

# Leverage Power of Machine Learning with ONNX

Ron Dagdag  
@rondagdag



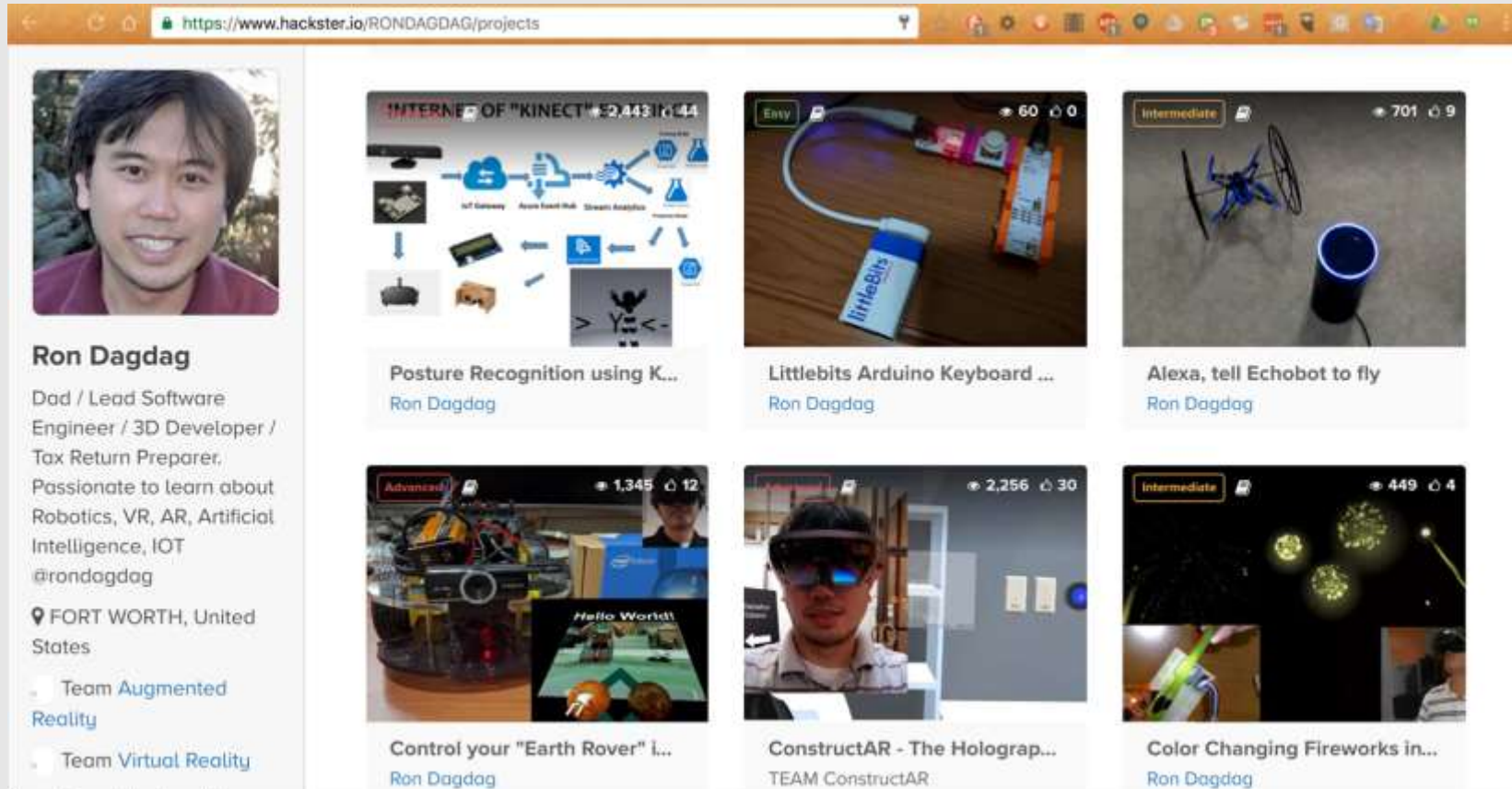
ONNX, Not ONIX



# Hackster Portfolio

[www.dagdag.net](http://www.dagdag.net)

@rondagdag

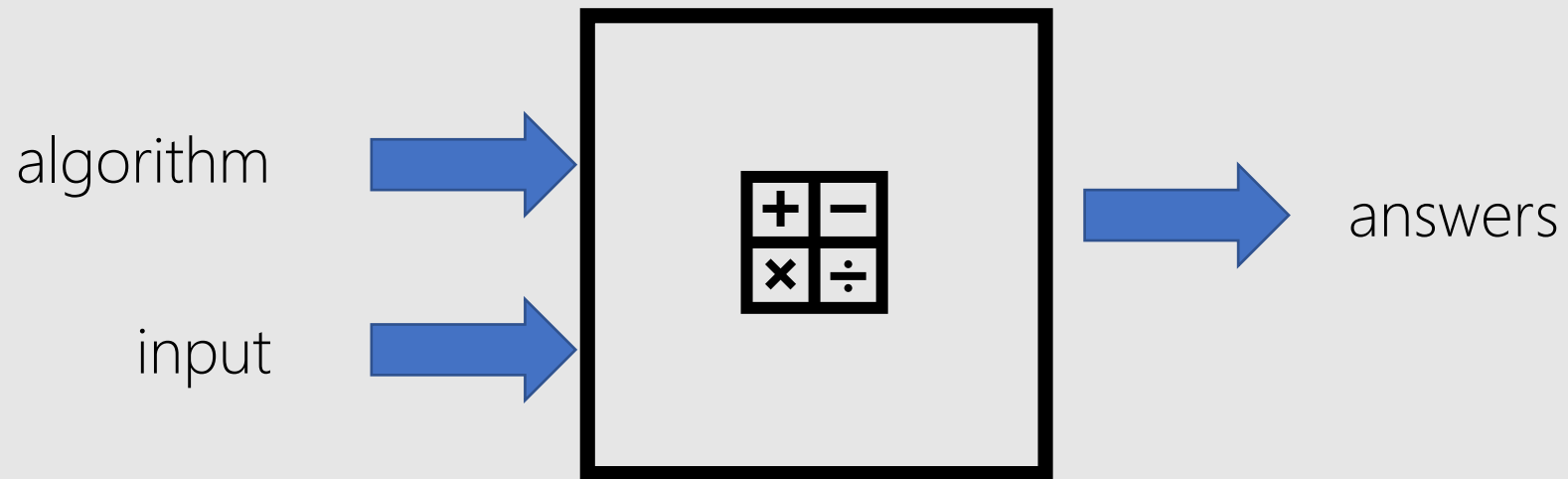


The screenshot displays a web browser window with the URL <https://www.hackster.io/RONDAGDAG/projects>. The profile of Ron Dagdag is shown on the left, including a profile picture, name, bio, location, and team affiliations. The main area features a grid of project thumbnails, each with a title, difficulty level, and engagement metrics.

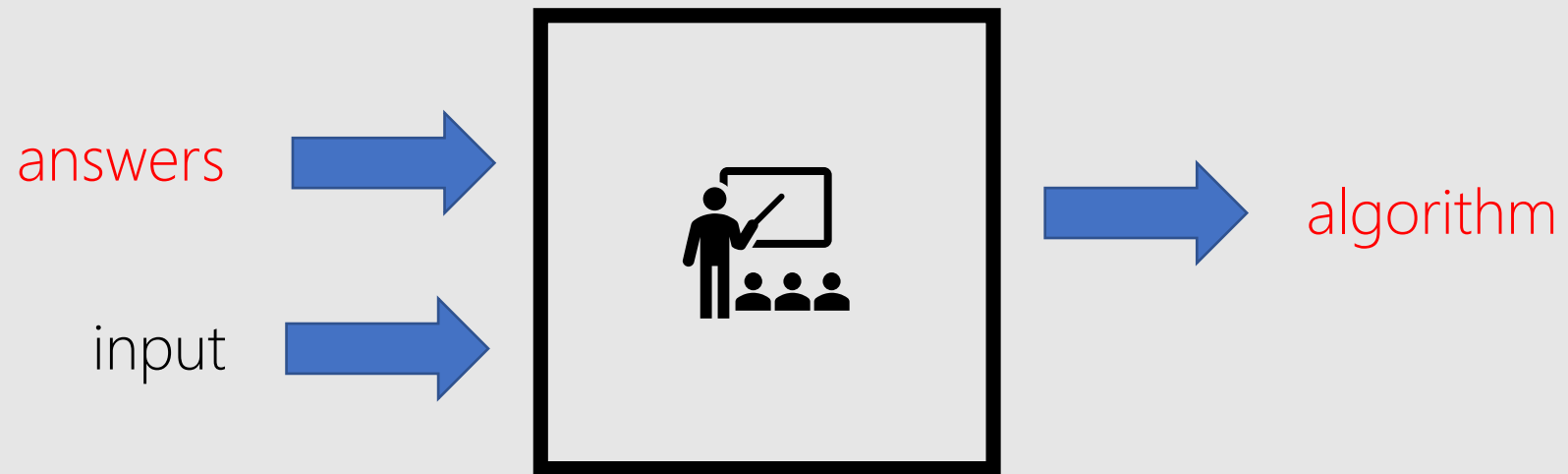
**Ron Dagdag**  
Dad / Lead Software Engineer / 3D Developer / Tax Return Preparer.  
Passionate to learn about Robotics, VR, AR, Artificial Intelligence, IOT  
@rondagdag  
FORT WORTH, United States  
Team [Augmented Reality](#)  
Team [Virtual Reality](#)

