

# ACCELERATE DEEP LEARNING INFERENCE USING INTEL TECHNOLOGIES

## INTRODUCTION: SMART VIDEO

April 2020

# SMART VIDEO WORKSHOP OVERVIEW

## INTRODUCTION

1. Introduction to Intel technologies for deep learning inference
2. Hardware acceleration techniques

Each module contains a hands-on lab exercise that introduces various Intel technologies to accelerate computer vision application with hardware heterogeneity.

INTEL® DISTRIBUTION OF  
OPENVINO™ 101

HARDWARE ACCELERATION ON LAPTOP  
AND DEVCLOUD

OPTIMIZATION

APPLICATION

CUSTOM LAYERS

2. Basic End-to-End Object Detection Example

3./4./5. Hardware Acceleration with CPU, Integrated GPU, Intel® Movidius™ NCS, FPGA

6. Optimization Tools and Techniques

7. Advanced Video Analytics

8. Custom layers



# OPTIMIZATION NOTICE

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 or 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. Refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice. Notice Revision #20110804.



# LEGAL NOTICES AND DISCLAIMERS (1 OF 2)

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software, or service activation. Performance varies depending on system configuration. No computer system can be absolutely secure. Check with your system manufacturer or retailer or learn more at [www.intel.com](http://www.intel.com).

Performance estimates were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown." Implementation of these updates may make these results inapplicable to your device or system.

Cost reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings. Circumstances will vary. Intel does not guarantee any costs or cost reduction.

This document contains information on products, services, and processes in development. All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest forecast, schedule, specifications and roadmaps.

Any forecasts of goods and services needed for Intel's operations are provided for discussion purposes only. Intel will have no liability to make any purchase in connection with forecasts published in this document.

Arduino\* 101 and the Arduino infinity logo are trademarks or registered trademarks of Arduino, LLC.

Altera, Arria, the Arria logo, Intel, the Intel logo, Intel Atom, Intel Core, Intel Nervana, Intel Xeon Phi, Movidius, Saffron, and Xeon are trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

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

Copyright © 2018 Intel Corporation. All rights reserved.



# LEGAL NOTICES AND DISCLAIMERS (2 OF 2)

This document contains information on products, services, and/or processes in development. All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest forecast, schedule, specifications, and roadmaps. Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software, or service activation. Learn more at [intel.com](https://www.intel.com), or from the OEM or retailer.

No computer system can be absolutely secure.

Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. Consult other sources of information to evaluate performance as you consider your purchase. For more complete information about performance and benchmark results, visit [www.intel.com/performance](https://www.intel.com/performance).

Cost-reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings. Circumstances will vary. Intel does not guarantee any costs or cost reduction.

Statements in this document that refer to Intel's plans and expectations for the quarter, the year, and the future are forward-looking statements that involve a number of risks and uncertainties.

A detailed discussion of the factors that could affect Intel's results and plans is included in Intel's SEC filings, including the annual report on Form 10-K.

The products described may contain design defects or errors, known as *errata*, which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Performance estimates were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown." Implementation of these updates may make these results inapplicable to your device or system.

No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.

Intel does not control or audit third-party benchmark data or the web sites referenced in this document. You should visit the referenced web site and confirm whether referenced data are accurate.

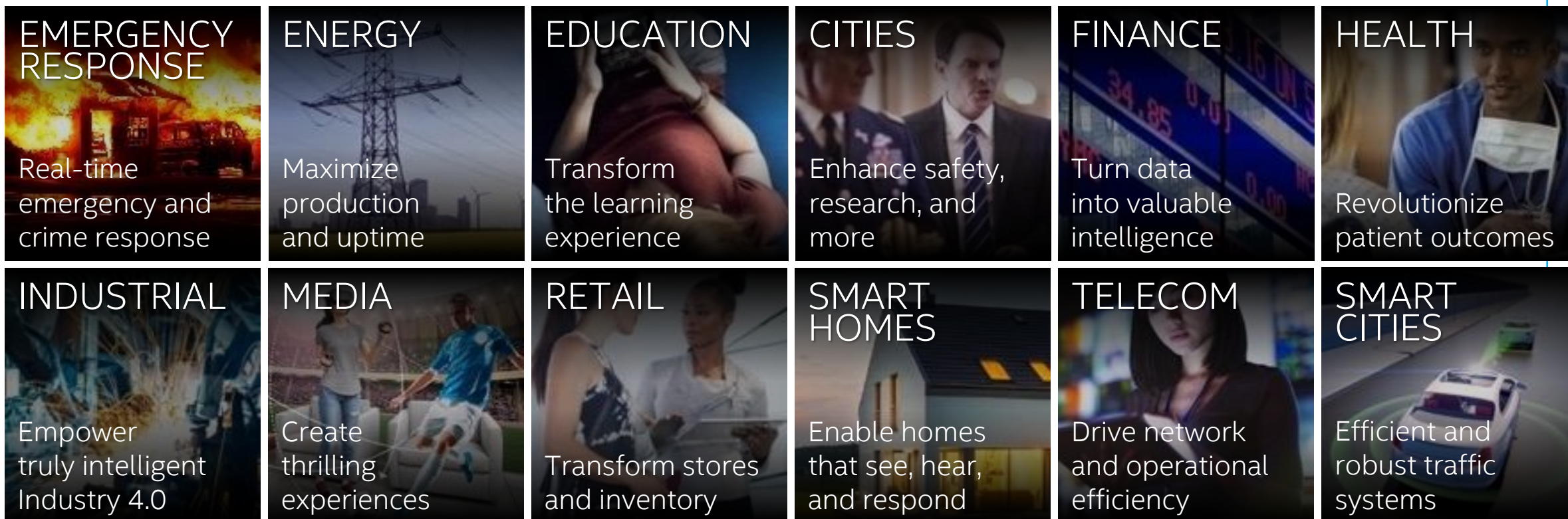
Intel, the Intel logo, Pentium, Celeron, Atom, Core, Xeon, Movidius, Saffron, and others are trademarks of Intel Corporation in the United States and other countries.

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

Copyright © 2018, Intel Corporation. All rights reserved.



