# arm

Machine Learning on MCUs

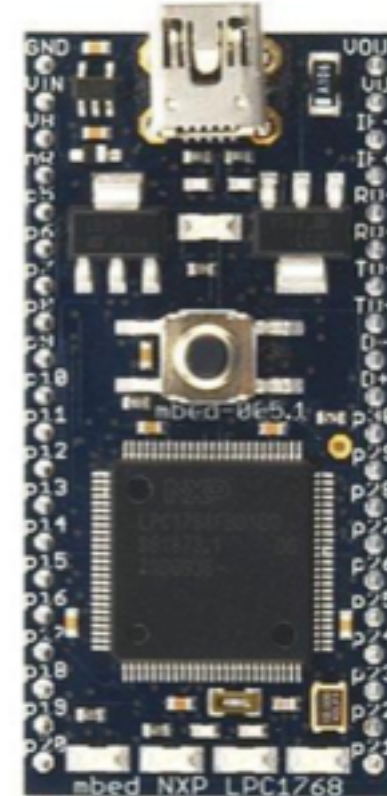# uTensor

A Mbed Labs Project

Neil Tan

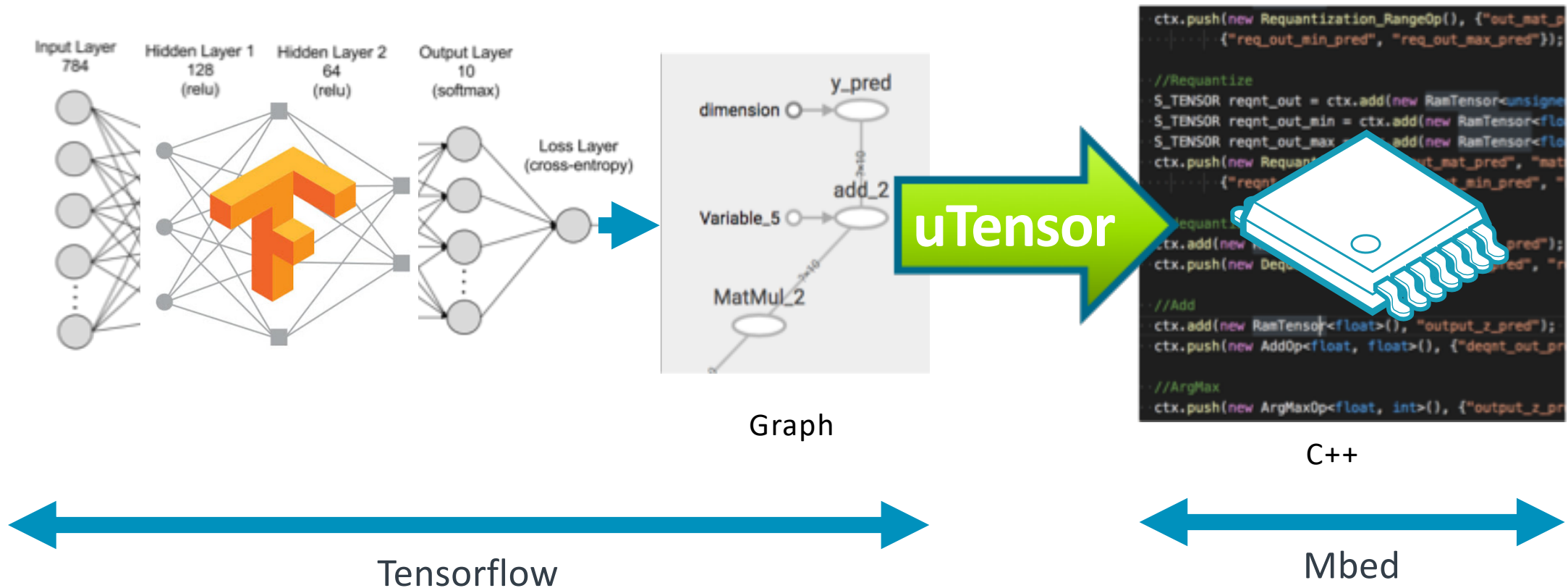July 2018

Machine learning

# uTensor

## Machine Learning for Microcontrollers

- Runs in <256K RAM

- Runs at ~100 MHz

- TensorFlow Compatible

- Inference Only

- Open source, Apache 2.0 license

arm

# uTensor



Graph

C++

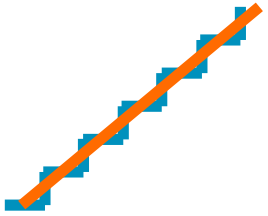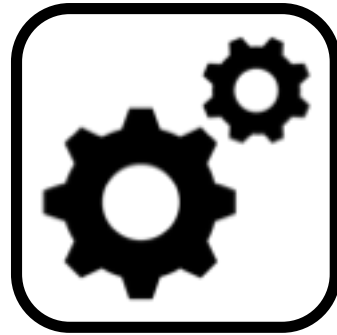Tensorflow

