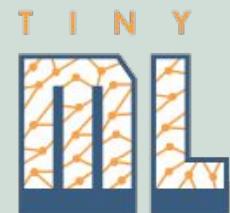


# TinyML Challenges



Hardware



Software

# Compute



# Memory

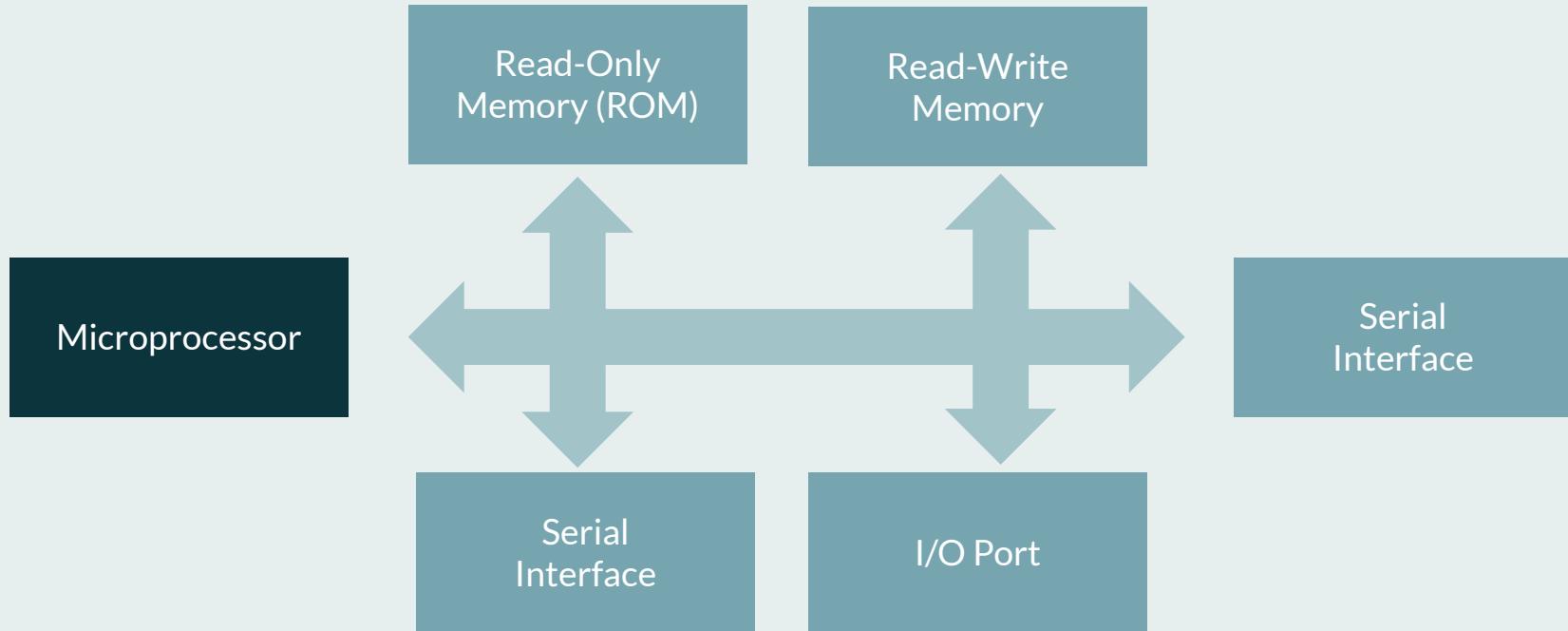


# Storage

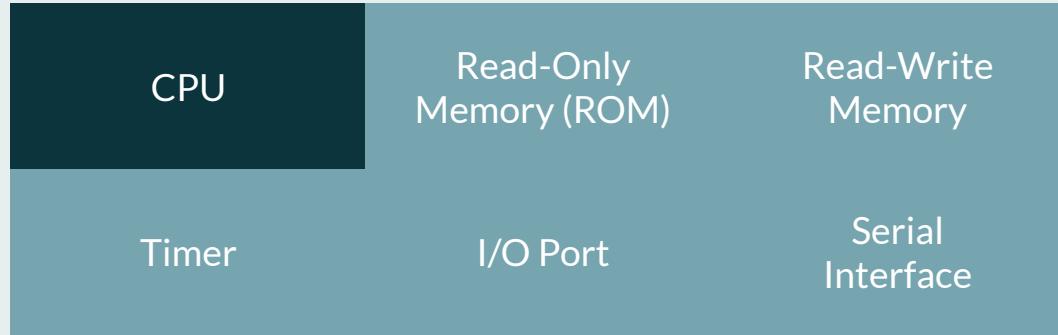


# Microprocessor Vs Microcontroller

# Microprocessor: only **one part** of the puzzle



# Microcontroller



# Microprocessor

- Heart of a **computer system**
- Just the processor, memory and storage are **external**
- Mainly used in general **purpose systems** like laptops, desktops and server
- **Offers flexibility** in design
- System size is **big**

# Microcontroller

- Heart of an **embedded system**
- Memory and storage are all **internal** to the system
- Mainly used in **specialized, fixed function system** like phones, MP3 player, etc
- **Limited flexibility** in design
- System size is **tiny**

# Orders of Magnitude Difference

	Microprocessor	>	Microcontroller
Platform			
Compute	1GHz - 4GHz	10x	1MHz - 400MHz
Memory	512MB - 64GB	10Kx	2KB - 512KB
Storage	64GB - 4TB	100Kx	32KB - 2MB
Power	30W - 100W	1Kx	150µW - 23.5mW

Microcontroller



1MHz - 400MHz

2KB - 512KB

32KB - 2MB

150 $\mu$ W - 23.5mW

# Implications

- How complicated is the running task?
- How much memory does it need to have?
- How long does the job have to perform?

