# ACCELERATE AI INFERENCE FROM CLOUD TO EDGE

## WITH ONNX* RUNTIME + OPENVINO™ TOOLKIT

DATE: JUNE, 2020

# NOTICES AND DISCLAIMER

INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH INTEL PRODUCTS. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO SALE AND/OR USE OF INTEL PRODUCTS INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

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.

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 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.

# ONNX\* EXCHANGE

## What it is

ONNX is an open ecosystem that empowers AI developers to choose the right tools as their project evolves. ONNX provides an open source format for AI models, both deep learning and traditional ML. It defines an extensible computation graph model, as well as definitions of built-in operators and standard data types. ONNX is currently focused on the capabilities needed for inferencing (scoring).

## Target audience

- Computer vision, machine learning and deep learning software developers
- Data scientists
- OEMs, ISVs, System Integrators

## Usages

AI for robotics, retail, healthcare, security surveillance, office automation, transportation, non-vision use cases (speech, NLP, Audio, text) & more.

**AI FRAMEWORK INTEROPERABILITY – COMMON FORMAT**

**TOOLS TO CONVERT MODEL FORMATS TO ONNX**

**MODEL CATALOG THROUGH ONNX MODEL ZOO**

**STREAMLINING PATH FROM PROTOTYPE TO PRODUCTION**

**Homepage** ▶ onnx.ai
**Github** ▶ github.com/onnx/onnx
**ONNX Model Zoo Github** ▶ https://github.com/onnx/models
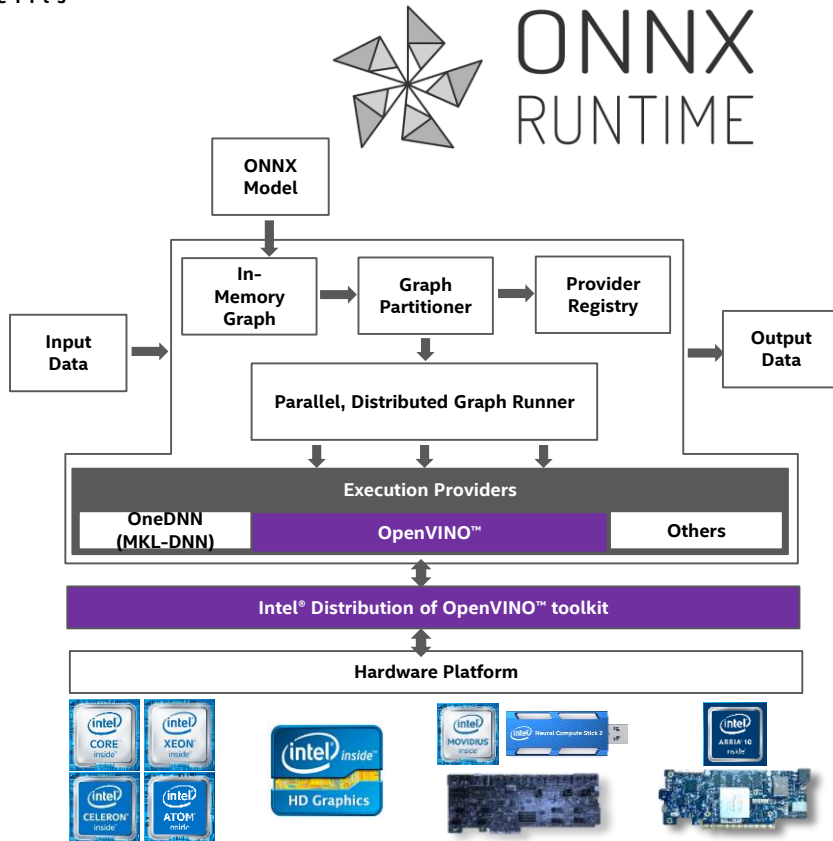
# ONNX* RUNTIME

## What it is

ONNX Runtime is a performance-focused complete scoring engine for Open Neural Network Exchange (ONNX) models, with an open extensible architecture to continually address the latest developments in AI and Deep Learning. ONNX Runtime stays up to date with the ONNX standard and supports all operators from the ONNX v1.2+ spec with both forwards and backwards compatibility. Execution Provider plugin allows the support of ONNX RT for Intel® Distribution of OpenVINO™ toolkit.

