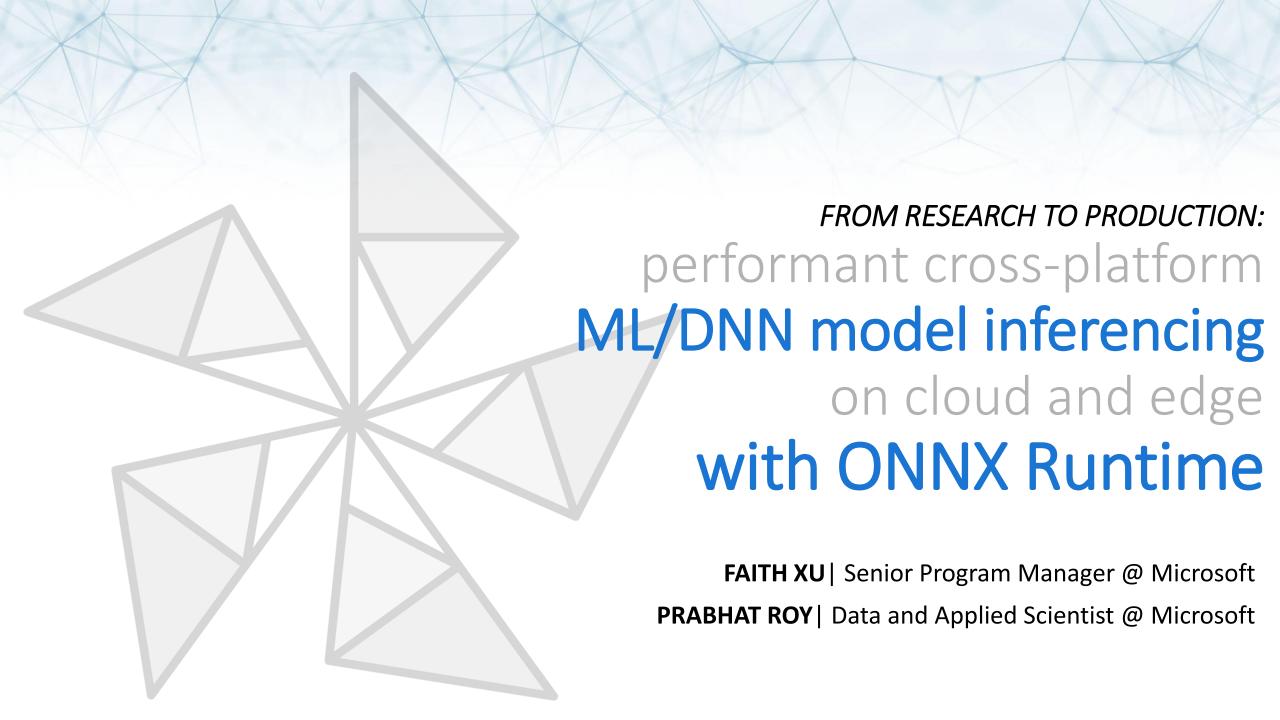
# OPEN DATA SCIENCE CONFERENCE

London | Nov. 19 - Nov. 22 2019





#### Agenda – What we'll cover today

- INTRODUCTION TO ONNX
- ACTIVITY A: Train an image classification model in PyTorch and convert to ONNX format for inferencing
- ACTIVITY B: Train a PyTorch model and deploy for production usage