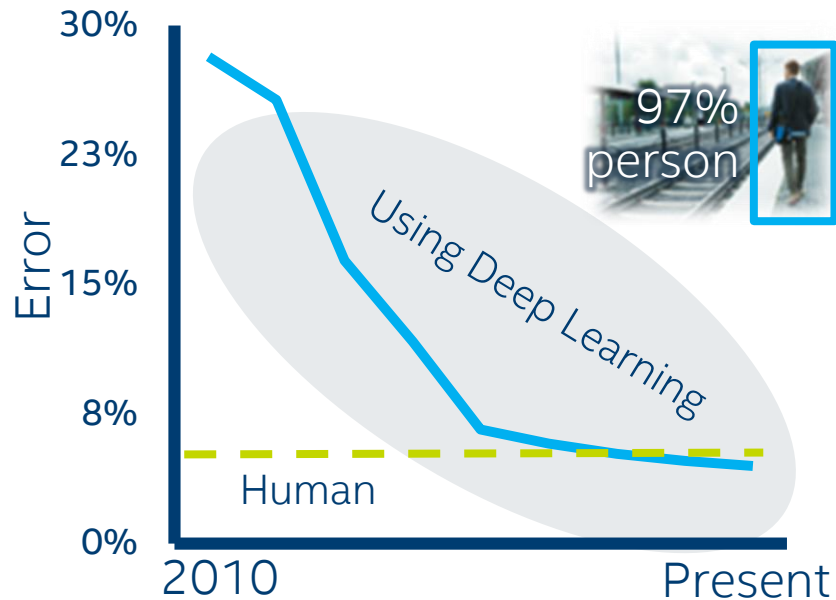
# AI IS CHANGING EVERY MARKET



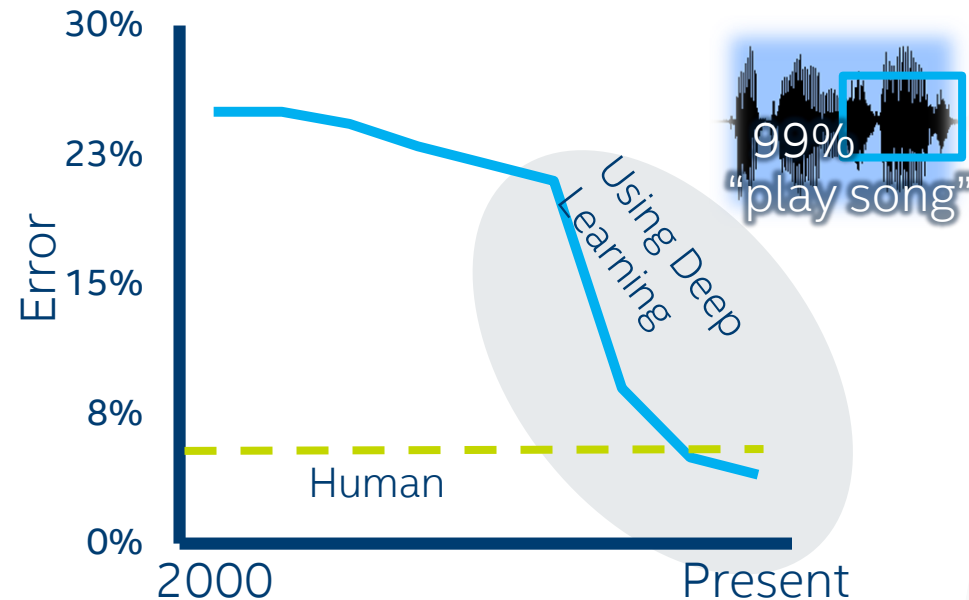
# DEEP LEARNING BREAKTHROUGHS AND OPPORTUNITIES

