



Well..  
that's  
ONNX-pected

“...there are no wrong turns, only  
unexpected paths.”




<http://bit.ly/onnxpected>

# Hackster Portfolio


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

@rondagdag


https://www.hackster.io/RONDAGDAG/projects




**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](#)




**Posture Recognition using K...**  
Ron Dagdag




**Littlebits Arduino Keyboard ...**  
Ron Dagdag




**Alexa, tell Echobot to fly**  
Ron Dagdag



**Control your "Earth Rover" i...**  
Ron Dagdag



**ConstructAR - The Holograp...**  
TEAM ConstructAR



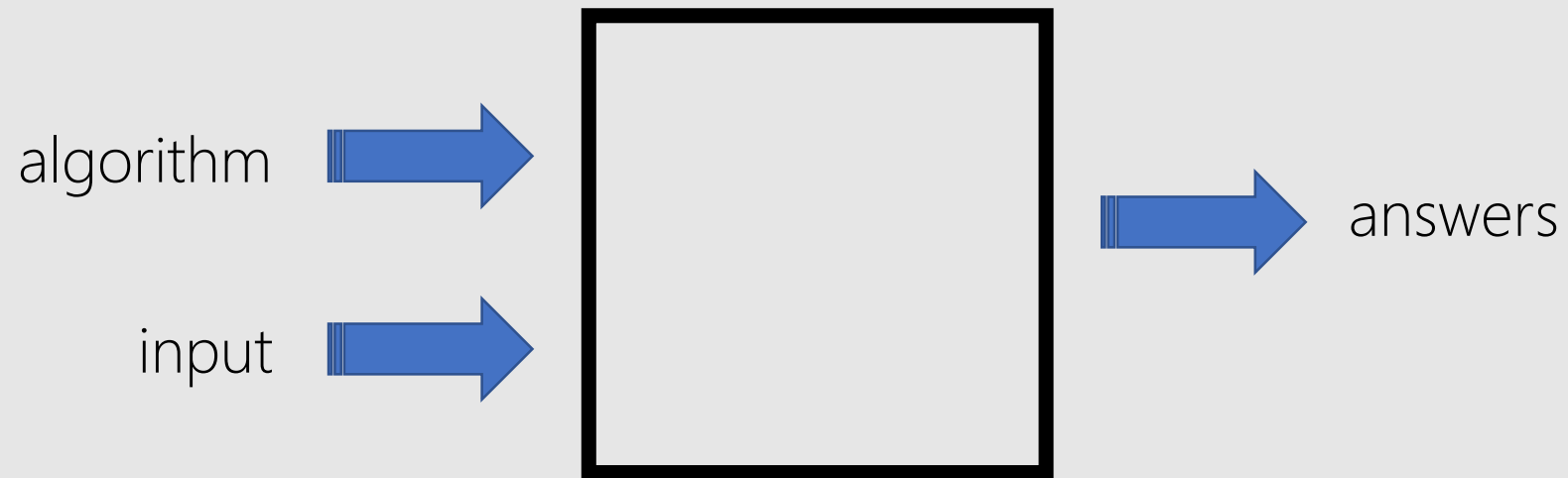
**Color Changing Fireworks in...**  
Ron Dagdag



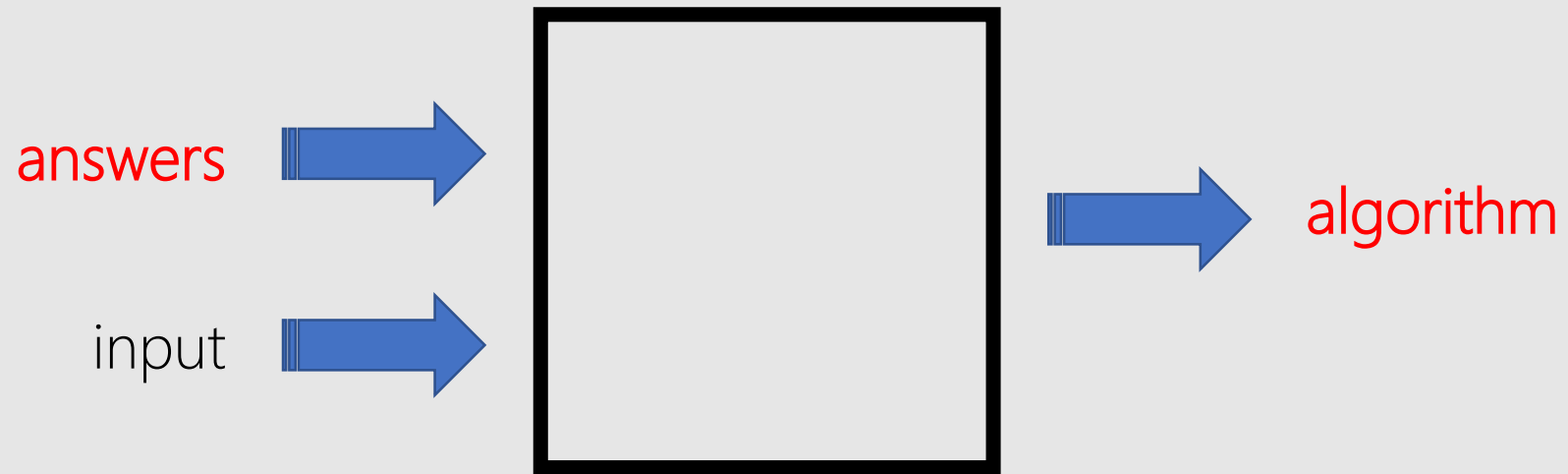
## Audience Survey:

- .NET Developers
- JavaScript Developers
- Python Developers
- Data Scientists
- Data Engineers
- Pokemon Trainers?

