# DELL AND INTEL AI ALIGNMENT

Gretchen Stewart,
Chief Data Scientist, Public Sector

Al Ford, Dell Technologies
Director AI Alliances

Key Take ways
- #S2MV – Joint AI strategy
- Democratizing AI – citizen data scientists
- Code Optimization and Code reuse - OpenSource
- Resources and Next steps

# Federal AI Initiatives

Department of Energy – AI for Science - https://anl.app.box.com/s/bpp2xokglo8z...0

"On the server side, storage systems, such as those that support Lustre and bu...t designed for and often perform poorly for these read-heavy, random access wo... realization has led to pursuit of alternative memory technologies including NAND flashand 3D Xpoint (e.g.,Intel Optane), because they offer superior energy efficiency and density."

NIST – US leadership in AI - https://www.nist.gov/system/files/documents/2019/08/1...ngagement_plan_9aug2019.pdf

Department of Defense – AI strategy - https://media.de...Feb/12/2002088963/-1/-1/1/SUMMARY-OF-DOD-AI-STRATEGY-PDF

- USAF AI Annex to DoD st...//www.af.mil/Portals/1/documents/5/USAF-AI-Annex-to-DoD-AI-Strategy.pdf
- Transformation of DoD... ng up the Joint AI Center - https://www.ai.mil/

NOAA – https://nrc.noaa.go...x?fileticket=0I2p2-Gu3rA%3D&...jd=0

"As data exploitation capabilities continue to increase exponentially with ...ystems and architecture, unmanned systems and commercial data sources, AI met... transformative advancements in the quality and timeliness of NOAA science... services."

AI in Government Act 2020

WH Executive Order 2019

National AI Research Institutes - NSF

Centers of Excellence - GSA

# Speed to Mission Value #S2MV



Can you quickly deploy the solution?

IT specialists with knowledge on AI infrastructure?

Developers to work on chosen accelerators?

Data Scientist available to manipulate data?

Deliver mission objective

What AI infrastructure is needed?

Developers to deliver and test applications?

Data scientist available to develop algorithms?

What are the data sets you have to work with?

DELL EMC    intel

# AI Strategy Powered by Dell and Intel

## Lead with Software



Leverage Software portfolio
Optimize agency software
Evangelize to developers

## Deliver the best AI Platforms



Silicon alternatives

TRAINING
CPUs and GPUs, limited FPGAs,
ASICs under investigation

EVALUATION
CPUs and FPGAs,
ASICs under investigation

CPUs    GPUs    FPGAs    ASICs

FLEXIBILITY ← → EFFICIENCY

Extend the CPU
Align compute and memory
Drive "Fit for Purpose"
compute

## Ignite the Ecosystem



Train the Team
Leverage >1000 Intel use cases
Pioneer AIoT deployments

**DELL**EMC    **intel**

# The AI Buyer

**1** **Agency CIO / CDO**



**2** **IT System Architect**



**3** **AI Developer, Programmer Data Scientist**



**Key Concern(s):** Agency mission effectiveness

**Key Messaging:**
- Time to results
- Alignment with Enterprise strategy

**Key Concern(s):** TCO, security & manageability

**Key Messaging:**
- Infrastructure integration
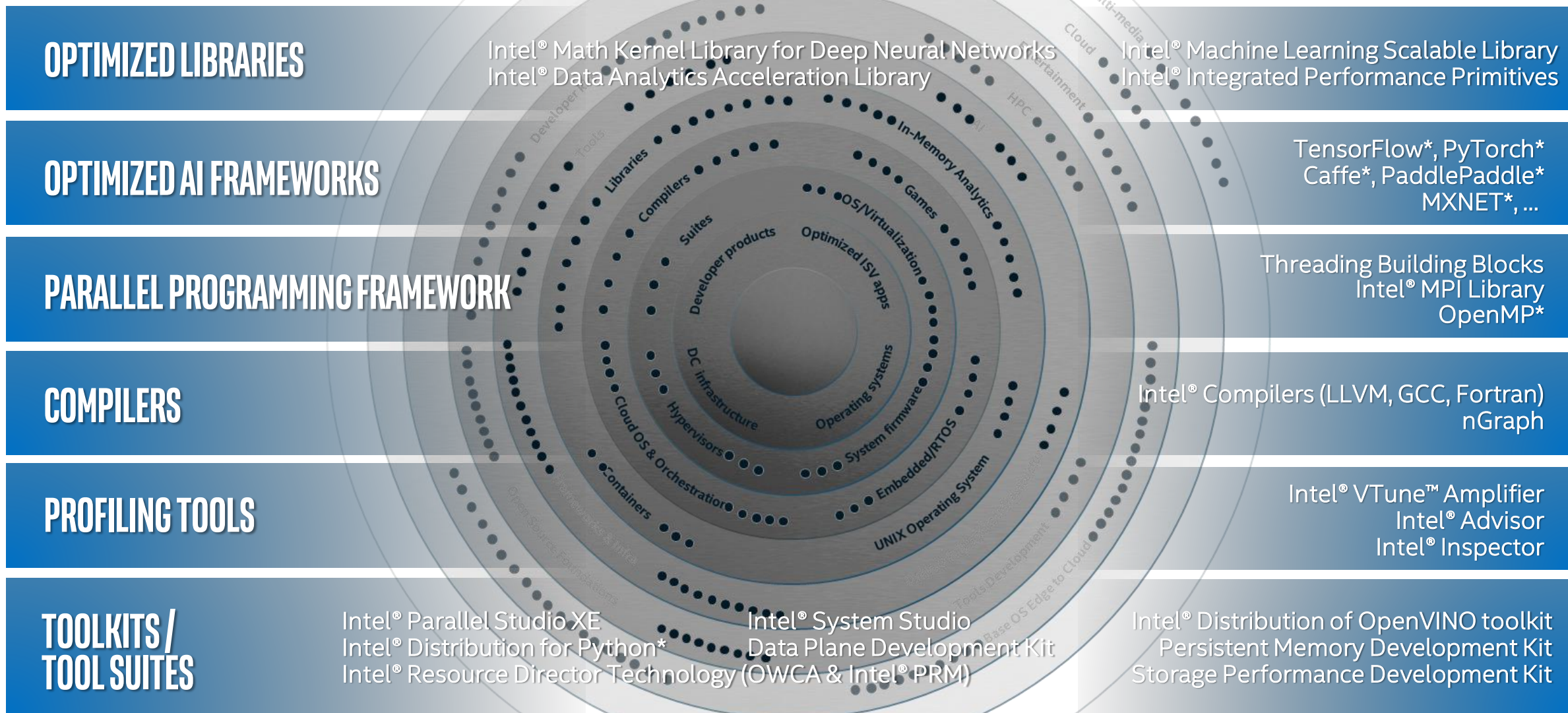- Create sharable platforms
- Time to results

