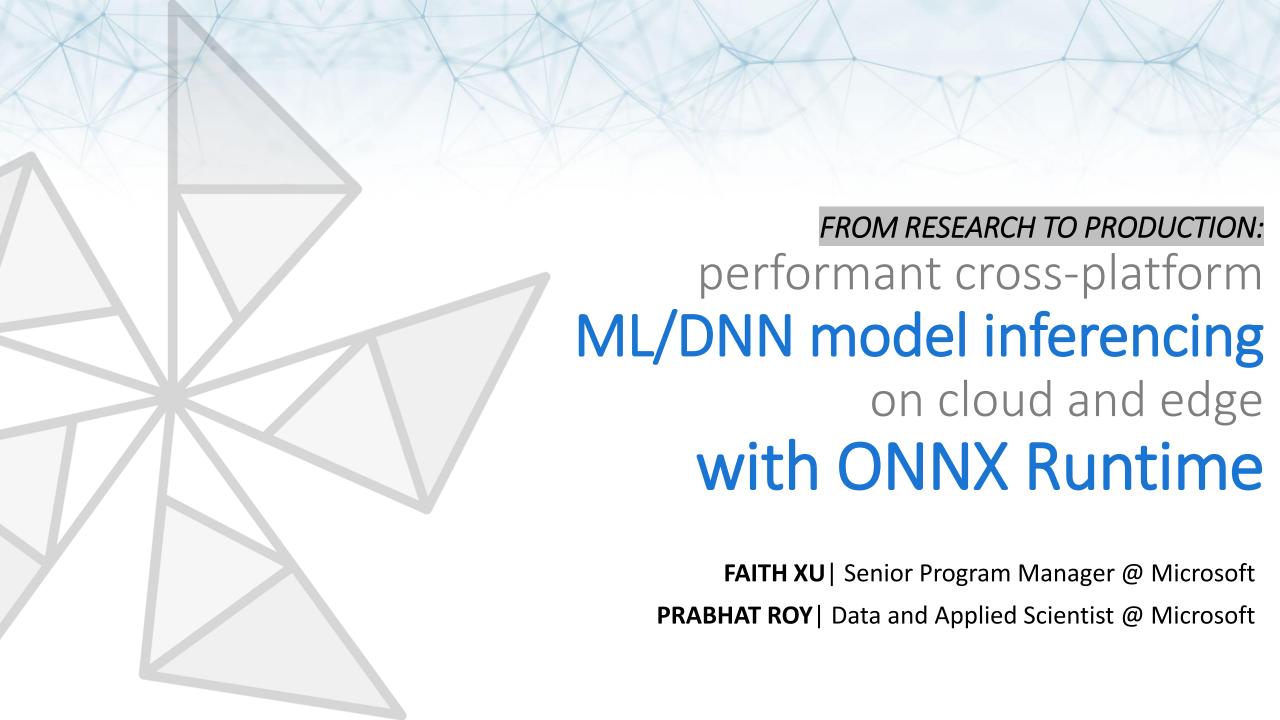
aka.ms/odsc-onnx

Please get started with the pre-requisite

OPEN DATA SCIENCE CONFERENCE

Virtual | Apr. 14 - Apr. 17 2020





Agenda – What we'll cover today

INTRODUCTION TO ONNX

 ACTIVITY: Convert the huggingface/transformers BERT model (trained with PyTorch) to ONNX format for accelerated inferencing

Optional extension: Deploy to a hosted web service using Azure Machine Learning Services

Trends and Growth Areas

Research -> Industry

- Automated Machine Learning services
- Startups applied AI
- Hosted services for cloud compute
- Hardware investments

