1 つの TILE は 16rows x 64-byte (1KB) なので、8TILE で 8KB
- TILE のサイズは変更可能

### アクセラレーター

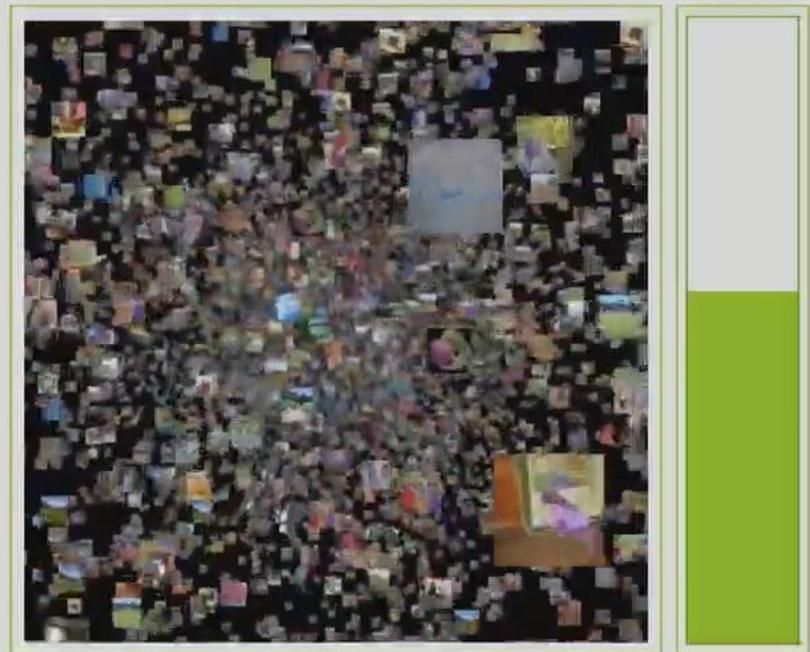
- TILE のデータを操作するオペレーター
- 現時点では Tile matrix multiply unit (TMUL) のみが定義
- TMUL は INT8 と BF16 に対応

+

AMXドキュメント: <https://software.intel.com/content/dam/develop/public/us/en/documents/architecture-instruction-set-extensions-programming-reference.pdf>

# Productivity & Performance

Resnet-50 v1.5 (Batchsize =504)