**Key Concern(s):** Time-to-solution

**Key Messaging:**
- Scalability of Models
- Reduction in training times
- Development and Deployment Framework Support

**DELL**EMC  **intel.**

# WORKLOAD PERFORMANCE THROUGH OPTIMIZED SOFTWARE

**OPTIMIZED LIBRARIES**

Intel® Math Kernel Library for Deep Neural Networks
Intel® Data Analytics Acceleration Library

Intel® Machine Learning Scalable Library
Intel® Integrated Performance Primitives

**OPTIMIZED AI FRAMEWORKS**

TensorFlow*, PyTorch*
Caffe*, PaddlePaddle*
MXNET*, ...

**PARALLEL PROGRAMMING FRAMEWORK**

Threading Building Blocks
Intel® MPI Library
OpenMP*

**COMPILERS**

Intel® Compilers (LLVM, GCC, Fortran)
nGraph

**PROFILING TOOLS**

Intel® VTune™ Amplifier
Intel® Advisor
Intel® Inspector

**TOOLKITS / TOOL SUITES**

Intel® Parallel Studio XE
Intel® Distribution for Python*
Intel® Resource Director Technology (OWCA & Intel® PRM)

Intel® System Studio
Data Plane Development Kit

Intel® Distribution of OpenVINO toolkit
Persistent Memory Development Kit
Storage Performance Development Kit

# Intel Distribution for Python

**software.intel.com/intel-distribution-for-python**

## FOR DEVELOPERS USING THE MOST POPULAR AND FASTEST-GROWING PROGRAMMING LANGUAGE FOR AI

| EASY, OUT-OF-THE-BOX ACCESS TO HIGH-PERFORMANCE PYTHON | DRIVE PERFORMANCE WITH MULTIPLE OPTIMIZATION TECHNIQUES | FASTER ACCESS TO LATEST OPTIMIZATIONS FOR INTEL ARCHITECTURE |
|---|---|---|
| ▪ Prebuilt, optimized for numerical computing, data analytics, HPC<br><br>▪ Drop-in replacement for your existing Python (no code changes required) | ▪ Accelerated NumPy/SciPy/Scikit-Learn with Intel Math Kernel Library (Intel MKL)<br><br>▪ Data analytics with pyDAAL, enhanced thread scheduling with TBB, Jupyter Notebook interface, Numba, Cython<br><br>▪ Scale easily with optimized MPI4Py and Jupyter notebooks | ▪ Distribution and individual optimized packages available through conda and Anaconda Cloud<br><br>▪ Optimizations upstreamed back to main Python trunk |

## ADVANCING PYTHON PERFORMANCE CLOSER TO NATIVE SPEEDS

# Optimized Deep Learning Frameworks and Toolkits

## Gen on gen performance gains for ResNet-50 with Intel DL Boost

**2S Intel Xeon Platinum 8280 Processor  vs  2S Intel Xeon Platinum 8180 Processor**

| Intel Xeon Scalable Processor | 2nd Gen Intel Xeon Scalable Processor | mxnet | PyTorch | TensorFlow | Caffe | OpenVINO |
|---|---|---|---|---|---|---|
| FP32 | → INT8 w/ Intel DL Boost | 3.0x | 3.7x | 3.9x | 4.0x | 3.9x |
| INT8 | → INT8 w/ Intel DL Boost | 1.8x | 2.1x | 1.8x | 2.3x | 1.9x |

# Intel's oneAPI Ecosystem

## Built on Intel's Rich Heritage of CPU Tools Expanded to XPUs

oneAPI

A cross-architecture language based on C++ and SYCL standards

Powerful libraries designed for acceleration of domain-specific functions

A complete set of advanced compilers, libraries, and porting, analysis and debugger tools

Powered by oneAPI

Frameworks and middleware that are built using one or more of the oneAPI industry specification elements, the DPC++ language, and libraries listed on oneapi.com.

---

**Application Workloads Need Diverse Hardware**

**Middleware & Frameworks (Powered by oneAPI)**

TensorFlow  PyTorch  MODIN  learn  NumPy  XGBoost  OpenVINO  ...

**1 oneAPI**  Intel® oneAPI Product

| Compatibility Tool | Languages | Libraries | | Analysis & Debug Tools |
|---|---|---|---|---|
| | | oneMKL oneTBB oneVPL oneDPL | oneDAL oneDNN oneCCL | |

Low-Level Hardware Interface

**XPUs**

CPU    GPU    FPGA    Other accelerators

Available Now

# oneAPI Ecosystem Support



These organizations support the oneAPI initiative 'concept' for a single, unified programming model for cross-architecture development. It does not indicate any agreement to purchase or use of Intel's products.
*Other names and brands may be claimed as the property of others.

