

2020年9月29日

# Intel AIのこれまでとこれから

インテル株式会社

APJデータセンター・グループ・セールス

大内山 浩



intel®

# AI “も” 動かすCPU

- あらゆるワークロードに対応できる汎用性と柔軟性がCPUの特徴

統合ワークフロー: sas, SAP, Microsoft, IBM, ORACLE, amazon web services, Google Cloud, TERADATA, KubeFlow, Spark, ANACONDA, IOT, DOMINO, インテル® DL Studio など...

**収集、統合、ETL、ELT**  
オープン, オープン (管理対象), 自社開発  
kafka, pentaho, talend, TIBCO Jaspersoft, INFORMATICA など...

**メタデータの管理**  
collibra, Waterline Data, Blue River ANALYTICS, ZALONI など...

**ディープラーニング** (Red Circle)  
TensorFlow, mxnet, Caffe, Spark, BIGDL, PaddlePaddle, PYTORCH, ONNX など...

**推論の導入**  
OpenVINO, plaidML など...

**データの保存と管理**  
オープン, オープン (管理対象)  
HBASE, cloudera, MAPR, mongoDB, Flink, HORTONWORKS, QUbole, STORM, Lightbend, 自社開発: Paxata など...

**データの前処理**  
DataRobot, Alation, Datameer, Lavastorm, DATAWATCH, Paxata, unifi, DataKitchen, TRIFACTA, panoply, tamr, alteryx, ClearStory, composable, Lore IO など...

**マシンラーニングとアナリティクス**  
pandas, Spark, MLib, H2O.ai, presto, XGBoost, NumPy, AMATLAB, CognitiveScale, feedzai, Amenity Analytics, ARCADIA DATA, gamalon, DataRobot, sentient, FIS, alteryx, Palantir, ATASCIENCE.COM, KNIME など...

**視覚化**  
ggplot2, mld3, IP[y]: IPython Interactive Computing, Bkeh, kibana, matplotlib, Gephi, Grafana, Data-Driven Documents, tableau, TIBCO Spotfire, Qlik など...

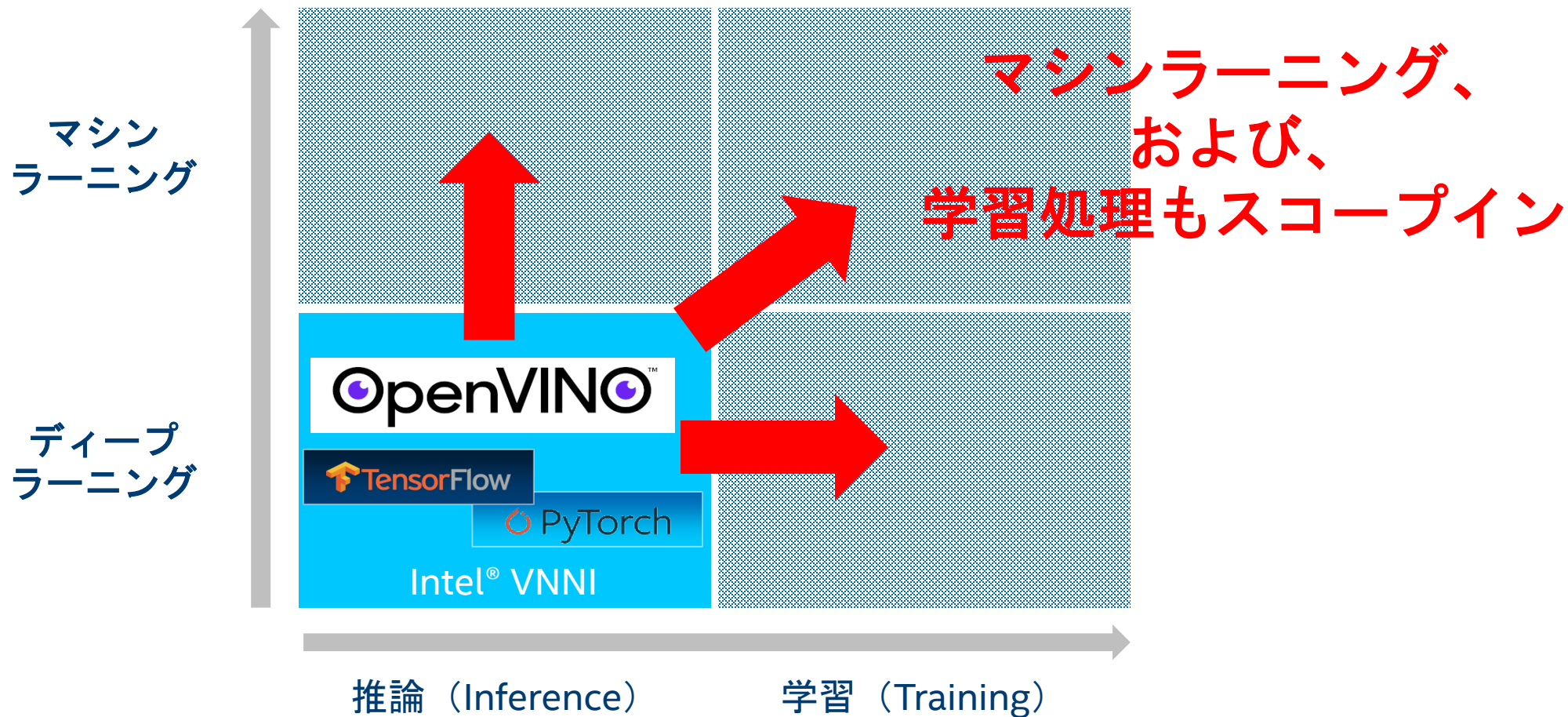
**IT システム管理**  
Vagrant, puppet, CHEF, STRATOSCALE, katacontainers, docker, APCERA, MESOS, Jira Software, XEN, MIRANTIS, CLOUDFOUNDRY, New Relic, Jenkins, bluedata など...

