



AI for Edge with Keras

K Agenda



What is EdgeAI?



Cloud to Edge



Key Metrics



Model Optimization
Techniques



Case Studies



Q & A

Edge AI

edge computing

artificial intelligence

K What is EdgeAI?

AI on the Edge refers to the practice of running artificial intelligence (AI) algorithms directly on endpoint devices, such as IoT devices, smartphones, drones, or sensors, rather than in a centralized data center or cloud environment.

The **Edge** refers to the computational capabilities embedded in these devices at the "edge" of the network, as opposed to centralized servers.



The benefits of On-Device Machine Learning



Low Latency

Unlock new user experiences by processing text, audio and video in real-time



Keep data on-device

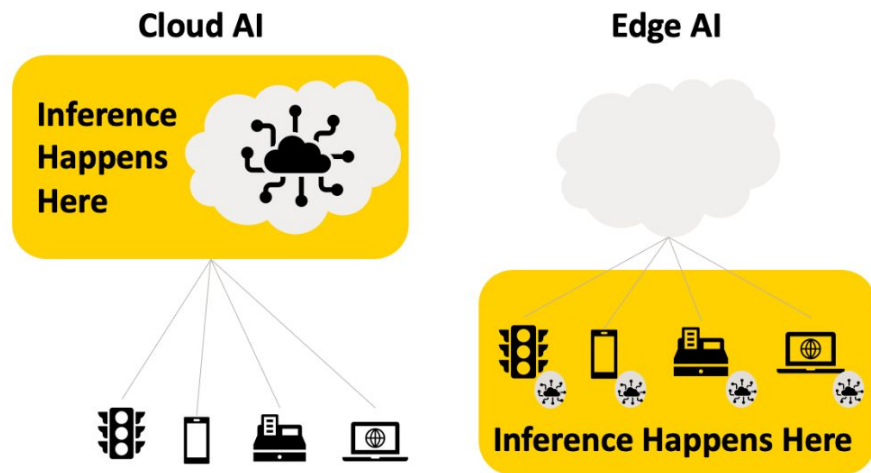
Perform inference locally without sending user data to the cloud



Works offline

No need for a network connection or running a service in the cloud

K Cloud to Edge



	Date of original paper	Energy consumption (kWh)	Carbon footprint (lbs of CO2e)	Cloud compute cost (USD)
Transformer (65M parameters)	Jun, 2017	27	26	\$41-\$140
Transformer (213M parameters)	Jun, 2017	201	192	\$289-\$981
ELMo	Feb, 2018	275	262	\$433-\$1,472
BERT (110M parameters)	Oct, 2018	1,507	1,438	\$3,751-\$12,571
Transformer (213M parameters) w/ neural architecture search	Jan, 2019	656,347	626,155	\$942,973-\$3,201,722
GPT-2	Feb, 2019	-	-	\$12,902-\$43,008

Note: Because of a lack of power draw data on GPT-2's training hardware, the researchers weren't able to calculate its carbon footprint.

Training a single AI model can emit as much carbon as five cars in their lifetimes