TensorFlow

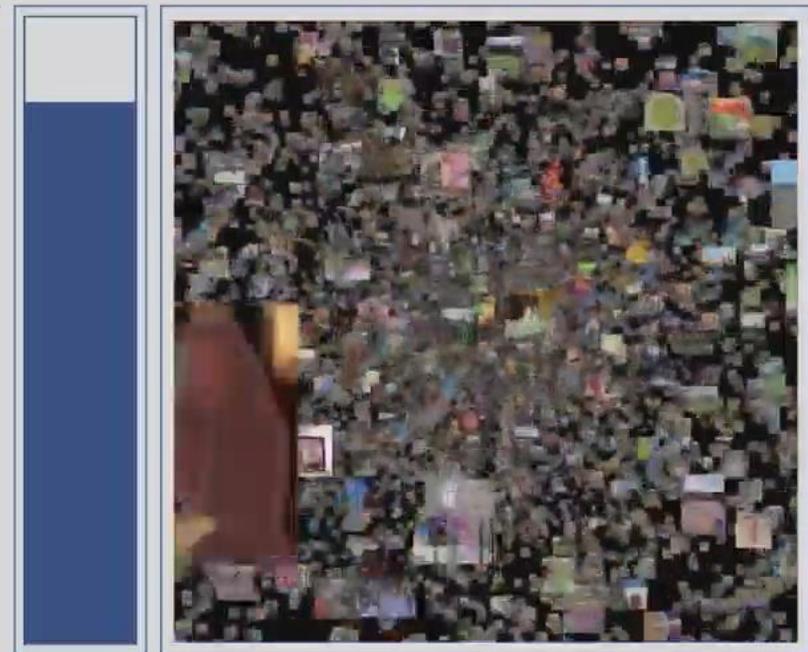
15,677 Img/sec

Competitor

1.54x

Speed Up

AMXの活用により他社製品を上回る性能を実現



24,085 Img/sec

TensorFlow

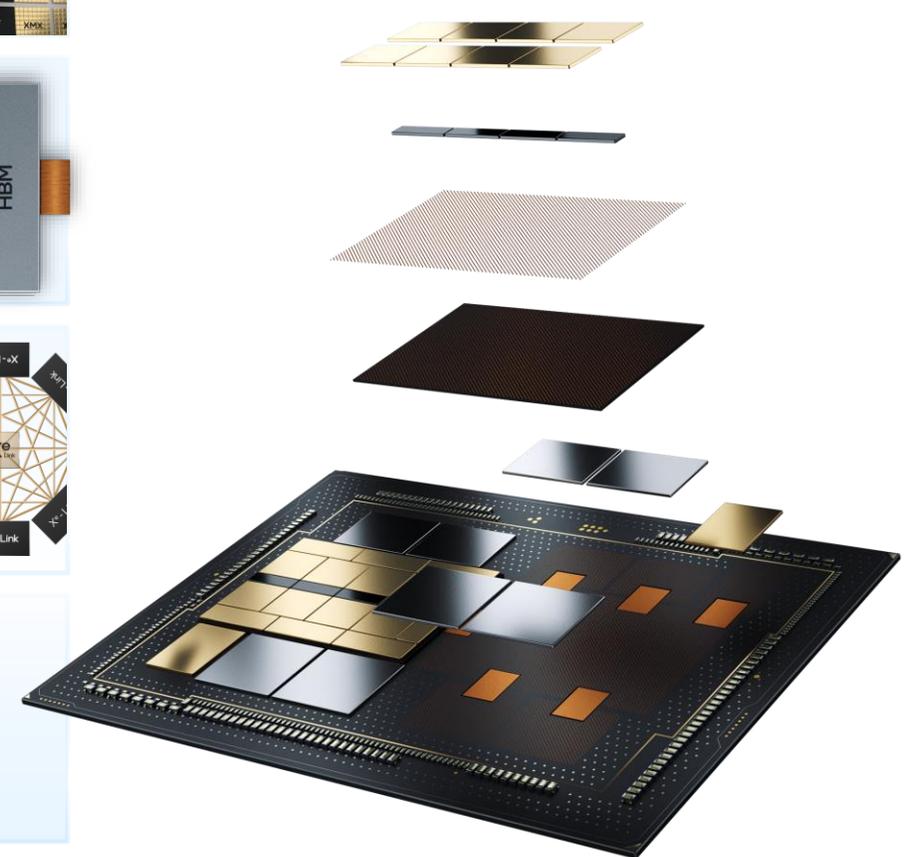
2S Next Gen Xeon Scalable Processor

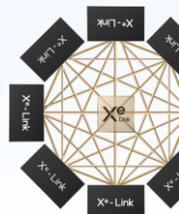
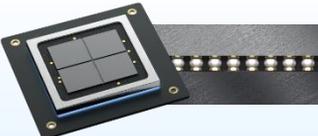
詳細はこちらへ↓

<https://newsroom.intel.la/news/intel-innovation-artificial-intelligence-newsenglish-only>

# Ponte Vecchio

X<sup>e</sup> HPC based GPU



<b>Compute</b>	最大 <b>128</b> Ray tracing Units	Highest Compute <b>Density</b> socket & node	<b>128 X<sup>e</sup> Cores</b> 
<b>Memory</b>	最大 <b>64MB</b> L1 cache in 2 Stacks	最大 <b>408MB</b> L2 Cache in 2 Stacks	<b>HBM2e</b> 
<b>I/O</b>	最大 <b>8</b> Fully Connected GPUs	<b>PCIe Gen 5</b>	<b>X<sup>e</sup> Link</b> High-Speed Coherent Unified Fabric 
<b>Technology</b>	 <b>EMIB</b>	 <b>Foveros</b>	Intel 7 TSMC N5 TSMC N7