**API**

**エンタープライズ・アプリケーション**

Intel Core i7 inside, Intel Xeon Platinum inside

# ここ最近のインテルAI ～より広いAIワークロードに対応～



# ディープラーニング高速化の要 ～AVX-512 & DL Boost～

- AVX-512 (SIMD) が搭載され、並列演算性能の向上に寄与しております。更に、Deep Learning Boostという専用命令により更なるアクセラレーションが期待できます。

## Intel® AVX-512

(Intel® Advanced Vector Extensions 512)

+

## Intel® DL Boost

(Intel® Deep Learning Boost)

はいつてる



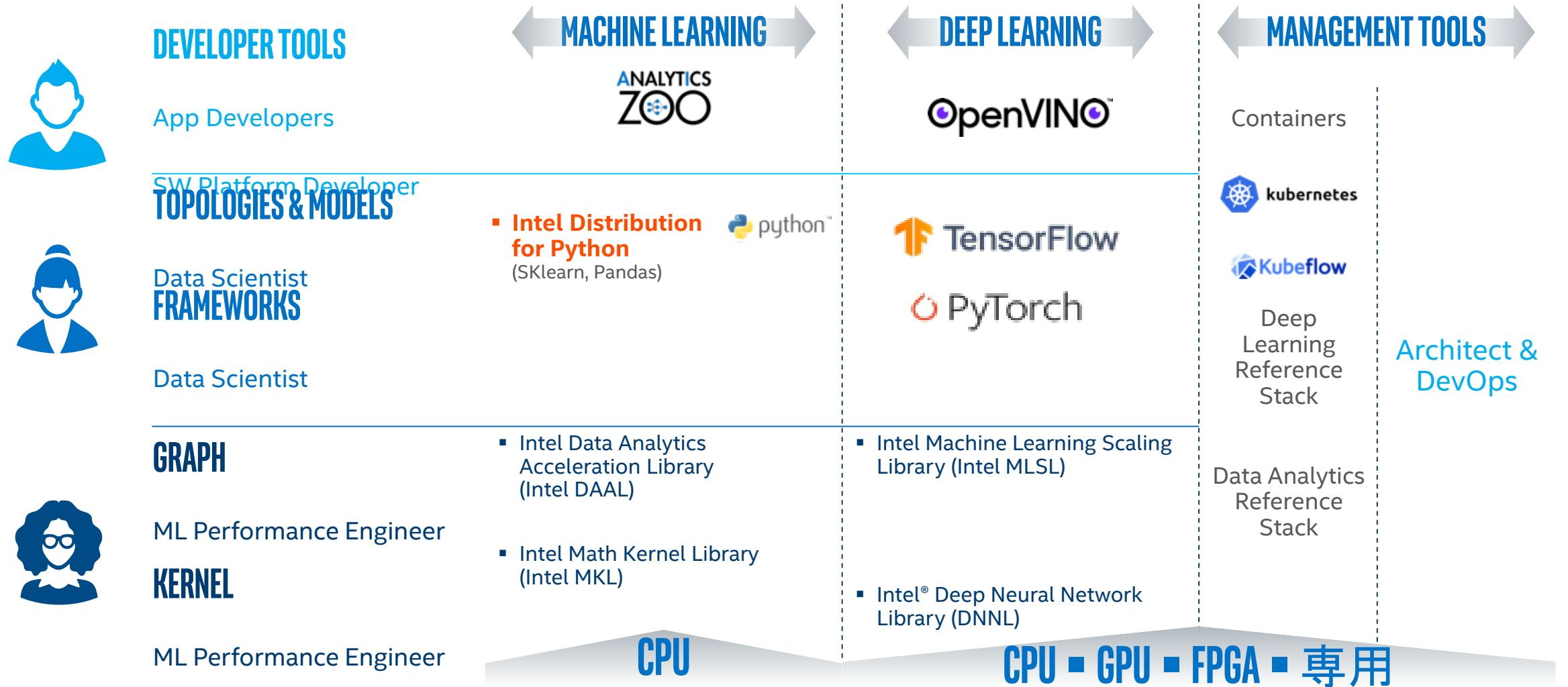
第10世代Ice Lakeから搭載

はいつてる



Skylake世代からAVX-512搭載  
Cascade Lake世代からDL Boost搭載

# インテル® AI ソフトウェア: マシンラーニングとディープラーニング



Red font products are the most broadly applicable SW products for AI users

# インテルによるディープラーニング・フレームワークの最適化



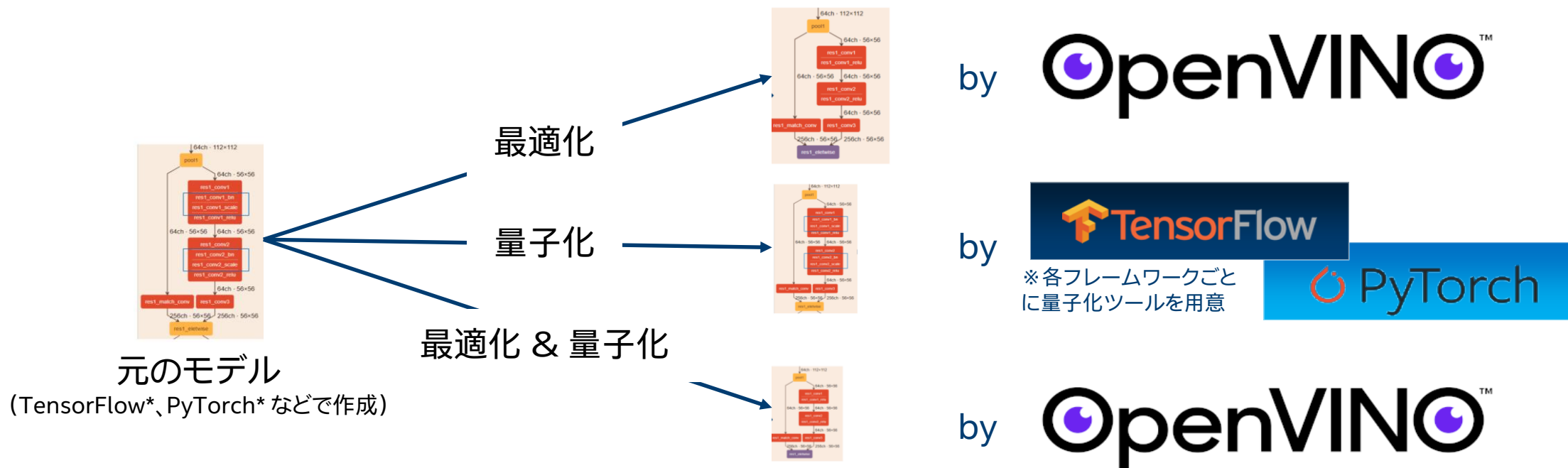
インストール・ガイドはこちら↓  
[ai.intel.com/framework-optimizations/](https://ai.intel.com/framework-optimizations/)