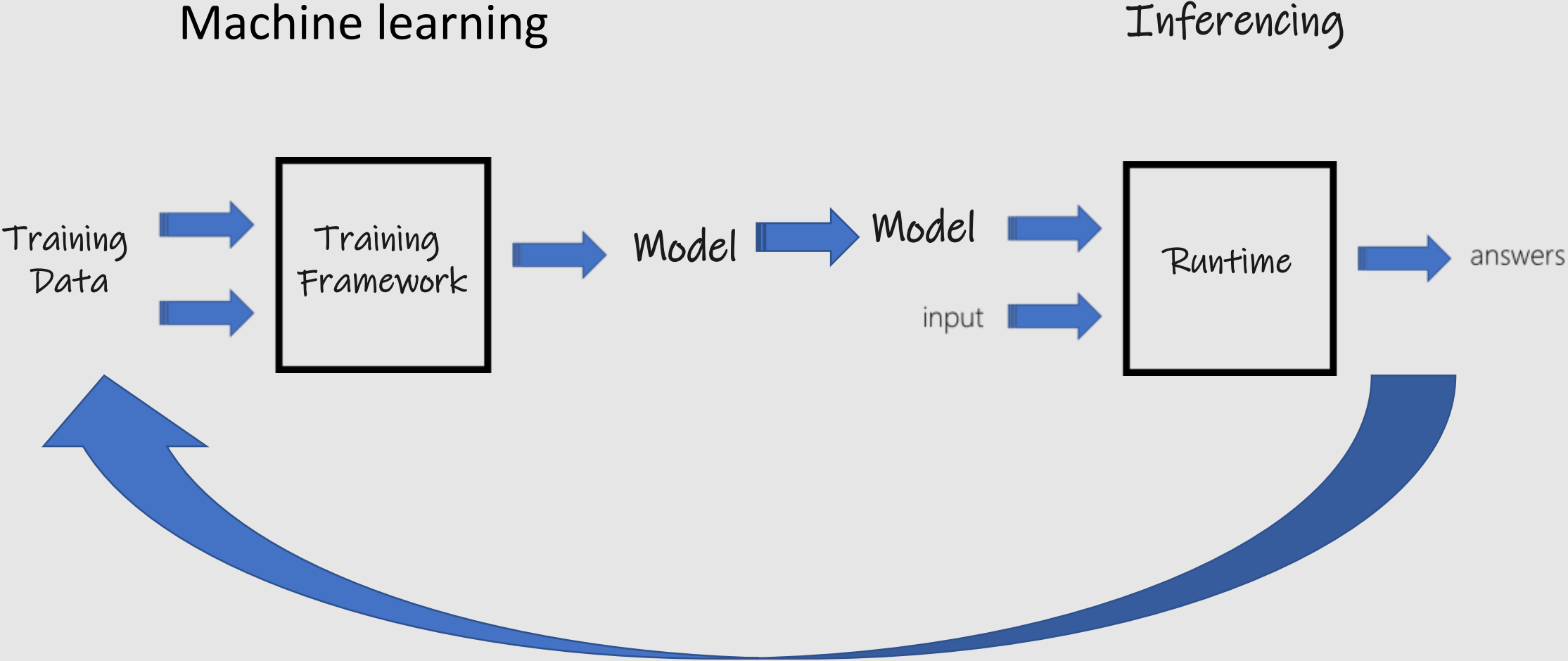
# programming



# machine learning

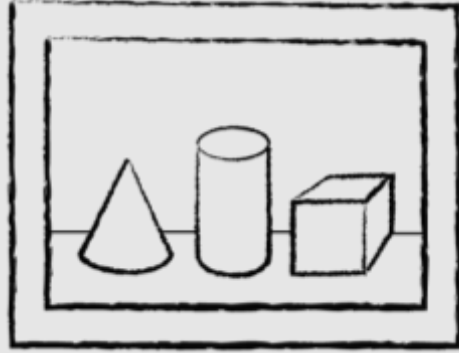


# ML Primer



# process

identify



1

predict



4

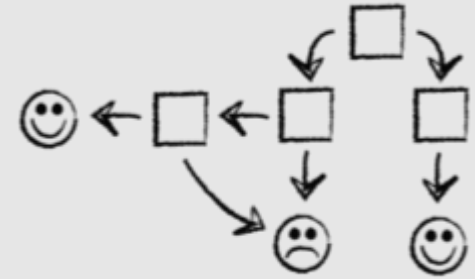
explore  
analyze  
encode

4	3	1	2	6
4	3	2	3	8
5	3	1	2	7
2	1	5	8	0
1	0	8	6	4
8	6	5	7	1

$$\begin{bmatrix} X + Y = Z \\ A + B = C \\ X + A = Z \\ B + Y = ? \end{bmatrix}$$

2

model



3





# Open and Interoperable AI





Open Neural Network Exchange

Open format for ML models

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



# ONNX Partners

---



Facebook  
Open Source



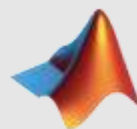
Microsoft



BITMAIN



Hewlett Packard  
Enterprise



MathWorks®



NVIDIA®



Preferred  
Networks



QUALCOMM®



unity

# Agenda

- ✓ What is ONNX?
- ☐ Create ONNX models
- ☐ Deploy ONNX models

# Create

## Frameworks



Native support

Converters

## Services



Native support

ONNX Model

# Deploy

## Azure

Azure Machine Learning services

Ubuntu VM

Windows Server 2019 VM

Windows 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



Windows Devices

**Step 2:  
Deploy**

Other Devices  
(iOS, etc)

Native  
support

Converters

A still life composition featuring several brown eggs nestled in a grey cardboard egg carton. The carton is positioned in the lower-left foreground. In the background, a wire mesh structure, possibly a decorative object or a piece of equipment, is visible on the right. To the left of the mesh, there are some papers or documents, one of which appears to be a recipe book or a notebook. The overall lighting is soft and warm, creating a cozy and rustic atmosphere. The text "Secret Recipe" is overlaid in a white, serif font, centered horizontally and slightly to the right of the eggs.

# Secret Recipe

# 4 ways to get an ONNX model



ONNX Model Zoo



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 and less computation. 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%

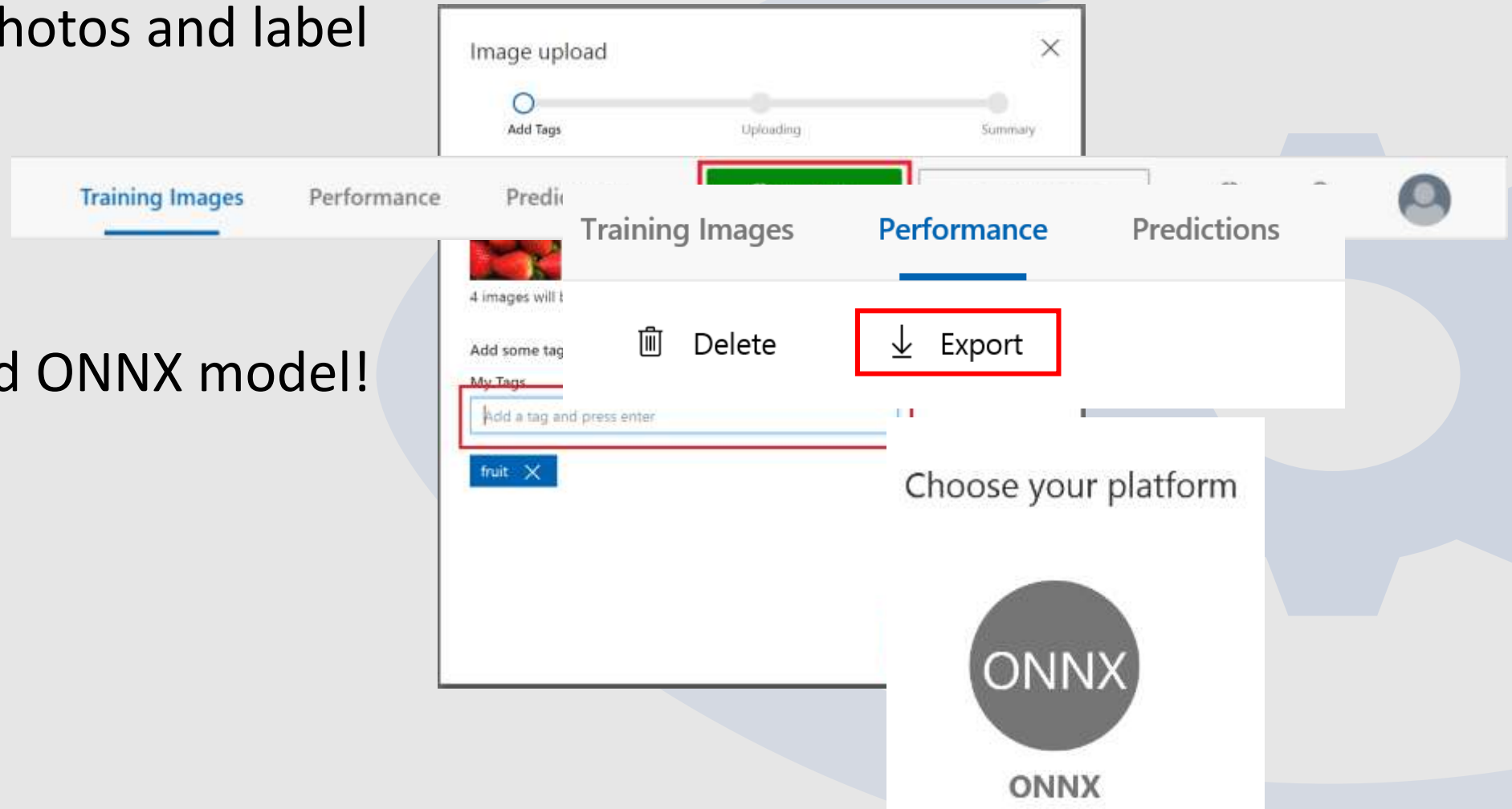
Model	Download	Checksum	Download (with sample test data)	ONNX version	Opset version	Top-1 accuracy (%)	Top-5 accuracy (%)
ResNet-18	<a href="#">44.6 MB</a>	<a href="#">MD5</a>	<a href="#">42.9 MB</a>	1.2.1	7	69.70	89.49
ResNet-34	<a href="#">83.2 MB</a>	<a href="#">MD5</a>	<a href="#">78.6 MB</a>	1.2.1	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!