Connectivity, compute, and resources

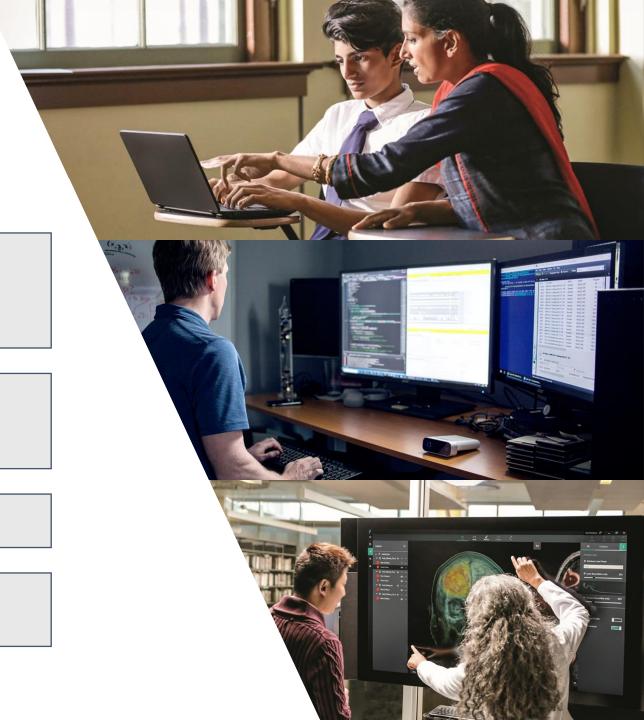
- Infinite storage and compute in the cloud
- CPU, GPUs for training
- LOTS of data