Hardware

Software



# Software

Applications

Libraries

Operating System

Hardware

# Widely Used Operating Systems



Mobile OS

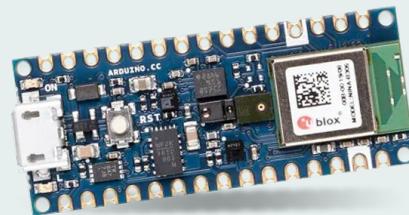


# The Android Software Stack

# Widely Used Operating Systems



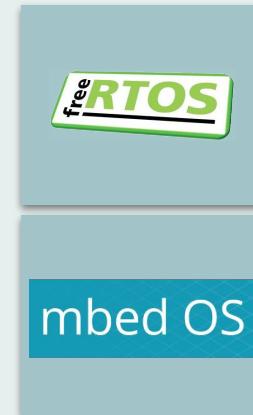
## Embedded Systems



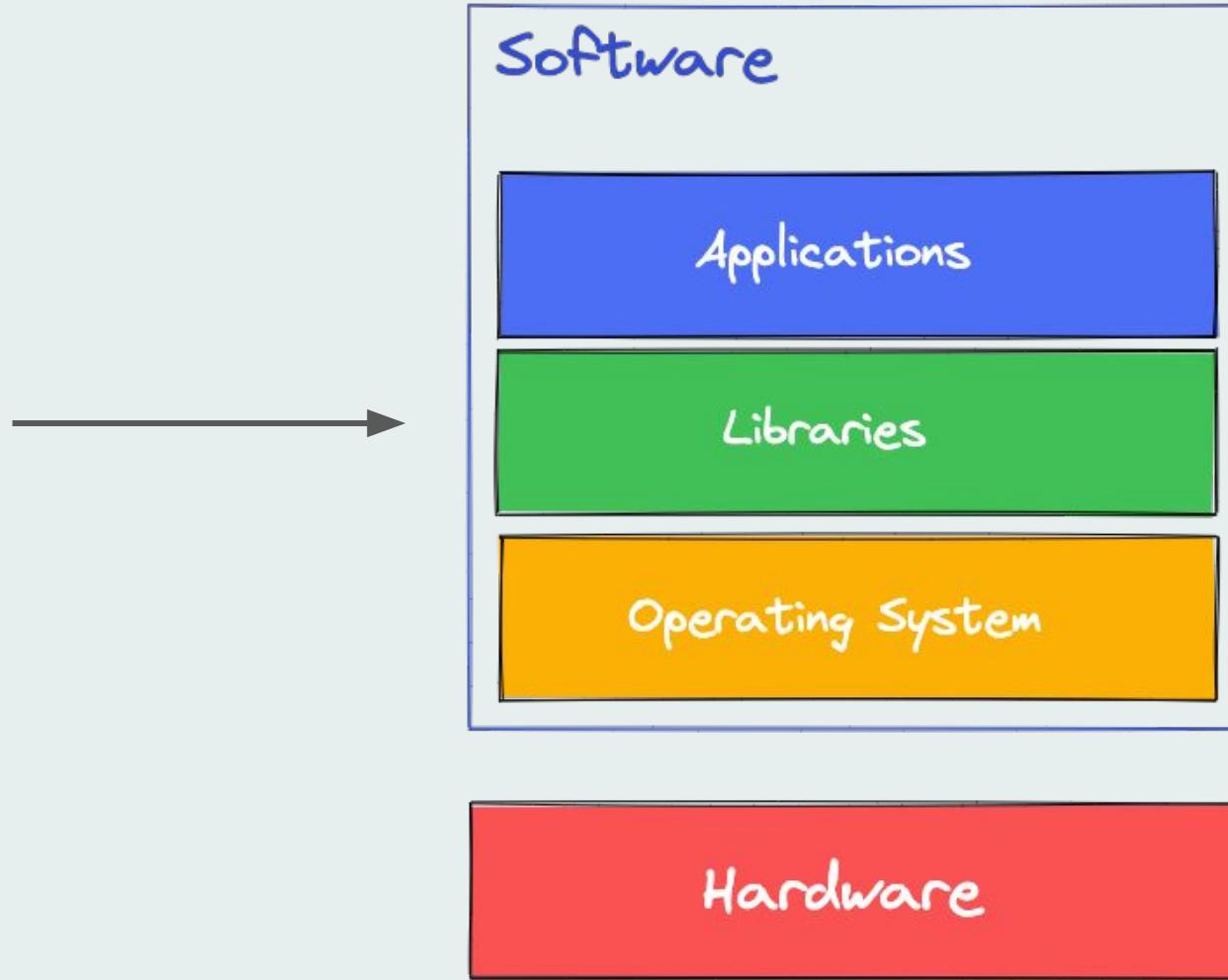
# Widely Used Operating Systems



Mobile OS



Embedded Sys.



# Software

Applications

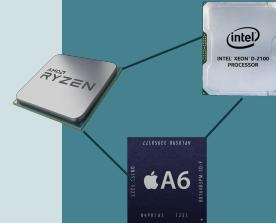
Libraries

Operating System

Hardware

```
import numpy as np
```

```
for x in range(10):  
    np.SaveTheWorld()
```



Software

Applications

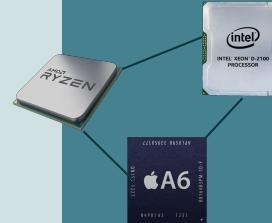
Libraries

Operating System

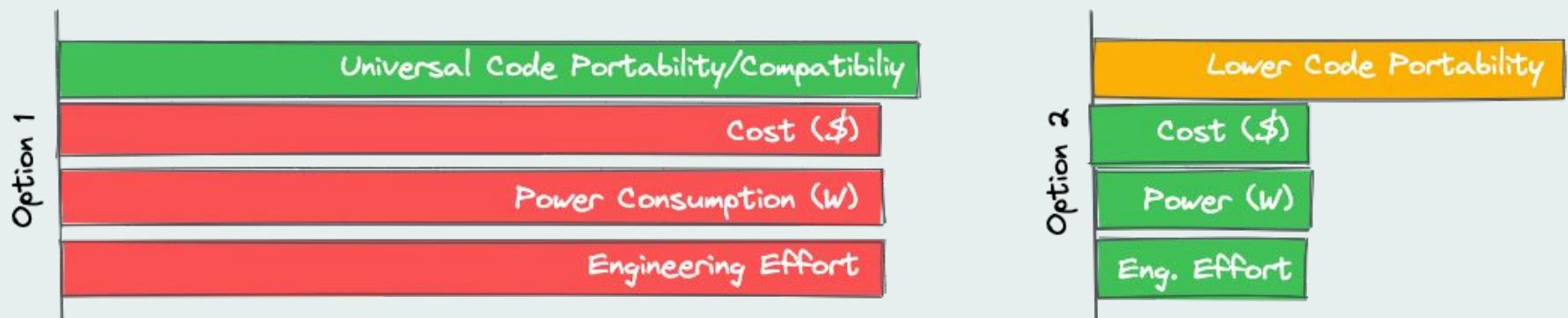
Hardware

# Portability Opportunity

Able to execute the same code on different microprocessor hardware and architecture

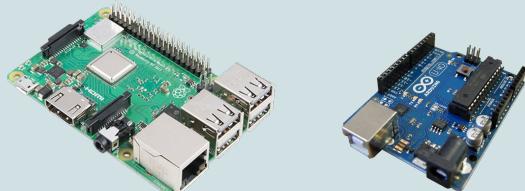


# Portability Trade-offs



# Portability Trade-offs

Sacrifice portability across systems for efficiency in system performance and power efficiency

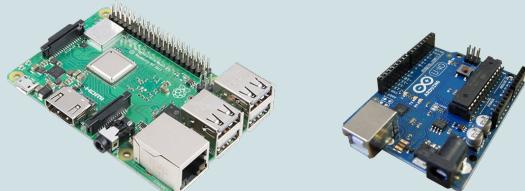


Specific HW Implementation of a Library



# Portability Trade-offs

Sacrifice portability across systems for efficiency in system performance and power efficiency



Specific HW Implementation of a Library

## Question:

How do we enable TinyML uniformly across these different systems if there is lower platform portability?