Project Title	Difficulty	Views	Upvotes	Comments
Posture Recognition using K...	Easy	2,443	44	
Littlebits Arduino Keyboard ...	Easy	60	0	
Alexa, tell Echobot to fly	Intermediate	701	9	
Control your "Earth Rover" i...	Advanced	1,345	12	
ConstructAR - The Holograp...	Advanced	2,256	30	
Color Changing Fireworks in...	Intermediate	449	4	

# programming



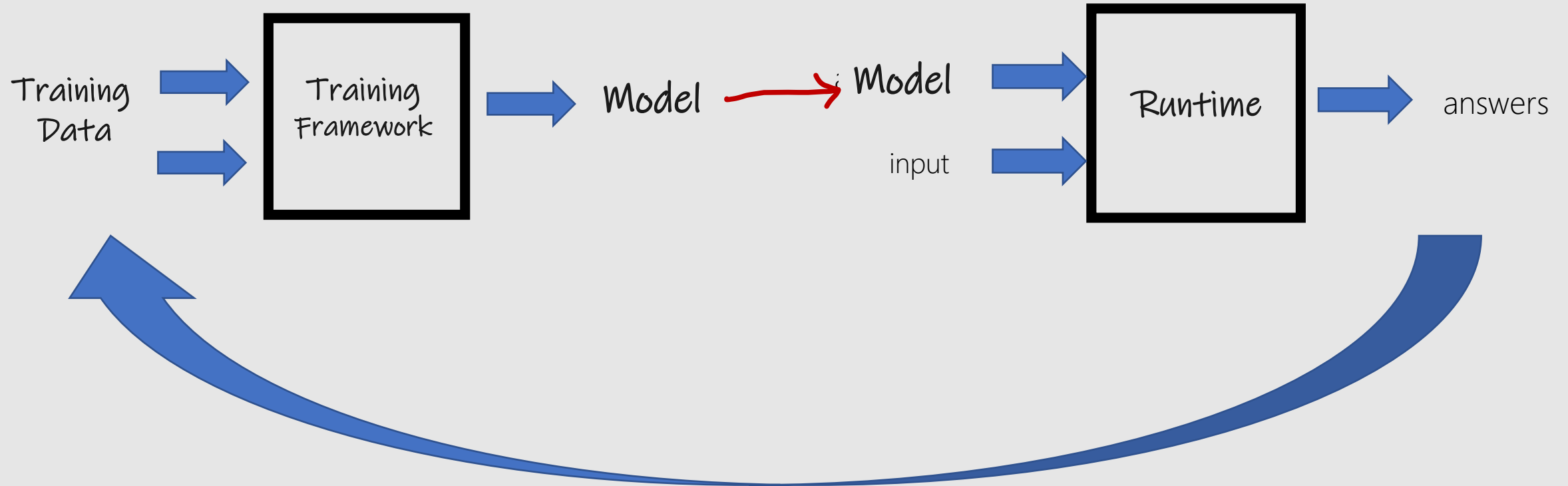
# machine learning



# ML Primer

Machine learning

Inferencing



# Open and Interoperable AI





# ONNX

Open Neural Network Exchange

## Open format for ML models

[github.com/onnx](https://github.com/onnx)

[onnx.ai/](https://onnx.ai/)





# ONNX Partners

---



# Agenda

- ✓ What is ONNX
- ☐ How to create ONNX models
- ☐ How to deploy ONNX models

# Create

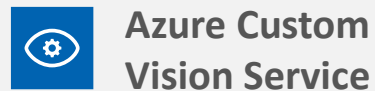
## Frameworks



Native support

Converters

## Services



Native support

ONNX Model

# Deploy

## Cloud Services

Azure Machine Learning services

Ubuntu VM

Windows Server 2019 VM

Windows Devices

IoT Edge Devices

Other Devices  
(iOS, Android, etc)

Native support

Converters

## Frameworks



**Step 1:  
Create**

Services



Azure Custom  
Vision Service

Native  
support

Converters

Native  
support



**ONNX Model**

## Azure

Azure Machine Learning services

Ubuntu VM

Windows Server 2019 VM



**Step 2:  
Deploy**

Other Devices  
(iOS, etc)

Native  
support

Converters

A still life composition featuring a cardboard egg carton with several brown eggs. A pair of wire-rimmed glasses is placed on a stack of papers in the background. The scene is dimly lit, creating a moody atmosphere. The text "Secret Recipe" is overlaid in white, with a vertical line separating the words.

# Secret Recipe

# 4 ways to get an ONNX model



ONNX Model Zoo



Azure Custom Vision Service



Convert existing models



Train models in Azure Machine Learning

Automated Machine Learning

# ONNX Model Zoo: [github.com/onnx/models](https://github.com/onnx/models)

## Image Classification

This collection of models take images as input, then classifies the major objects in the images into a set of predefined classes.

Model Class	Reference	Description
<a href="#">MobileNet</a>	<a href="#">Sandler et al.</a>	Efficient CNN model for mobile and embedded vision applications. Top-5 error from paper - ~10%
<a href="#">ResNet</a>	<a href="#">He et al., He et al.</a>	Very deep CNN model (up to 152 layers), won the ImageNet Challenge in 2015. Top-5 error from paper - ~3.6%
<a href="#">SqueezeNet</a>	<a href="#">Iandola et al.</a>	A lightweight CNN model with fewer parameters than AlexNet. Top-5 error from paper - ~4.8%
<a href="#">VGG</a>	<a href="#">Simonyan et al.</a>	Deep CNN model, won the ImageNet Challenge in 2014. Top-5 error from paper - ~7.4%

Model	Download	Checksum	Download (with sample test data)	ONNX version	Opset version	Top-1 accuracy (%)	Top-5 accuracy (%)
ResNet-18	<a href="#">44.6 MB</a>	<a href="#">MD5</a>	<a href="#">42.9 MB</a>	1.2.1	7	69.70	89.49
ResNet-34	<a href="#">83.2 MB</a>	<a href="#">MD5</a>	<a href="#">78.6 MB</a>	1.2.1	7	73.36	91.43
ResNet-50	<a href="#">97.7 MB</a>	<a href="#">MD5</a>	<a href="#">92.0 MB</a>	1.2.1	7	75.81	92.82
ResNet-101	<a href="#">170.4 MB</a>	<a href="#">MD5</a>	<a href="#">159.4 MB</a>	1.2.1	7	77.42	93.61
ResNet-152	<a href="#">230.3 MB</a>	<a href="#">MD5</a>	<a href="#">216.0 MB</a>	1.2.1	7	78.20	94.21

