Experts

# Edge ML Explosion

- Lower latency & close knit interactions

Experts

# Edge ML Explosion

- Lower latency & close knit interactions
- Network Connectivity

# Edge ML Explosion

- Lower latency & close knit interactions
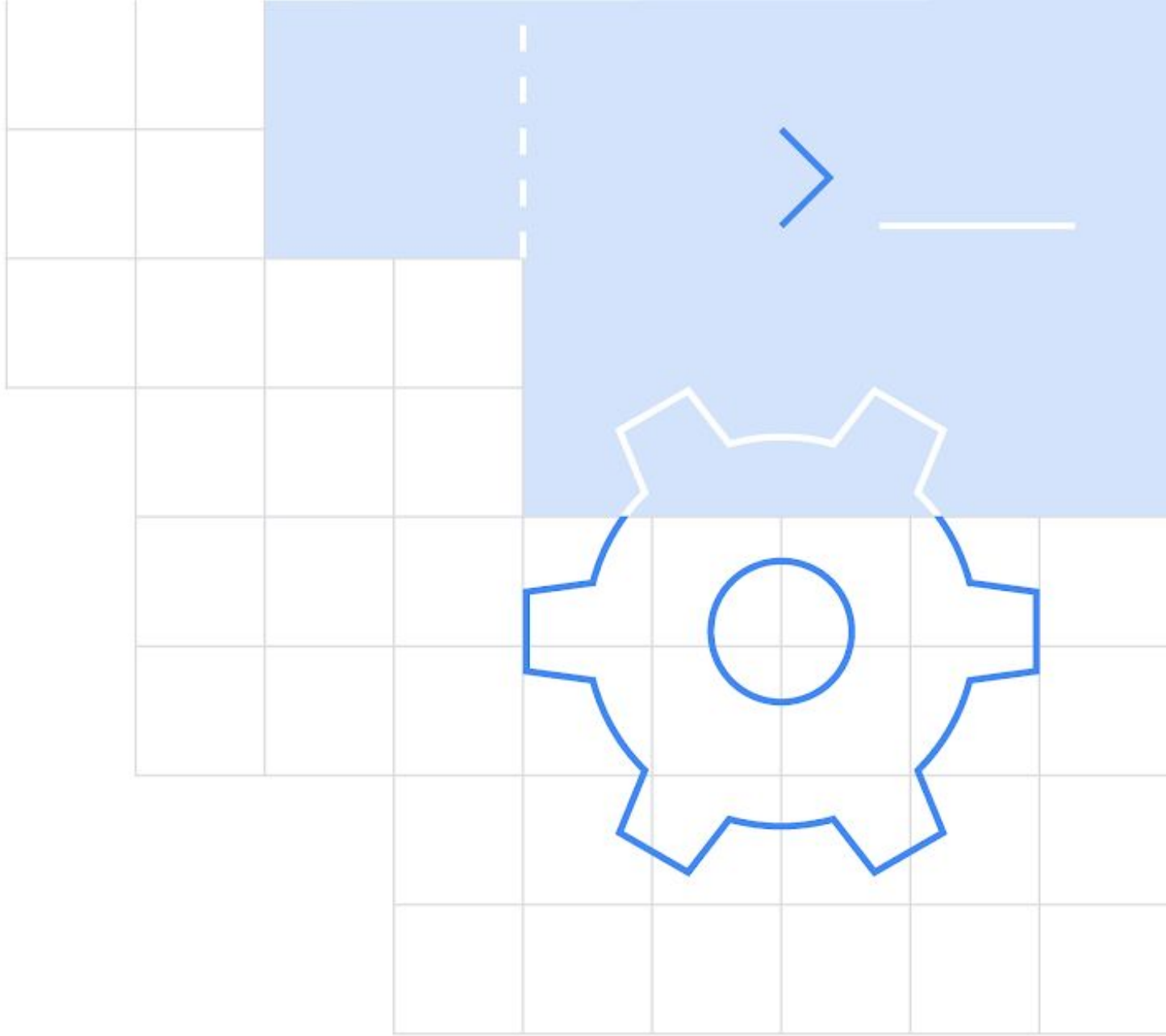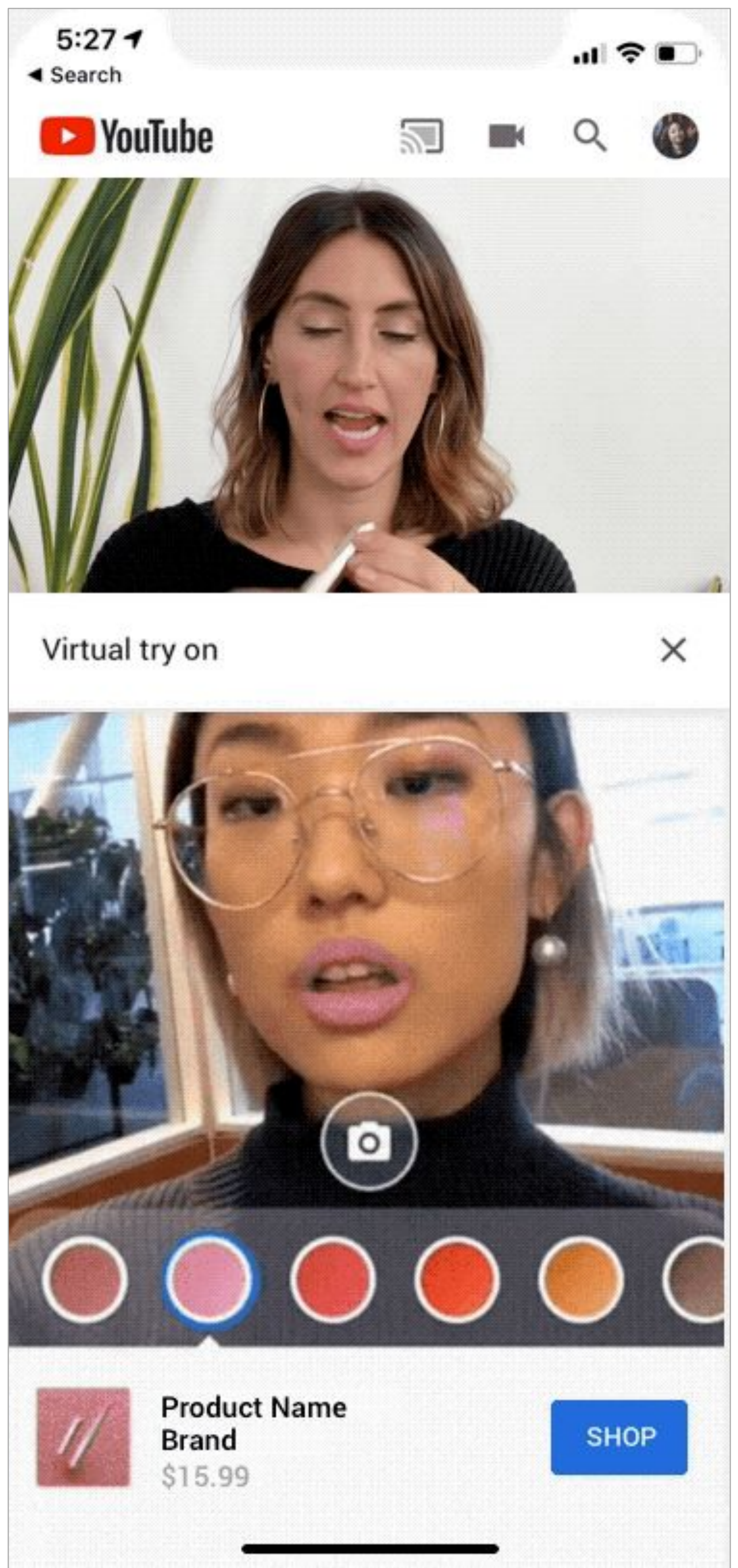- Network Connectivity
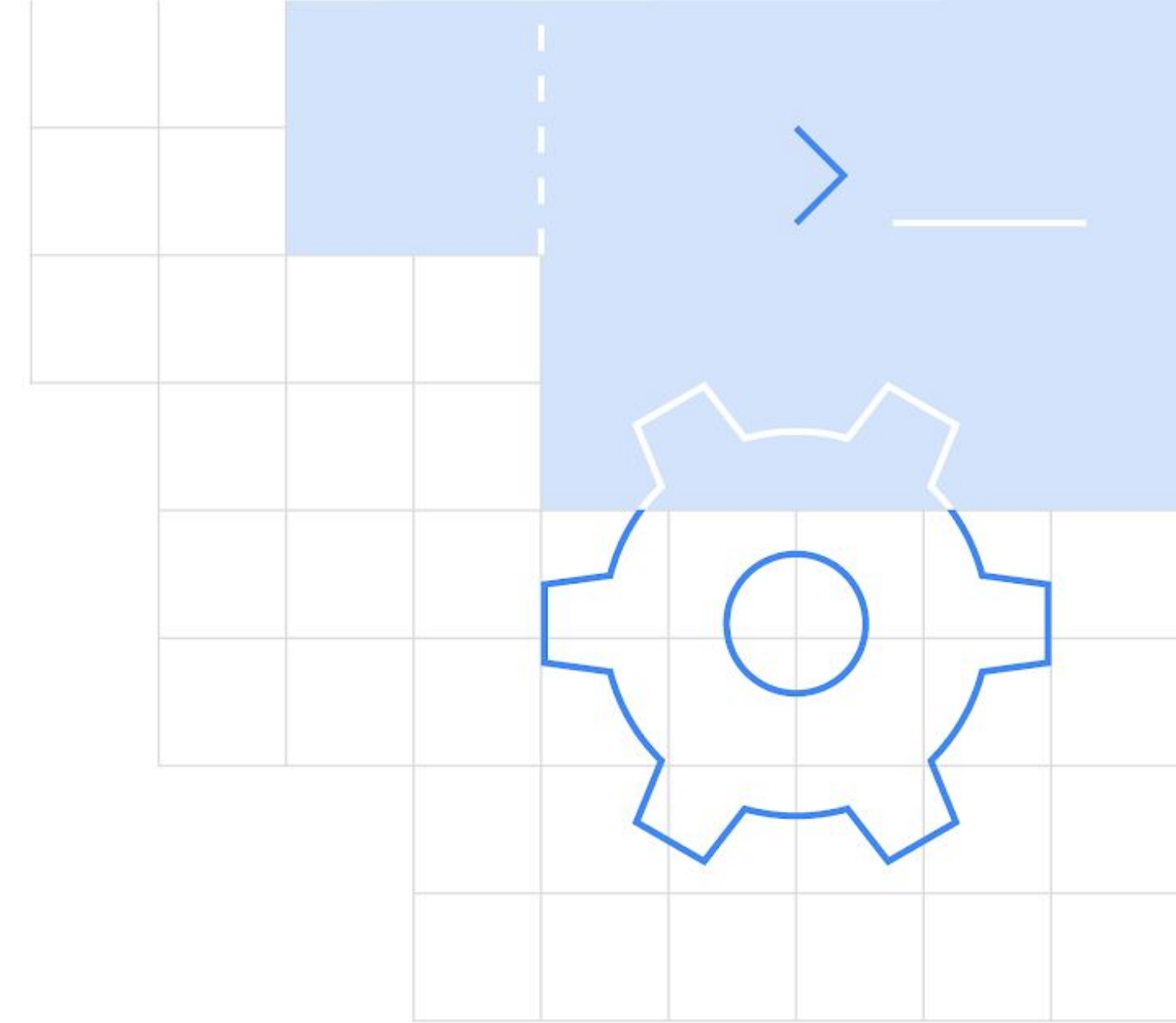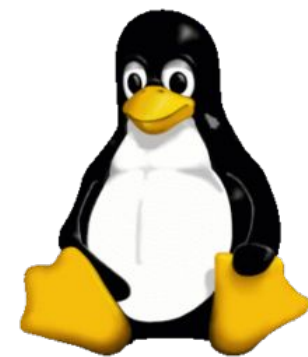- Privacy preserving