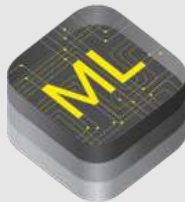
# Convert models



ML.NET



dmlc  
**XGBoost**



# Convert models

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

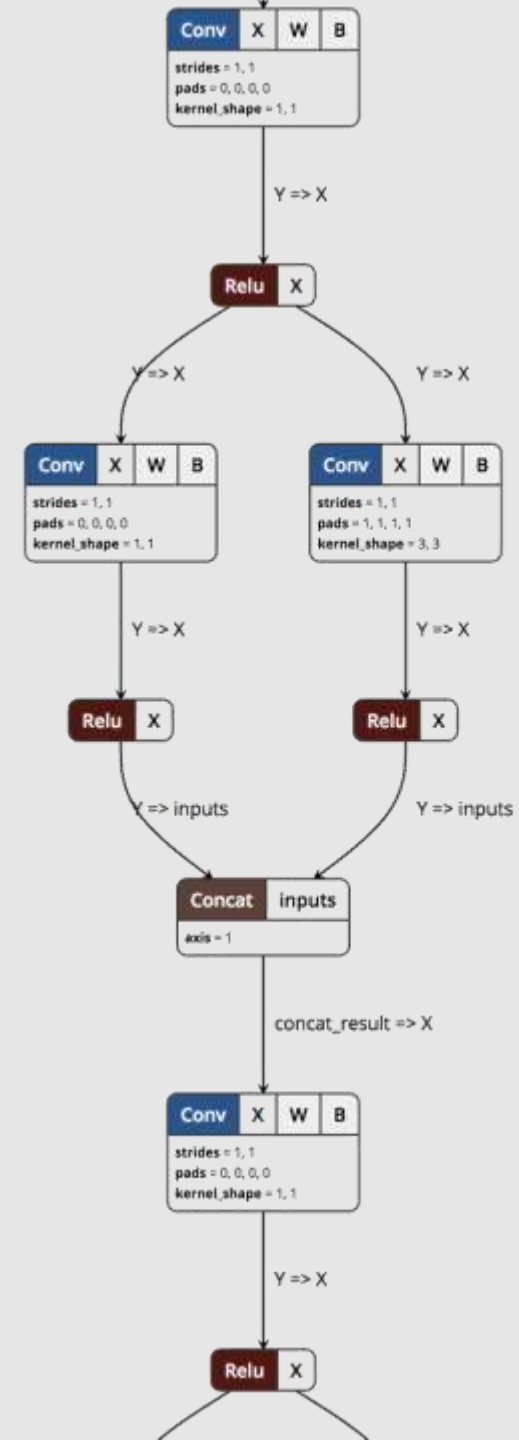


# ONNX Models

Graph of operations

WinML Dashboard

<https://github.com/microsoft/Windows-Machine-Learning/releases>



# ONNX Ecosystem Docker Image

- [onnx-ecosystem](#): Jupyter notebook environment for getting started quickly with ONNX models, ONNX converters, and inference using ONNX Runtime.

```
> docker pull onnx/onnx-ecosystem
```

```
> docker run -p 8888:8888 onnx/onnx-ecosystem
```

```
http://127.0.0.1:8888/?token={tokenId}
```

```
http://127.0.0.1:8888/tree/converter_scripts
```



ONNX

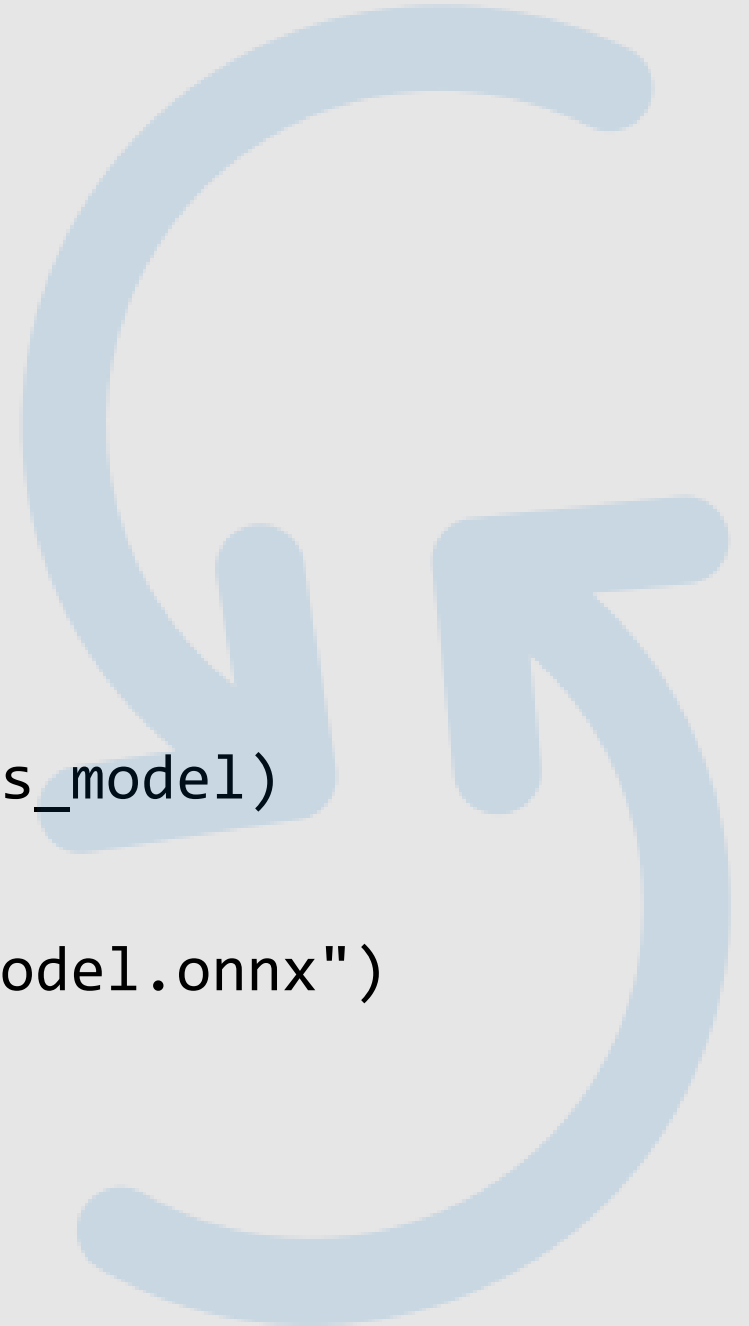
# Convert models: Keras

```
from keras.models import load_model  
import winmltools
```

```
keras_model = load_model("model.h5")
```

```
onnx_model = winmltools.convert_keras(keras_model)
```

```
winmltools.utils.save_model(onnx_model, "model.onnx")
```

