

ACCELERATE DEEP LEARNING INFERENCE USING INTEL TECHNOLOGIES

OBJECT DETECTION USING INTEL[®] DISTRIBUTION OF OPENVINO[™] TOOLKIT

February 2020

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

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.

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/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.

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

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.

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.



SMART VIDEO WORKSHOP OVERVIEW

INTRODUCTION

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

**Intel®
Distribution of
OpenVINO™ 101**

**Hardware
Acceleration on
laptop and DevCloud**

2. Basic End-to-End Object Detection Example

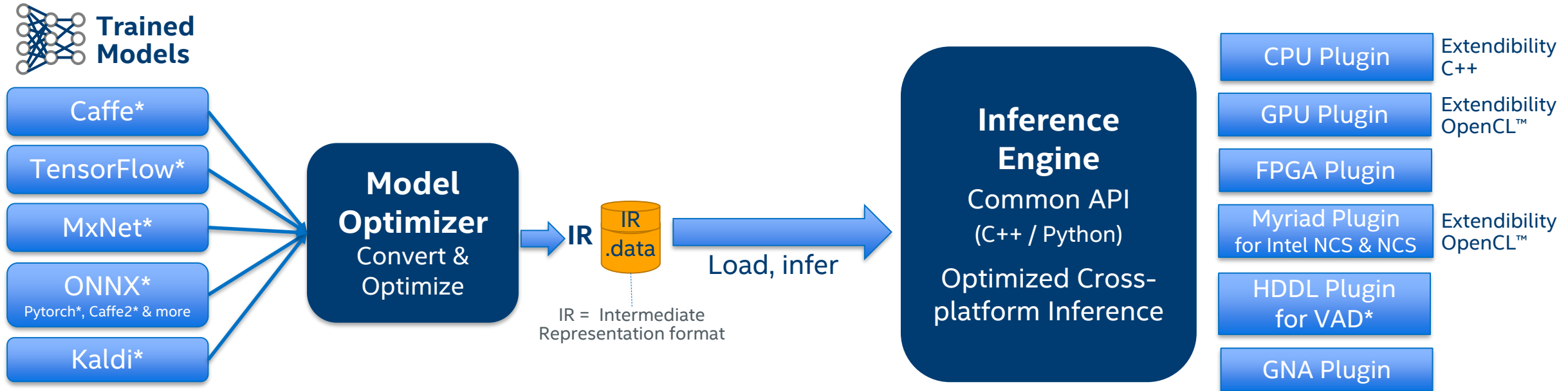
3./4./5. Hardware Acceleration with CPU, Integrated GPU, Intel® Movidius™ NCS, FPGA

INTEL[®] DEEP LEARNING DEPLOYMENT TOOLKIT

FOR DEEP LEARNING INFERENCE

Model Optimizer

- **What it is:** A Python*-based tool to import trained models and convert them to Intermediate representation.
- **Why important:** Optimizes for performance/space with conservative topology transformations; biggest boost is from conversion to data types matching hardware.



Inference Engine

- **What it is:** High-level inference API
- **Why important:** 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.

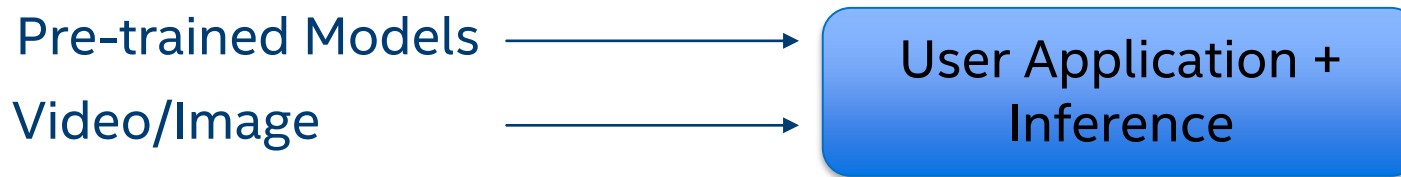
GPU = Intel CPU with integrated GPU/Intel[®] Processor Graphics, Intel[®] NCS = Intel[®] Neural Compute Stick (VPU)

*VAD = Intel[®] Vision Accelerator Design Products (HDDL-R)

