arm

Deep learning
**on microcontrollers**

**Jan Jongboom**

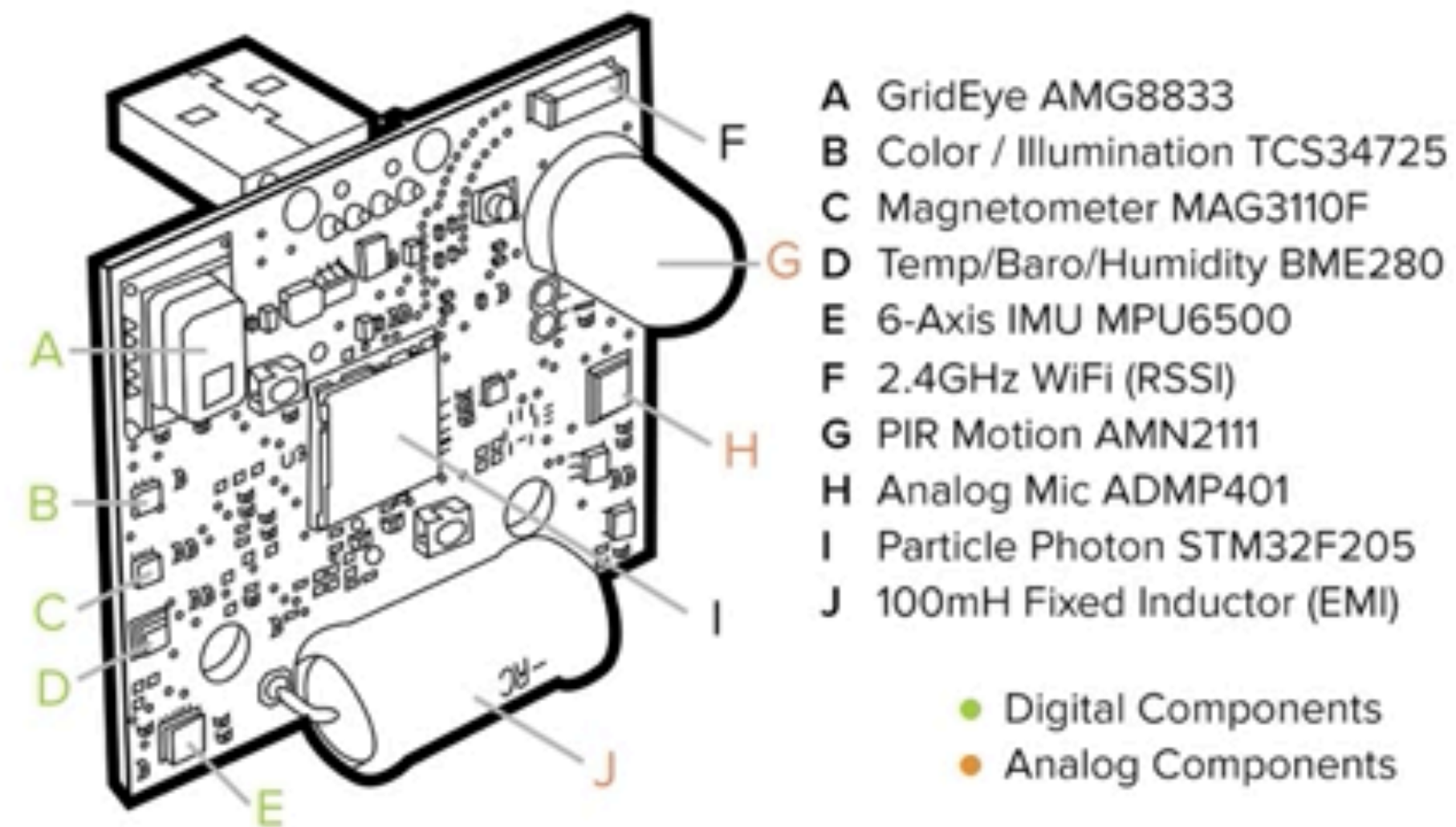IETF 101, London

22 March 2018

# Machine learning

arm

# Why machine learning on the edge?



A  GridEye AMG8833
B  Color / Illumination TCS34725
C  Magnetometer MAG3110F
D  Temp/Baro/Humidity BME280
E  6-Axis IMU MPU6500
F  2.4GHz WiFi (RSSI)
G  PIR Motion AMN2111
H  Analog Mic ADMP401
I  Particle Photon STM32F205
J  100mH Fixed Inductor (EMI)