✓ **BUILT SPECIFICALLY FOR ONNX FORMAT MODELS**

▢ **SUPPORTS EXECUTION ON MANY TYPES OF HARDWARE**

◎ **COMPLETELY OPEN SOURCED ON GITHUB**



**Homepage ▶ onnx.ai**

**Github ▶ github.com/Microsoft/onnxruntime**

4

# INTEL® DISTRIBUTION OF OPENVINO™ TOOLKIT

## What it is

A toolkit to accelerate development of **high performance computer vision** & **deep learning inference into vision/AI applications** used from edge to cloud. It enables deep learning on hardware accelerators and easy deployment across multiple types of Intel® platforms.

## Target audience

- Computer vision, deep learning software developers
- Data scientists
- OEMs, ISVs, System Integrators

## Usages

AI for robotics, retail, healthcare, security surveillance, office automation, transportation, non-vision use cases (speech, NLP, Audio, text) & more.

OpenVINO™

HIGH PERFORMANCE, PERFORM AI AT THE EDGE

STREAMLINED & OPTIMIZED DEEP LEARNING INFERENCE

HETEROGENEOUS, CROSS-PLATFORM FLEXIBILITY

**Free Download** ▶ software.intel.com/openvino-toolkit

**Open Source version** ▶ 01.org/openvinotoolkit

# INTEL® DISTRIBUTION OF OPENVINO™ TOOLKIT

## Deep Learning

### Intel® Deep Learning Deployment Toolkit

**Model Optimizer**
Convert & Optimize

IR

**Inference Engine**
Optimized Inference

IR = Intermediate Representation file

### Open Model Zoo

**40+ Pretrained Models**

**Samples**

**Model Downloader**

### Deep Learning Workbench

| Calibration Tool | Model Analyzer | Benchmark App | Accuracy Checker | Aux. Capabilities |

## Traditional Computer Vision

### Optimized Libraries & Code Samples

**OpenCV***

**OpenVX***

**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 XEON inside | intel CORE inside | intel CELERON inside | intel ATOM inside | intel ARRIA 10 inside | intel MOVIDIUS inside

intel IRIS Pro GRAPHICS

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

An open source version is available at 01.org/openvinotoolkit (deep learning functions support for Intel CPU/GPU/NCS/GNA).
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

6

# AT THE EDGE

# ONNX* ECOSYSTEMS

| Frameworks | Caffe2  Chainer  mxnet  ML.NET  PaddlePaddle  PyTorch  MATLAB  Microsoft Cognitive Toolkit |
|---|---|
| Converters | ML  LibSVM  Keras  TensorFlow  scikit learn  dmlc XGBoost  XG |
| Runtimes | ONNX RUNTIME  OpenVINO™  And others… |
| Visualizers | NETRON  Visual DL |

# GET STARTED

# Accelerate Time to Production with Intel® DevCloud for the Edge

## See immediate AI Performance Across Intel's Array of Edge Solutions



**Instant, Global Access**
Run AI applications from anywhere in the world

**Prototype on the Latest Hardware and Software**
Develop knowing you're using the latest Intel technology

**Benchmark your Customized AI Application**
Immediate feedback – frames per second, performance

**Reduce Development Time and Cost**
Quickly find the right compute for your edge solution

**Sign up now for access**

# 3 SETUP OPTIONS IN GITHUB

## BUILD FROM SOURCE

- ➤ (ONNX RT + OV) *
- ➤ RUN NATIVELY, COMPILE FROM SCRATCH – PROVIDES MAXIMUM FLEXIBILITY
- ➤ DEPLOY NATIVELY AT THE EDGE
- ➤ **Github Readme** ▶
  https://github.com/microsoft/onnxruntime/blob/master/docs/execution_providers/OpenVINO-ExecutionProvider.md

## DEPLOY CONTAINERS FROM AZURE IOT EDGE

- ➤ [ONNX RT + OV + AZURE IOT EDGE] *
- ➤ PROVIDES FLEXIBILITY AS WELL AS CONVENIENCE THROUGH CONTAINER SUPPORT
- ➤ DEPLOY CUSTOM APPLICATIONS IN CONTAINERS FROM AZURE IOT EDGE
- ➤ **Github Azure IoT Hub Instructions** ▶
  https://github.com/intel/Edge-Analytics-FaaS/tree/master/Azure-IoT-Edge/OnnxRuntime