# AI Software Stack for Intel® XPUs

Intel offers a Robust Software Stack to Maximize Performance of Diverse Workloads



| E2E Workloads (Census, NYTaxi, Mortgage…) | Intel® Low Precision Optimization Tool | Model Zoo for Intel® Architecture | Open Model Zoo | DL/ML Tools |

numba, pandas, numpy, Scikit-learn, Modin, scipy, daal4Py, xgboost | TensorFlow | PyTorch | Model Optimizer Inference Engine | DL/ML Middleware & Frameworks

DPC++ / DPPY | oneMKL | oneDAL | oneTBB | oneCCL | oneDNN | oneVPL | Libraries & Compiler
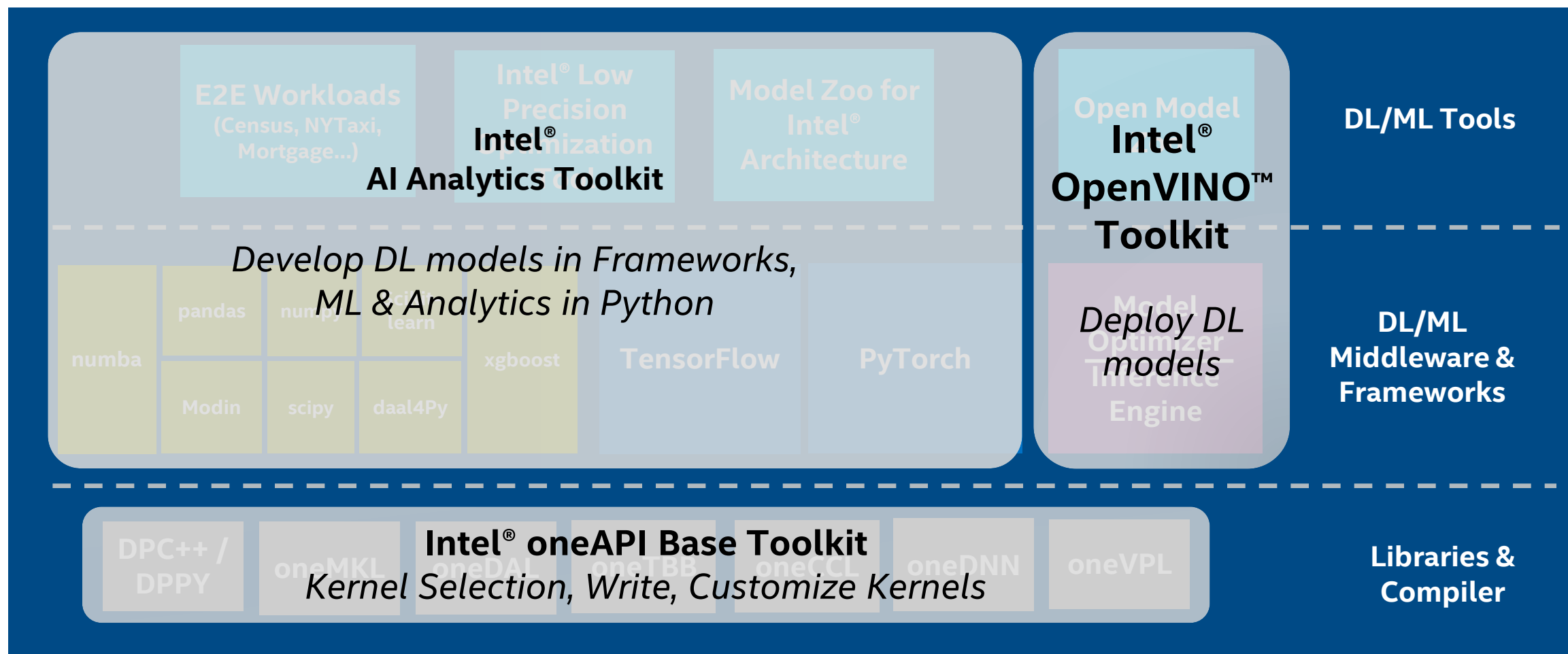
# AI Software Stack for Intel® XPUs

## Intel offers a robust software stack to maximize performance of diverse workloads

E2E Workloads
(Census, NYTaxi, Mortgage...)