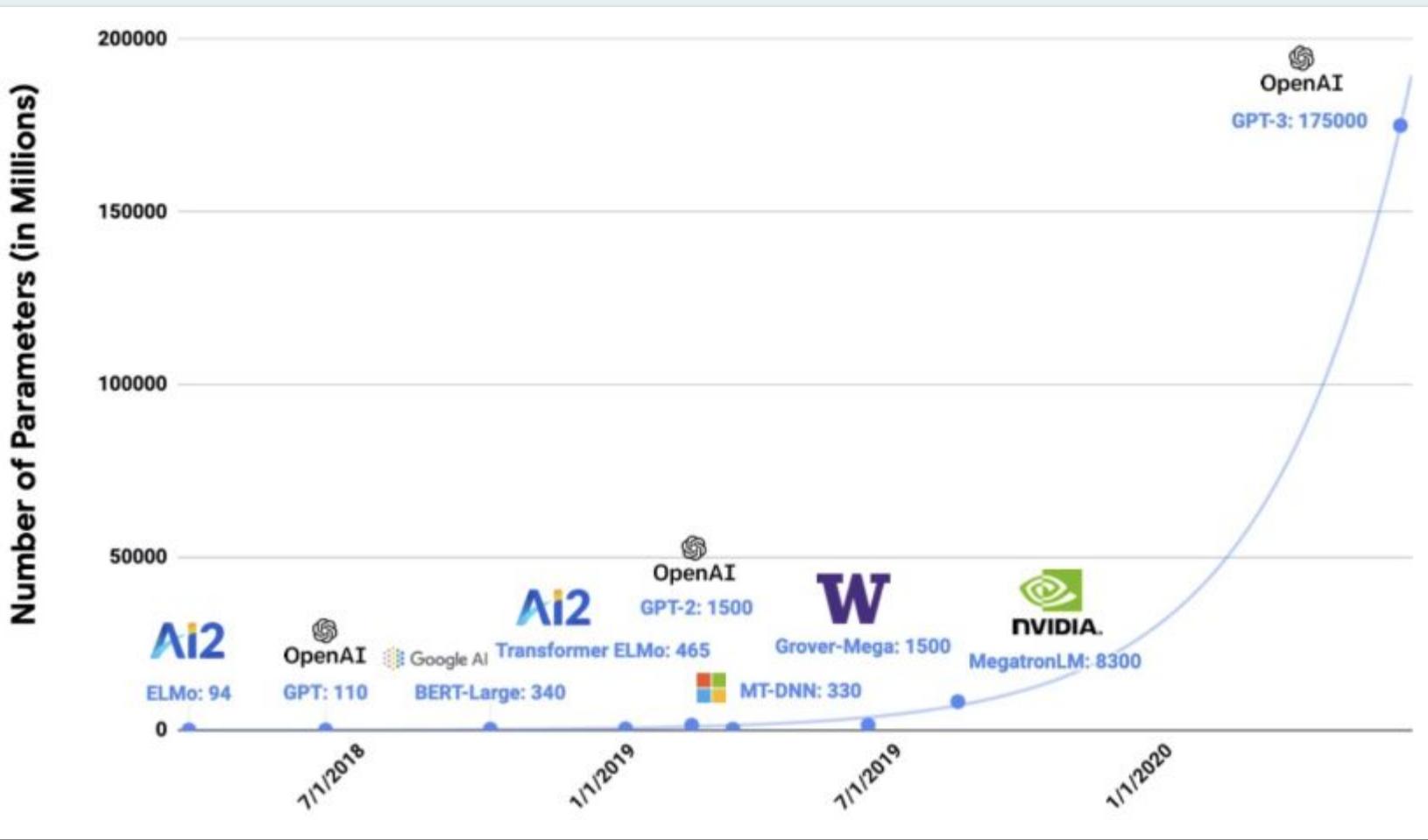
# Summary

**Embedded hardware** is extremely limited in performance, power consumption and storage

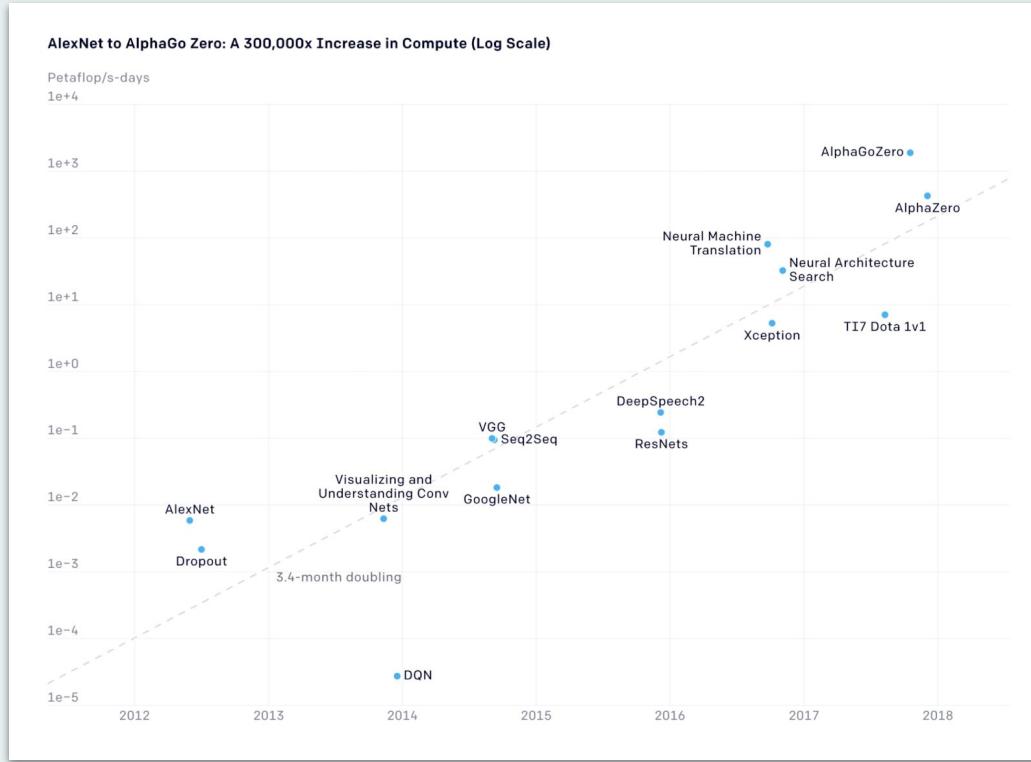
**Embedded software** is not as portable and flexible as mainstream computing



