



# THAT<sup>®</sup>

## CONFERENCE



# THANK YOU, THAT CONFERENCE PARTNERS!

***Unspecified***

SOFTWARE CO



#RonDagdag



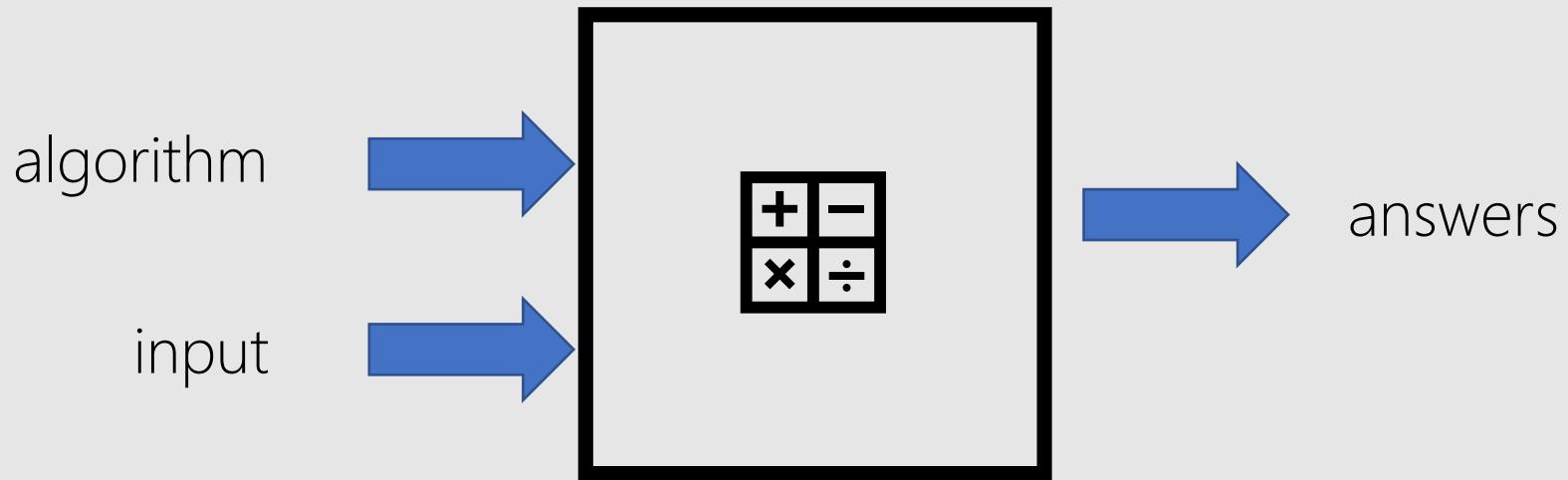
# RON DAGDAG

Director of Software Engineering  
and Microsoft MVP

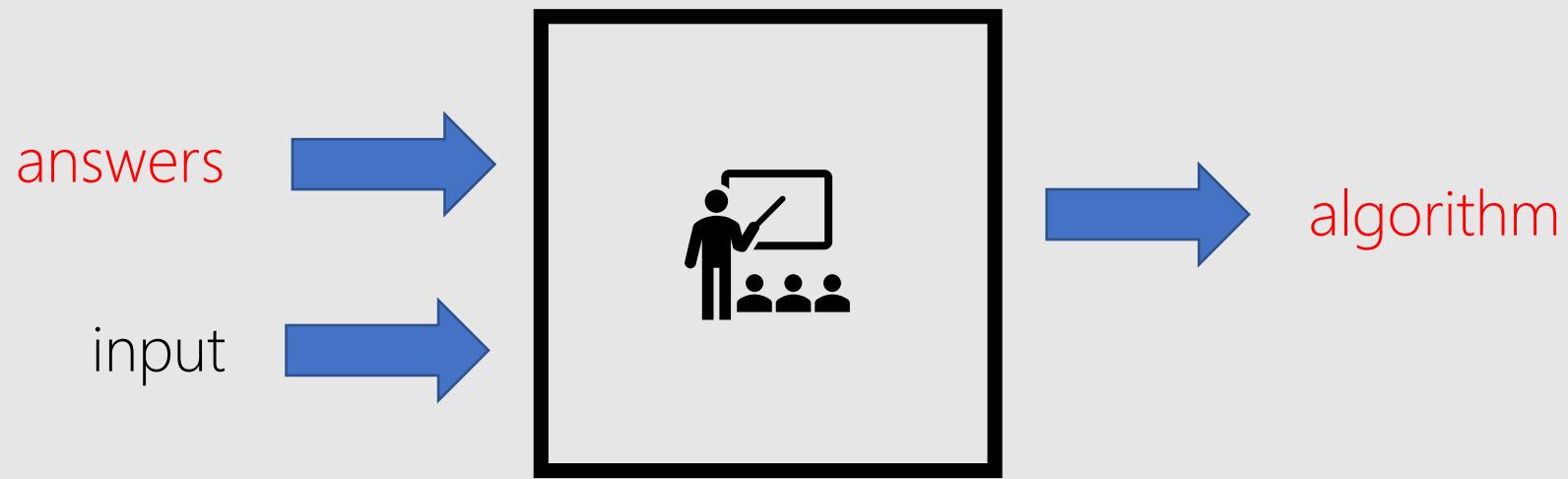
ONNX  
Not ONIX  
Not ONYX



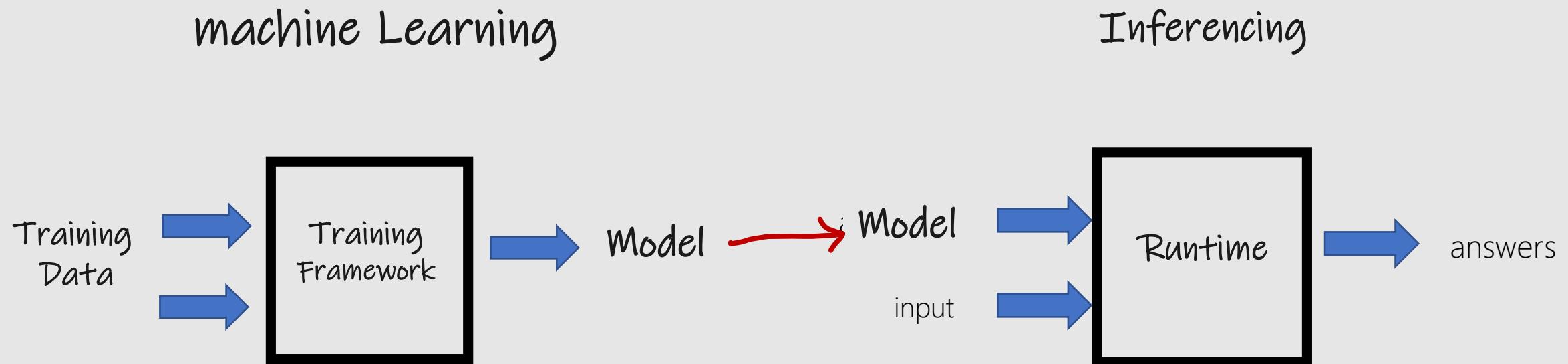
# programming



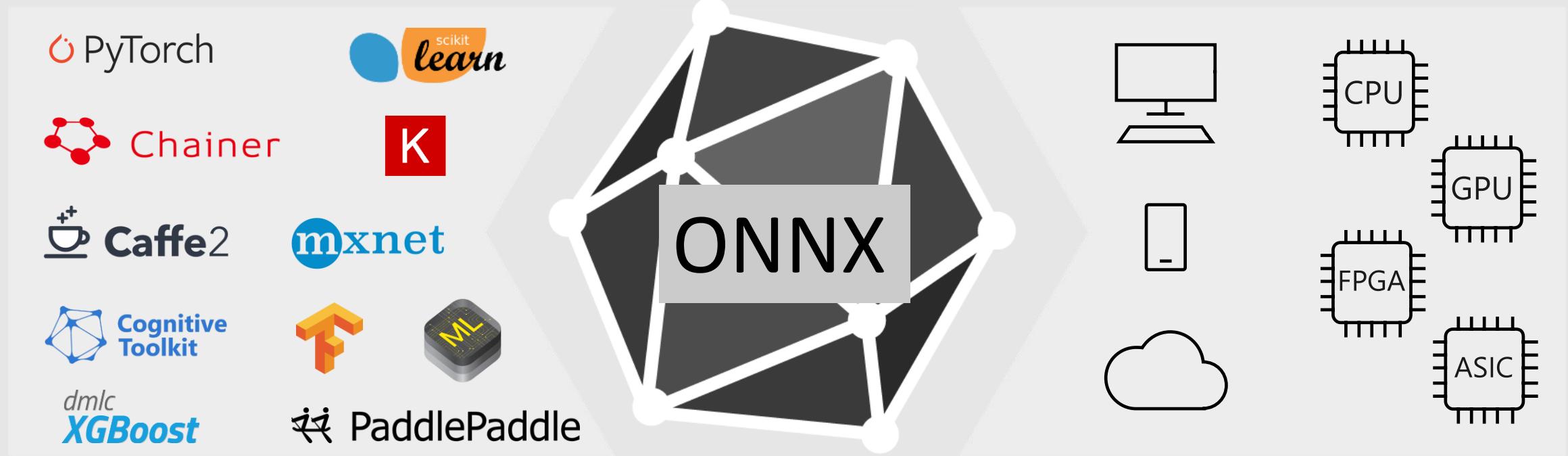
# machine learning

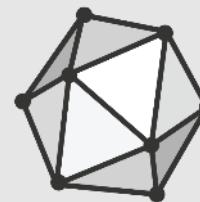


# ML Primer



# 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

**ABBYY**

 Alibaba Group  
阿里巴巴集团

**AMD**

**arm**

 aws

 Baidu 百度

**BECKHOFF**

**BITMAIN**



**CEVA**

 Facebook  
Open Source

**GRAPHCORE**





 Hewlett Packard  
Enterprise

 HUAWEI



 Idein Inc



 MathWorks<sup>®</sup>







 Microsoft

 NVIDIA.



 OctoML



**OPEN AI LAB**  
开放智能

 Preferred  
Networks



 SONY



























12.5k  
Github Stars



2100+  
Pull  
Requests

220+  
Contributors



2.7k Github  
Forks



40+ Model  
Zoo



**DLF AI**  
GRADUATE  
PROJECT

# When to use ONNX?

- Trained in Python or ML.NET - deploy into a C#/Java/Javascript app
- High Inferencing latency for production use
- Model to run resource on IoT/edge devices
- Model to run on different OS or Hardware
- Combine running models created from different frameworks
- Training takes too long (transformer models)

# Agenda

- ✓ What is ONNX, When to use ONNX
- ❑ How to create ONNX models
- ❑ How to deploy ONNX models