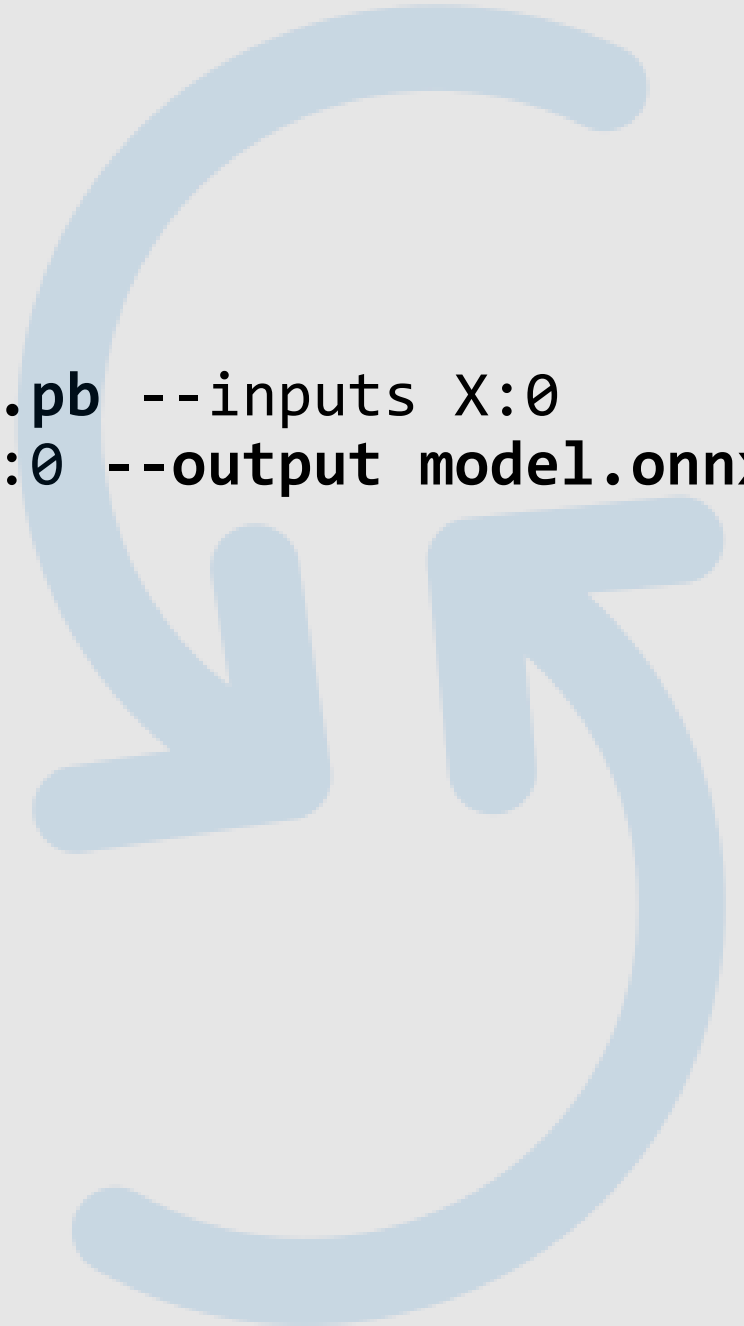


# Convert models: TensorFlow

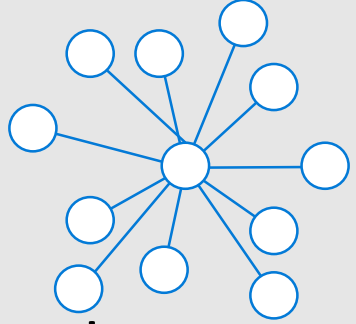
```
> pip install -U tf2onnx  
> python -m tf2onnx.convert --input frozen.pb --inputs X:0  
                                --outputs output:0 --output model.onnx
```

Learn more at

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





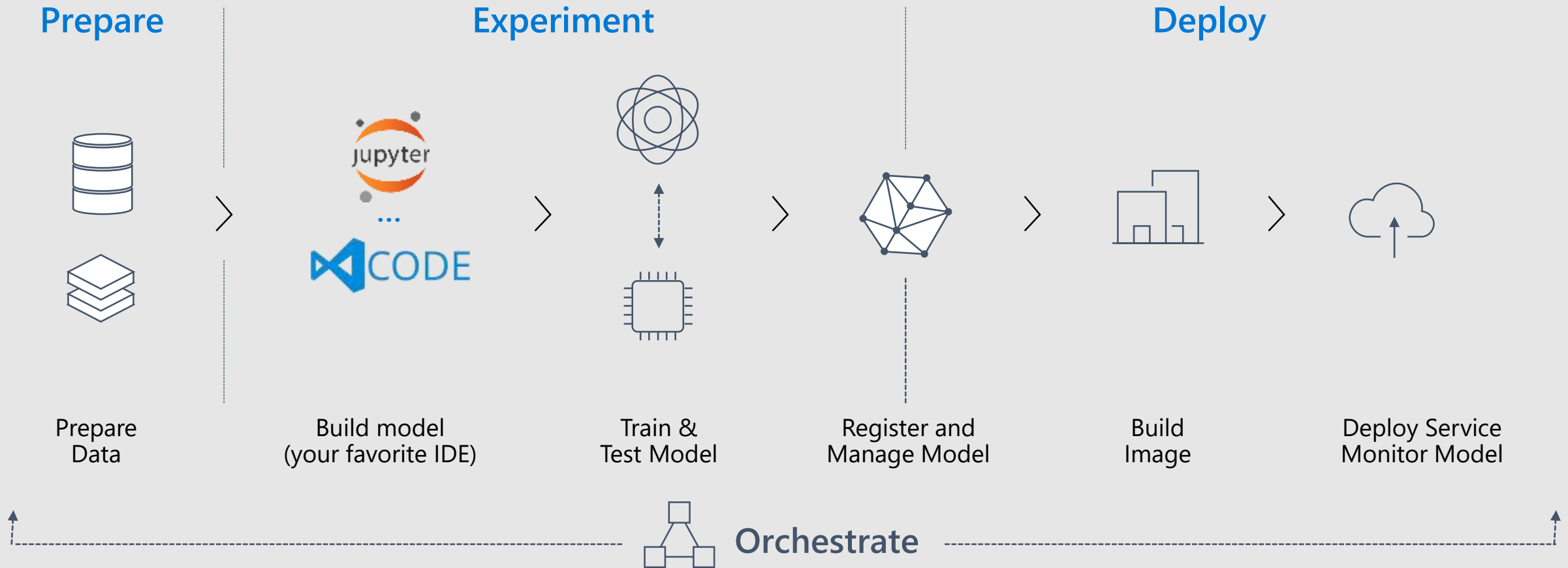


# Train models in Azure Machine Learning

- 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
- Clone Azure ML Gallery Samples
- <https://notebooks.azure.com/import/gh/Azure/MachineLearningNotebooks/>

# 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]

## Frameworks

Caffe2 Chainer Cognitive Toolkit

mxnet PyTorch PaddlePaddle

ML.NET Microsoft ML.NET 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


Windows Devices

**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

Native support

Windows Devices

Other Devices  
(iOS, etc)