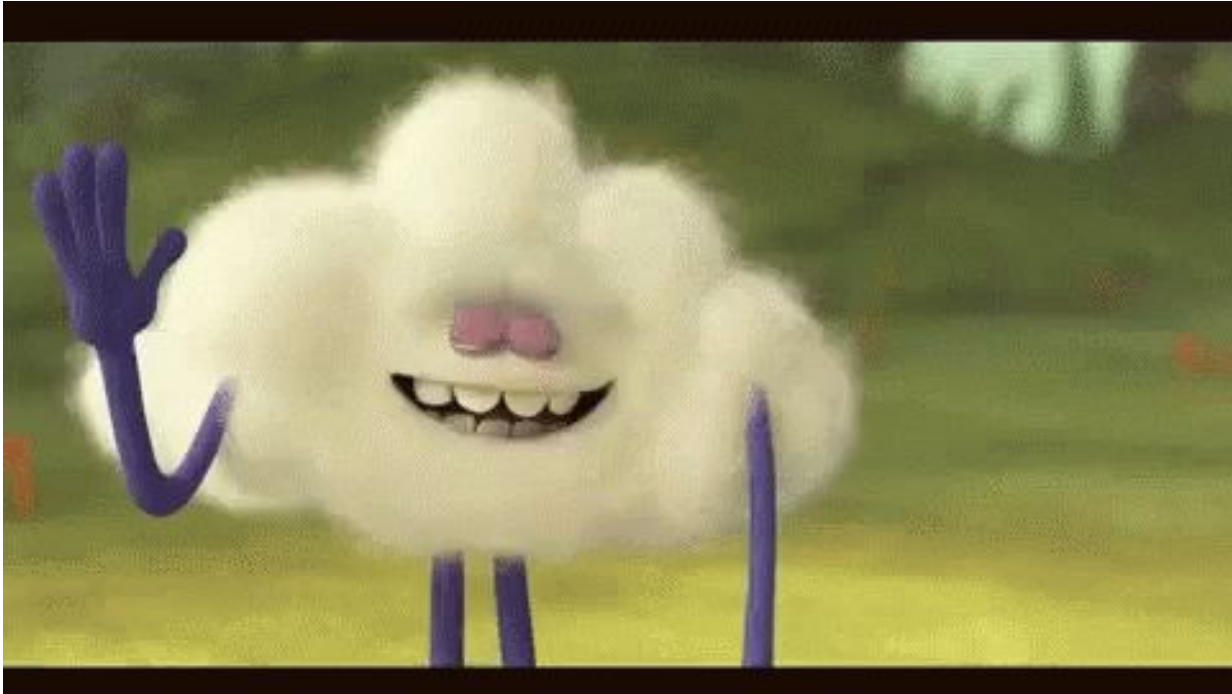
K A Simple Keras Model

```
from tensorflow import keras
```

```
model = keras.Sequential([  
    keras.layers.Dense(units=16, activation='relu', input_shape=(2,)),  
    keras.layers.Dropout(rate=0.2),  
    keras.layers.Dense(units=8, activation='relu'),  
    keras.layers.Dropout(rate=0.2),  
    keras.layers.Dense(units=3, activation='softmax')  
])
```

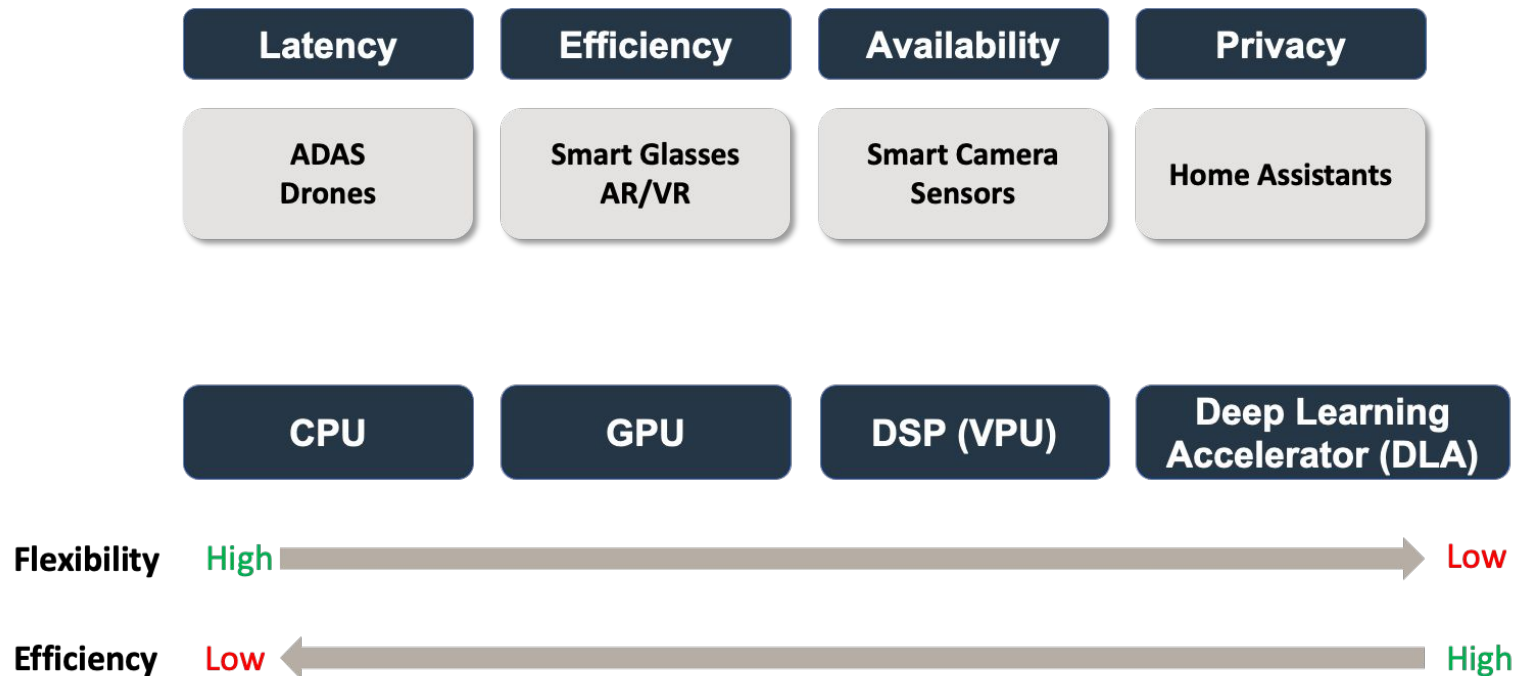
Train it! Deploy it!

```
model.save('my_model.h5')
```