更なるフレームワークの最適化が進行中  
 (例、PaddlePaddle\*、CNTK\* など)

SEE ALSO: Machine Learning Libraries for Python (Scikit-learn, Pandas, NumPy), R (Cart, randomForest, e1071), Distributed (MLlib on Spark, Mahout)  
 \*Limited availability today  
 Optimization Notice

# 推論処理の更なる性能向上のための ディープラーニング・モデルの最適化と量子化

- 最適化：不要な Ops の除去、複数の Ops の統合などによりモデルをスマート化
- 量子化\*：モデル内部の数値表現を FP32→INT8 に変換することでスリム化



\* 2020年5月現在、インテル® ディープラーニング・ブースト (VNNI) が搭載された 第 2 世代インテル® Xeon® スケーラブル・プロセッサ以降、第 10 世代 インテル® Core™ プロセッサ・ファミリー (Ice Lake† のみ)以降 にてより効力を発揮する

# ディープラーニング推論処理ベンチマーク

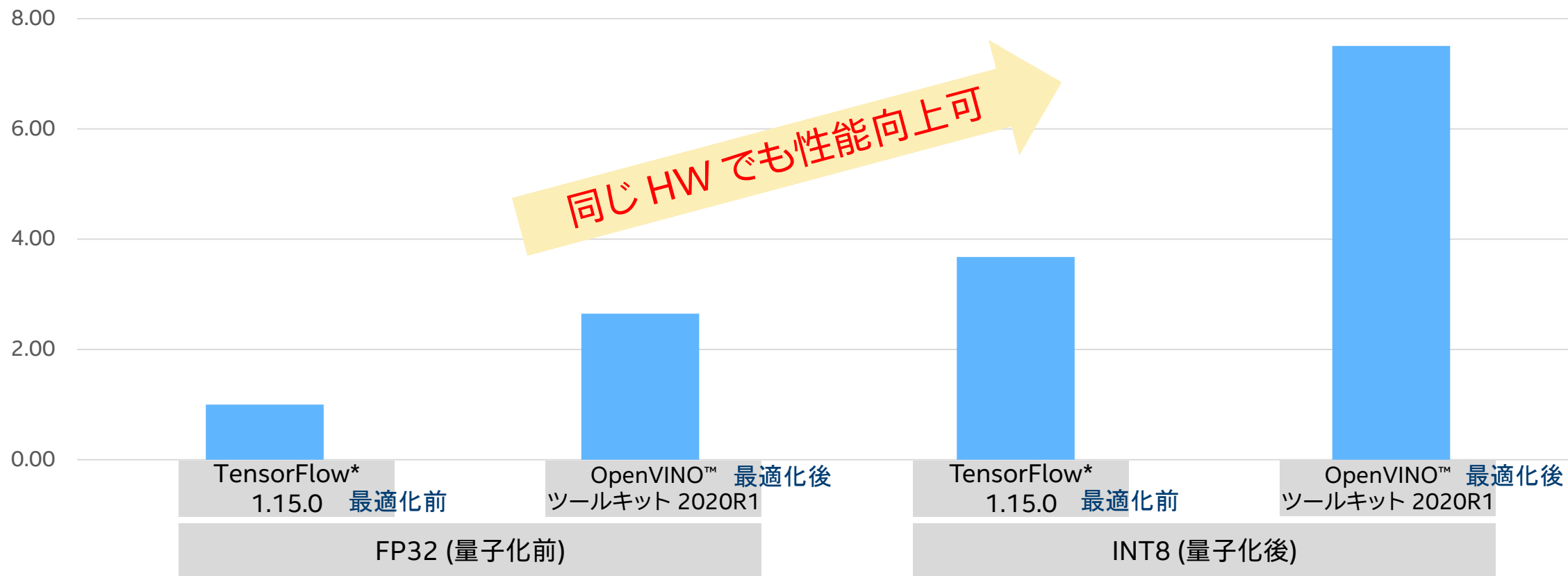
インテル® Xeon® Gold 6254 プロセッサ @ 2.10GHz (18 cores × 1 sockets)

性能比 (倍)

Resnet50 推論スループット (FPS)