Machines able to meet or exceed human image and speech recognition

## Image Recognition



## Speech Recognition

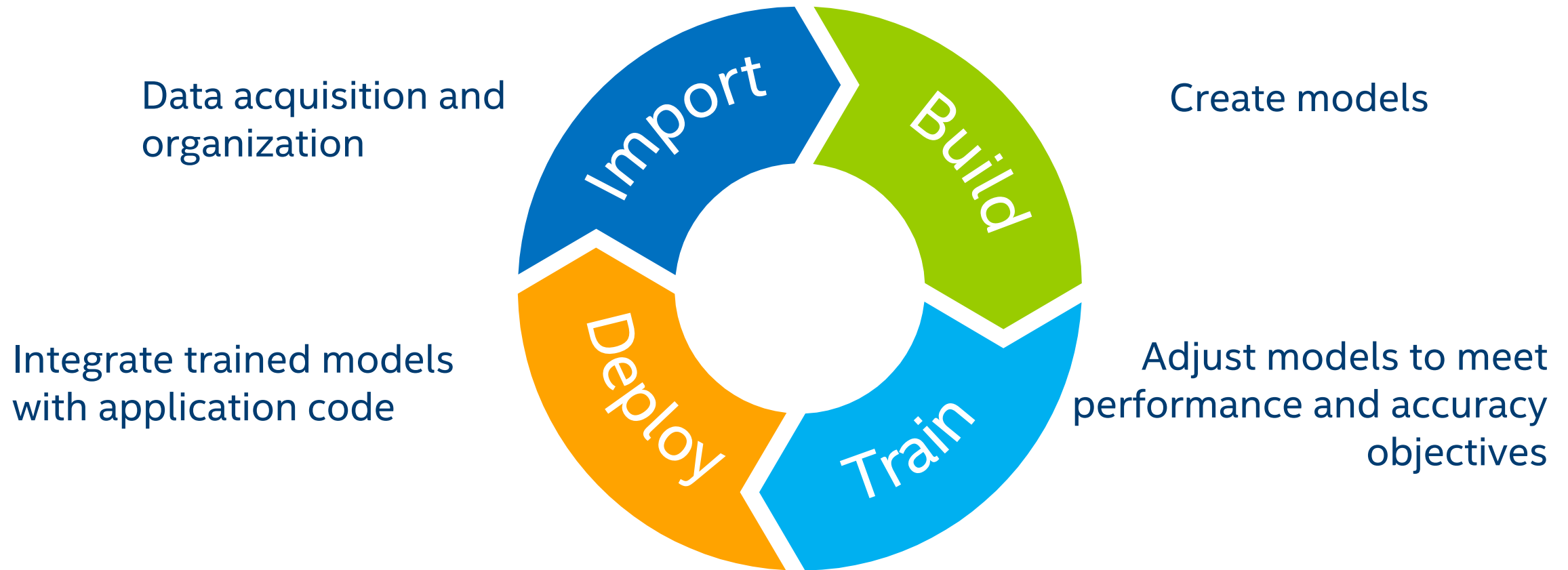


**ADDITIONAL ECONOMIC  
IMPACT DRIVEN BY AI  
\$13 TRILLION IN 2030**



Source: ILSVRC ImageNet winning entry classification error rate each year 2010-2016 (Left), <https://www.microsoft.com/en-us/research/blog/microsoft-researchers-achieve-new-conversational-speech-recognition-milestone/> (Right)  
Source: <https://www.mckinsey.com/featured-insights/artificial-intelligence/notes-from-the-ai-frontier-applications-and-value-of-deep-learning>

# DEEP LEARNING DEVELOPMENT CYCLE



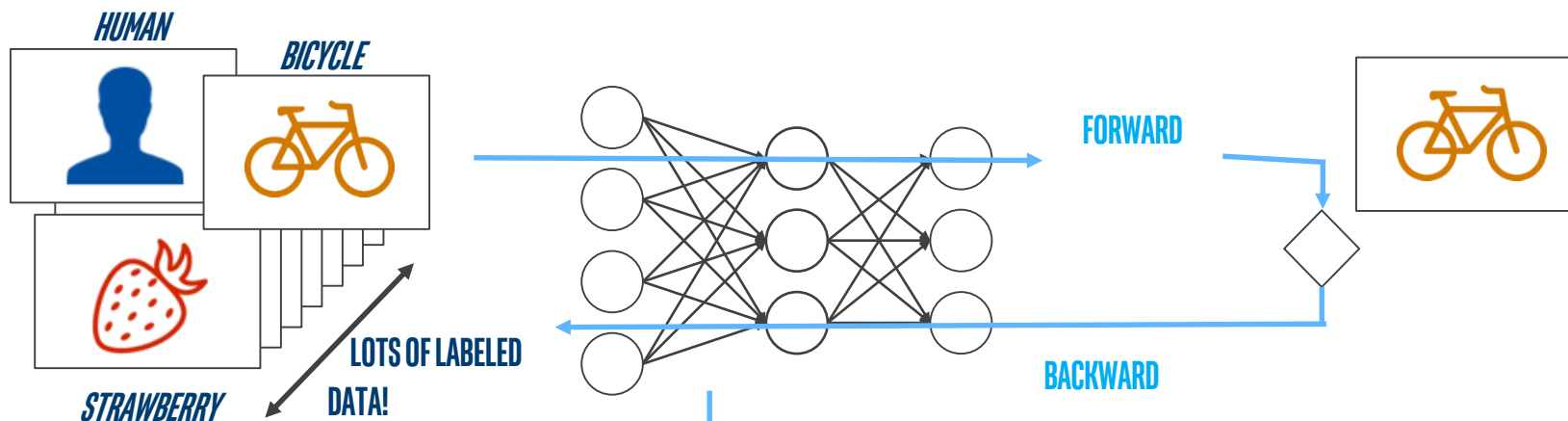
Intel® Distribution OpenVINO™ Toolkit Provides Deployment from Intel® Edge to Cloud



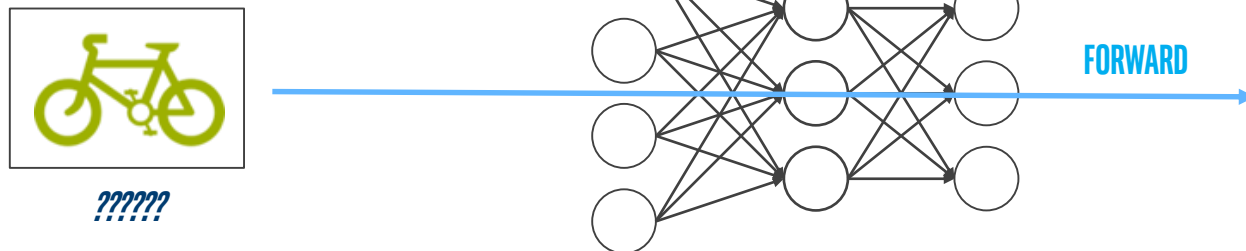


# DEEP LEARNING: TRAINING VS. INFERENCE

## TRAINING

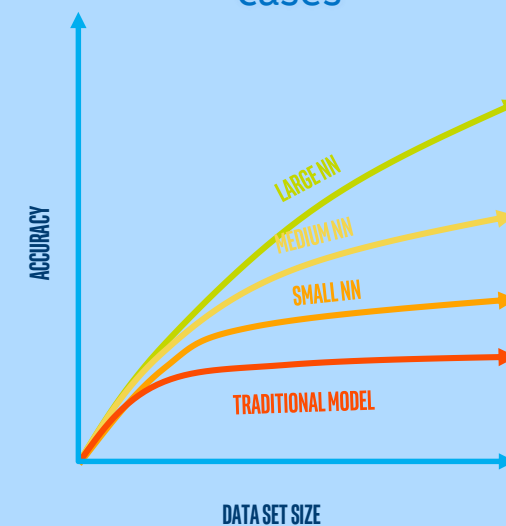


## INFERENCE



## DID YOU KNOW?

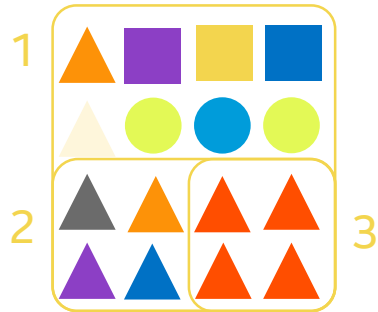
Training requires a very large data set and deep neural network (many layers) to achieve the highest accuracy in most cases



# AI COMPUTE CONSIDERATIONS

How do you determine the right computing for your AI needs?