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



# Deploy with Azure Machine Learning

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

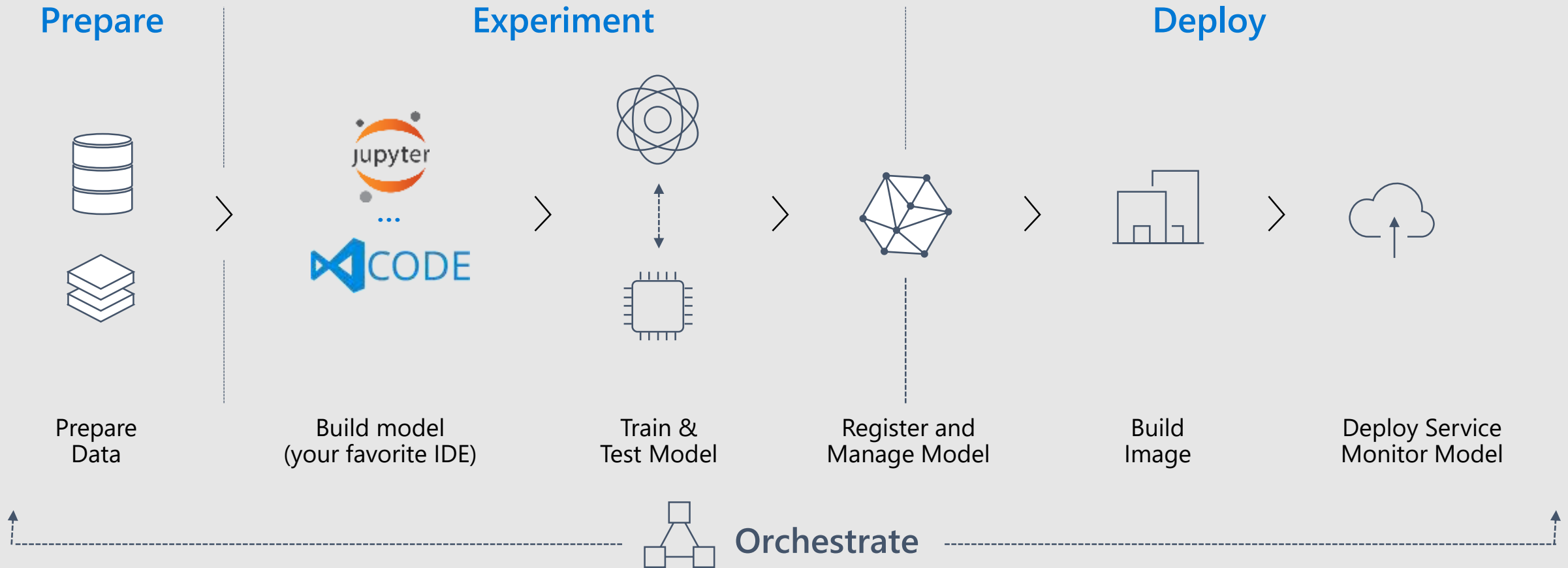


Azure  
Machine Learning

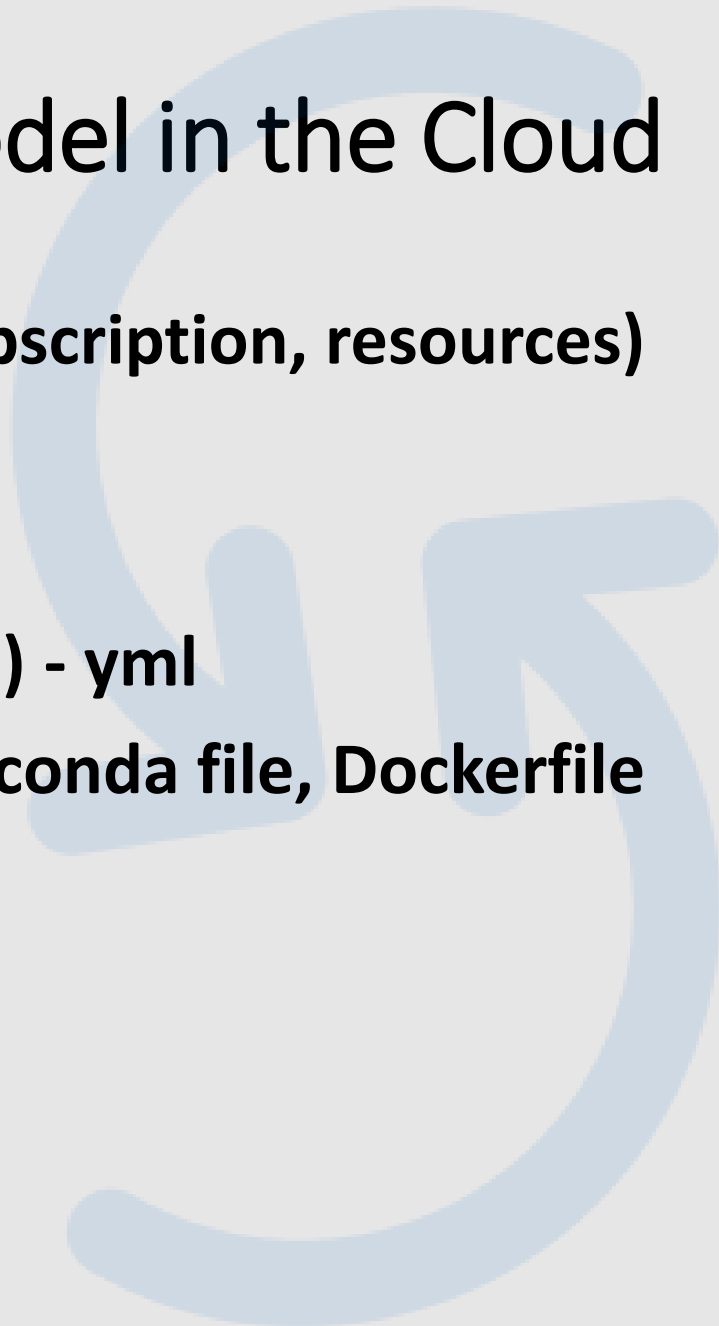


# Machine Learning

Typical E2E Process



# Deploy a VM with your ONNX model in the Cloud

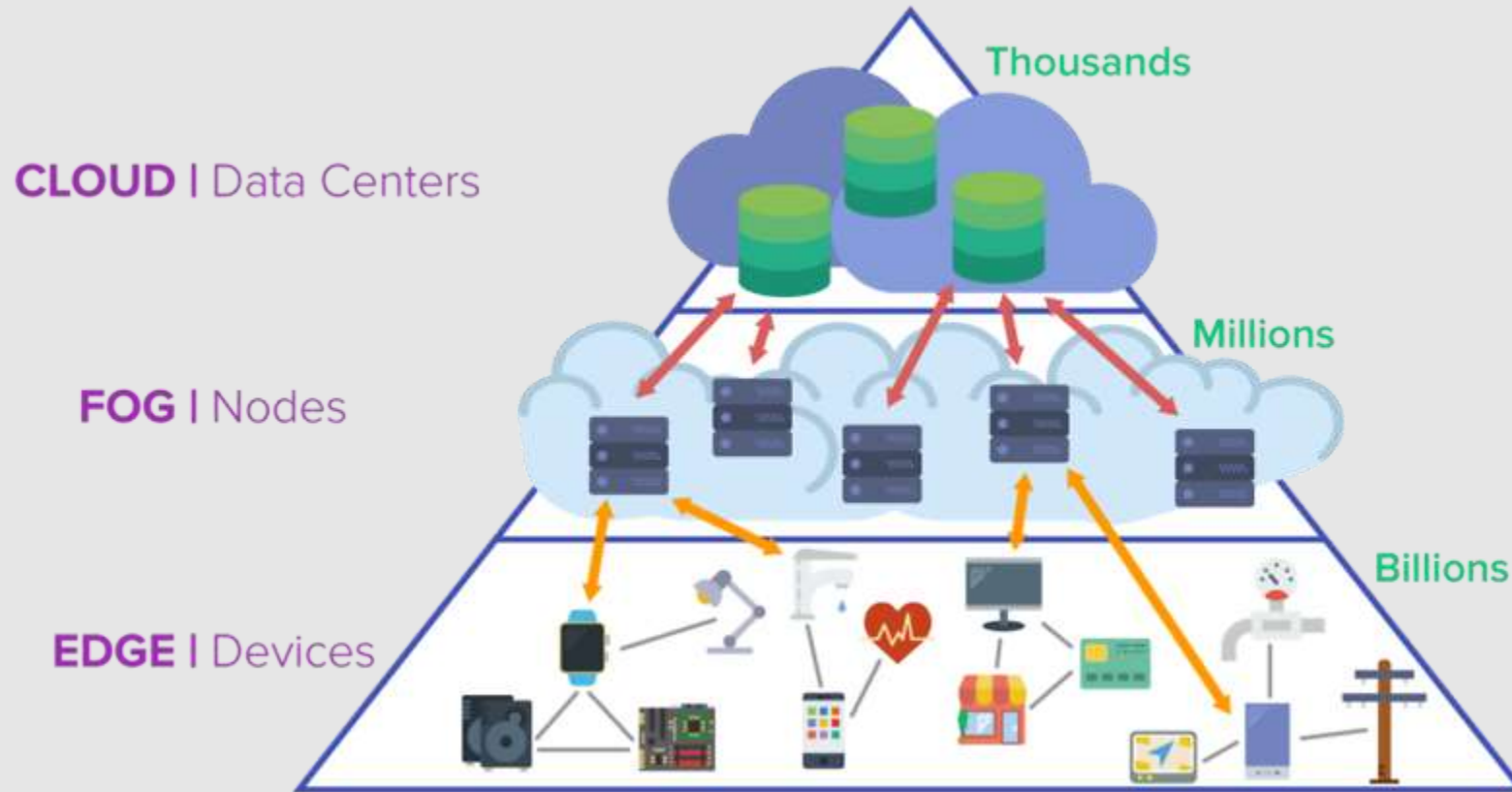
1. **Load Azure ML workspace – (config.json - subscription, resources)**
  2. **Registering your model with Azure ML**
  3. **Write Score File - python**
  4. **Write Environment File (conda dependencies) - yml**
  5. **Setup inference configuration – entry script, conda file, Dockerfile**
  6. **Deploy the model to ACI – type of cpu, mem**
- 



