

**Choose the Best Accelerated Technology**

# Intel Performance optimizations for Deep Learning

Shailen Sobhee – Deep Learning Engineer

[shailen.sobhee@intel.com](mailto:shailen.sobhee@intel.com)

15 June 2021



# Agenda

- Quick recap of oneAPI
- Overview of oneDNN
- Training:
  - Overview of performance-optimized DL frameworks
    - Tensorflow
    - PyTorch
- Inferencing:
  - Intel® Low Precision Optimization Tool
  - Primer to Intel® Distribution of OpenVINO™
    - (Deeper presentation in next session)

# Intel's oneAPI Ecosystem

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

### oneAPI

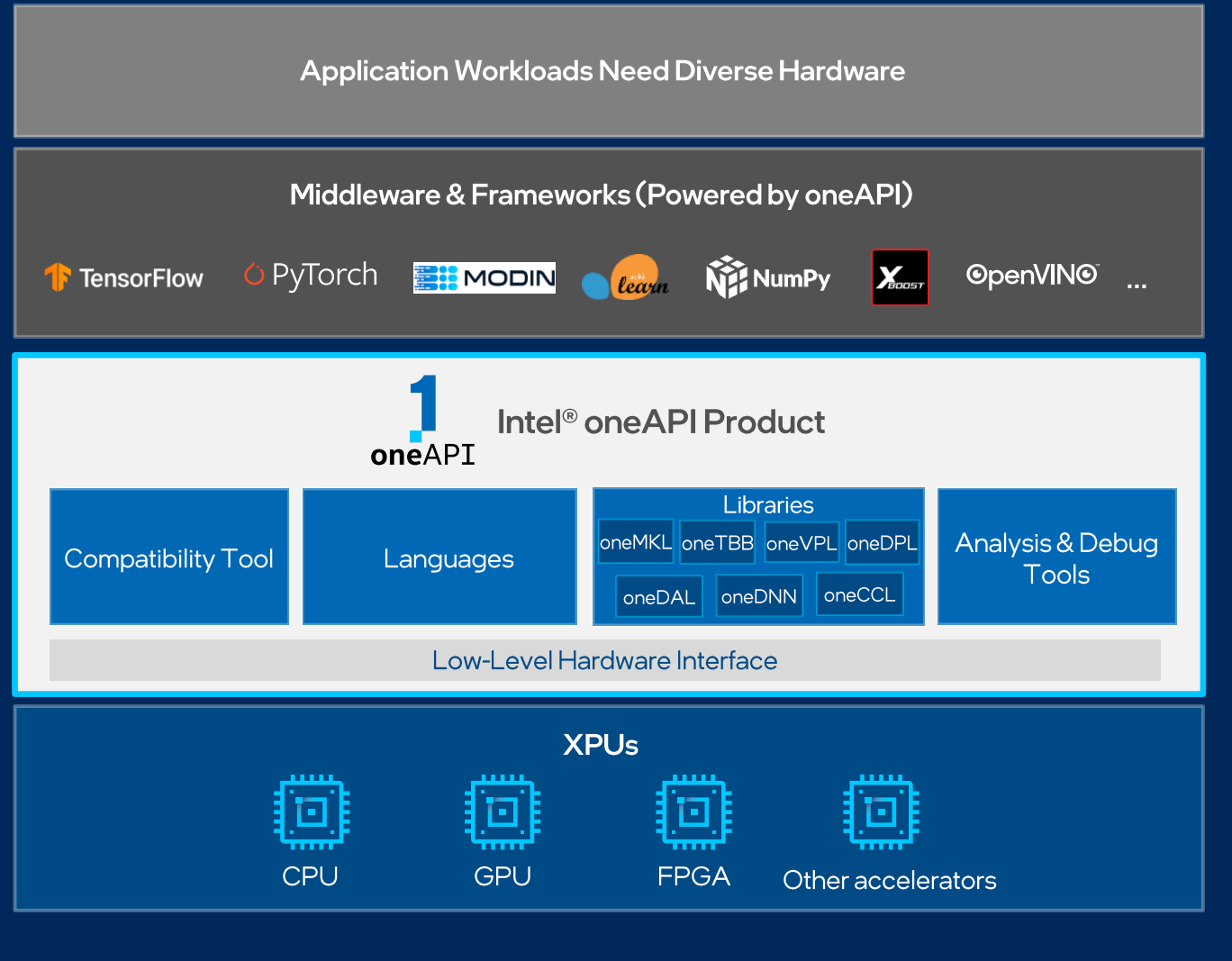
A cross-architecture language based on C++ and SYCL standards

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](https://oneapi.com).



[Available Now](#)

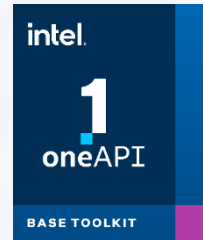
# Intel® oneAPI Toolkits

A complete set of proven developer tools expanded from CPU to XPU



## Intel® oneAPI Base Toolkit

Native Code Developers



A core set of high-performance tools for building C++, Data Parallel C++ applications & oneAPI library-based applications

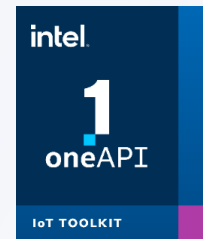
## Add-on Domain-specific Toolkits

Specialized Workloads



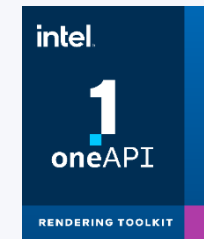
### Intel® oneAPI Tools for HPC

Deliver fast Fortran, OpenMP & MPI applications that scale



### Intel® oneAPI Tools for IoT

Build efficient, reliable solutions that run at network's edge



### Intel® oneAPI Rendering Toolkit

Create performant, high-fidelity visualization applications

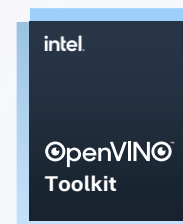
## Toolkits powered by oneAPI

Data Scientists & AI Developers



### Intel® AI Analytics Toolkit

Accelerate machine learning & data science pipelines with optimized DL frameworks & high-performing Python libraries



### Intel® Distribution of OpenVINO™ Toolkit

Deploy high performance inference & applications from edge to cloud

Latest version is 2021.1

# Intel® oneAPI AI Analytics Toolkit

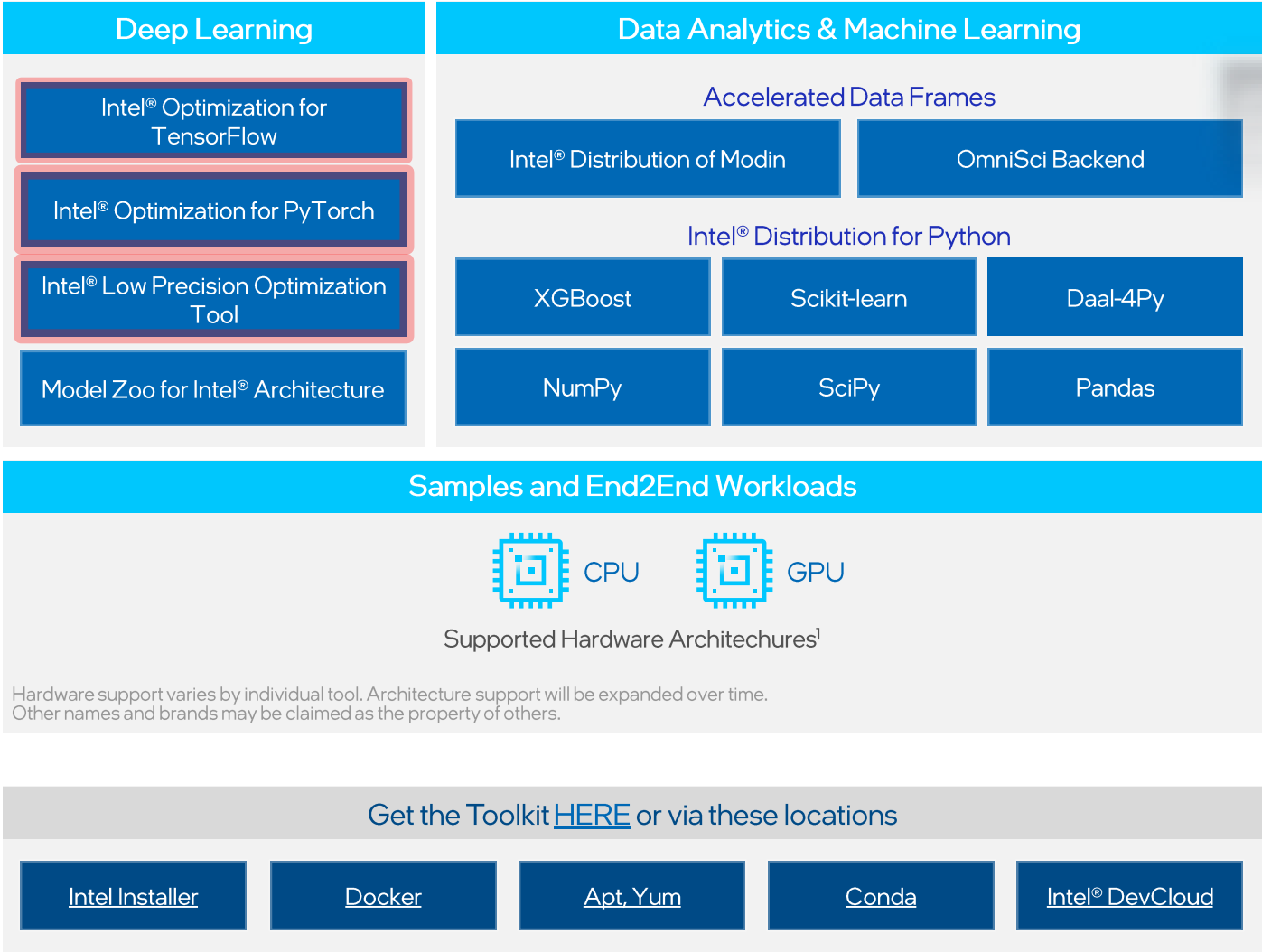
Accelerate end-to-end AI and data analytics pipelines with libraries optimized for Intel® architectures

## Who Uses It?

Data scientists, AI researchers, ML and DL developers, AI application developers

## Top Features/Benefits

- Deep learning performance for training and inference with Intel optimized DL frameworks and tools
- Drop-in acceleration for data analytics and machine learning workflows with compute-intensive Python packages



Learn More: [software.intel.com/oneapi/ai-kit](https://software.intel.com/oneapi/ai-kit)

**Develop Fast Neural Networks on Intel® CPUs & GPUs**

**with Performance-optimized Building Blocks**

# Intel® oneAPI Deep Neural Network Library (oneDNN)



# Intel® oneAPI Deep Neural Network Library (oneDNN)

An **open-source cross-platform** performance library for deep learning applications

- Helps developers create high performance deep learning frameworks
- Abstracts out instruction set and other complexities of performance optimizations
- **Same API for both Intel CPUs and GPUs, use the best technology for the job**
- Supports Linux, Windows and macOS
- Open source for community contributions

More information as well as sources:

<https://github.com/oneapi-src/oneDNN>

# Intel® oneAPI Deep Neural Network Library

## Basic Information

- Features
  - API: C, C++, SYCL
  - Training: float32, bfloat16<sup>(1)</sup>
  - Inference: float32, bfloat16<sup>(1)</sup>, float16<sup>(1)</sup>, and int8<sup>(1)</sup>
  - MLPs, CNNs (1D, 2D and 3D), RNNs (plain, LSTM, GRU)
- Support Matrix
  - Compilers: Intel, GCC, CLANG, MSVC, DPC++
  - OS: Linux, Windows, macOS
  - CPU
    - Hardware: Intel® Atom, Intel® Core™, Intel® Xeon™
    - Runtimes: OpenMP, TBB, DPC++
  - GPU
    - Hardware: Intel HD Graphics, Intel® Iris® Plus Graphics
    - Runtimes: OpenCL, DPC++

	Intel® oneDNN
<b>Convolution</b>	2D/3D Direct Convolution/Deconvolution, Depthwise separable convolution 2D Winograd convolution
<b>Inner Product</b>	2D/3D Inner Production
<b>Pooling</b>	2D/3D Maximum 2D/3D Average (include/exclude padding)
<b>Normalization</b>	2D/3D LRN across/within channel, 2D/3D Batch normalization
<b>Eltwise (Loss/activation)</b>	ReLU(bounded/soft), ELU, Tanh; Softmax, Logistic, linear; square, sqrt, abs, exp, gelu, swish
<b>Data manipulation</b>	Reorder, sum, concat, View
<b>RNN cell</b>	RNN cell, LSTM cell, GRU cell
<b>Fused primitive</b>	Conv+ReLU+sum, BatchNorm+ReLU
<b>Data type</b>	f32, bfloat16, s8, u8

(1) Low precision data types are supported only for platforms where hardware acceleration is available



# Overview of Intel-optimizations for TensorFlow\*



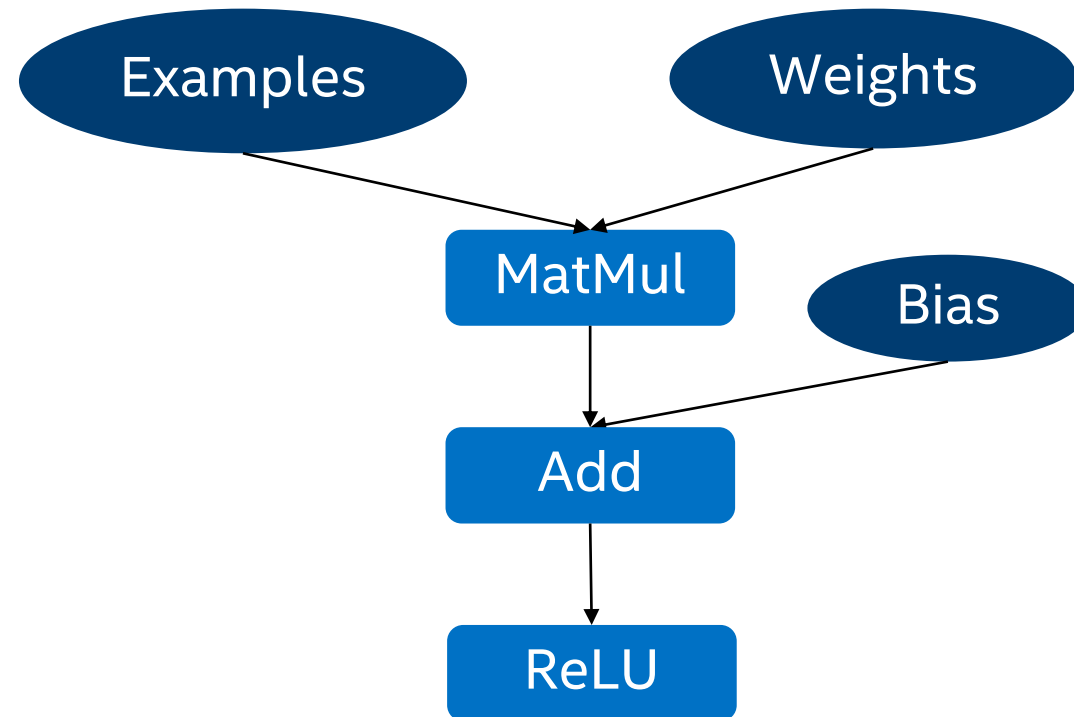
# Intel® TensorFlow\* optimizations

1. Operator optimizations: Replace default (Eigen) kernels by highly-optimized kernels (using Intel® oneDNN)
2. Graph optimizations: Fusion, Layout Propagation
3. System optimizations: Threading model

# Run TensorFlow\* benchmark

# Operator optimizations

In TensorFlow, computation graph is a data-flow graph.

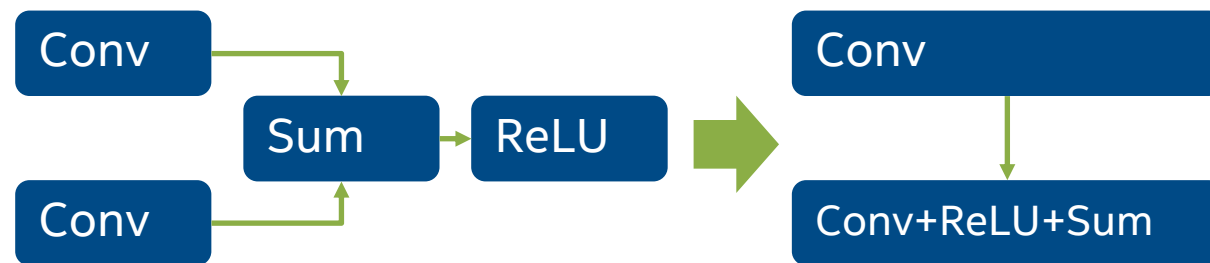


# Operator optimizations

- Replace default (Eigen) kernels by highly-optimized kernels (using Intel® oneDNN)
- Intel® oneDNN has optimized a set of TensorFlow operations.
- Library is open-source (<https://github.com/oneapi-src/oneDNN>) and downloaded automatically when building TensorFlow.