# Create

## Frameworks



## Services



Native support

Converters

Native support

## ONNX Model



# Deploy

## Cloud Services



## Windows Devices

## IoT/Edge Devices

Native support

Converters

## Other Devices (iOS, Android, etc)

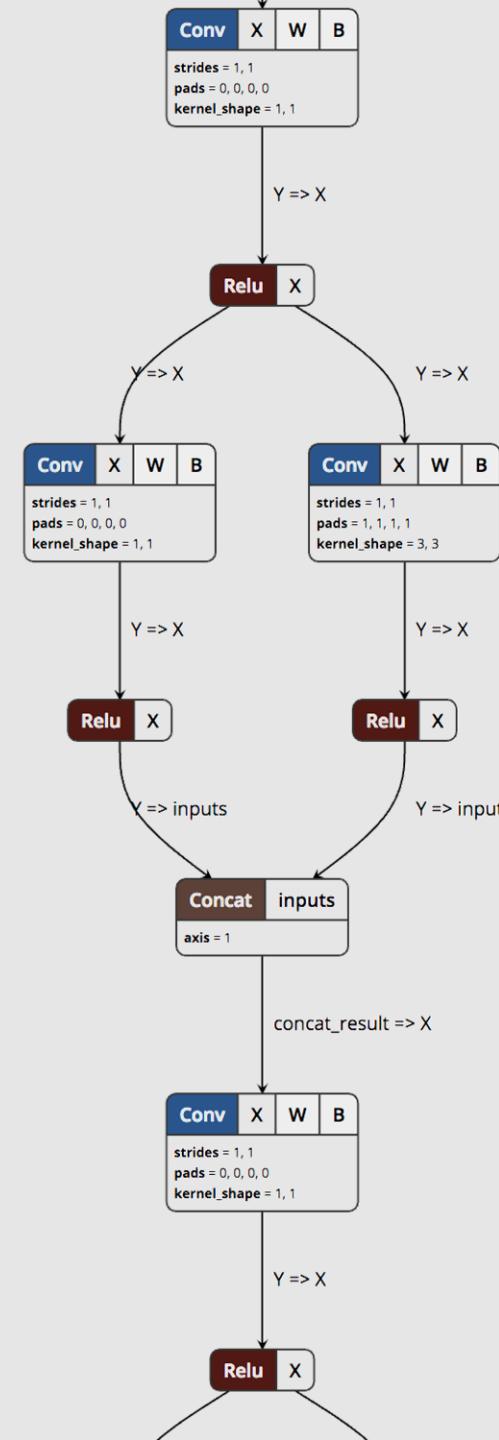
# ONNX Models

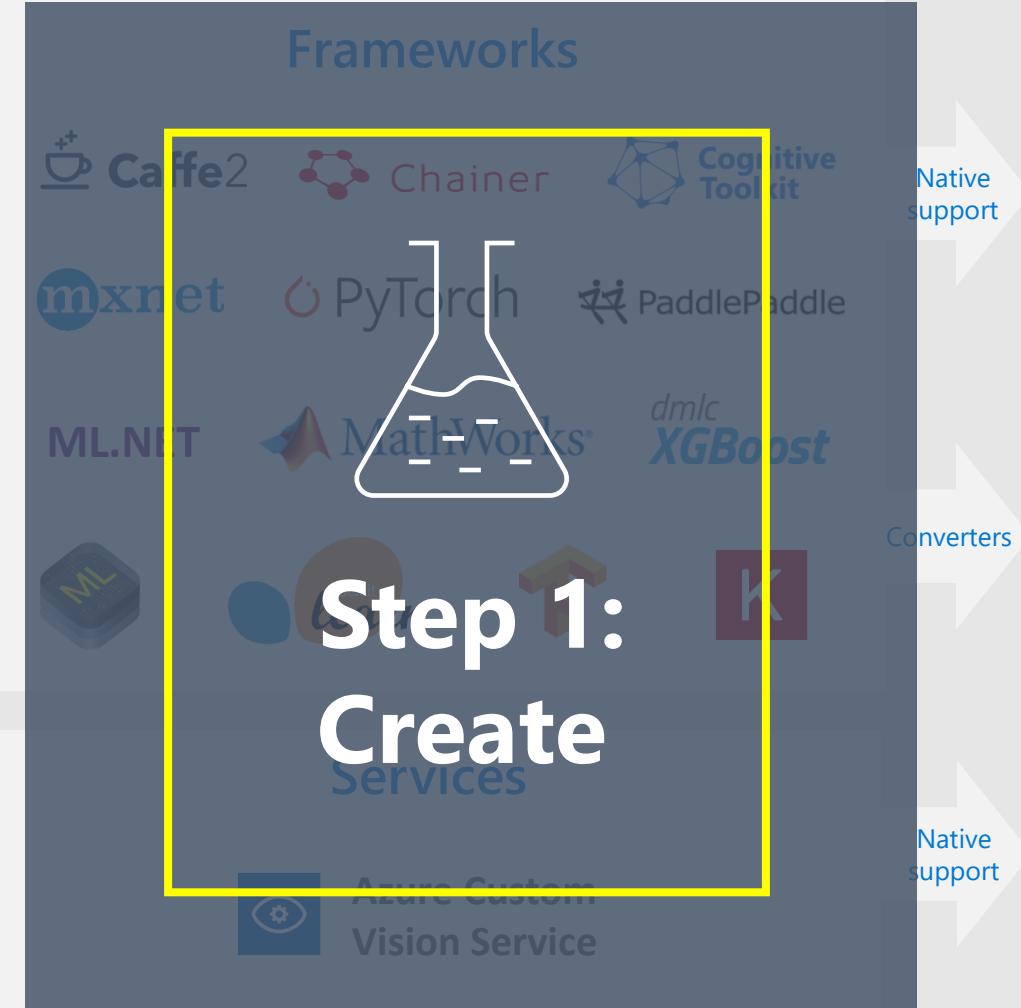
## Graph of operations

Netron

<https://netron.app/>

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





ONNX Model





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.

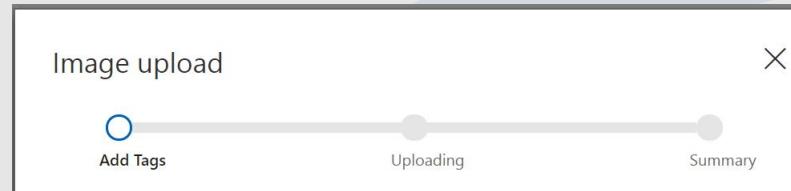
Model Class	Reference	Description
MobileNet	Sandler et al.	Efficient CNN model for mobile and embedded vision applications. Top-5 error from paper - ~10%
ResNet	He et al., He et al.	Very deep CNN model (up to 152 layers), won the ImageNet Challenge in 2015. Top-5 error from paper - ~3.6%
SqueezeNet	Iandola et al.	A light network with fewer parameters. Top-5 error from paper - ~1.2%
VGG	Simonyan et al.	Deep Convolutional Neural Networks for Visual Recognition. Top-5 error from paper - ~16.4%

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

A screenshot of the 'Training Images' tab. It shows a strawberry image with the text '4 images will be used'. Below it is a 'My Tags' section with a text input field containing 'fruit' and a delete button. To the right are 'Delete' and 'Export' buttons, with 'Export' highlighted by a red box. The 'Performance' tab is also visible.

3. Download ONNX model!

A screenshot of the 'Performance' tab. It features a large 'Choose your platform' button. Below it is a large grey circle with the word 'ONNX' repeated twice. At the bottom of the page, there's a footer with the text 'GET STARTED' and 'TRY IT FREE'.