Mbed

arm

# Technologies

**Quantization**



Float to 8-bit
75% memory saving
Faster Computation

**Code Generation**



Copy and Paste
Easy Integration

**Intermediate Representation**
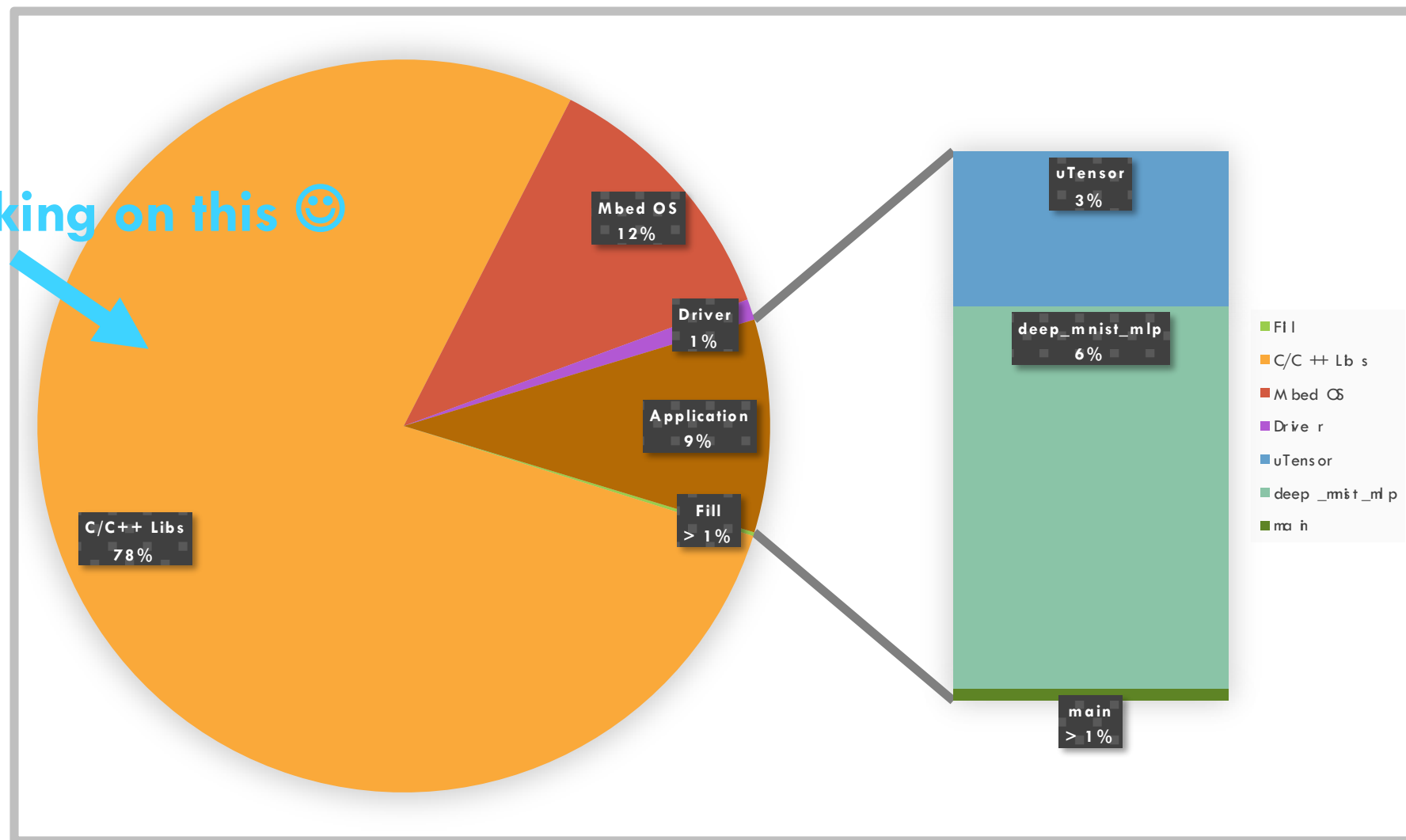


IR

Cross-Framework
Multi-Language

**Mbed**



arm
MBED

RTOS
CMSIS-NN
Connectivity
Production Ready

arm

# Binary



We're working on this ☺

Mbed OS
12%

Driver
1%

Application
9%

C/C++ Libs
78%

Fill
> 1%

uTensor
3%

deep_mnist_mlp
6%

main
> 1%

Legend:
- Fill
- C/C++ Libs
- Mbed OS
- Driver
- uTensor
- deep_mnist_mlp
- main

arm

# This is a Neural Network



Matrix multiplication (weight), bias, then activation

Input Layer
784

Hidden Layer 1
128
(relu)

Hidden Layer 2
64
(relu)

Output Layer
10
(softmax)

Output Classes

Loss Layer
(cross-entropy)

Output

Neuron

arm

# Why Matrix Multiplication

$$\text{ReLu}\left( \begin{bmatrix} \text{Input}_0 & \text{Input}_1 & - & - & - & \text{Input}_n \end{bmatrix} * \begin{bmatrix} W_{0,0} & W_{0,1} & & & W_{0,m} \\ W_{1,0} & W_{1,1} & & & W_{1,m} \\ W_{2,0} & W_{2,1} & - & - & W_{2,m} \\ W_{3,0} & W_{3,1} & & & W_{3,m} \\ W_{4,0} & W_{4,1} & & & W_{4,m} \\ & & & & \\ W_{n,0} & W_{n,1} & - & - & W_{n,m} \end{bmatrix} + \begin{bmatrix} B_{0,m} \\ B_{1,m} \\ B_{2,m} \\ B_{3,m} \\ B_{4,m} \\ \\ B_{n,m} \end{bmatrix} \right)$$

Input$_0$ x **W$_0$**
Input$_1$ x **W$_1$**
Input$_2$ x **W$_2$**
Input$_3$ x **W$_3$**
Input$_4$ x **W$_4$**

$\Sigma$ $z$

**Bias**

$$\textbf{Output} = \begin{bmatrix} \text{Output}_0 & \text{Output}_1 & - & - & \text{Output}_m \end{bmatrix}$$

**arm**

# Tensors and Operators

Tensors

Operators

$$\text{ReLu}\left(\begin{bmatrix}\text{Input}_0 & \text{Input}_1 & \text{---} & \text{Input}_n\end{bmatrix} * \begin{bmatrix} W_{0,0} & W_{0,1} & & W_{0,m} \\ W_{1,0} & W_{1,1} & & W_{1,m} \\ W_{2,0} & W_{2,1} & \text{---} & W_{2,m} \\ W_{3,0} & W_{3,1} & & W_{3,m} \\ W_{4,0} & W_{4,1} & & W_{4} \\ & & & \\ W_{n,0} & W_{n,1} & \text{---} & W_{n,m} \end{bmatrix} + \begin{bmatrix} B_{0,m} \\ B_{1,m} \\ B_{2,m} \\ B_{3,m} \\ B_{4,m} \\ \\ B_{n,m} \end{bmatrix}\right)$$

$$= \begin{bmatrix} \text{Output}_0 & \text{Output}_1 & \text{---} & \text{Output}_m \end{bmatrix}$$