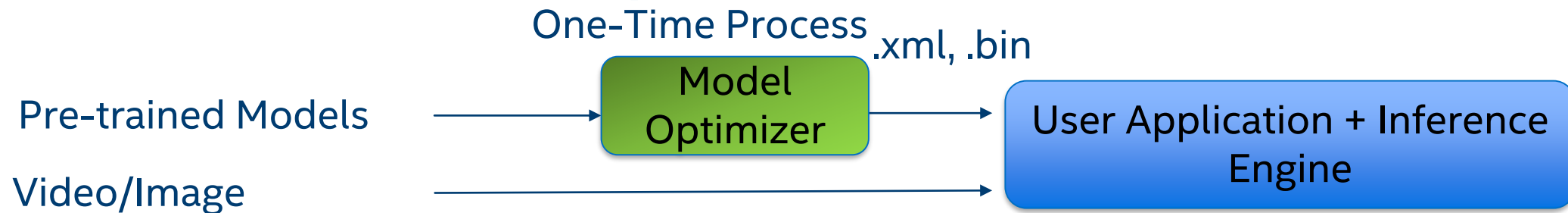


DEEP LEARNING APPLICATION DEPLOYMENT

Traditional

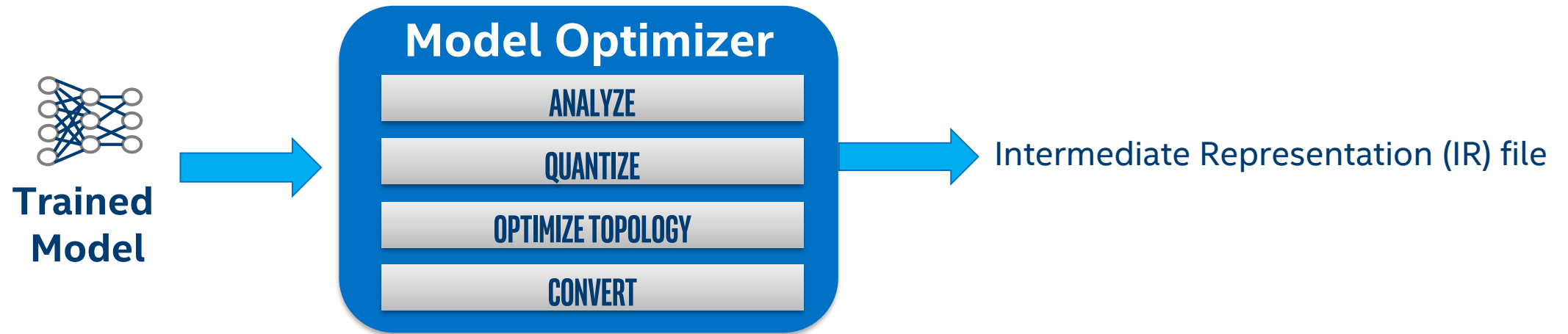


With Intel® Distribution of OpenVINO™ Toolkit



MODEL OPTIMIZER

IMPROVE PERFORMANCE WITH MODEL OPTIMIZER



- Easy to use, Python*-based workflow does not require rebuilding frameworks.
- Import Models from many supported frameworks: Caffe*, TensorFlow*, MXNet*, Kaldi*, exchange formats like ONNX* (Pytorch*, Caffe2* and others through ONNX).
- 100+ models for Caffe, MXNet, TensorFlow validated. Supports all ONNX* model zoo public models.
- Extends inferencing for non-vision networks with support of LSTM, Bert, GNMT, TDNN-LSTM, ESPNet and more.
- IR files for models using standard layers or user-provided custom layers do not require Caffe.
- Fallback to original framework is possible in cases of unsupported layers, but requires original framework.

MODEL OPTIMIZER

Model optimizer performs generic optimization:

- Node merging
- Horizontal fusion
- Batch normalization to scale shift
- Fold scale shift with convolution
- Drop unused layers (dropout)
- FP16/Int8 quantization
- Model optimizer can add normalization and mean operations, so some preprocessing is 'added' to the IR

`--mean_values (104.006, 116.66, 122.67)`

`--scale_values (0.07, 0.075, 0.084)`

PLUGIN	FP32	FP16	INT8
CPU plugin	Supported and preferred	Supported	Supported
GPU plugin	Supported	Supported and preferred	Supported*
FPGA plugin	Supported	Supported	Not supported
VPU plugins	Not supported	Supported	Not supported
GNA plugin	Supported	Supported	Not supported

MODEL OPTIMIZATION TECHNIQUES

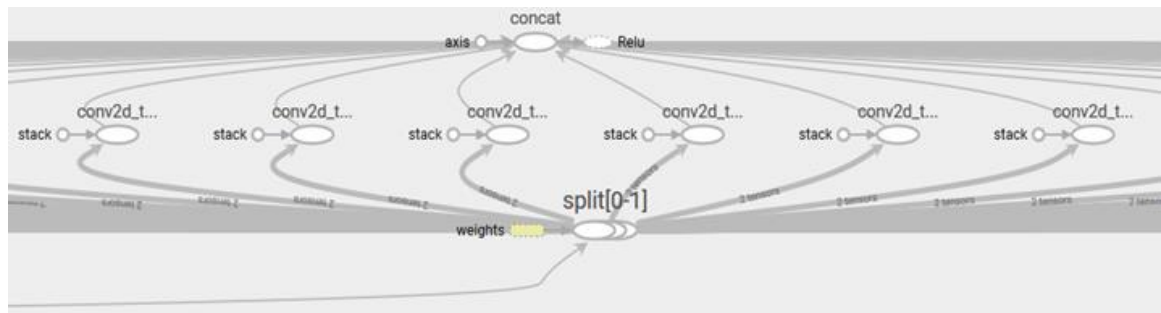
LINEAR OPERATION FUSING & GROUPED CONVOLUTIONS FUSING

Linear Operation Fusing: 3 stages

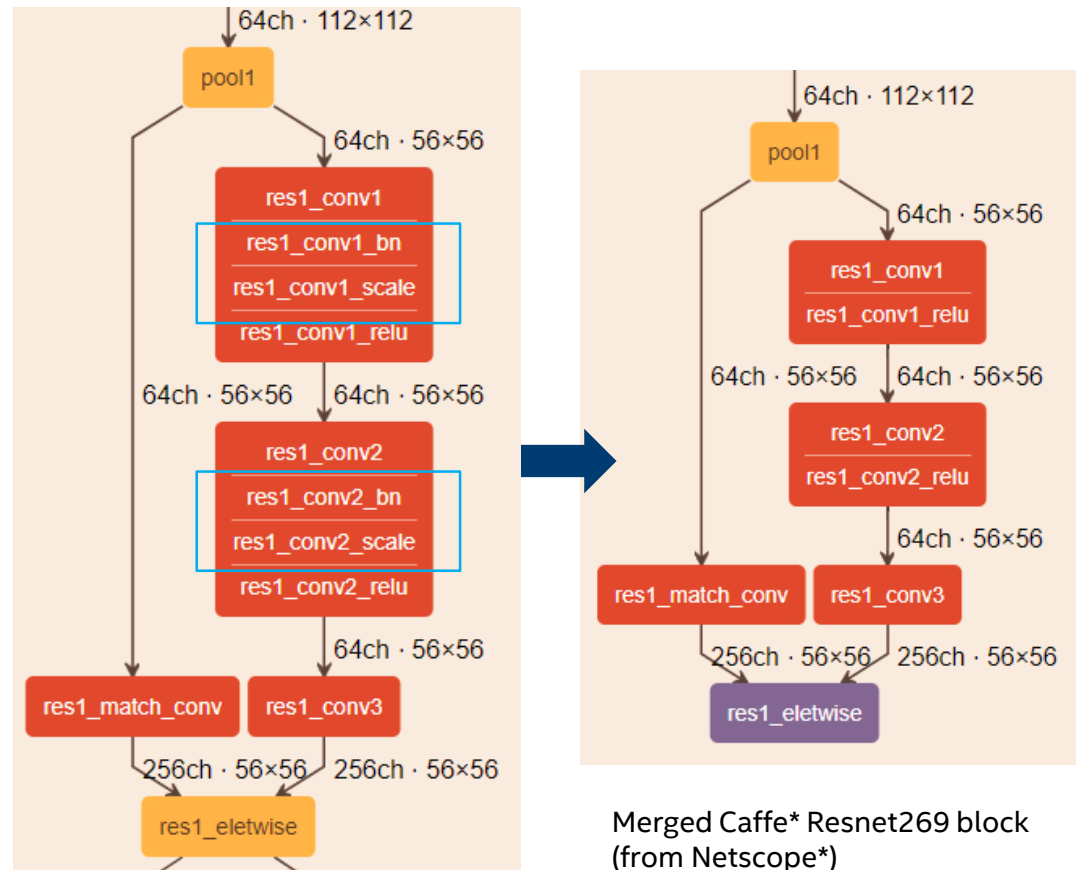
1. **BatchNorm and ScaleShift decomposition:** BN layers decomposes to *Mul->Add->Mul->Add* sequence; ScaleShift layers decomposes to *Mul->Add* sequence.
2. **Linear operations merge:** Merges sequences of Mul and Add operations to the **single** Mul->Add instance.
3. **Linear operations fusion:** Fuses Mul and Add operations to Convolution or FullybConnected layers.