# Convert models



Keras



NCNN



# Convert models

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



# Convert models: PyTorch

```
import torch  
import torch.onnx  
  
model = torch.load("model.pt")  
  
sample_input = torch.randn(1, 3, 224, 224)  
  
torch.onnx.export(model, sample_input, "model.onnx")
```

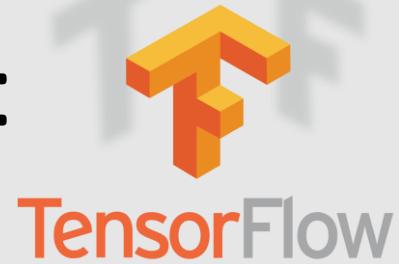
# Convert models: Keras

```
In [ ]: import onnxmltools  
from keras.models import load_model
```

```
In [ ]: # Update the input name and path for your Keras model  
input_keras_model = 'model.h5'  
  
# Change this path to the output name and path for the ONNX model  
output_onnx_model = 'model.onnx'
```

```
In [ ]: # Load your Keras model  
keras_model = load_model(input_keras_model)  
  
# Convert the Keras model into ONNX  
onnx_model = onnxmltools.convert_keras(keras_model)  
  
# Save as protobuf  
onnxmltools.utils.save_model(onnx_model, output_onnx_model)
```

# Convert models:



```
> python -m tf2onnx.convert  
    --saved-model tensorflow-model-path  
    --output model.onnx
```

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

# Convert models:

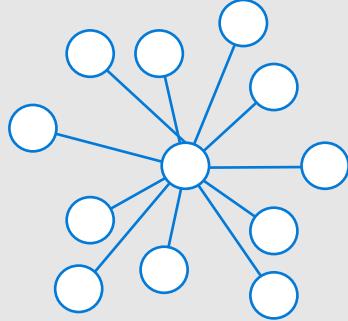


```
# Train a model.  
from sklearn.datasets import load_iris  
from sklearn.model_selection import train_test_split  
from sklearn.ensemble import RandomForestClassifier  
iris = load_iris()  
X, y = iris.data, iris.target  
X_train, X_test, y_train, y_test = train_test_split(X, y)  
clr = RandomForestClassifier()  
clr.fit(X_train, y_train)  
  
# Convert into ONNX format  
from skl2onnx import convert_sklearn  
from skl2onnx.common.data_types import FloatTensorType  
initial_type = [('float_input', FloatTensorType([None, 4]))]  
onx = convert_sklearn(clr, initial_types=initial_type)  
with open("rf_iris.onnx", "wb") as f:  
    f.write(onx.SerializeToString())
```

A large, light blue, abstract graphic consisting of several thick, rounded, horizontal strokes that curve and overlap each other, creating a sense of motion or a stylized logo.

<https://github.com/onnx/tutorials>

# 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

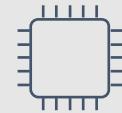
Typical E2E Process

## Prepare



Prepare Data

## Experiment



Build model  
(your favorite IDE)

Train &  
Test Model



Register and  
Manage Model

## Deploy



Deploy Service  
Monitor Model

## Orchestrate



**ONNX Model**



A black and white photograph showing a baker from the waist down, wearing a white apron over a dark shirt. The baker is standing at a light-colored wooden counter, using their hands to knead a large, round loaf of bread. The dough has a textured, slightly cracked surface. The background is a plain, light-colored wall.

# Baker vs Starting a Bakery

# Create

## Frameworks

 Caffe2  Chainer  Cognitive Toolkit

 mxnet  PyTorch  PaddlePaddle

 ML.NET  MathWorks  XGBoost

## Services

 Azure Custom Vision Service

Native support

Converters

Native support

## ONNX Model



# Deploy

## Azure

Azure Machine Learning services  
Ubuntu VM  
Windows Server 2019 VM

Windows/Linux Devices

IoT Edge Devices

Other Devices  
(iOS, etc)

Native support

Converters

A close-up photograph of a person's hands holding a large, round loaf of bread. The bread has a thick, cracked, and golden-brown crust. The person is wearing a blue and white striped cloth around their wrists. The background is blurred.

Cloud  
or  
Edge

# Deploy with Azure Machine Learning

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



Azure  
Machine Learning

# Machine Learning

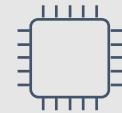
Typical E2E Process

## Prepare



Prepare Data

## Experiment



Build model  
(your favorite IDE)

