

# ACCELERATE DEEP LEARNING INFERENCE USING INTEL TECHNOLOGIES

## INTRODUCTION: SMART VIDEO

September 2018

Core and Visual Computing Group

# 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 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. 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](http://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](http://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.



EMERGENCY RESPONSE



FINANCIAL SERVICES



MACHINE VISION



CITIES/TRANSPORTATION

# VIDEO: THE “EYE OF IOT”

USE OF VIDEO, COMPUTER VISION AND DEEP LEARNING IS GROWING RAPIDLY



AUTONOMOUS VEHICLES



RESPONSIVE RETAIL



MANUFACTURING

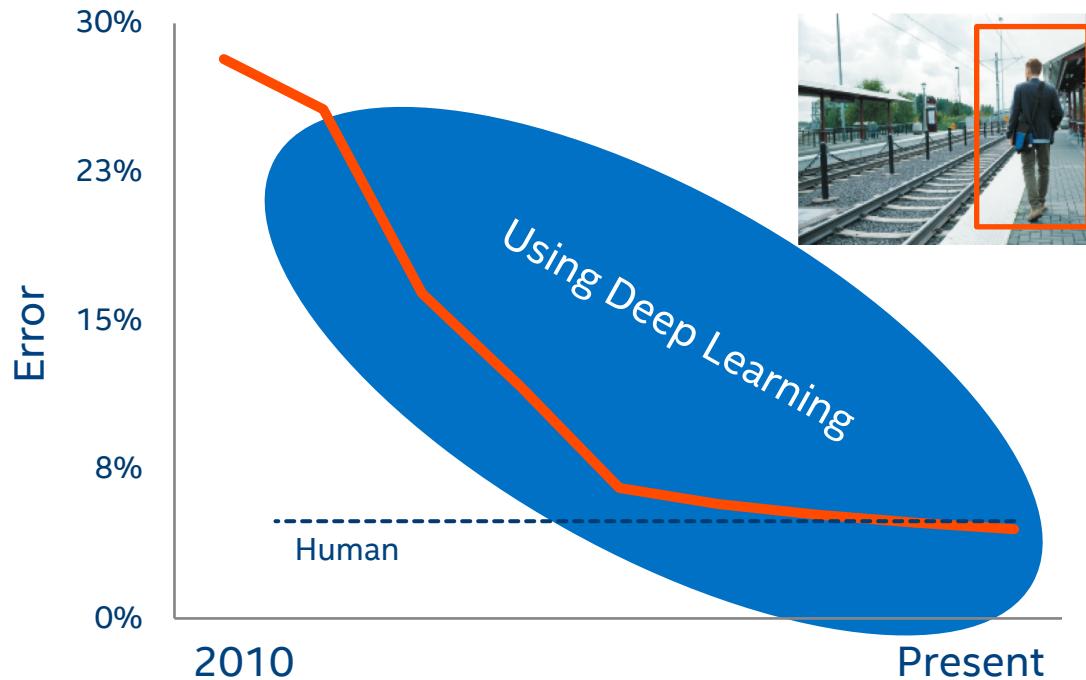


PUBLIC SECTOR

# Deep Learning Usage Is Increasing

Deep learning revenue is estimated to grow from \$655M in 2016 to **\$35B** by 2025<sup>1</sup>.

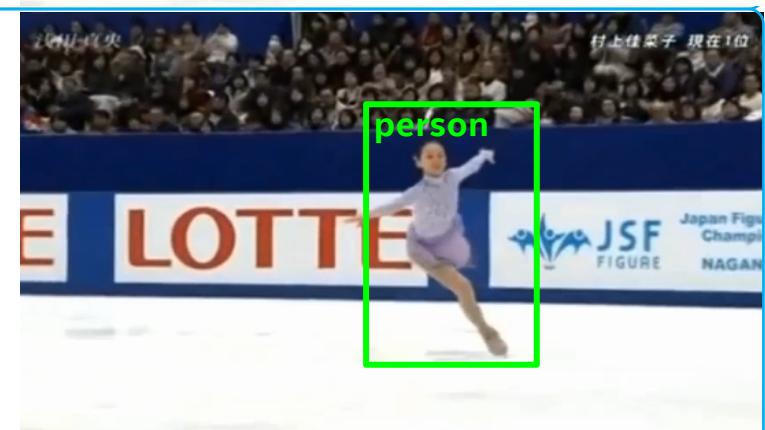
## Image Recognition



## Traditional Computer Vision Object Detection



## Deep Learning Computer Vision Person Recognition



Market Opportunities + Advanced Technologies Have Accelerated Deep Learning Adoption

<sup>1</sup>Tractica\* 2Q 2017

Activities

Terminal

Mon 15:19 ●

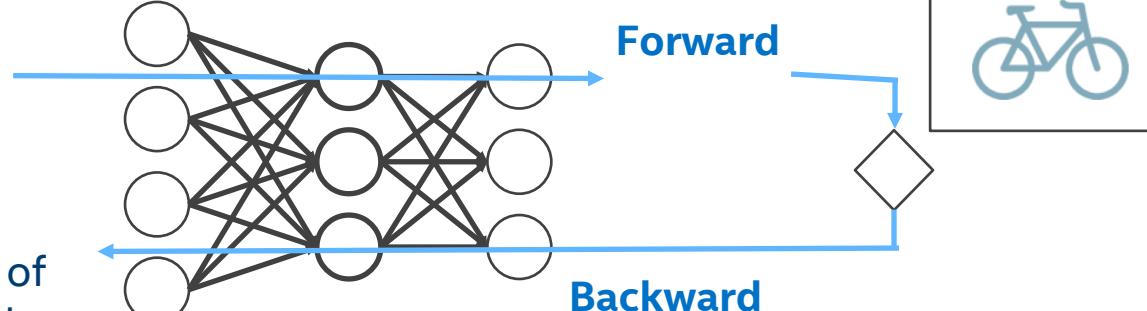
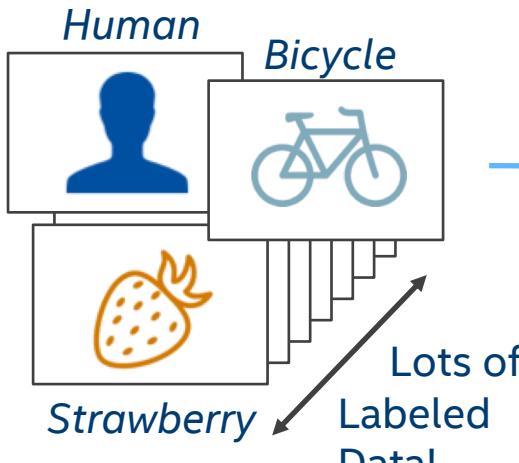


```
jeff@KBL: ~/evs2018_standalone
File Edit View Search Terminal Help
Gtk-Message: 15:18:10.900: Failed to load module "canberra-gtk-module"
Preprocess: 7.40115 ms/frame
Inference: 37.3826 ms/frame
Postprocess:34.8848 ms/frame
inputDims=300 300 3 1
outputDims=1 1 200 7
nuseclasses=20
Gtk-Message: 15:18:31.679: Failed to load module "canberra-gtk-module"
Preprocess: 4.7722 ms/frame
Inference: 40.2658 ms/frame
Postprocess:32.954 ms/frame
inputDims=300 300 3 1
outputDims=1 1 200 7
nuseclasses=20
Gtk-Message: 15:18:52.148: Failed to load module "canberra-gtk-module"
Preprocess: 5.40599 ms/frame
Inference: 44.5869 ms/frame
Postprocess:34.0411 ms/frame
jeff@KBL:~/evs2018_standalone$ ./run_demo2.sh
inputDims=227 227 3 1
outputDims=1 1000 1 1
nuseclasses=1000
Gtk-Message: 15:19:50.392: Failed to load module "canberra-gtk-module"
frame: 41
```