# Demo

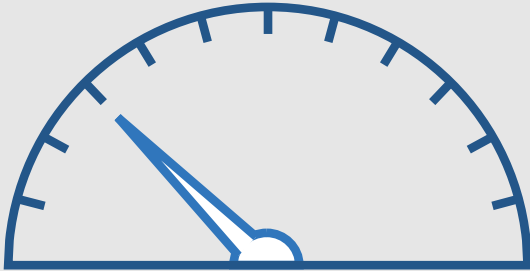
<https://github.com/rondagdag/onnx-pected/tree/master/ONNX-AML>

# What is the Edge?

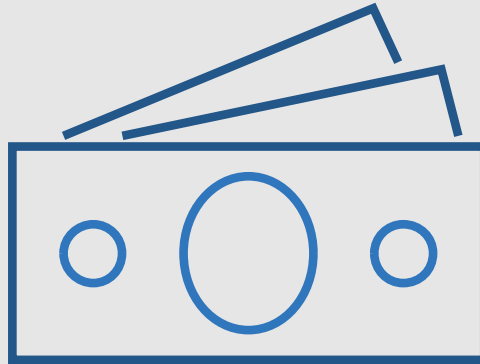


[Imagimob AB](#)

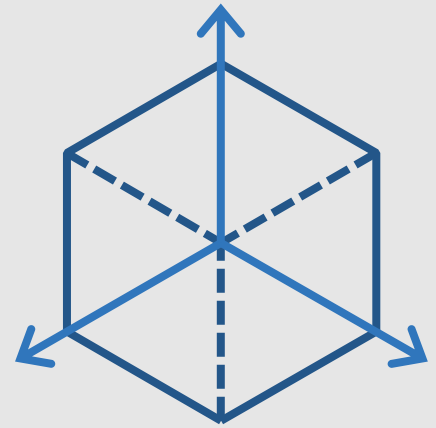
# AI on the edge



Low latency



Scalability



Flexibility

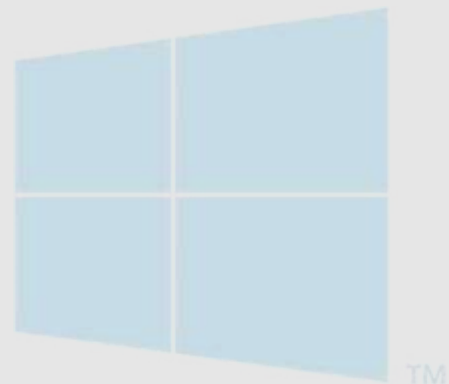
# ONNX as an intermediary format

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

# Deploy to Windows Devices

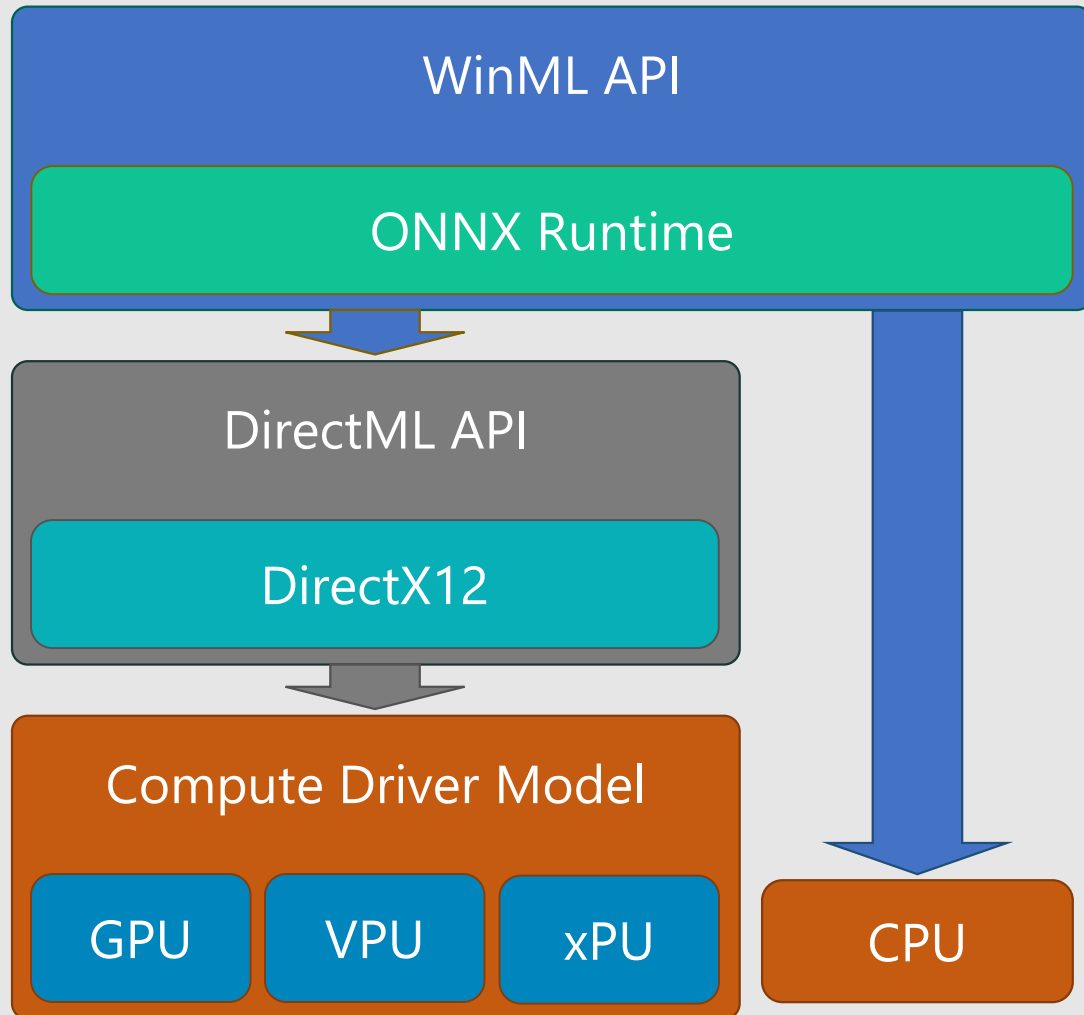
## Windows ML

- 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





# Demo

<https://github.com/rondagdag/onnx-pected/tree/master/GenerateONNX-AutoML>

# 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



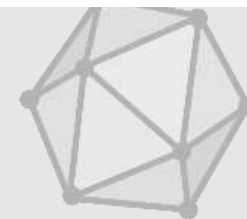
ONNX

# ONNX Runtime

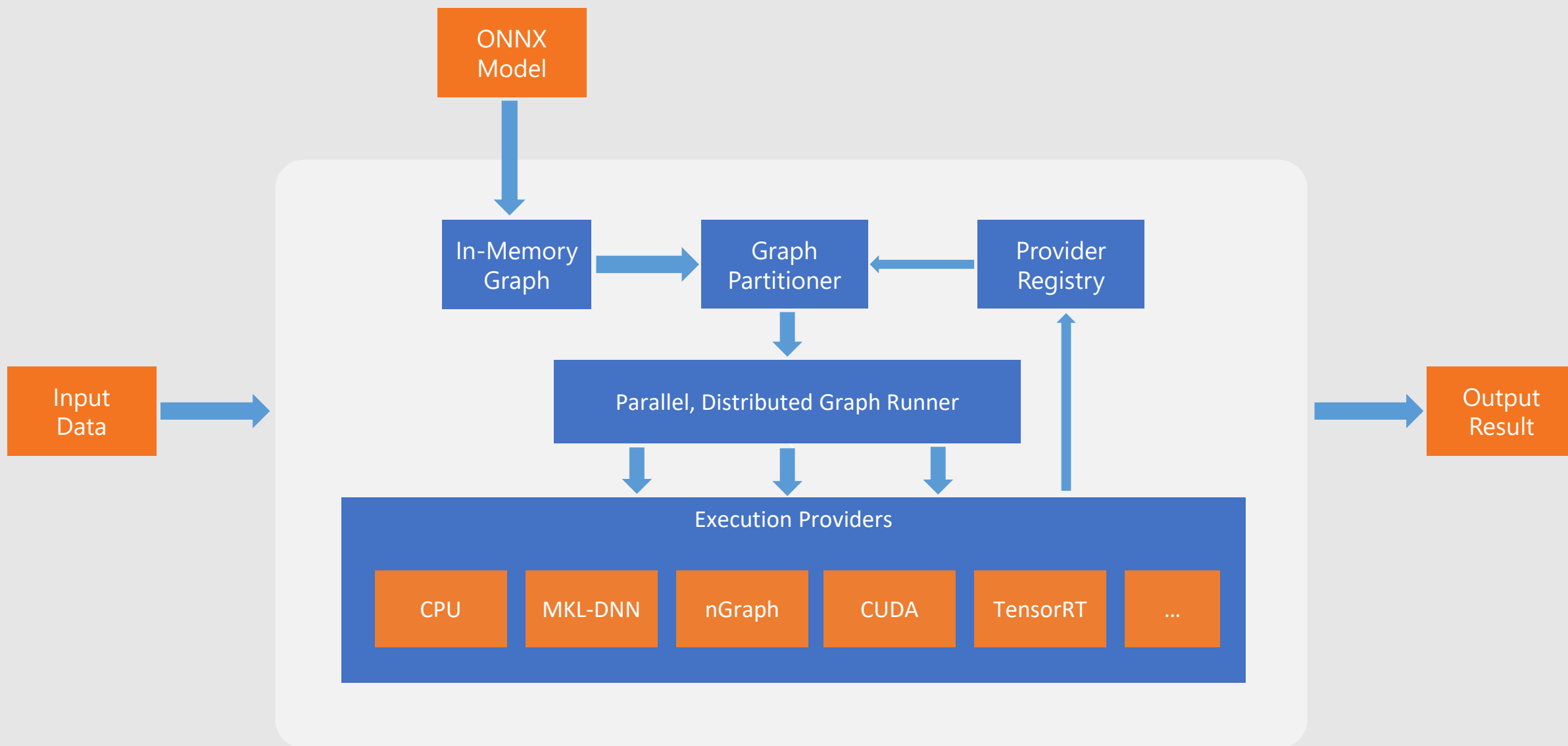
## Get Started Easily

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</a>				