### Why now?

#### Trends and Growth Areas

#### Research -> Industry

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

#### Connectivity, compute, and resources

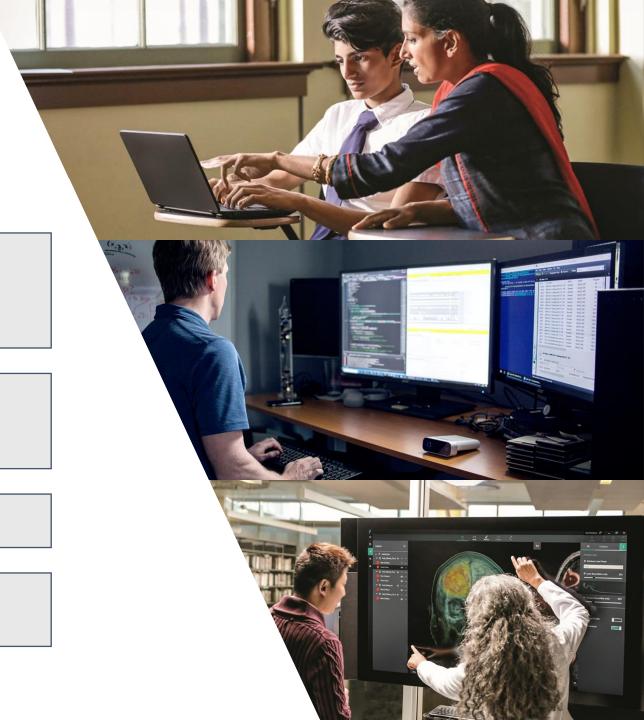
- Infinite storage and compute in the cloud
- CPU, GPUs for training
- LOTS of data