# Deep Learning: Training vs. Inference

## Training

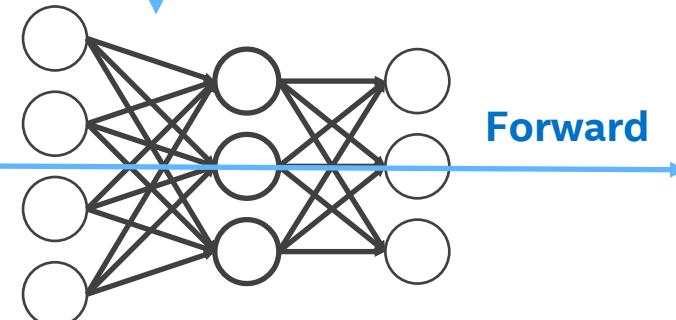


Model Weights

## Inference

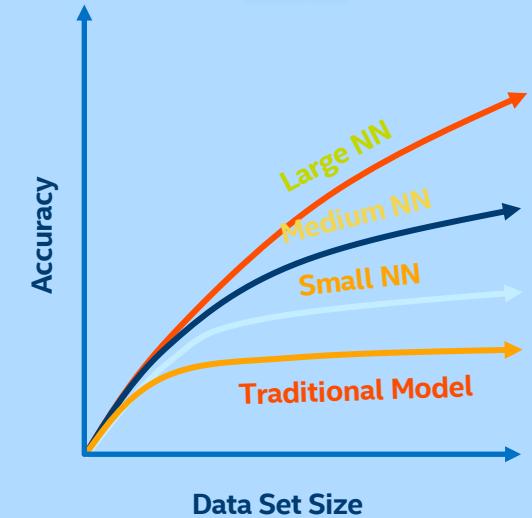


??????