# Challenges

- Limited compute power

- Limited memory

- Battery consumption

- App size

**TensorFlow Lite** is a production ready, cross-platform framework for deploying ML on **mobile devices and embedded systems**
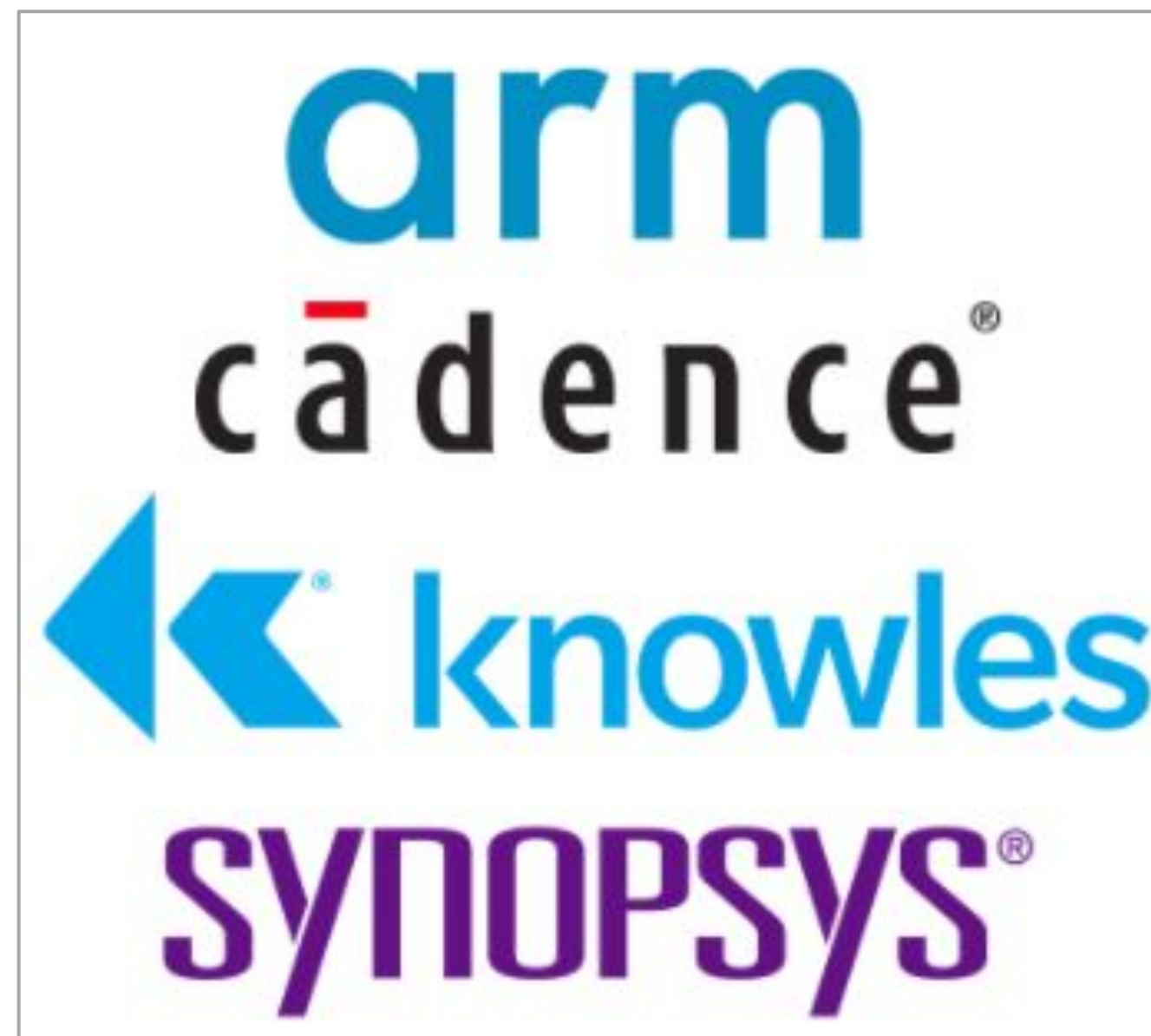
# Folks, hacking with embedded stuff & microcontroller stuff

- [Build TensorFlow Lite for ARM64 boards](#)
- [Build TensorFlow Lite for Raspberry Pi](#)