# Ponte Vecchio

Execution Progress

## A0 Silicon Current Status

**> 45 TFLOPS**

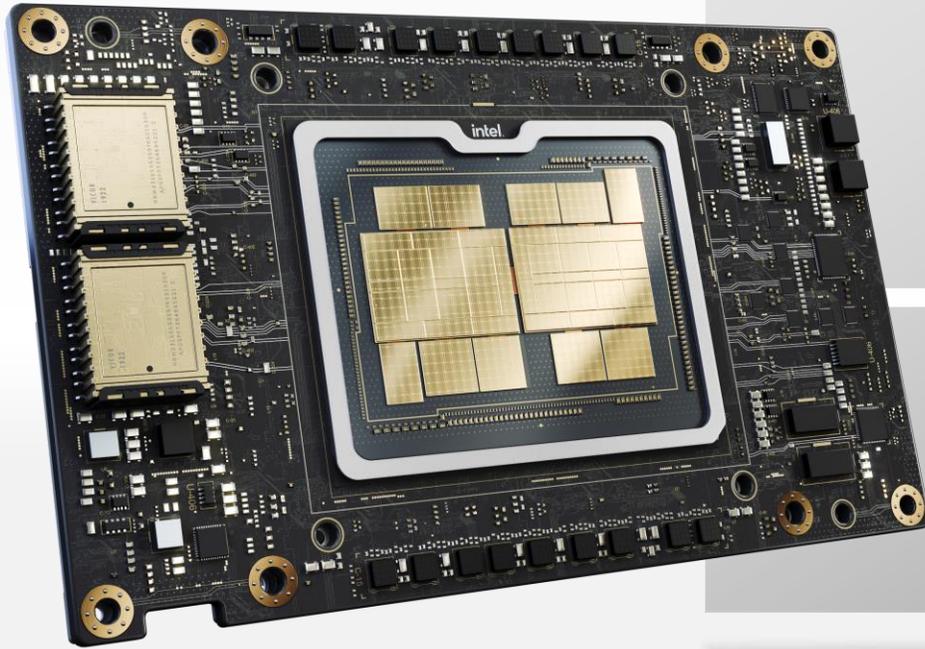
FP32 Throughput

**> 5 TBps**

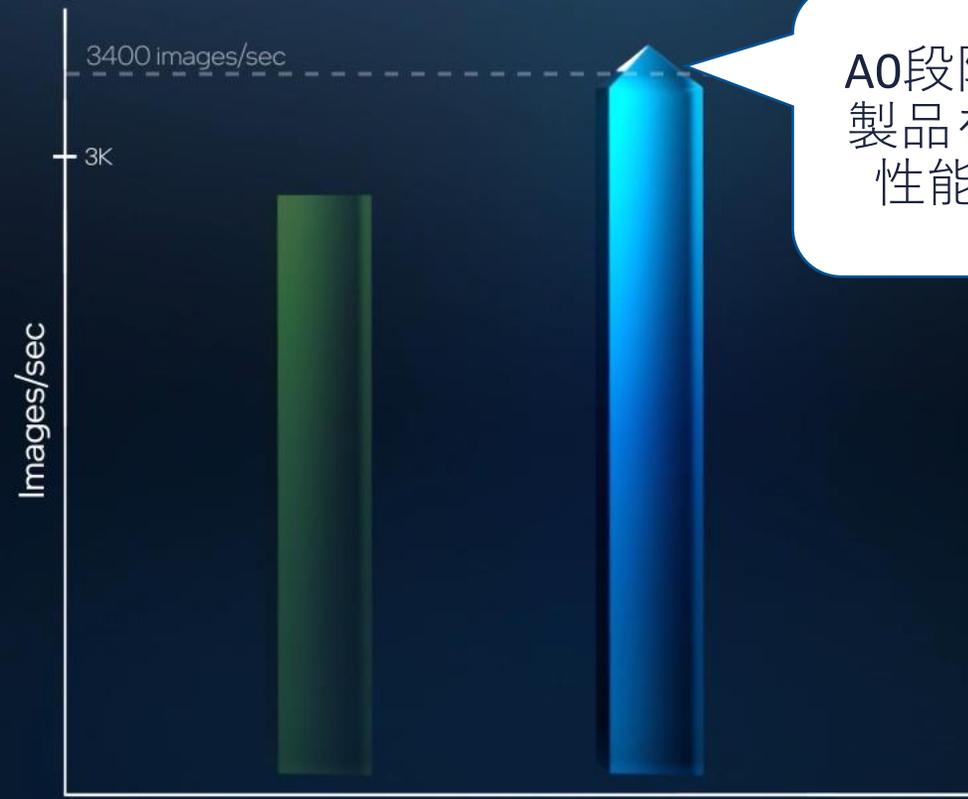
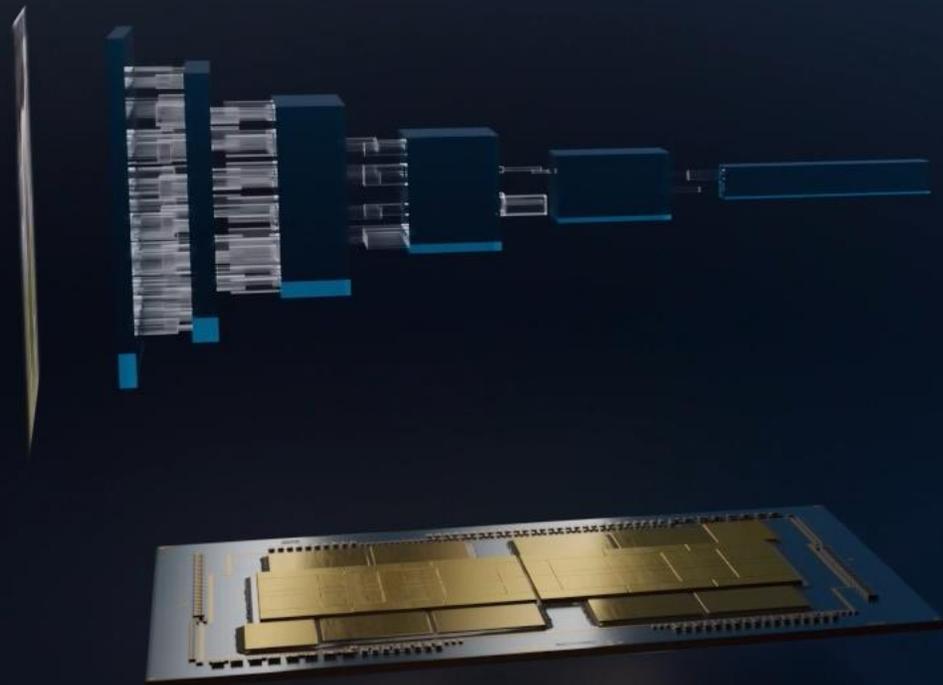
Memory Fabric Bandwidth