# TinyML

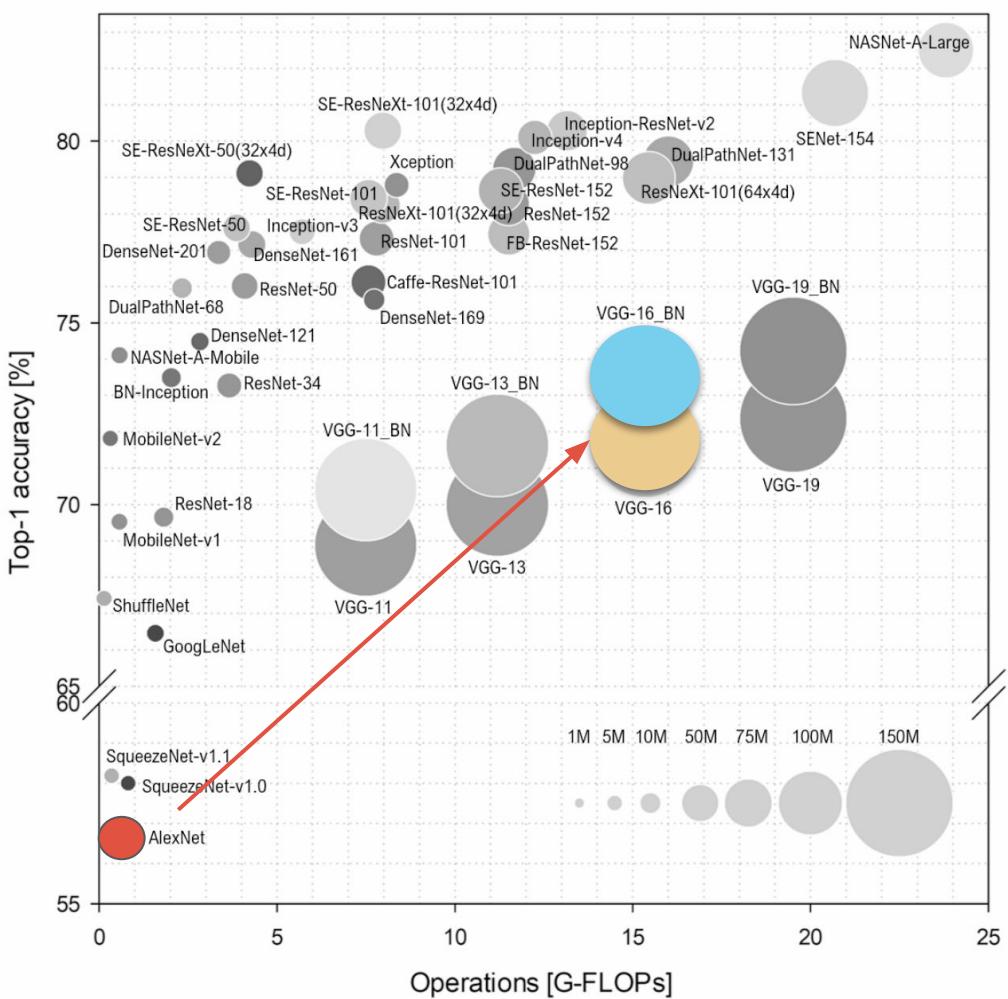


# ML Compute Needs (2012 to Present Day)



In recent years,  
**computing needs grew  
by 300,000x** to train  
the machine learning  
models that are widely  
deployed in the  
industry

# ML Model Evolution

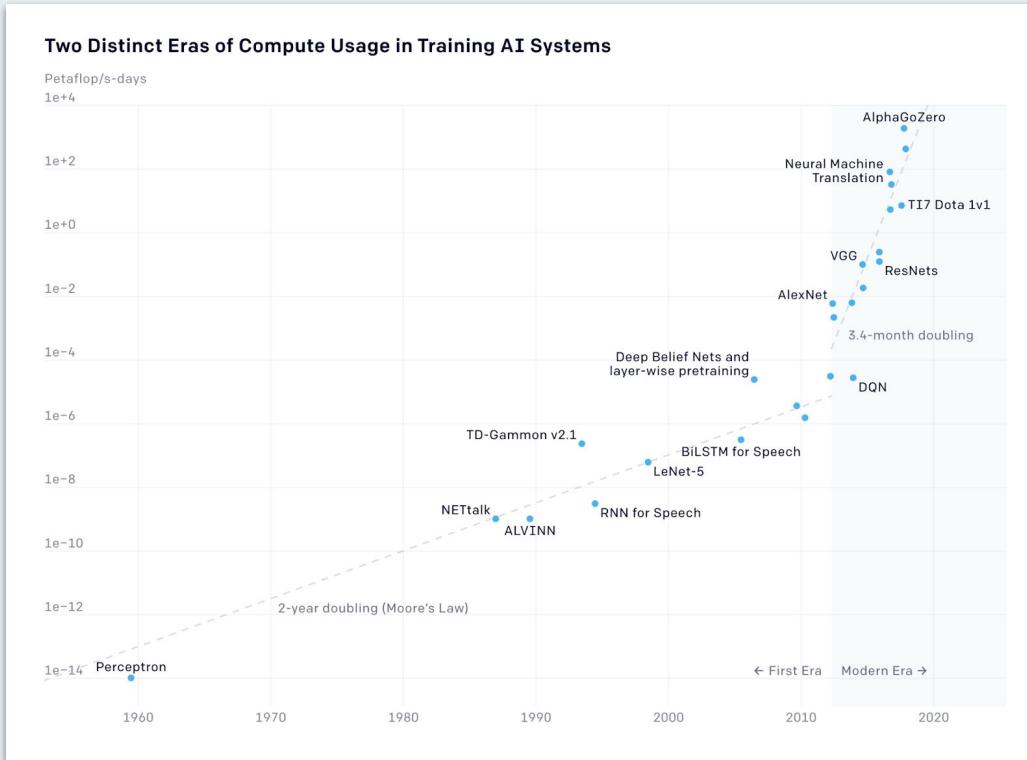


- VGGNet (2014)
  - 71.5% accuracy
  - 528 MB in size

Source: S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures". IEEE Access, vol. 6, 2018.



# ML Compute Needs (from the 1960s)



In recent years, the amount of computing needed has grown remarkably fast.

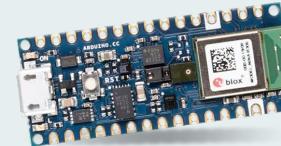
Computer requirements are **doubling nearly every 3 to 4 months.**