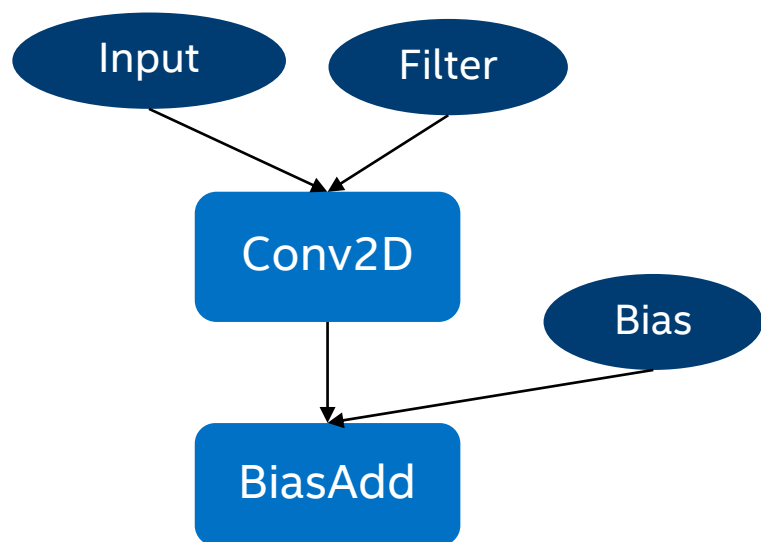
Forward	Backward
Conv2D	Conv2DGrad
Relu, TanH, ELU	ReLUGrad, TanHGrad, ELUGrad
MaxPooling	MaxPoolingGrad
AvgPooling	AvgPoolingGrad
BatchNorm	BatchNormGrad
LRN	LRNGrad
MatMul, Concat	

# Fusing computations

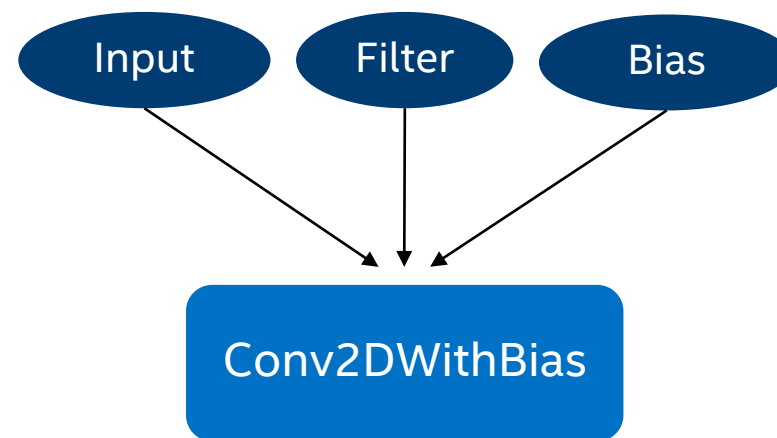


- On Intel processors a high percentation of time is typically spent in BW-limited ops
  - ~40% of ResNet-50, even higher for inference
- The solution is to fuse BW-limited ops with convolutions or one with another to reduce the # of memory accesses
  - Conv+ReLU+Sum, BatchNorm+ReLU, etc
- The frameworks are expected to be able to detect fusion opportunities
  - IntelCaffe already supports this

# Graph optimizations: fusion

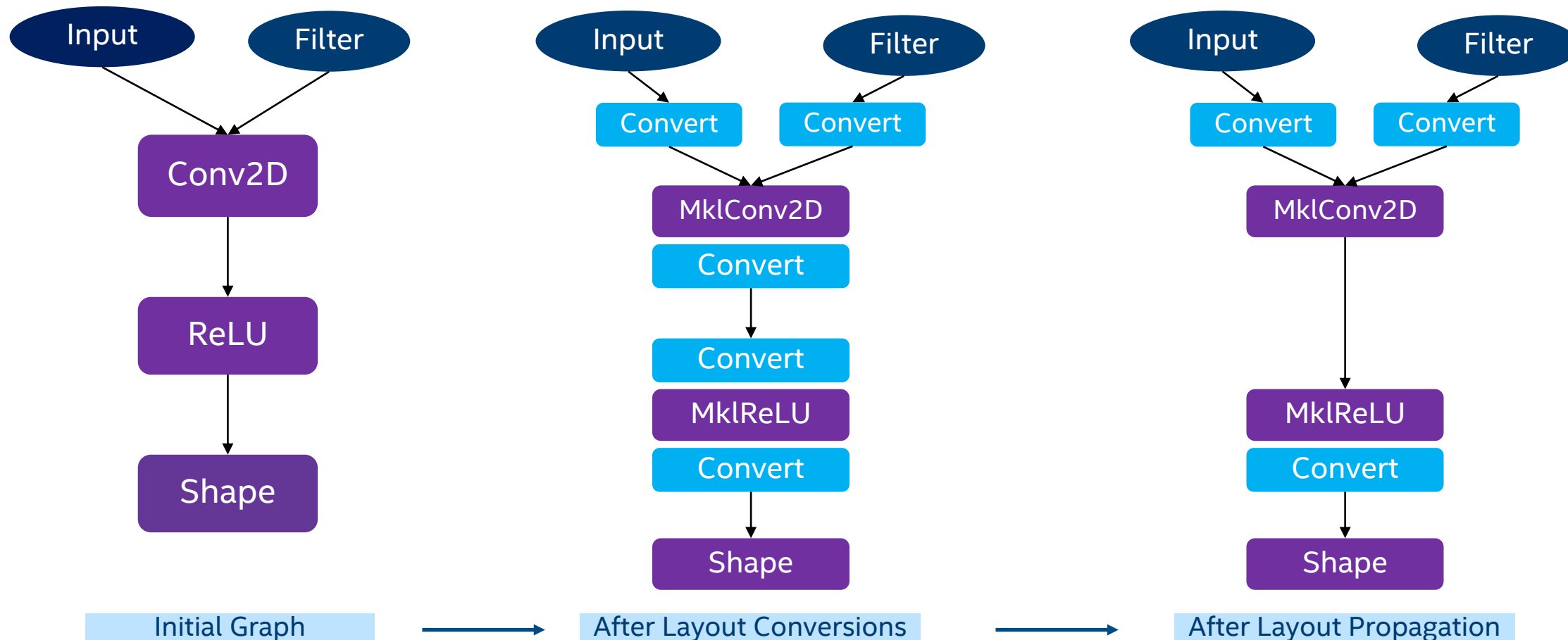


Before Merge



After Merge

# Graph optimizations: layout propagation



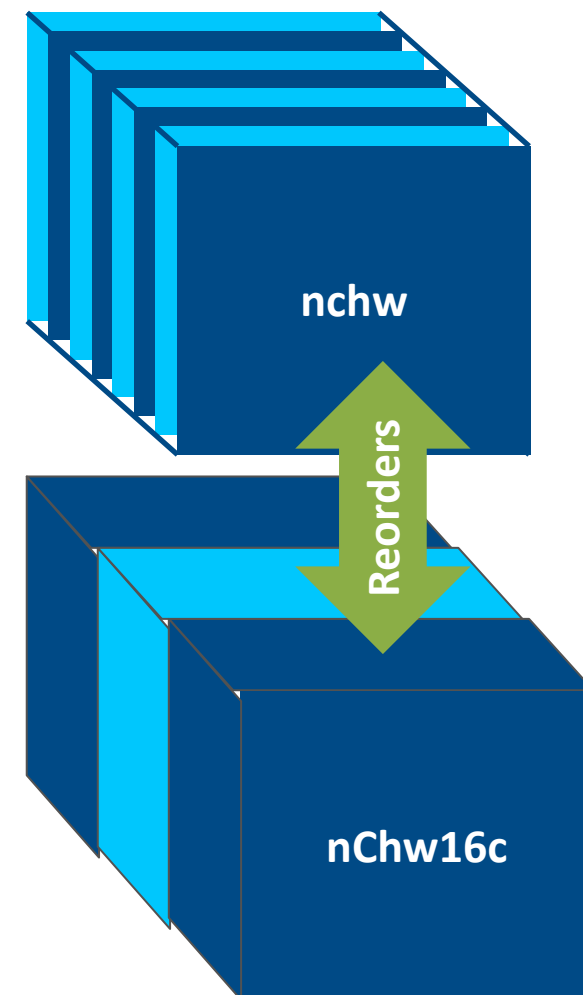
All oneDNN operators use highly-optimized layouts for TensorFlow tensors.



# More on memory channels: Memory layouts

- Most popular memory layouts for image recognition are **nhwc** and **nchw**
  - Challenging for Intel processors either for vectorization or for memory accesses (cache thrashing)
- Intel oneDNN convolutions use blocked layouts
  - Example: **nhwc** with channels blocked by 16 – **nChw16c**
  - Convolutions define which layouts are to be used by other primitives
  - Optimized frameworks track memory layouts and perform reorders **only** when necessary

More details: [https://oneapi-src.github.io/oneDNN/understanding\\_memory\\_formats.html](https://oneapi-src.github.io/oneDNN/understanding_memory_formats.html)



# Data Layout has a BIG Impact

- Continuous access to avoid gather/scatter
- Have iterations in inner most loop to ensure high vector utilization
- Maximize data reuse; e.g. weights in a convolution layer
- Overhead of layout conversion is sometimes negligible, compared with operating on unoptimized layout

21	18	32	6	3	
1	8	92	37	29	44
40	11	9	22	3	26
23	3	47	29	88	1
5	15	16	22	46	12
	29	9	13	11	1

21	18	...	1	..	8	92	..
----	----	-----	---	----	---	----	----

Channel based  
(NCHW)

21	8	18	92	..	1	11	..
----	---	----	----	----	---	----	----

Pixel based  
(NHWC)

```
for i= 1 to N # batch size
  for j = 1 to C # number of channels, image RGB = 3 channels
    for k = 1 to H # height
      for l = 1 to W # width
        dot_product( ...)
```

# System optimizations: load balancing

- TensorFlow graphs offer opportunities for parallel execution.
- Threading model
  1. **inter\_op\_parallelism\_threads** = max number of operators that can be executed in parallel
  2. **intra\_op\_parallelism\_threads** = max number of threads to use for executing an operator
  3. **OMP\_NUM\_THREADS** = oneDNN equivalent of **intra\_op\_parallelism\_threads**

# Performance Guide

- Maximize TensorFlow\* Performance on CPU: Considerations and Recommendations for Inference Workloads: <https://software.intel.com/en-us/articles/maximize-tensorflow-performance-on-cpu-considerations-and-recommendations-for-inference>

Example setting system environment variables with python `os.environ` :

```
os.environ["KMP_AFFINITY"] = "granularity=fine,compact,1,0"
```

```
os.environ["KMP_SETTINGS"] = "0"
```

## Tuning MKL for the best performance

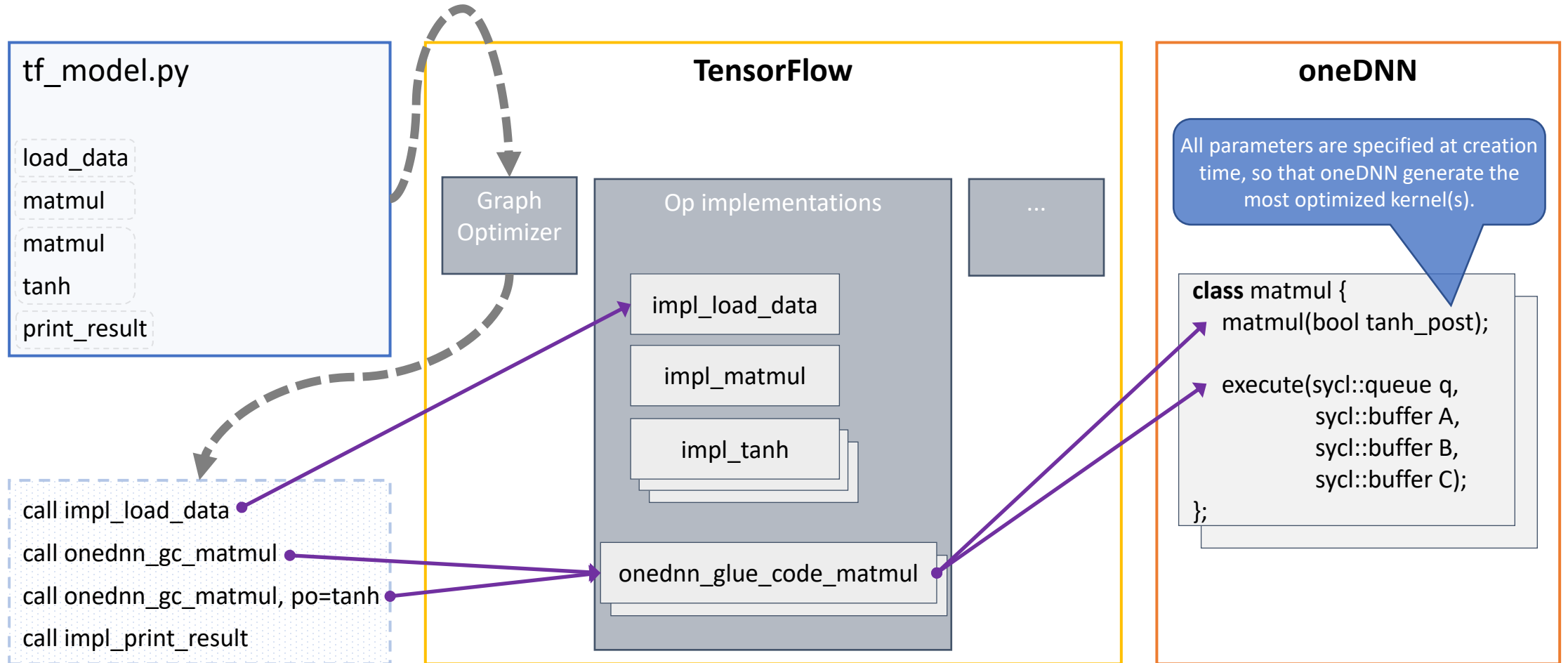
This section details the different configurations and environment variables that can be used to tune the MKL to get optimal performance. Before tweaking various environment variables make sure the model is using the `NCHW` ( `channels_first` ) [data format](#). The MKL is optimized for `NCHW` and Intel is working to get near performance parity when using `NHWC`.

MKL uses the following environment variables to tune performance:

- `KMP_BLOCKTIME` - Sets the time, in milliseconds, that a thread should wait, after completing the execution of a parallel region, before sleeping.
- `KMP_AFFINITY` - Enables the run-time library to bind threads to physical processing units.
- `KMP_SETTINGS` - Enables (true) or disables (false) the printing of OpenMP\* run-time library environment variables during program execution.
- `OMP_NUM_THREADS` - Specifies the number of threads to use.