arm

## Execution

**A + B => C**

**D * C => E**

**A + E => F**

Shared Pointer

arm

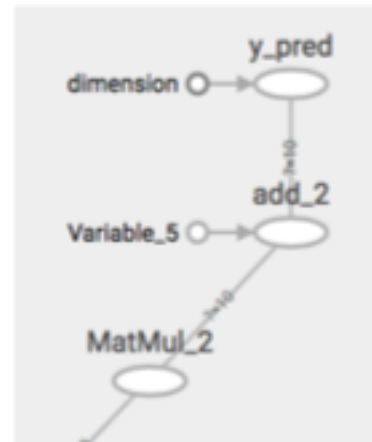# Code generator

A Python Tool

Turns Graph into C++ source

- Graph Traversal

- Optimizer

- Template Engine

TensorFlow



Graph

Mbed



C++

arm

# Graph

Operator Class

Tensor Class

init

Ops
In
Out

Kernel

de-init

Reference C

Remote
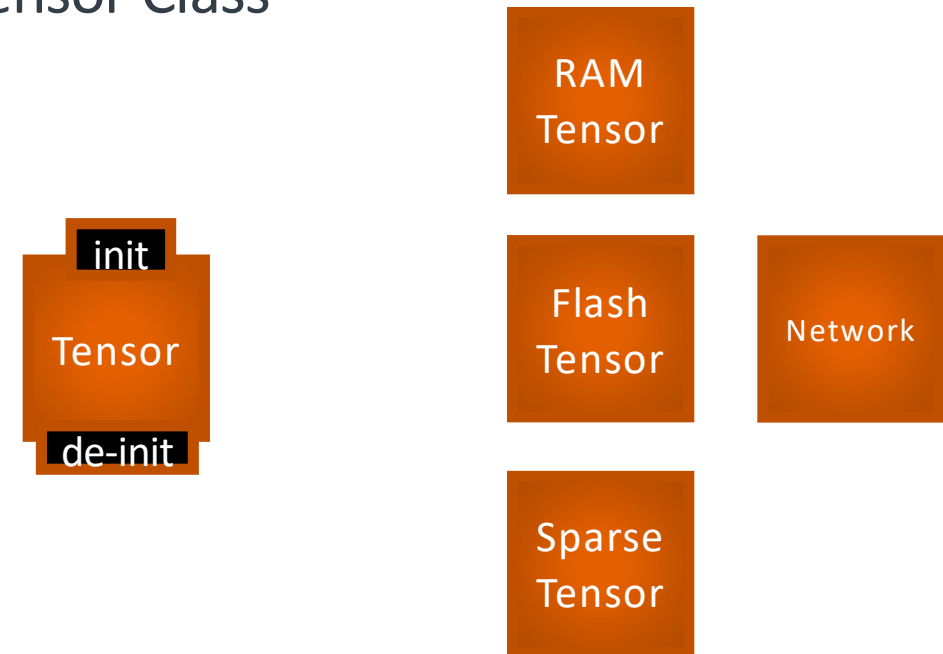