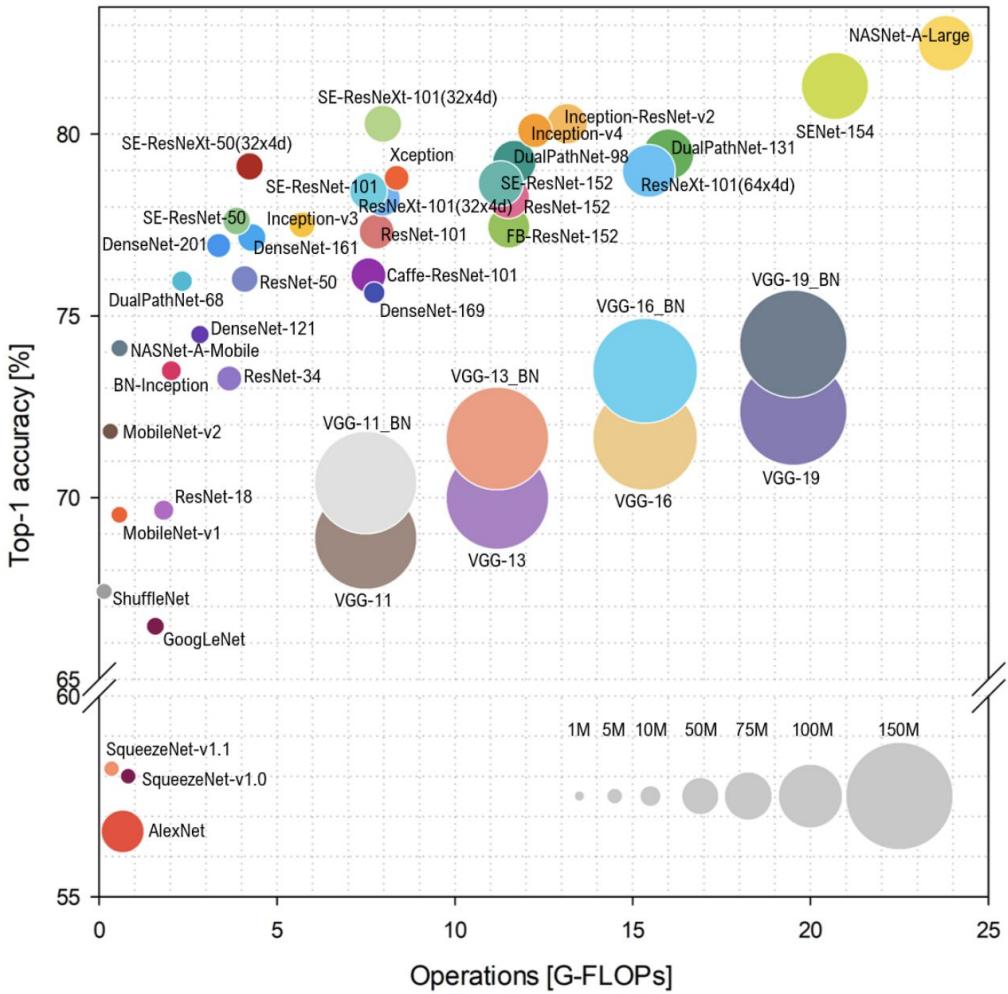
Cloud TPU



TinyML



# ML Model Evolution

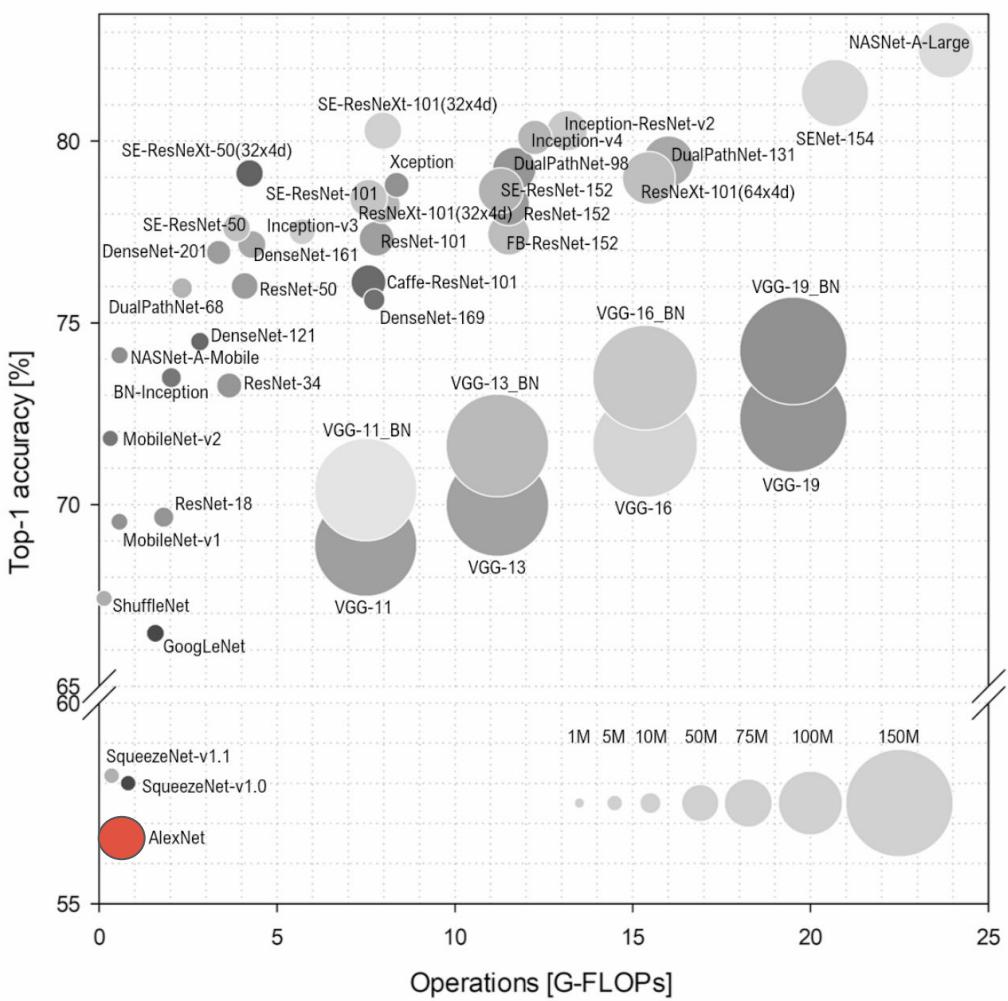


Source: S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures". IEEE Access, vol. 6, 2018.