Application spans across all industries

• Healthcare, finance, farming, manufacturing, consumer products, and more

Investments in AI education and jobs

- Universities
- ML Engineer



Al everywhere

Microsoft 365















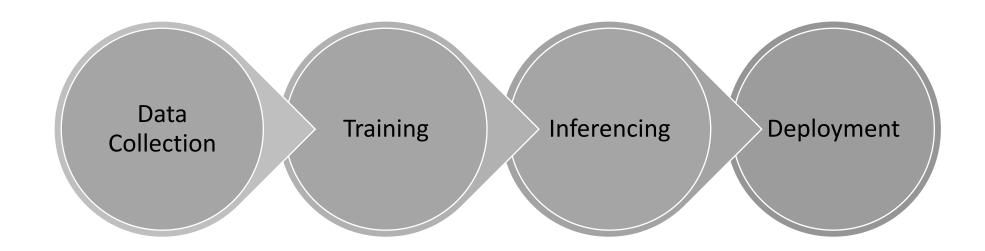




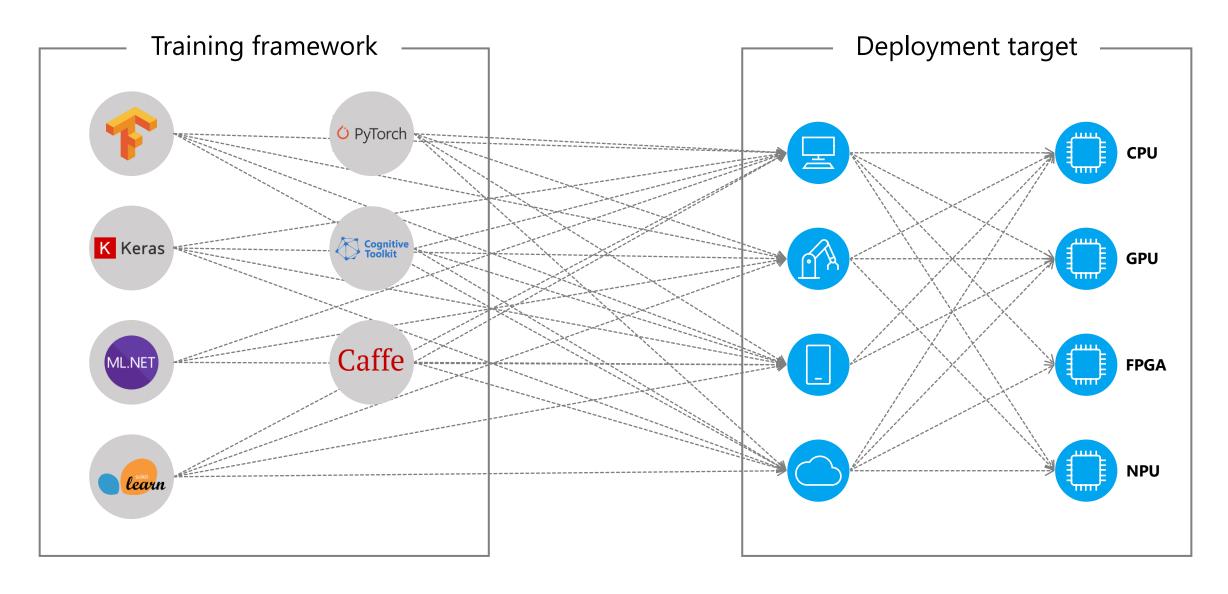




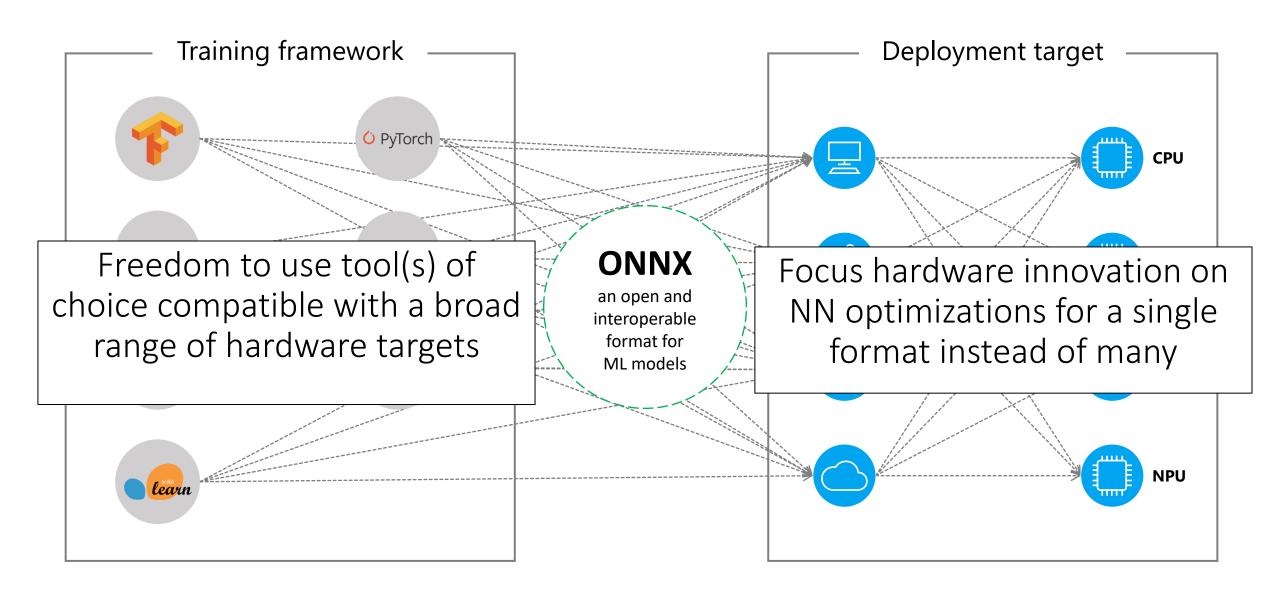
ML Models: Research to Production

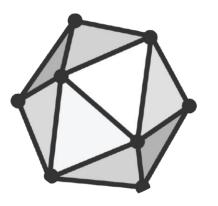


Training frameworks **x** Deployment targets



ONNX: an open and interoperable format for ML models





ONNX

github.com/onnx onnx.ai

OPEN NEURAL NETWORK EXCHANGE











































Neural Network Libraries



















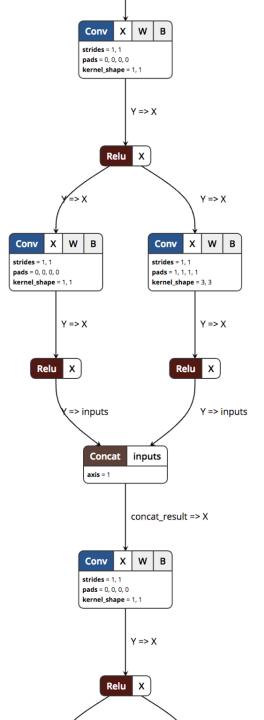




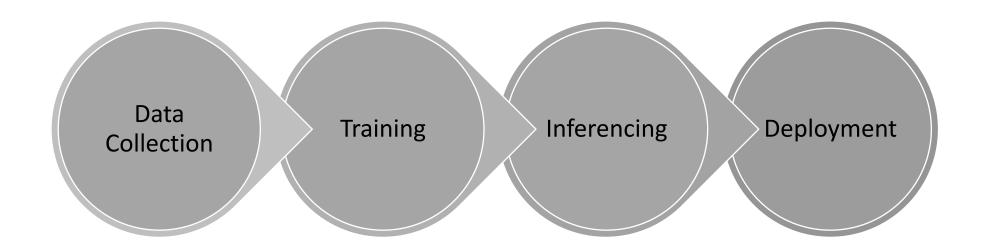


ONNX

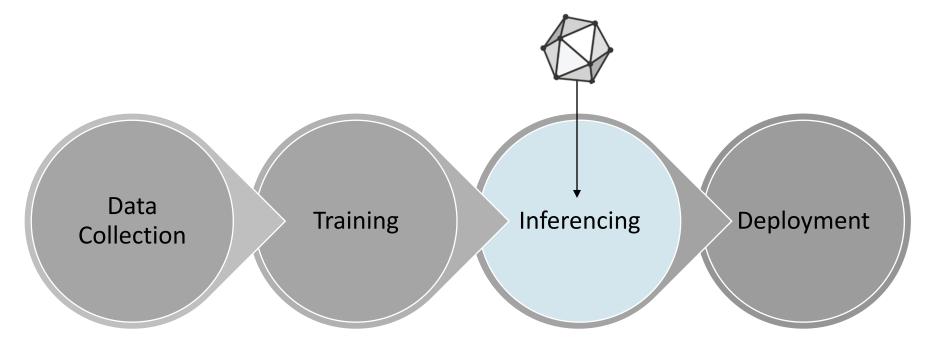
- Standard format for ML models consisting of:
 - common Intermediate Representation