**> 2 TBps**

Connectivity Bandwidth



# ResNet-50 v1.5 Training (Single GPU)

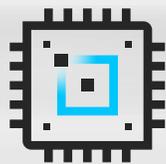
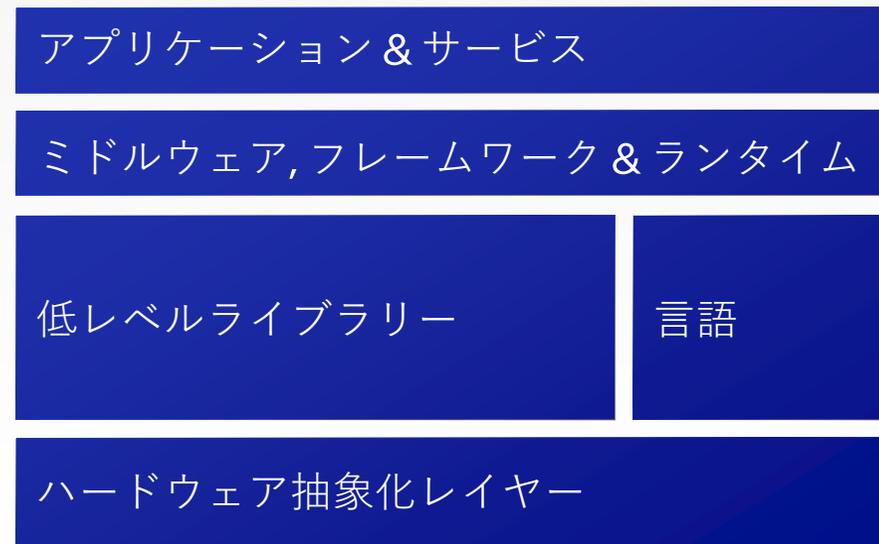
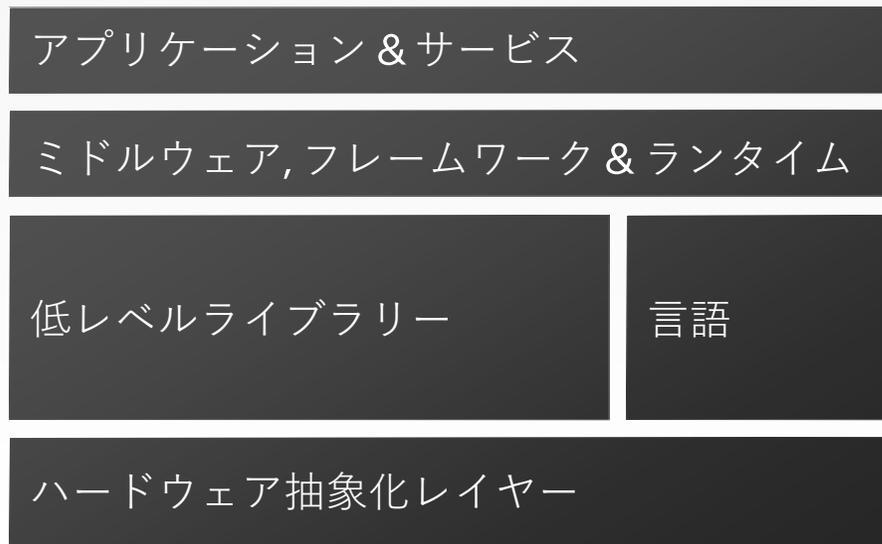