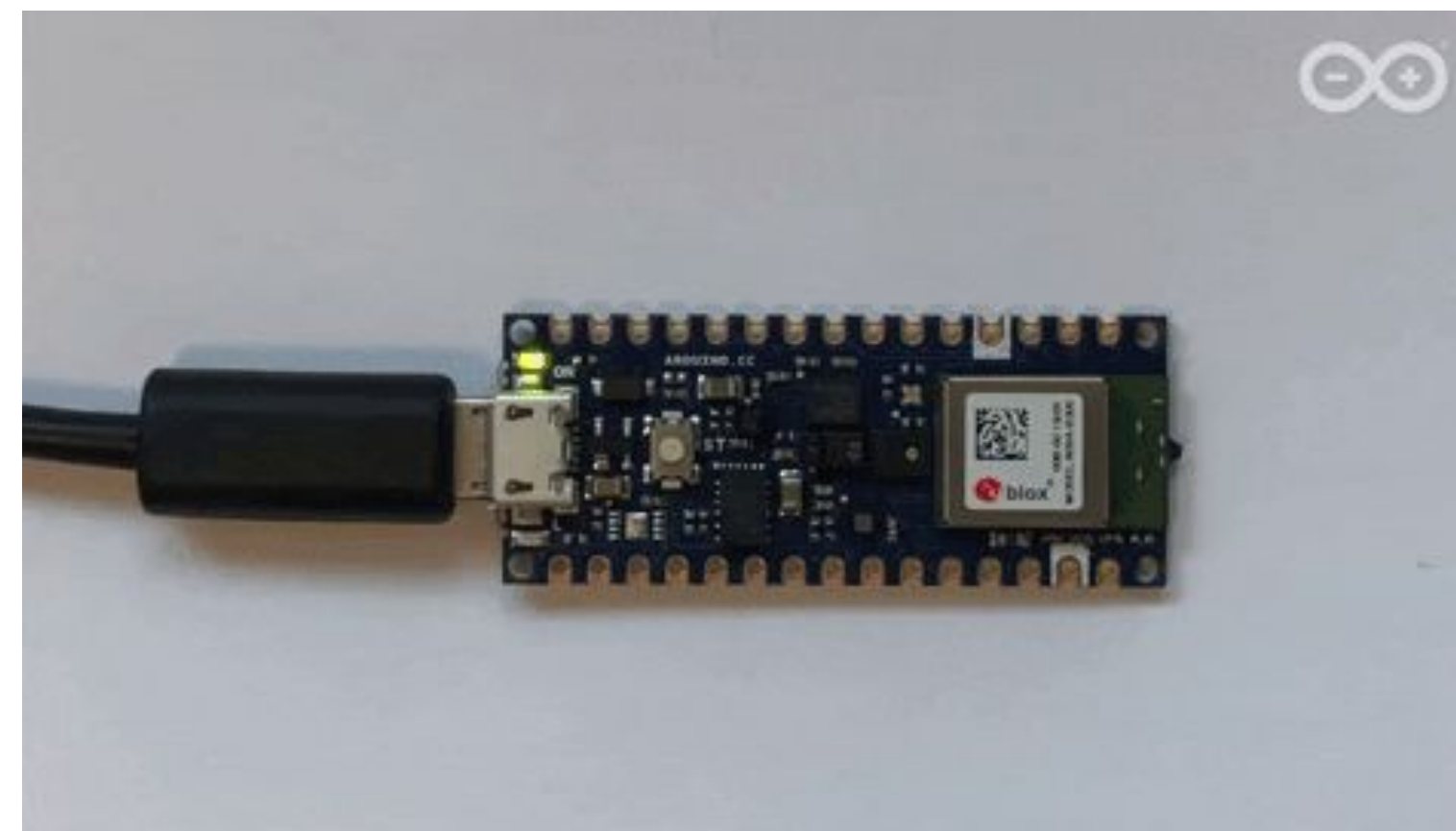# Folks, hacking with embedded stuff & microcontroller stuff

- Tremendous speed up with Edge TPU compatible TF Lite models



Check out here: [Edge TPU performance benchmarks | Coral](#)

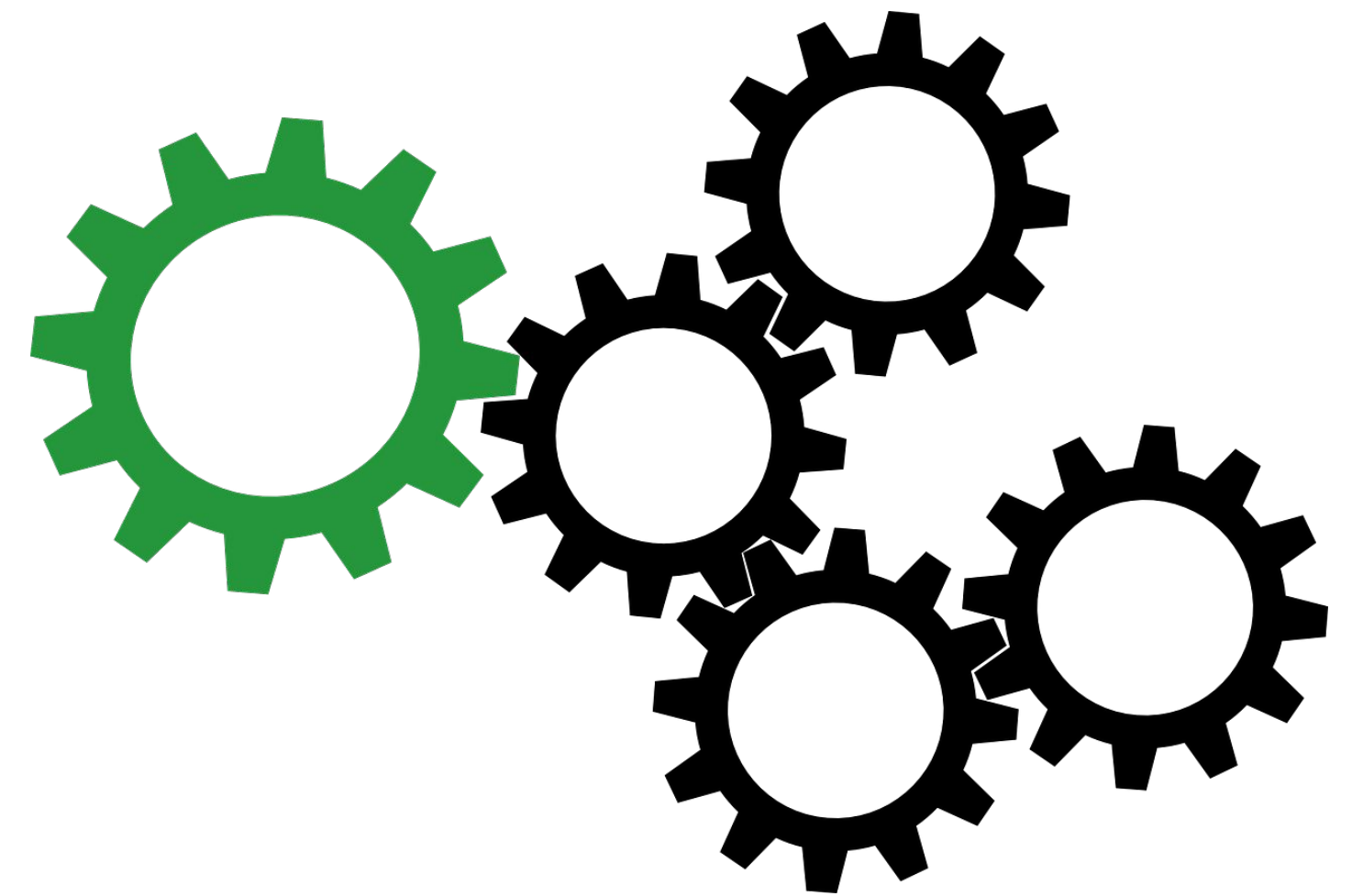# Folks, hacking with embedded stuff & microcontroller stuff