● Digital Components
● Analog Components

# Sensor fusion

http://www.gierad.com/projects/supersensor/

arm

# Why machine learning on the edge?



Federated learning

https://research.googleblog.com/2017/04/federated-learning-collaborative.html

arm

# Why machine learning on the edge?



LPWANs

**arm**

# Why machine learning on the edge?



Offline self-contained systems
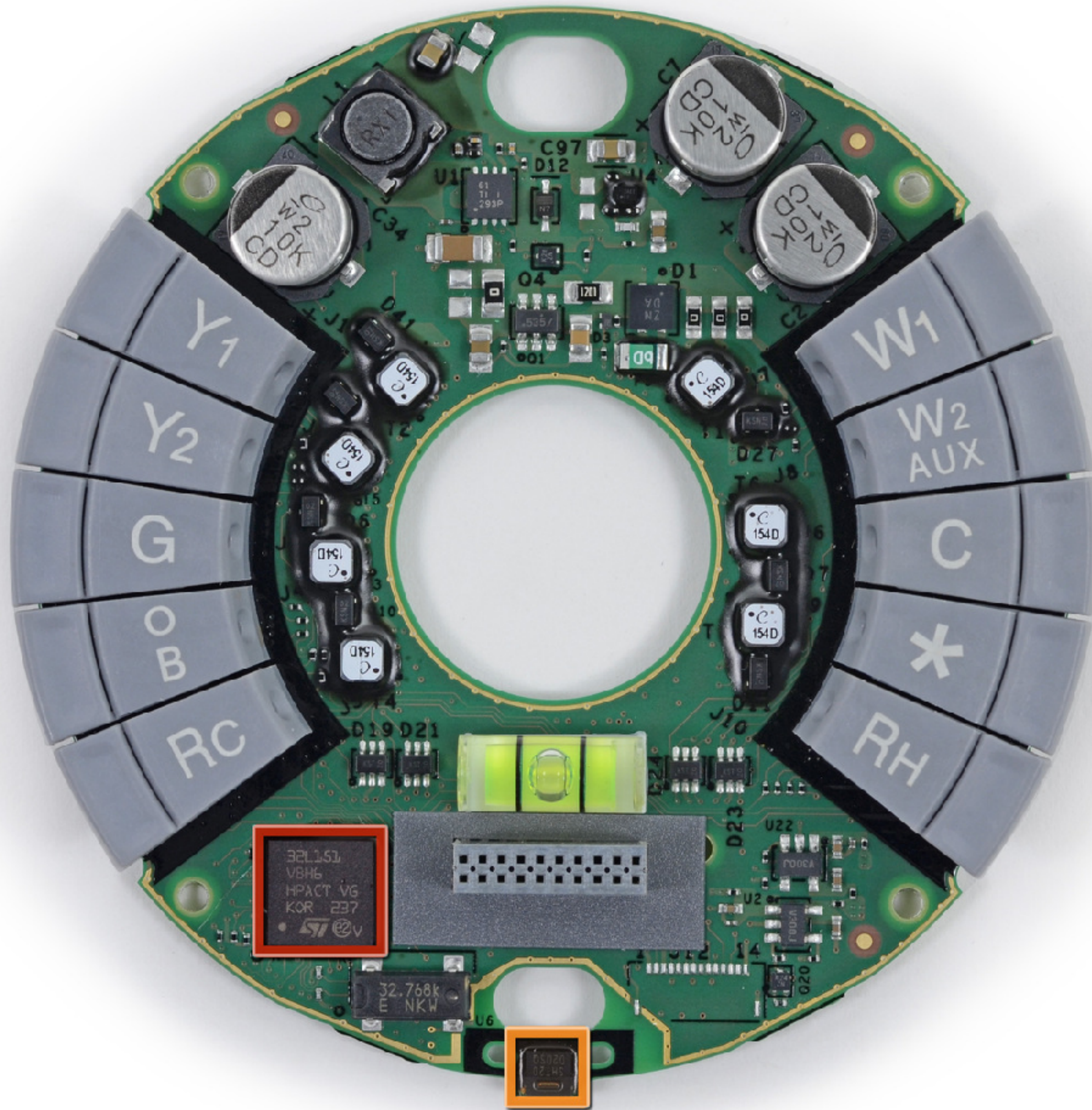
https://os.mbed.com/blog/entry/streaming-data-cows-dsa2017/

arm

# Edge vs. Cloud

arm

# Microcontrollers

Small (1cm²)