A0段階で他社製品を上回る性能を達成

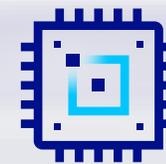
Competition Ponte Vecchio  
Early A0 Si Measurements

# CPUとGPUのソフトウェアスタック分離を克服するために

CPUに最適化されたスタック ← → GPUに最適化されたスタック



CPU



GPU



オープンな標準規格に基づく  
統合ソフトウェアスタック

独自のプログラミングモデルからの解放

ハードウェアの性能をフルに発揮

開発者の安心感

## CPU & XPU – 最適化されたスタック

アプリケーション & サービス

ミドルウェア, フレームワーク & ランタイム

TensorFlow PyTorch mxnet *learn* NumPy *dmlc XGBoost* OpenVIN® ...

### 低レベルライブラリー

oneMKL

oneDNN

oneDAL

oneVPL

oneTBB

oneCCL

oneDPL

Other  
Libraries

### 言語

DPC++

他の言語

ハードウェア抽象化層

Level Zero

演算ハードウェア



CPU



GPU

**>200K** 開発者

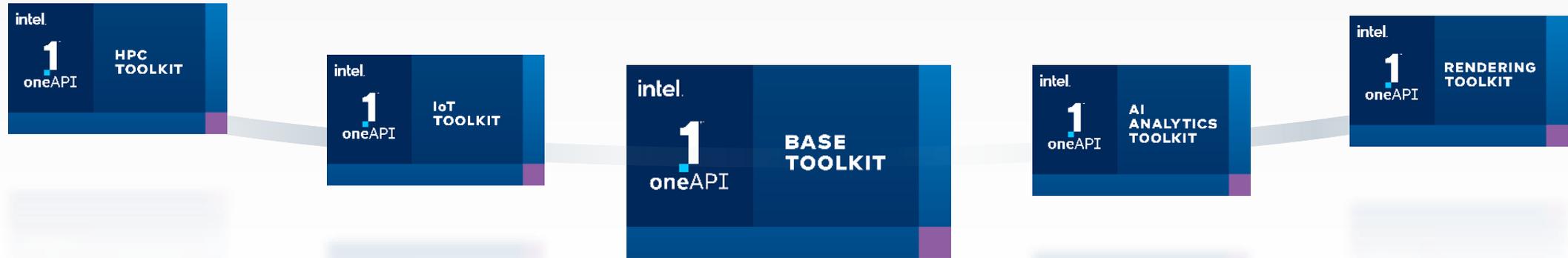
2020年12月のリリース以来、インテル®  
oneAPI 製品のユニークなインストール数

**>300** アプリケーション

インテル® oneAPI 言語とライブラリを使用  
した市場展開

**>80** HPC & AI アプリケーション

インテル® oneAPIを使用したインテル®  
Xe HPCでの機能性



**1**  
**oneAPI**

# 業界を牽引する勢い

## End Users



## National Labs



## ISVs & OSVs



## OEMs & SIs



## Universities & Research Institutes



## CSPs & Frameworks





# Aurora Blade

Building Block for the ExaScale Supercomputer

1  
oneAPI

Argonne   
NATIONAL LABORATORY

 U.S. DEPARTMENT OF  
**ENERGY**

  
Hewlett Packard  
Enterprise

intel.