Intel Tensorflow\* install guide is available →  
<https://software.intel.com/en-us/articles/intel-optimization-for-tensorflow-installation-guide>

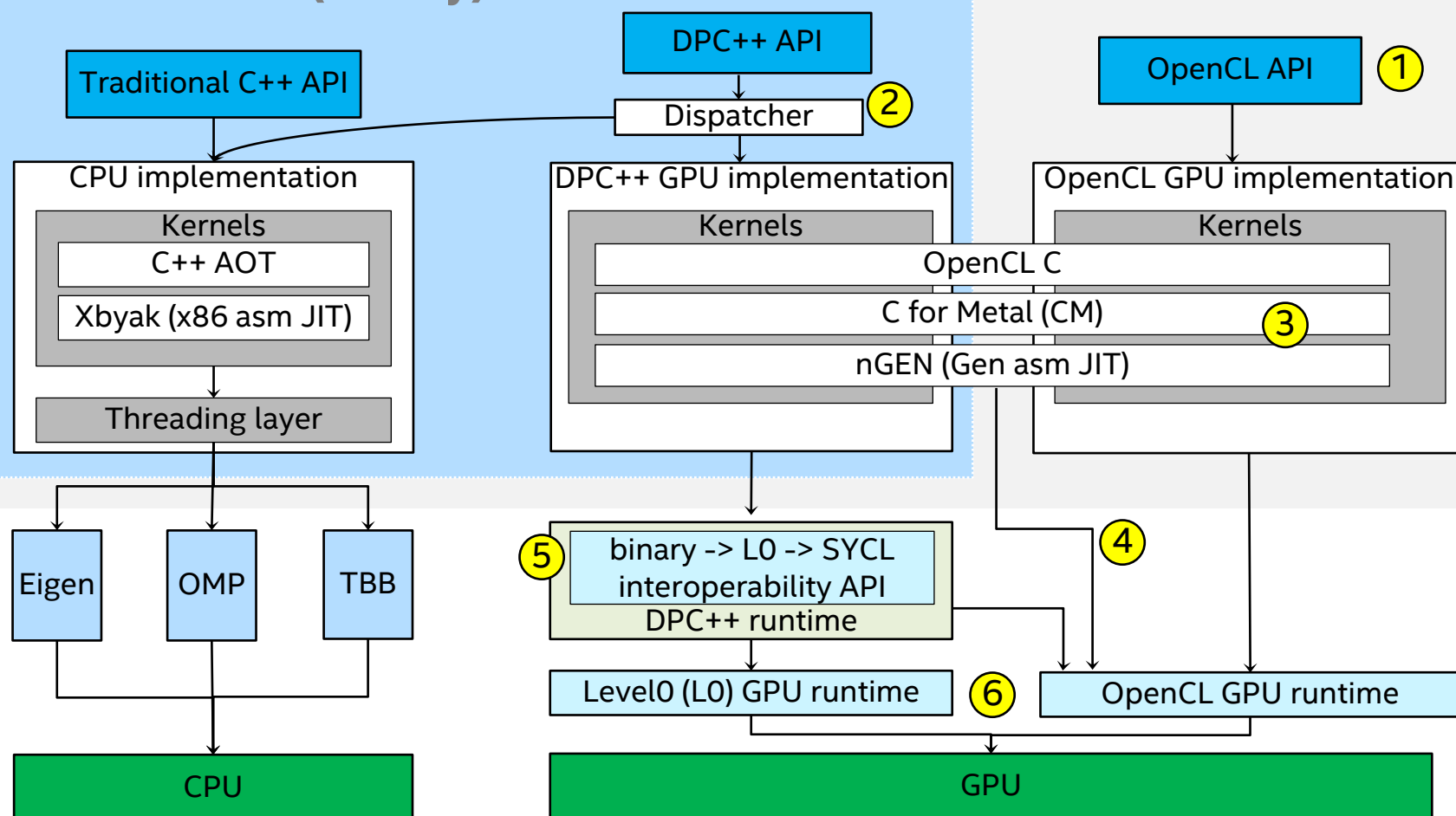
# oneDNN <-> Frameworks interaction



# oneDNN architecture overview

oneDNN (open source)

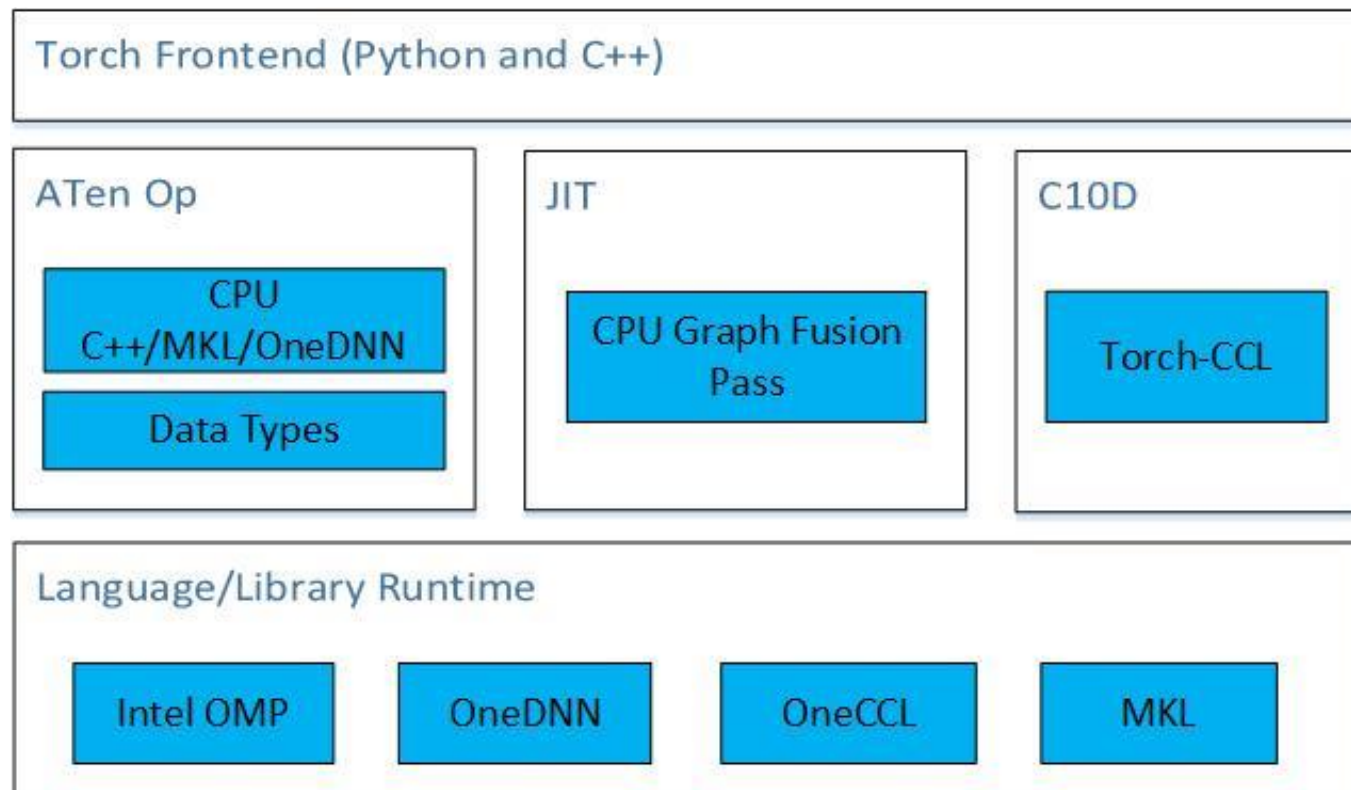
## Intel oneDNN (binary)



- 1 OpenCL API is not available as part of Intel oneAPI binary distribution
- 2 Dispatching between CPU and GPU is based on the kind of device associated with the DPC++ queue
- 3 All GPU kernels are compiled in runtime. CM and nGEN support is not available publicly yet. Adding/migrating to DPC++ kernels is under consideration
- 4 OpenCL GPU RT is always needed to compile OpenCL C and CM kernels
- 5 In case of DPC++ and LO, binary kernels need to be wrapped to LO modules to create SYCL kernels eventually
- 6 Under DPC++ API/runtime, users can run on GPU via either OpenCL or LO GPU runtime: it should be specified in compile time, but can be checked during execution time

# Intel Optimizations for PyTorch

- Accelerated operators
- Graph optimization
- Accelerated communications



# Motivation for Intel Extension for PyTorch (IPEX)

- Provide customers with the up-to-date Intel software/hardware features
- Streamline the work to enable Intel accelerated library



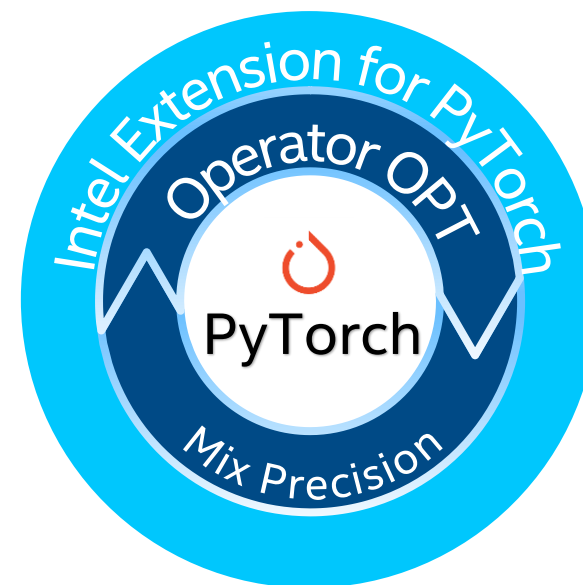
## Operator Optimization

- Auto dispatch the operators optimized by the extension backend
- Auto operator fusion via PyTorch graph mode



## Mix Precision

- Accelerate PyTorch operator by bfloat16
- Automatic mixed precision





# PyTorch-IPEX Demo

# How to get IPEX

1. oneAPI AI Analytics Toolkit
2. Install from source

# IPEX from the oneAPI AI Analytics Toolkit

## Intel Optimizations for PyTorch

### Intel-Optimized PyTorch

- PyTorch back-end optimizations
- Up-streamed to regular PyTorch
- Same front-end code as regular PyTorch

### Intel Extension for PyTorch (IPEX)

- Additional optimizations and Mixed Precision support
- Different front-end

### Torch-CCL

- For distributed learning
- PyTorch bindings for oneCCL

# Installing IPEX from source

<https://github.com/intel/intel-extension-for-pytorch>

License - Apache 2.0

## Build and install

1. Install PyTorch from source
2. Download and install Intel PyTorch Extension source
3. Add new backend for Intel Extension for PyTorch
4. Install Intel Extension for PyTorch

# Automatic Mixed Precision Feature (FP32 + BF16)

```
import torch
import intel_pytorch_extension as ipex
ipex.enable_auto_optimization(mixed_dtype = torch.bfloat16, train = True)

EPOCH = 20
BATCH_SIZE = 128
LR = 0.001

def main():
    train_loader = ...
    test_loader = ...
    net = topology()
    net = net.to(ipex.DEVICE)
    criterion = torch.nn.CrossEntropyLoss()
    optimizer = torch.optim.SGD(net.parameters(), lr = LR, momentum=0.9)
    for epoch in range(EPOCH):
        net.train()
        for batch_idx, (data, target) in enumerate(train_loader):
            data = data.to(ipex.DEVICE)
            target = target.to(ipex.DEVICE)
            optimizer.zero_grad()
            output = net(data)
            loss = criterion(output, target)
            loss.backward()
            optimizer.step()

        net.eval()
        test_loss = 0
        correct = 0
        with torch.no_grad():
            for data, target in test_loader:
                data = data.to(ipex.DEVICE)
                target = target.to(ipex.DEVICE)
                output = net(data)
                test_loss += criterion(output, target, reduction='sum').item()
                pred = output.argmax(dim=1, keepdim=True)
                correct += pred.eq(target.view_as(pred)).sum().item()
        test_loss /= len(test_loader.dataset)

    if __name__ == '__main__':
        main()
```

1. import ipex

2. Enable Auto-Mix-Precision by API

\* Subject to change

3. Convert the input tensors to the extension device

4. Convert the model to the extension device

# Data types

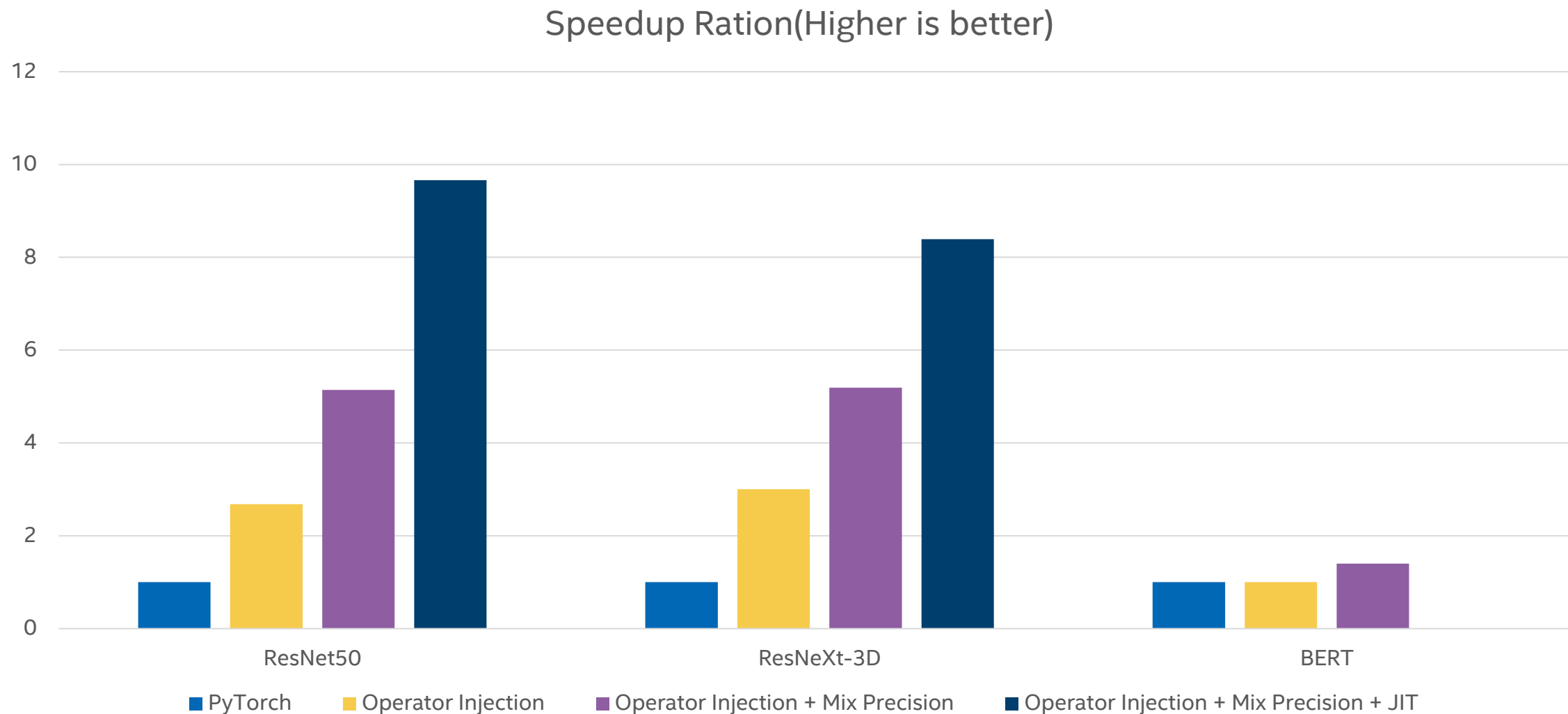


<https://software.intel.com/sites/default/files/managed/40/8b/bf16-hardware-numerics-definition-white-paper.pdf?source=techstories.org>

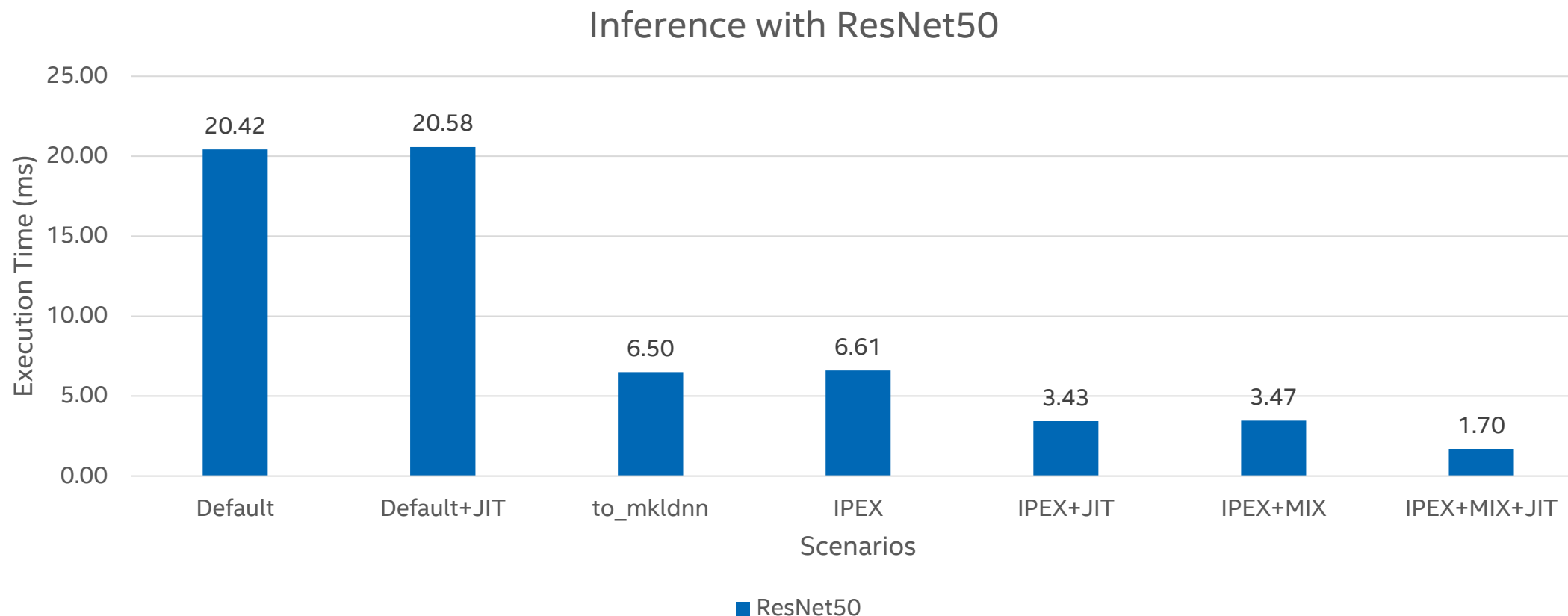
- Benefit of bfloat16
  - Performance 2x up
  - Comparable accuracy loss against fp32
  - No loss scaling, compared to fp16