## WORKLOADS



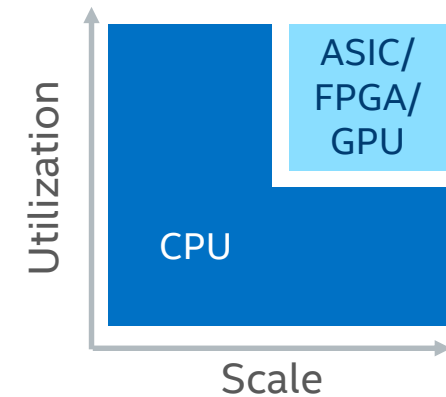
What is my workload profile?

## REQUIREMENTS



What are my use case requirements?

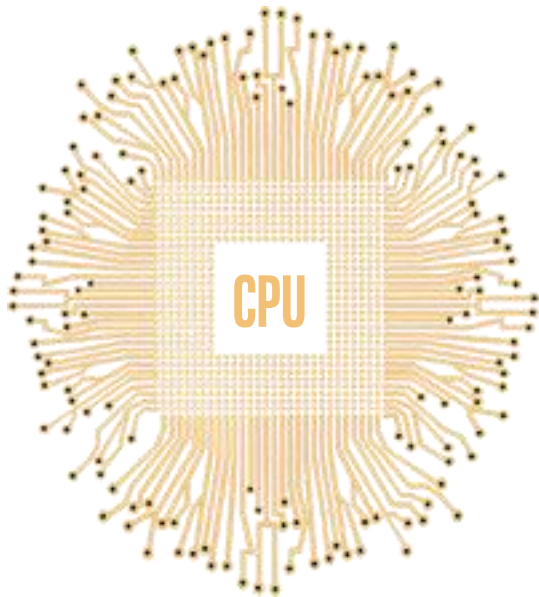
## DEMAND



How prevalent is AI in my environment?

# WHY INTEL AI COMPUTE?

## MAXIMIZE



Get the most out of the foundation for AI from the CPU leader

## OPTIMIZE



Choose the right compute for you from the one with all the options

## SIMPLIFY

OPTIMIZED SW  
DATA PIPELINE  
ANALYTICS & AI  
SUPPORT  
MOVE/STORE



Reduce “moving parts” by building on an optimized AI platform

## LEAD



Lead your industry by aligning with the builder of next-gen AI solutions

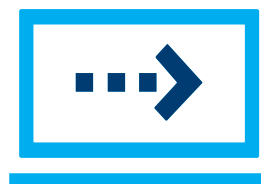
# INTEL® DISTRIBUTION OF OPENVINO™ TOOLKIT

## Tool Suite for High-Performance, Deep Learning Inference

Faster, more accurate real-world results using high-performance, AI and computer vision inference deployed into production across Intel® architecture from edge to cloud



High-Performance,  
Deep Learning Inference



Streamlined Development,  
Ease of Use



Write Once,  
Deploy Anywhere



# INSIDE INTEL® DISTRIBUTION OF OPENVINO™ TOOLKIT

## Deep Learning

### Intel® Deep Learning Deployment Toolkit

**Model Optimizer**  
Convert & Optimize



**Inference Engine**  
Optimized Inference

+ samples

IR = Intermediate Representation file

### Open Model Zoo

**Intel & Public  
Pretrained Models**

**Demos**

**Model  
Downloader**

**Accuracy  
Checker**

**Deployment  
Manager**

**Post Training Optimization Toolkit**

**Benchmark  
App**

**DL Workbench**

**DL Streamer**

New

## Traditional Computer Vision

**OpenCV\***

**Samples**

For Intel® CPU & GPU/Intel® Processor Graphics

## Tools & Libraries

### Increase Media/Video/Graphics Performance

**Intel® Media SDK**  
Open Source version

**OpenCL™  
Drivers & Runtimes**

For GPU/Intel® Processor Graphics

### Optimize Intel® FPGA (Linux\* only)

**FPGA RunTime  
Environment**  
(from Intel® FPGA SDK for OpenCL™)

**Bitstreams**

**OS Support:** CentOS\* 7.4 (64 bit), Ubuntu\* 16.04.3 LTS (64 bit), Microsoft Windows\* 10 (64 bit), Yocto Project\* version Poky Jethro v2.0.3 (64 bit), macOS\* 10.13 & 10.14 (64 bit)

Intel® Architecture-Based  
Platforms Support



Intel® Vision Accelerator  
Design Products &  
AI in Production/  
Developer Kits

An open source version is available at [01.org/openvinotoolkit](https://01.org/openvinotoolkit) (some deep learning functions support Intel CPU/GPU only).



OpenVX and the OpenVX logo are trademarks of the Khronos Group Inc.  
OpenCL and the OpenCL logo are trademarks of Apple Inc. used by permission by Khronos



# DEPLOY DEEP LEARNING SOLUTIONS WITH INTEL<sup>®</sup> DISTRIBUTION OF OPENVINO<sup>™</sup> TOOLKIT

**1. BUILD**

**2. OPTIMIZE**

**3. DEPLOY**





**1. BUILD**



**2. OPTIMIZE**



**3. DEPLOY**

# BREADTH OF SUPPORTED FRAMEWORKS MAXIMIZES DEVELOPMENT

011010110110  
110101101011  
001011010100  
011010110110  
110101101011  
001011010100  
011010110110  
110101101011  
001011010100  
011010110110  
110101101011  
001011010100  
011010110110  
110101101011  
001011010100  
011010110110  
110101101011  
001011010100  
011010110110  
110101101011  
001011010100  
011010110110



mxnet



Caffe



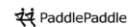
OpenVINO™

(and other tools via  
ONNX\* conversion)



LibSVM

MATLAB®



SIEMENS



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)



[Optimization Notice](#)

Copyright © 2020, Intel Corporation. All rights reserved.