 - full operator spec
- Model = graph composed of computational nodes
- Supports both DNN and traditional ML



Where does ONNX fit in?



Focus on inferencing



Framework Compatibility



















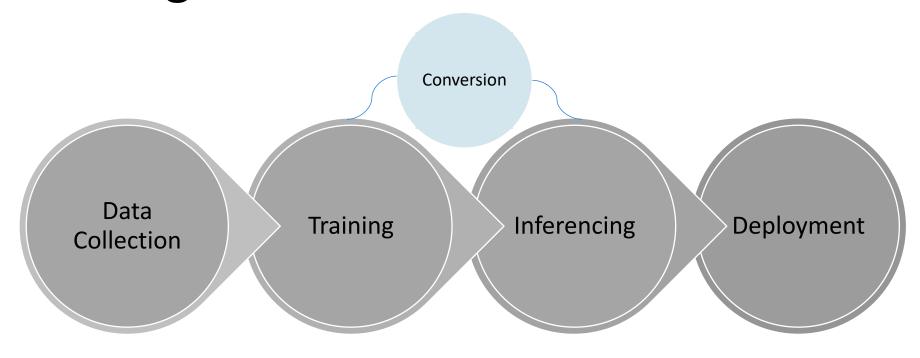








How do I get an ONNX model?



How do I get an ONNX model?

• Get a pre-trained ready to use model from the ONNX Model Zoo

Use a model builder service that supports export to the ONNX format

Convert an existing model from another framework

Open Source converters for popular ML frameworks

```
Tensorflow: onnx/tensorflow-onnx
PyTorch (native export)
Keras: onnx/keras-onnx
Scikit-learn: onnx/sklearn-onnx
CoreML: onnx/onnxmltools
LightGBM: onnx/onnxmltools
LibSVM: onnx/onnxmltools
XGBoost: onnx/onnxmltools
H2O: onnx/onnxmltools
ML.NET (native export)
SparkML (experimental): onnx/onnxmltools
CNTK (native export)
```