\* bfloat16 intrinsic support starts from 3rd Generation Intel® Xeon® Scalable Processors

# Extension Performance comparison



# Inference with IPEX for ResNet50



Worker11 (CPX)

LD\_PRELOAD=/root/anaconda3/lib/libiomp5.so OMP\_NUM\_THREADS=26 KMP\_AFFINITY=granularity=fine,compact,1,0 numactl -N 0 -m 0 python resnet50.py



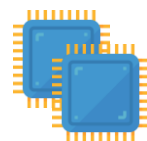
# Intel Low Precision Optimization Tool Tutorial



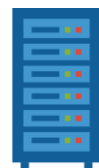
# The motivation for low precision



**Lower  
Power**



**Lower memory  
bandwidth**



**Lower  
storage**



**Higher  
performance**

**Important:  
Acceptable  
accuracy loss**

The key term:

- Quantization

# Quantization in a nutshell

**Floating Point**

96.1924



**Integer**

96

32 -bit

8 bit

10110110

10110110



10110110

10110110

10110110

# Challenge & Solution of Low Precision Optimization Tool (for Inferencing in Deep Learning)

- Low Precision Inference can speed up the performance by reducing the computing, memory and storage of AI model.
- Intel provides solution to cover the challenge of it:

Challenge	Intel Solution	How
Hardware support	Intel® Deep Learning Boost supported by the Second-Generation Intel® Xeon® Scalable Processors and later.	VNNI intrinsic. Support INT8 MulAdd.
Complex to convert the FP32 model to INT8/BF16 model	Intel® Low Precision Optimization Tool (LPOT)	Unified quantization API
Accuracy loss in converting to INT8 model	Intel® Low Precision Optimization Tool (LPOT)	Auto tuning

# Product Definition

- Convert the FP32 model to INT8/BF16 model. Optimize the model in same time.
- Support multiple Intel optimized DL frameworks (TensorFlow, PyTorch, MXNet) on both CPU and GPU.
- Support automatic accuracy-driven tuning, along with additional custom objectives like performance, model size, or memory footprint
- Provide the easy extension capability for new backends (e.g., PDPD, ONNX RT) and new tuning strategies/metrics (e.g., HAWQ from UCB)

# Tuning Zoo

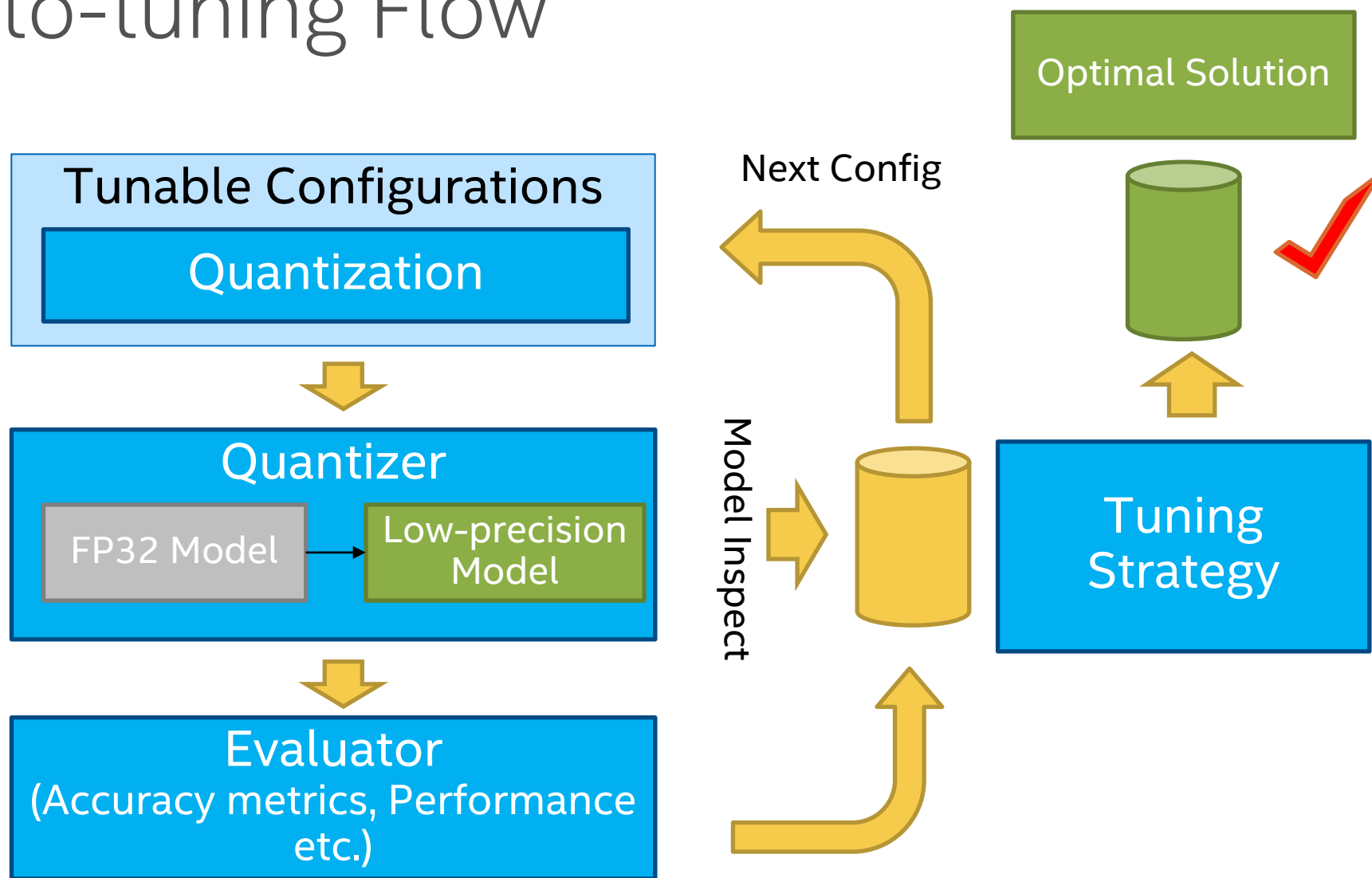
The followings are the models supported by Intel® Low Precision Optimization Tool for auto tuning.

TensorFlow Model	Category
<a href="#">ResNet50 V1</a>	Image Recognition
<a href="#">ResNet50 V1.5</a>	Image Recognition
<a href="#">ResNet101</a>	Image Recognition
<a href="#">Inception V1</a>	Image Recognition
<a href="#">Inception V2</a>	Image Recognition
<a href="#">Inception V3</a>	Image Recognition
<a href="#">Inception V4</a>	Image Recognition
<a href="#">ResNetV2_50</a>	Image Recognition
<a href="#">ResNetV2_101</a>	Image Recognition
<a href="#">ResNetV2_152</a>	Image Recognition
<a href="#">Inception ResNet V2</a>	Image Recognition
<a href="#">SSD ResNet50 V1</a>	Object Detection
<a href="#">Wide &amp; Deep</a>	Recommendation
<a href="#">VGG16</a>	Image Recognition
<a href="#">VGG19</a>	Image Recognition
<a href="#">Style transfer</a>	Style Transfer

PyTorch Model	Category
<a href="#">BERT-Large RTE</a>	Language Translation
<a href="#">BERT-Large QNLI</a>	Language Translation
<a href="#">BERT-Large CoLA</a>	Language Translation
<a href="#">BERT-Base SST-2</a>	Language Translation
<a href="#">BERT-Base RTE</a>	Language Translation
<a href="#">BERT-Base STS-B</a>	Language Translation
<a href="#">BERT-Base CoLA</a>	Language Translation
<a href="#">BERT-Base MRPC</a>	Language Translation
<a href="#">DLRM</a>	Recommendation
<a href="#">BERT-Large MRPC</a>	Language Translation
<a href="#">ResNext101_32x8d</a>	Image Recognition
<a href="#">BERT-Large SQUAD</a>	Language Translation
<a href="#">ResNet50 V1.5</a>	Image Recognition
<a href="#">ResNet18</a>	Image Recognition
<a href="#">Inception V3</a>	Image Recognition
<a href="#">YOLO V3</a>	Object Detection
<a href="#">Peleenet</a>	Image Recognition
<a href="#">ResNest50</a>	Image Recognition
<a href="#">SE_ResNext50_32x4d</a>	Image Recognition
<a href="#">ResNet50 V1.5 QAT</a>	Image Recognition

MxNet Model	Category
<a href="#">ResNet50 V1</a>	Image Recognition
<a href="#">MobileNet V1</a>	Image Recognition
<a href="#">MobileNet V2</a>	Image Recognition
<a href="#">SSD-ResNet50</a>	Object Detection
<a href="#">SqueezeNet V1</a>	Image Recognition
<a href="#">ResNet18</a>	Image Recognition
<a href="#">Inception V3</a>	Image Recognition

# Auto-tuning Flow





# System Requirements

## ■ Hardware

Intel® Low Precision Optimization Tool supports systems based on Intel 64 architecture or compatible processors.

The quantization model could get acceleration by Intel® Deep Learning Boost if running on the Second-Generation Intel® Xeon® Scalable Processors and later:

Verified:

- Cascade Lake & Cooper Lake, with Intel DL Boost VNNI
- Skylake, with AVX-512 INT8

## ■ OS: Linux

Verified: CentOS 7.3 & Ubuntu 18.04

## ■ Software

Intel® Low Precision Optimization Tool requires to install Intel optimized framework version for TensorFlow, PyTorch, and MXNet.

Verified Release	Installation Example
Intel Optimization for TensorFlow: v1.15 (up1), v2.1, v2.2, v2.3	<code>pip install intel-tensorflow==2.3.0</code>
PyTorch: v1.5	<code>pip install torch==1.5.0+cpu*****</code>
MXNet: v1.6, v1.7	<code>pip install mxnet-mkl==1.6.0</code>

# Installation

- **Install from Intel AI Analytics Toolkit (Recommended)**

```
source /opt/intel/oneapi/setvars.sh
```

```
conda activate tensorflow
```

```
cd /opt/intel/oneapi/iLiT/latest
```

```
sudo ./install_iLiT.sh
```

- **Install from source**

```
git clone https://github.com/intel/lpot.git
```

```
cd lpot
```

```
python setup.py install
```

- **Install from binary**

```
# install from pip
```

```
pip install lpot
```

```
# install from conda
```

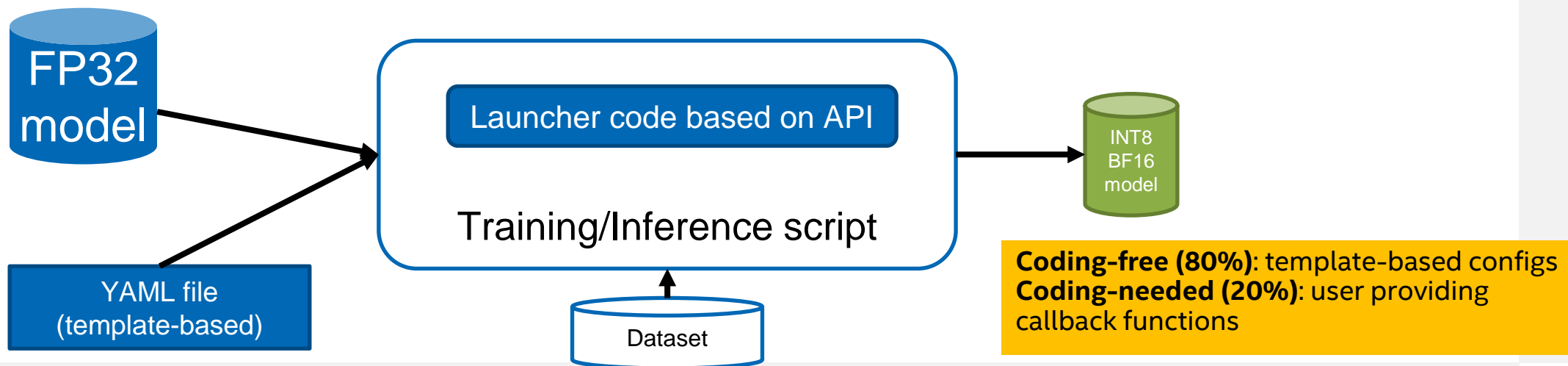
```
conda install lpot -c intel -c conda-forge
```

For more detailed installation info, please refer to <https://github.com/intel/lpot>

# Usage: Simple Python API + YAML config

LPOT is designed to reduce the workload of the user and keep the flexibility.

Python API	YAML
<ul style="list-style-type: none"><li>Simple API is easy to integrated in original training/inference script.</li></ul>	<ul style="list-style-type: none"><li>Common functions are integrated and controlled by parameters;</li><li>Templates are easy to refer;</li><li>Lots of advance parameters provide powerful tuning capability.</li></ul>



# Python API

- Core User-facing API:

- Quantization()

- Follow a specified tuning strategy to tune a low precision model through QAT or PTQ which can meet pre-defined accuracy goal and objective.

```
class Quantization(object):  
    def __init__(self, conf_fname):  
        ...  
  
    def __call__(self, model, q_data_loader=None, q_func=None,  
                 eval_data_loader=None, eval_func=None):  
        ...
```

# Intel LPOT YAML Configure

Intel LPOT YAML config consists of 6 building blocks:

- ❑ model
- ❑ device
- ❑ quantization
- ❑ evaluation
- ❑ tuning

```
# ilit yaml building block
model:          # model specific info, such as model name, framework,
                 input/output node name required for tensorflow.
    ...

device: ...     # the device ilit runs at, cpu or gpu. default is cpu.

quantization:   # the setting of calibration/quantization behavior. only
                 required for PTQ and QAT.
    ...

evaluation:     # the setting of how to evaluate a model.
    ...

tuning:         # the tuning behavior, such as strategy, objective, accuracy
                 criterion.
    ...
```

# Easy: TensorFlow ResNet50

model:

```
name: resnet50_v1_5
framework: tensorflow
inputs: input_tensor
outputs: softmax_tensor
```

YAML config

quantization:

```
calibration:
  sampling_size: 50, 100
dataloader:
  batch_size: 10
  dataset:
    Imagenet:
      root: /path/to/calibration/dataset
  transform:
    ParseDecodeImagenet:
    ResizeCropImagenet:
      height: 224
      width: 224
      mean_value: [123.68, 116.78, 103.94]
```

evaluation:

```
accuracy:
  metric:
    topk: 1
dataloader:
  batch_size: 32
  dataset:
    Imagenet:
      root: /path/to/evaluation/dataset
  transform:
    ParseDecodeImagenet:
    ResizeCropImagenet:
      height: 224
      width: 224
      mean_value: [123.68, 116.78, 103.94]
```

tuning:

```
accuracy_criterion:
  relative: 0.01
exit_policy:
  timeout: 0
random_seed: 9527
```

```
from lpot import Quantization
quantizer = Quantization("./conf.yaml")
q_model = quantizer(model)
```

Code change

Full example:

[https://github.com/intel/lpot/tree/master/examples/tensorflow/image\\_recognition](https://github.com/intel/lpot/tree/master/examples/tensorflow/image_recognition)