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





**1. BUILD**



**2. OPTIMIZE**



**3. DEPLOY**





## 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 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)

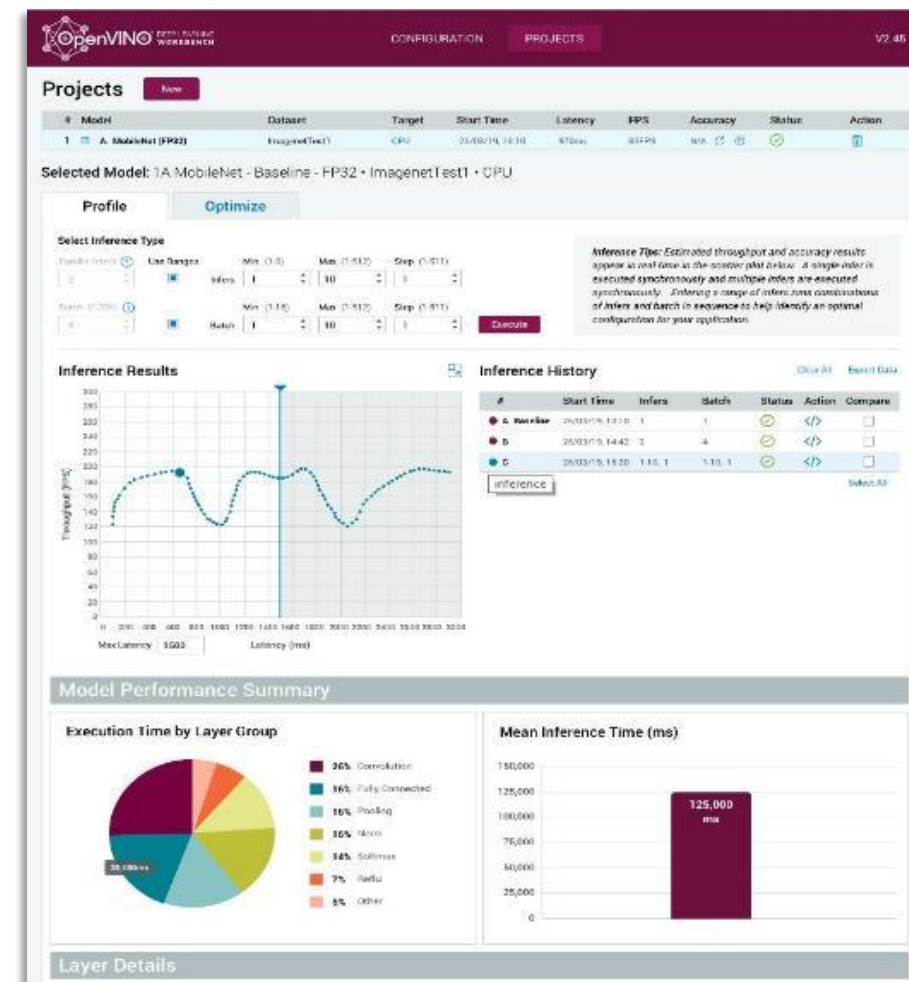
# Deep Learning Workbench



- Web-based, UI extension tool of the Intel® Distribution of OpenVINO™ toolkit
- 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
- Provide accuracy information through accuracy checker
- Direct access to models from public set of Open Model Zoo

## Development Guide ►

[https://docs.openvino toolkit.org/latest/\\_docs\\_Workbench\\_DG\\_Introduction.html](https://docs.openvino toolkit.org/latest/_docs_Workbench_DG_Introduction.html)



## Optimization Notice

Copyright © 2020, Intel Corporation. All rights reserved.

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

# TOOLS TO SPEED UP TEST CYCLES AND DEVELOPMENT



## [NEW] Post-training Optimization

- Reduce model size into low precision data types, such as INT8
- Reduces model size while also improving latency



## Deployment Manager

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



## Model Analyzer

- Provides theoretical data on models: computational complexity (flops), number of neurons, memory consumption



## Accuracy Checker

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



## Benchmark App

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



## Model Downloader

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

Get Started ► [https://docs.openvino toolkit.org/latest/\\_docs\\_IE\\_DG\\_Tools\\_Overview.html](https://docs.openvino toolkit.org/latest/_docs_IE_DG_Tools_Overview.html) –or– by using the [Deep Learning Workbench](#)



[Optimization Notice](#)

Copyright © 2020, Intel Corporation. All rights reserved.