- **Launch of official Arduino library** - run example code directly from desktop and web IDEs onto Arduino hardware

- **Speech detection in 5 minutes** - open source models available to get started quickly on Arduino
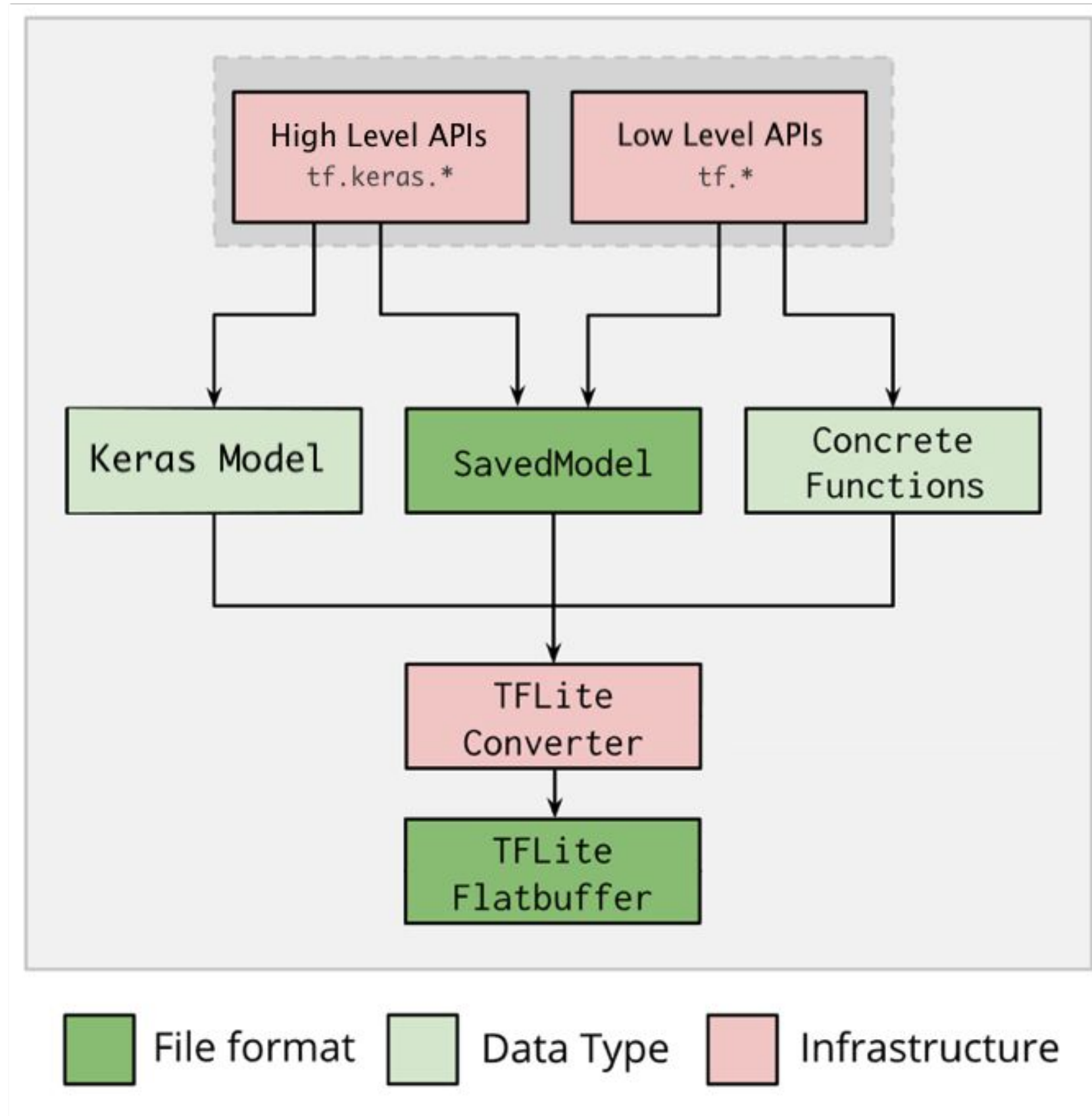