# DEMO

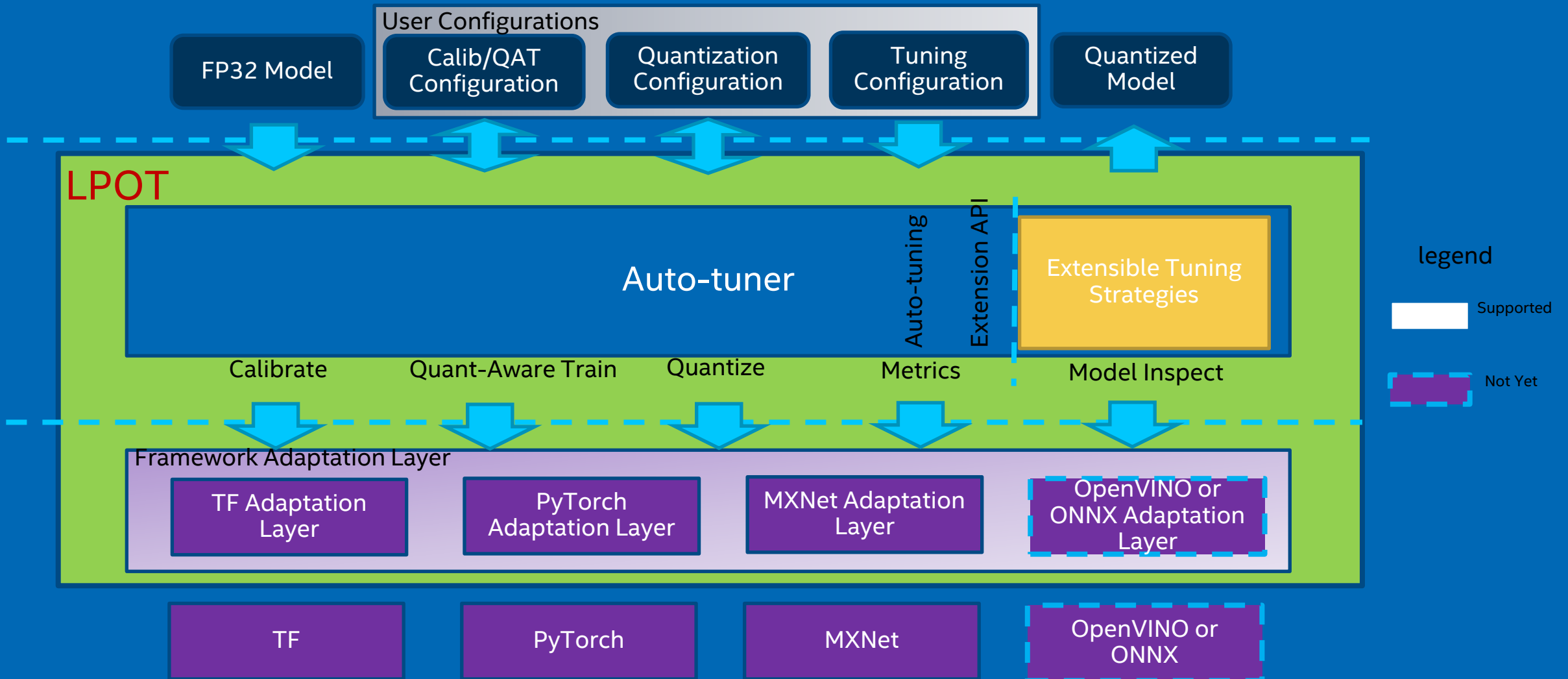
# Demo

- Intel AI Analytics Toolkit Samples:
- <https://github.com/oneapi-src/oneAPI-samples/tree/master/AI-and-Analytics>
- Intel LPOT Sample for Tensorflow: [-samples](#)
- <https://github.com/oneapi-src/oneAPI-samples/tree/master/AI-and-Analytics/Getting-Started-Samples/LPOT-Sample-for-Tensorflow>

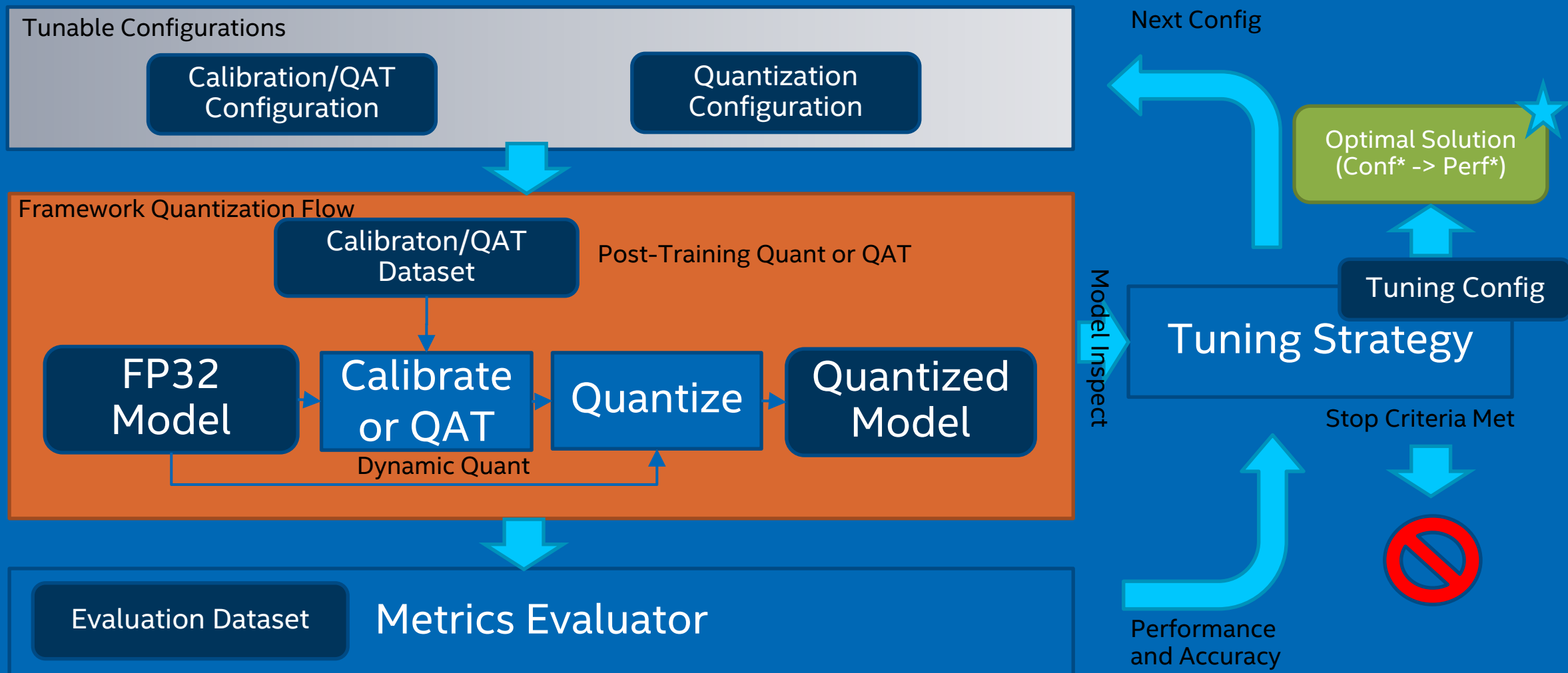




# Infrastructure



# Working Flow





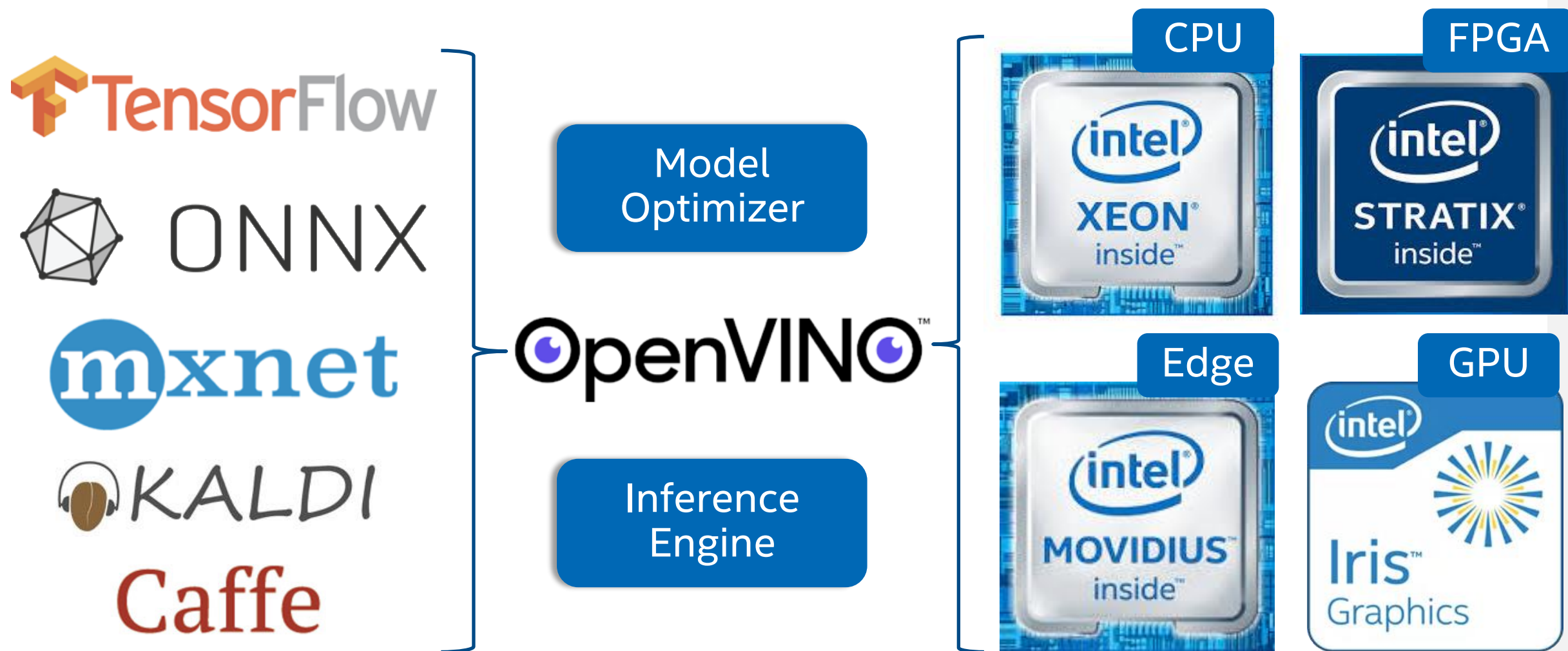
Product Overview

# Intel® Distribution of OpenVINO™ toolkit



2019 Developer Tool of the Year  
Awarded by the Edge AI and Vision Alliance

# WRITE once, deploy & scale diversely



\*Other names and brands may be claimed as the property of others.

# From a bird's eye-view

Advanced capabilities to streamline deep learning deployments

## 1. Build

### Trained Model

TensorFlow Caffe  
KALDI mxnet  
ONNX



### Open Model Zoo

100+ open sourced and optimized pre-trained models;  
80+ supported public models

## 2. Optimize



### Model Optimizer

Converts and optimizes trained model using a supported framework

Read, Load, Infer



Intermediate Representation  
(.xml, .bin)

Post-Training Optimization Tool

Deep Learning Streamer

OpenCV

OpenCL™

Deep Learning Workbench

Code Samples & Demos

(e.g. Benchmark app, Accuracy Checker, Model Downloader)

## 3. Deploy



### Inference Engine

Common API that abstracts low-level programming for each hardware

Deployment Manager

CPU Plugin

GPU Plugin

GNA Plugin

Myriad Plugin

For Intel® NCS2 & NCS

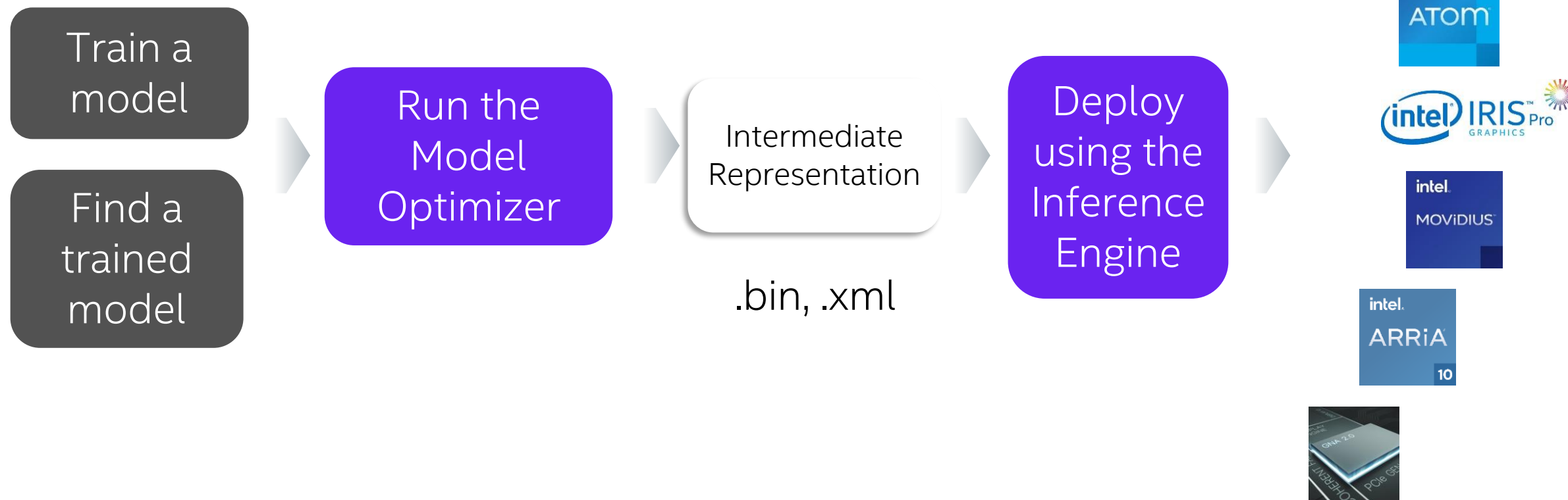
HDDL Plugin

FGPA Plugin



# Get Started

Typical workflow from development to deployment



# Supported Frameworks

Breadth of supported frameworks to enable developers with flexibility



**Supported Frameworks and Formats** ▶ [https://docs.openvinotoolkit.org/latest/\\_docs\\_IE\\_DG\\_Introduction.html#SupportedFW](https://docs.openvinotoolkit.org/latest/_docs_IE_DG_Introduction.html#SupportedFW)

**Configure the Model Optimizer for your Framework** ▶ [https://docs.openvinotoolkit.org/latest/\\_docs\\_MO\\_DG\\_prepare\\_model\\_Config\\_Model\\_Optimizer.html](https://docs.openvinotoolkit.org/latest/_docs_MO_DG_prepare_model_Config_Model_Optimizer.html)



# Core Components

## Model optimization to deployment

### Model Optimizer



- A Python-based tool to **import** trained models and **convert** them to Intermediate Representation
- **Optimizes for performance** or space with conservative topology transformations
- **Hardware-agnostic** optimizations

#### Development Guide ►

[https://docs.openvino toolkit.org/latest/\\_docs\\_MO\\_DG\\_Deep\\_Learning\\_Model\\_Optimizer\\_DevGuide.html](https://docs.openvino toolkit.org/latest/_docs_MO_DG_Deep_Learning_Model_Optimizer_DevGuide.html)

### Inference Engine



- High-level, C, C++ and Python, inference **runtime API**
- Interface is implemented as **dynamically loaded plugins** for each hardware type
- Delivers best performance for each type **without requiring** users to implement and maintain multiple code pathways

#### Development Guide ►

[https://docs.openvino toolkit.org/latest/\\_docs\\_IE\\_DG\\_Deep\\_Learning\\_Inference\\_Engine\\_DevGuide.html](https://docs.openvino toolkit.org/latest/_docs_IE_DG_Deep_Learning_Inference_Engine_DevGuide.html)

# Model Optimization

Breadth of supported frameworks to enable developers with flexibility

**Model Optimizer** loads a model into memory, reads it, builds the internal representation of the model, optimizes it, and produces the **Intermediate Representation**.

Optimization techniques available are:

- Linear operation fusing
- Stride optimizations
- Group convolutions fusing

*Note:* Except for ONNX (.onnx model formats), all models have to be converted to an IR format to use as input to the Inference Engine



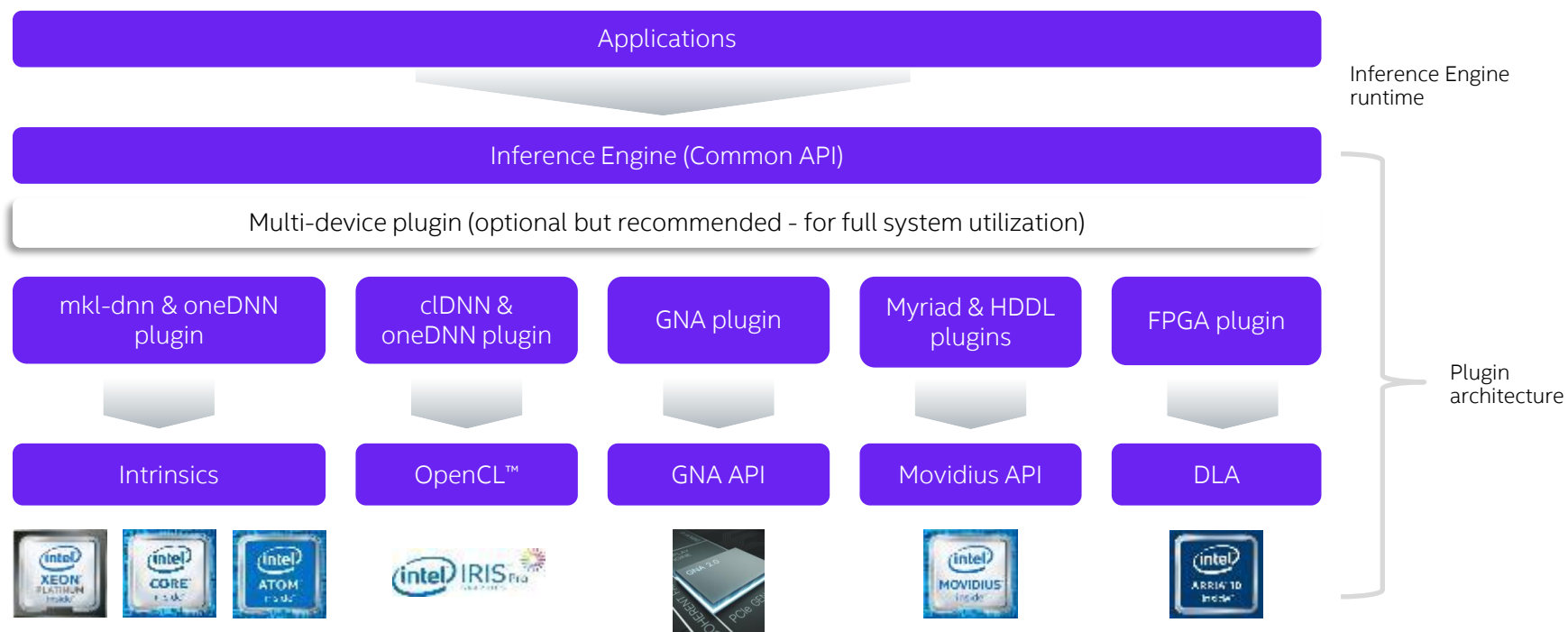
**.xml** – describes the network topology