#### **Application spans across all industries**

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

#### Investments in AI education and jobs

- Universities
- ML Engineer



#### Product teams want to incorporate ML

Microsoft 365









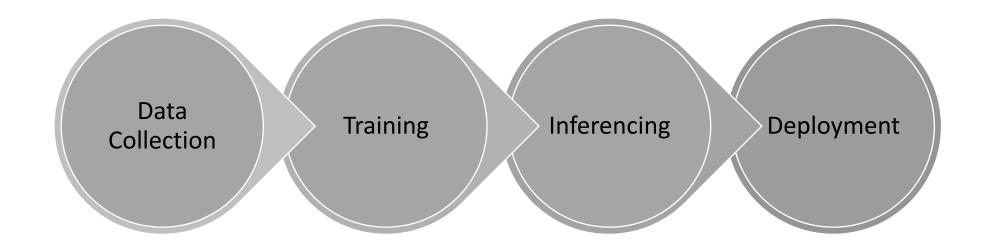




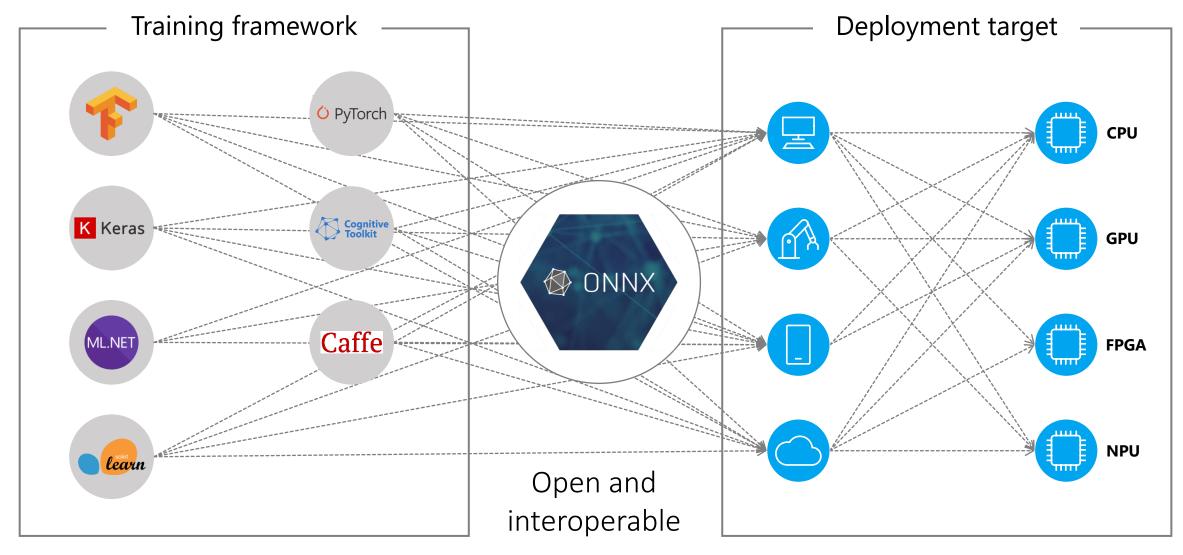


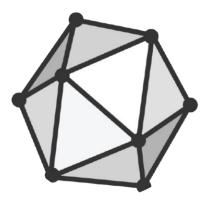


#### ML Models: Research to Production



#### Reality



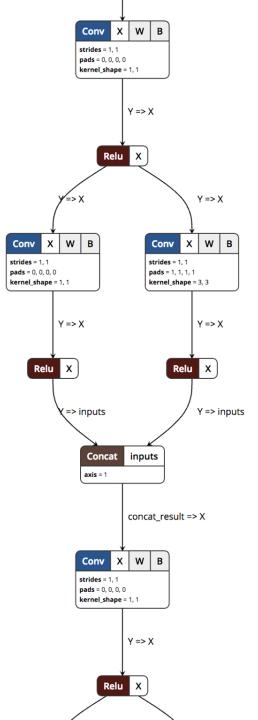


## ONNX https://github.com/onnx

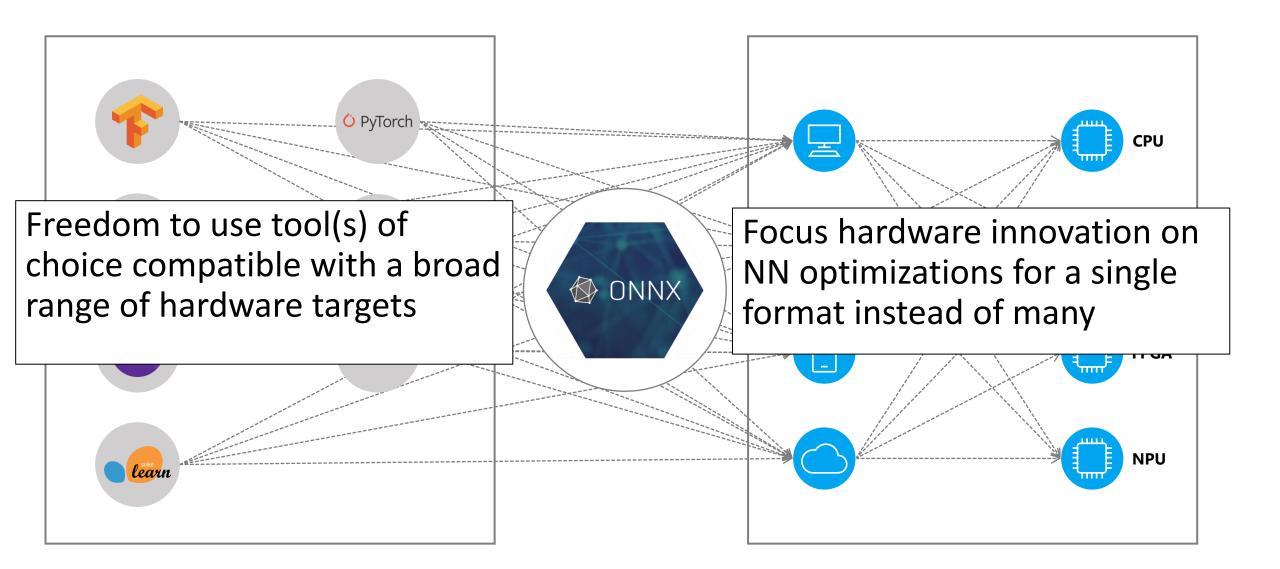
#### OPEN NEURAL NETWORK EXCHANGE

#### What is ONNX?

- Interoperable standard format for AI models consisting of:
  - common Intermediate Representation (IR)
  - full operator spec
- Model = graph composed of computational nodes, based on Google protobuf
- Graph = Compact and cross-platform representation for serialization
- Supports both DNN and traditional ML
- Backward compatible with comprehensive versioning