Examples: Model Conversion

```
from keras.models import load_model
import keras2onnx
import onnx

keras_model = load_model("model.h5")

onnx_model = keras2onnx.convert_keras(keras_model, keras_model.name)

onnx.save_model(onnx_model, 'model.onnx')
```

```
import torch
import torch.onnx

O PyTorch

model = torch.load("model.pt")

sample_input = torch.randn(1, 3, 224, 224)

torch.onnx.export(model, sample_input, "model.onnx")
```

```
import numpy as np
import chainer
from chainer import serializers
import onnx_chainer

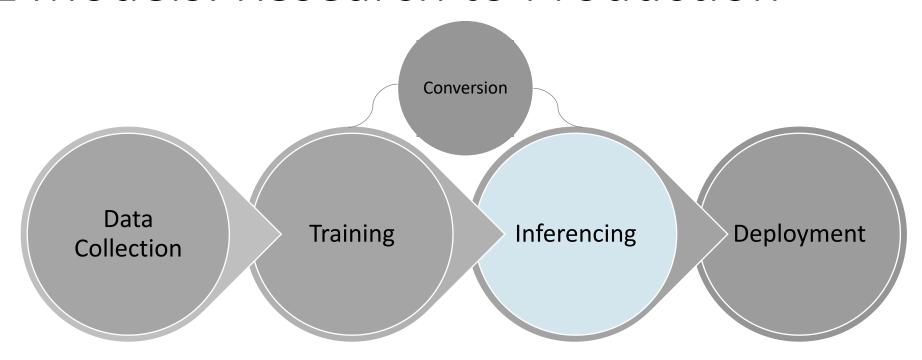
serializers.load_npz("my.model", model)

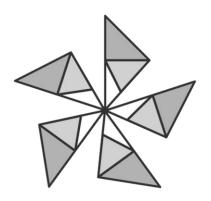
sample_input = np.zeros((1, 3, 224, 224), dtype=np.float32)
chainer.config.train = False

onnx_chainer.export(model, sample_input, filename="my.onnx")
```

Inferencing ONNX models

ML Models: Research to Production





ONNX Runtime

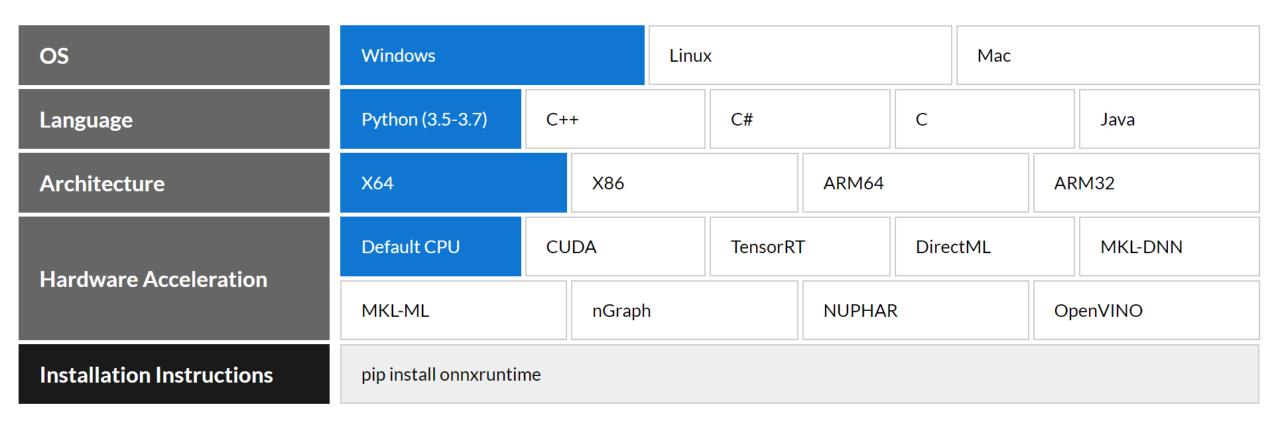
aka.ms/onnxruntime

github.com/microsoft/onnxruntime

ONNX Runtime: open source high performance Inference Engine

- Performance focused design
- Full ONNX operator support
- Flexibility for custom operators not in the spec
- Backwards compatible to minimize versioning issues with software or model updates