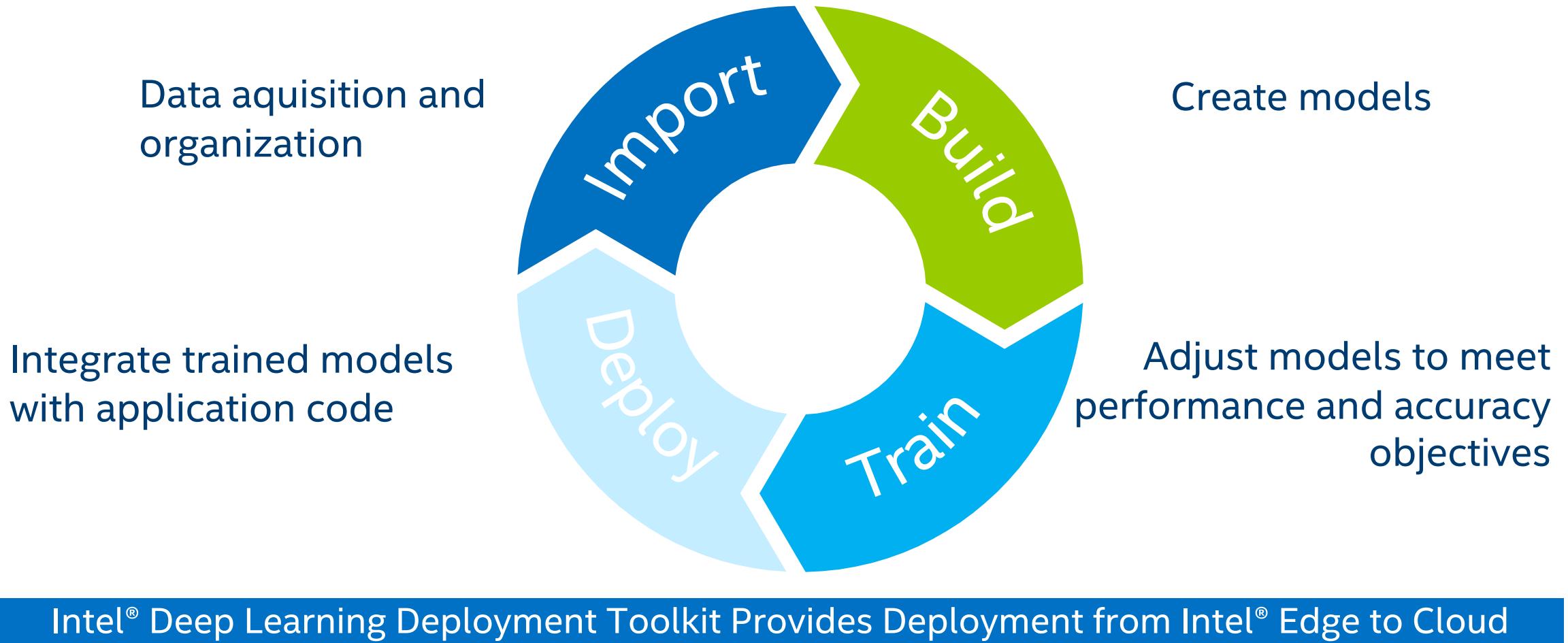


## Did You Know?

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



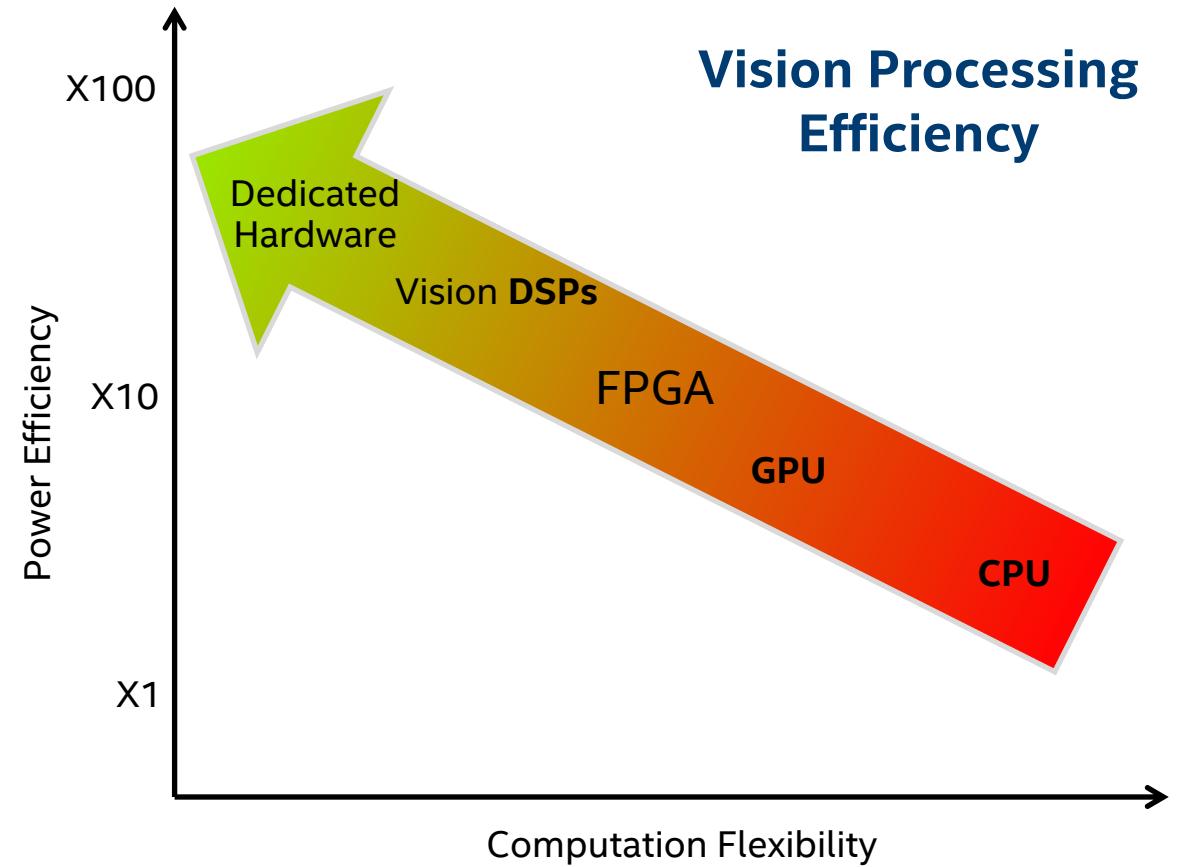
# Artificial Intelligence Development Cycle



# Choosing the “Right” Hardware

## Power/Performance Efficiency Varies

- Running the right workload on the right piece of hardware → higher efficiency
- Hardware acceleration is a must
- Heterogeneous computing?