#### What does this provide?



#### Framework Compatibility



























#### **ONNX** Community











































Neural Network Libraries

























#### Open Governance



#### **Steering Committee**

<u>Prasanth</u> <u>Pulavarthi</u> (Microsoft)

Joe Spisak (Facebook)

Vin Sharma (Amazon)

Harry Kim (Intel)

Dilip Sequeira (NVIDIA)



#### **SIG** (special interest group)

#### **Architecture/Infrastructure**

<u>Lu Fang</u> (Facebook)

Ke Zhang (Microsoft)

#### **Operators**

Michał Karzyński (Intel)

**Emad Barsoum** (Microsoft)

#### **Converters**

Chin Huang (IBM)

**Guenther Schmuelling (Microsoft)** 

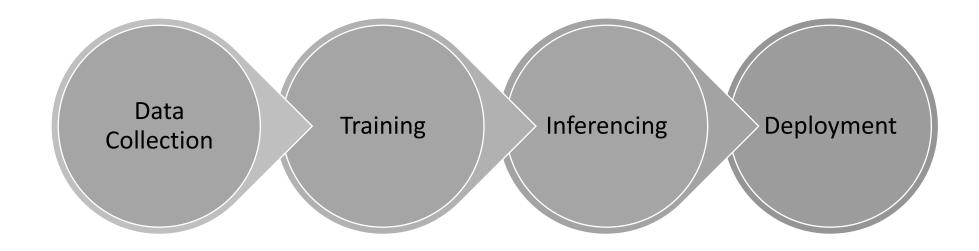


#### **Working Groups**

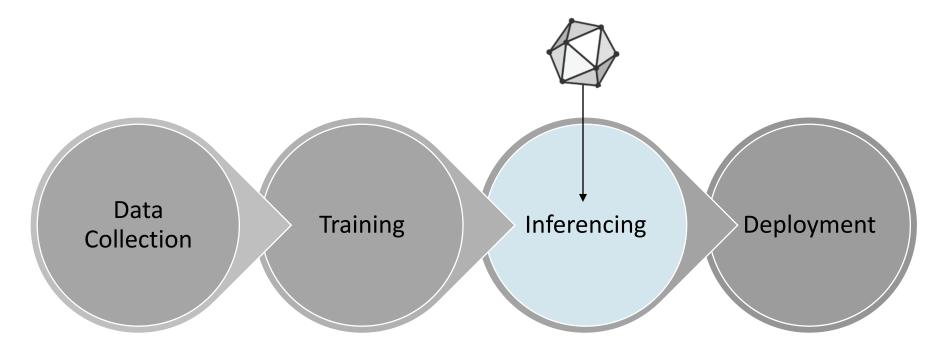
Training

Edge/Mobile

#### ML Models: Research to Production



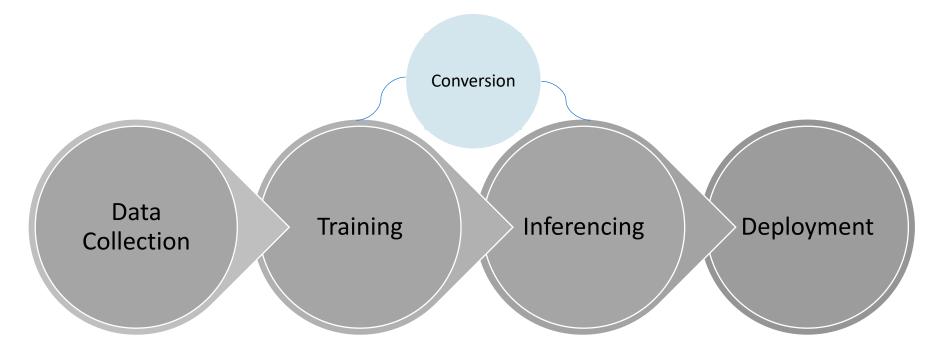
#### ML Models: Research to Production



#### How do I get an ONNX model?

- Get a pre-trained ready to use model from the <u>ONNX Model Zoo</u>
- Use a model builder service that supports export to the ONNX format
- Convert an existing model from another framework