# インテル® FPGAデバイス 最新技術情報



intel®

# インテル® FPGA

## 急激に変化する世界へフレキシビリティを提供



データ・セントリック  
時代へ



最適化された  
バンド幅



ミッドレンジ  
FPGA & SoC



エッジ・セントリック  
FPGA



# FPGAとは？

- 動的に再プログラム可能なハードウェア回路を形成するシリコンデバイス
- 様々なワークロードに対応可能なデータパスを備え、処理速度が高く電力効率に優れた、低レイテンシーのサービスを提供

データ分析

映像処理

ネットワーク

符号化処理

組み込み・アプライアンス向けのチップ製品

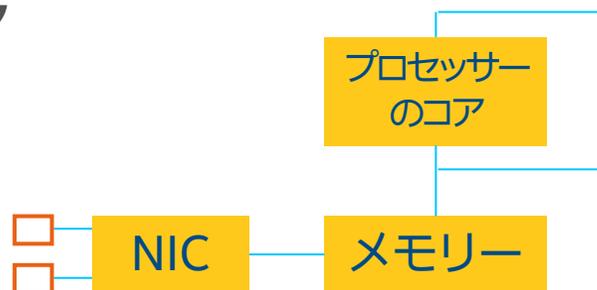


コンピューティング向けのカード製品



# FPGA: 低遅延・確定的レイテンシー

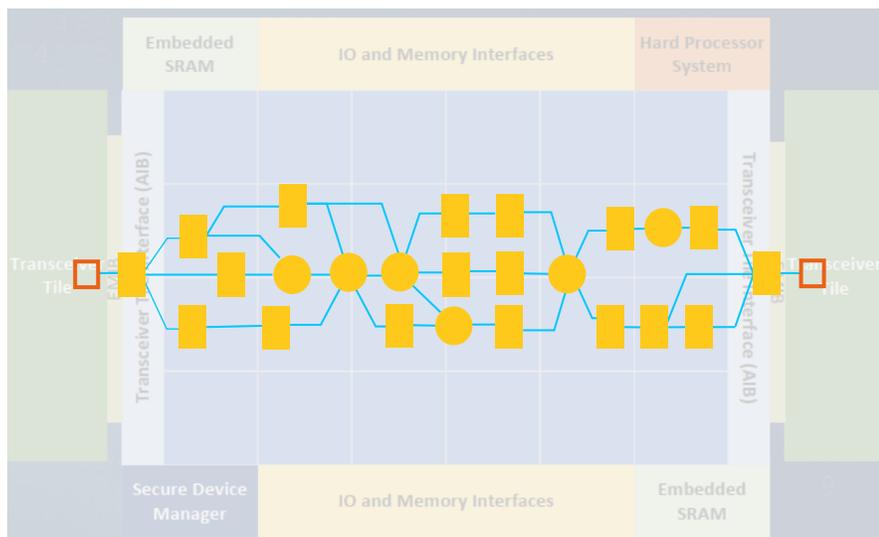
## ソフトウェア



## シーケンシャル・アーキテクチャー

- NIC とメモリー間でパケットを転送
- プロセッサのコアが順次処理

## FPGA



## 空間アーキテクチャー

- プロセッシング・パイプラインへデータを直接フィード
- 複数ステージでデータを同時処理

# ヘテロジニアス・コンピューティングでの FPGA の役割

FPGA はヘテロジニアス・システムにおいて他のテクノロジーの補完的役割

- 他のアーキテクチャーの性能上のボトルネックにアプローチ

## FPGA のターゲット

リアルタイム性が要求される処理

データが再利用される繰り返し処理

混合精度やベクターが使われる処理

サブブロックと I/O の間のストリーミング・データ処理

演算密度や電力効率を考慮すべき処理

# HPC向けFPGAカード製品

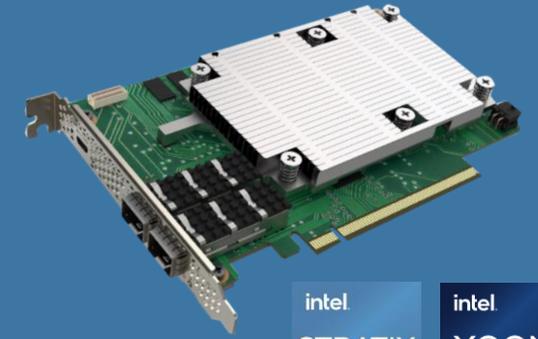
ワークロードオフロード・ネットワーク・ストレージなど用途に応じたカードをご提案