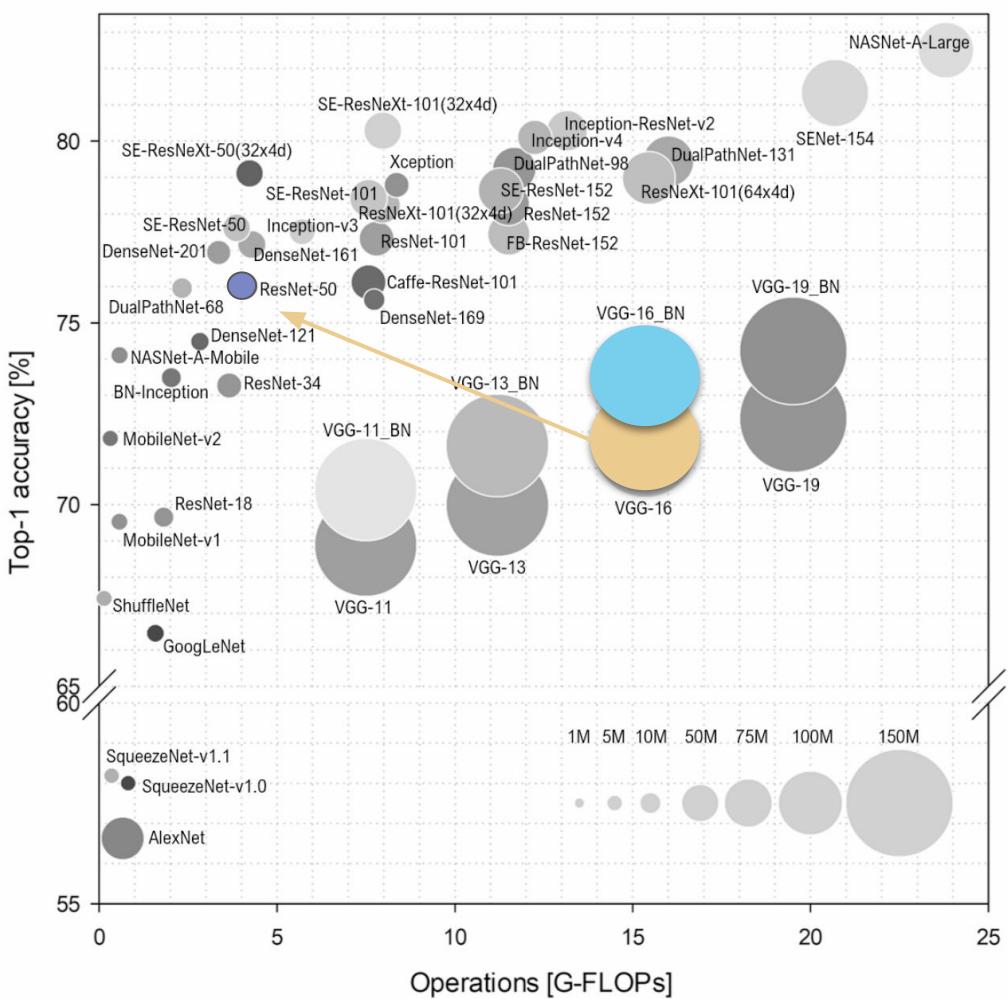
# ML Model Evolution

- **AlexNet (2012)**
  - 57.1% accuracy
  - 61 MB in size



Source: S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures". IEEE Access, vol. 6, 2018.

# ML Model Evolution

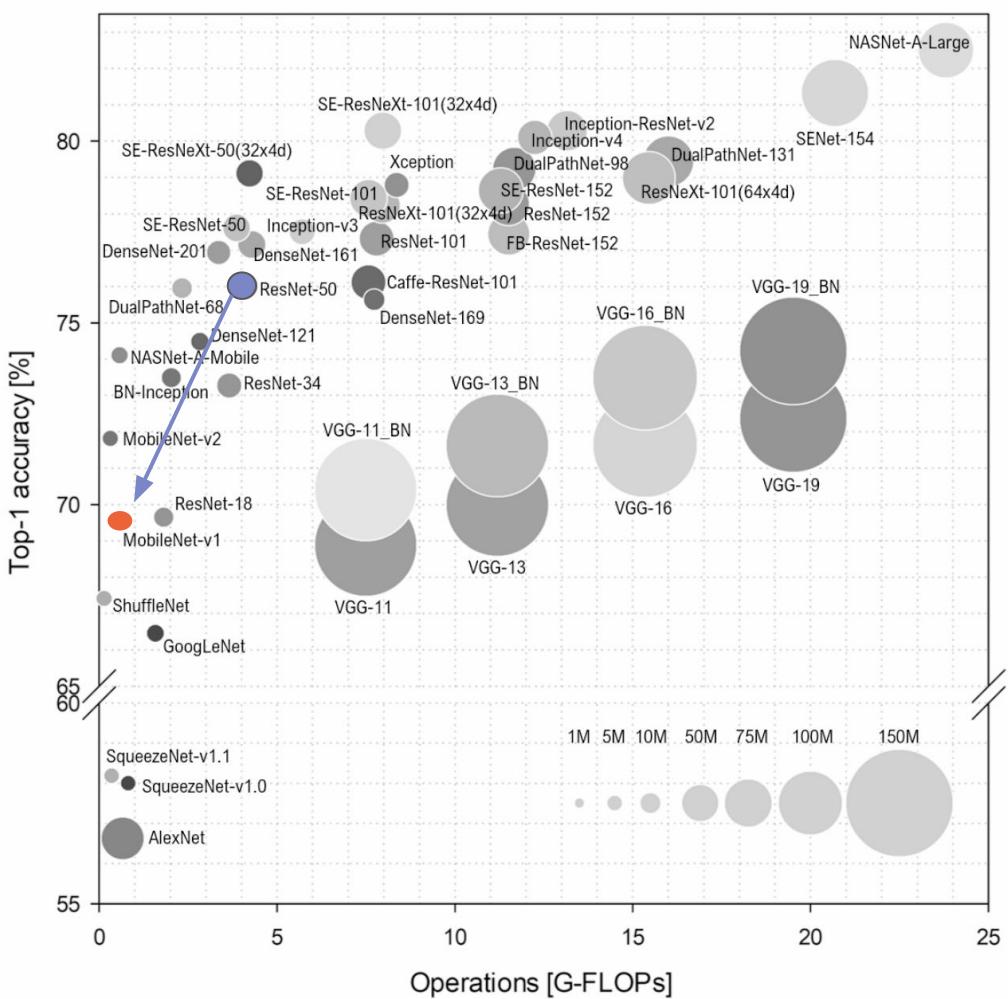


- ResNet (2015)
  - 75.8% accuracy
  - 22.7 MB in size

Source: S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures". IEEE Access, vol. 6, 2018.



# ML Model Evolution



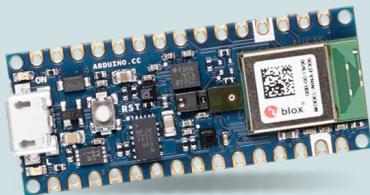
- **MobileNet (2015)**
  - 70.6% accuracy
  - 16.9 MB in size

Source: S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures". IEEE Access, vol. 6, 2018.



# Problem:

Our board only has **256 KB** of RAM (memory) yet MobileNetv1 needs **16.9MB!!!**



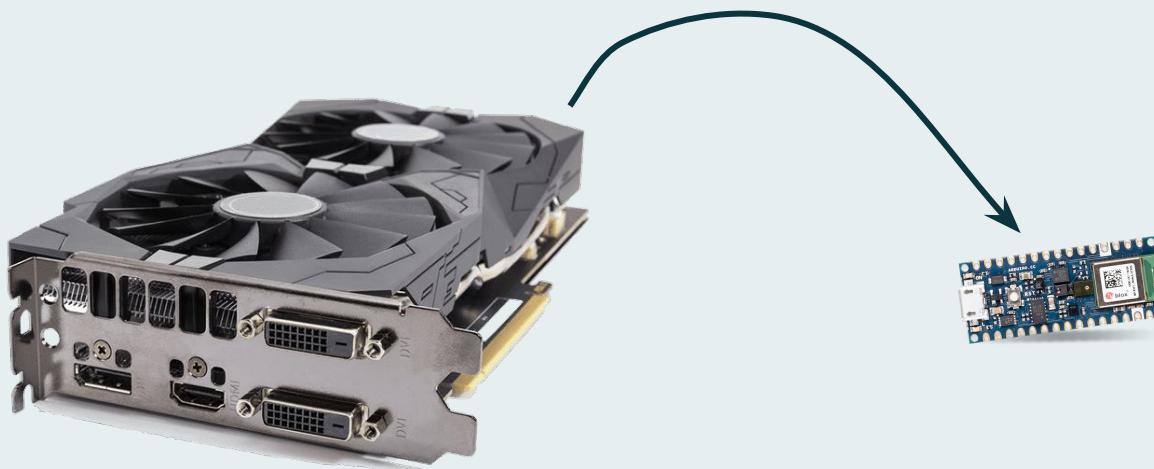
# ML Model Evolution

MobileNet (**2015**)

- 70.6% accuracy
- 16.9 MB in size

Source: S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures". IEEE Access, vol. 6, 2018.