Train &  
Test Model



Register and  
Manage Model

## Deploy



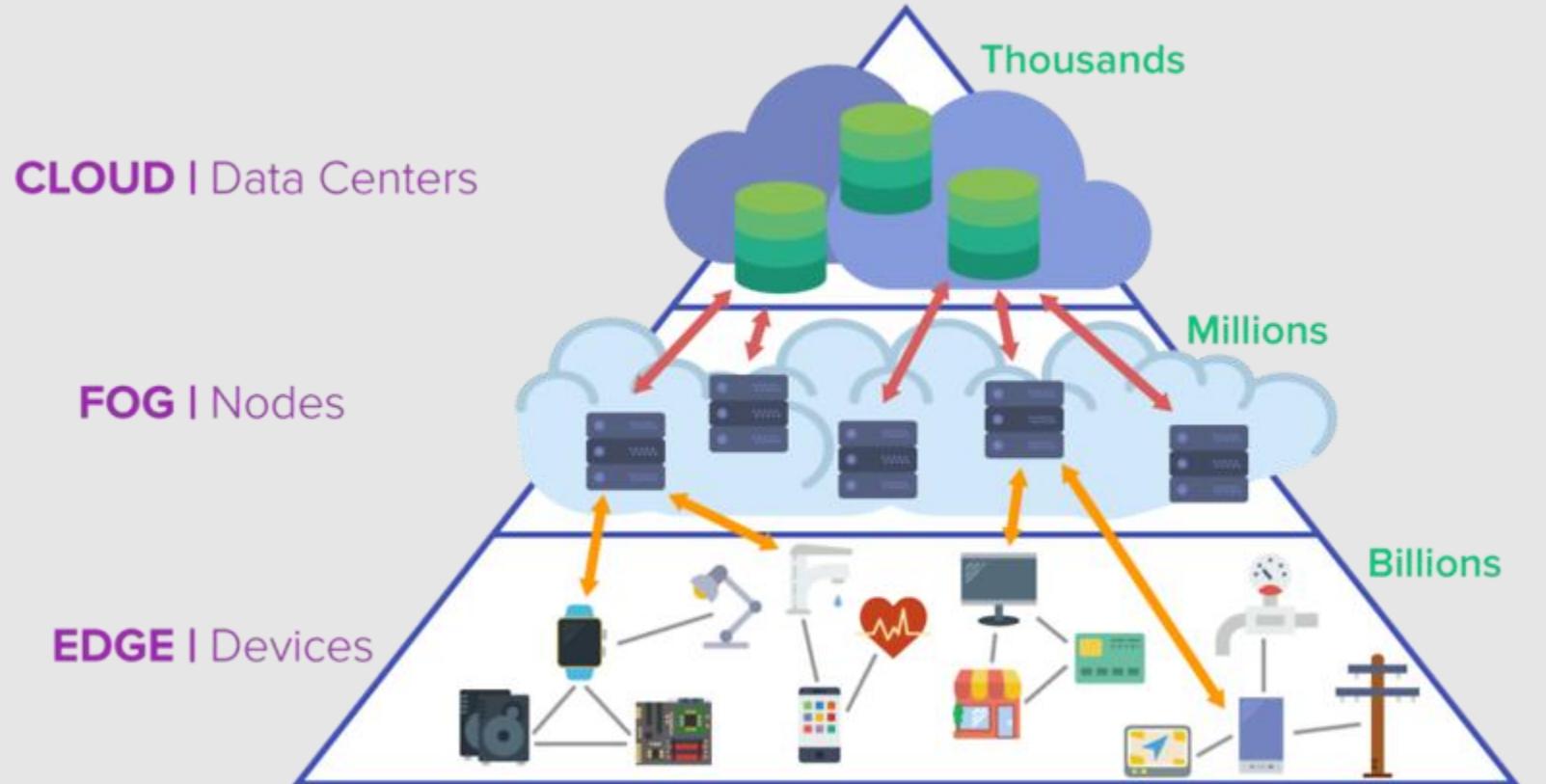
Deploy Service  
Monitor Model



Orchestrate



# What is the Edge?



[Imagimob AB](#)

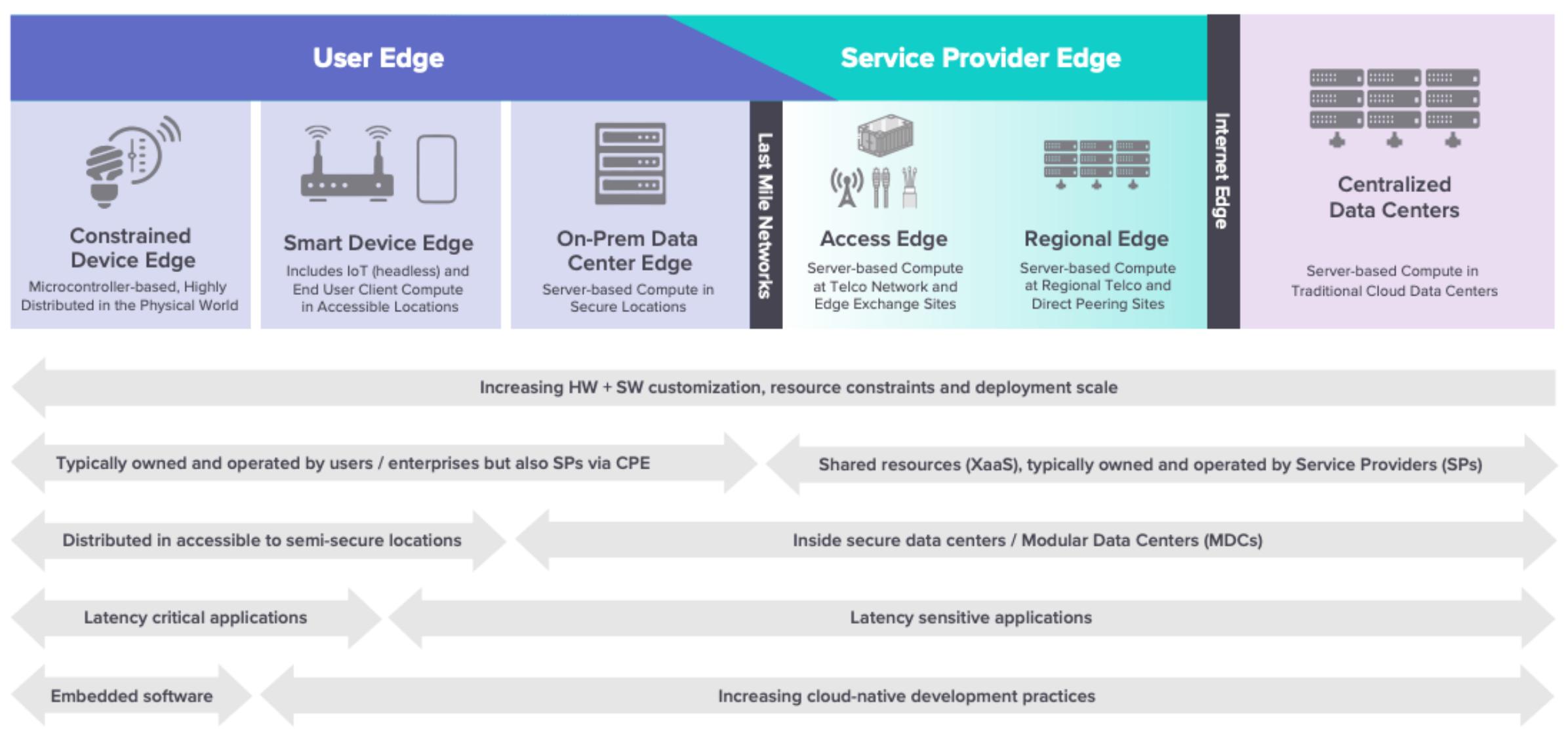


Image source: [LF Edge](#)

# ONNX Runtime

- High performance inference engine for ONNX models
- Founded and Open Sourced by Microsoft under MIT License
- Supports full ONNX-ML spec
- Extensible architecture to plug-in hardware accelerators
- Ships with Windows 10 as WinML
- [onnxruntime.ai](http://onnxruntime.ai)



ONNX

# ONNX Runtime

## Get Started Easily

Optimize Inferencing

Optimize Training

Platform

Windows

Linux

Mac

Android

iOS

Web Browser  
(Preview)

API

Python

C++

C#

C

Java

JS

Obj-C

WinRT

Architecture

X64

X86

ARM64

ARM32

IBM Power

Hardware Acceleration

Default CPU

CUDA

DirectML

oneDNN

OpenVINO

TensorRT

NNAPI

ACL (Preview)

ArmNN  
(Preview)

CoreML  
(Preview)

MIGraphX  
(Preview)

NUPHAR  
(Preview)