Intel® Low Precision Optimization Tool

Model Zoo for Intel® Architecture

**Intel® AI Analytics Toolkit**

Open Model

**Intel® OpenVINO™ Toolkit**

**DL/ML Tools**

*Develop DL models in Frameworks, ML & Analytics in Python*

numba

pandas

numpy

sci- learn

xgboost

Modin

scipy

daal4Py

TensorFlow

PyTorch

Model Optimizer

Inference Engine

*Deploy DL models*

**DL/ML Middleware & Frameworks**

DPC++ / DPPY

oneMKL

oneDAL

**Intel® oneAPI Base Toolkit**
*Kernel Selection, Write, Customize Kernels*

oneTBB

oneCCL

oneDNN

oneVPL

**Libraries & Compiler**

**Full Set of AI ML and DL Software Solutions Delivered with Intel's oneAPI Ecosystem**

# Intel® Distribution of OpenVINO™ toolkit

## Powered by oneAPI

A toolkit for faster, more accurate real-world results using high-performance, AI and computer vision inference deployed into production across Intel® architecture from edge to cloud
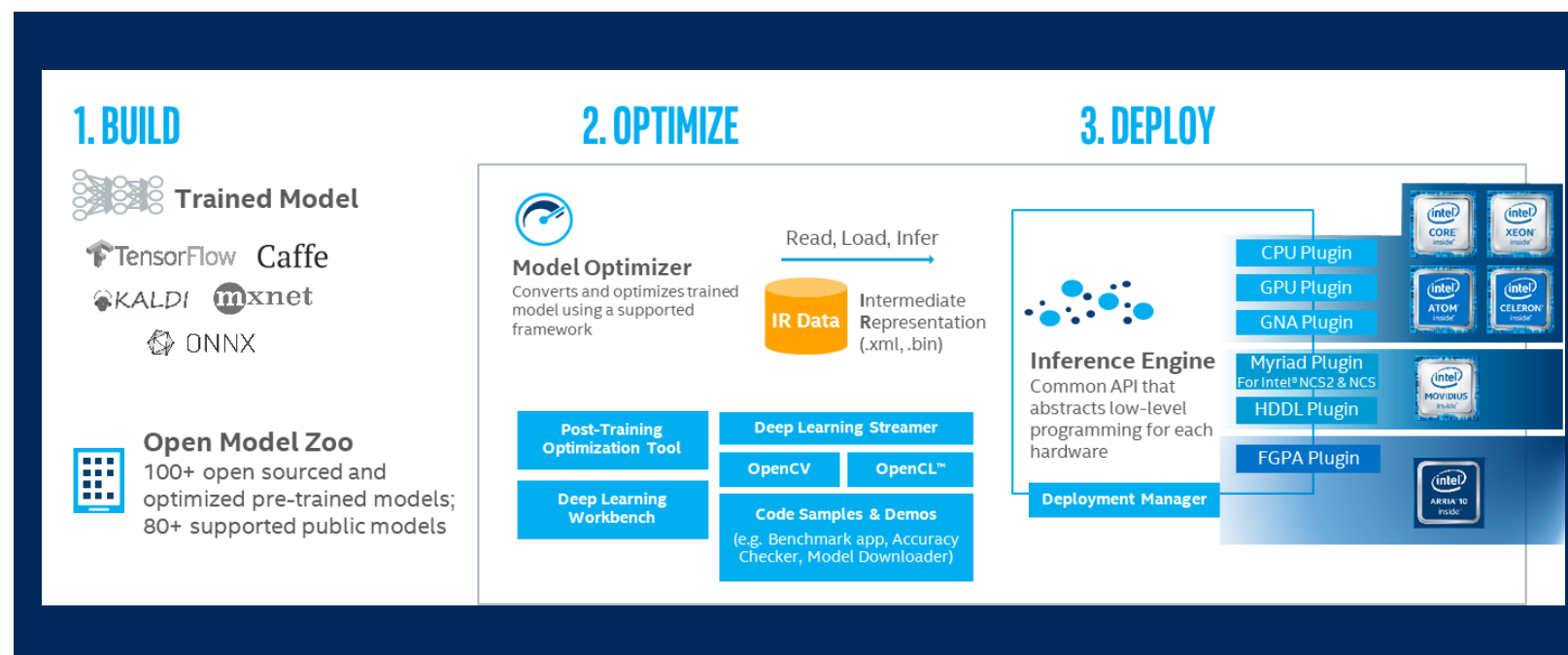
**Who needs this product?**

AI application developers, OEMs, ISVs, System Integrators, Vision and Media developers

**Top Features/Benefits**

High-performance, deep learning inference deployment

Streamlined development; ease of use

Write once, deploy anywhere

Proven, industry-leading accelerated technology



https://software.intel.com/en-us/openvino-toolkit

# Getting Started with Intel® AI Analytics Toolkit

| Overview | Installation | Hands on | Learning | Support |
|---|---|---|---|---|
| <ul><li>Visit Intel® AI Analytics Toolkit (AI Kit) for more details and up-to-date product information</li><li>Release Notes</li></ul> | <ul><li>Download the AI Kit from Intel, Anaconda or any of your favorite package managers</li><li>Get started quickly with the AI Kit Docker Container</li><li>Installation Guide</li><li>Utilize the Getting Started Guide</li></ul> | <ul><li>Code Samples</li><li>Build, test and remotely run workloads on the Intel® DevCloud for free. No software downloads. No configuration steps. No installations.</li></ul> | <ul><li>Machine Learning & Analytics Blogs at Intel Medium</li><li>Intel AI Blog site</li><li>Webinars and Articles at Intel® Tech Decoded</li></ul> | <ul><li>Ask questions and share information with others through the Community Forum</li><li>Discuss with experts at AI Frameworks Forum</li></ul> |

## Download Now

# Which Toolkit to Use When ?

| | Intel® AI Analytics Toolkit | OpenVINO™ Toolkit |
|---|---|---|
| Key Value Prop | • Provide performance and easy integration across end-to-end data science pipeline for efficient AI model development<br>• Maximum compatibility with opensource FWKs and Libs with drop-in acceleration that require minimal to no code changes<br>• Audience: Data Scientists; AI Researchers; DL/ML Developers | • Provide leading performance and efficiency for DL inference solutions to deploy across any Intel HW (cloud to edge).<br>• Optimized package size for deployment based on memory requirements<br>• Audience: AI Application Developers; Media and Vision Developers |
| Use Cases | • Data Ingestion, Data pre-processing, ETL operations<br>• Model training and inference<br>• Scaling to multi-core / multi-nodes / clusters | • Inference apps for vision, Speech, Text, NLP<br>• Media streaming / encode, decode<br>• Scale across HW architectures – edge, cloud, datacenter, device |
| HW Support | • CPUs - Datacenter and Server segments – Xeons, Workstations<br>• GPU - ATS and PVC (in future) | • CPU - Xeons, Client CPUs and Atom processors<br>• GPU - Gen Graphics; DG1 (current), ATS, PVC (in future)<br>• VPU - NCS & Vision Accelerator Design Products,<br>• FPGA<br>• GNA |
| Low Precision Support | **Use Intel® Low Precision Optimization Tool when using AI Analytics Toolkit**<br>• Supports BF16 for training and FP16, Int8 and BF16 for Inference<br>• Seamlessly integrates with Intel optimized frameworks<br>• Available in the AI toolkit and independently | **Use Post Training Optimization Tool when using OpenVINO**<br>• Supports FP16, Int8 and BF16 for inference<br>• Directly works with Intermediate Representation Format<br>• Available in the Intel Distribution of OpenVINO toolkit<br>• Provides Training extension via NNCF for PyTorch with FP16, Int8 |

**Exception**: If a model is not supported by OpenVINO™ toolkit for Inference deployment, build custom layers for OV or fall back to the AI Analytics Toolkit and use optimized DL frameworks for inference.