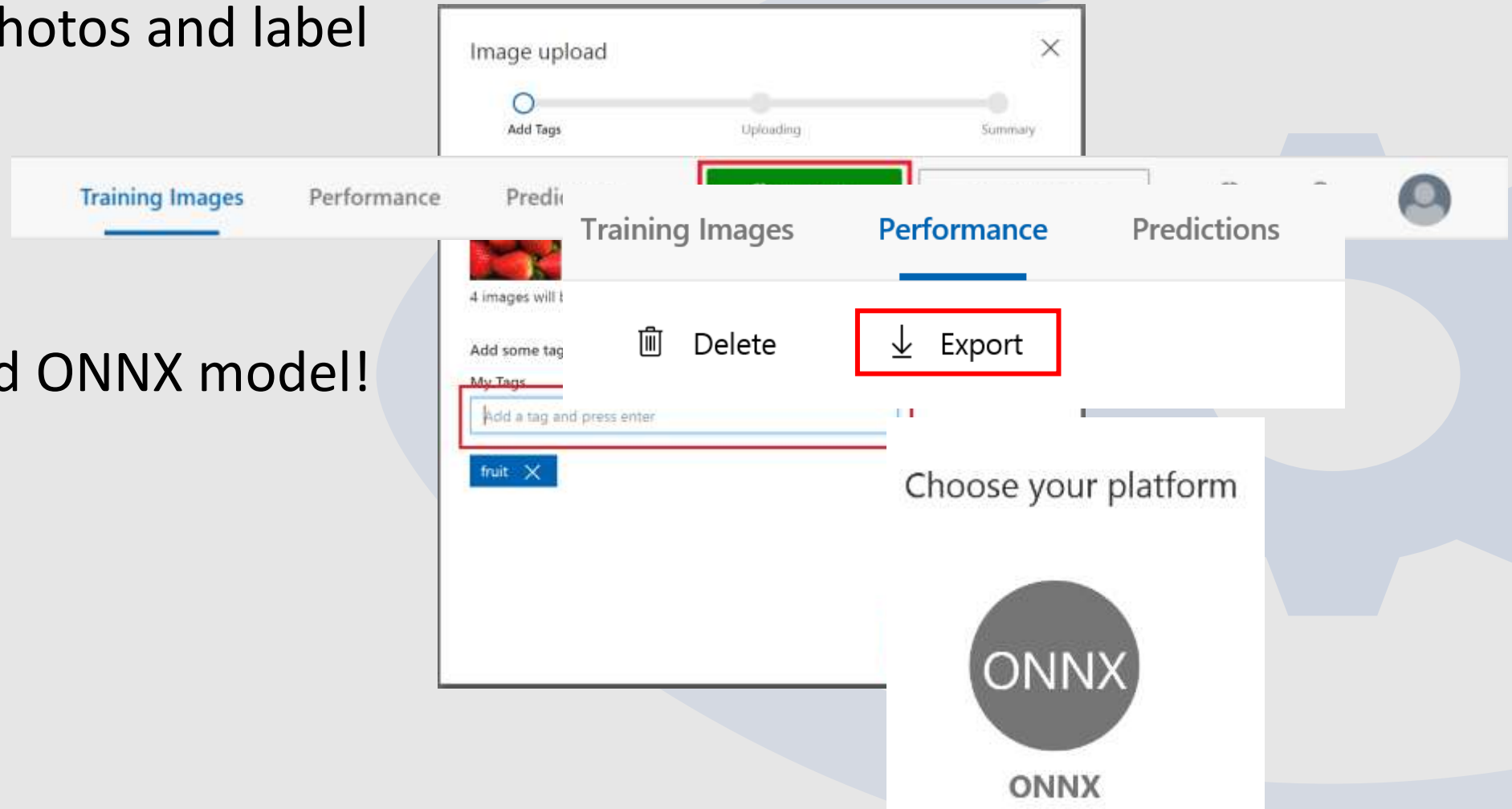


# Custom Vision Service: [customvision.ai](https://customvision.ai)

1. Upload photos and label

2. Train

3. Download ONNX model!





Convert  
models



# Convert models

1. Load existing model
2. (Convert to ONNX)
3. Save ONNX model

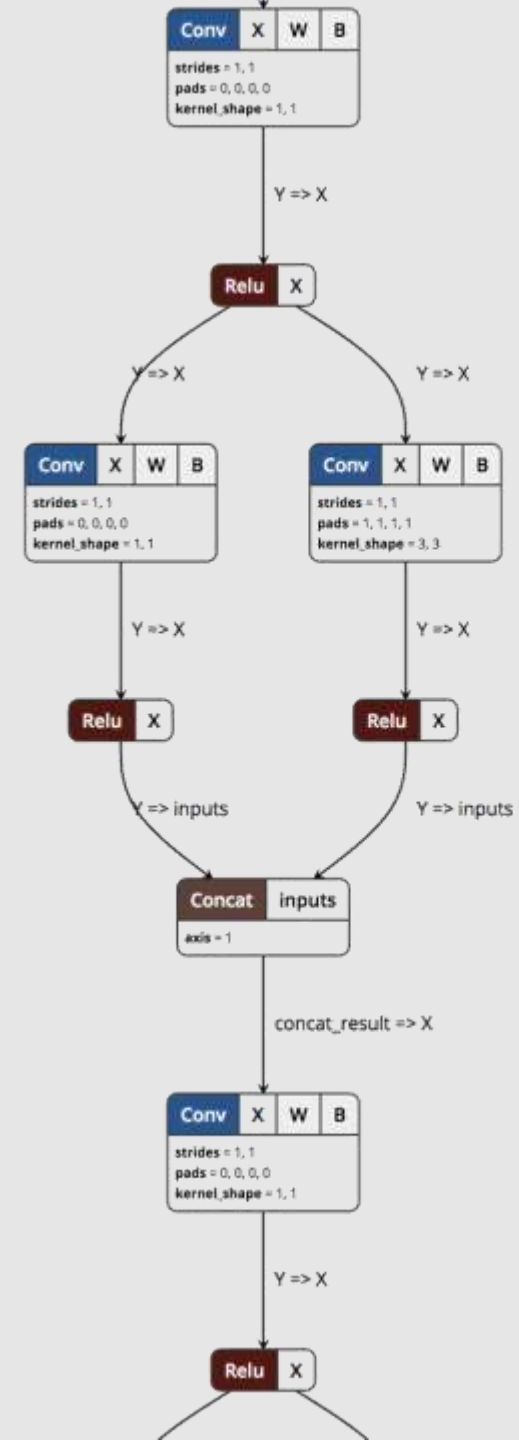


# ONNX Models

## Graph of operations

Netron

<https://lutzroeder.github.io/netron/>

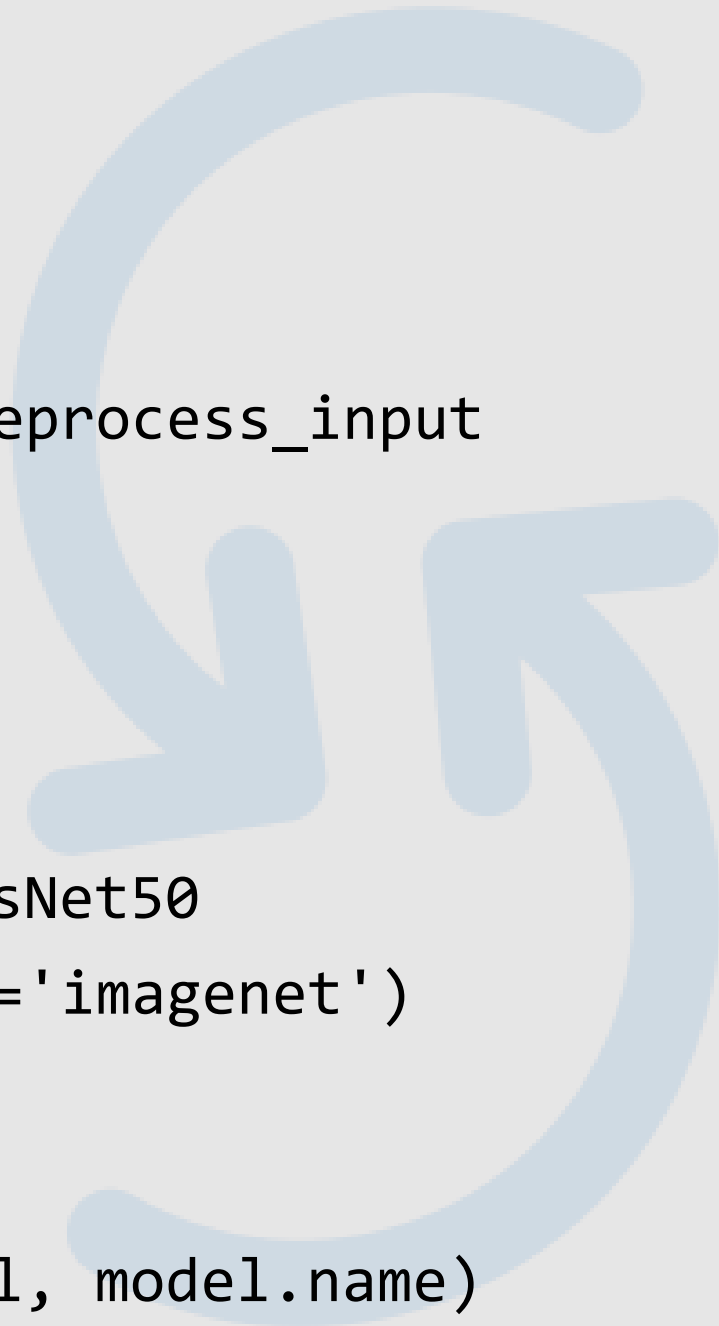


# Convert models: Keras

```
import numpy as np
from keras.preprocessing import image
from keras.applications.resnet50 import preprocess_input
import keras2onnx
import onnxruntime

# load keras model
from keras.applications.resnet50 import ResNet50
model = ResNet50(include_top=True, weights='imagenet')

# convert to onnx model
onnx_model = keras2onnx.convert_keras(model, model.name)
```