## DEPLOY CONTAINERS FROM AZURE ML

- ➤ [ONNX RT + OV + AZURE IOT EDGE]
- ➤ MORE AUTOMATED, AZURE ML CONSTRUCTS THE CONTAINER FROM PRE-DEFINED AZURE ML FORMAT APPLICATIONS
- ➤ DEPLOY AZURE ML APPLICATIONS IN CONTAINERS FROM AZURE ML SERVICES
- ➤ **Github Azure ML Container Dockerfiles** ▶
  https://github.com/microsoft/onnxruntime/tree/master/dockerfiles

---

*Note: Download the Intel® Distribution of OpenVINO™ toolkit installer(tgz) before building the above Docker image.

Additional Github Resource: Azure ML Instructions ▶

**Cloud to Edge Deployment flow using Azure ML and Azure IoT Edge**

*Using Azure ML to deploy Azure ML container applications*

https://github.com/Azure-Samples/onnxruntime-iot-edge/tree/master/AzureML-OpenVINO

# HOW IT WORKS (RUNTIME)

```
import onnxruntime                    Simple runtime call pointing to model location


session = onnxruntime.InferenceSession("model.onnx")

x = GetInputData()

y = session.run([session.get_outputs()[0].name],

        {session.get_inputs()[0].name : x})
```

# CSP MARKETPLACE OFFERS





- [Intel DevCloud](#) link on Azure Marketplace

- Use Intel maintained HW Devices to quickly deploy models

- Understand HW device performance trade-off before you purchase your device(s).

- [Vision Ready-to-deploy app](#) link Azure Marketplace

- Leverages customvision.ai and OpenVINO with enhanced DX Training-to-Inference

- Bring your own Intel device(s) to experience the Cloud to Edge inference and quick relaunch of vision application

# ONNX* MODEL ZOO



> ➤ **Github** ▶ https://github.com/onnx/models

Sources: Microsoft

# ONNX* TUTORIALS

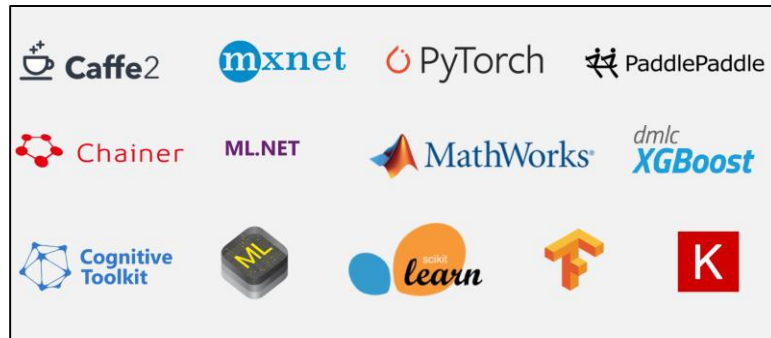## Get started with ONNX and tutorials

Docker image for ONNX and Caffe2/PyTorch
Docker image for ONNX, ONNX Runtime, and various converters

- **Getting ONNX models –** ONNX Model Zoo

- **Services –** Output ONNX models customized for your data
    Azure Custom Vision service
    Azure Machine Learning automated ML

- **Converting to ONNX format**

- **Scoring ONNX models –** Score accuracy



Sources: Microsoft

# FEATURES SET

- **Compute and Accelerator support:**
  - **Intel® CPU, integrated GPU**
  - **Intel® Movidius™ Myriad™ X VPU (USB and embedded)**
  - **Intel® Vision Accelerator Design Products with Intel® Movidius™ Myriad™ X VPU (2x, 4x & 8x)**
  - **Intel® Vision Accelerator Design Products with Intel® Arria® 10 FPGA**

- **Quantization support: Full precision (32 bit) and Half precision (16 bit) floating point**

- **Operator coverage: Majority models from ONNX Model Zoo github.com/onnx/models**

- **OS Support: Linux* and Win10***

- **Docker container support: Linux* only**

- **Azure ML integration: Train model on Azure* ML and deploy on connected edge devices**