# Use both !

**Toolkits are complimentary to each other and recommendation is to use them both based on your current phase of AI Journey**

- I am **exploring and analyzing data;** I am **developing models**

- I want **performance and compatibility** with frameworks and libraries I use

- I would like to have **drop-in acceleration** with little to no additional code changes

- I prefer **not to learn any new tools** or languages

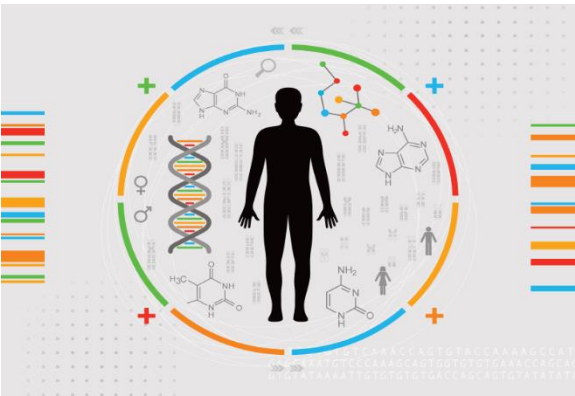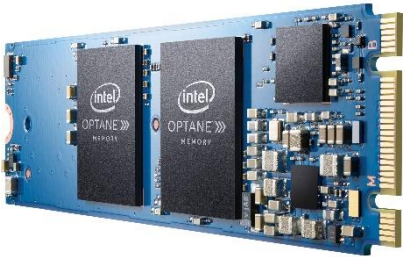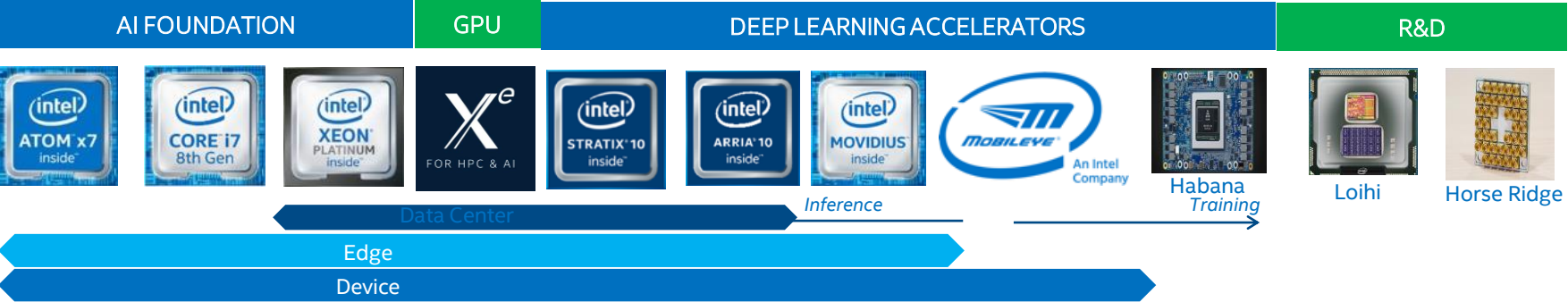**Data Scientist/ML Developer**
Intel® AI Analytics Toolkit

- I am **deploying models**

- I want **leading performance and efficiency** across multiple target HW

- I'm concerned about **having lower memory footprint**, which is critical for deployment

- I am **comfortable with learning and adopting a new tool or API** to do so

**App Developer**
Intel® Distribution of OpenVINO™ toolkit

*If you prefer working on primitives and optimize kernels and algorithms directly using oneAPI libraries (oneDNN, oneCCL & oneDAL), then use Intel® Base Toolkit*

intel.

# Fit for Purpose AI Compute



| AI FOUNDATION | | | GPU | DEEP LEARNING ACCELERATORS | | | | | R&D | |
|---|---|---|---|---|---|---|---|---|---|---|



intel ATOM x7 inside · intel CORE i7 8th Gen · intel XEON PLATINUM inside · X^e FOR HPC & AI · intel STRATIX 10 inside · intel ARRIA 10 inside · intel MOVIDIUS inside · MOBILEYE An Intel Company · Habana *Training* · Loihi · Horse Ridge

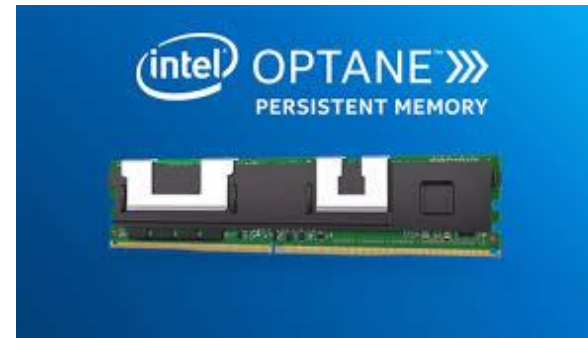Data Center

*Inference*

Edge

Device

# Align Compute and Memory

Hard Problem: Traditional AI architectures don't address the challenge of accessing large and small data sets and the data latency from intensive ETL (Extract, Transform, Load) processes



- **Supercharges AI Application Performance**

- **Increases Speed 2 Mission Value**

Key Trend: Application developers want to run AI applications in persistent memory to eliminate bottlenecks and accelerate performance

# Software Driven Market Ready Solutions
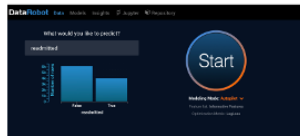
## IGUAZIO Reference Architecture



| ToR Switch | Management: PowerSwitch S3148 (1GbE) |
|---|---|
| | Cluster: Mellanox SN2700 Open Ethernet Switch 100GbE |