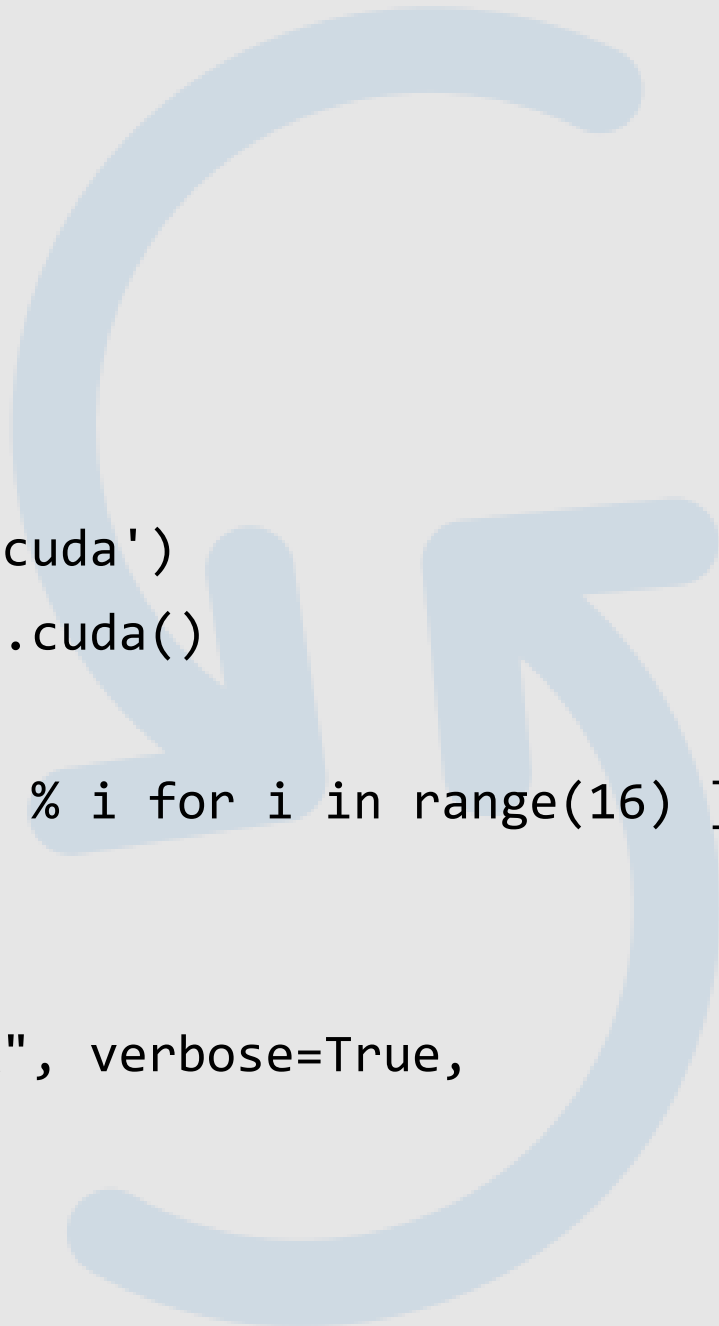
# Convert models: Pytorch

```
import torch
import torchvision

dummy_input = torch.randn(10, 3, 224, 224, device='cuda')
model = torchvision.models.alexnet(pretrained=True).cuda()

input_names = [ "actual_input_1" ] + [ "learned_%d" % i for i in range(16) ]
output_names = [ "output1" ]

torch.onnx.export(model, dummy_input, "alexnet.onnx", verbose=True,
input_names=input_names, output_names=output_names)
```

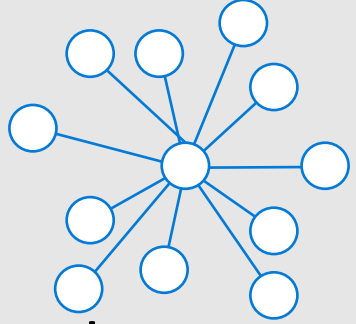


# Convert models: TensorFlow

```
python -m tf2onnx.convert --input frozen.pb --inputs X:0  
--outputs output:0 --output model.onnx
```

<https://github.com/onnx/tensorflow-onnx>



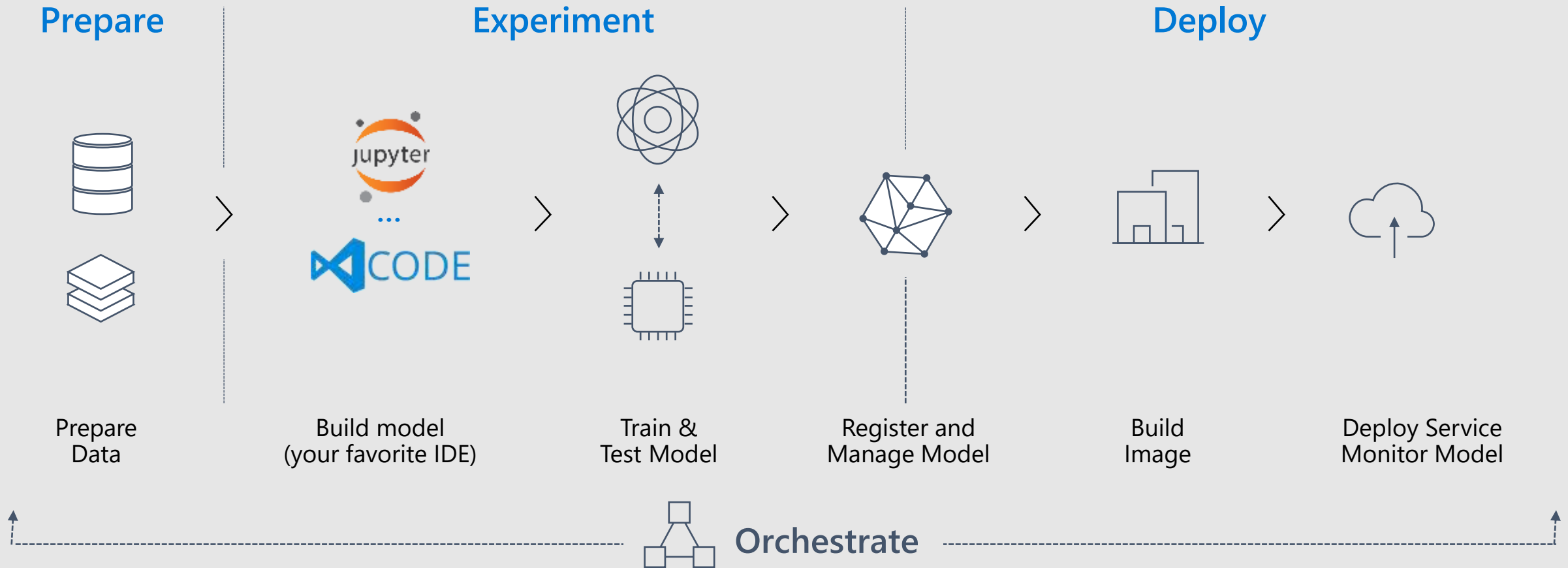


# Train models in Azure Machine Learning

- Experiment locally then quickly scale with GPU clusters in the cloud
- Use automated machine learning and hyper-parameter tuning.
- Keeping Track of experiments, manage models, and easily deploy with integrated CI/CD tooling

# Machine Learning

Typical E2E Process





high  
dimensional  
matrices

# tensor

't'
'e'
'n'
's'
'o'
'r'

tensor of dimensions [6]  
(vector of dimension 6)

3	1	4	1
5	9	2	6
5	3	5	8
9	7	9	3
2	3	8	4
6	2	6	4

tensor of dimensions [6,4]  
(matrix 6 by 4)

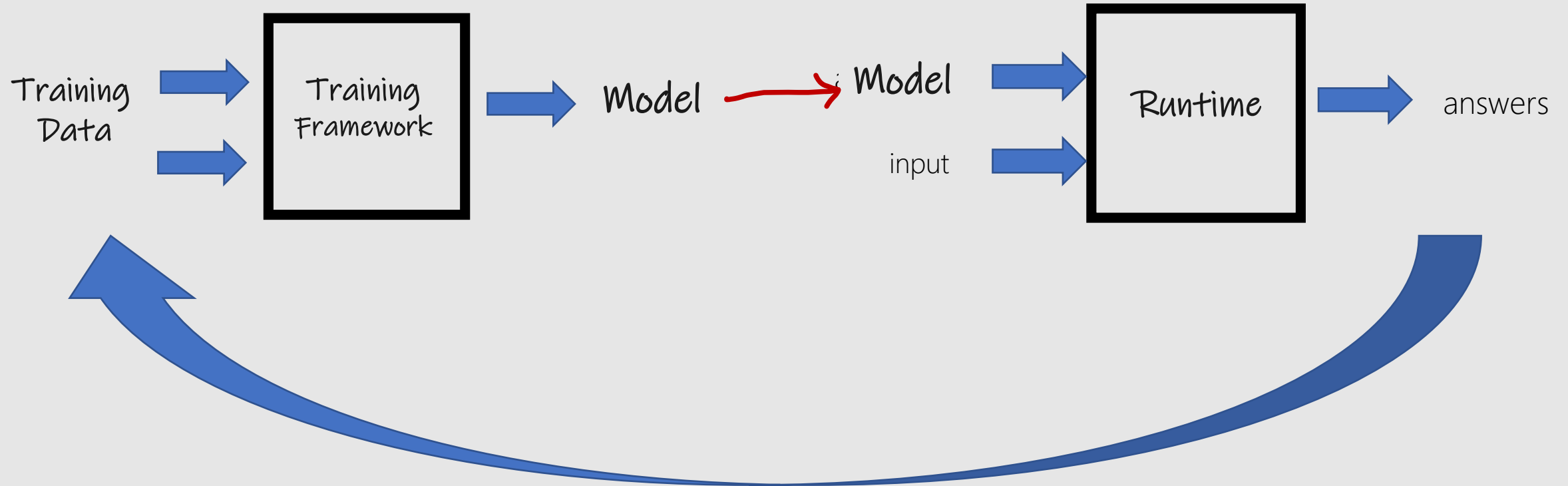
2	1	2	1
2	4	9	4
2	5	6	2
7	7	3	2

tensor of dimensions [4,4,2]