# What are the Challenges for TinyML?



Machine Learning Models



Machine Learning Runtimes

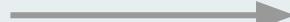
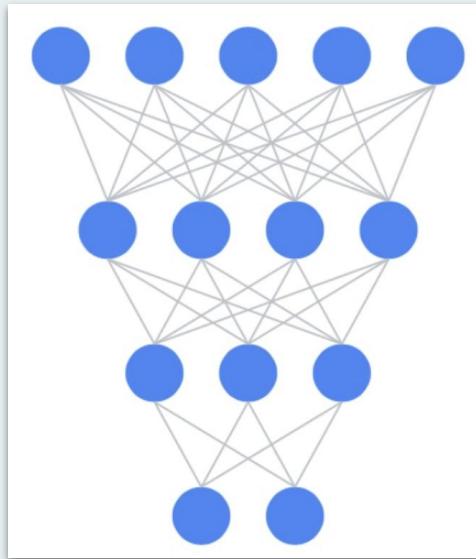


Machine Learning Hardware

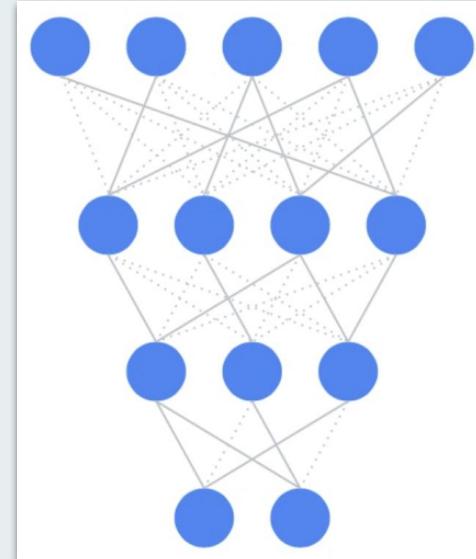
# Model Compression Techniques

Pruning  
Quantization  
Knowledge Distillation  
...

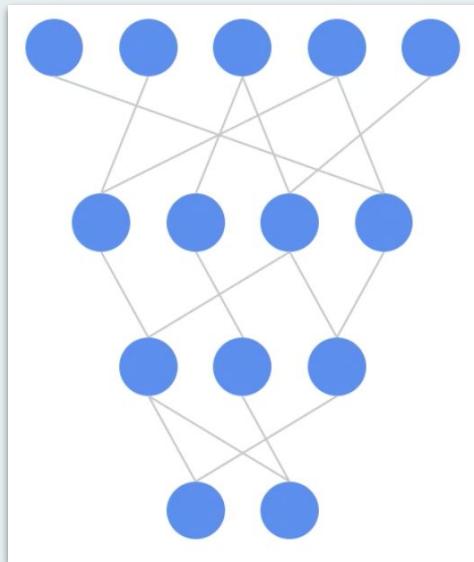
# Pruning



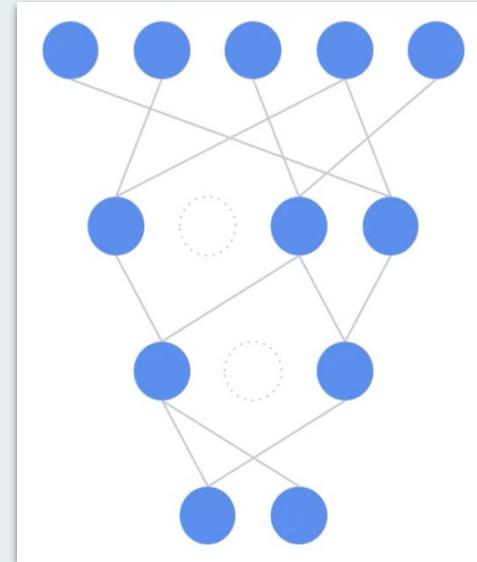
Pruning  
Synapses



# Pruning



→  
Pruning  
Neurons



Machine Learning Models



Machine Learning Runtimes

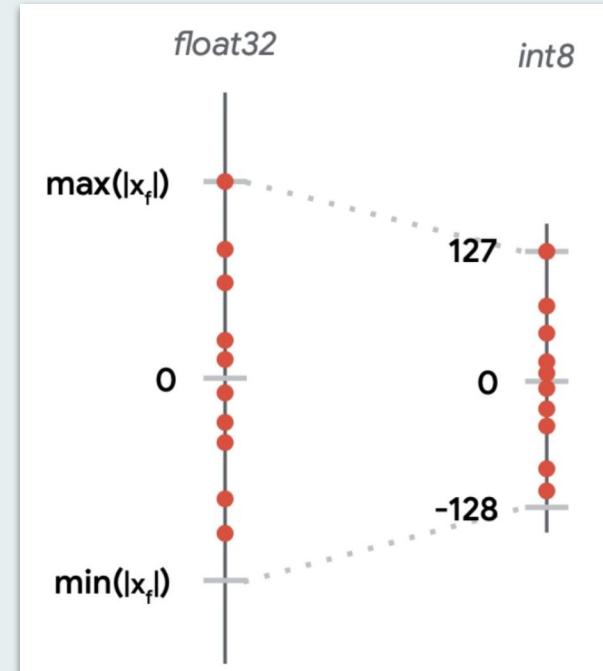


Machine Learning Hardware

# Model Compression Techniques

Pruning  
Quantization  
Knowledge Distillation  
...