Source: Tenor

K Challenges in EdgeAI Deployment



K Understanding Key Metrics

Size Reduction

Accuracy

Throughput

Latency

Memory Bandwidth

**Power & Energy
Consumption**

K Model Optimization Techniques

“Model optimization refers to the process of improving the performance, efficiency, and resource utilization of a trained machine learning model. This involves various techniques aimed at achieving better results while using fewer computational resources.”

We'll talk about

- Knowledge Distillation
- Quantization

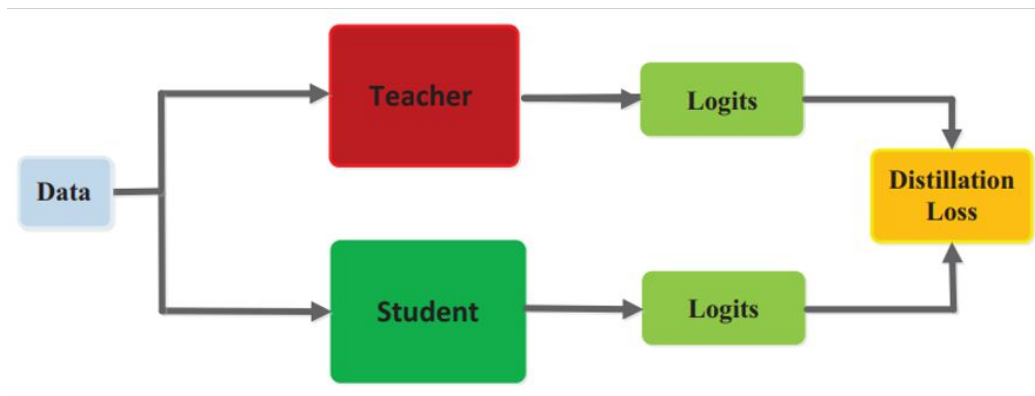


Optimize models to reduce size, latency and power for negligible loss in accuracy

Source: TensorFlow Blog

K Knowledge Distillation

“A technique where a smaller model (student) is trained to reproduce the behavior of a larger model (teacher) or an ensemble of models, often leading to a compact model with comparable performance”