Grouped Convolutions Fusing

Specific optimization that applies for TensorFlow* topologies. (Xception*)



Split->Convolutions->Concat block from TensorBoard*



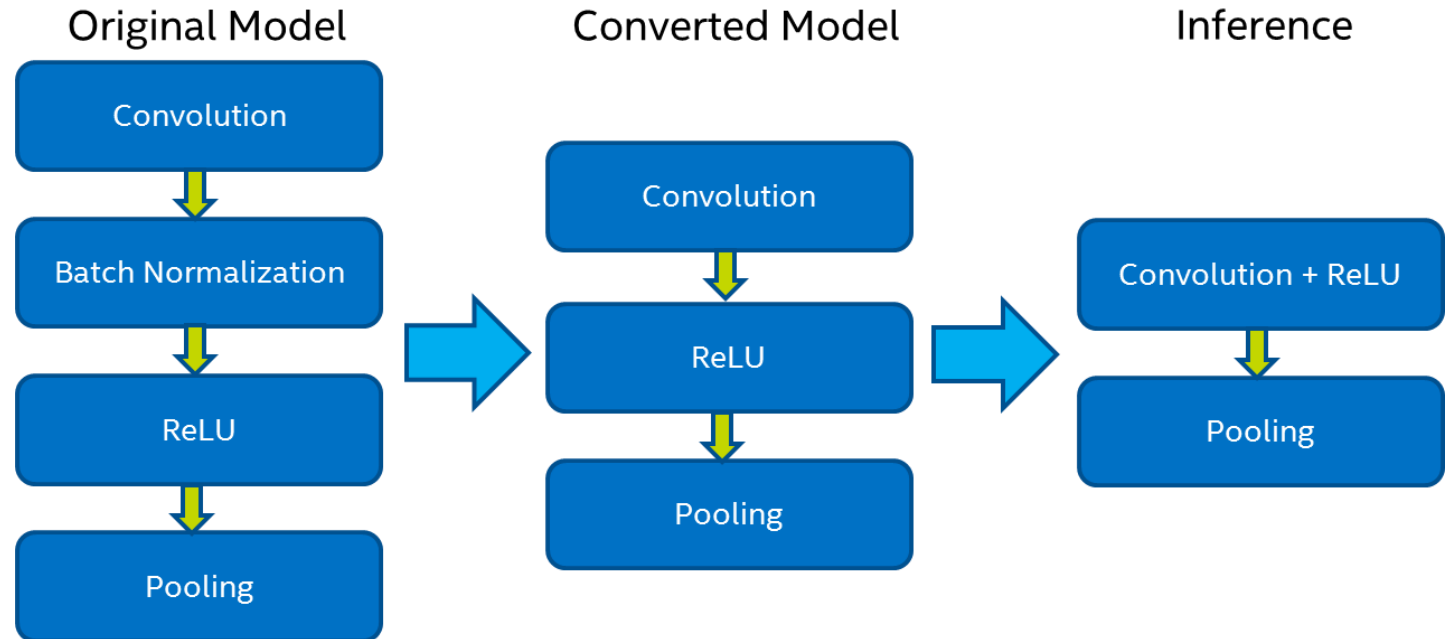
Caffe Resnet269 block (from Netscope)

Merged Caffe* Resnet269 block
(from Netscope*)

MODEL OPTIMIZER: LINEAR OPERATION FUSING

Example

1. Remove Batch normalization stage.
2. Recalculate the weights to 'include' the operation.
3. Merge Convolution and ReLU into one optimized kernel.



MODEL OPTIMIZER: CUTTING OFF PARTS OF A MODEL

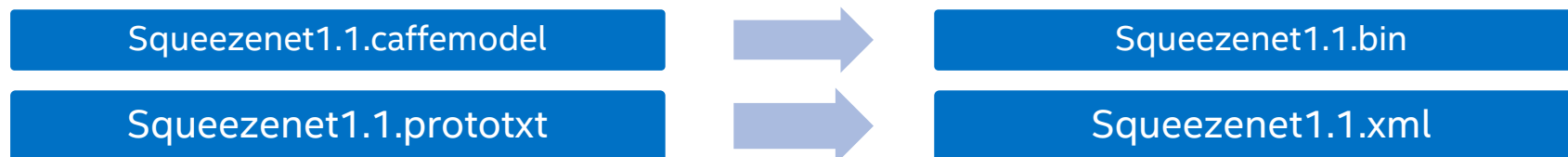
Model optimizer can cut out a portion of the network:

- Model has pre/post-processing parts that cannot be mapped to existing layers.
- Model has a training part that is not used during inference.
- Model is too complex and cannot be converted in one shot.

Command line options

- Model Optimizer provides command line options --input and --output to specify new entry and exit nodes ignoring the rest of the model:
- --input option accepts a comma-separated list of layer names of the input model that should be treated as new entry points to the model;
- --output option accepts a comma-separated list of layer names of the input model that should be treated as new exit points from the model.