#### ML Models: Research to Production



# Open Source converters for popular 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
SparkML (alpha): 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

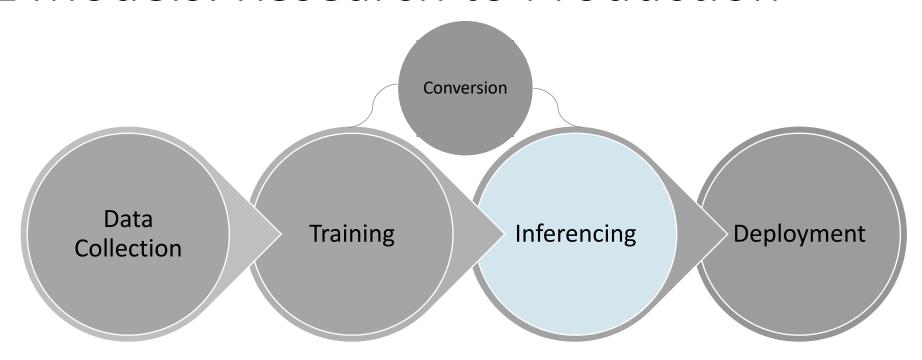
onnx_chainer.export(model, sample_input, filename="my.onnx")
```

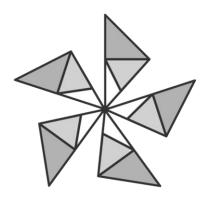
#### **ACTIVITY** A

Train an image classification model in PyTorch and convert to ONNX format for inferencing

### Inferencing ONNX models

#### ML Models: Research to Production





#### **ONNX** Runtime

aka.ms/onnxruntime

github.com/microsoft/onnxruntime

# ONNX Runtime is an open source high performance Inference Engine for ONNX models

## ONNX Runtime can run all operators defined in the ONNX spec

- ONNX domain (DNN) and ONNX-ML (traditional)
- Backwards and forwards compatible to minimize versioning issues with software or model upgrades
- Flexibility for custom operators not in the spec

#### Cross platform, multi language API

Windows, Linux, Mac X64, X86, ARM CPU, GPU Python, C, C++, C#, Ruby, Java (future)

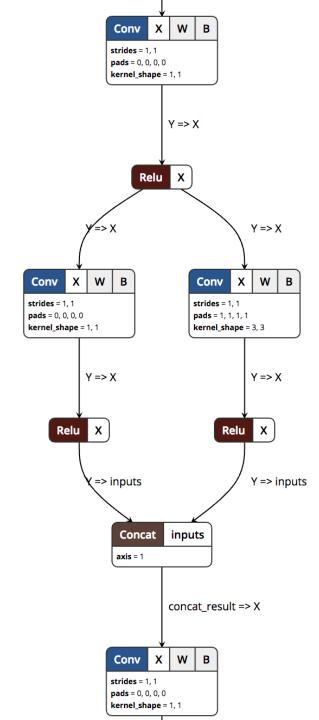
#### Inferencing with ONNX Runtime

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

```
using Microsoft.ML.OnnxRuntime;