https://www.tensorflow.org/lite/microcontrollers

# Workflow

1. Train a TensorFlow model

2. Convert to TensorFlow Lite format

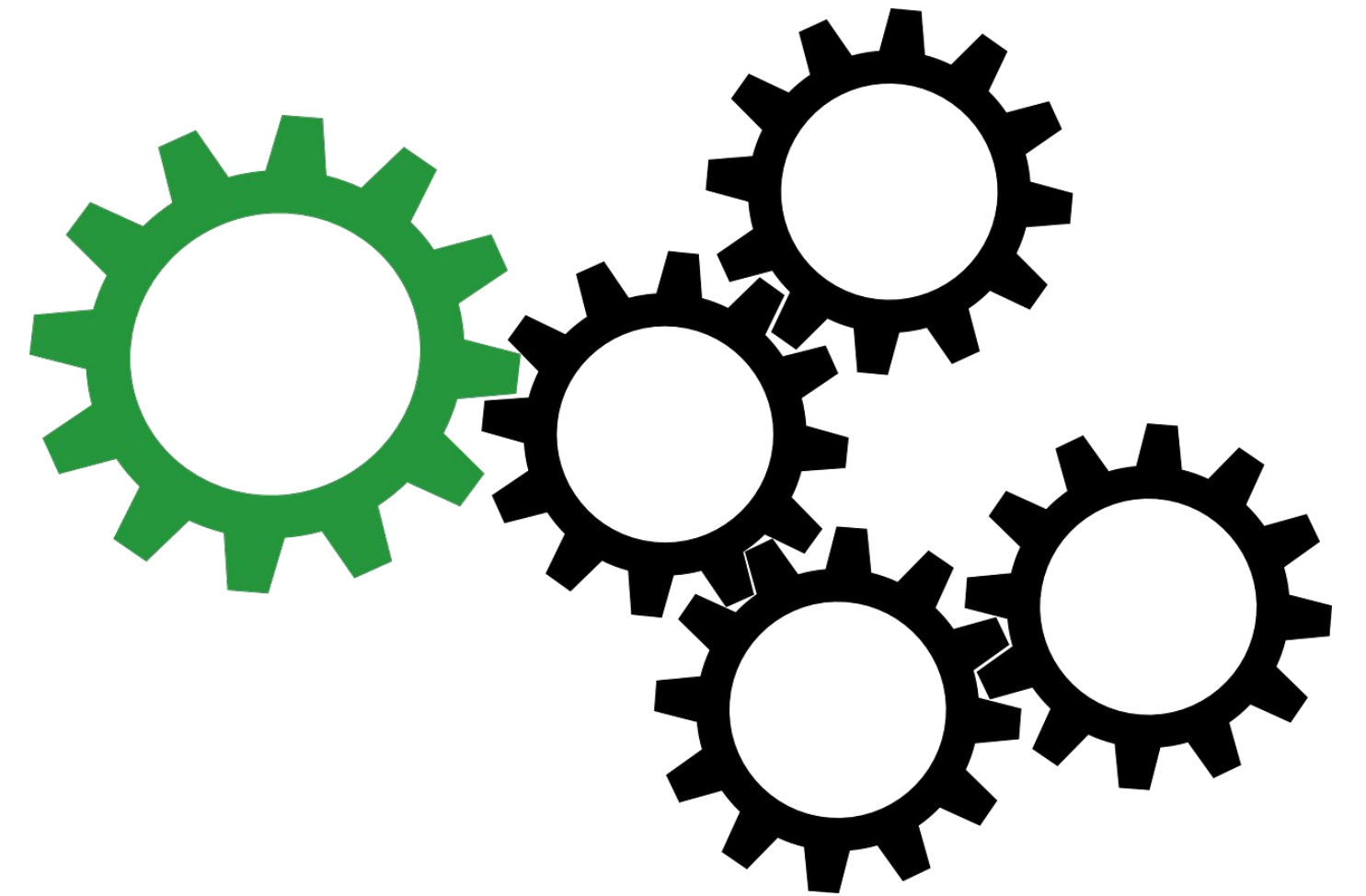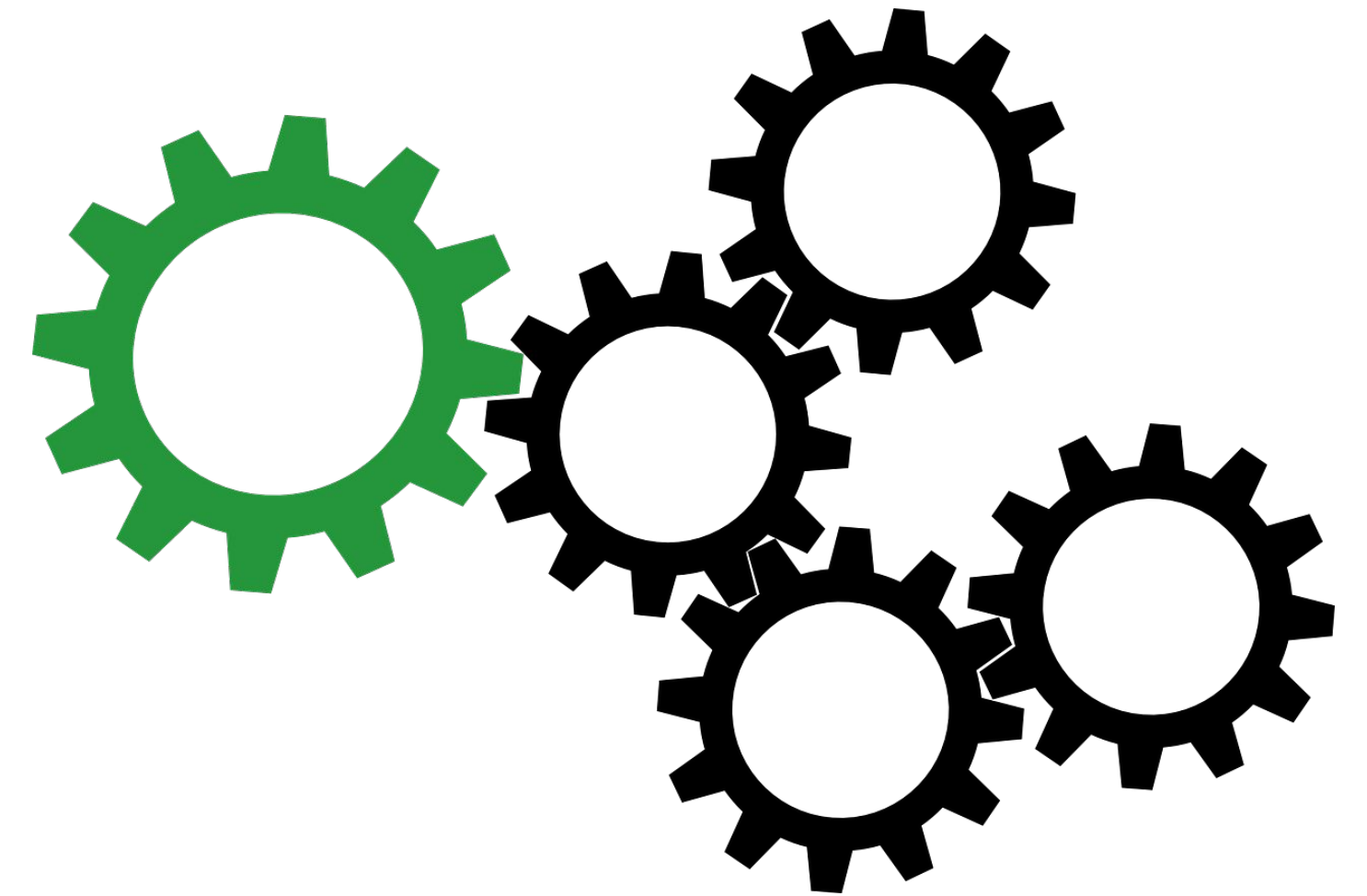3. Deploy and run on edge device