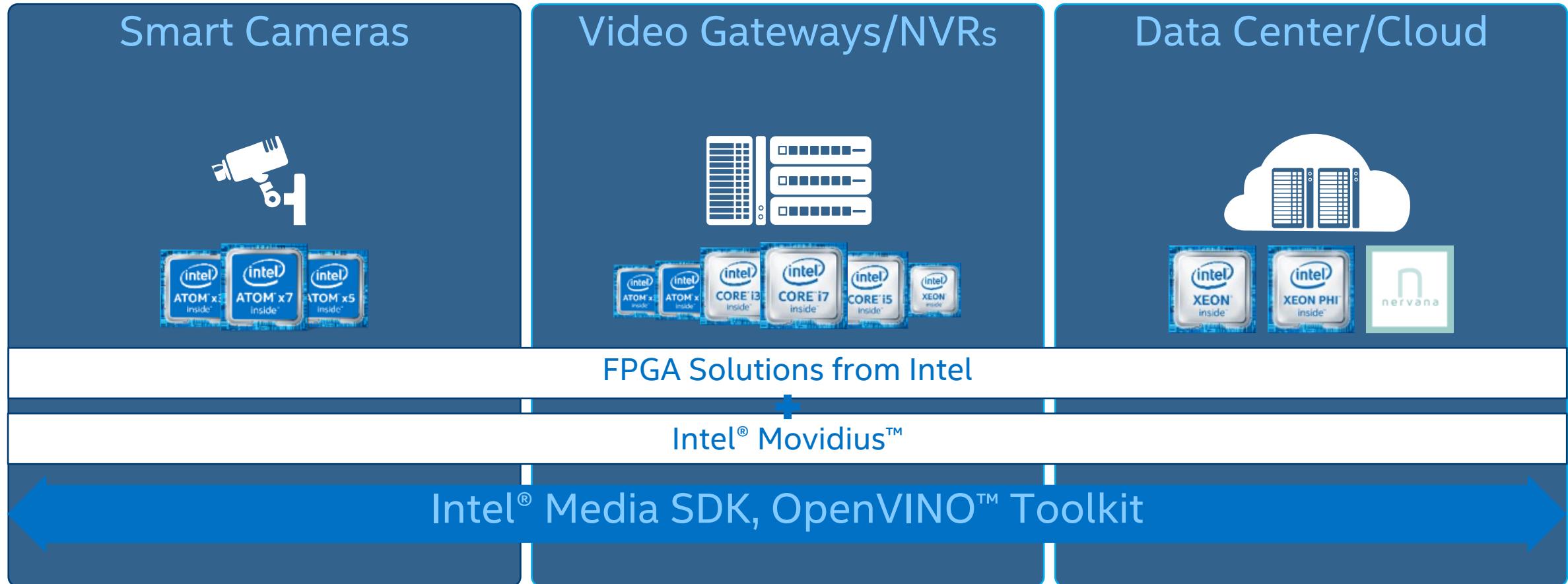


## Tradeoffs

- Power/performance
- Price
- Software flexibility, portability

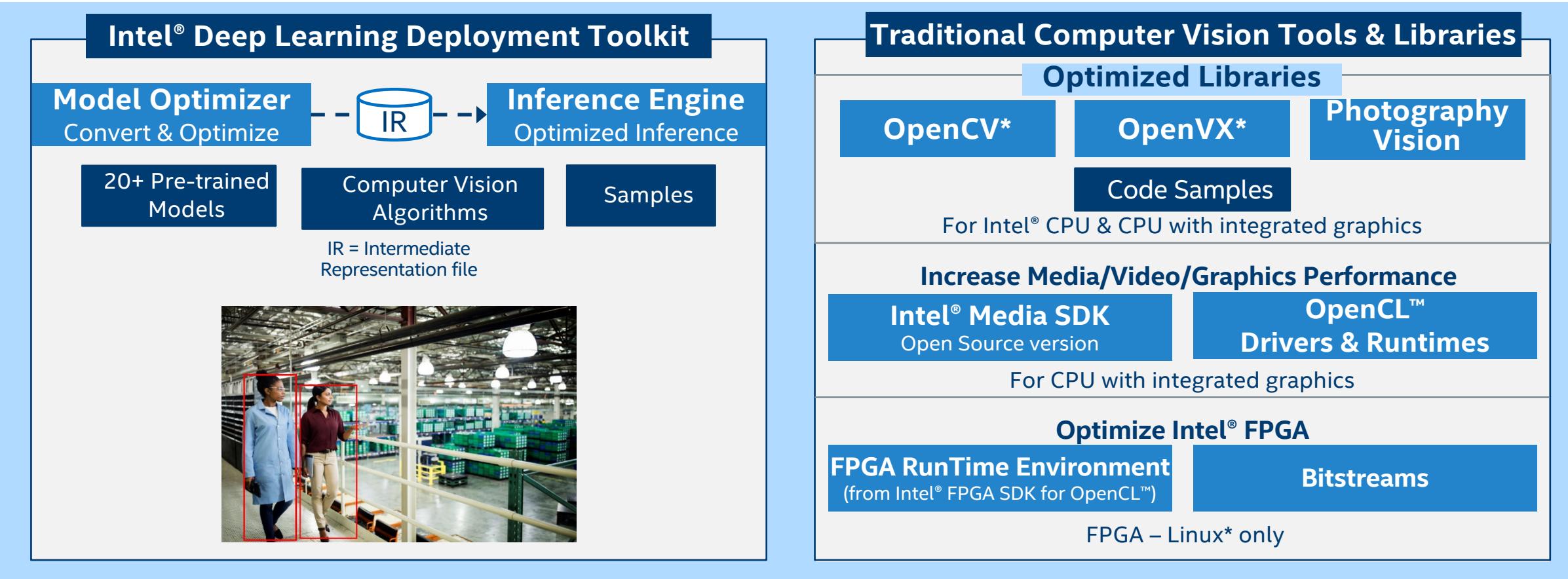
# Intel Internet of Things (IoT) Video Portfolio

## Intel Invests in AI, Computer Vision, and Deep Learning for IoT



Industry's Broadest Media and Computer Vision and Deep Learning Portfolio

# Open Visual Inference & Neural network Optimization (OpenVINO™) toolkit & Components

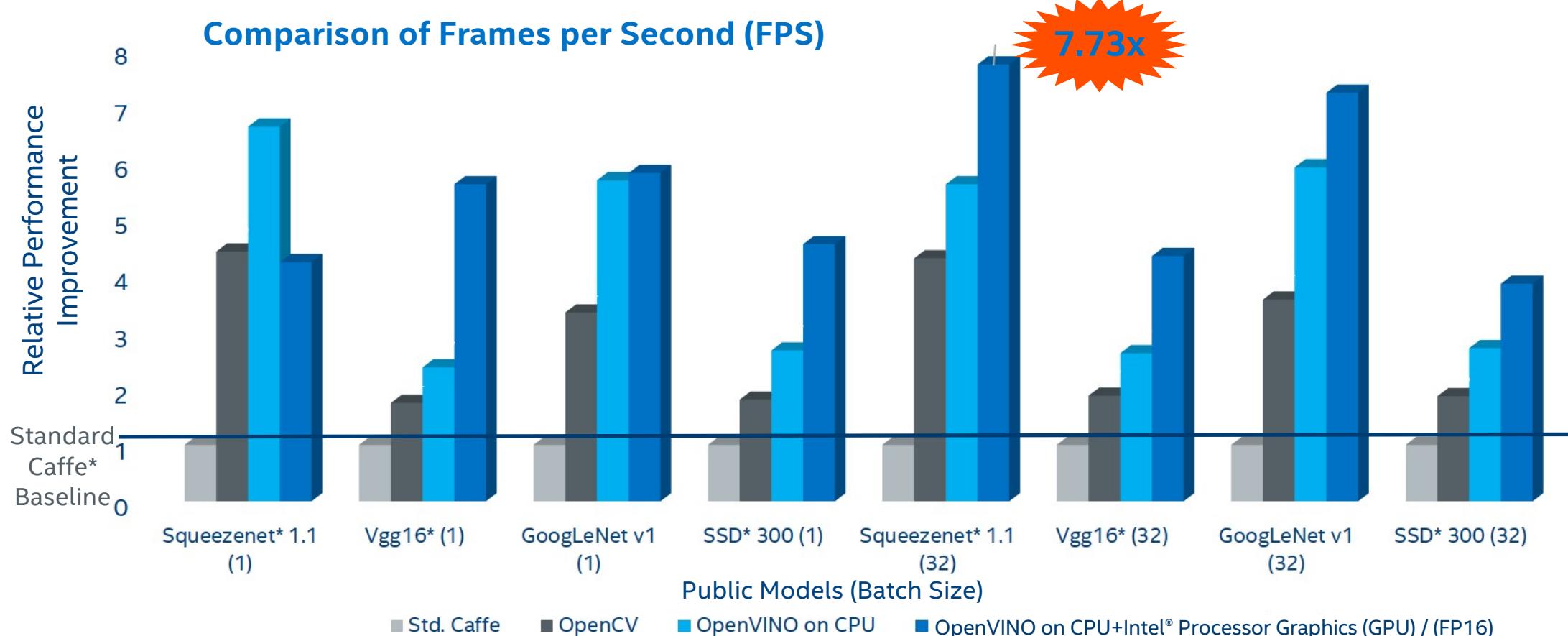


**OS Support** CentOS\* 7.4 (64 bit) Ubuntu\* 16.04.3 LTS (64 bit) Microsoft Windows\* 10 (64 bit) Yocto Project\* version Poky Jethro v2.0.3 (64 bit)