Cross platform, multi language API



Inferencing with ONNX Runtime

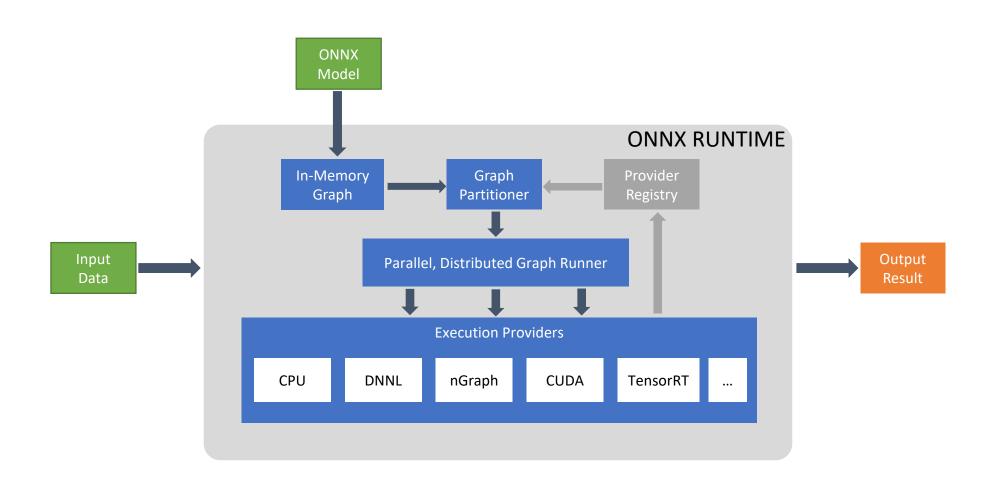
```
import onnxruntime
session = onnxruntime.InferenceSession("mymodel.onnx")
results = session.run([], {"input": input_data})
```

```
using Microsoft.ML.OnnxRuntime;

var session = new InferenceSession("model.onnx");

var results = session.Run(input);
```

Leverages and abstracts hardware accelerators



W Conv strides = 1, 1 pads = 0, 0, 0, 0kernel_shape = 1, 1 Y => XRelu X Y => XW Conv Conv strides = 1, 1 strides = 1.1 pads = 0, 0, 0, 0pads = 1, 1, 1, 1 kernel_shape = 1, 1 kernel_shape = 3, 3 Y => XY => XRelu X Relu Y => inputs => inputs inputs Concat axis = 1 concat_result => X W Conv kernel_shape = 1, 1

Graph Optimizations

- Constant folding
- Node eliminations
- Simple and complex node fusions
- Layout optimizations (e.g. NCHWc vs NCHW)
- Extendible and pluggable to add new optimizations

Accelerators for a range of hardware

Base CPU

Microsoft Linear

Algebra Subprograms

NVIDIA CUDA

NVIDIA TensorRT

DirectML (Windows)

preview

Intel nGraph

Intel DNNL

Intel OpenVINO

NUPHAR
TVM/LLVM-based
model compiler

NN API (Android)

preview

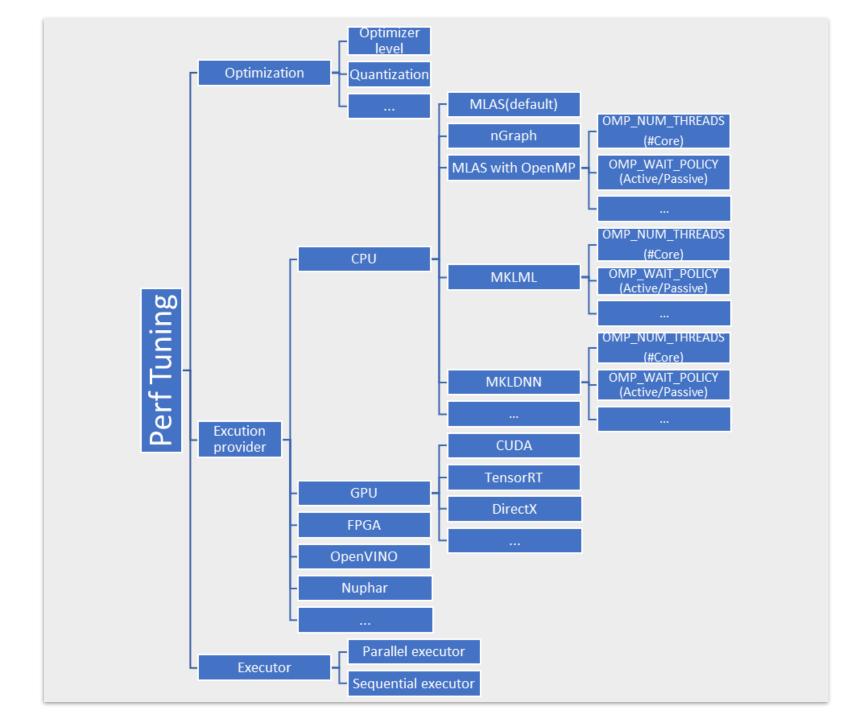
NXP ARM Compute Library preview

Perf Tuning in ONNX Runtime

ONNX Go Live (OLive)

https://github.com/microsoft/OLive

- Automates the process of ONNX model shipping by integrating model conversion, correctness test, and performance tuning into a single pipeline
- Outputs a production ready ONNX model with ONNX Runtime configurations (execution provider + optimization options)