Cheap (~1$)

Efficient (standby: 0.3 µA)

# Downsides

Slow (max. 100 MHz)

Limited memory (max. 256K RAM)
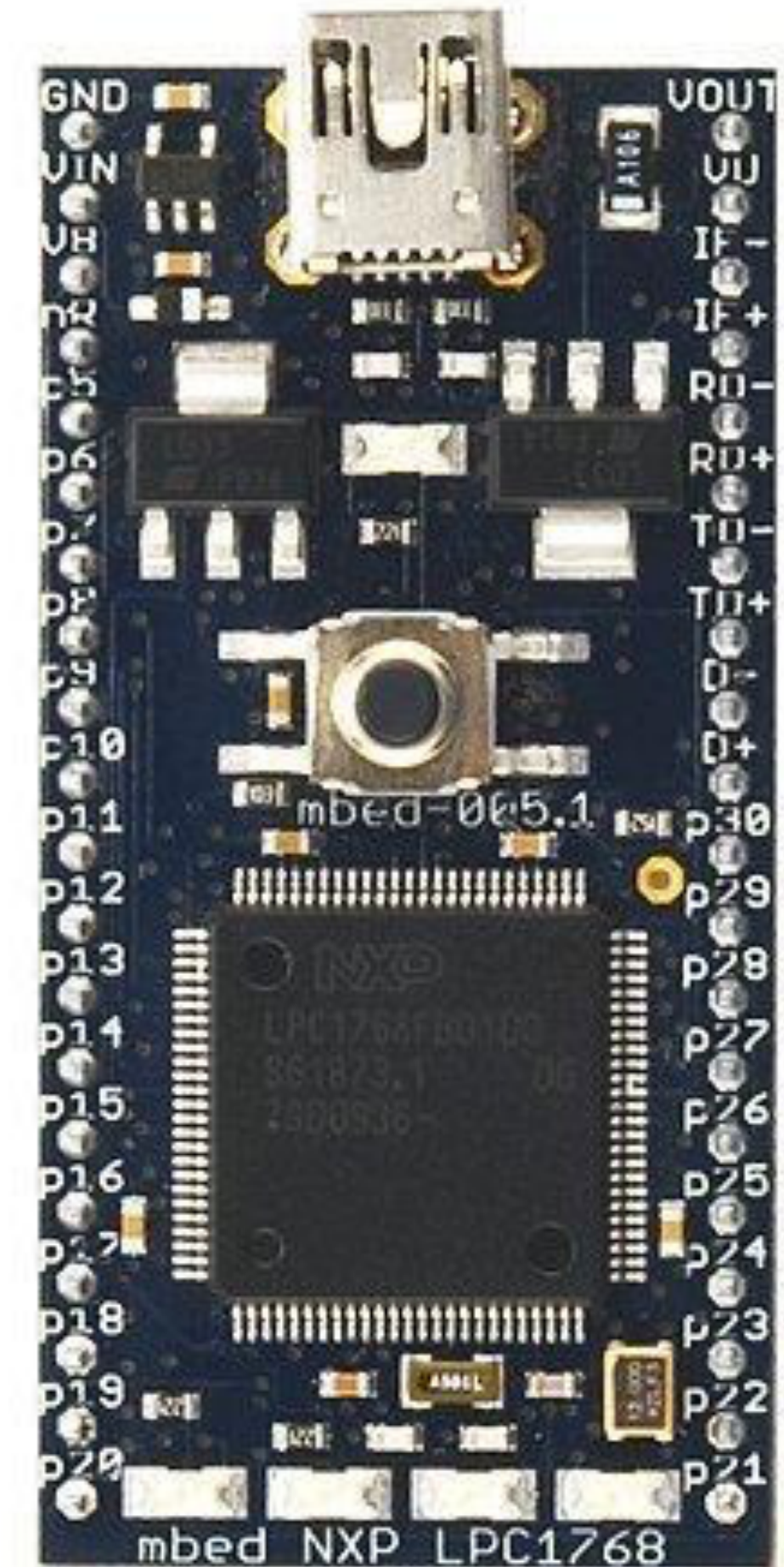
arm

# uTensor

Machine learning for microcontrollers