# ML Primer

Machine learning

Inferencing



## Frameworks

Caffe2 Chainer Cognitive Toolkit

mxnet PyTorch PaddlePaddle

ML.NET MathWorks dmlc XGBoost



# Step 1: Create

Services



Azure Custom  
Vision Service

Native  
support

Converters

Native  
support



## ONNX Model

## Azure

Azure Machine Learning services

Ubuntu VM

Windows Server 2019 VM




# Step 2: Deploy

Other Devices  
(iOS, etc)

Native  
support

Converters

A baker in a white shirt is shown from the chest down, working on a wooden table. They are shaping a large, round loaf of bread. The table is covered with a layer of white flour. The background is a plain, light-colored wall.

# Baker vs Starting a Bakery

# Create

## Frameworks



Native support

Converters

## Services



Native support

ONNX Model

# Deploy

## Azure

Azure Machine Learning services

Ubuntu VM

Windows Server 2019 VM

Windows Devices

IoT Edge Devices

Other Devices  
(iOS, etc)

Native support

Converters

A person's hands are visible, holding a large, round, rustic loaf of bread. The bread has a thick, golden-brown crust with some darker, caramelized spots. It is wrapped in a blue and white striped cloth. The background is a blurred wooden surface.

# Cloud or Edge



Cloud  
or  
Edge



# Deploy with Azure Machine Learning

- Model management services
- Deploy as web service to ACI or AKS
- Capture model telemetry