# Workflow

# Workflow

FlatBuffers is an efficient cross platform serialization library for C++, C#, C, Go, Java, Kotlin, JavaScript, Lobster, Lua, TypeScript, PHP, Python, Rust and Swift. It was originally created at Google for game development and other performance-critical applications.

# Why not use Protocol Buffers?

Protocol Buffers is indeed relatively similar to FlatBuffers, with the primary difference being that FlatBuffers does not need a parsing/unpacking step to a secondary representation before you can access data, often coupled with per-object memory allocation. The code is an order of magnitude bigger, too.
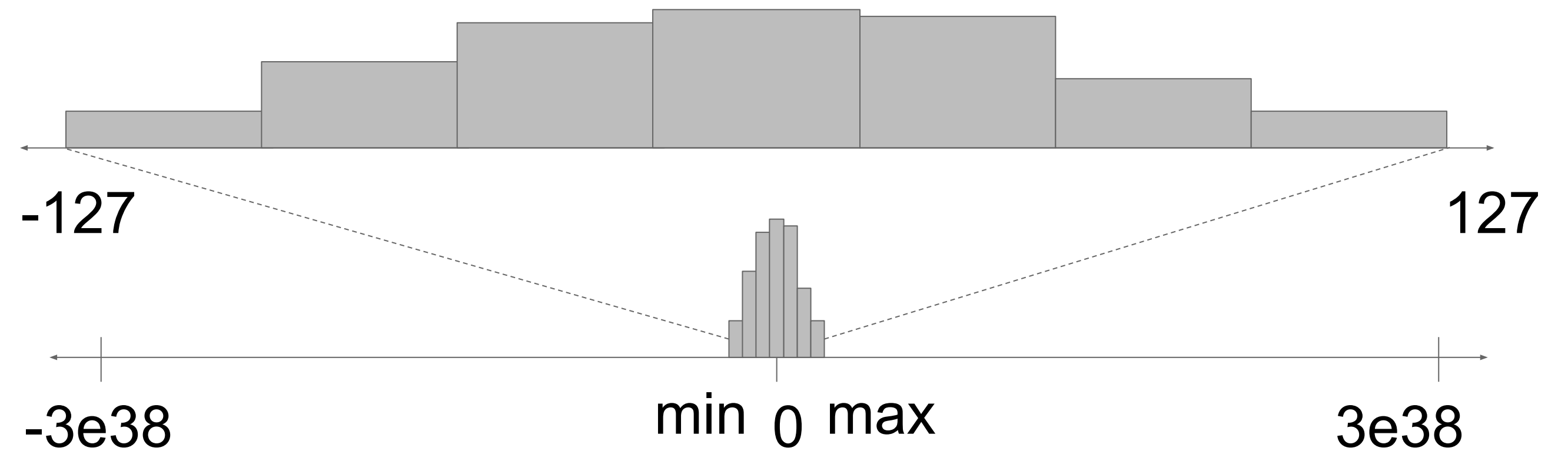
# TensorFlow Lite