K Knowledge Distillation

Size Reduction

Up to 10% – 50%
w.r.t teacher model

Throughput

2x – 10x higher
w.r.t teacher model

Memory Bandwidth

50% – 75% reduction,
depending on student model
design

Accuracy

1% - 5% drop,
w.r.t teacher model

Latency

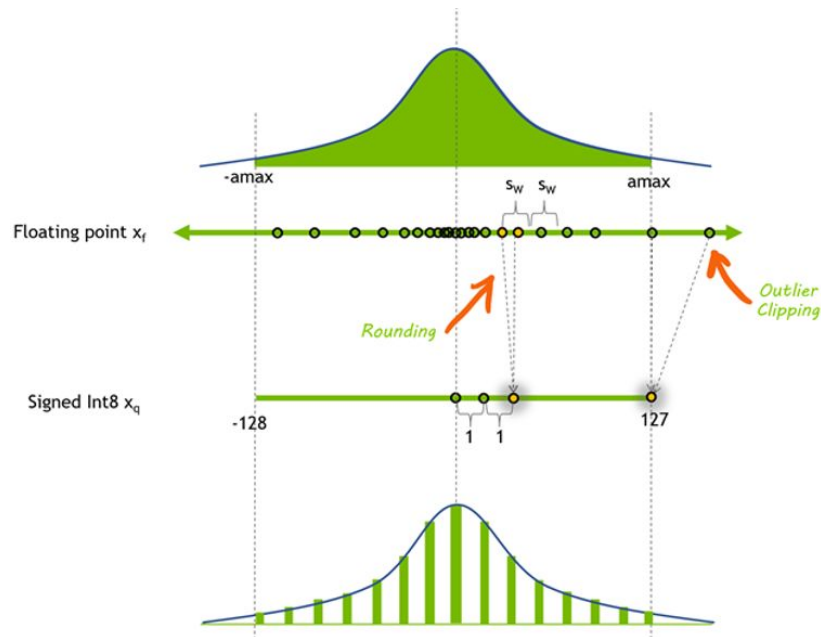
2x – 10x reduction

Power& Energy Efficiency

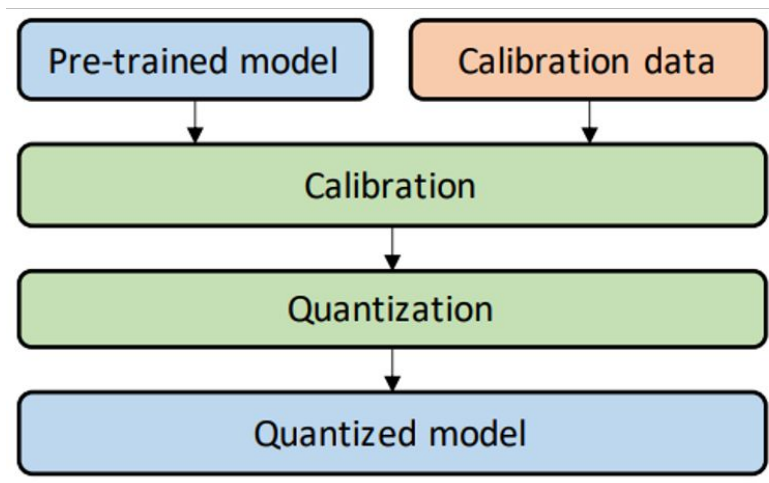
2x – 10x lower consumption

K Quantization

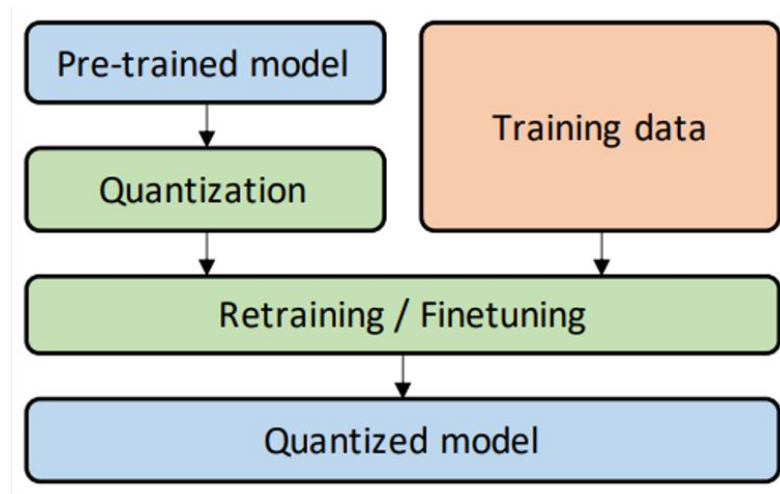
“The process of reducing the numerical precision of model parameters by mapping it from a large number of possible values to a reduced set of values”