2020年3月20日に計測

Input=224x224, BS=1, 1 stream

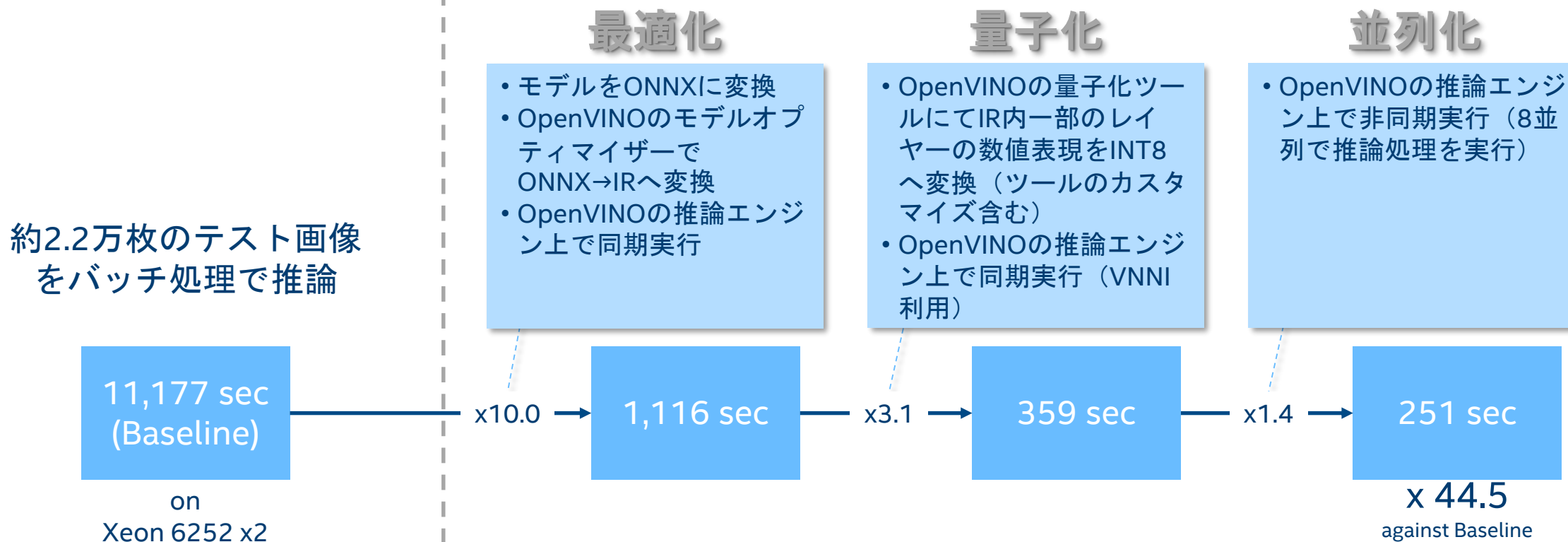


注) インテル社員による性能確認のための個人的なベンチマーク結果であり、インテルの公式結果ではありません。



# 事例：理化学研究所様 CheXNet の推論性能改善

Before Optimization ← → After Optimization



※オリジナルモデルは  
PyTorch 1.2.0にて実装

上記対応内容は下記Githubを参照  
<https://github.com/taneishi/CheXNet>  
 （計算科学研究機構 種石様のレポジトリ）

# Training with Huge Memory ~U-Net Training by NUS~



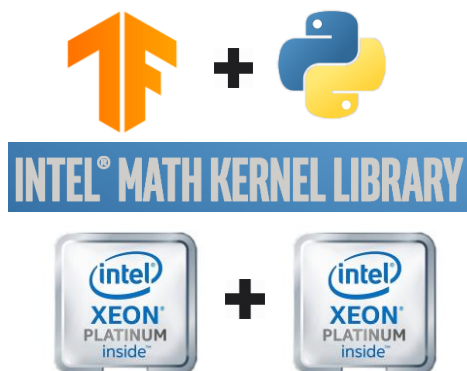
Saw Swee Hock  
School of Public Health

GPU-based env



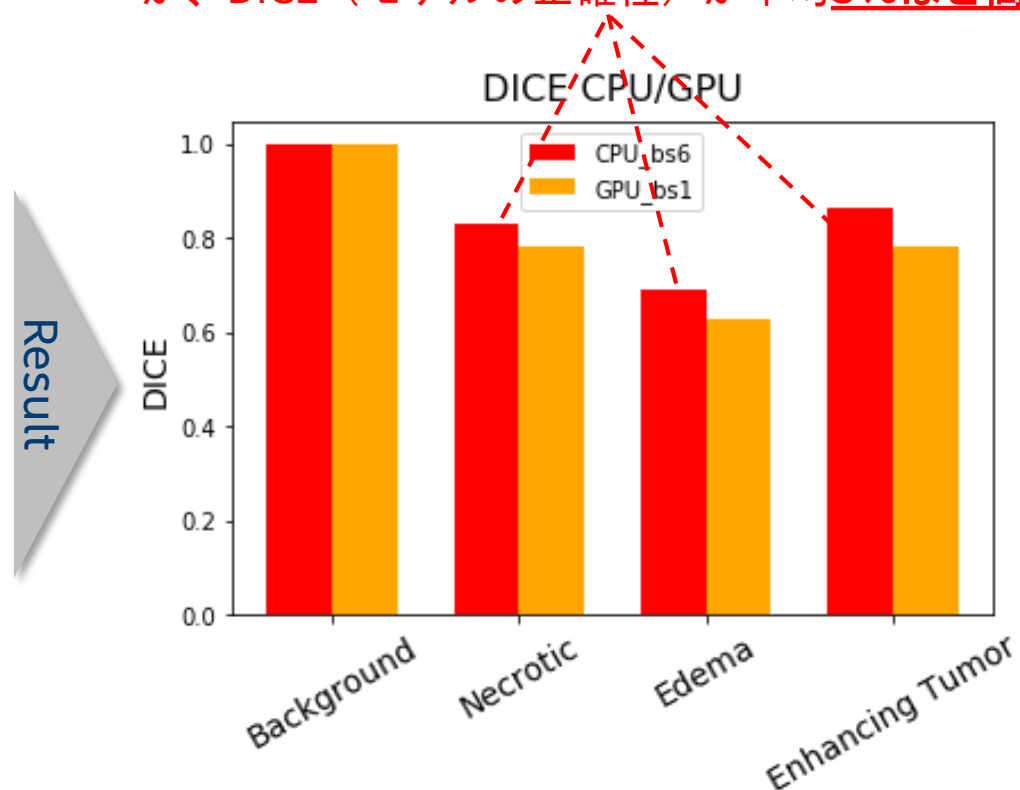
- V100 GPU (32GB memory)
- 10 CPU cores
- 126GB RAM
- Batch size of 1