Intel® Architecture-Based Platforms Support



# Increase Deep Learning Workload Performance on Public Models using OpenVINO™ toolkit & Intel® Architecture

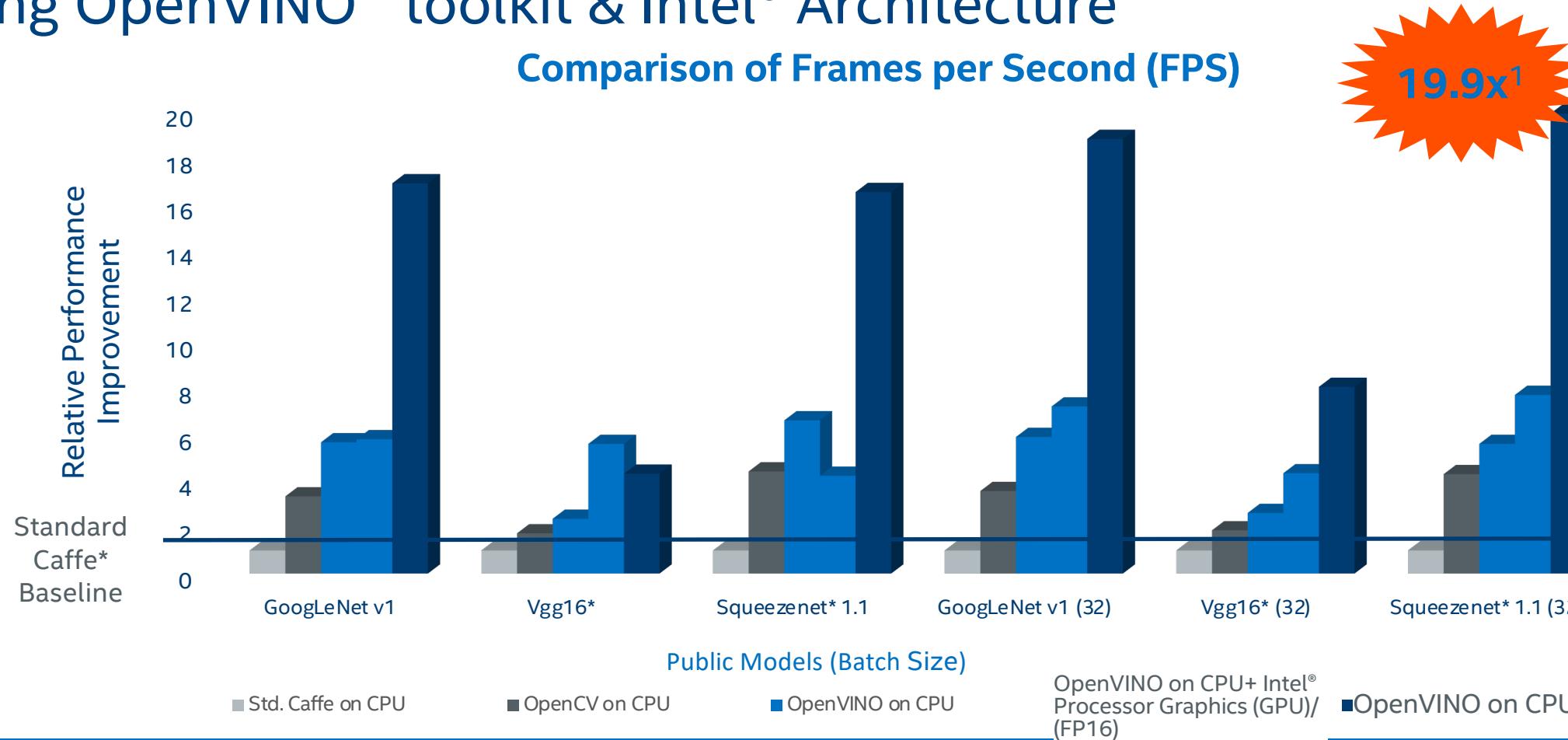


Fast Results on Intel Hardware, even before using Accelerators

1Depending on workload, quality/resolution for FP16 may be marginally impacted. A performance/quality tradeoff from FP32 to FP16 can affect accuracy; customers are encouraged to experiment to find what works best for their situation. The benchmark results reported in this deck may need to be revised as additional testing is conducted. Performance results are based on testing as of April 10, 2018 and may not reflect all publicly available security updates. See configuration disclosure for details. No product can be absolutely secure. For more complete information about performance and benchmark results, visit [www.intel.com/benchmarks](http://www.intel.com/benchmarks).

**Configuration:** Testing by Intel as of April 10, 2018. Intel® Core™ i7-6700K CPU @ 2.90GHz fixed, GPU GT2 @ 1.00GHz fixed Internal ONLY testing, Test v312.30 – Ubuntu\* 16.04, OpenVINO™ 2018 RC4. Tests were based on various parameters such as model used (these are public), batch size, and other factors. Different models can be accelerated with different Intel hardware solutions, yet use the same Intel software tools.