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





**1. BUILD**



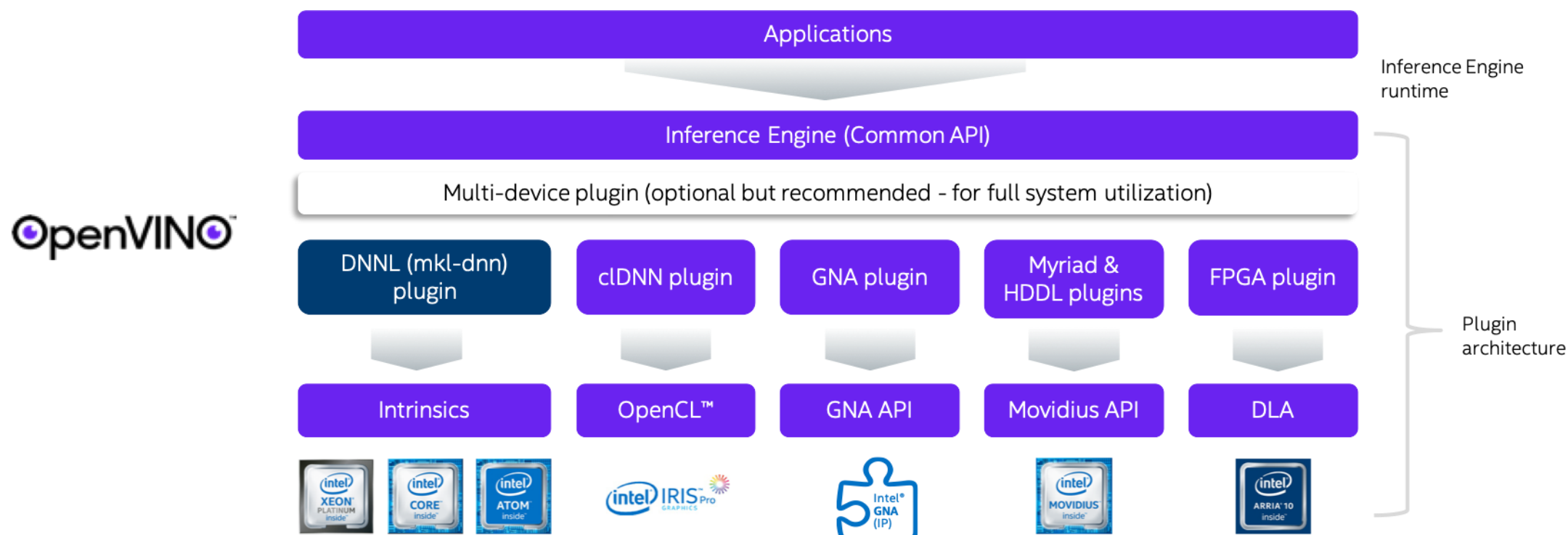
**2. OPTIMIZE**



**3. DEPLOY**

# INTEL® DISTRIBUTION OF OPENVINO™ TOOLKIT

## WRITE ONCE, DEPLOY ANYWHERE

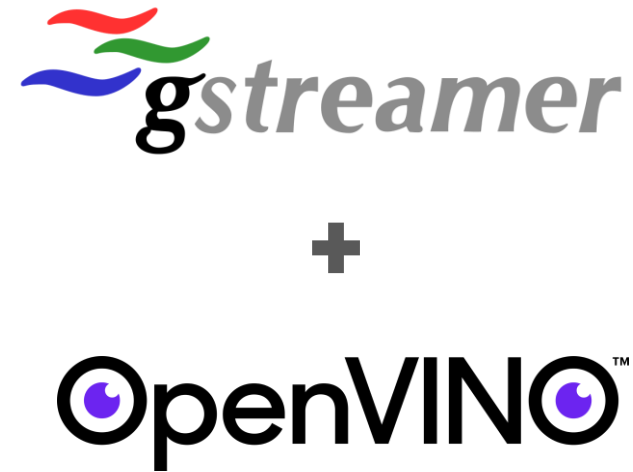


## Deep Learning Streamer

- Intel® Distribution of OpenVINO™ toolkit **Deep Learning (DL) Streamer**, now part of the default installation package
- Enables developers to **create and deploy** optimized streaming media analytics pipelines across Intel® architecture from edge to cloud
- Optimal pipeline interoperability with a familiar developer experience built using the **GStreamer\*** multimedia framework

**Learn More** ►

[https://docs.openvino toolkit.org/latest/index.html#toolkit\\_components](https://docs.openvino toolkit.org/latest/index.html#toolkit_components)



# SPEED UP DEVELOPMENT USING THE OPEN MODEL ZOO

Open source resources with pre-trained models, samples and demos



## Computer Vision

[Object detection](#)

[Object recognition](#)

[Reidentification](#)

[Semantic segmentation](#)

[Instance segmentation](#)

[Human pose estimation](#)

[Image processing](#)



## Audio, Speech, Language

[Text detection](#)

[Text recognition](#)



## Recommender

[Action recognition](#)



## Other

(Data Generation,  
Reinforcement Learning)

[Compression models](#)

[Image retrieval](#)

And more..

## PRE-TRAINED MODELS

[https://github.com/opencv/open\\_model\\_zoo](https://github.com/opencv/open_model_zoo)



[Optimization Notice](#)

Copyright © 2020, Intel Corporation. All rights reserved.

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



# SPEED UP DEVELOPMENT USING THE OPEN MODEL ZOO

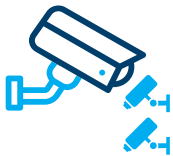
## Open source resources with pre-trained models, demos, and tools

The Open Model Zoo demo applications are console applications that demonstrate how you can use your applications to solve specific use-cases.



### Smart Classroom

Recognition and action detection demo for classroom settings



### Multi-Camera, Multi-Person

Tracking multiple people on multiple cameras for public safety use cases



### Gaze Estimation

Face detection followed by gaze estimation, head pose estimation and facial landmarks regression.



### Super Resolution

Enhances the resolution of the input image



### Action Recognition

Classifies actions that are being performed on input video

*And more..*

## DEMO APPLICATIONS

[https://github.com/opencv/open\\_model\\_zoo](https://github.com/opencv/open_model_zoo)