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

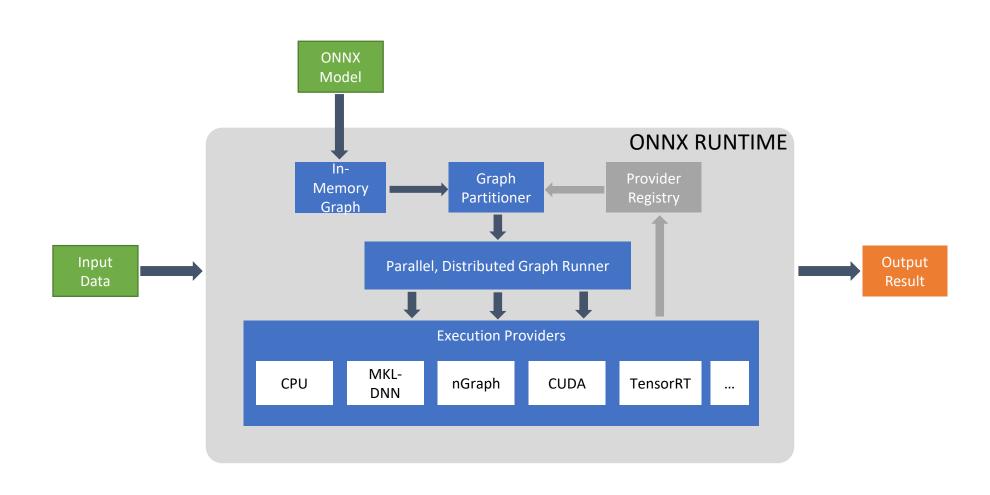
var results = session.Run(input);
```



## Graph optimizations for performance

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

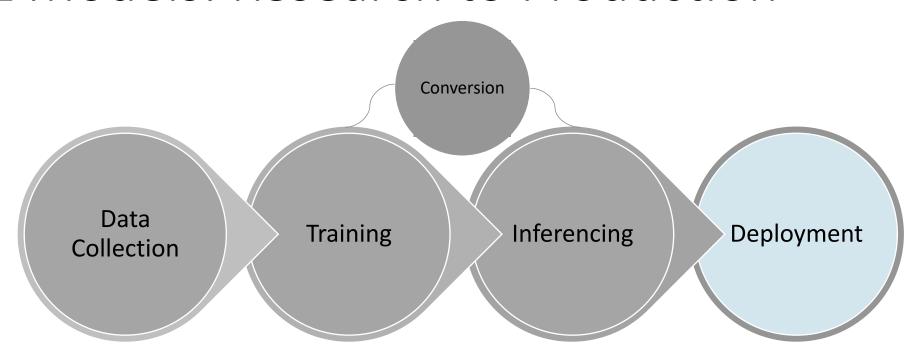
#### Leverages and abstracts hardware accelerators



#### Accelerators for a range of hardware

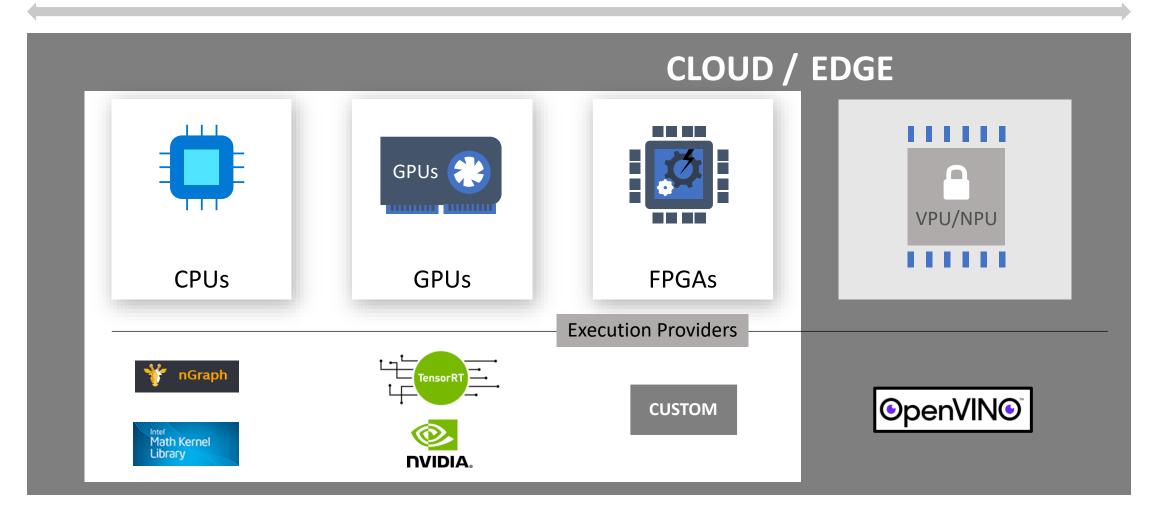
Base CPU **NVIDIA CUDA NVIDIA TensorRT** Microsoft Linear Algebra Subprograms Intel OpenVINO Intel MKL-DNN Intel nGraph **NUPHAR** NN API for Android DirectML TVM/LLVM-based (future) model compiler

#### ML Models: Research to Production

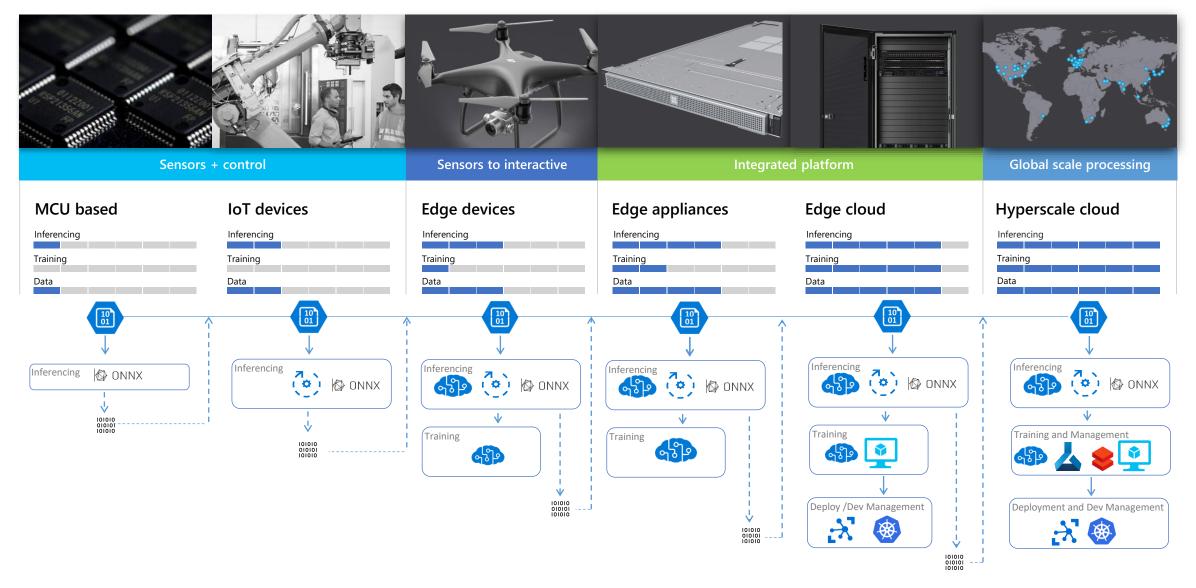


#### Variety of Deployment Options

FLEXIBILITY EFFICIENCY



#### Deployment targets with varying compute power



#### ACTIVITY B

Train a PyTorch model and deploy for production usage

#### References

- ONNX: https://github.com/onnx/onnx
- ONNX Converters: <a href="https://github.com/onnx/onnxmltools/tree/master/onnxmltools">https://github.com/onnx/onnxmltools/tree/master/onnxmltools</a>
- ONNX Tutorials: https://github.com/onnx/tutorials
- ONNX Runtime: https://github.com/microsoft/onnxruntime
- ONNX Runtime Tutorials: <a href="https://github.com/microsoft/onnxruntime#examples-and-tutorials">https://github.com/microsoft/onnxruntime#examples-and-tutorials</a>
- Performance Tuning with ONNX Runtime: <a href="https://github.com/microsoft/onnxruntime/blob/master/docs/ONNX Runtime Perf Tuning.md">https://github.com/microsoft/onnxruntime/blob/master/docs/ONNX Runtime Perf Tuning.md</a>
- Training, Inferencing, and deployment in AzureML with ONNX models: <a href="https://aka.ms/onnxnotebooks">https://aka.ms/onnxnotebooks</a>
- AzureML resources: <a href="https://azure.microsoft.com/en-us/services/machine-learning/">https://azure.microsoft.com/en-us/services/machine-learning/</a>
- 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: <a href="https://docs.microsoft.com/en-us/windows/ai/windows-ml/">https://docs.microsoft.com/en-us/windows/ai/windows-ml/</a>

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