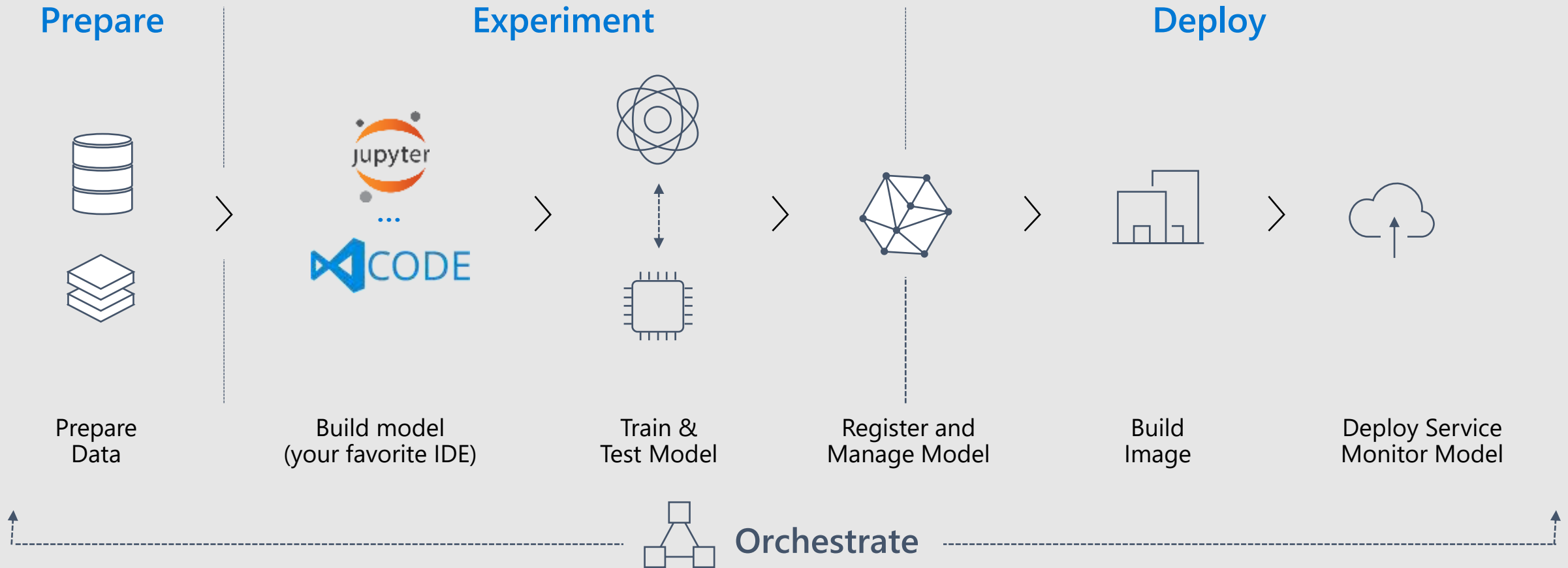


Azure  
Machine Learning

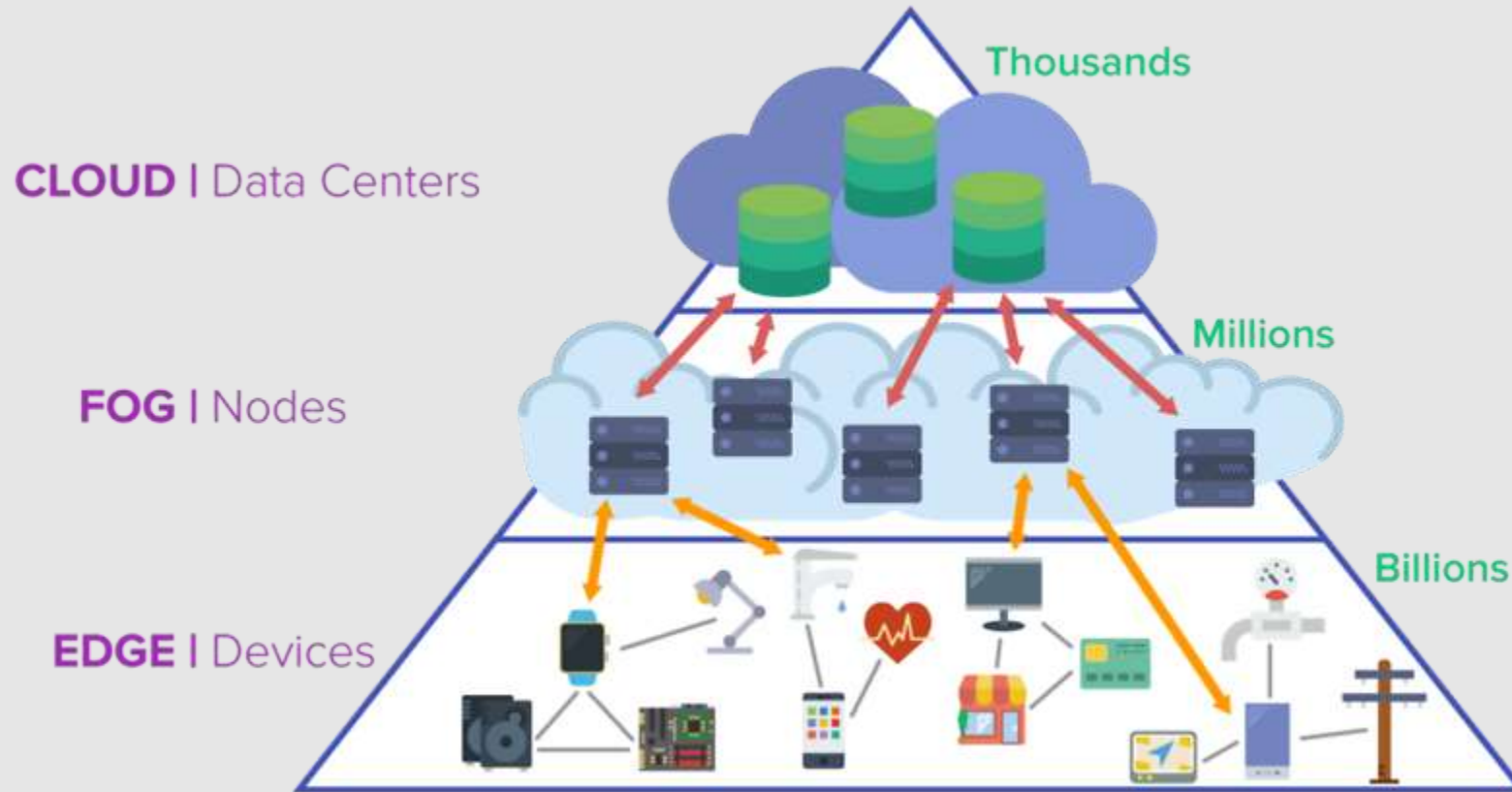


# Machine Learning

Typical E2E Process

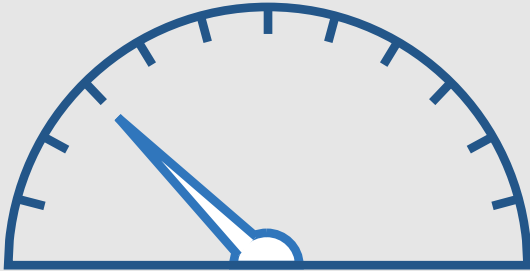


# What is the Edge?

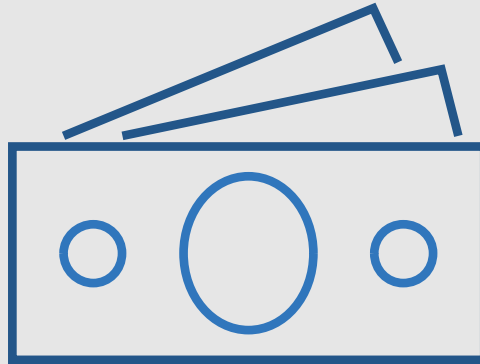


[Imagimob AB](#)

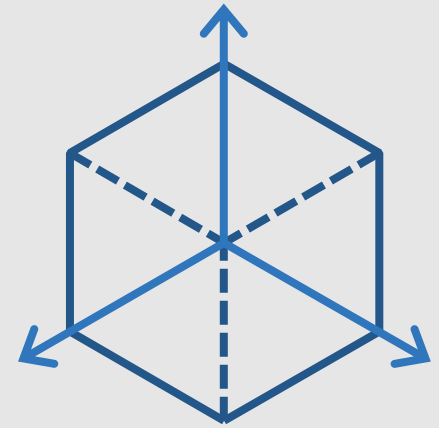
# AI on the edge



Low latency



Scalability



Flexibility

# ONNX as an intermediary format

- **Convert to Tensorflow for Android**
  - [Convert a PyTorch model to Tensorflow using ONNX](#)
- **Convert to CoreML for iOS**
  - <https://github.com/onnx/onnx-coreml>
- **Fine-tuning an ONNX model with MXNet/Gluon**
  - [https://mxnet.apache.org/versions/master/tutorials/onnx/fine\\_tuning\\_gluon.html](https://mxnet.apache.org/versions/master/tutorials/onnx/fine_tuning_gluon.html)

# ONNX Runtime

- High performance runtime for ONNX models
- Supports full ONNX-ML spec (currently v1.2+)
- Extensible architecture to plug-in hardware accelerators
- Simple Python API



ONNX

# ONNX Runtime



# ONNX

## Get Started Easily

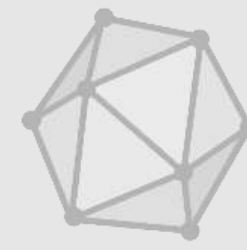
Select your requirements and use the resources provided to get started quickly

OS	Windows		Linux		Mac	
Language	Python (3.5-3.7)	C++	C#	C	Java	
Architecture	X64	X86		ARM64		ARM32
Hardware Acceleration	Default CPU	CUDA		TensorRT	DirectML	MKL-DNN
	MKL-ML	nGraph		NUPHAR		OpenVINO
Installation Instructions	Install Nuget package <a href="#">Microsoft.ML.OnnxRuntime.Gpu</a>					



# Demo

<https://github.com/rondagdag/LeverageONNX>



ONNX

# ONNX Docker Image

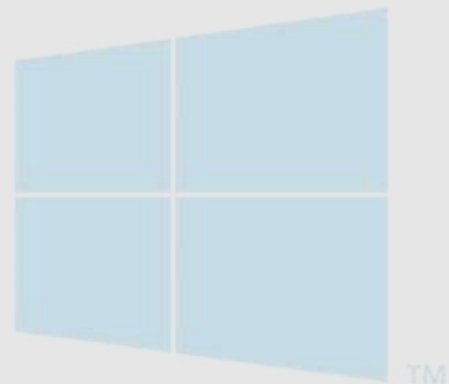
- [onnx-base](#): Use published ONNX package from PyPi with minimal dependencies.
- [onnx-dev](#): Build ONNX from source with minimal dependencies.
- [onnx-ecosystem](#): Jupyter notebook environment for getting started quickly with ONNX models, ONNX converters, and inference using ONNX Runtime.



# Deploy to Windows Devices

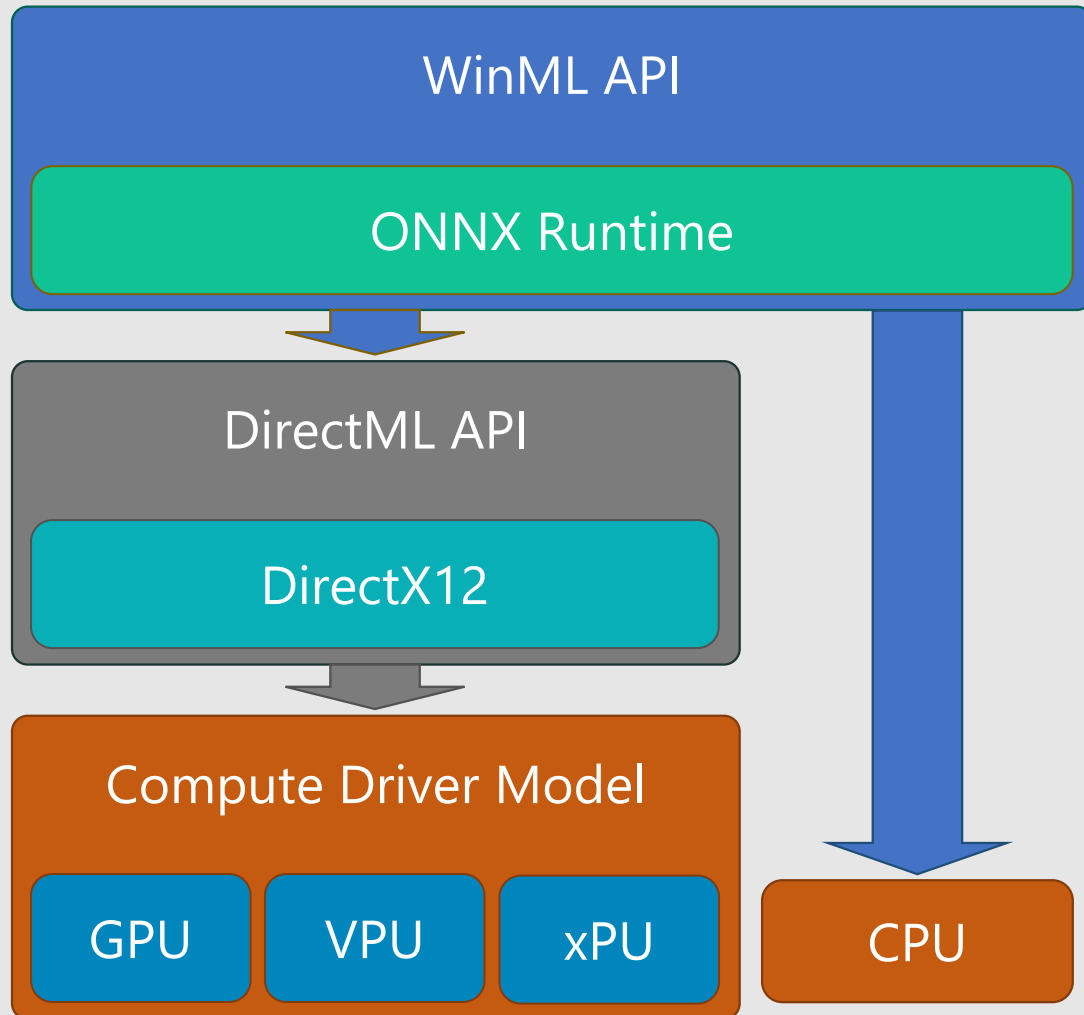
## Windows ML

- Available across Windows family of devices
- Hardware abstraction via DirectML
- Unified API for Win32 and WinRT
- Optimized for performance
- Virtualization ready



Windows®

# Windows AI platform



- WinML
  - **Practical**, simple model-based API for ML inferencing on Windows
- DirectML
  - **Realtime, high control** ML operator API; part of DirectX family
- Compute Driver Model
  - Robust **hardware reach**/abstraction layer for compute and graphics silicon

# ONNX Runtime

- Microsoft services have seen an **average 2x performance gain on CPU**
- Office team saw a **14.6x reduction in latency** for a grammar checking model (thousands of queries per minute)
- Azure Cognitive Services saw a **3.5x reduction in latency** for an optical character recognition (OCR) model
- Bing QnA saw a **2.8x reduction in latency** for a model that generates answers to questions
- Bing Visual Search saw a **2x reduction in latency** for a model that helps identify similar images

# ONNX Runtime - Python API

```
import onnxruntime
```

```
session = onnxruntime.InferenceSession("mymodel.onnx")
```

```
results = session.run([], {"input": input_data})
```



ONNX

# Reference implementation to use ONNX Runtime with Azure IoT Edge



- <https://github.com/Azure-Samples/onnxruntime-iot-edge>



ONNX

# ONNX.js

- ONNX.js is a JavaScript library for running ONNX models on browsers and on Node.js.
- ONNX.js has adopted Web Assembly and WebGL technologies
- optimized ONNX model inference runtime for both CPUs and GPUs.