SIMD

SPI

init

Tensor

de-init

RAM Tensor

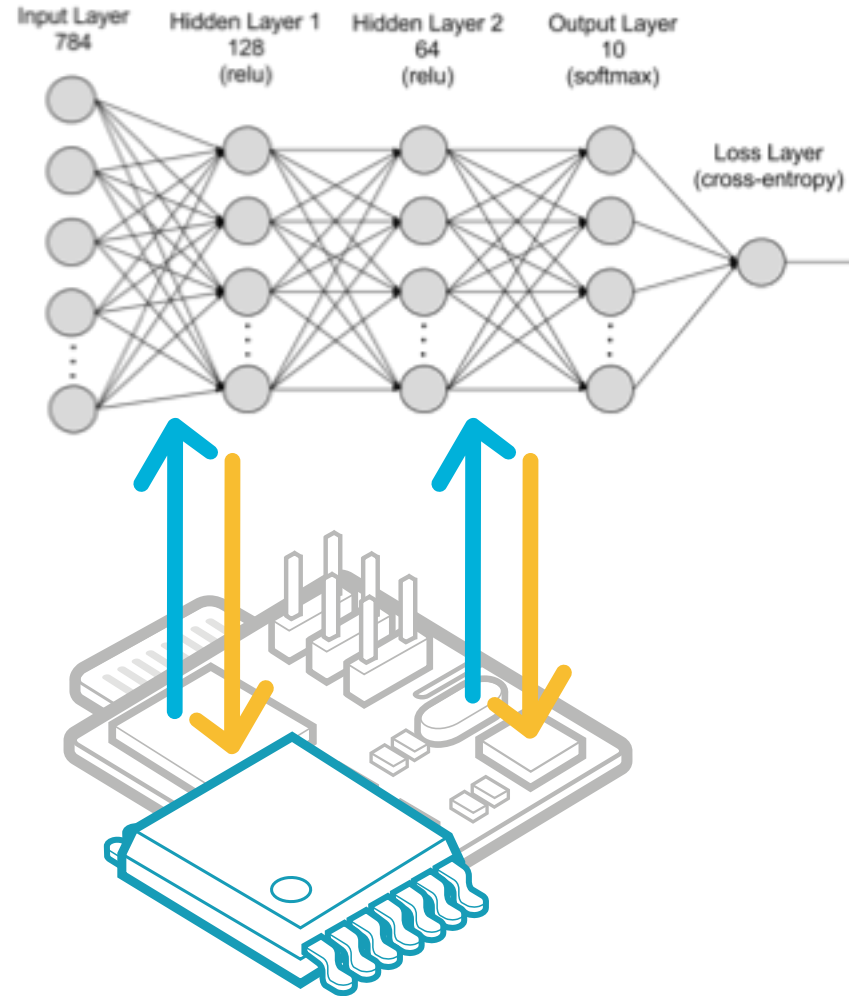Flash Tensor

Network

Sparse Tensor

arm

# FOTA Graph Update

uTensor's design allows the graph to be embedded in the binary

- Graph is in Firmware

- Firmware Over The Air

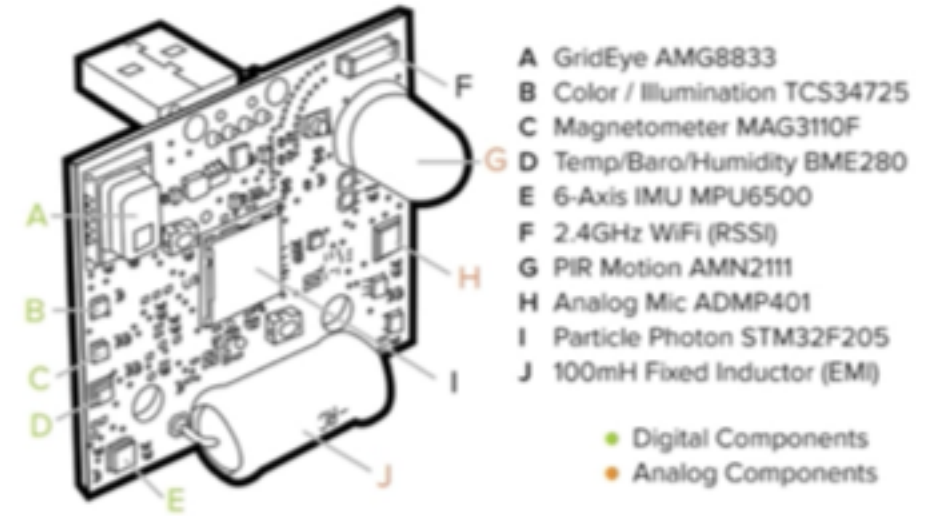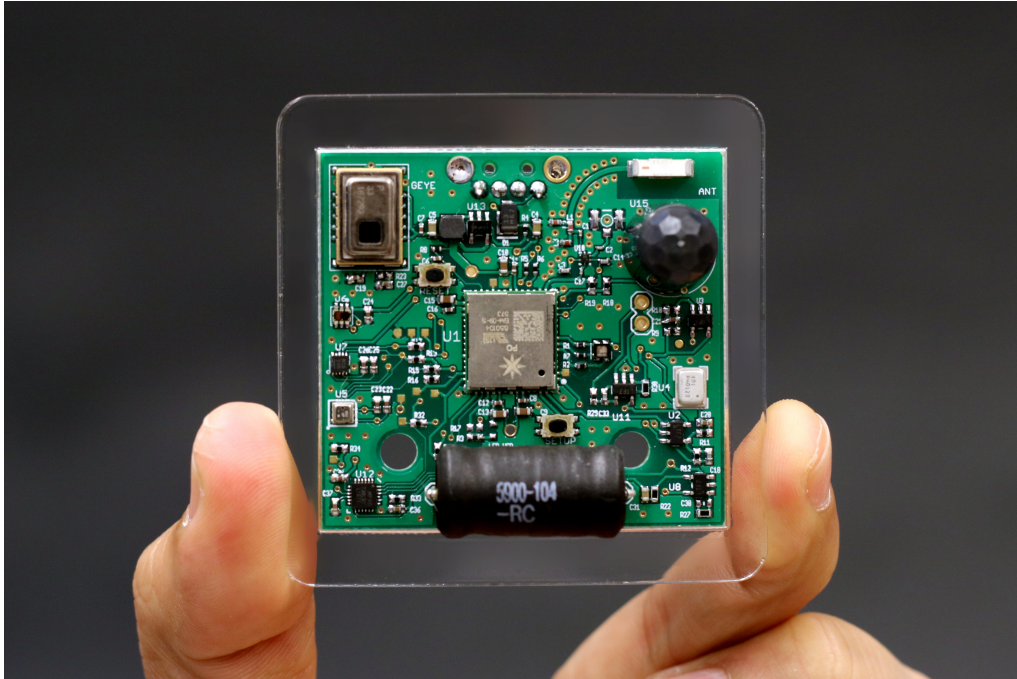- Weights stored in SD or Firmware