## H2O.ai DAI Reference Architecture

**DRIVERLESSAI**
Automatic feature engineering, machine learning and interpretability

## DataRobot Reference Architecture

| Configuration | Data Robot compute node |
|---|---|
| Server | 4x PE-C6420 |
| Processor | 2 x Intel Xeon® Gold Scalable 6248 processor |
| Memory | 512GB DDR4 @ 2,667 MHz |
| Drives | OS: Boss Card with 2 X 480G SSD<br>Data: 12 x Dell Express Flash NVMe P4610 1.6TB SFF |
| Networking | Intel(R) Ethernet 10G 4P X710 SFP+ rNDC |
| ToR Switch | Management: PowerSwitch S3048-ON (1GbE)<br>Cluster: PowerSwitch S5248F-ON (10/25GbE) |

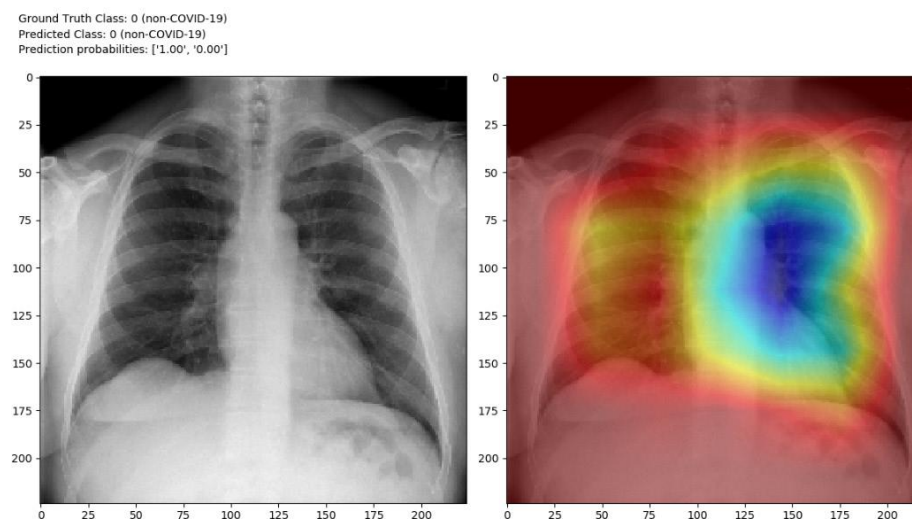| | Stand Alone (VM) | | | Hadoop |
|---|---|---|---|---|
| Data Robot Configs | Small Config 12Modeler+2Pred | Medium Config 20Modeler+3Pred | Large Config 30Modeler+6Pred | 12Modeler+2Pred |

Demonstrating Customer Obsession

- 1x software ~10x hardware pull thru
- Reference Arch accelerates AI Adoption
- Software optimized for AI performance

**DELL**EMC    **intel.**

# Case Study - Accrad

## Better together – using both toolkits

CheXRad is a machine learning at the edge application that helps Radiologists and Physicians to identify COVID-19, viral pneumonia and other diseases on chest X-ray images and predict the need for ventilators

- CheXRad comes pre-configured with a COVID-19 and viral pneumonia classification neural network

- To architect, train and validate the neural network, Accrad used Intel Tensorflow from AI Analytics Toolkit and the infrastructure provided by Intel oneAPI DevCloud and developed the model.

- To further optimize its model for deployment, they used OpenVINO Toolkit and Intel DevCloud for Edge

- CheXRad could label pathologies in 140 chest x-rays in just 90 seconds—up to **160x faster** than radiologists, at comparable levels of accuracy, sensitivity and specificity.



Ground Truth Class: 0 (non-COVID-19)
Predicted Class: 0 (non-COVID-19)
Prediction probabilities: ['1.00', '0.00']

*"With the help of Intel, we were able to **train, optimize, and deploy** a machine learning model in **less time and at a lower operational cost** than available alternatives, enabling us to get to market fast with a powerful solution that's optimized for Intel® architecture." – Moloti Nakampe, R&D Director*

# Case Study - AbbVie

## Better together – using both toolkits

Drop-in acceleration

- AbbVie is a research-based biopharmaceutical company using Xeons

- Abbelfish Machine Translation uses Intel® Optimization for TensorFlow of AI Analytics Toolkit
  - Custom model was used to provide more accurate translations than commercially available ones. Model includes 24 layers and over 500 million parameters, which took over four months to train;
  - Intel TF provided greater performance boost while the customer did not have to change their code / api's from standard TF

- AbbVie Search uses Intel® Distribution of OpenVINO™ toolkit

Deployment

  - OpenVINO provided great speed up to answer questions from a scientific article or clinical report when compared to standard TF.
  - Requires scaling across the company, so uses OpenVINO Model Server to serve inferences



**Abbelfish Model Performance**

Intel® Xeon® Gold 6252N Processor 1.9x speedup. Higher is better.

Sentences per second — DNNL Disabled: 2.6, DNNL Enabled: 5.0
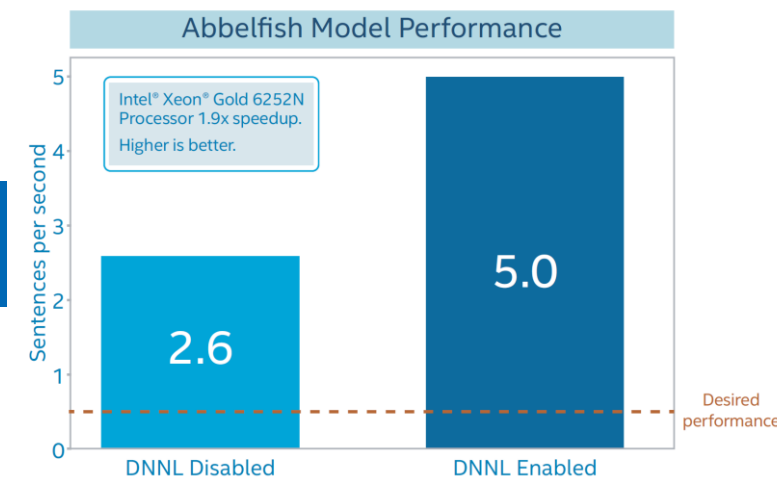
Desired performance