<https://github.com/microsoft/onnxjs>



ONNX

# ONNX.js

## Compatibility

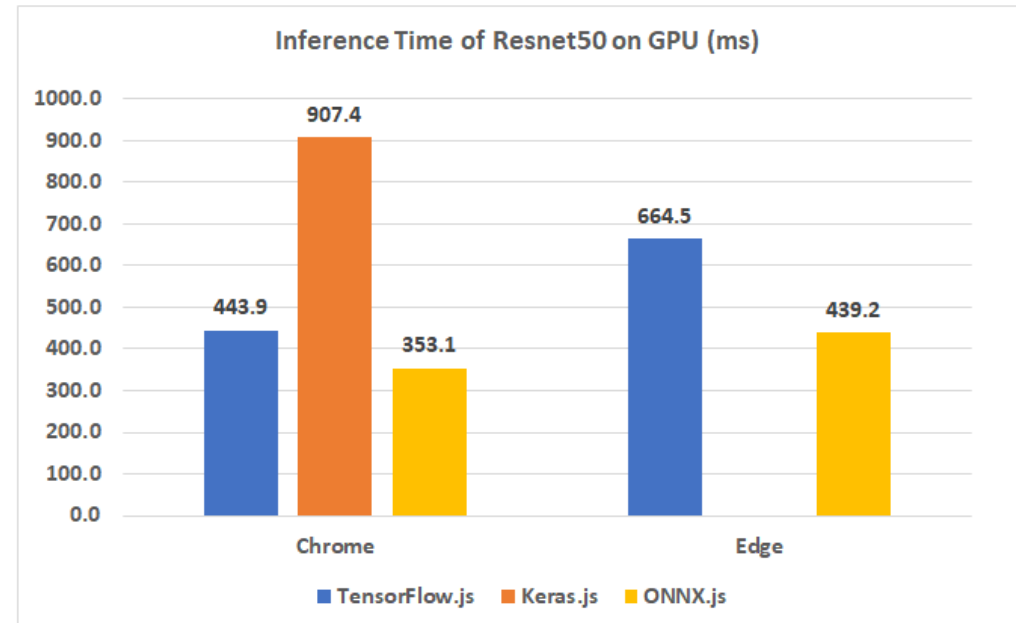
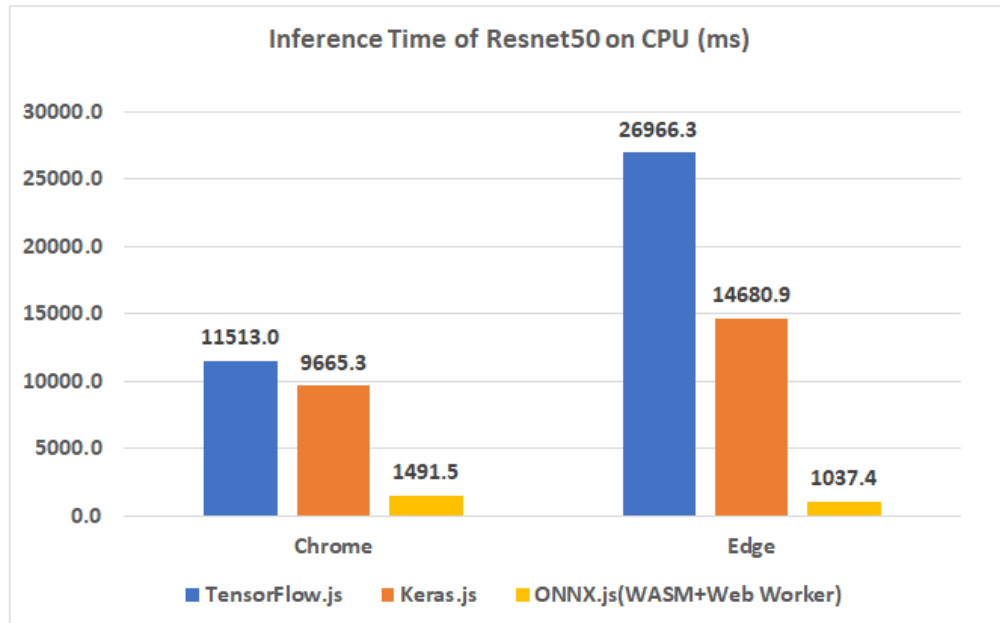
### Desktop Platforms

OS/Browser	Chrome	Edge	FireFox	Safari	Opera	Electron	Node.js
Windows 10	✓	✓	✓	-	✓	✓	✓
macOS	✓	-	✓	✓	✓	✓	✓
Ubuntu LTS 18.04	✓	-	✓	-	✓	✓	✓

### Mobile Platforms

OS/Browser	Chrome	Edge	FireFox	Safari	Opera
iOS	✓	✓	✓	✓	✓
Android	✓	✓	Coming soon	-	✓