# Increase Deep Learning Workload Performance on Public Models using OpenVINO™ toolkit & Intel® Architecture

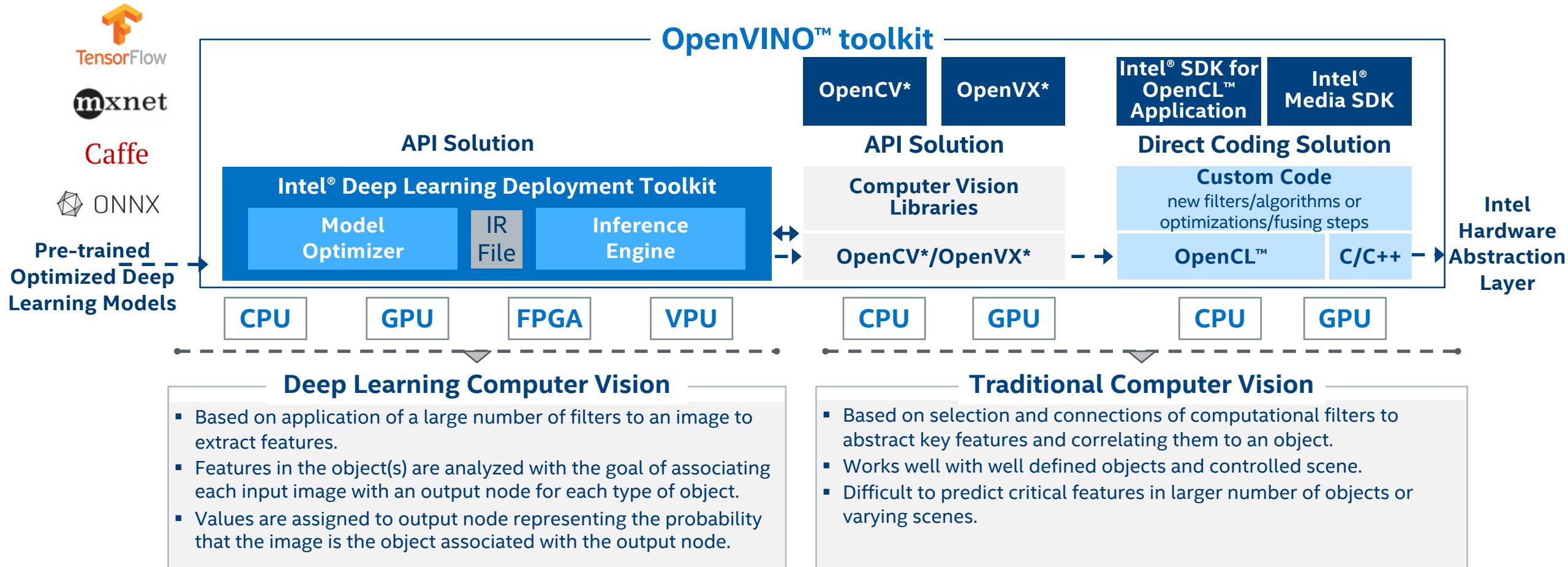


Get an even Bigger Performance Boost with Intel® FPGA

<sup>1</sup>Depending on workload, quality/resolution for FP16 may be marginally impacted. A performance/quality tradeoff from FP32 to FP16 can affect accuracy; customers are encouraged to experiment to find what works best for their situation. Performance results are based on testing as of June 13, 2018 and may not reflect all publicly available security updates. See configuration disclosure for details. No product can be absolutely secure. For more complete information about performance and benchmark results, visit [www.intel.com/benchmarks](http://www.intel.com/benchmarks). Configuration: Testing by Intel as of June 13, 2018. Intel® Core™ i7-6700K CPU @ 2.90GHz fixed, GPU GT2 @ 1.00GHz fixed Internal ONLY testing, Test v3.15.21 – Ubuntu\* 16.04, OpenVINO 2018 RC4, Intel® Arria® 10 FPGA 1150GX. Tests were based on various parameters such as model used (these are public), batch size, and other factors. Different models can be accelerated with different Intel hardware solutions, yet use the same Intel software tools.

# Deep Learning vs. Traditional Computer Vision

OpenVINO™ toolkit has Tools for an End-to-End Vision Pipeline



IR = Intermediate Representation File

GPU = Intel CPU with integrated graphics processing unit/Intel® Processor Graphics

VPU = Intel® Movidius™ Vision Processing Unit

## Optimization Notice

Copyright © 2018, Intel Corporation. All rights reserved.

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

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



# Application development with OpenVINO™ Toolkit

## Train

Train a DL model.  
Currently supports:

- Caffe\*
- Mxnet\*
- TensorFlow\*
- ONNX\*



Caffe

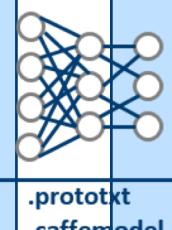


## Prepare Optimize

Model optimizer:

- Converting
- Optimizing
- Preparing to inference

(device agnostic,  
generic optimization)



## Inference

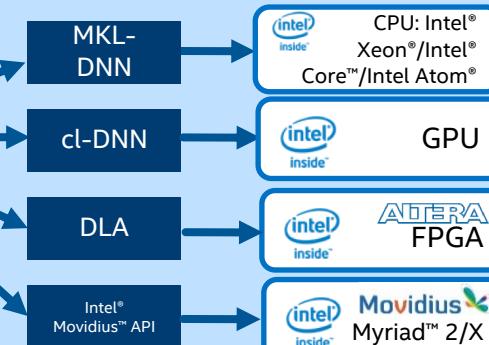
Inference engine lightweight API to use in applications for inference.

User Application

Inference Engine

## Optimize/ Heterogeneous

Inference engine supports multiple devices for heterogeneous flows.  
(device-level optimization)

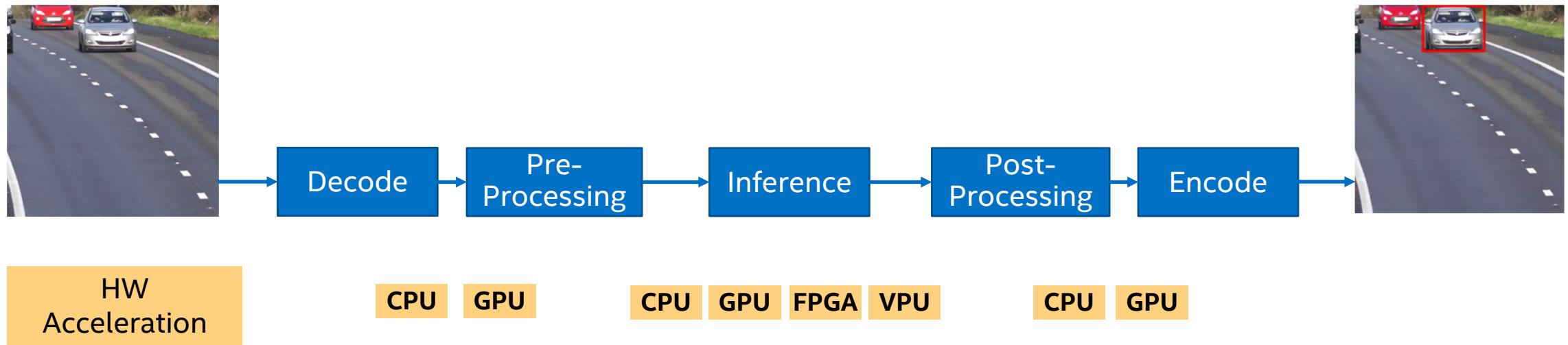


## Extend

Inference engine supports extensibility and allows custom kernels for various devices.