**Figure 3.** AbbVie's Abbelfish translated over five sentences per second using Intel Optimization for TensorFlow with oneAPI Deep Neural Network Library (oneDNN).[1]
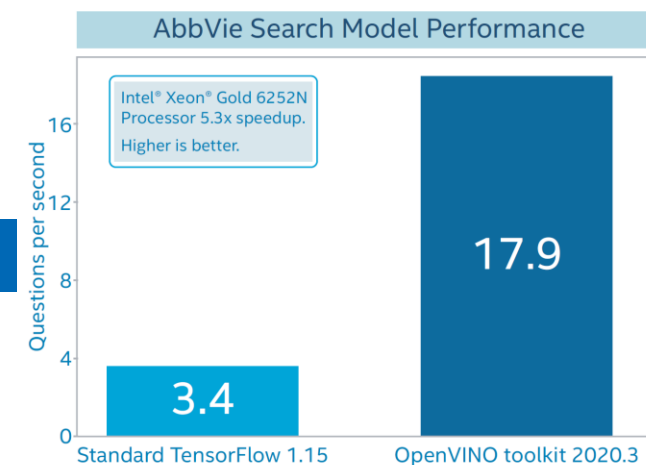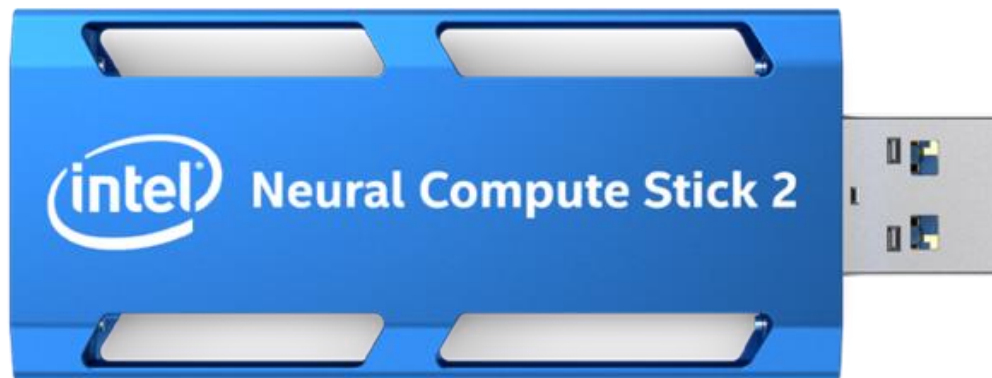
### AbbVie Search



**AbbVie Search Model Performance**

Intel® Xeon® Gold 6252N Processor 5.3x speedup. Higher is better.

Questions per second — Standard TensorFlow 1.15: 3.4, OpenVINO toolkit 2020.3: 17.9

**Figure 5.** Comparison of AbbVie Search inference between unoptimized TensorFlow 1.15 (oneDNN disabled) and OpenVINO toolkit 2020.3.[1]

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

intel. 21

# Intel Neural Compute Stick 2

**USB STICK FORM FACTOR**
for neural network acceleration

**REAL-TIME ON-DEVICE INFERENCE**
no cloud connectivity required

**NO ADDITIONAL PERIPHERALS**
needed to start deploying solutions

**ACCELERATE DEVELOPMENT**
with Intel Distribution of OpenVINO toolkit

**INDUSTRY-LEADING PERFORMANCE[1]**
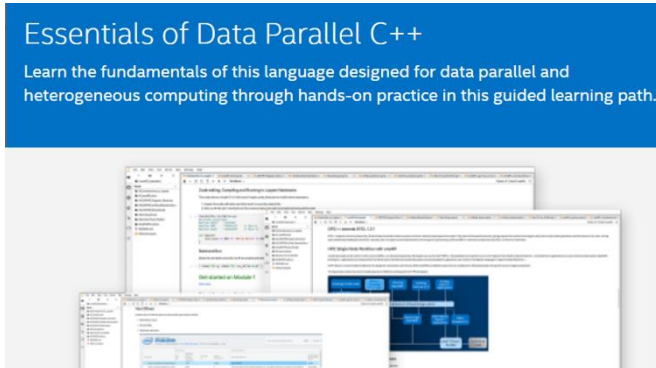with Intel Movidius Myriad™ X Vision Processing Unit (VPU)

**HIGHER PERFORMANCE[1]**
on deep neural networks compared to Intel Movidius Myriad 2 VPU

**8X**

[1]Testing by Intel as of October 12th, 2018
Deep Learning Workload Configuration. Comparing Intel® Movidius™ Neural Compute Stick based on Intel® Movidius™ Myriad™ 2 VPU vs. Intel® Neural Compute Stick 2 based on Intel® Movidius™ Myriad™ X VPU with Asynchronous Plug-in enabled for (2xNCE engines). As measured by images per second across GoogleNetV1. Base System Configuration: Intel® Core™ I7-8700K 95W TDP (6C12T at 3.7GHz base freq and 4.7GHz max turbo freq), Graphics: Intel® UHD Graphics 630 Total Memory 65830088 kB Storage: INTEL SSDSC2BB24 (240GB), Ubuntu 16.04.5 Linux-4.15.0-36-generic-x86_64-with-Ubuntu-16.04-xenial, deeplearning_deploymenttoolkit_2018.0.14348.0, API version 1.2, Build 14348, myriadPlugin, FP16, Batch Size = 1
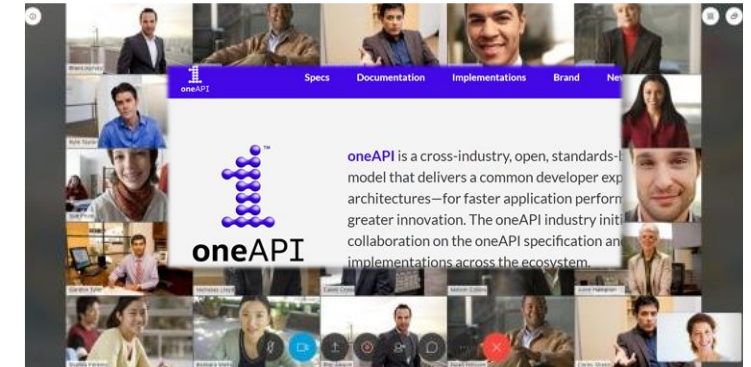
# Ecosystem Adoption & Support