# Quantization



Machine Learning Models



Machine Learning Runtimes



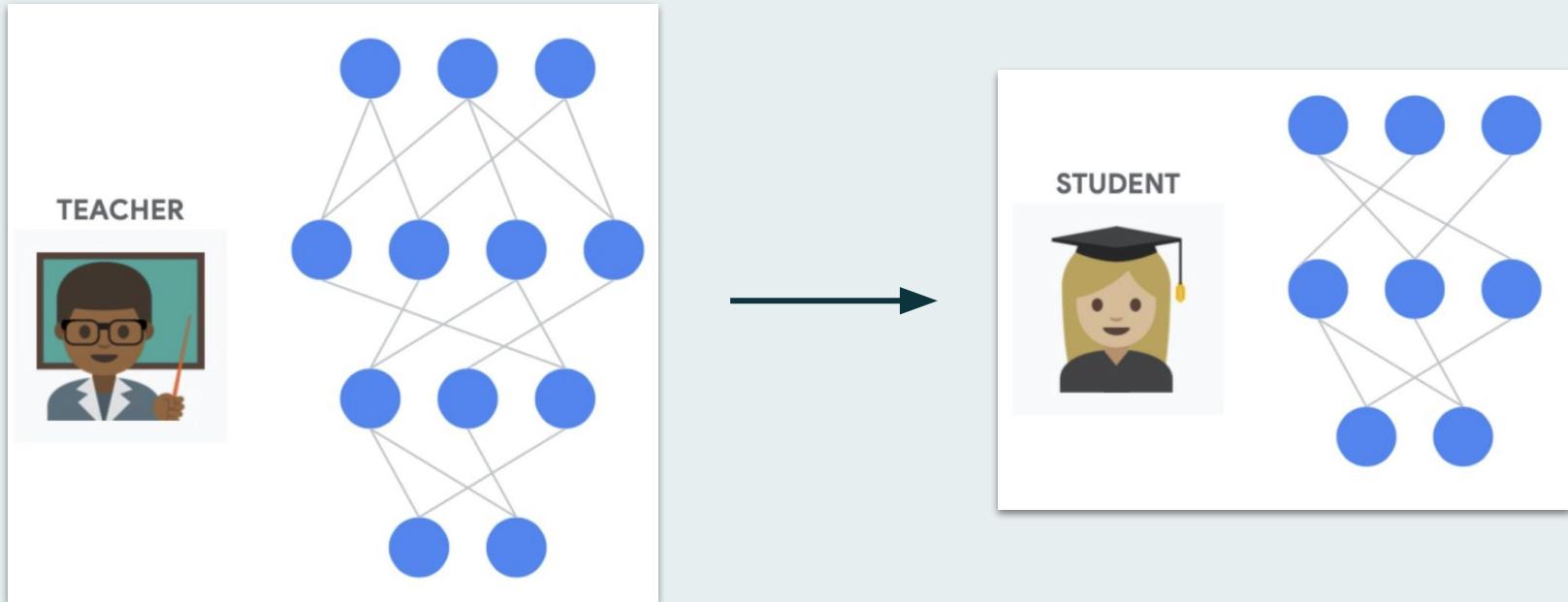
Machine Learning Hardware

# Model Compression Techniques

Pruning  
Quantization  
Knowledge Distillation

...

# Knowledge Distillation



Machine Learning Models



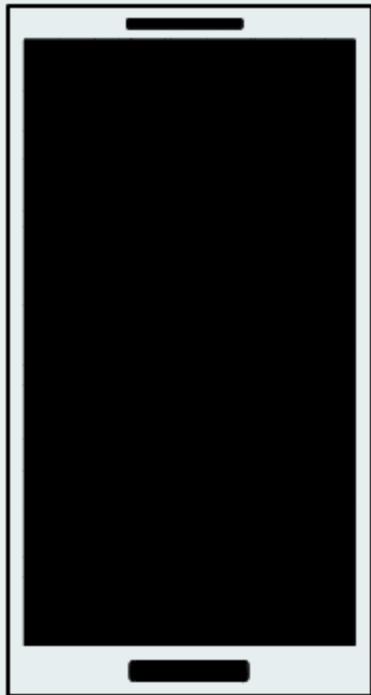
Machine Learning Runtimes



Machine Learning Hardware



TensorFlow



- Less Memory
- Less computer power
- Only focused on inference

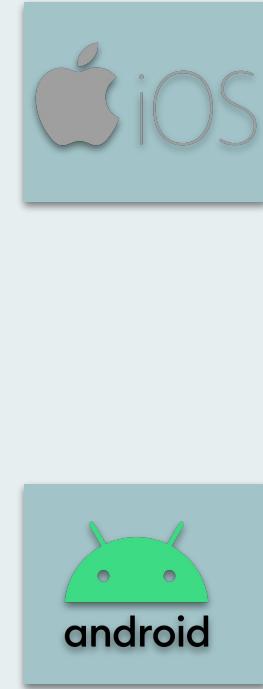
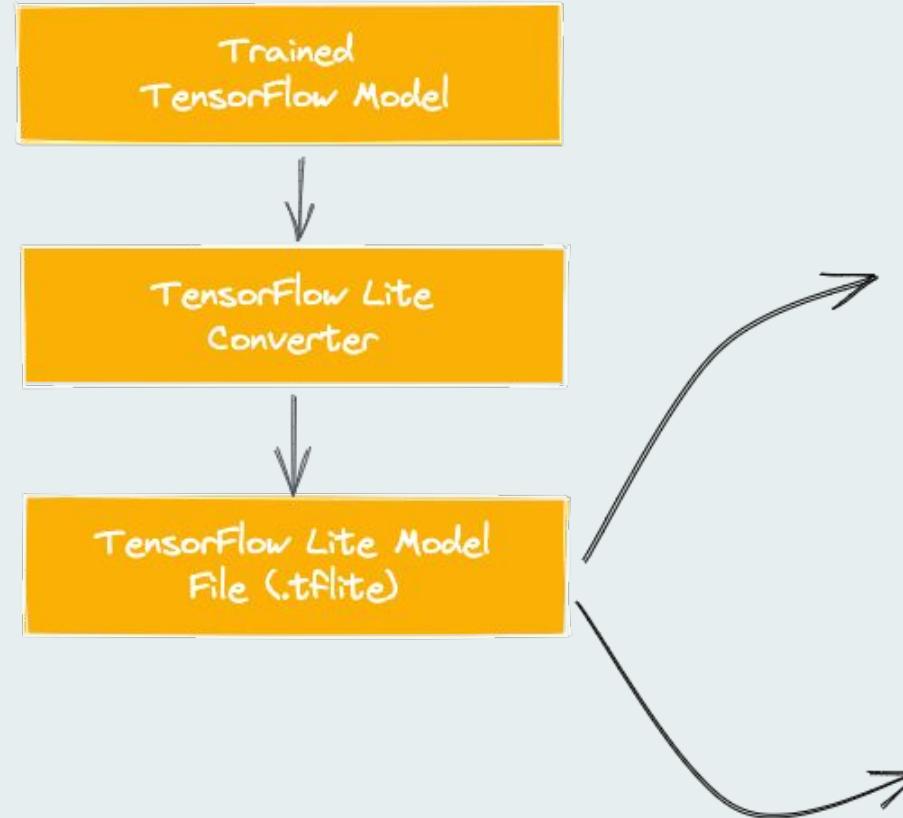
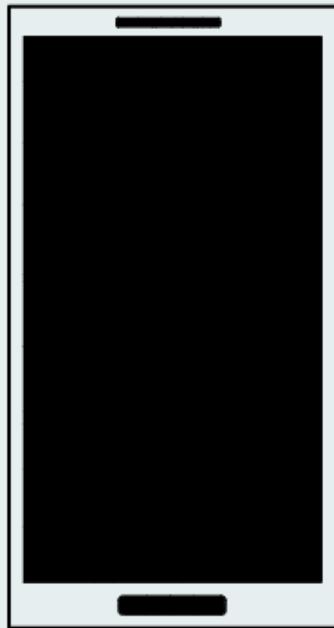


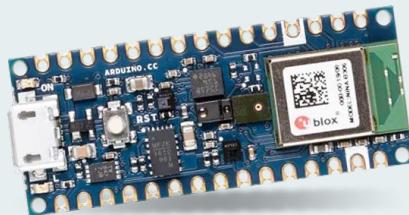
TensorFlow Lite

# Key Differences

	 TensorFlow	 TensorFlow Lite
<b>Topology</b>	<b>Variable</b>	<b>Fixed</b>
<b>Weights</b>	<b>Variable</b>	<b>Fixed</b>
<b>Binary Size</b>	<b>Unimportant</b>	<b>High Priority</b>
<b>Distributed Compute</b>	<b>Needed</b>	<b>Not Needed</b>
<b>Developer Background</b>	<b>ML Researcher</b>	<b>Application Developer</b>

# Architecture





Even less memory

Even less computer power

Also, only focused on inference