INTERMEDIATE REPRESENTATION (IR)



```
layer {
  name: "data"
  type: "Input"
  top: "data"
  input_param { shape: { dim: 1 dim: 3 dim: 227 dim: 227 } }
}
layer {
  name: "conv1"
  type: "Convolution"
  bottom: "data"
  top: "conv1"
  convolution_param {
    num_output: 64
    kernel_size: 3
    stride: 2
  }
}
```

```
<net batch="1" name="model" version="2">
  <layers>
    <layer id="100" name="data" precision="FP32" type="Input">
      <output>
        <port id="0">
          <dim>1</dim>
          <dim>3</dim>
          <dim>227</dim>
          <dim>227</dim>
        </port>
      </output>
    </layer>
    <layer id="129" name="conv1" precision="FP32" type="Convolution">
      <data dilation-x="1" dilation-y="1" group="1" kernel-x="3" kernel-y="3" output="64" pa
      <input>
        <port id="0">
          <dim>1</dim>
          <dim>3</dim>
          <dim>227</dim>
          <dim>227</dim>
        </port>
      </input>
      <output>
        <port id="3">
          <dim>1</dim>
          <dim>64</dim>
          <dim>113</dim>
          <dim>113</dim>
        </port>
      </output>
      <blobs>
        <weights offset="2275104" size="6912"/>
        <biases offset="4805920" size="256"/>
      </blobs>
    </layer>
  </layers>
</net>
```

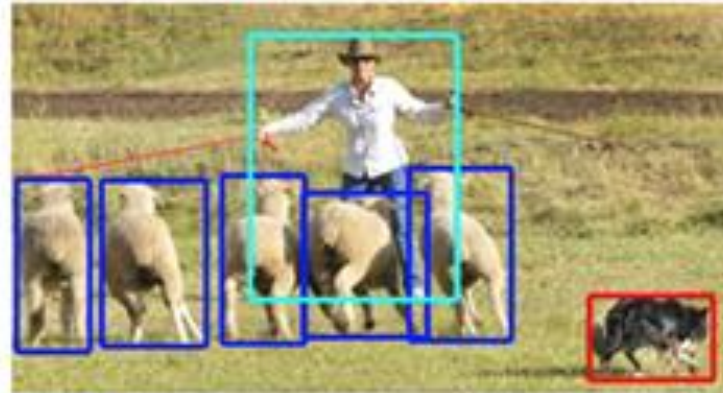
INFERENCE ENGINE

INFERENCE ON AN INTEL® EDGE SYSTEMS

Many deep learning networks are available—choose the one you need.



(a) classification



(b) detection

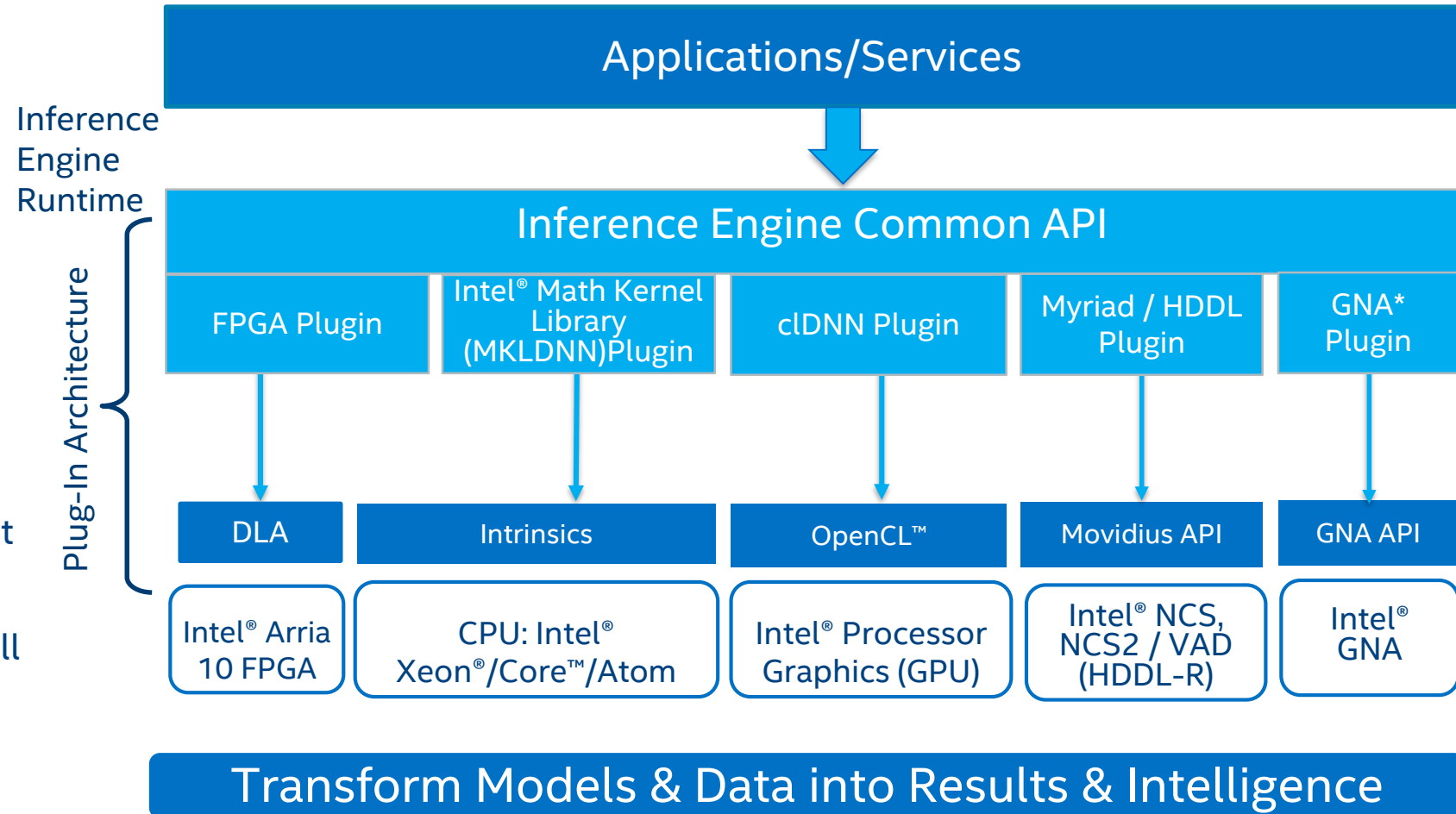


(c) segmentation

The complexity of the problem (data set) dictates the network structure. The more complex the problem, the more 'features' required, the deeper the network.