## Training



Online webinars & courses, developer guides, sample code
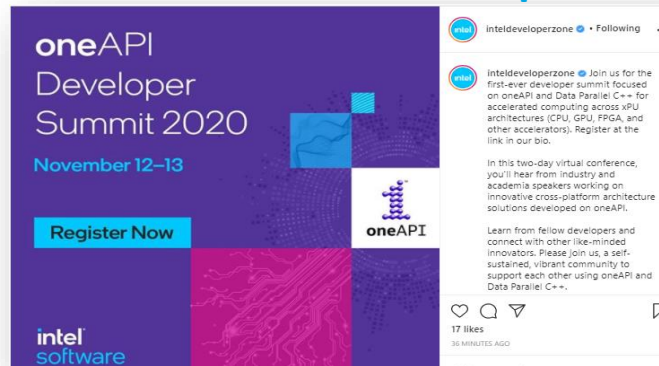
## Academia



oneAPI Centers of Excellence: research, enabling code, curriculum, teaching

## Community



oneAPI open specification, DevMesh innovators, community support forums

## Summits & Workshops



Live & on-demand virtual workshops, community-led sessions

## Industry Experts



Training by leading technical training companies worldwide

## Intel® DevCloud



State-of-the-art software and hardware
Intel® oneAPI Toolkits + latest Intel® Xeon® processors, GPUs (integrated & discrete), FPGAs

# Key Takeaways & Call to Action

$3-4 SW generate $10 in HW

➤ Create 3 – 5 proof of concepts. Ready to partner and we can aid with customer analytics meetings.

➤ The toolkits and software are optimized, and we can help with development and deployment to achieve performance and efficiency across different stages of AI Journey

➤ Recommend the toolkits based on current phase of customer pipeline

Your customers can download the toolkits for free. Intel partnered with DELL we can deliver workshops, coding sessions and POC's/

Intel® AI Analytics Toolkit

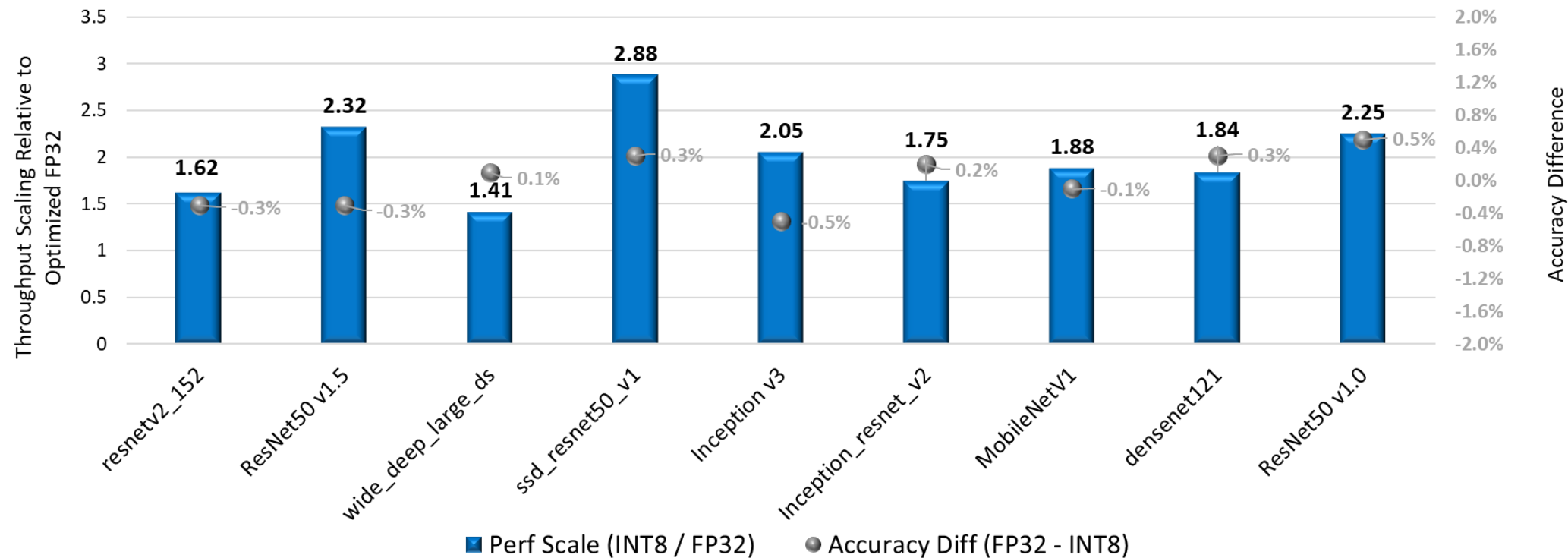Intel® Distribution of OpenVINO™ toolkit

Intel® oneAPI Base Toolkit

# INT8 Quantized Inference Performance

## Uses Intel® Optimization for Tensorflow and Intel® Low Precision Optimization Tool



INT8 Inference Throughput Scaling up to 3.8x and Accuracy Drop within 0.6%

For more complete information about performance and benchmark results, visit www.intel.com/benchmarks.
See backup for configuration details.

intel.