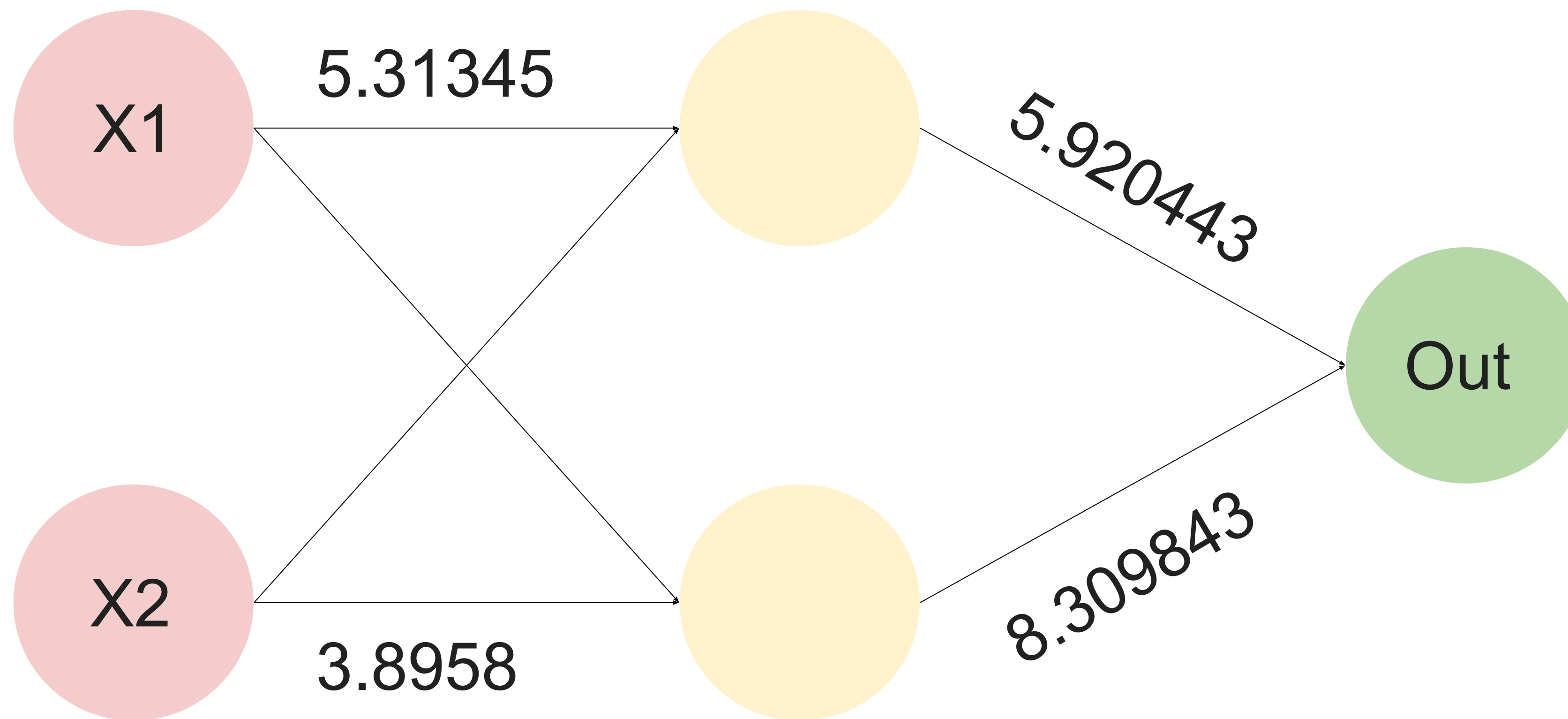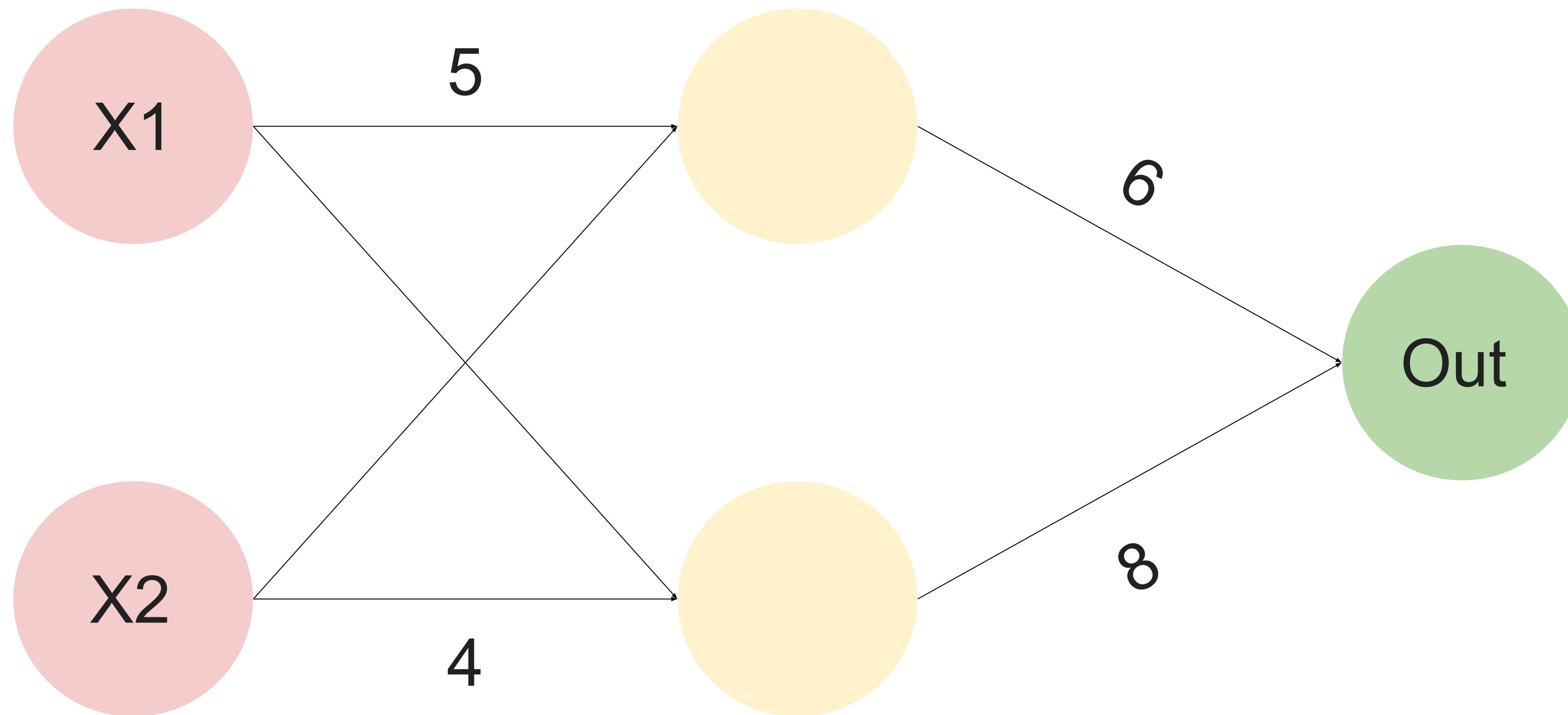Lets Go Forward & witness the
magic of TFLite

# Quantization



-127                                                     127

-3e38                          min 0 max                   3e38

Reduce number of bits for model weights and activations

# Q&A



THANK YOU FOR ATTENTION

Q&A SESSION

# Scan to access Slides & Code



Experts

# Thank You!

Bhavesh Bhatt
@_bhaveshbhatt