# Full Pipeline Optimization



# Intel® Media SDK

## API to Access Intel® Quick Sync Video: Hardware Accelerated Encoding, Decoding, and Processing

- H.265 (HEVC)
- H.264 (AVC)
- MPEG-2 and more
- Resize, scale, deinterlace
- Color conversion, composition
- Denoise, sharpen, and more

## Benefits

- Outstanding performance
- Rich API to tune encoding pipeline
- Future proofed: support new processor without code changes

## Targeting Digital Security and Surveillance, Connected Car Applications, and More



Smart Camera

Car Infotainment and Cluster Display

using



Intel Atom®, Pentium®, and Celeron®<sup>1</sup>

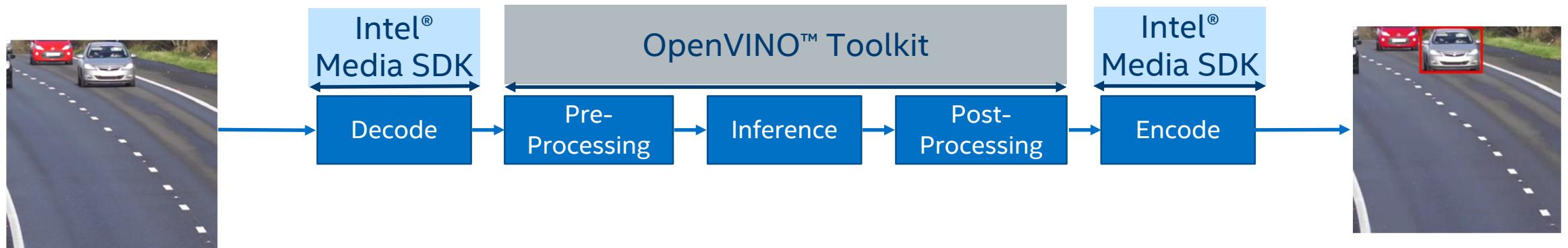
Embedded Linux\*



<sup>1</sup> Intel® Celeron® Processor N3350, Intel® Pentium® Processor N4200, Intel Atom® E3930, E3940, E3950 processors

# Accelerate Streaming Performance, Integrate Video Analytics Computer Vision Needs Intel® Media SDK

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

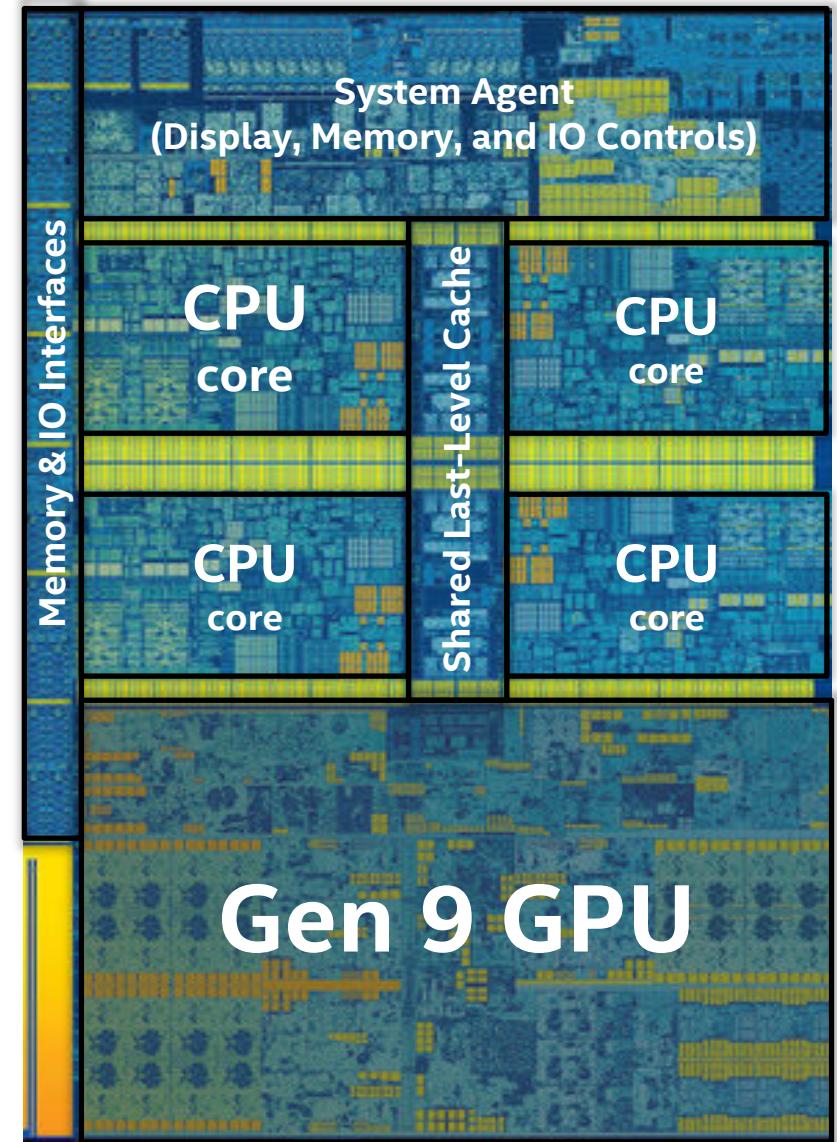


# Intel Integrated Graphics

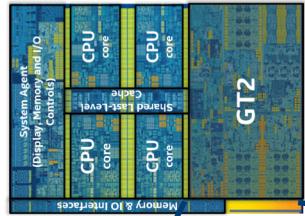
**Gen** is the internal name for Intel's on-die GPU solution. It's a hardware ingredient with various configurations.

- Intel® Core™ Processors include Gen hardware.
- Gen GPUs can be used for graphics and also as general compute resources.
- Libraries contained in the OpenVINO™ toolkit (and many others) support Gen offload using OpenCL™.