OPTIMAL MODEL PERFORMANCE USING THE INFERENCE ENGINE

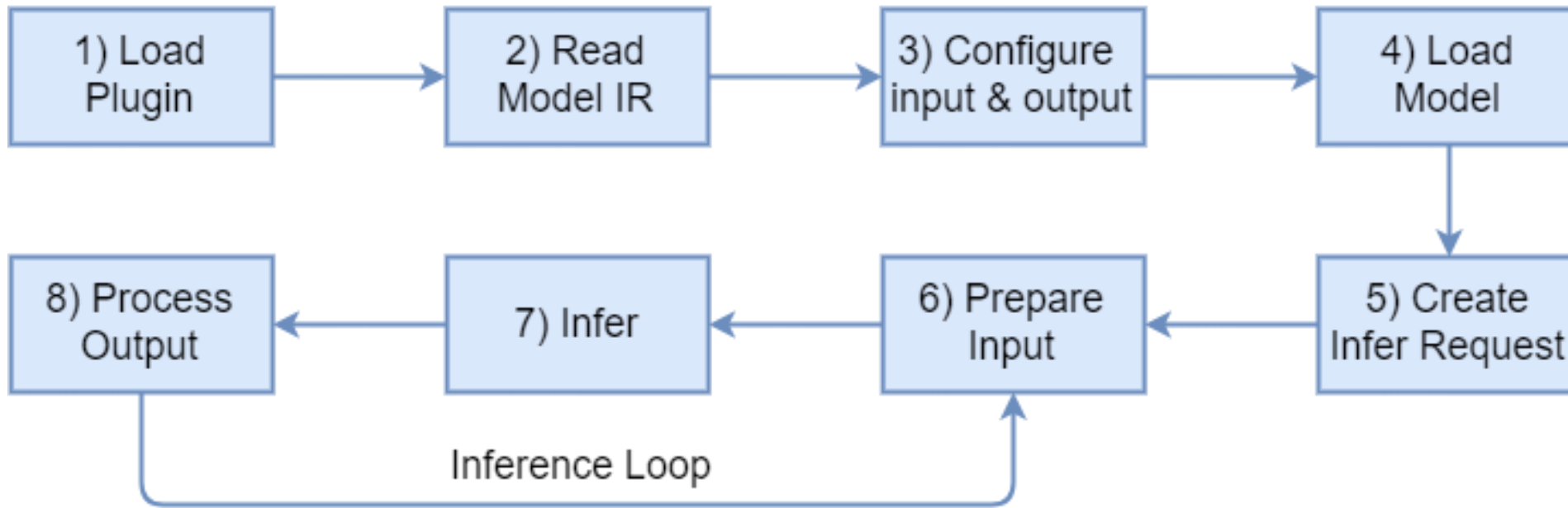
- Simple & unified API for inference across all Intel® architecture
- Optimized inference on large IA hardware targets (CPU/GEN/FPGA)
- Heterogeneity support allows execution of layers across hardware types
- Asynchronous execution improves performance
- Futureproof/scale your development for future Intel® processors
- Supports serialized FP16 IR across all plugins / platforms (CPU inference remains at FP32)



GPU = Intel CPU with integrated graphics/Intel® Processor Graphics/GEN
GNA = Gaussian mixture model and Neural Network Accelerator



APPLICATION WORKFLOW



http://docs.openvinotoolkit.org/latest/_docs_IE_DG_Integrate_with_customer_application_new_API.html

INFERENCE ENGINE ASYNC API

- Improves overall frame-rate of the application.
- Executes a request asynchronously (in the background) and waits until ready, when the result is actually needed.
- Continues doing things on the host, while the accelerator is busy.
- The demo keeps two parallel infer requests, and, while the current is processed, the input frame for the next is captured. This essentially hides the latency of capturing, so the overall framerate is determined by the `MAXIMUM(detection time, input capturing time)` and not the `SUM(detection time, input capturing time)`.

PROCESS OUTPUT

Developers are responsible for parsing inference output.

Many output formats are available. Some examples include:

- **Simple Classification (alexnet*)**: an array of float confidence scores, # of elements=# of classes in the model
- **SSD**: Many “boxes” with a confidence score, label #, xmin,ymin, xmax,ymax

Unless a model is well documented, the output pattern may not be immediately obvious.



NETWORK LOADING OPTIMIZATIONS

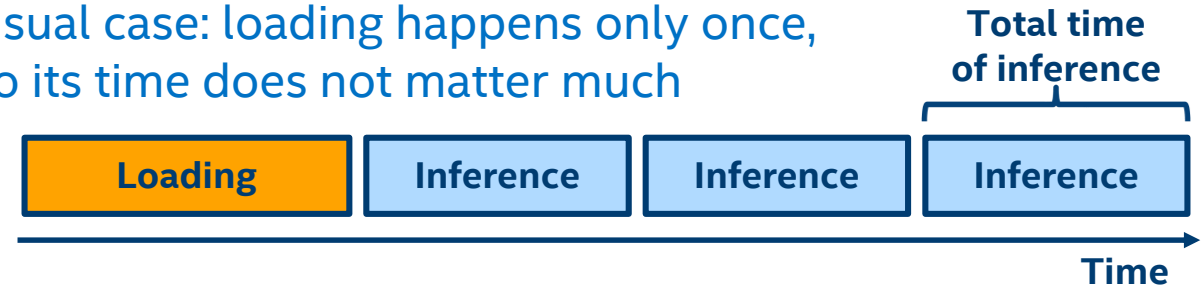
REDUCED MODEL LOAD TIMES FOR FASTER PERFORMANCE

Audience: Helpful when shape size changes from inference to inference, and resizing is undesirable (e.g., leads to accuracy degradation)

Problem: Shape change requires reloading of the model which can be slow

UseCase: Input shape is defined by a previous network in the pipeline (i.e., in the case of object detection/classification), or ROI is defined by operator (common case in medical applications)