CPU-based env



- 2 x Intel Platinum CPUs.
- 2 x 24 CPU cores
- 384GB RAM
- Batch size of 6

インテル® CPU上でトレーニングしたモデルの方が、DICE（モデルの正確性）が平均5%ほど高い。



パフォーマンスが欲しい場合はどうすればいい？



複数の CPU を束ねて使いましょう

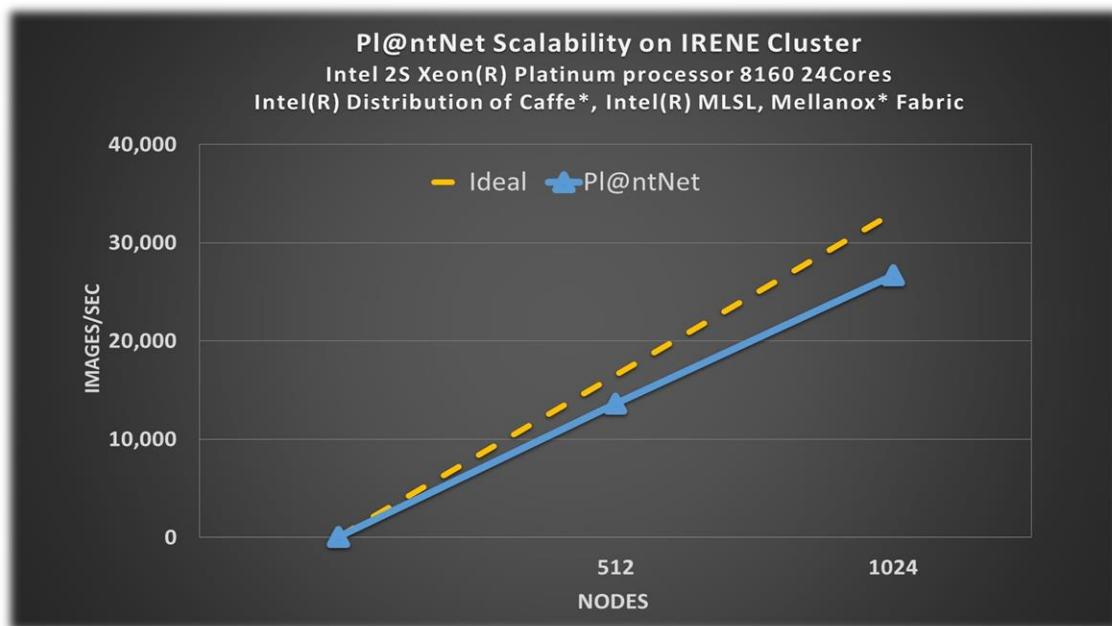
つまり、分散学習 (Distributed Training) です 

# 既存インフラ上で効率的な深層学習のスケール ~GENCI と CERN の事例~

## GENCI

French research institute focused on numerical simulation and HPC across all scientific and industrial fields

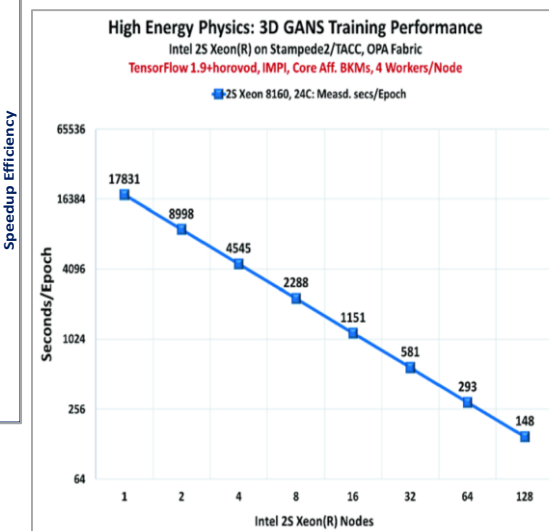
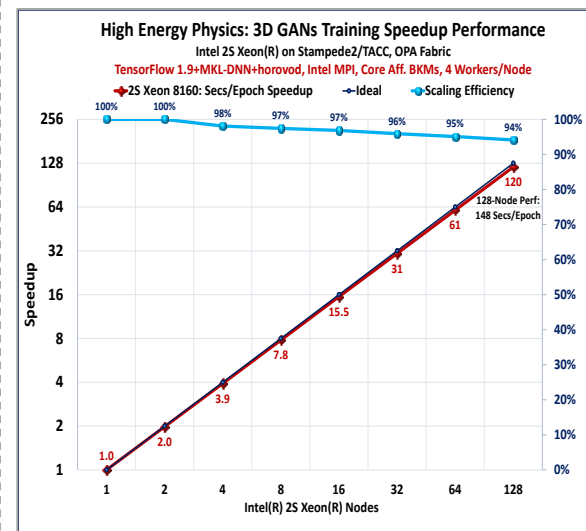
*Succeeded in training a plant classification model for 300K species, 1.5TByte dataset of 12 million images on 1024 2S Intel® Xeon® Nodes with Resnet50.*



## CERN

the European Organization for Nuclear Research, which operates the Large Hadron Collider (LHC), the world's largest particle accelerator

*94% scaling efficiency up to 128 nodes, with a significant reduction in training time per epoch for 3D-GANs*