K Quantization



Post Training Quantization



Quantization Aware Training

K Quantization

Size Reduction

Up to 50% – 75%
w.r.t FP32 model

Throughput

2x – 4x higher

Memory Bandwidth

50% – 75% reduction,
depending on bit-width

Accuracy

1% - 5% drop,
depending on bit-width and
quantization technique

Latency

2x – 3x reduction

Power& Energy Efficiency

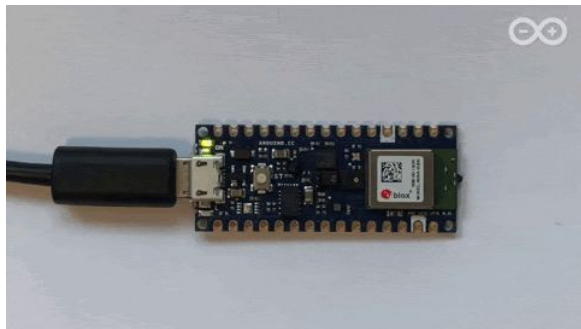
2x – 3x lower consumption

K

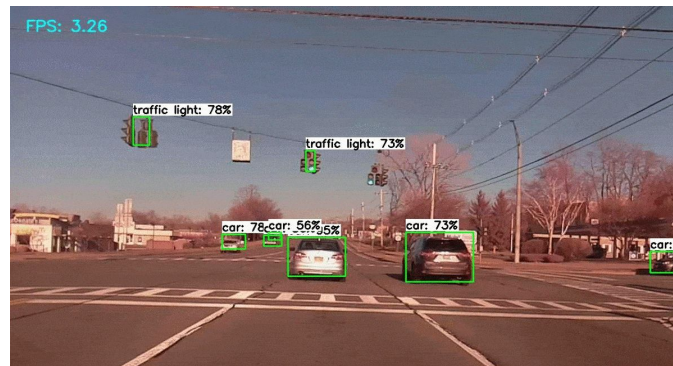
**Few of us after looking
at Model Optimization
techniques:**