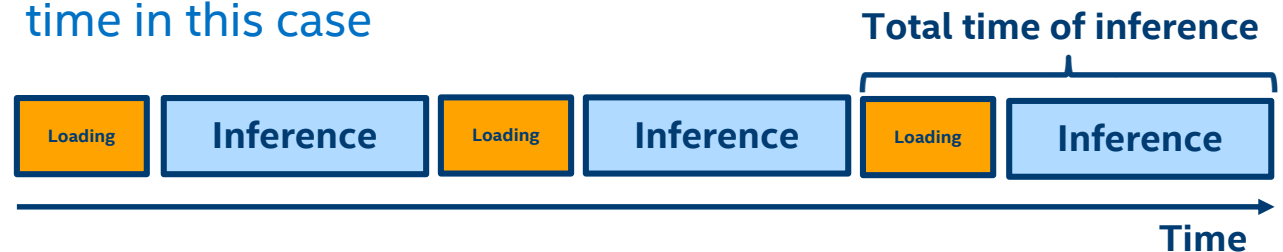
Usual case: loading happens only once, so its time does not matter much



Changing shape case: loading happens before each inference, so its time matters



Loading time optimization decreases total inference time in this case



LAB1 - BASIC END TO END OBJECT DETECTION EXAMPLE

URL: <https://github.com/intel-iot-devkit/smart-video-workshop/blob/master/object-detection/README.md>

Objective: This Lab uses a Single Shot MultiBox Detector (SSD) on a trained mobilenet-ssd* Caffe model to walk you through the basic steps of using two key components of the Intel® Distribution of OpenVINO™ toolkit: Model Optimizer and Inference Engine.

Estimated Complete Time: 30min



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

2. Basic End-to-End Object
Detection Example

HARDWARE ACCELERATION ON LAPTOP
AND DEVCLOUD

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

OPTIMIZATION

6. Optimization Tools and
Techniques

APPLICATION

7. Advanced Video Analytics

CUSTOM LAYERS

8. Custom layers

CPU/GPU/FPGA/Intel® Movidius™ Neural Compute Stick

HARDWARE HETEROGENITY

HETEROGENEOUS DICHOTOMIES

Data vs. Task Parallelism

- Data parallelism between devices is usually explicit and extremely error-prone (like in OpenCL™).
- Inference engine heterogeneity support is inherently **task-oriented** (node-level).

Dynamic vs. Static

- **User-defined work split** vs. adaptive schemes (like load-balancing).
- For both, there is a question on the granularity (communications costs, next).

Heterogeneity is basically Fallback.

INTRODUCING HETEROGENEOUS PLUGIN

The heterogeneous plugin enables computing for inference on one network on several devices. Purposes to execute networks in heterogeneous mode

- To utilize accelerators power and calculate heaviest parts of network on accelerator and execute not supported layers on fallback devices like CPU
- To utilize all available hardware more efficiently during one inference

The execution through heterogeneous plugin can be divided to two independent steps:

- Setting of affinity to layers (binding them to devices in `InferenceEngine::ICNNNetwork`)
- Loading a network to the Heterogeneous plugin, splitting the network to parts, and executing them through the plugin

These steps are decoupled. The setting of affinity can be done automatically using fallback policy or in manual mode.

APPLY DEVICE AFFINITIES TO LAYERS AUTOMATICALLY USING FALLBACK POLICY (1 OF 2)

```
$ ./object_detection_sample_ssd -m <path_to_model>/Model.xml -i  
<path_to_pictures>/picture.jpg -d HETERO:FPGA,CPU
```

The “**priorities**” defines a greedy behavior:

- Keeps all layers that can be executed on the device (FPGA)
- Carefully respects topological and other limitations
- Follows priorities when searching (for example, CPU)

APPLY DEVICE AFFINITIES TO LAYERS AUTOMATICALLY USING FALLBACK POLICY (2 OF 2)

```
1. InferenceEngine::PluginDispatcher dispatcher({ FLAGS_pp,
    archPath , "" });
2. InferenceEngine::InferenceEnginePluginPtr enginePtr;
3. enginePtr =
    dispatcher.getPluginByDevice("HETERO:FPGA,CPU");
4. InferencePlugin plugin(enginePtr);
5. CNNNetReader reader;
6. reader.ReadNetwork("Model.xml");
7. reader.ReadWeights("Model.bin");
8. auto executable_network = plugin.LoadNetwork(network,
    {});
```

APPLY DEVICE AFFINITIES TO LAYERS MANUALLY

```
1.InferenceEngine::PluginDispatcher dispatcher({ FLAGS_pp,
    archPath , "" });
2.InferenceEngine::InferenceEnginePluginPtr enginePtr;
3.enginePtr =
    dispatcher.getPluginByDevice("HETERO:FPGA,CPU");
4.HeteroPluginPtr hetero(enginePtr);
5.hetero->SetAffinity(network, { }, &resp);
6.network.getLayerByName("qqq")->affinity = "CPU";
7.InferencePlugin plugin(enginePtr);
8.CNNNetReader reader;
9. reader.ReadNetwork("Model.xml");
10.reader.ReadWeights("Model.bin");
11.auto executable_network = plugin.LoadNetwork(network,
    {});
```



LAB2 - HARDWARE HETEROGENEITY

URL: <https://github.com/intel-iot-devkit/smart-video-workshop/blob/master/hardware-heterogeneity/README.md>

Objective: This example shows how to use hetero plugin to define preferences to run different network layers on different hardware types, then use option -pc to get performance data on each subgraph.