6<sup>th</sup> Generation Intel® Core™ i7 (Skylake) Processor



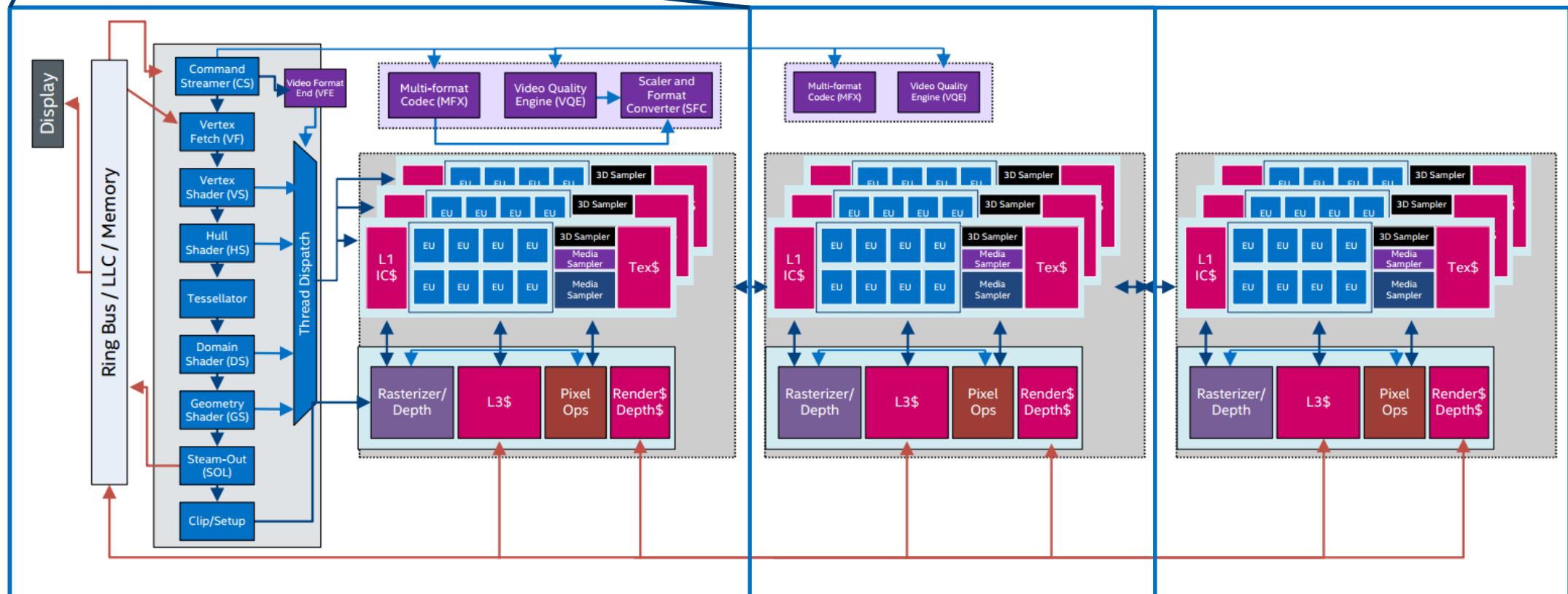
# Intel GPU Configurations



GT2  
**Intel® HD Graphics**  
24 EUs, 1 MFX

GT3  
**Intel® Iris® Graphics**  
48 EUs, 2 MFX

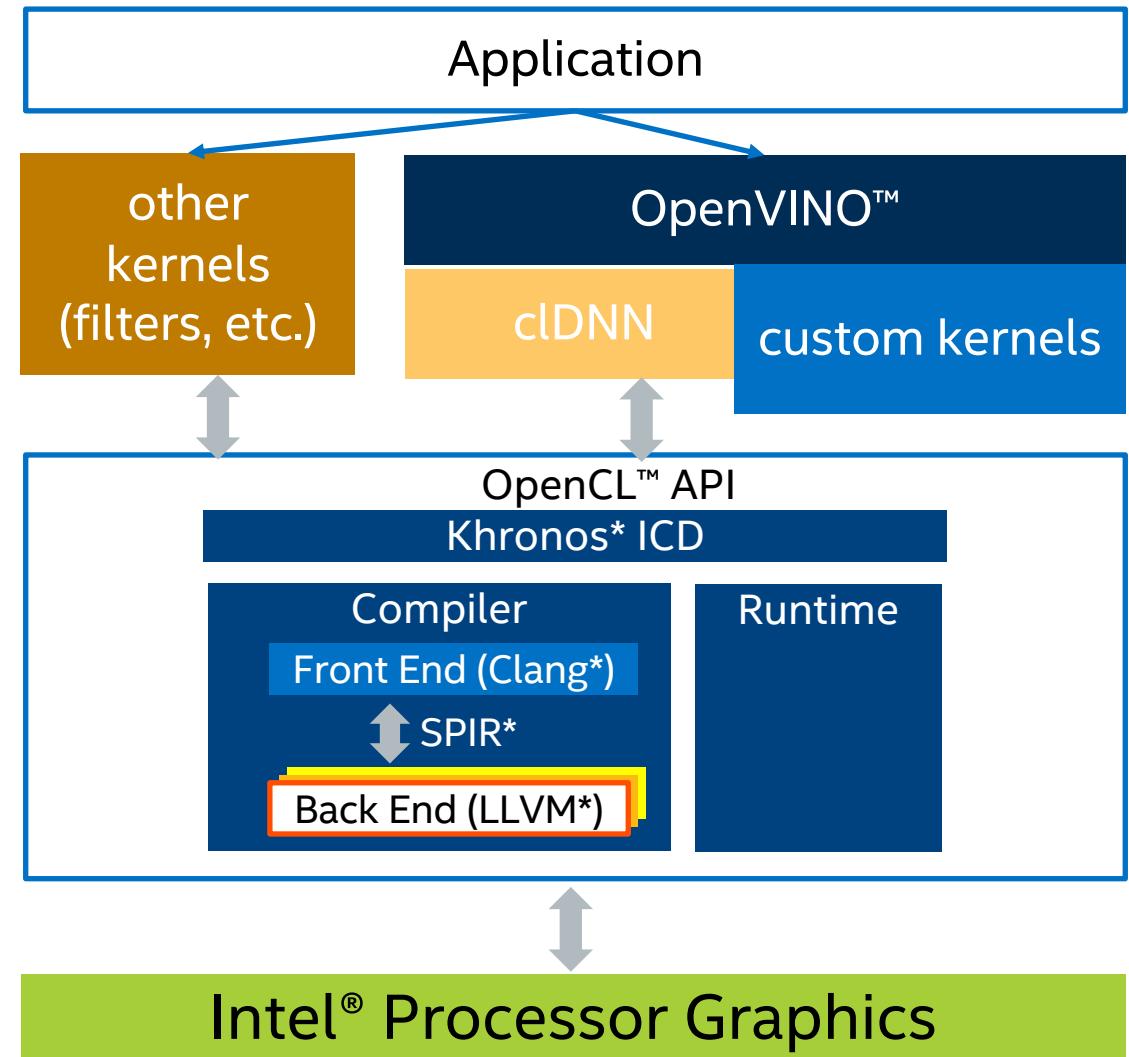
GT4  
**Intel® Iris® Pro Graphics**  
72 EUs, 2 MFX



# OpenCL™

## OpenCL™:

- Required to run with a GPU target (clDNN) using Intel® Processor Graphics
- Custom kernels
- Other kernels can be used for other non-inference pipeline stages, such as color conversions



# Putting It All Together

- A major challenge is to get all these tool and libraries to work together in the best possible way to minimize development time and optimize system power/performance.
- A good way to abstract that workload is using an end-to-end pipeline

## Computer Vision



## Deep Learning



## Media



### SDKs



Optimized CV  
Capabilities



Intel® Deep Learning  
Deployment Toolkit

*OpenVINO™ Toolkit*



Intel® Media SDK

### Tools

Compiler, Analyzers, Debuggers



### Libraries

IPP



TBB



Intel®  
MKL-DNN



Intel® MKL  
DAAL



# 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.

## OpenVINO™ 101

### Hardware Acceleration

### Optimization

### Application

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