**.bin** – describes the weights and biases binary data

# Inference Engine

Common high-level inference runtime for cross-platform flexibility

OpenVINO™

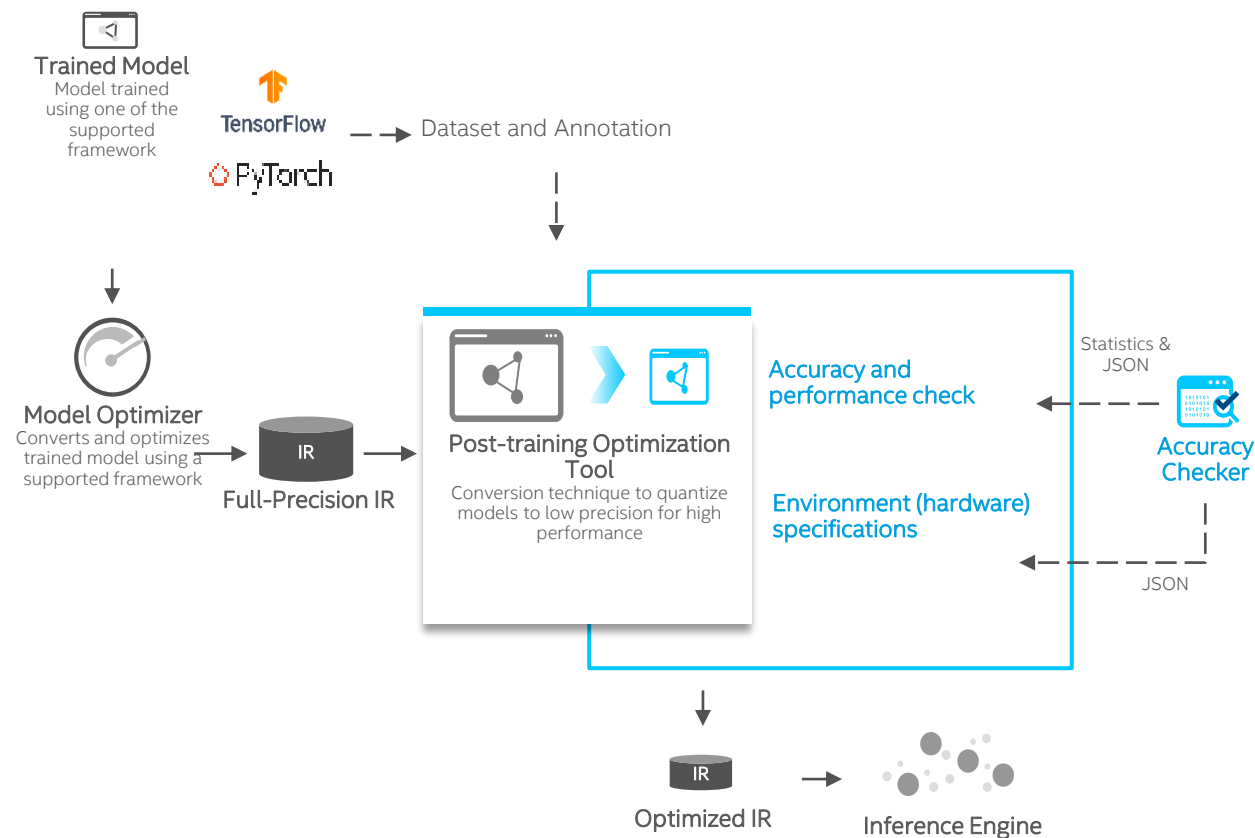


# Post-Training Optimization Tool

Conversion technique that reduces model size into low-precision without re-training

Reduces model size **while also improving latency**, with **little degradation** in model accuracy and without model re-training.

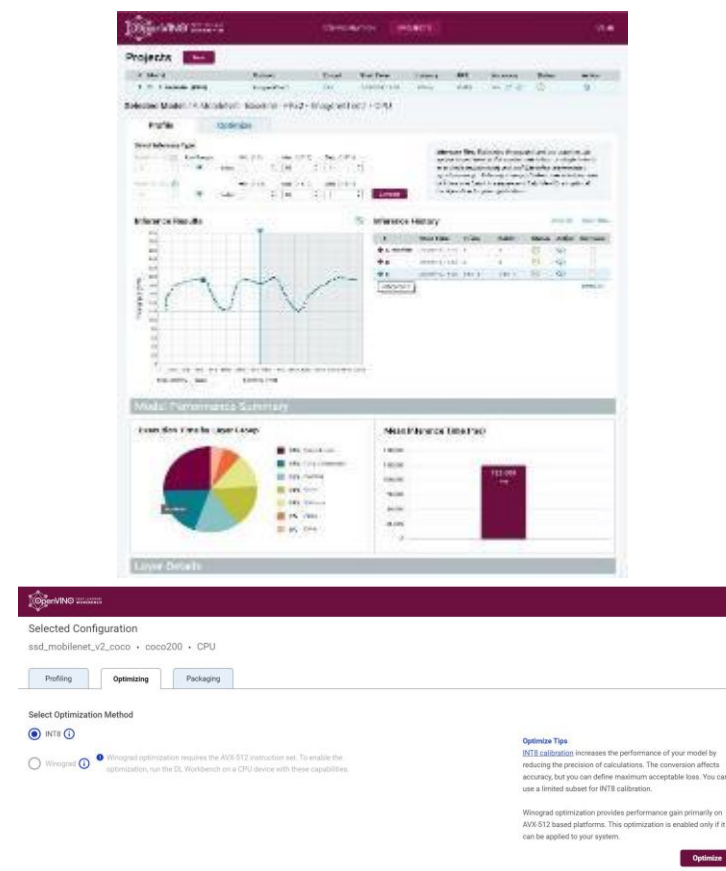
Different optimization approaches are supported: quantization algorithms, sparsity, etc.



# Deep Learning Workbench

Web-based UI extension tool for model analyses and graphical measurements

- Visualizes performance data for topologies and layers to aid in model analysis
- Automates analysis for optimal performance configuration (streams, batches, latency)
- Experiment with INT8 or Winograd calibration for optimal tuning using the Post Training Optimization Tool
- Provide accuracy information through accuracy checker
- Direct access to models from public set of Open Model Zoo
- Enables **remote profiling**, allowing the collection of performance data from multiple different machines without any additional set-up.



# Additional Tools and Add-ons

Streamlined development experience and ease of use



Model Downloader

- Provides an easy way of accessing a number of public models as well as a set of pre-trained Intel models



Deployment Manager

- Generate an optimal, minimized runtime package for deployment
- Deploy with smaller footprint compared to development package



Benchmark App

- Measure performance (throughput, latency) of a model
- Get performance metrics per layer and overall basis



Accuracy Checker

- Check for accuracy of the model (original and after conversion) to IR file using a known data set

## Computer Vision Annotation Tool

This web-based tool helps annotate videos and images before training a model

## Deep Learning Streamer

Streaming analytics framework to create and deploy complex media analytics pipelines

## OpenVINO™ Model Server

Scalable inference server for serving optimized models and applications

## Dataset Management Framework

Use this add-on to build, transform and analyze datasets

## Neural Network Compression Framework

Training framework based on PyTorch\* for quantization-aware training

## Training Extensions

Trainable deep learning models for training with custom data

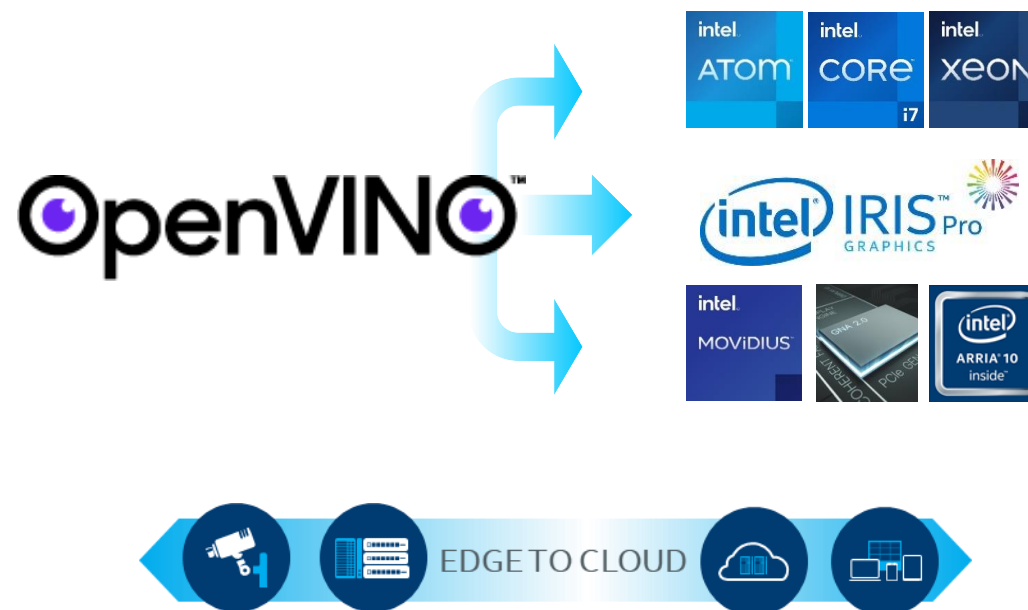
# Write Once, Deploy Anywhere

Common high-level inference runtime for cross-platform flexibility

Write once, deploy across different platforms with the same API and framework-independent execution

Consistent accuracy, performance and functionality across all target devices with no re-training required

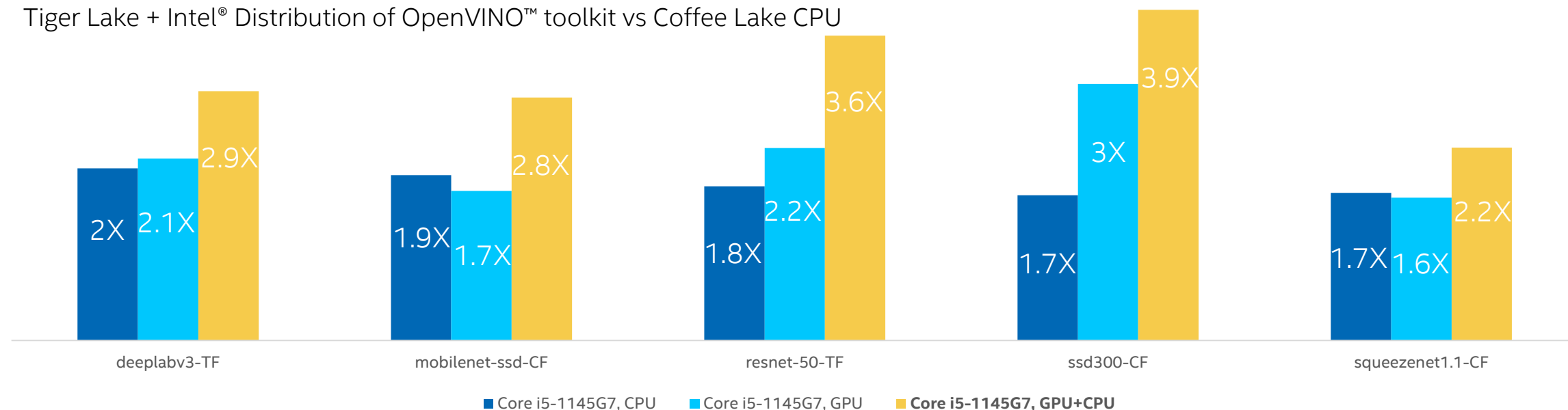
Full environment utilization, or multi-device plugin, across available hardware for greater performance results



# Compounding Effect of Hardware and Software

Use Intel® Xe Graphics + CPU combined for maximum inferencing

Tiger Lake + Intel® Distribution of OpenVINO™ toolkit vs Coffee Lake CPU



11<sup>th</sup> Gen Intel® Core™ (Tiger Lake) Core i5-1145G7 relative inference FPS compared to Coffee Lake, Core i5-8500

Using the Multi-device plugin

The above is preliminary performance data based on pre-production components. For more complete information about performance and benchmark results, visit [www.intel.com/benchmarks](https://www.intel.com/benchmarks). See backup for configuration details.



# Pre-Trained Models and Public Models

Open-sourced repository of pre-trained models and support for public models

Use free **Pre-trained Models** to speed up development and deployment

Take advantage of the **Model Downloader** and other automation tools to quickly get started

Iterate with the **Accuracy Checker** to validate the accuracy of your models

## 100+ Pre-trained Models

*Common AI tasks*

Object Detection  
Object Recognition  
Reidentification  
Semantic Segmentation  
Instance Segmentation  
Human Pose Estimation  
Image Processing  
Text Detection  
Text Recognition  
Text Spotting  
Action Recognition  
Image Retrieval  
Compressed Models  
Question Answering

## 100+ Public Models

*Pre-optimized external models*

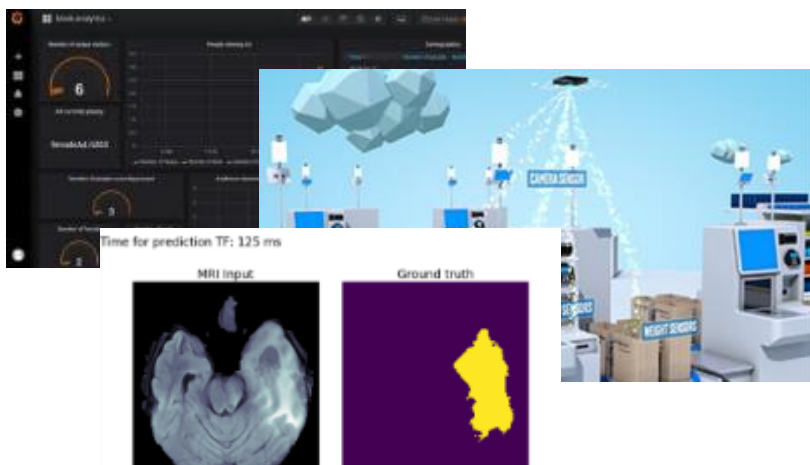
Classification  
Segmentation  
Object Detection  
Human Pose Estimation  
Monocular Depth Estimation  
Image Inpainting  
Style Transfer  
Action Recognition  
Colorization

# DEMO

# Demos and Reference Implementations

Quickly get started with example demo applications and reference implementations

Take advantage of **pre-built, open-sourced** example implementations with step-by-step guidance and required components list



[Face Access Control - C++](#)

[Parking Lot Counter - Go](#)

[Intruder Detector - C++](#)

[People Counter - C++](#)

[Machine Operator Monitor - C++](#)

[Restricted Zone Notifier - Go](#)

[Machine Operator Monitor - Go](#)

[Shopper Gaze Monitor - C++](#)

[Motor Defect Detector - Python](#)

[Shopper Mood Monitor - Go](#)

[Object Flaw Detector - C++](#)

[Store Aisle Monitor - C++](#)

[Object Size Detector - C++](#)

[Store Traffic Monitor - C++](#)

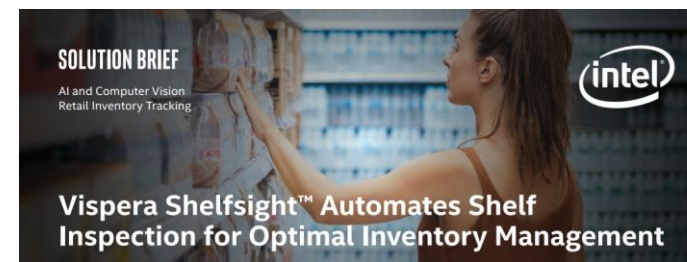
[Object Size Detector - Go](#)

[Store Traffic Monitor - Python](#)

[Parking Lot Counter - C++](#)

# Case Studies

Use cases and successful implementation across a variety of industries powered by the Intel® Distribution of OpenVINO™ toolkit



**JLK INSPECTION Healthcare Access and Quality** [Solution Brief](#)

Reduced average inference time on Intel® NUC (with no GPU) from **4.23 seconds to just 2.81 seconds**, which helps medical professionals reach more people, accelerate screening and help improve quality of care.

**ZEROFOX Security Against Social and Digital Attacks** [Solution Brief](#)

Performance improvements of up to **2.3x faster**, reducing latency by up to **50 percent** for threat detection and remediation to protect businesses against targeted social and digital attacks.