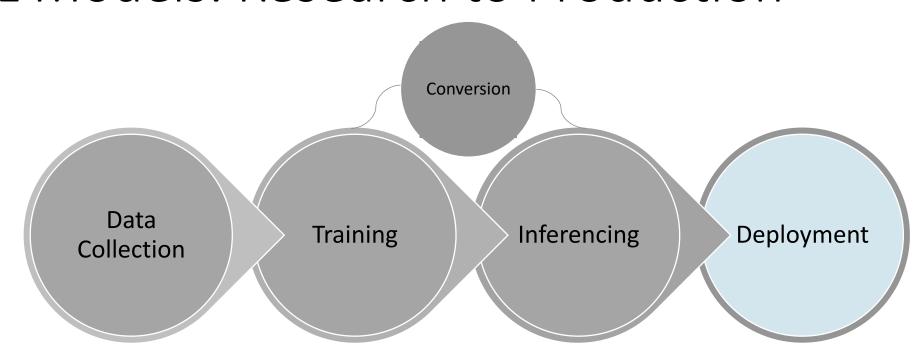


Focus areas and further investments

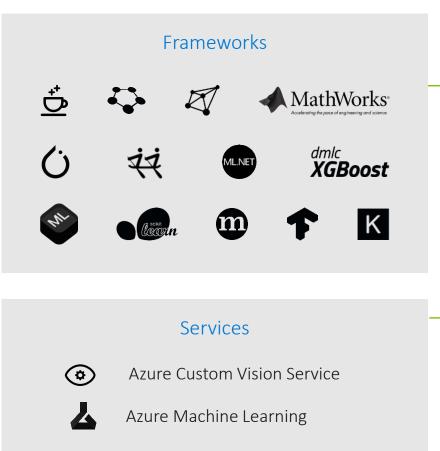
- ONNX spec alignment
- Performance optimizations
- Hardware ecosystem and mobile/edge/IoT devices
- Production models
- Training acceleration