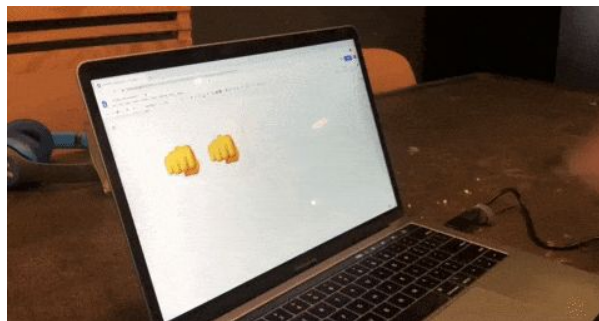
K Case Study 1: AI on embedded



Voice Activation

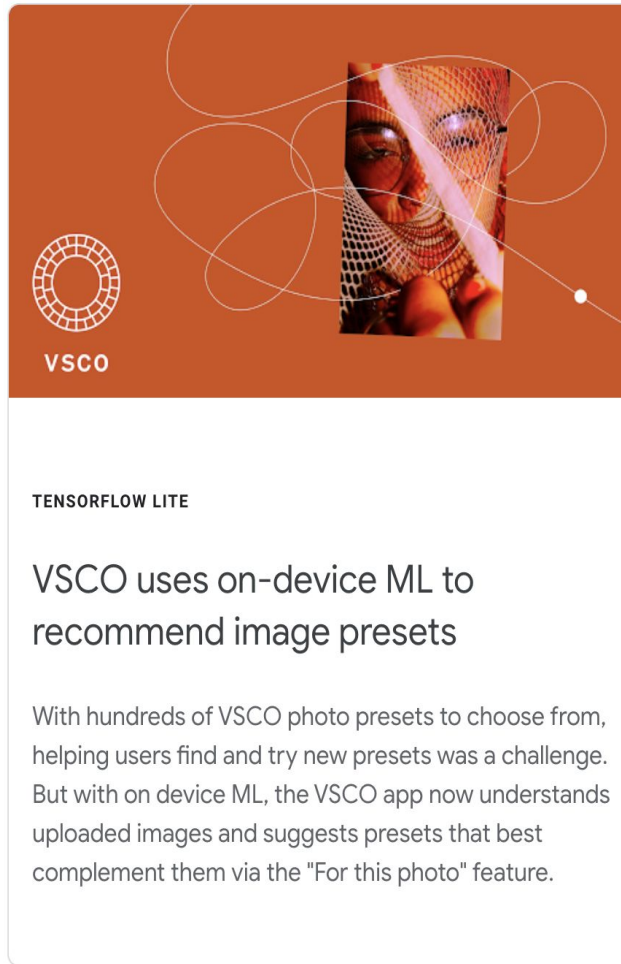


Object detection

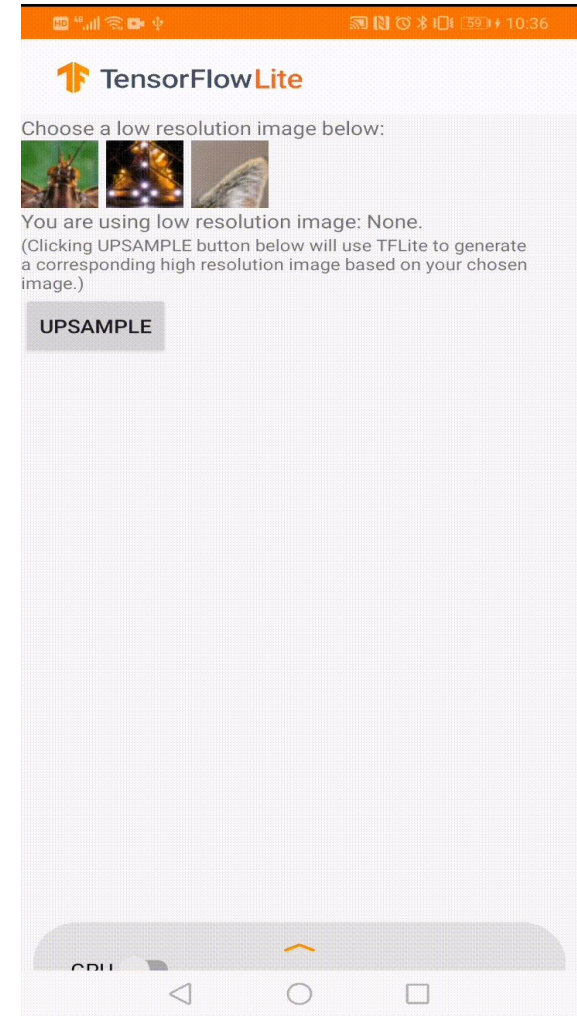


Gesture Recognition

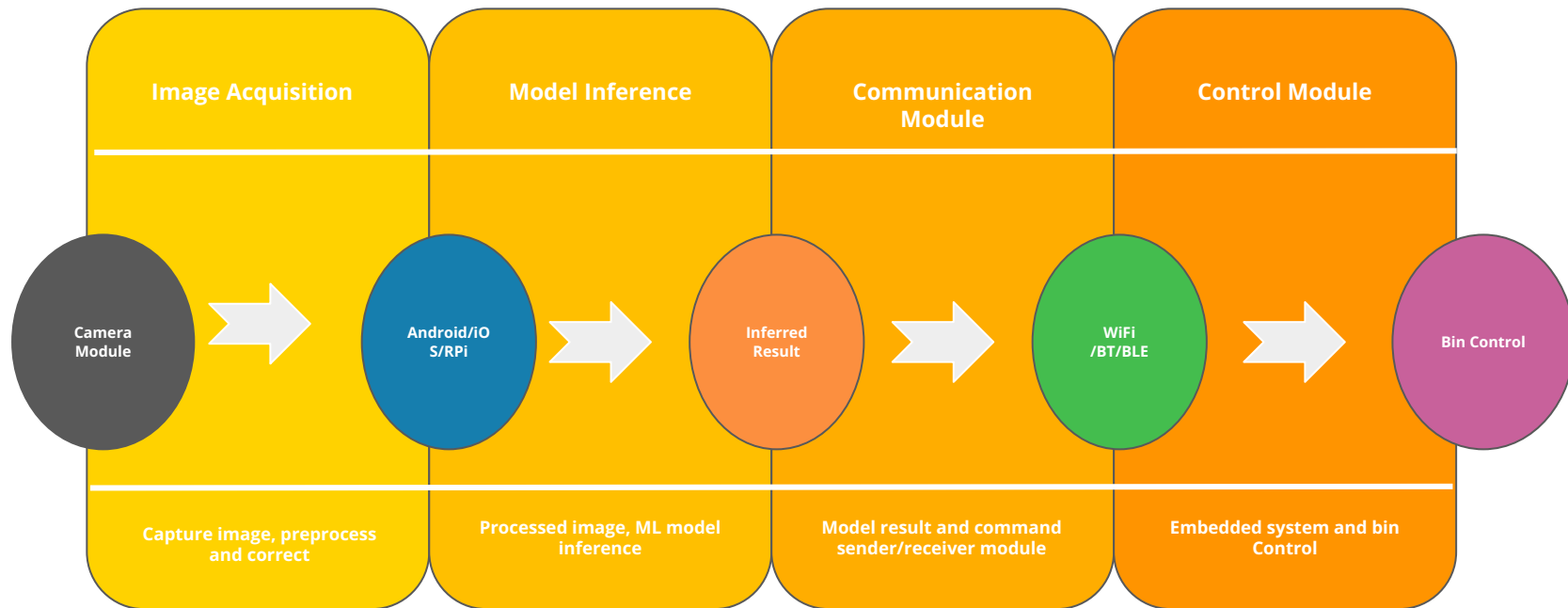
Case Study 2: On Mobile



Source: TF.org



K Let's club the two: Mobile & Embedded





CAPTURE IMAGE

Waste type:	Biodegradable
Waste Subclass:	food waste, 98%
Motor 1 Status:	ON
Motor 2 Status:	OFF




K Check out the complete project here 



K Resources for you







New to ML?



Train a computer to recognize your own images, sounds, & poses.

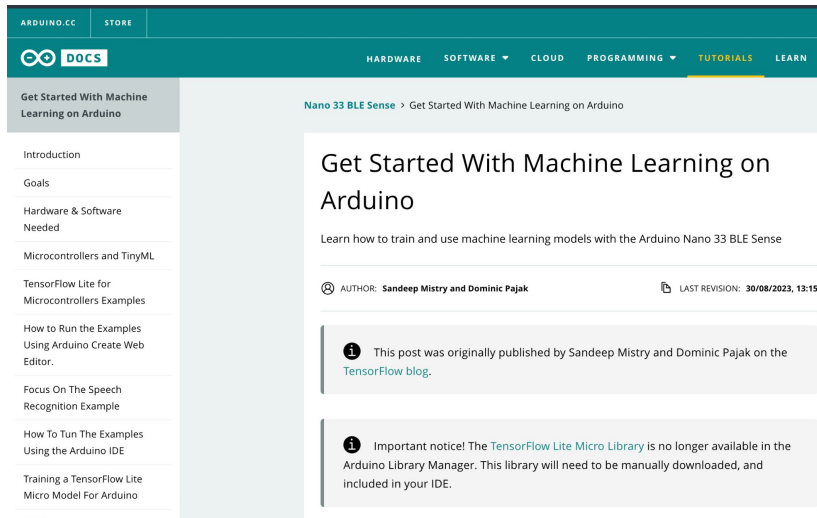
A fast, easy way to create machine learning models for your sites, apps, and more – no expertise or coding required.

[Get Started](#)

<https://teachablemachine.withgoogle.com/>

New to ML on edge?



The screenshot shows the Arduino.cc documentation page for "Get Started With Machine Learning on Arduino". The page is part of the "NANO 33 BLE SENSE" tutorial series. It includes a table of contents on the left with links to Introduction, Goals, Hardware & Software Needed, Microcontrollers and TinyML, TensorFlow Lite for Microcontrollers Examples, How to Run the Examples Using Arduino Create Web Editor, Focus On The Speech Recognition Example, How To Tun The Examples Using the Arduino IDE, and Training a TensorFlow Lite Micro Model For Arduino. The main content area has the title "Get Started With Machine Learning on Arduino" and a subtitle "Learn how to train and use machine learning models with the Arduino Nano 33 BLE Sense". It also features a note about the original author (Sandeep Mistry and Dominic Pajak) and a warning that the TensorFlow Lite Micro Library is no longer available in the Arduino Library Manager.

<https://docs.arduino.cc/tutorials/nano-33-ble-sense/get-started-with-machine-learning#introduction>



Aashi Dutt

Organizer, TensorFlow User
Group Chandigarh

Connect With Me 🖐️



@AashiDutt

Join TFUG Chandigarh Chapter 🎉



@TFUGChandigarh



@tfugchandigarh