**dc water is life® Sewer pipe inspection analysis** [Solution Brief](#)

Inference time was improved with a reduction of up to **80%** using Intel Xeon processors with the OpenVINO toolkit, while not producing significant loss in model precision or accuracy.

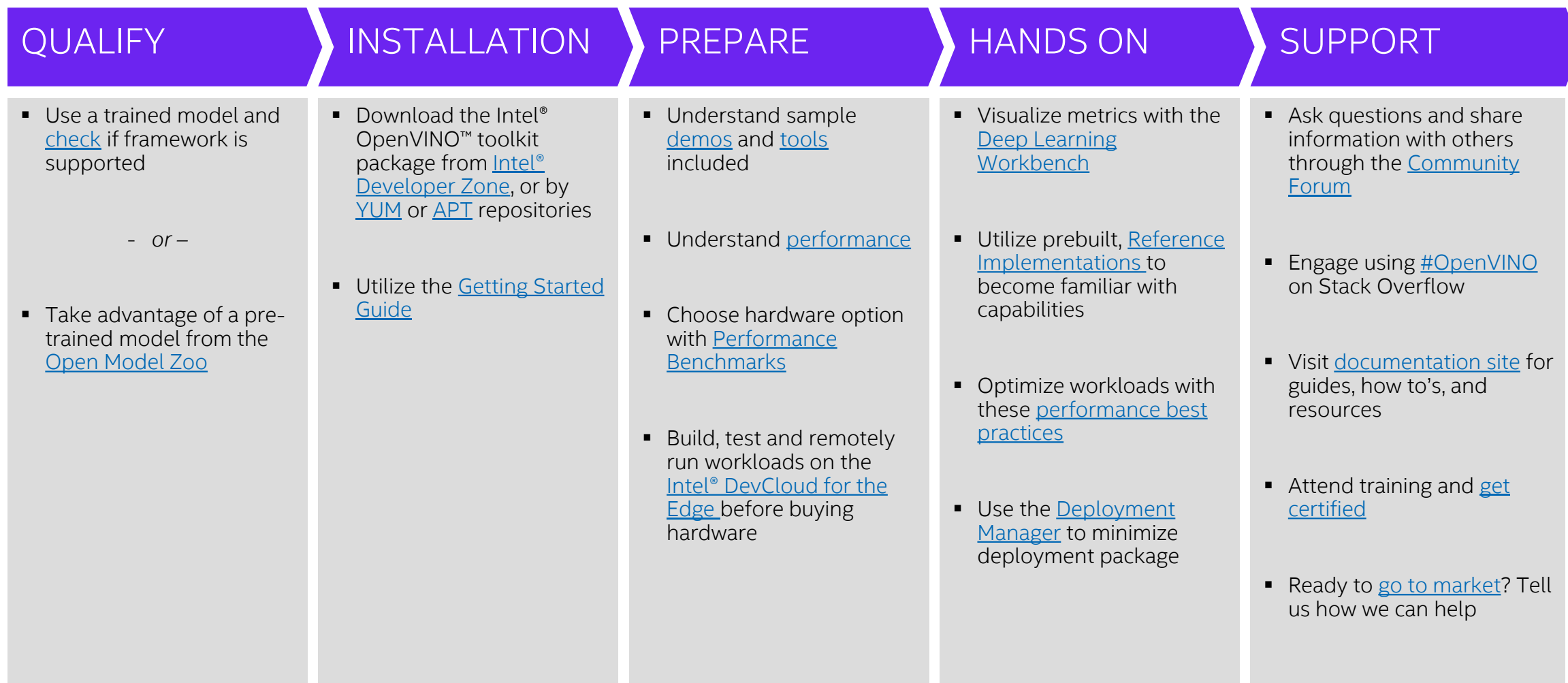
**Retail Business Services Frictionless retail checkout** [Solution Brief](#)

Using **existing Intel-based point-of-sale systems**, automatic inventory and shopper tracking with cashier-less checkout at a physical retail store was deployed at Quincy, Massachusetts.

Success Stories ► <https://intel.com/openvino-success-stories>

# Resources and Community Support

Vibrant community of developers, enterprises and skills builders



# Ready to get started?

Download directly from Intel for free

[Intel® Distribution of OpenVINO™ toolkit](#)  
(Recommended)

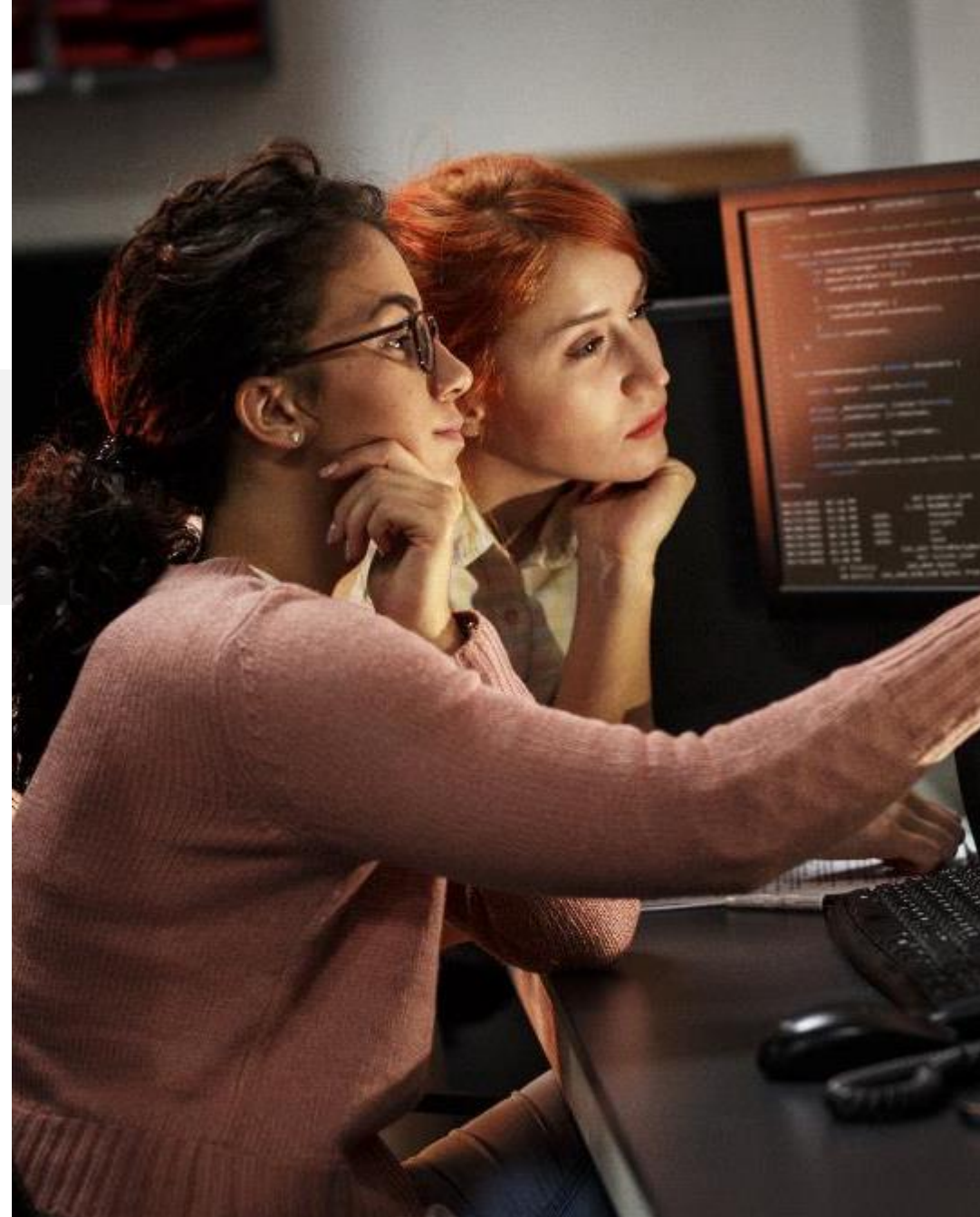
*Also available from*

Intel's Edge Software Hub | Intel® DevCloud for the Edge | PIP |  
DockerHub | Dockerfile | Anaconda Cloud | YUM | APT

*Build from source*

GitHub | Gitee (for China)

[Choose & Download](#)







# Choose between Distributions

Tool/Component	Intel® Distribution of OpenVINO™ toolkit	OpenVINO™ toolkit (open source)	Open Source Directory
Installer (including necessary drivers)	✓		
Model Optimizer	✓	✓	<a href="https://github.com/openvinotoolkit/openvino/tree/master/model-optimizer">https://github.com/openvinotoolkit/openvino/tree/master/model-optimizer</a>
Inference Engine - Core	✓	✓	<a href="https://github.com/openvinotoolkit/openvino/tree/master/inference-engine">https://github.com/openvinotoolkit/openvino/tree/master/inference-engine</a>
Intel CPU plug-in	✓ Intel® Math Kernel Library (Intel® MKL) only <sup>1</sup>	✓ BLAS, Intel® MKL <sup>1</sup> , jit (Intel MKL)	<a href="https://github.com/openvinotoolkit/openvino/tree/master/inference-engine">https://github.com/openvinotoolkit/openvino/tree/master/inference-engine</a>
Intel GPU (Intel® Processor Graphics) plug-in	✓	✓	<a href="https://github.com/openvinotoolkit/openvino/tree/master/inference-engine">https://github.com/openvinotoolkit/openvino/tree/master/inference-engine</a>
Heterogeneous plug-in	✓	✓	<a href="https://github.com/openvinotoolkit/openvino/tree/master/inference-engine">https://github.com/openvinotoolkit/openvino/tree/master/inference-engine</a>
Intel GNA plug-in	✓	✓	<a href="https://github.com/openvinotoolkit/openvino/tree/master/inference-engine">https://github.com/openvinotoolkit/openvino/tree/master/inference-engine</a>
Intel® FPGA plug-in	✓		
Intel® Neural Compute Stick (1 & 2) VPU plug-in	✓	✓	<a href="https://github.com/openvinotoolkit/openvino/tree/master/inference-engine">https://github.com/openvinotoolkit/openvino/tree/master/inference-engine</a>
Intel® Vision Accelerator based on Movidius plug-in	✓		
Multi-device & hetero plug-ins	✓	✓	
Public and Pretrained Models - incl. Open Model Zoo (IR models that run in IE + open sources models)	✓	✓	<a href="https://github.com/openvinotoolkit/open_model_zoo">https://github.com/openvinotoolkit/open_model_zoo</a>
Samples (APIs)	✓	✓	<a href="https://github.com/openvinotoolkit/openvino/tree/master/inference-engine">https://github.com/openvinotoolkit/openvino/tree/master/inference-engine</a>
Demos	✓	✓	<a href="https://github.com/openvinotoolkit/open_model_zoo">https://github.com/openvinotoolkit/open_model_zoo</a>
<b>Traditional Computer Vision</b>			
OpenCV*	✓	✓	<a href="https://github.com/opencv/opencv">https://github.com/opencv/opencv</a>
Intel® Media SDK	✓	✓ <sup>2</sup>	<a href="https://github.com/Intel-Media-SDK/MediaSDK">https://github.com/Intel-Media-SDK/MediaSDK</a>
OpenCL™ Drivers & Runtimes	✓	✓ <sup>2</sup>	<a href="https://github.com/intel/compute-runtime">https://github.com/intel/compute-runtime</a>
FPGA Runtime Environment, Deep Learning Acceleration & Bitstreams (Linux* only)	✓		



# System Requirements

Target Solution Platforms	<b>Intel® Platforms</b> <ul style="list-style-type: none"> <li><b>CPU</b> <ul style="list-style-type: none"> <li>6<sup>th</sup>-10<sup>th</sup> generation Intel® Core™ and Xeon® processors</li> <li>1<sup>st</sup> and 2<sup>nd</sup> generation Intel® Xeon® Scalable processors</li> </ul> </li> <li>Intel® Pentium® processor N4200/5, N3350/5, N3450/5 with Intel® HD Graphics</li> <li><b>Iris® Pro &amp; Intel® HD Graphics</b> <ul style="list-style-type: none"> <li>6<sup>th</sup>-10<sup>th</sup> generation Intel® Core™ processor with Intel® Iris™ Pro graphics &amp; Intel® HD Graphics</li> <li>Intel® Xeon® processor with Intel® Iris™ Pro Graphics &amp; Intel® HD Graphics (excluding E5 product family, which does not have graphics<sup>1</sup>)</li> </ul> </li> <li><b>FPGA</b> <ul style="list-style-type: none"> <li>Intel® Arria® FPGA 10 GX development kit</li> <li>Intel® Programmable Acceleration Card with Intel® Arria® 10 GX FPGA operating systems</li> <li>OpenCV* &amp; OpenVX* functions must be run against the CPU or Intel® Processor Graphics (GPU)</li> </ul> </li> <li><b>VPU:</b> Intel Movidius™ Neural Compute Stick; Intel® Neural Compute Stick2</li> <li><b>Intel® Vision Accelerator Design Products</b> <ul style="list-style-type: none"> <li>Intel® Vision Accelerator Design with Intel® Arria10 FPGA</li> <li>Intel® Vision Accelerator Design with Intel® Movidius™ VPUs</li> </ul> </li> </ul>	<b>Compatible Operating Systems</b> <ul style="list-style-type: none"> <li>Ubuntu* 18.04.3 LTS (64 bit)</li> <li>Microsoft Windows* 10 (64 bit)</li> <li>CentOS* 7.4 (64 bit)</li> <li>macOS* 10.13 &amp; 10.14 (64 bit)</li> <li>Yocto Project* Poky Jethro v2.0.3 (64 bit)</li> <li>Ubuntu 18.04.3 LTS (64 bit)</li> <li>Windows 10 (64 bit)</li> <li>CentOS 7.4 (64 bit)</li> </ul>
	<b>Development Platforms</b> <ul style="list-style-type: none"> <li>6<sup>th</sup>-10<sup>th</sup> generation Intel® Core™ and Intel® Xeon® processors</li> <li>1<sup>st</sup> and 2<sup>nd</sup> generation Intel® Xeon® Scalable processors</li> </ul>	<ul style="list-style-type: none"> <li>Ubuntu 18.04.2 LTS (64 bit)</li> <li>CentOS 7.4 (64 bit)</li> <li>Ubuntu 18.04.3 LTS (64 bit)      CentOS 7.4 (64 bit)</li> <li>Windows 10 (64 bit)      macOS* (64 bit)      Raspbian (target only)</li> <li>Ubuntu 18.04.2 LTS (64 bit)</li> <li>Ubuntu 8.04.3 LTS (64 bit)</li> <li>Windows 10 (64 bit)</li> <li>Ubuntu* 18.04.3 LTS (64 bit)</li> <li>Windows® 10 (64 bit)</li> <li>CentOS* 7.4 (64 bit)</li> <li>macOS* 10.13 &amp; 10.14 (64 bit)</li> </ul>
	<b>Additional Software Requirements</b> <p>Linux* build environment required components</p> <ul style="list-style-type: none"> <li><a href="#">OpenCV 3.4</a> or higher</li> <li><a href="#">CMake* 2.8</a> or higher</li> <li><a href="#">GNU Compiler Collection (GCC) 3.4</a> or higher</li> <li><a href="#">Python* 3.4</a> or higher</li> </ul> <p>Microsoft Windows* build environment required components</p> <ul style="list-style-type: none"> <li><a href="#">Intel® HD Graphics Driver</a> (latest version)<sup>†</sup></li> <li><a href="#">Intel® C++ Compiler 2017 Update 4</a></li> <li><a href="#">Python 3.4</a> or higher</li> <li><a href="#">OpenCV 3.4</a> or higher</li> <li><a href="#">CMake 2.8</a> or higher</li> <li><a href="#">Microsoft Visual Studio* 2015</a></li> </ul>	
External Dependencies/Additional Software		View Product Site, detailed System Requirements

# Commonly Asked Questions

**Can I use the Intel® Distribution of OpenVINO™ toolkit for commercial usage?** Yes, the Intel® Distribution of OpenVINO™ toolkit is licensed under [Intel's End User License Agreements](#) and the open-sourced OpenVINO™ toolkit is licensed under [Apache License 2.0](#). For information, review the licensing directory inside the package.

**Is the Intel® Distribution of OpenVINO™ toolkit subject to export control?** Yes, the ECCN is EAR99.

**How often does the software get updated?** Standard releases are updated 3-4 times a year, while LTS releases are updated once a year.

**What is the difference between Standard and LTS releases?** Standard Releases are recommended for new users and users currently prototyping. It offers new features, tools and support to stay current with deep learning advancements. LTS Releases are recommended for experienced users that are ready to take their application into production and who do not require new features and capabilities for their application.

**For technical questions,** visit the [Model Optimizer FAQ](#) and [Performance Benchmarks FAQ](#). If you don't find an answer, please visit the following community and support links.