## インテル® FPGA SmartNIC N6000 プラットフォーム



- QSF56 x2 インターフェイスで最大 100GbE x2 のネットワーキングをサポート
  - SyncE、CPRI、eCPRI をサポート
  - O-RAN LLS-C1/-C2/-C3 のサポート
- 4G & 5G vRAN  
イネーブルメント・パッケージ
- OVS、Contrail、SRv6、vFW アクセラレーションのサポート

## インテル® FPGA IPU C5000X-PL (OVS, ストレージ)



- インテル® Xeon® D プロセッサ + FPGA プラットフォーム、ハードウェア・プログラマブル・データパスを提供
- ストレージ / 仮想スイッチのワークロードオフロード

## BittWare IA-840F インテル® Agilex™ F シリーズ FPGA 搭載 FPGA アクセラレーター PCIe カード



- インテル® Agilex™ AGF027 FPGA (2.6M LE)
- インテル® oneAPI サポート
- ハイパフォーマンス I/O
  - 3x QSFP-DD インターフェイスポート
- PCIe Gen4 x16 ホストインターフェイス
- さまざまなアプリケーション向けのMCIO 拡張ポート

## Silicom FPGA SmartNIC N5010 (HBM搭載, ネットワークオフロード)



インテル® イーサネット

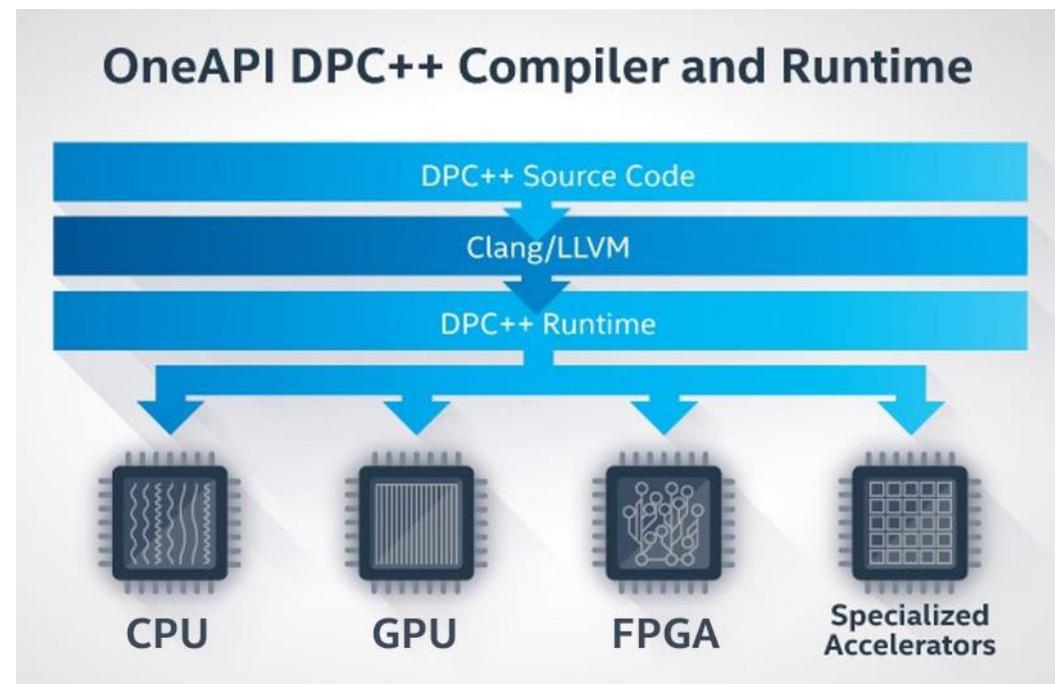
- ハードウェア・プログラマブル 4x100GE FPGA アクセラレーションを実装する初の SmartNIC
- 5G コア・ネットワーク (UPF)、アクセス・ゲートウェイ (BNG、AGF)、セキュリティ機能 (ファイアウォール、IPsec) のパフォーマンスとスケーリングのニーズに対応。
- 機能 / アクセス・ゲートウェイ機能、その他のワークロード