# ONNX



# ONNX Runtime

- 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

# ONNX Docker Image

- [onnx-docker-cpu](#): Image with ONNX, PyTorch, Tensorflow support
- [onnx-docker-gpu](#): Image with ONNX, PyTorch (CUDA), Caffe2 support
- [onnx-ecosystem](#): Jupyter notebook environment for getting started quickly with ONNX models, ONNX converters, and inference using ONNX Runtime.

## ONNX Runtime Server

- **ONNX Runtime Server**
  - <https://github.com/onnx/tutorials/blob/master/tutorials/OnnxRuntimeServerSSDModel.ipynb>



ONNX

# Reference implementation to use ONNX Runtime with Azure IoT Edge



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





# 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

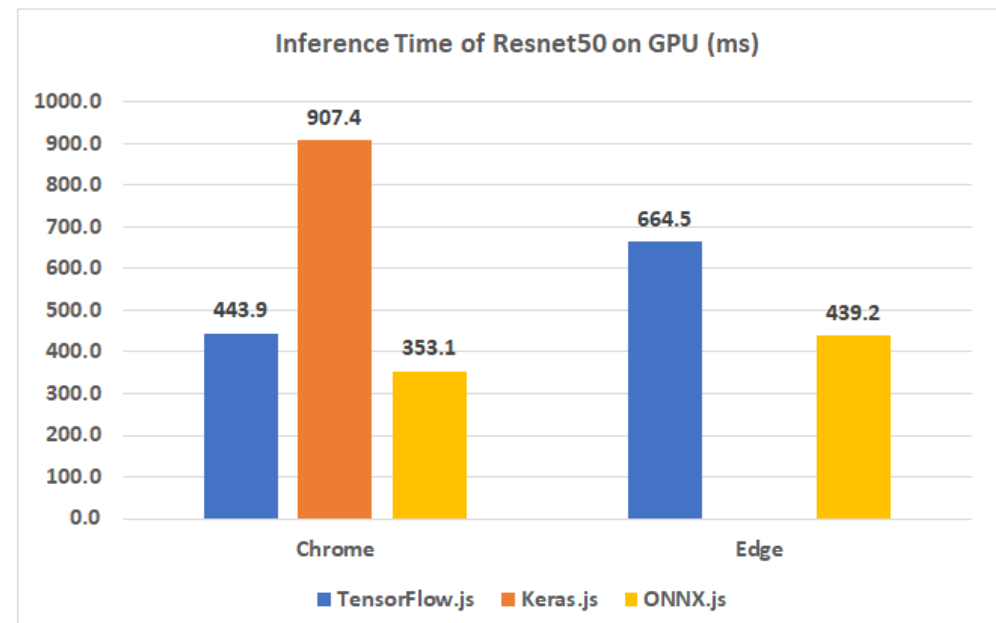
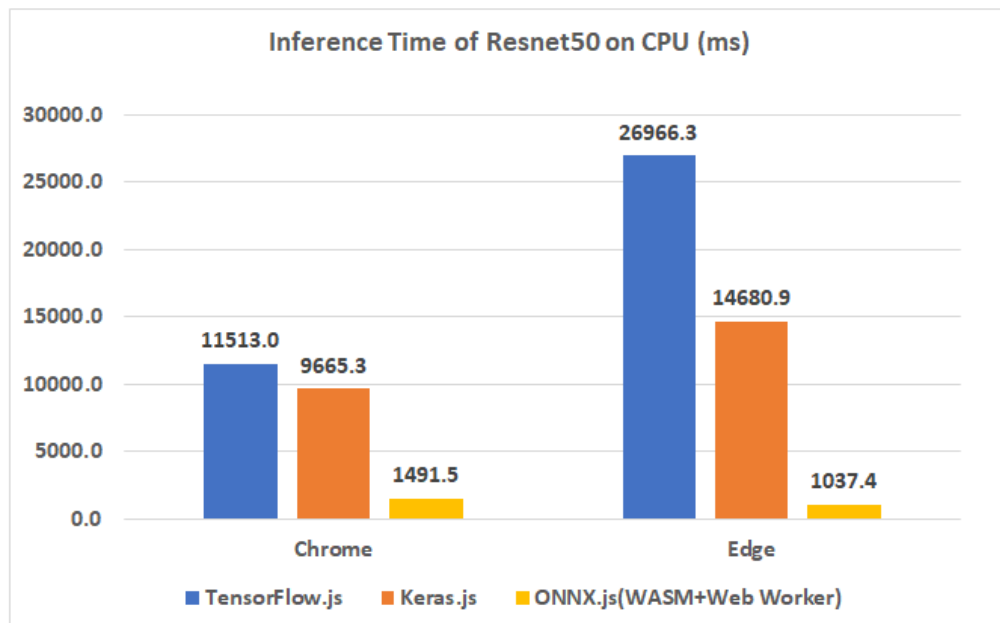
### Desktop Platforms