## Get Help

- [Ask on the Community Forum](#)
- [Contact Intel Support](#)
- [File an Issue on GitHub\\*](#)
- [Get Answers on StackOverflow\\*](#)

## Get Involved

- [Contribute to the Code Base](#)
- [Contribute to Documentation](#)

## Stay Informed

- [Join the Mailing List](#)
- [Read the Documentation](#)
- [Read the Knowledge Base](#)
- [Read the Blog](#)

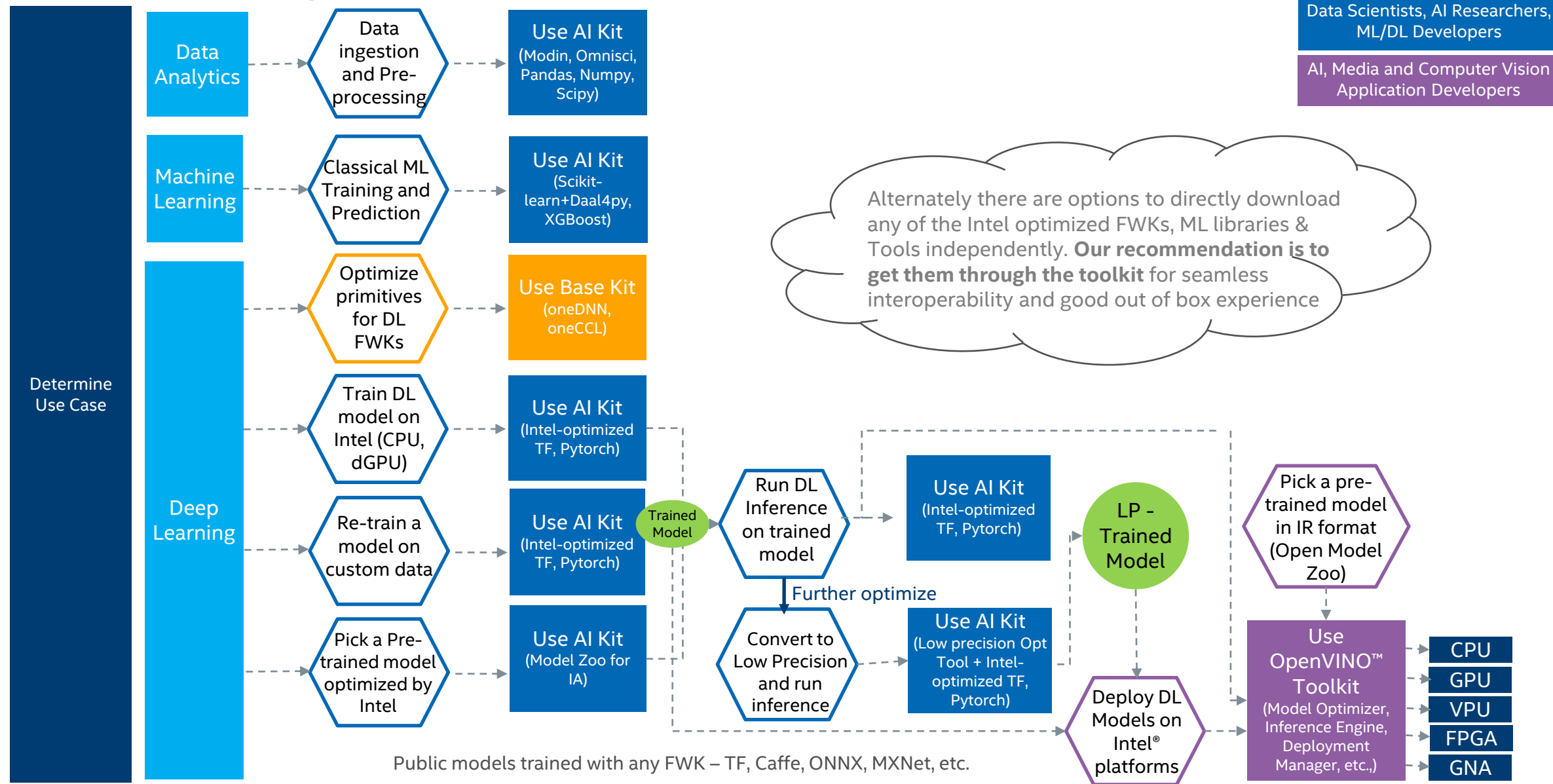
# Which Toolkit should I use

# Which Toolkit to Use When ?

	Intel® AI Analytics Toolkit	OpenVINO™ Toolkit
Key Value Prop	<ul style="list-style-type: none"> <li>• Provide performance and easy integration across end-to-end data science pipeline for efficient AI model development</li> <li>• Maximum compatibility with opensource FWKs and Libs with drop-in acceleration that require minimal to no code changes</li> <li>• Audience: Data Scientists; AI Researchers; DL/ML Developers</li> </ul>	<ul style="list-style-type: none"> <li>• Provide leading performance and efficiency for DL inference solutions to deploy across any Intel HW (cloud to edge).</li> <li>• Optimized package size for deployment based on memory requirements</li> <li>• Audience: AI Application Developers; Media and Vision Developers</li> </ul>
Use Cases	<ul style="list-style-type: none"> <li>• Data Ingestion, Data pre-processing, ETL operations</li> <li>• Model training and inference</li> <li>• Scaling to multi-core / multi-nodes / clusters</li> </ul>	<ul style="list-style-type: none"> <li>• Inference apps for vision, Speech, Text, NLP</li> <li>• Media streaming / encode, decode</li> <li>• Scale across HW architectures – edge, cloud, datacenter, device</li> </ul>
HW Support	<ul style="list-style-type: none"> <li>• CPUs - Datacenter and Server segments – Xeons, Workstations</li> <li>• GPU - ATS and PVC (in future)</li> </ul>	<ul style="list-style-type: none"> <li>• CPU - Xeons, Client CPUs and Atom processors</li> <li>• GPU - Gen Graphics; DG1 (current), ATS, PVC (in future)</li> <li>• VPU - NCS &amp; Vision Accelerator Design Products,</li> <li>• FPGA</li> <li>• GNA</li> </ul>
	<b>Use Intel® Low Precision Optimization Tool when using AI Analytics Toolkit</b>	<b>Use Post Training Optimization Tool when using OpenVINO</b>
Low Precision Support	<ul style="list-style-type: none"> <li>• Supports BF16 for training and FP16, Int8 and BF16 for Inference</li> <li>• Seamlessly integrates with Intel optimized frameworks</li> <li>• Available in the AI toolkit and independently</li> </ul>	<ul style="list-style-type: none"> <li>• Supports FP16, Int8 and BF16 for inference</li> <li>• Directly works with Intermediate Representation Format</li> <li>• Available in the Intel Distribution of OpenVINO toolkit</li> <li>• Provides Training extension via NNCF for PyTorch with FP16, Int8</li> </ul>

**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.

# AI Development Workflow



Native Code developers,  
Framework Developers

Data Scientists, AI Researchers,  
ML/DL Developers

AI, Media and Computer Vision  
Application Developers

# AI Model Deployment Workflow

Data Scientists, AI Researchers,  
ML/DL Developers

AI, Media and Computer Vision  
Application Developers

## Use AI Kit

### 0. PLAN

Determine model type & framework needed for your usage

### 1. BUILD

Find a model

Build a model (ie Train for accuracy)

### 1A. TRAIN

Bring a DL model

Is the model accurate enough?

Train/  
Re-train model

## Use OpenVINO™ Toolkit

### 3. OPTIMIZE

Run Model Optimizer

Did the model convert?

3A. MODIFY THE MODEL  
POST-TRAINING

Use Custom Layers

### 4. TEST

Test the model with IE & Benchmark

Do you prefer GUI

Use DL Workbench

### 5. PACKAGE

Use Deployment Manager

### 6. INTEGRATE

Integrate Model to Application

### 7. DEPLOY

Get Intel HW for deployment

Deploy App & Model

Fast enough with acceptable accuracy hit?

4A. ADVANCE MODEL TUNING

Use POT to change attributes of IR

Use multi-device plugin

Use Dev Cloud to find the right HW

Did you quantize the model?

Did you use multi-device?

Can you use different HW?

MO	Model Optimizer
OMZ	Open Model Zoo
DL	Deep Learning
IR	Intermediate Representation
IE	Inference Engine
HW	Hardware
POT	Post Training Optimization Tool

A comprehensive workflow to optimize your DL model for the Intel Hardware that will be used for running inference

- 1) We run the demo on DC
  - TF demo
  - PyTorch demo
  - future: (ATS demo)
- Slide on how to access DevCloud
- 2) What's behind the DC

# Intel DevCloud: Getting started with oneAPI

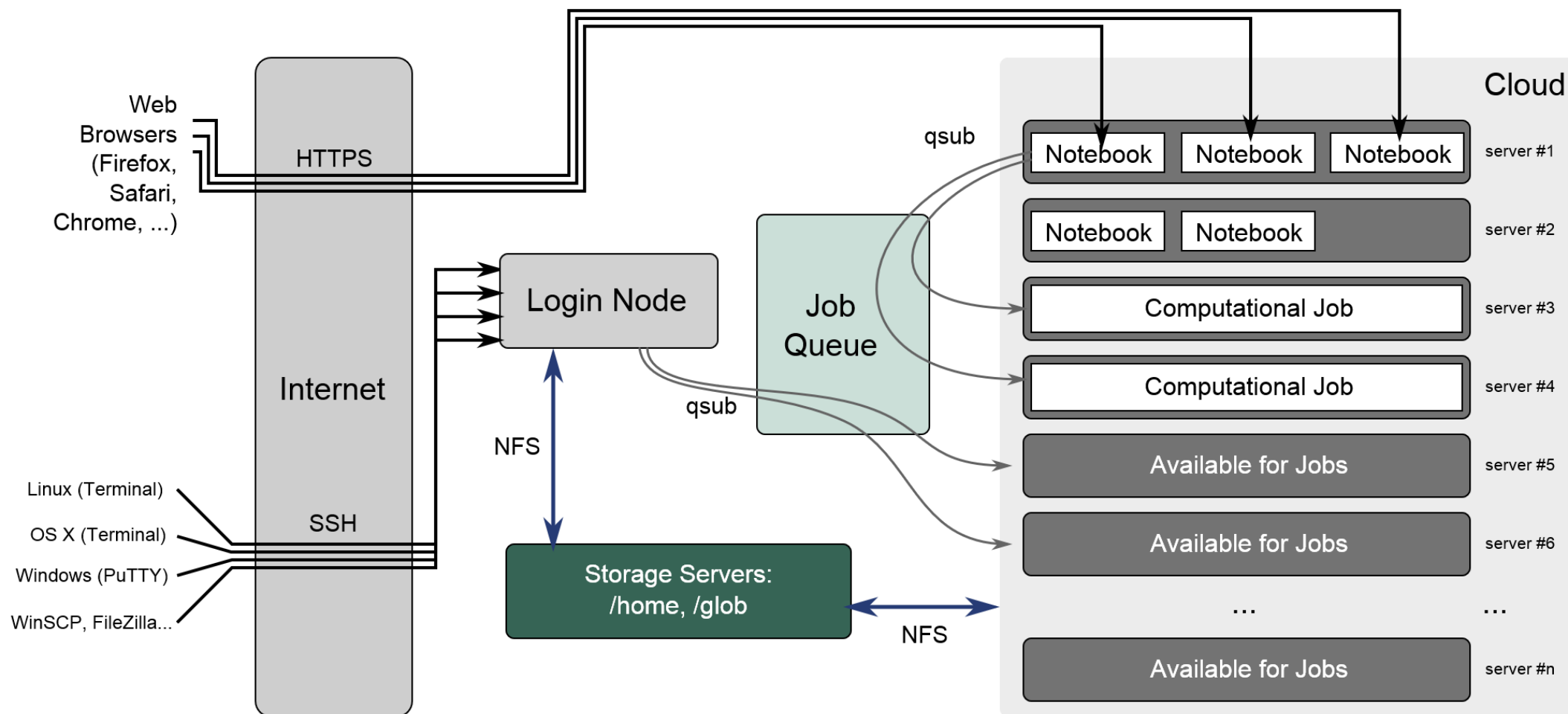




# Objectives of the External DevCloud Strategy

1. Demonstrate the promise of oneAPI.
2. Provide developers easy access to oneAPI h/w & s/w environment
3. Get high value feedback on oneAPI tools, libraries, language.
4. Seed research, papers, curriculum, lighthouse apps (the metrics output).
5. Support two tracks with web front end for consistent experience:
  - oneAPI production hardware/software
  - NDA SDP hardware/software

# Network Arch



New Additions for One API  
Under  
NDA





## Notices and Disclaimers

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors.

Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit [www.intel.com/benchmarks](http://www.intel.com/benchmarks).

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 technologies may require enabled hardware, software or service activation.

© 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.

### Optimization Notice

<sup>1</sup> Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice. Notice revision #20110804.

<sup>2</sup> Software and workloads used in performance tests may have been optimized for performance only on microprocessors from Intel. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations, and functions. Any change to any of those factors may cause the results to vary. Consult other information and performance tests while evaluating potential purchases, including performance when combined with other products. For more information, see Performance Benchmark Test Disclosure. Source: Intel measurements, as of June 2017.

## Notices and Disclaimers

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. Check with your system manufacturer or retailer or learn more at [www.intel.com](http://www.intel.com).

Intel, the Intel logo, Xeon™, Arria™ and Movidius™ are trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries.

\*Other names and brands may be claimed as the property of others.

© Intel Corporation.

Slide Reference	1	2	3
<b>System Board</b>	Intel® Server S2600 (Dual socket)	Supermicro / X11SPL-F	Supermicro / X11SPL-F
<b>Product</b>	Xeon Silver 4216	Intel(R) Xeon(R) Silver 4112	Intel(R) Xeon(R) Silver 4112
<b>CPU sockets</b>	2	-	1
<b>Physical cores</b>	2 x 16	4	4
<b>Processor Base Frequency</b>	2.10 GHz	2.60GHz	2.60GHz
<b>HyperThreading</b>	enabled	-	enabled
<b>Turbo</b>	On	-	On
<b>Power-Performance Mode</b>	Performance Mode	-	-
<b>Total System Memory size</b>	12 x 64GB	16384	16384
<b>Memory speed</b>	2400MHz	2400MHz	2400MHz
<b>Software OS</b>	Ubuntu 18.04	Ubuntu 16.04.3 LTS	Ubuntu 16.04.6 LTS
<b>Software Kernel</b>	4.15.0-66-generic x86_64	4.13.0-36-generic	4.15.0-29-generic
<b>Test Date</b>	27 September 2019	25 May 2018	18 April 2019
<b>Precision (IntMode)</b>	Int 8 (Throughput Mode)	FP32	Int 8 (Throughput Mode)
<b>Power (TDP)</b>	200W	85W	85W
<b>Price Link on 30 Sep 2019 (Prices may vary)</b>	\$2,024	\$483	\$483
<b>Network</b>	Mobilenet SSD	Mobilenet SSD	Mobilenet SSD

## Notices and Disclaimers

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. Check with your system manufacturer or retailer or learn more at [www.intel.com](http://www.intel.com).

Intel, the Intel logo, Xeon™, Arria™ and Movidius™ are trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries.

\*Other names and brands may be claimed as the property of others.

© Intel Corporation.

System Board	Intel prototype, TGL U DDR4 SODIMM RVP	ASUSTeK COMPUTER INC. / PRIME Z370-A
CPU	11 <sup>th</sup> Gen Intel® Core™ -5-1145G7E @ 2.6 GHz.	8 <sup>th</sup> Gen Intel® Core™ i5-8500T @ 3.0 GHz.
Sockets / Physical cores	1 / 4	1 / 6
HyperThreading / Turbo Setting	Enabled / On	Na / On
Memory	2 x 8198 MB 3200 MT/s DDR4	2 x 16384 MB 2667 MT/s DDR4
OS	Ubuntu* 18.04 LTS	Ubuntu* 18.04 LTS
Kernel	5.8.0-050800-generic	5.3.0-24-generic
Software	Intel® Distribution of OpenVINO™ toolkit 2021.1.075	Intel® Distribution of OpenVINO™ toolkit 2021.1.075
BIOS	Intel TGLIFUI1.R00.3243.A04.2006302148	AMI, version 2401
BIOS release date	Release Date: 06/30/2021	7/12/2019
BIOS Setting	Load default settings	Load default settings, set XMP to 2667
Test Date	9/9/2021	9/9/2021
Precision and Batch Size	CPU: INT8, GPU: FP16-INT8, batch size: 1	CPU: INT8, GPU: FP16-INT8, batch size: 1
Number of Inference Requests	4	6
Number of Execution Streams	4	6
Power (TDP Link)	<u>28 W</u>	<u>35W</u>
Price (USD) Link on Sep 22,2021 Prices may vary	<u>\$309</u>	<u>\$192</u>

- 1): Memory is installed such that all primary memory slots are populated.
- 2): Testing by Intel as of September 9, 2021