# インテル® oneAPI DPC++ コンパイラー

## 並列プログラミングの生産性とパフォーマンス

- CPUとアクセラレーターの両方で並列プログラミングの生産性を引き出すことが可能
  - 単一アーキテクチャーの独自言語に代わる、業界間で共通のオープン言語
- ベースは最新規格のC++とSYCL\*
  - 広く採用されている慣れ親しんだCやC++の構造により、C++の生産性メリットを活用
- アーキテクチャーと高性能コンパイラーにおけるインテルの数十年にわたる経験を基に構築



# インフラストラクチャー・プロセッシング・ユニット (IPU)

データセンターの新たな価値を提供



高度にインテリジェントな**インフラストラクチャー・アクセラレーション**

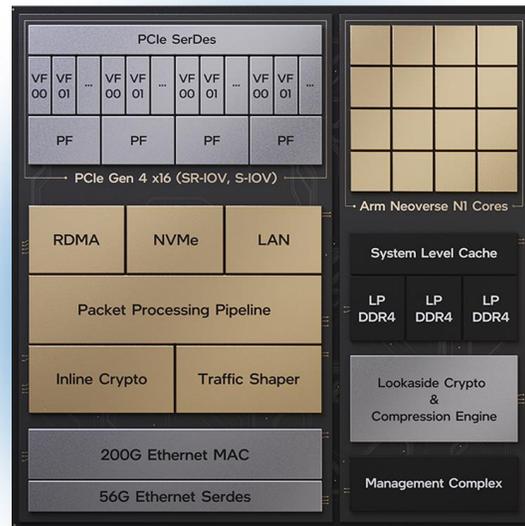
システムレベルの**セキュリティ、制御、分離**

**共通の**ソフトウェア・フレームワーク

ハードウェア / ソフトウェア・**プログラマブル**、顧客のニーズに合わせて構築

# インフラストラクチャー・プロセッシング・ユニット (IPU)

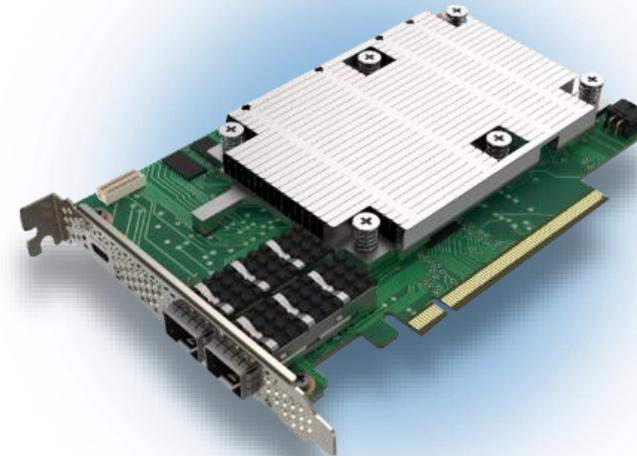
## ASIC IPU SoC



Mount Evans<sup>†</sup>

## FPGA IPU

プラットフォームとアダプター



インテル® FPGA IPU  
C5000X-PL

現在提供中

Oak Springs Canyon<sup>†</sup>

インテル® FPGA IPU C6001X-PL

<sup>†</sup>開発コード名

intel.

# インテル® XPU ポートフォリオ

ここから内蔵 AI  
アクセラレーション  
を開始

AI、HPC、  
グラフィックス、  
リアルタイム・  
メディアが中心

ディープ  
ラーニング・  
ワークロードが  
中心



オープンな標準  
ベースの一体型  
ソフトウェア・  
スタック\*



Thank You



intel®