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



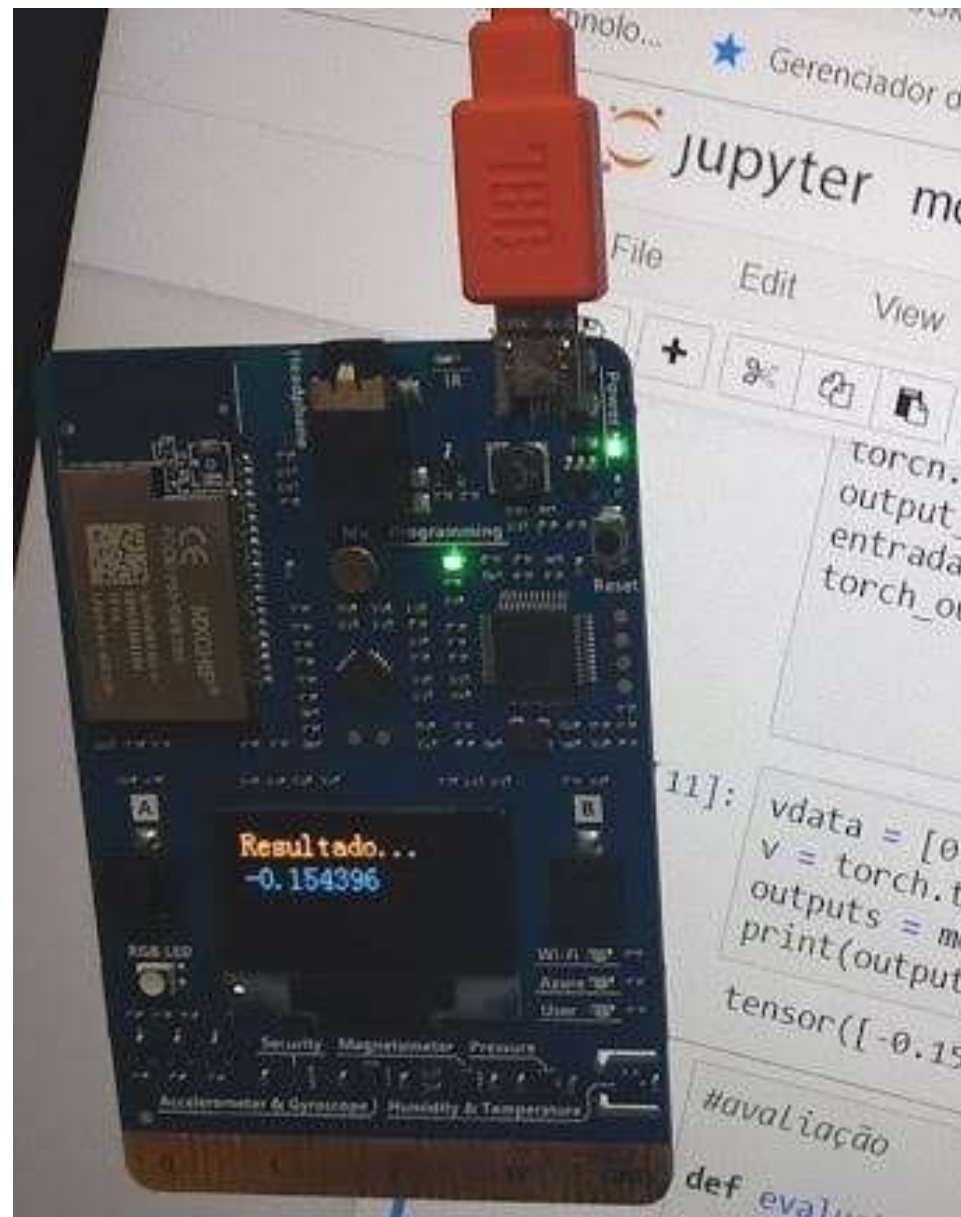
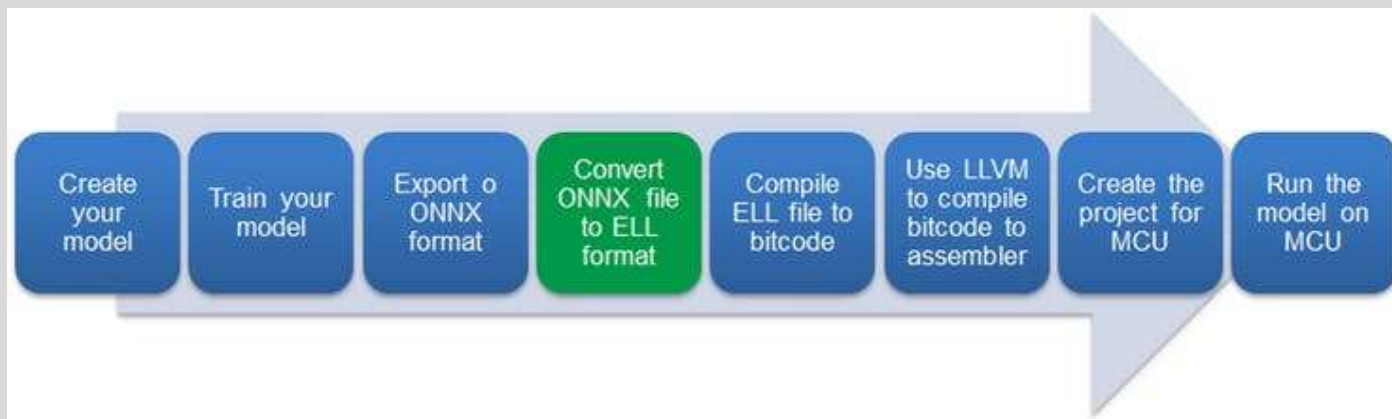


# Demo

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

# 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/onnxpected>

# Ron Dagdag

- Sr. Software Engineer/Voice AI Assistant Specialist at Crestron Electronics
- Microsoft MVP award
- Hackster DFW Ambassador [meetup.com/Hackster-DFW](https://meetup.com/Hackster-DFW)
- Dallas Littlebits Chapter Leader [meetup.com/amRobotics](https://meetup.com/amRobotics)
- Dallas AR/VR Development meetup [meetup.com/Dallas-Virtual-Reality](https://meetup.com/Dallas-Virtual-Reality)
- *"Opinions expressed are solely my own and do not express the views or opinions of my employer."*



**Ron Lyle Dagdag**

Immersive Experience Developer

Cell: 682-560-3988

ron@dagdag.net



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

@rondagdag

<http://ron.dagdag.net>

Experience AR