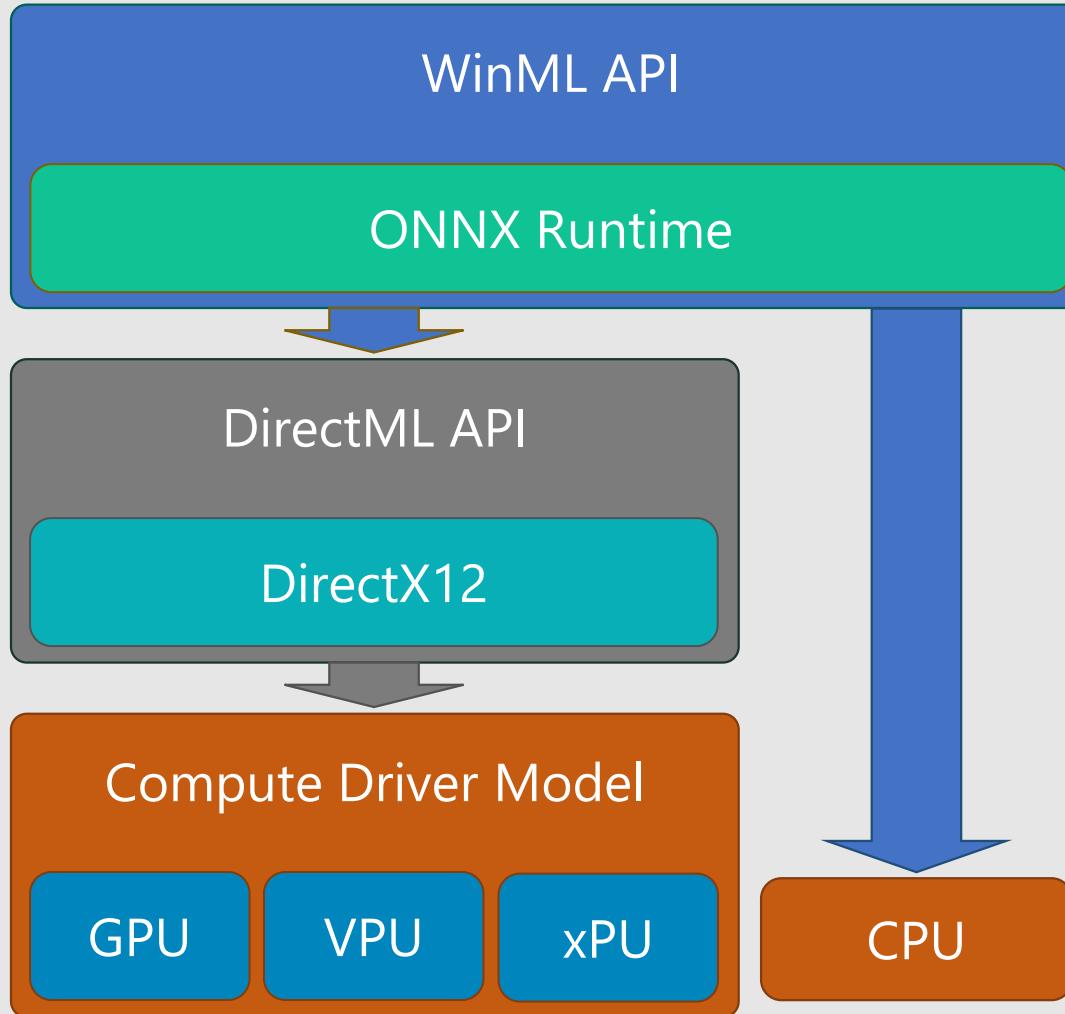
Rockchip NPU  
(Preview)

Vitis AI (Preview)

Installation Instructions

Install Nuget package [Microsoft.ML.OnnxRuntime.Gpu](#)  
Refer to [docs](#) for requirements.

# 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  
JavaScript

---

Node.js binding

---

Web

---

React Native

# ONNX Runtime Node.js

- Node.js binding
- ONNX model inferencing
- Electron
- Uses web assembly

## Install

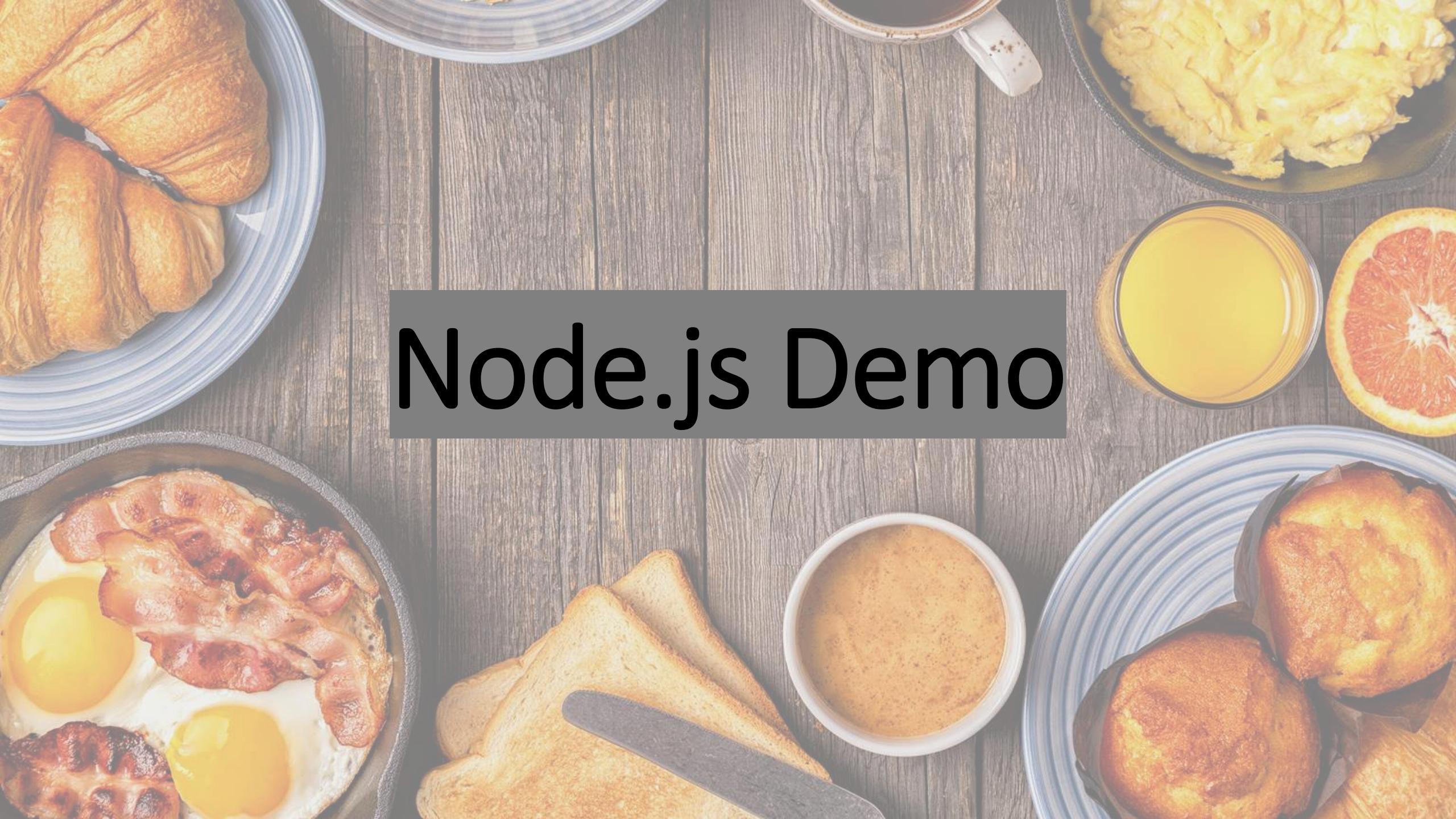
```
# install latest release version  
npm install onnxruntime-node
```