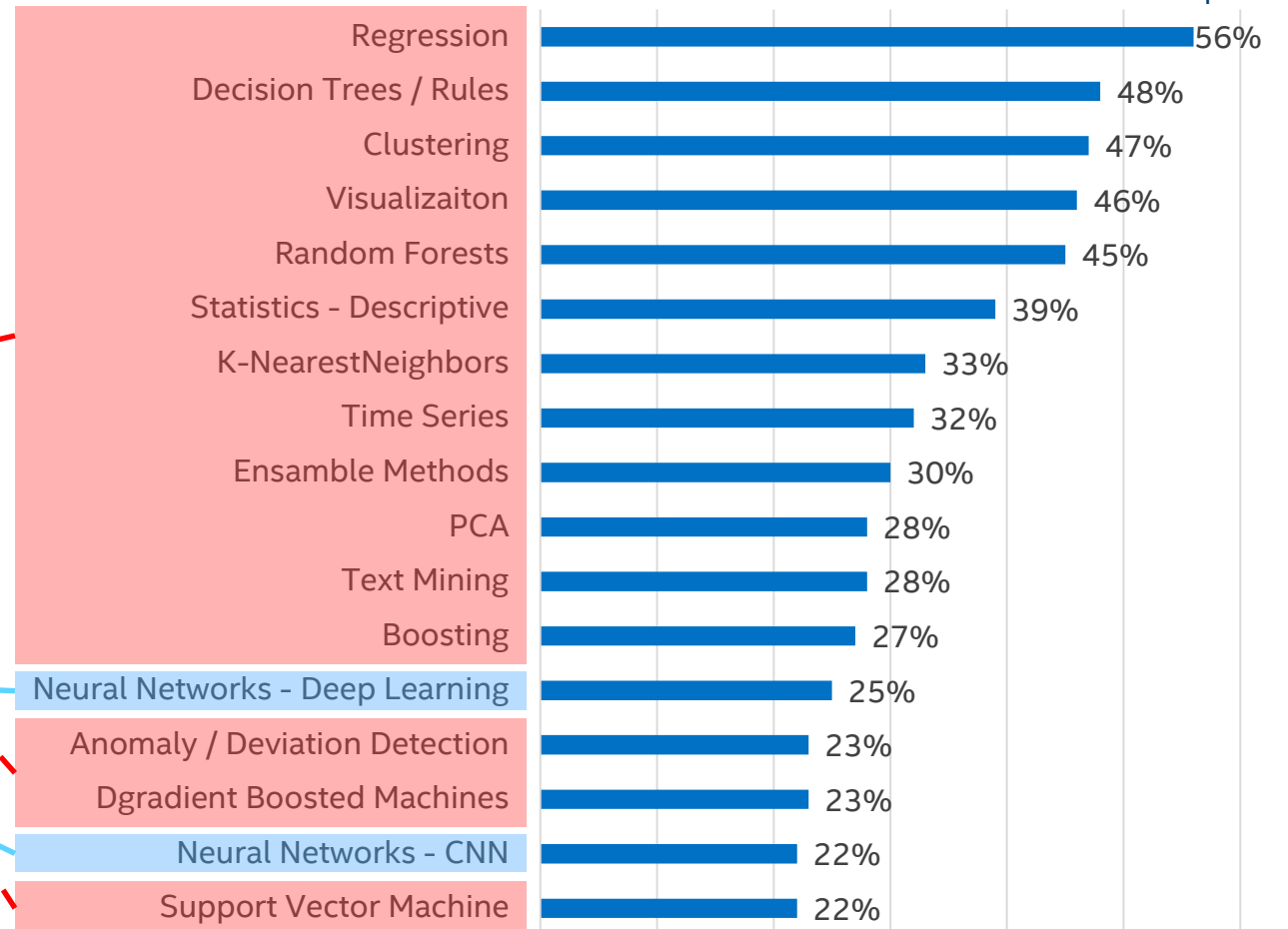
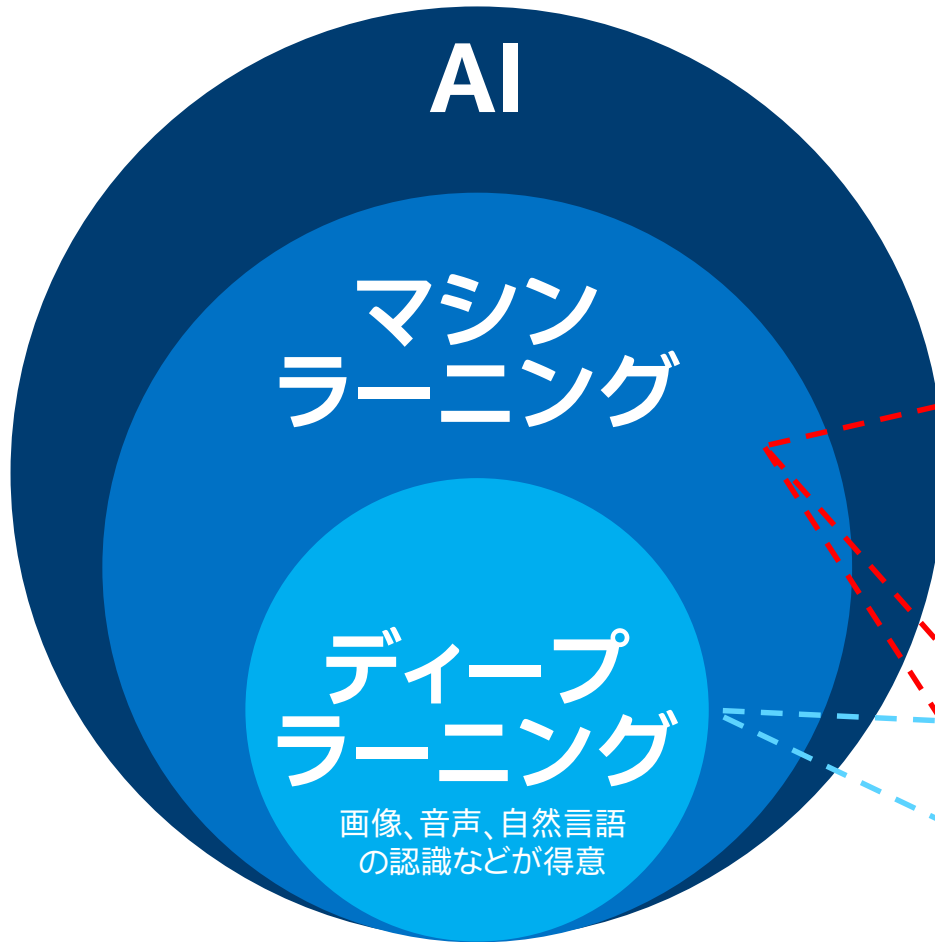