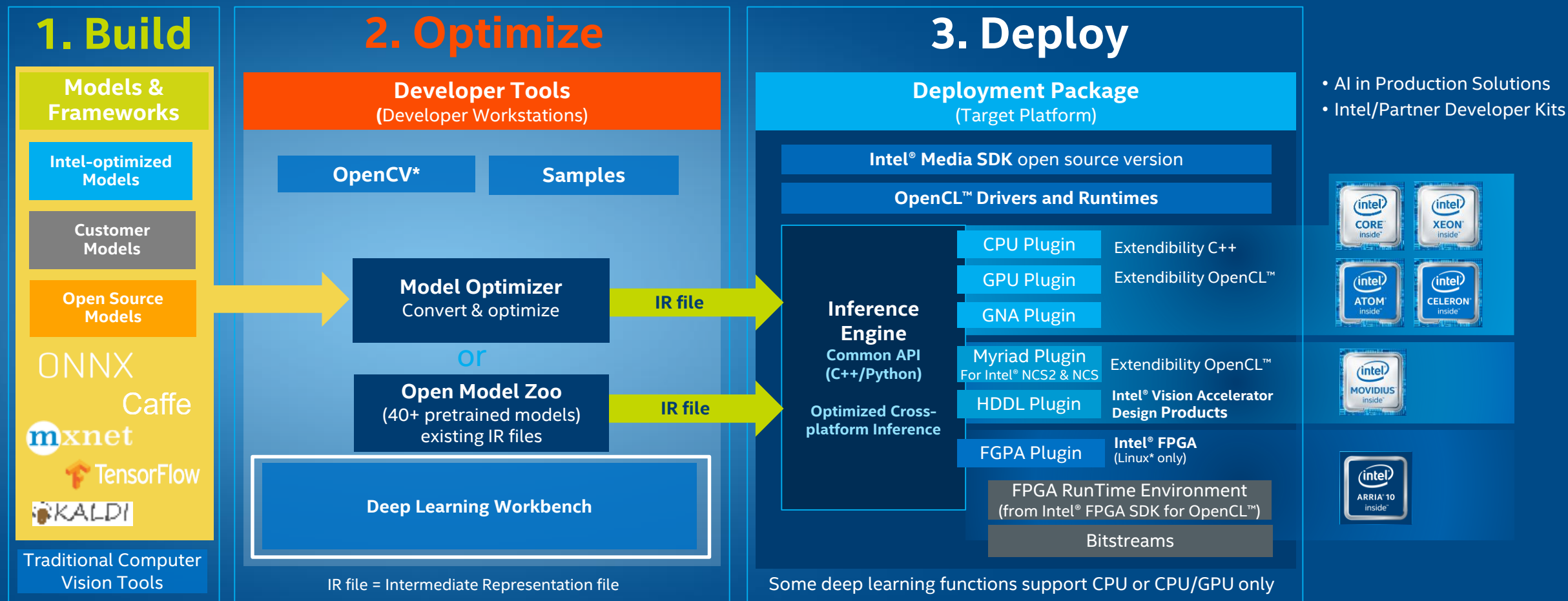
### Optimization Notice

Copyright © 2020, Intel Corporation. All rights reserved.

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

# USING THE INTEL® DISTRIBUTION OF OPENVINO™ TOOLKIT

ADVANCED CAPABILITIES TO STREAMLINE DEEP LEARNING DEPLOYMENT



Intel® NCS = Intel® Neural Compute Stick (VPU)

## Optimization Notice

Copyright © 2019, Intel Corporation. All rights reserved.

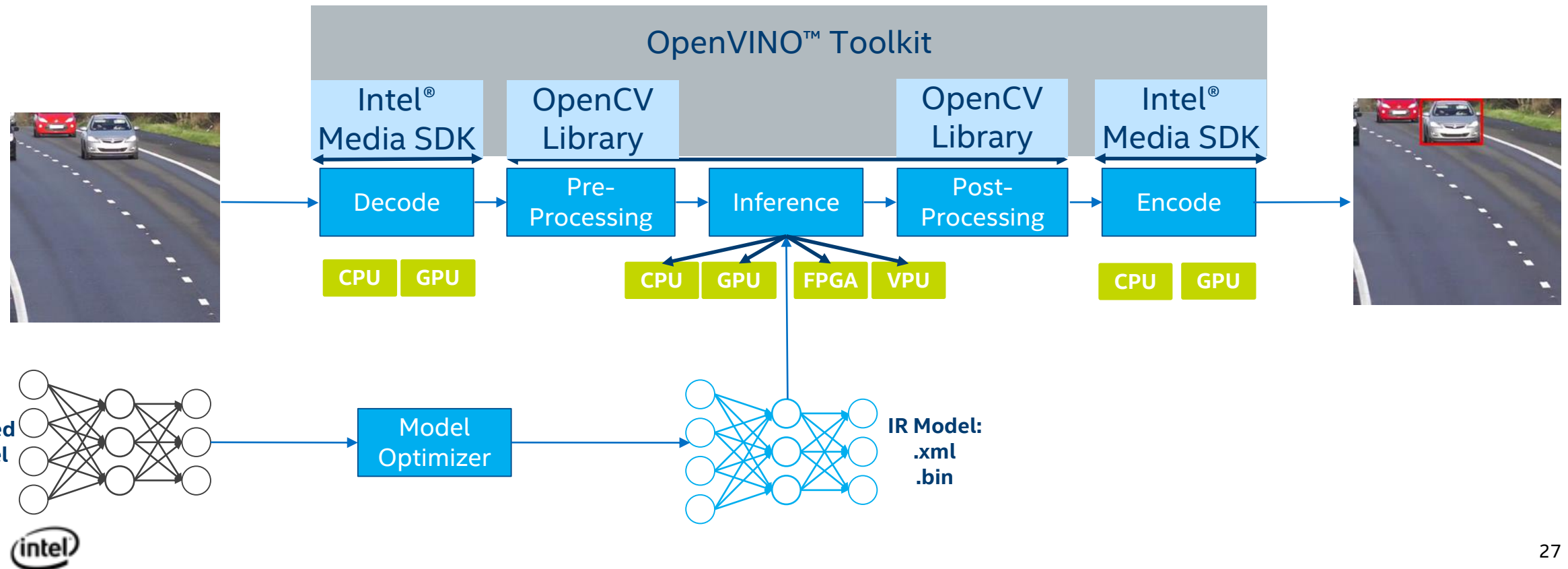
\*Other names and brands may be claimed as the property of others.

OpenCL and the OpenCL logo are trademarks of Apple Inc. used by permission by Khronos