Estimated Complete Time: 20min





INFERENCE ENGINE MULTI DEVICE

MULTI-DEVICE SUPPORT

OOB automatic load-balancing between devices (inference requests level)

- Fully general machinery: any devices combinations
 - **MULTI: HDDL,GPU,<optional>CPU**
 - **CPU+GPU**
 - Multiple MYX sticks, etc
- As easy as “-d **MULTI:HDDL,GPU**” for cmd-line option of your favorite sample
- Also C++ example (Python is similar)

```
//NEW IE-CENTRIC API
```

```
Core ie;
```

```
ExecutableNetwork exec = ie.LoadNetwork(network,{{“DEVICE_PRIORITIES”, “HDDL,GPU”}}, “MULTI”);
```

```
//Old (plugin-centric) API
```

```
auto plugin = PluginDispatcher().getPluginByDevice(“MULTI:CPU,GPU”);
```

```
ExecutableNetwork executable_network = plugin.LoadNetwork(network, config);
```

- Multi-Device implications: CPU is critical resource
 - Rule of thumb: every accelerator need a CPU core for scheduling, data transfers, etc



LAB3 - MULTI-DEVICE PLUGIN

URL: <https://github.com/intel-iot-devkit/smart-video-workshop/blob/master/hardware-heterogeneity/Multi-devices.md>

Objective: This example shows how to use MULTI plugin to share the inference burden onto different hardware types. Here, we will use the command line option to demonstrate MULTI plugin usage.

Estimated Complete Time: 20min



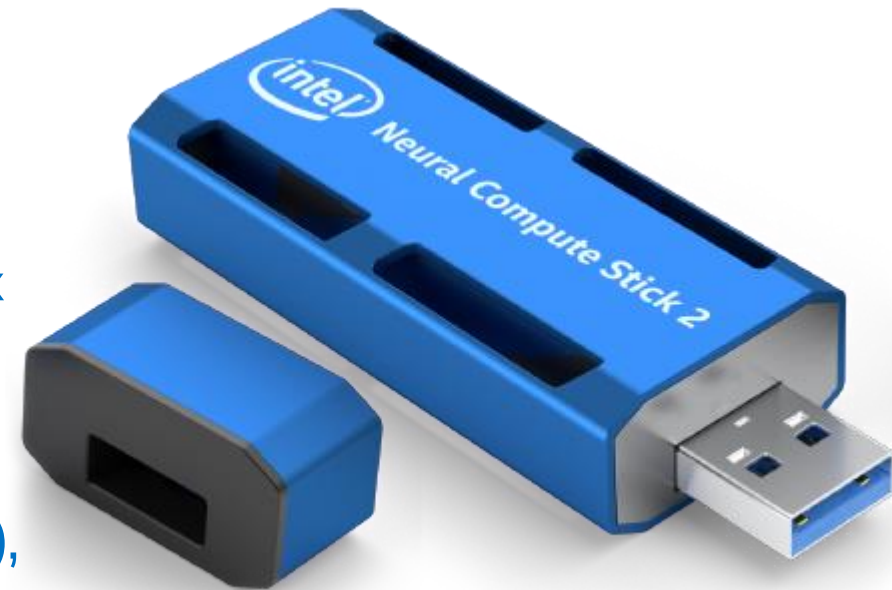
INTEL® MOVIDIUS™ NEURAL COMPUTE STICK 2

INTEL® MOVIDIUS™ NEURAL COMPUTE STICK 2

REDEFINING THE AI DEVELOPER KIT

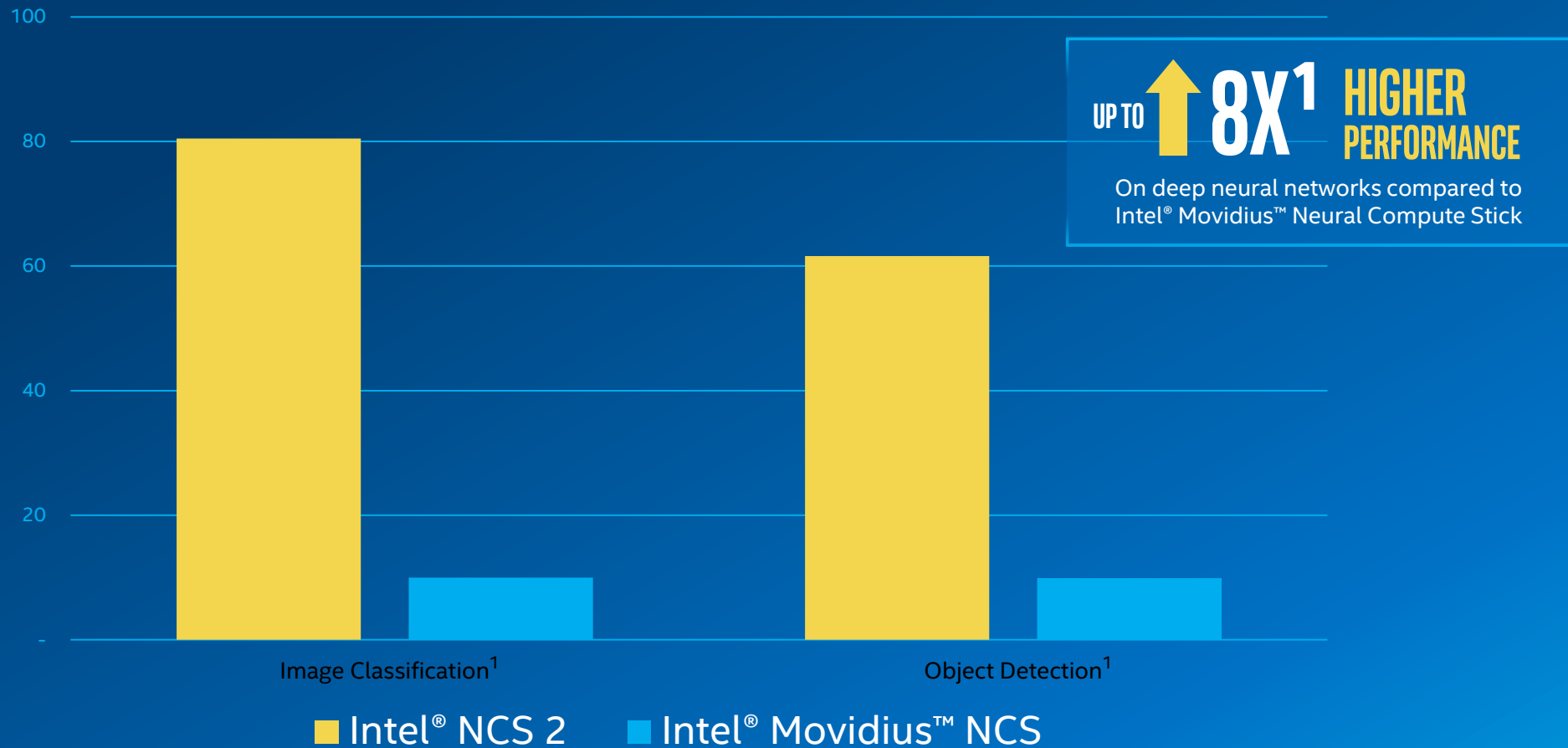
Technical Specifications

- Processor: Intel® Movidius™ Myriad™ X Vision Processing Unit (VPU)
- Supported frameworks: TensorFlow* and Caffe*
- Connectivity: USB 3.0 Type-A
- Dimensions: 2.85 in. x 1.06 in. x 0.55 in. (72.5 mm x 27 mm x 14 mm)
- Operating temperature: 0° C to 40° C
- Compatible operating systems: Ubuntu* 16.04.3 LTS (64 bit), CentOS* 7.4 (64 bit), and Windows® 10 (64 bit)