# まだまだマシンラーニングは重要

Top Data Science, Machine Learning

Methods used in 2018/2019

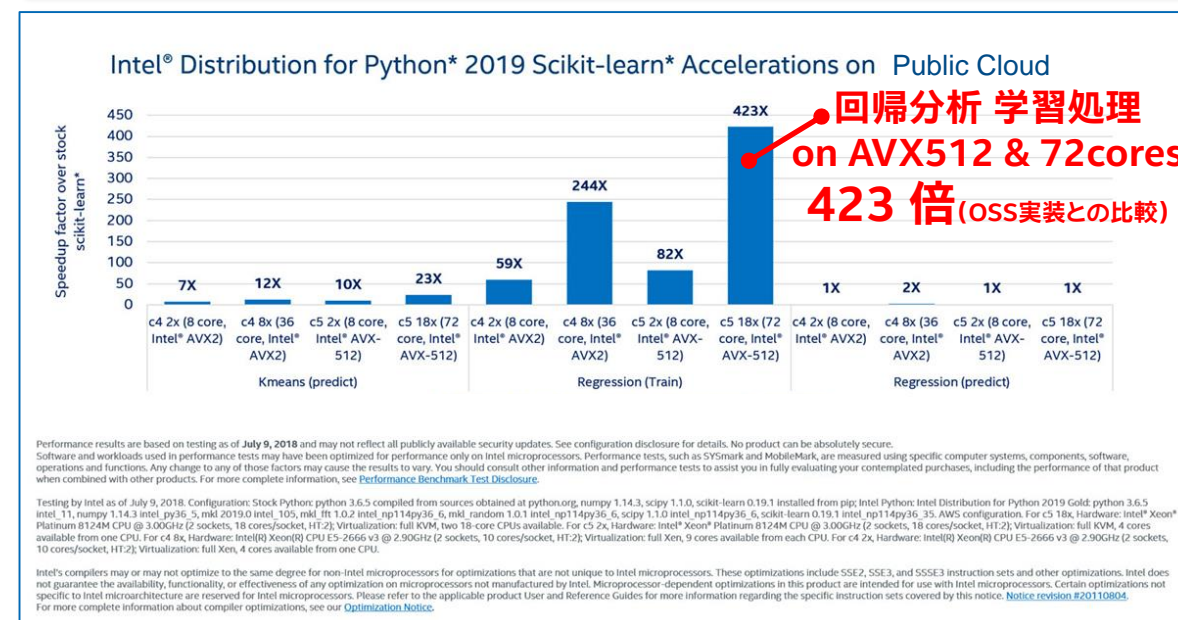
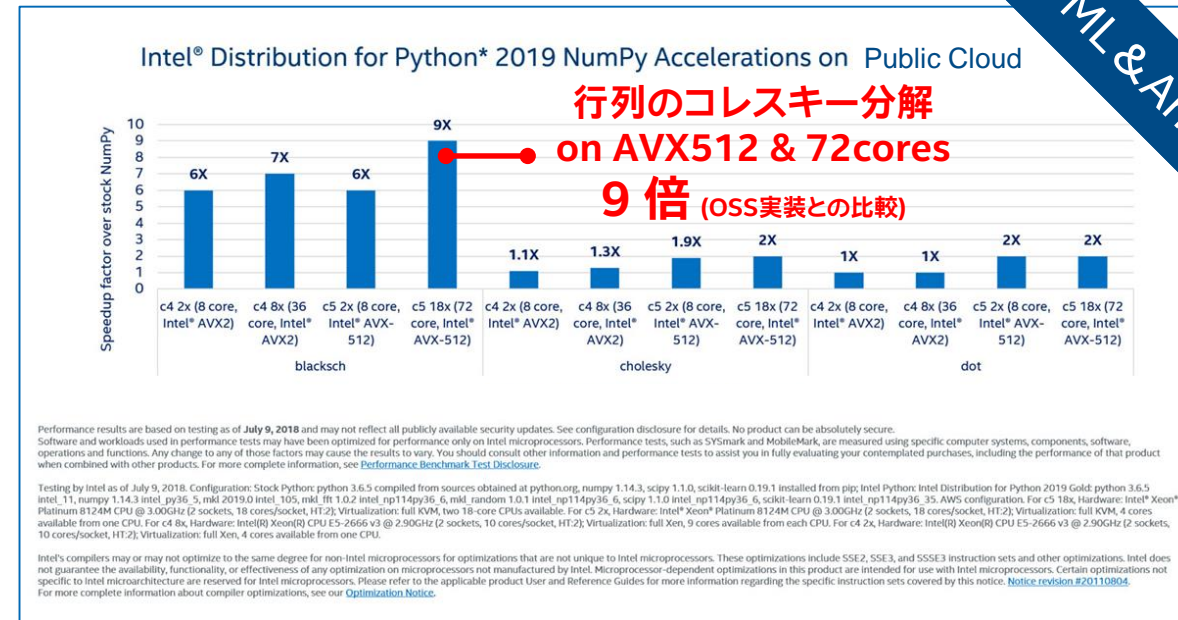
Share of Respondents



# Intel® Distribution for Python\*

インテルが実装、かつ、最適化した Python、および、周辺ライブラリ

- Numpy
- Pandas
- Scipy
- Scikit-learn
- XGBoost
- TensorFlow
- etc..



# インテル® のAI系ライブラリー & oneDALの使い方

数学	マシンラーニング/ データ分析	ディープラーニング	集団通信
インテル® oneAPI Math Kernel Library (oneMKL)	インテル® oneAPI Data Analytics Library (oneDAL)	インテル® oneAPI Deep Neural Network Library (oneDNN)	インテル® oneAPI Collective Communication Library (oneCCL)



daal4py



パートナー  
ソリューション

`pip install daal4py`

`pip install intel-scikit-learn`

Sparkへのインストール可能

<http://www.intel.com/analytics>

# 新たなデマンドと新たなテクノロジー

## Security

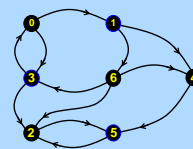
### PPML

(Privacy Preserving Machine Learning)

プライバシー情報保護に重きを置いた機械学習技術

## Data

### Graph



グラフデータに対する分析  
または機械学習を用いての  
パターン検出など

## Algorithm

### SLIDE

(Sub-Linear Deep learning Engine)

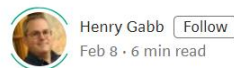
ライス大学との共同研究。  
ディープラーニング学習アルゴリズムを抜本的に見直すことでCPUにてGPUを上回る学習性能を実現



# グラフ分析に関するインテルの技術ブログ

## Measuring Graph Analytics Performance

The Diverse Landscape of Graph Analytics Requires a Comprehensive Benchmark

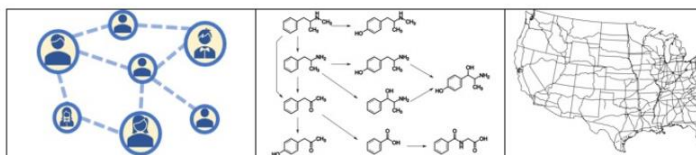


Henry Gabb [Follow](#)  
Feb 8 · 6 min read



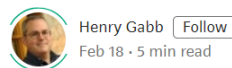
### What Is Graph Analytics And Why Does It Matter?