## Import

```
// use ES6 style import syntax (recommended)  
import * as ort from 'onnxruntime-node';
```

```
// or use CommonJS style import syntax  
const ort = require('onnxruntime-node');
```

# Node.js Demo



# ONNX Runtime Web (ORT-Web)

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

## Install

```
# install latest release version  
npm install onnxruntime-web
```

```
# install nightly build dev version  
npm install onnxruntime-web@dev
```

## Import

```
// use ES6 style import syntax (recommended)  
import * as ort from 'onnxruntime-web';
```

```
// or use CommonJS style import syntax  
const ort = require('onnxruntime-web');
```

# Why inference in the browser



**It's faster**



**It's safer** and helps with  
privacy



**It works offline**

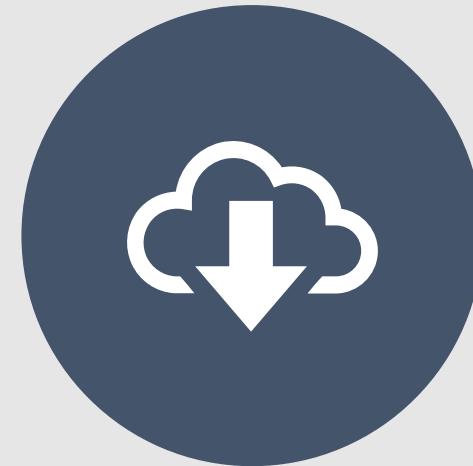


**It's cheaper**

# Why not in the browser?



THE MODEL IS TOO LARGE AND  
REQUIRES HIGHER HARDWARE SPECS.



DOWNLOADED ONTO THE DEVICE

A close-up photograph of a young girl with light brown hair and blue eyes. She is wearing a red long-sleeved shirt. She is holding a gingerbread man cookie in front of her face. The cookie has white icing for a bow tie and a smile. The background is slightly blurred.

# Web Browser Demo

# React Native

- score pre-trained ONNX models
- ONNX Runtime Mobile
- light-weight inference solution
- Android and iOS

## Install

```
# install latest release version  
npm install onnxruntime-react-native
```

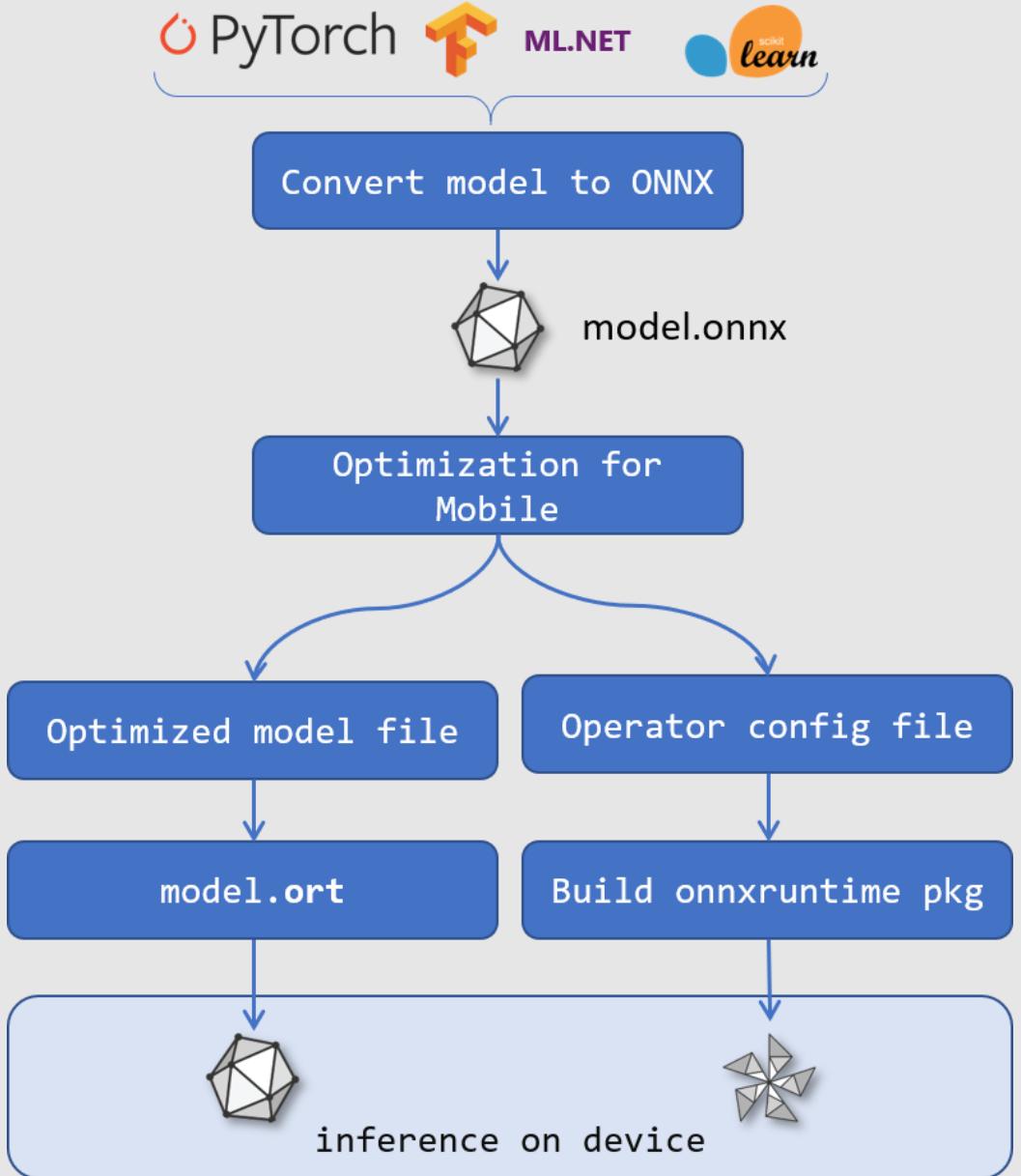
## Import

```
// use ES6 style import syntax (recommended)  
import * as ort from 'onnxruntime-react-native';
```

```
// or use CommonJS style import syntax  
const ort = require('onnxruntime-react-native');
```