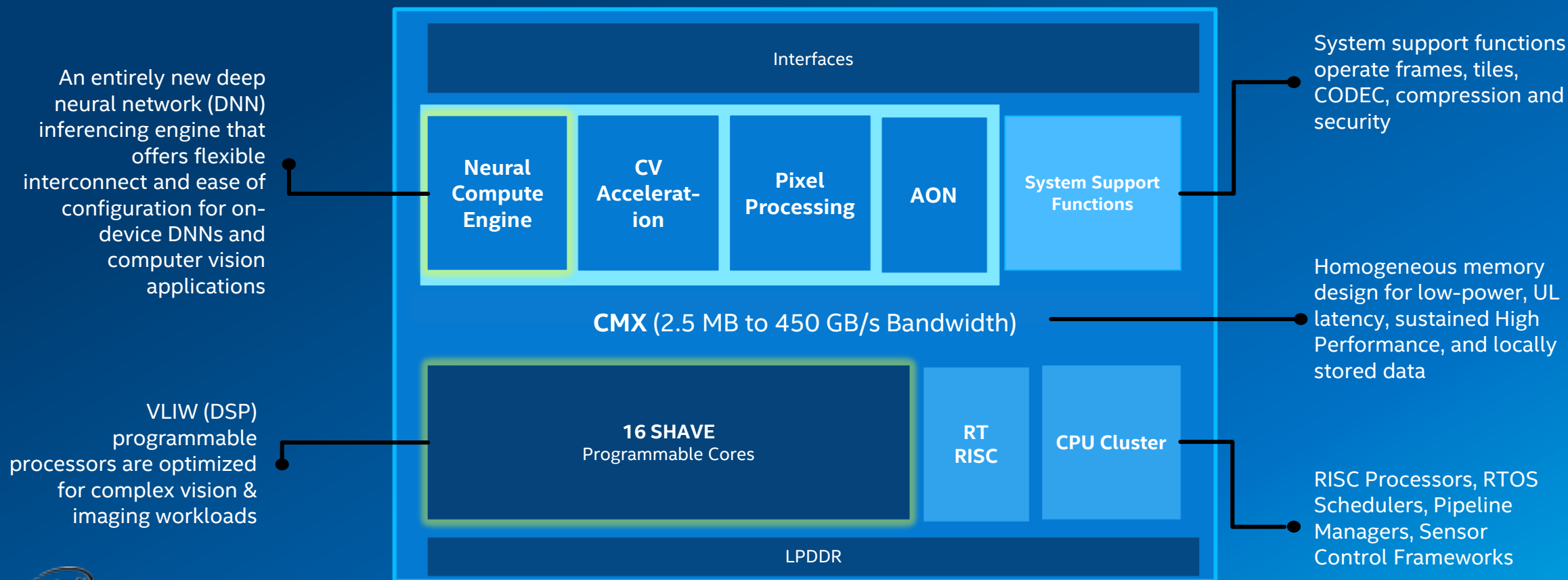
INTEL® NEURAL COMPUTE STICK 2

MORE CORES. MORE AI INFERENCE



INTEL® NEURAL COMPUTE STICK 2: FEATURING THE INTEL® MOVIDIUS™ MYRIAD™ X VPU

A self-sufficient, all-in-one processor that features the powerful **Neural Compute Engine** and **16 programmable SHAVE cores** that deliver class-leading performance for deep neural network inference applications.



INTEL[®] DISTRIBUTION OF OPENVINO[™] TOOLKIT

SUPPORTED LAYERS

View Documentation ► https://docs.openvino toolkit.org/latest/docs_IE_DG_supported_plugins_Supported_Devices.html

- Activation-Clamp
- Activation-ELU
- Activation-Leaky ReLU
- Active-PReLU
- Activation-ReLU
- **Activation-ReLU6**
- Activation-Sigmoid/Logistic
- Activation-TanH
- **ArgMax**
- BatchNormalization
- Concat
- **Const**
- Convolution-Dilated
- Convolution-Grouped
- Convolution-Ordinary
- Crop
- CTCGreedyDecoder*
- Deconvolution
- DetectionOutput*
- Eltwise-Max
- Eltwise-Mul
- Eltwise-Sum
- Flatten
- FullyConnected (Inner Product)
- GRN
- Interp
- LRN (Norm)
- MVN*
- Normalize*
- **Pad***
- Permute
- Pooling(AVG,MAX)*
- Power
- PriorBox
- PriorBoxClustered
- Proposal
- PSROIPooling
- **RegionYolo**
- ReorgYolo
- **Resample**
- Reshape
- RNN
- **ROIPooling**
- ScaleShift*
- Slice
- SoftMax
- Split
- Tile

* Support is limited to the specific parameters. Refer to "Known Layers Limitation" section for the device [from the list of support](#) 38

Changed since last update



LAB4 - HW ACCELERATION WITH INTEL® MOVIDIUS™ NEURAL COMPUTE STICK

URL: <https://github.com/intel-iot-devkit/smart-video-workshop/blob/master/HW-Acceleration-with-Movidious-NCS/README.md>

Objective: This lab shows how the Intel® Distribution of OpenVINO™ toolkit provides hardware abstraction to run the sample object detection application which was built in previous modules on Intel® Movidius™ Neural Compute Stick.

Estimated Complete Time: 20min