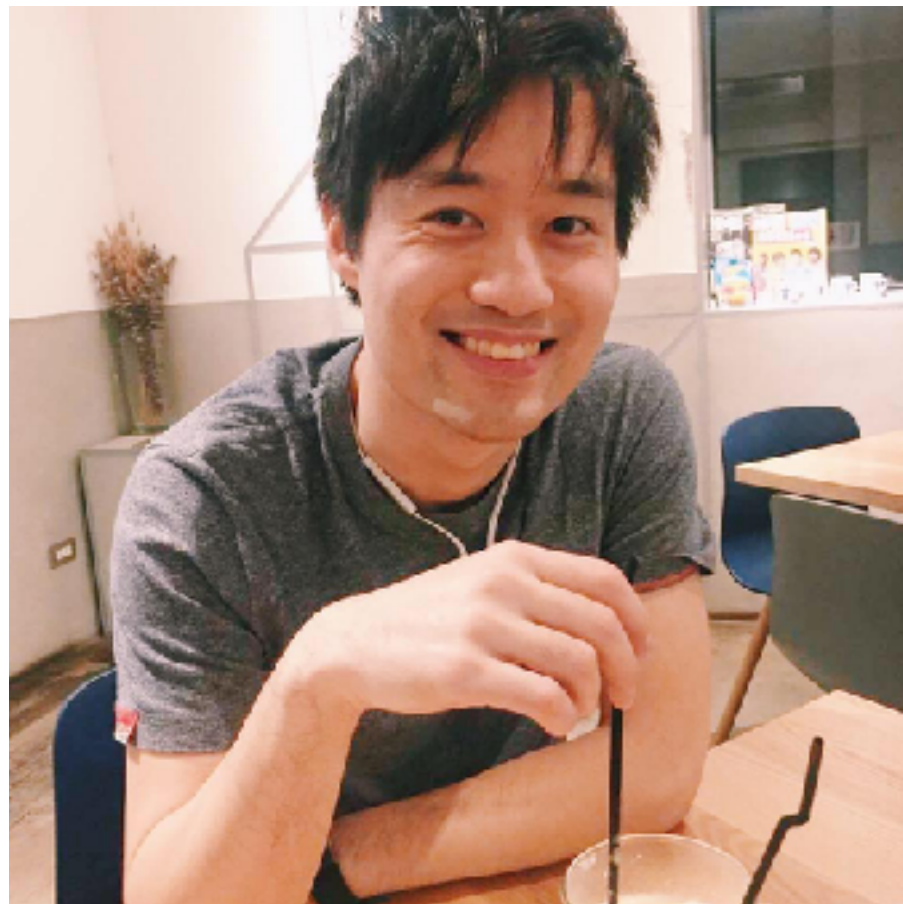Runs in <256K RAM

TensorFlow compatible

Built on top of Mbed OS 5
> (file systems, drivers, 150 boards compatible)