# DEVELOPER KITS, USE CASES, & CASE STUDIES

# DEVELOPER KITS

## INTEL® NEURAL COMPUTE STICK 2
*link*

Powered by the
Intel® Movidius™ Myriad™ X VPU

## IEI TANK* AIOT DEVELOPER KIT
*link*

**Intel® Vision Accelerator
Design Product Choices**

Powered by Intel® Movidius™ VPU (link)

Powered by Intel® Arria® 10 FPGA  (link)
*In Preview*

## UP SQUARED* AI VISION X DEV KIT
*link*

**Intel® Vision Accelerator
Design Product**

Powered by Intel® Movidius™ VPU (link)

# EQUIPMENT MAKER OFFERS

- **IEI\* TANK AIoT Developer Kit**
  **Intel® Core® i7/i5/i3 Processor & Intel® Xeon® Processor**
  **Use Case: Industrial**
- **IEI\* FLEX-BX200**
  **Intel® Core® i3/5/i7 Processor**
  **Use Cases: Public Safety, Parking Mgmt., License Plate Detection**

- **UP\* Squared; UP\* Core Plus**
  **Intel® Atom™ Processor ; Intel® Core® i7/i5/i3 Processor**
  **Use Cases: Retail, DSS**
- **Aaeon\* BOXER-6841M**
  **Intel® Core® i7/i5/i3 Processor**
  **Use Cases: Industrial, Smart Retail and Smart City**

- **Advantech\* ARK-1124 + VEGA-320**
  **Intel® Atom™ Processor + Intel® Movidius™ Myriad™ X VPU**
  **Use Cases: Age & Gender Recognition**

# SUPPORT & RESOURCES

intel

# SUPPORT

**Software Issues**
Software issues related to ONNX Runtime with OpenVINO Execution Provider code should be logged at: "Issues"
Tab https://github.com/Microsoft/onnxruntime with [OpenVINO-EP] tag.

**Hardware Issues**
Hardware issues should be routed towards your equipment maker suppliers, your Intel Representative, or Intel
Premier Support

**Supported Models**
Link to supported models for the ONNX Runtime with OpenVINO Execution Provider
https://github.com/microsoft/onnxruntime/blob/master/docs/execution_providers/OpenVINO-
ExecutionProvider.md.
All issues related to these models should be routed towards your Intel Representative

**Intel® Distribution of OpenVINO™ toolkit Support**
OpenVINO issues should be reported through the OpenVINO "Computer Vision" Forum
https://software.intel.com/en-us/forums/computer-vision

**ONNX Support**
All other ONNX model issues should be logged at "Issues" tab https://github.com/Microsoft/onnxruntime

# INTEL® IOT RFP READY KITS

**Check RFP Ready Kit [Playbook] for details on each kit**

| Retail | Industrial | Transportation | Smart Cities | Healthcare | Horizontal |
|--------|-----------|----------------|--------------|------------|------------|

**Retail**
→ Kiosk & digital signage
→ POS & mobile POS
→ Inventory Management

**Industrial**
→ Manufacturing
→ Building Management
→ Agriculture
→ Energy

**Transportation**
→ Fleet Management
→ Logistics

**Smart Cities**
→ Security surveillance
→ Smart lighting
→ Connected Transportation
→ Air quality management

**Healthcare**
→ Medical (in-hospital)
→ Remote health management

**Horizontal**
→ Security Surveillance Video
→ Connectivity

# ACCELERATE PROTOTYPE TO PRODUCTION & SOLUTION DEPLOYMENT

## SCALE

**Deploy solution & solve business problems, and scale with Intel® IoT Solution Aggregators & Ecosystem**

## DEVELOP USE CASE SPECIFIC OFFERS

**Developer optimization & use case specific applications**

**Intel® RFP Ready Kits**

## INCREASE PERFORMANCE

**Intel® Vision Accelerator Design Products**

Intel® Movidius™ Myriad™ X VPU
Intel® Arria® 10 FPGA

## USE VISION ACCELERATOR KITS

**Intel® Distribution of OpenVINO toolkit**

AAEON UP Squared AI Vision Developer Kit

IEI Tank AIoT Developer Kit

## DEVELOP ON HOST SYSTEM

& TEST on

INTEL® DEVCLOUD FOR THE EDGE

OpenVINO™

**Intel® AI: In Production** ▶
https://software.intel.com/ai-in-production