arm

# Sensor Fusion

http://www.aikidocenterla.com/blog/1212

# Synthetic Sensors



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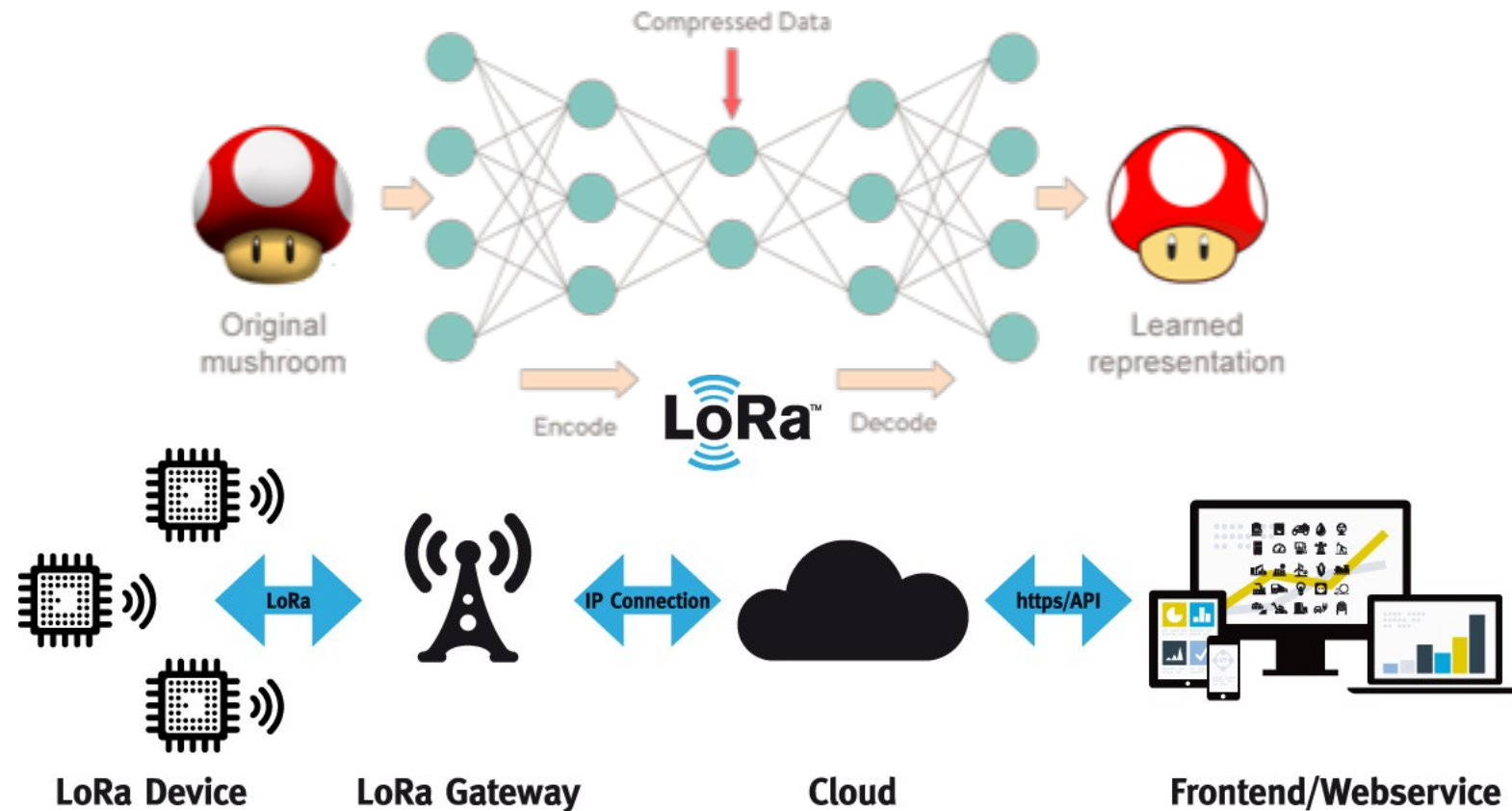
http://www.gierad.com/projects/supersensor/

arm

Recognized As:
Coffee Grinding

# Compression + LPWAN



http://curiousily.com/data-science/2017/02/02/what-to-do-when-data-is-missing-part-2.html
https://blog.microtronics.com/lora-and-2g-in-one-module/

# utensor.ai

## uTensor Timeline

- Test Release

  - Fully Connected, Dropout, Documentation

- Alpha Release

  - Convolution, SD Free

- Beta Release

  - CMSIS-NN



https://futurism.com/common-misunderstandings-of-evolution-part-2/

arm

# Thank You!

arm