Open source, Apache 2.0 license

arm

# uTensor Team

**Neil Tan**
Arm
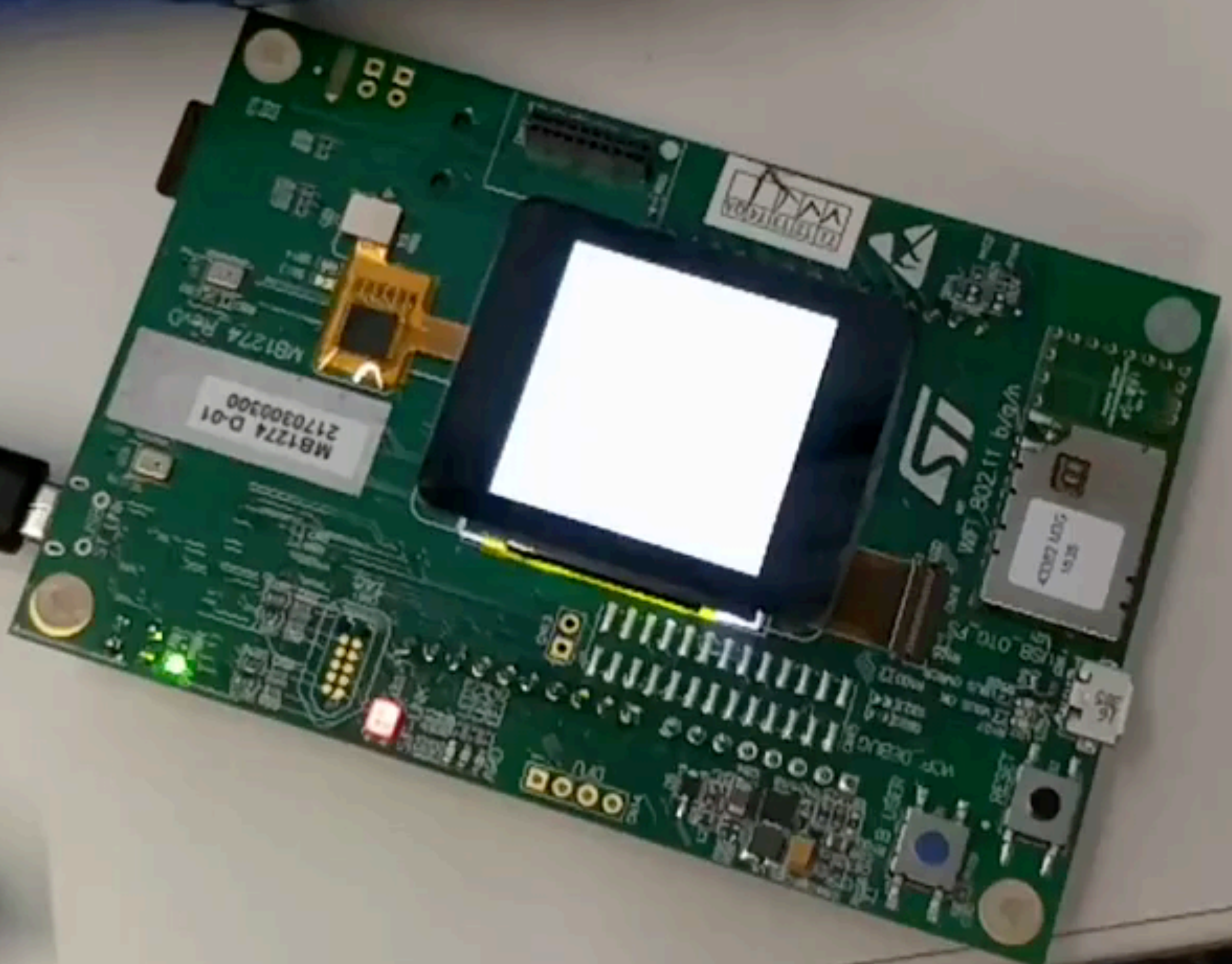
**Michael Bartling**
Arm

**Dboy Liao**
Piniko

**Kazami Hsieh**
Academia Sinica
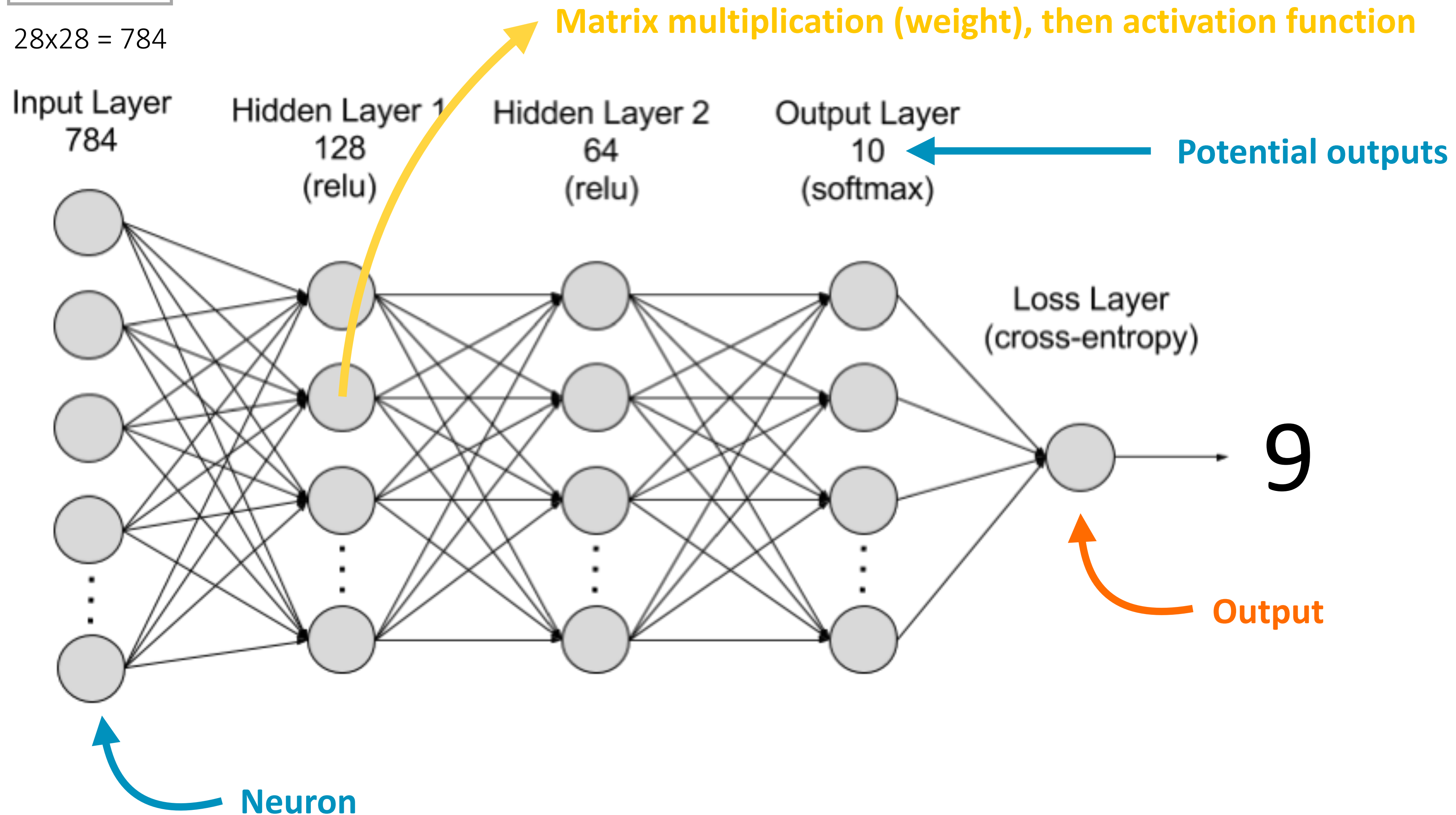
arm

# How?



MNIST data set

Training set: 60,000 images

Every drawing is downsampled to 28x28 pixels
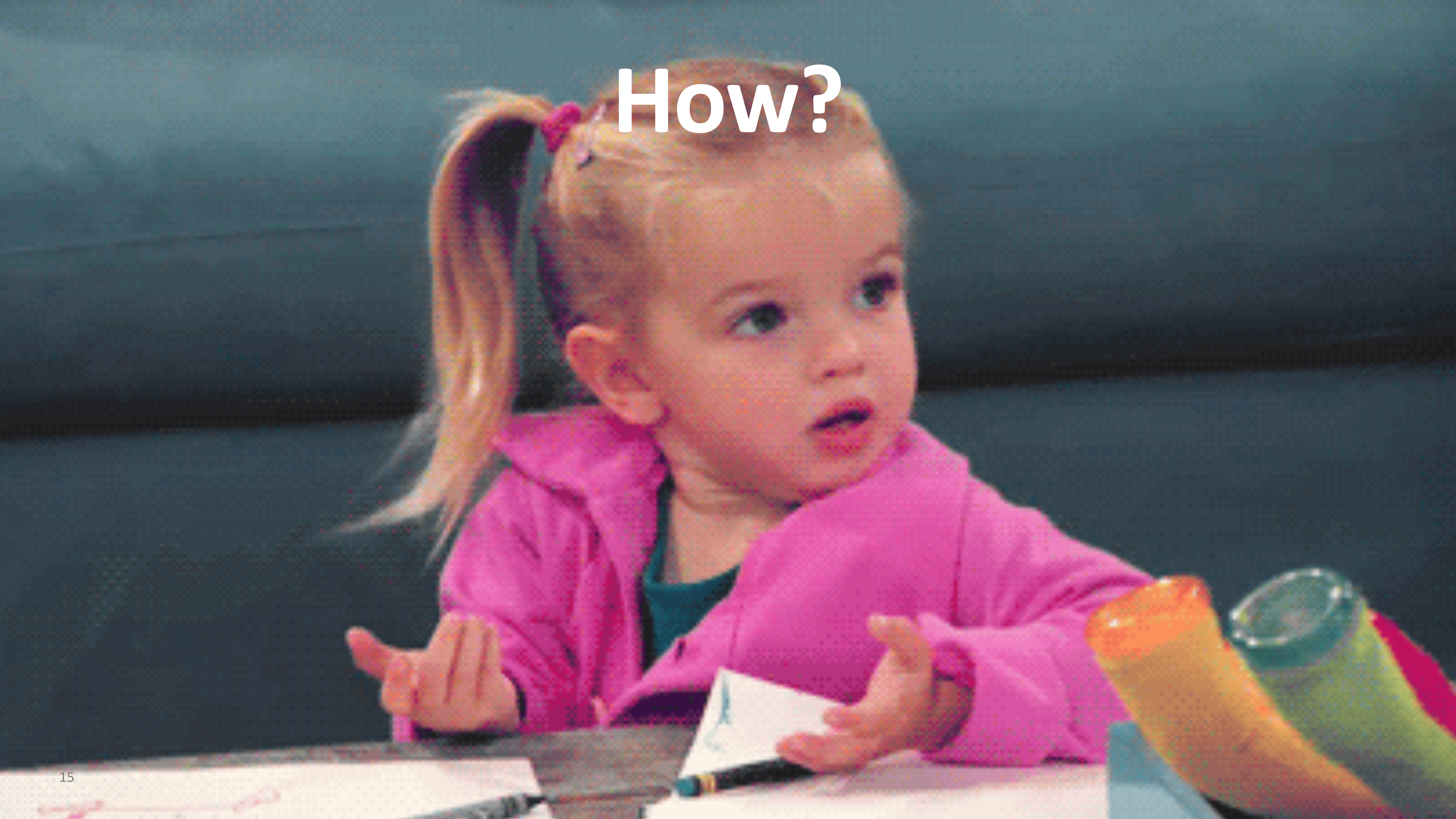
Supervised learning through backpropagation

**arm**

# Multi-layer perceptron (MLP) classification

28x28 = 784

**Matrix multiplication (weight), then activation function**

Input Layer
784

Hidden Layer 1
128
(relu)

Hidden Layer 2
64
(relu)

Output Layer
10
(softmax)

**Potential outputs**

Loss Layer
(cross-entropy)

9

**Output**

**Neuron**

arm

# How?

# Quantization

8-bit integers instead of 32-bit floats

Only during classification

79.9% accuracy vs. 80.3% accuracy (CIFAR-10)

TensorFlow requires floating-point de-quantization between layers

https://petewarden.com/2016/05/03/how-to-quantize-neural-networks-with-tensorflow/

arm

# Memory usage

Matrix multiplication in first hidden layer dominates RAM usage:

```
Input elements:                              784
Number of neurons (1st layer):               128
Number of weight (input to 1st layer):       128 * 784
Resulting values (Pre-activation function):  128
Data type:                                   8-bit integer (1 byte)
```

**1 byte * (784 + (128 * 784) + 128) = 98.891 kB**

**arm**

# Other tricks

Paging of memory for larger models (sacrifices speed)

Graph in ROM (requires pre-processing) (MNIST: 26K)