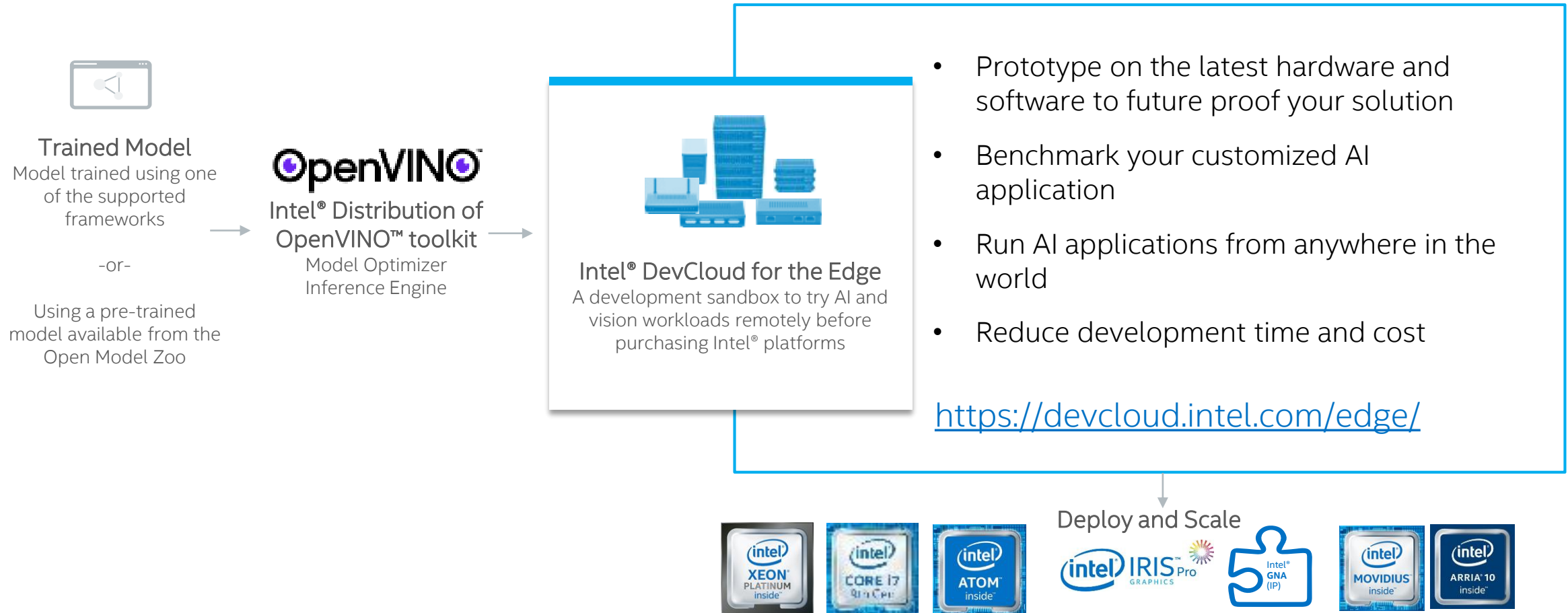
# Workflow of Applying OpenVINO™ in CV Applications, Accelerate Streaming Performance

Using Intel® Media SDK and the OpenVINO™ toolkit together enables customers to build high performance, intelligent vision solutions.



# TEST HARDWARE WITH THE INTEL® DEVCLOUD FOR THE EDGE

## Powered by Intel® Distribution of OpenVINO™ toolkit



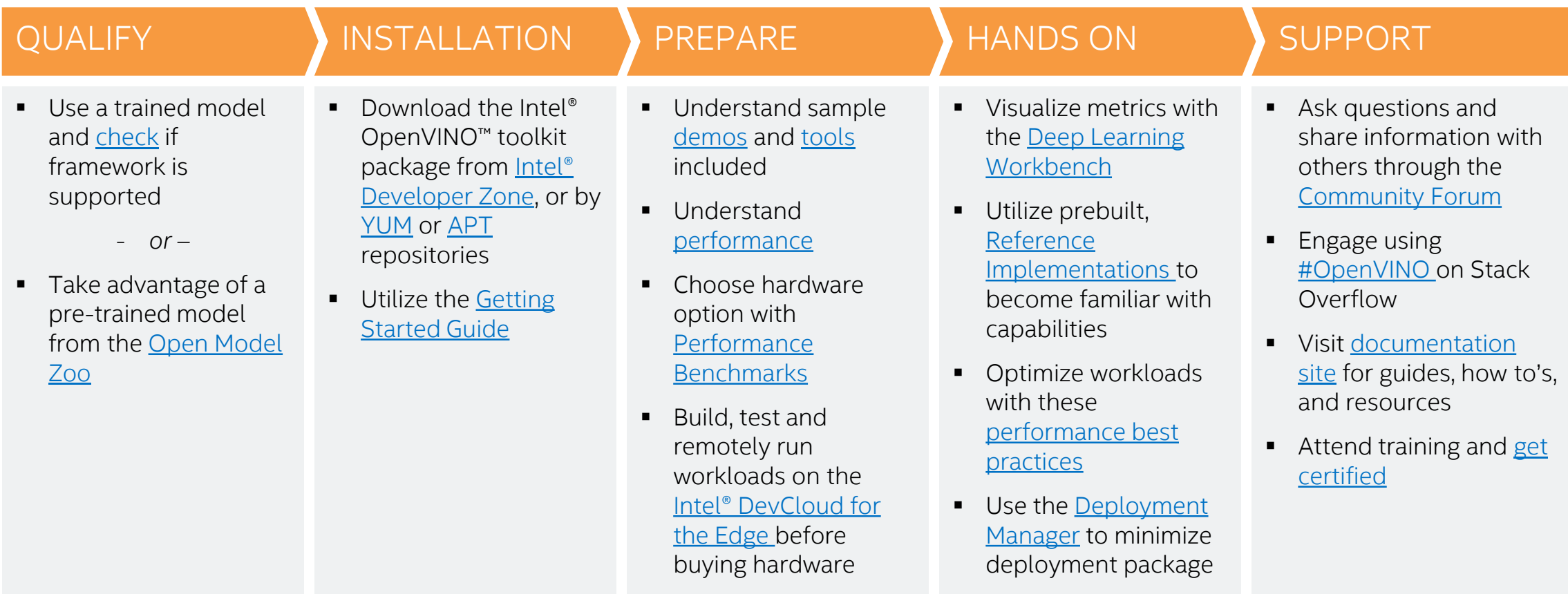
[Optimization Notice](#)

Copyright © 2020, Intel Corporation. All rights reserved.

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

# GETTING STARTED WITH INTEL® DISTRIBUTION OF OPENVINO™ TOOLKIT

Recommendations to the customer or developer



# JUMPSTART DEEP LEARNING TODAY!

Download Free ►

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

Also available from

[Docker](#) | [YUM](#) | [APT](#) | [\[NEW\] Anaconda Cloud](#)



[Optimization Notice](#)

Copyright © 2020, Intel Corporation. All rights reserved.

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

OpenVINO



For public use – OK for non-NDA disclosure