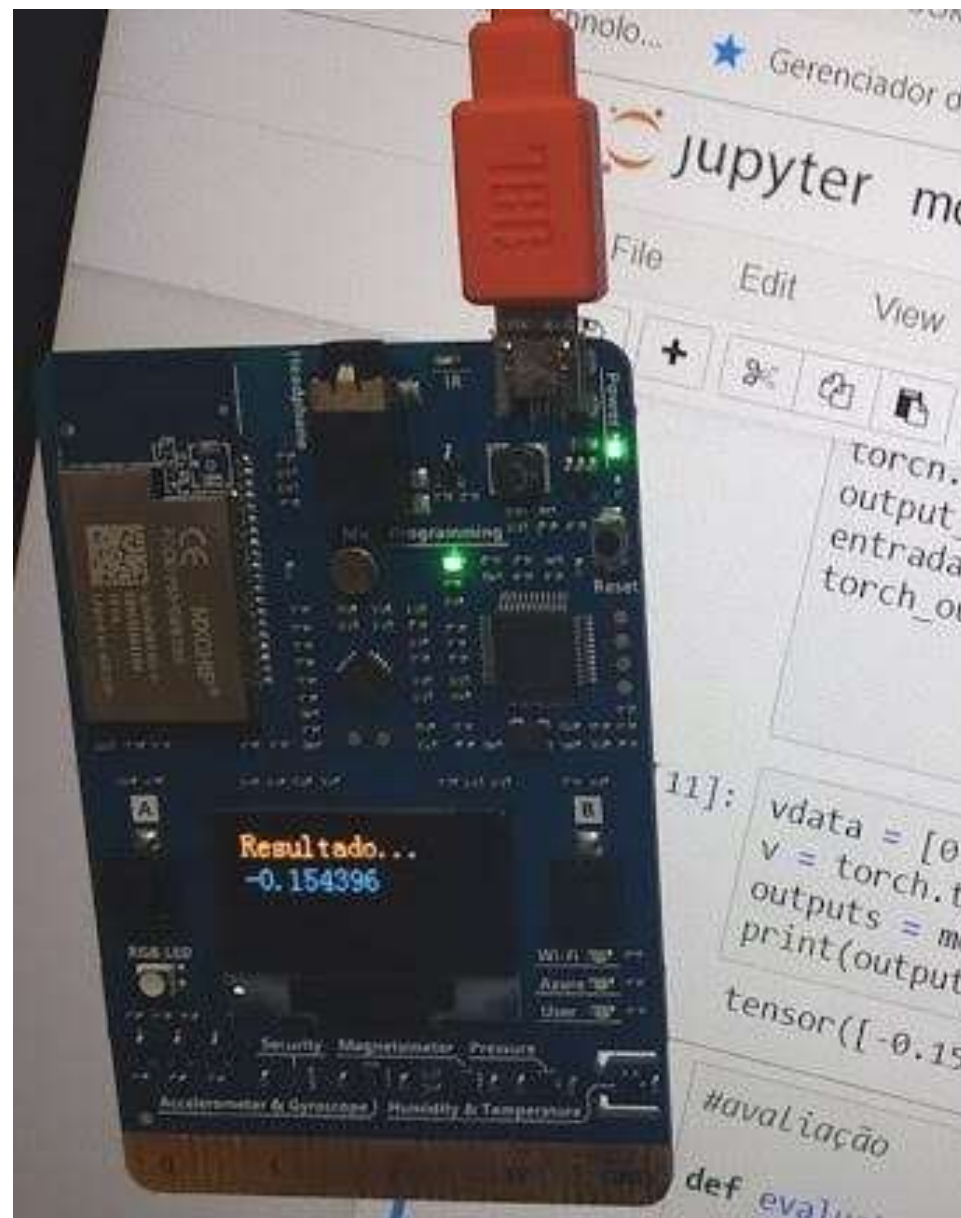
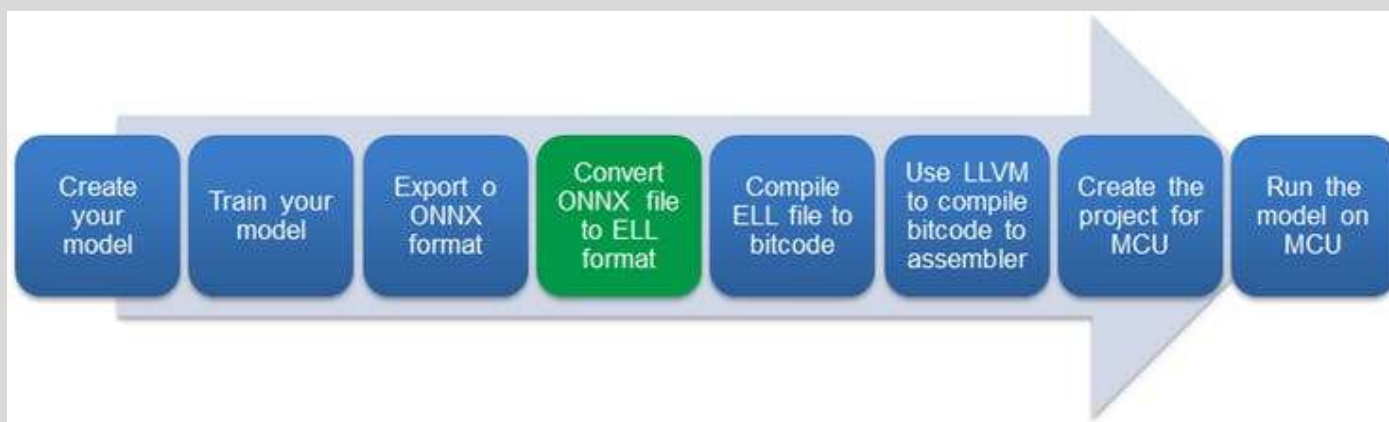
# ONNX.js





# Wait... there's more

- Embedded Learning Library
  - <https://github.com/microsoft/ELL>
- Machine Learning Model Running on Azure IoT Starter Kit
  - <https://www.hackster.io/waltercoan/machine-learning-model-running-on-azure-iot-starter-kit-f9608b>





# Recap

- ✓ What is ONNX

**ONNX is an open standard so you can use the right tools for the job and be confident your models will run efficiently on your target platforms**

- ✓ How to create ONNX models

**ONNX models can be created from many frameworks**

- ✓ How to deploy ONNX models

**ONNX models can be deployed with Windows ML, .NET/Javascript/Python and to the cloud with Azure ML and the high performance ONNX Runtime**

# Try it for yourself!

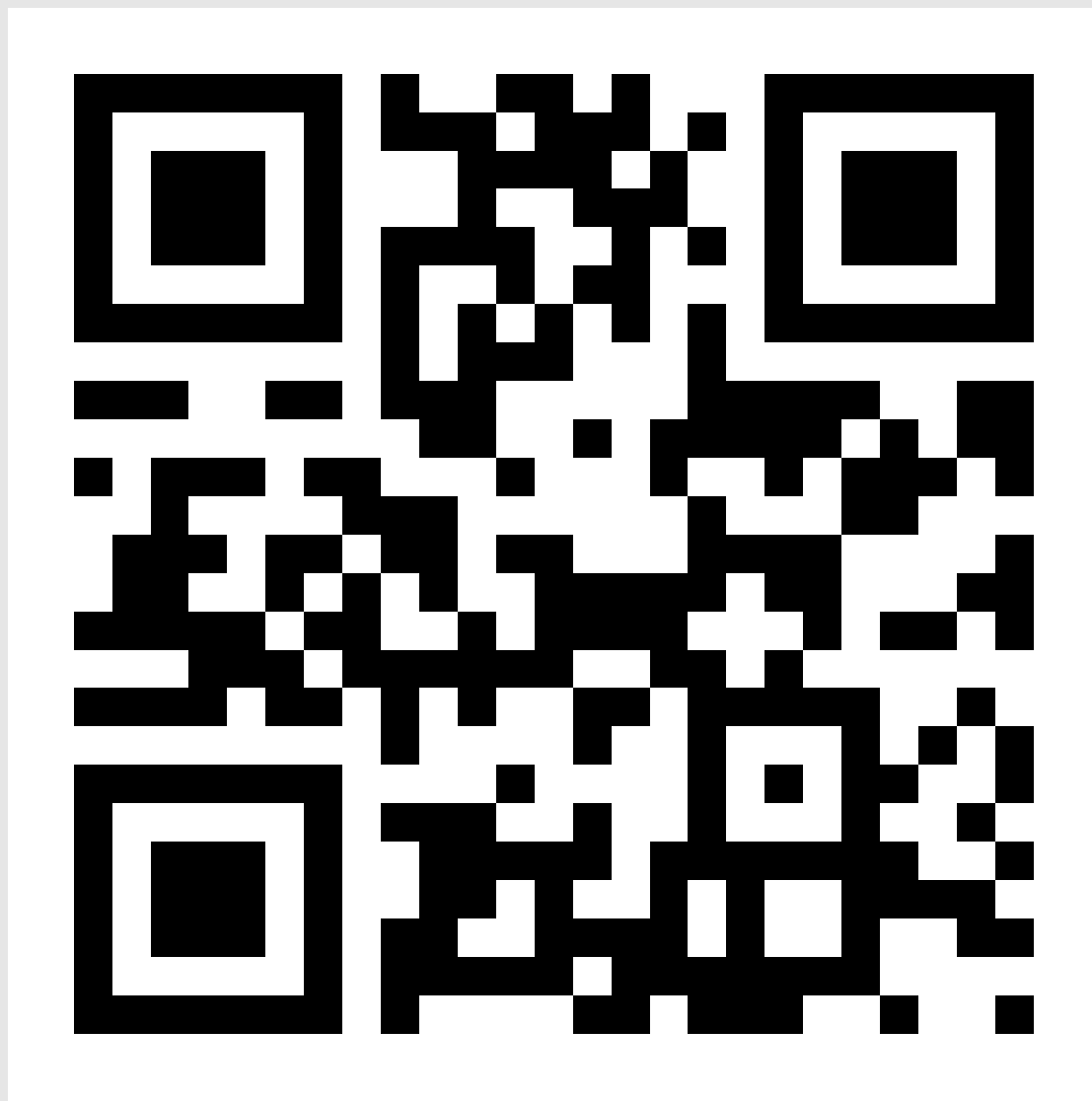
ONNX Runtime is available now!

```
pip install onnxruntime  
pip install onnxruntime-gpu
```

Documentation and samples at [aka.ms/onnxruntime](https://aka.ms/onnxruntime)

Source for Demo:

<https://github.com/rondagdag/onnx-pected>



<http://bit.ly/ml-onnx>

# About Me

## Ron Dagdag



**Ron Lyle Dagdag**

Immersive Experience Developer

Cell: 682-560-3988

[ron@dagdag.net](mailto:ron@dagdag.net)



Experience AR

[www.dagdag.net](http://www.dagdag.net)

[@rondagdag](https://twitter.com/rondagdag)

<http://ron.dagdag.net>

Sr. Software Developer/AI Edge Engineer at Spacee

Microsoft MVP award – Windows Development

Personal Projects  
[www.dagdag.net](http://www.dagdag.net)

Email: [ron@dagdag.net](mailto:ron@dagdag.net)  
Twitter [@rondagdag](https://twitter.com/rondagdag)

Connect me via Linked In  
[www.linkedin.com/in/rondagdag/](https://www.linkedin.com/in/rondagdag/)

Feedback appreciated, help improve my presentation skills