Take advantage in sparsity of data, sacrifice accuracy *(TBD)*

**arm**

# Operators

Add, Subtract

Min, Max, ArgMax

ReLU, Matrix multiplication, Reshape, Quantization
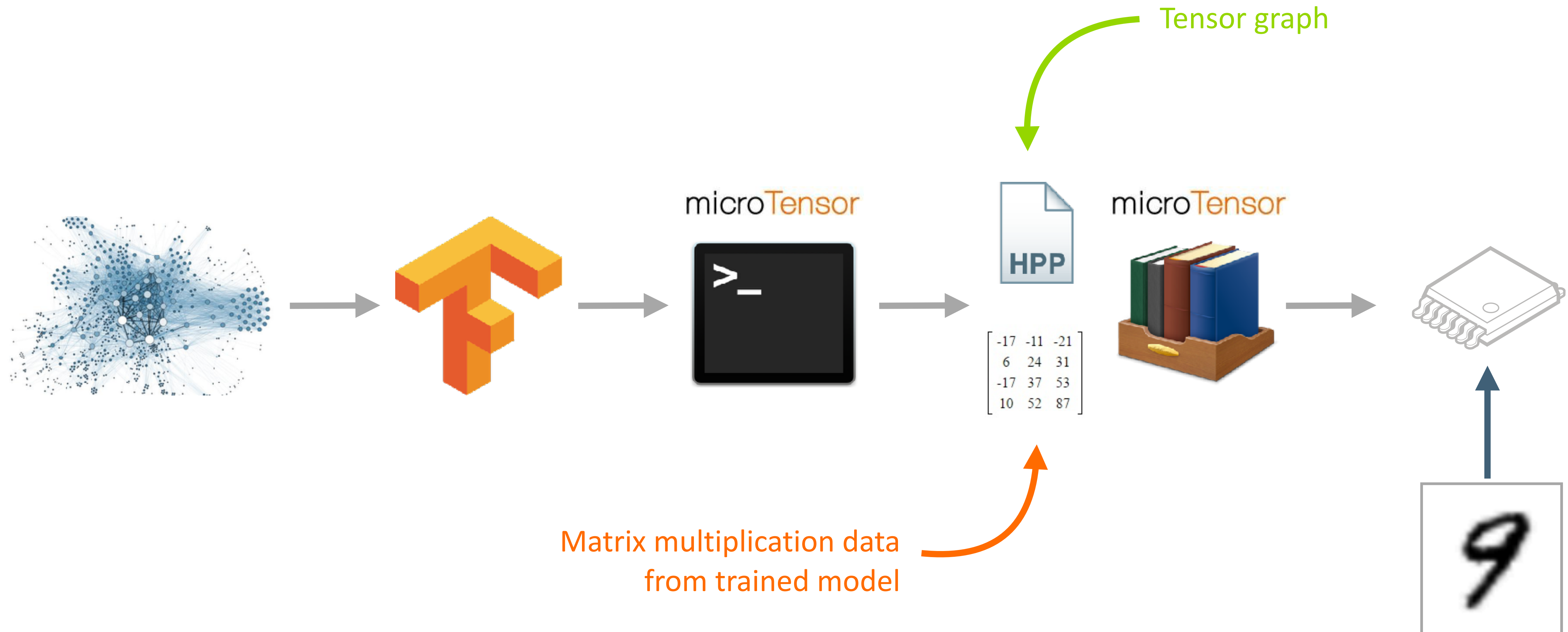
Convolution *(WIP)*

Pooling *(WIP)*

**arm**

# Tensors

RAM tensor

Flash tensor

Sparse tensor

Networked tensor


Tensors can be paged to fit larger networks

**arm**

# Workflow



Tensor graph

microTensor

HPP

microTensor

$$\begin{bmatrix} -17 & -11 & -21 \\ 6 & 24 & 31 \\ -17 & 37 & 53 \\ 10 & 52 & 87 \end{bmatrix}$$

Matrix multiplication data
from trained model

arm

# Developing using the simulator

# CMSIS-NN

New neural network kernel functions

Leverages the DSP/SIMD functions in silicon

See speedup of 4-5x

Hardware acceleration for convolution, pooling, etc.

uTensor will be built on top of CMSIS-NN

**arm**

# Recap

1. Buy a development board (http://os.mbed.com/platforms)

2. Clone uTensor (https://github.com/uTensor/uTensor)

3. ???

4. PROFIT!!!

# Thank you!

https://labs.mbed.com

Jan Jongboom, Arm