A graph is a good way to represent a set of objects and the relations between them (**Figure 1**). Graph analytics is the set of techniques to extract information from connections between entities.



## Adventures in Graph Analytics Benchmarking

It's Important to Use a Benchmark for Its Intended Purpose



Henry Gabb [Follow](#)  
Feb 18 · 5 min read



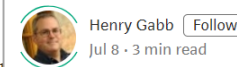
With all the attention graph analytics is getting lately, it's increasingly important to measure its performance in a comprehensive, objective, and reproducible way. I covered this in a [previous article](#), in which I recommended using an off-the-shelf benchmark like the [GAP Benchmark Suite](#) from the University of California, Berkeley. There are other graph benchmarks, of course, like [LDBC Graphalytics](#), but they can't beat GAP for ease of use. There's significant overlap between GAP and Graphalytics, but the latter is an industrial-strength benchmark that requires a special software configuration.

Measuring Graph Analytics Performance



## You Don't Have to Spend \$800,000 to Compute PageRank

There's a Better Way to Do Large-Scale Graph Analytics



Henry Gabb [Follow](#)  
Jul 8 · 3 min read



Benchmarking isn't my favorite topic, but I have a passing interest in graph analytics benchmarking:

Measuring Graph Analytics Performance

What Is Graph Analytics And Why Does It Matter?

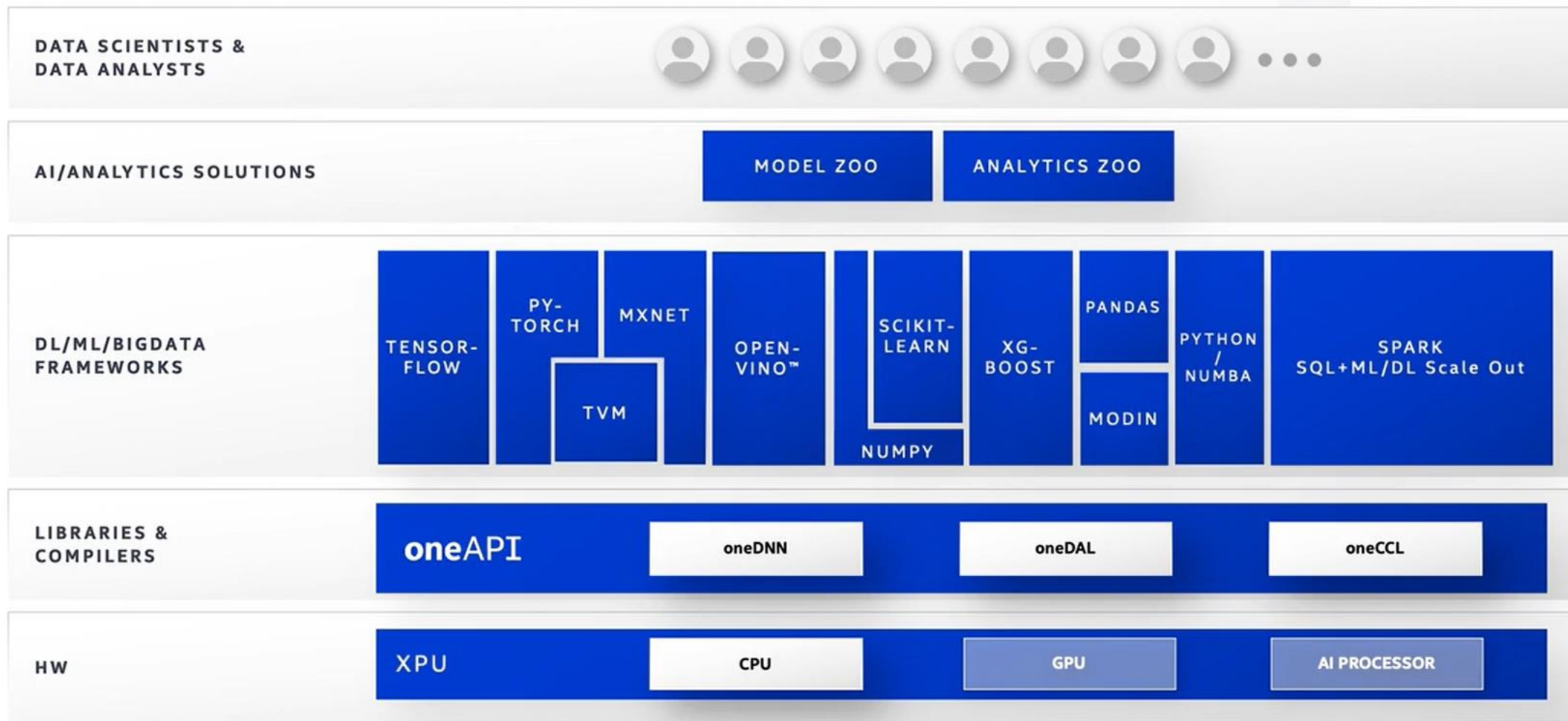
medium.com



I'll occasionally dissect benchmarks that I think are inaccurate or misleading:

<https://medium.com/intel-analytics-software>

# AI Software Ecosystem on Intel



Refer to <https://software.intel.com/articles/optimization-notice> for more information regarding performance and optimization choices in Intel software products.

# インテル製品で AI への取り組みを加速



発見

可能性と次のステップ



データ

セットアップ、取り込み、  
クリーニング



開発

分析 / AI を使用するモデル



導入

実稼動へ & 反復

## エコシステム

インテル® AI  
ビルダーズ・  
プログラム

100 以上の垂直的 / 水平的  
エコシステムのソリューション

最適化された  
クラウド

Amazon Web Services\*  
Baidu\* Cloud  
Google Cloud\* Platform  
Microsoft\* Azure\*  
など

AI に最適化  
された構成



## ソフトウェア

データ  
分析

50 以上の最適化された  
ソフトウェア・プラットフォーム

マシン  
ラーニング

Python\* 向け  
インテル®  
ディストリ  
ビューション



ディープ  
ラーニング



## ハードウェア

転送



格納



処理



すべての製品、コンピューター・システム、日付、および数値は、現在の予想に基づくものであり、予告なく変更されることがあります。最適化に関する注意事項

intel®