Deploying models with ONNX Runtime

ML Models: Research to Production







DEPLOY



Azure Machine Learning service

Ubuntu VM

Windows Server 2019 VM

Devices

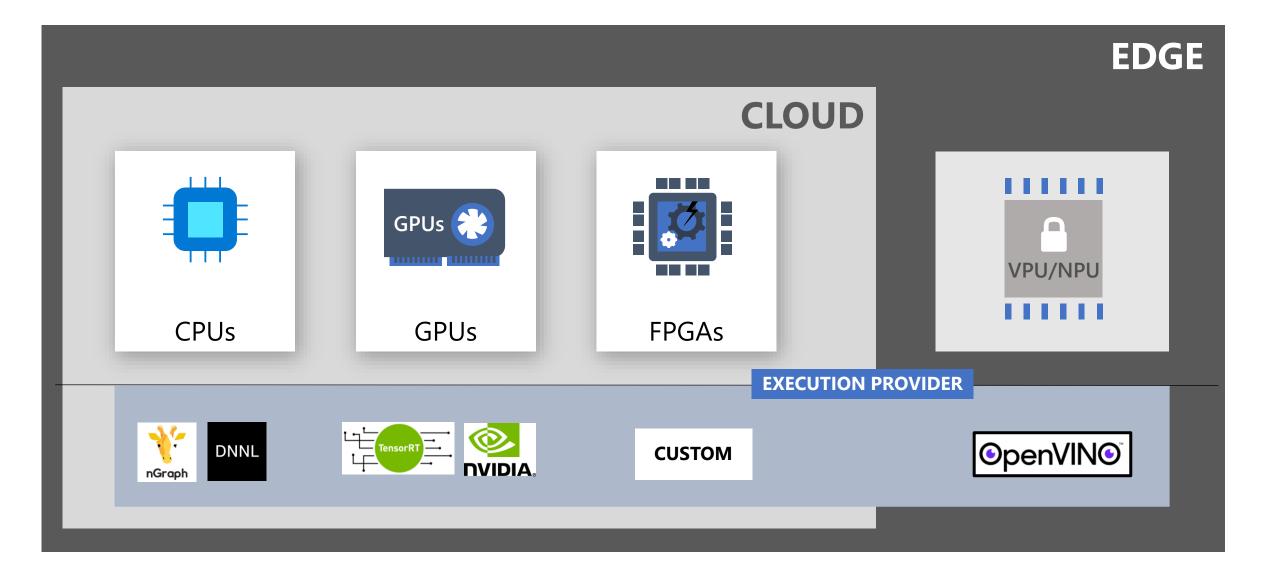
ONNX Model

Local applications

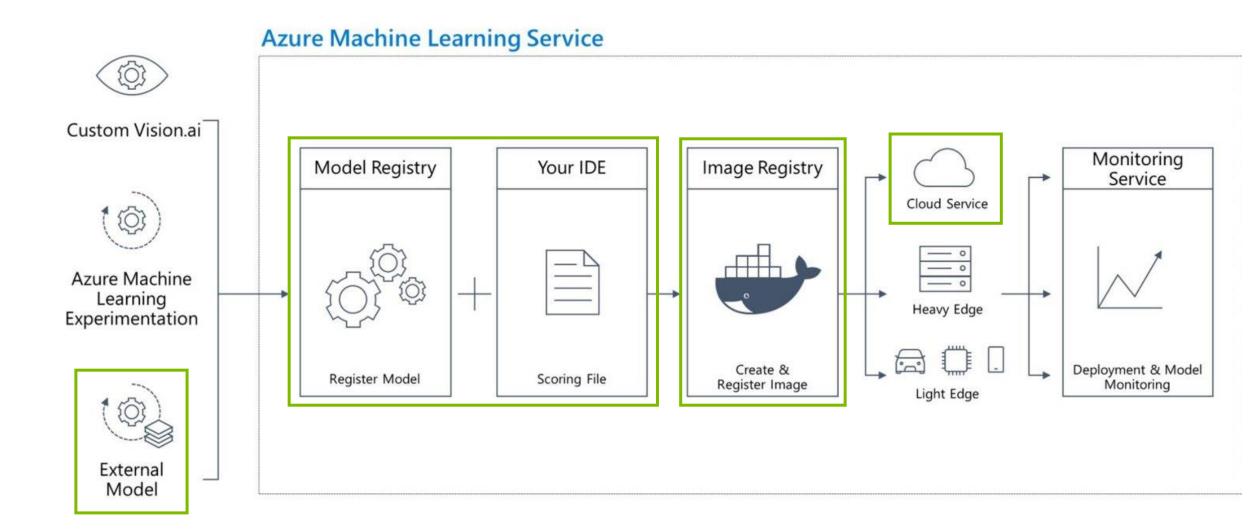
Edge Cloud & Appliances

Edge & IoT Devices

Variety of Deployment Options



Deployment through Azure ML



ACTIVITY

Convert the huggingface/transformers BERT model (trained with PyTorch) to ONNX format for accelerated inferencing

References

- ONNX: https://github.com/onnx/onnx
- ONNX Converters: https://github.com/onnx/onnxmltools/tree/master/onnxmltools
- ONNX Tutorials: https://github.com/onnx/tutorials
- ONNX Runtime: https://github.com/microsoft/onnxruntime
- ONNX Runtime Tutorials: https://github.com/microsoft/onnxruntime#examples-and-tutorials
- Performance Tuning with ONNX Runtime: https://github.com/microsoft/onnxruntime/blob/master/docs/ONNX Runtime Perf Tuning.md
- Training, Inferencing, and deployment in AzureML with ONNX models: https://aka.ms/onnxnotebooks
- AzureML resources: https://azure.microsoft.com/en-us/services/machine-learning/
- Deploying to Edge and IoT devices: <u>Deploying to Intel OpenVINO based devices</u>, <u>Deploying to NVIDIA Jetson</u> <u>Nano (ARM64)</u>
- Windows ML: https://docs.microsoft.com/en-us/windows/ai/windows-ml/

FAITH XU | faxu@microsoft.com PRABHAT ROY | Prabhat.Roy@microsoft.com