# ONNX Runtime Mobile

- minimizes the binary size
- pre-optimized ONNX model to an internal format ('ORT format model')

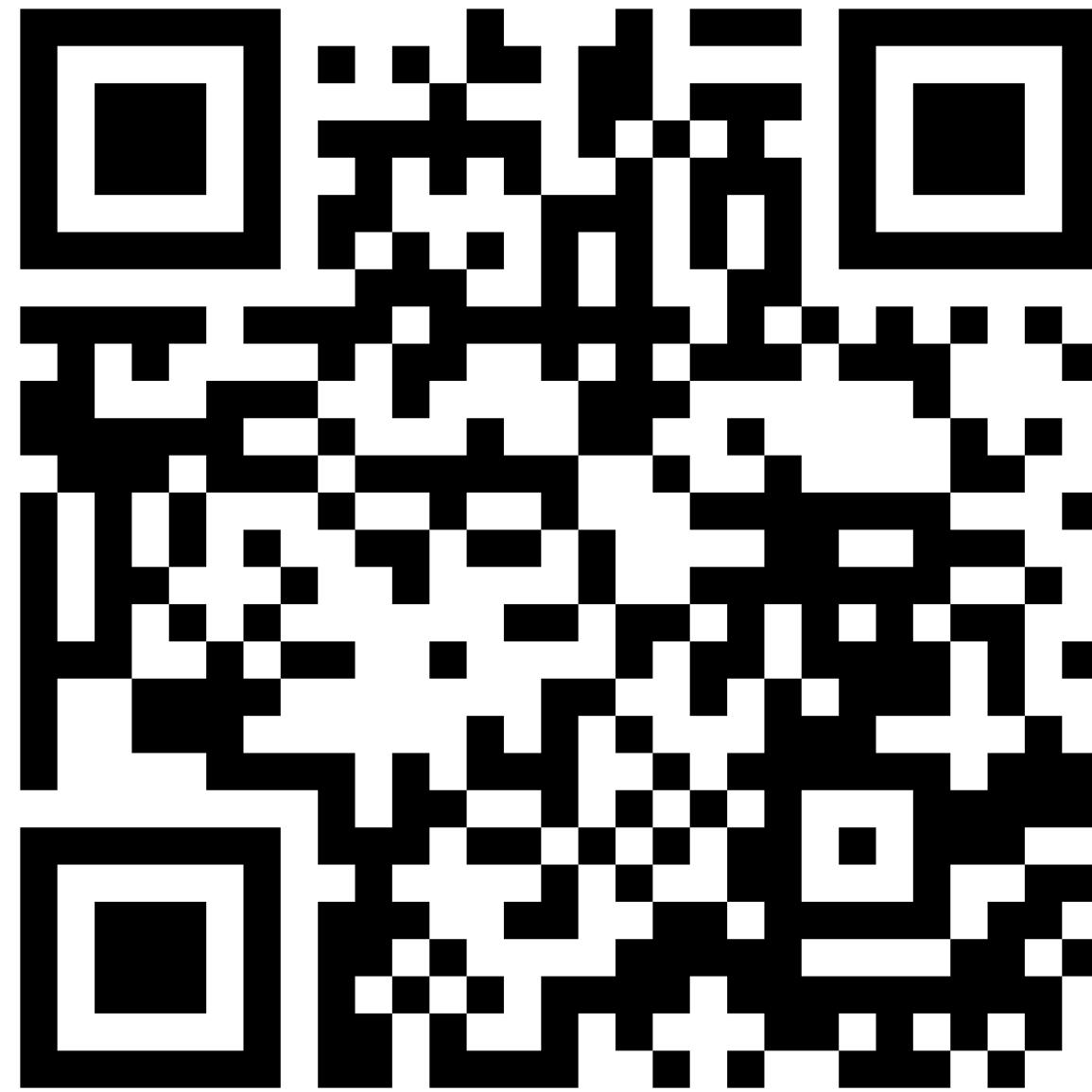


# Compatibility Chart

## Compatibility

OS/Browser	Chrome	Edge	Safari	Electron	Node.js
Windows 10	wasm, webgl	wasm, webgl	-	wasm, webgl	wasm
macOS	wasm, webgl	wasm, webgl	wasm, webgl	wasm, webgl	wasm
Ubuntu LTS 18.04	wasm, webgl	wasm, webgl	-	wasm, webgl	wasm
iOS	wasm, webgl	wasm, webgl	wasm, webgl	-	-
Android	wasm, webgl	wasm, webgl	-	-	-





<https://github.com/rondagdag/onnx-web-presentation>

# 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**

# About Me

Ron Dagdag



**Ron Lyle Dagdag**

Immersive Experience Developer

Cell: 682-560-3988

ron@dagdag.net



Experience AR

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

<http://ron.dagdag.net>

<https://linktr.ee/rondagdag>

Director of Software Engineering at Spaceee

5<sup>th</sup> year Microsoft MVP awardee

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

Email: [ron@dagdag.net](mailto:ron@dagdag.net)  
Twitter @rondagdag

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

Thanks for geeking out with me about ONNX



SEE YOU NEXT YEAR! JANUARY 2023

CALL FOR SPEAKERS STARTS JUNE 1, 2022