Aneel Damaraju    Abhishek Sharma    Jason Jabbour

# Tensorflow Lite Micro

## CS 249 TinyML

# **Roadmap**

- High-Level Tensorflow Lite Micro Sales Pitch

- Review of Important Concepts

- Under the hood of Tensorflow Lite Micro
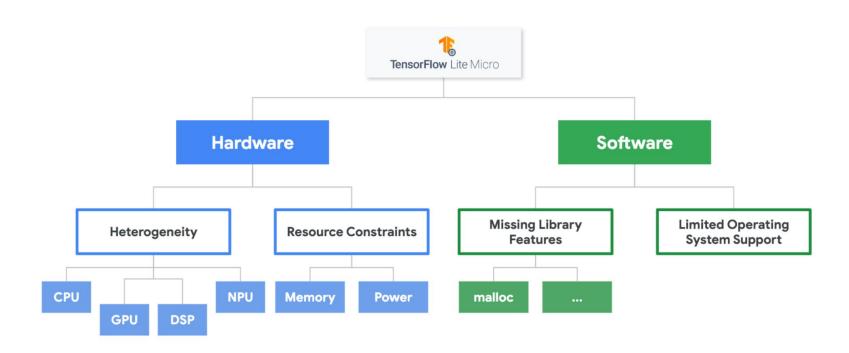
- Tensorflow Lite Micro Demo

# ▶ **Roadmap**

- High-Level Tensorflow Lite Micro Sales Pitch

- Review of Important Concepts

- Under the hood of Tensorflow Lite Micro
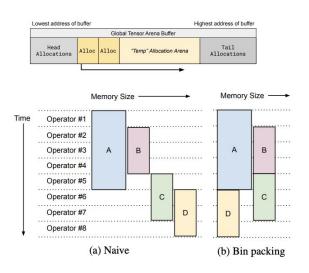
- Tensorflow Lite Micro Demo
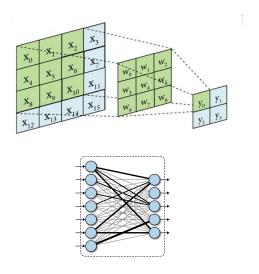
# What is TFLite Micro?

# What is TFLite Micro?

**Hardware agnostic**, with minimal external dependencies

**Minimal memory overhead** for trained models

Uses TFLite to cover a **wide variety of model architectures**



...and **many more!**

# TFLM in a Machine Learning Pipeline

| Collect Data | Preprocess Data | Design a Model | Train a Model | Evaluate Optimize | Convert Model | Deploy Model | Make Inferences |
|---|---|---|---|---|---|---|---|

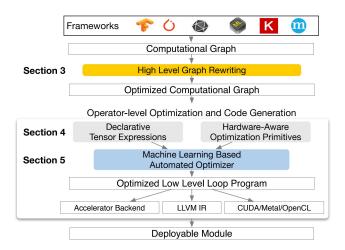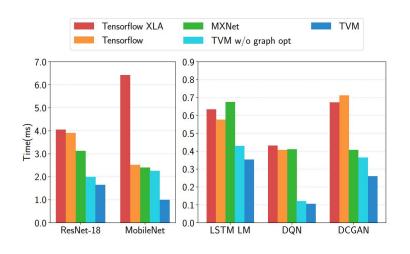| TensorFlow | TensorFlow Lite | TensorFlow Lite Micro |
|---|---|---|

# Ecosystem of TinyML Frameworks

# Ecosystem of TinyML Frameworks
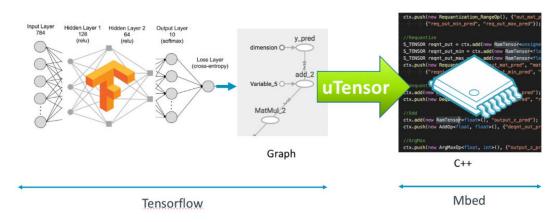


SVM System Overview



TVM Time Comparison

Compiler that exposes **graph-level** and **operator-level optimizations** to provide performance portability to deep learning workloads across diverse hardware back-ends
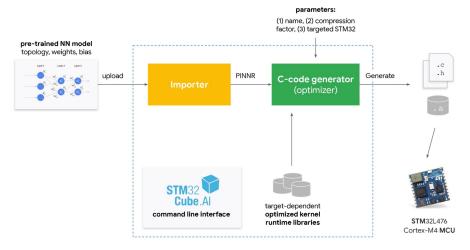
## Ecosystem of TinyML Frameworks



TVM Time Comparison

Framework that **converts ML models** to **C++ source files** to increase human-readability and ease of editing with a core runtime **size of 2KB**

# Ecosystem of TinyML Frameworks



STM32Cube.AI Workflow

Generates an internal representation of a network, **generates optimized code** using a **compiler** and **optimized kernel libraries**, generates executable code runnable on a microcontroller

# ➤ **Roadmap**

- High-Level Tensorflow Lite Micro Sales Pitch

- Review of Important Concepts

- Under the hood of Tensorflow Lite Micro

- Tensorflow Lite Micro Demo

# Embedded System

Computer hardware system with software that:
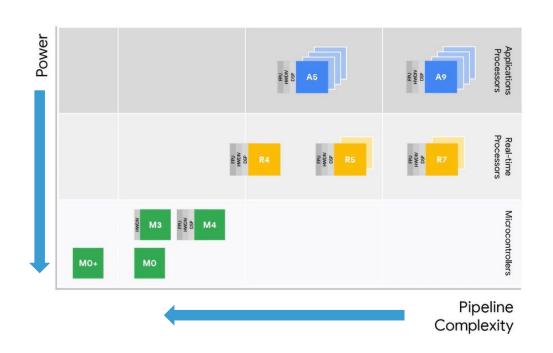
- is managed by **microcontrollers**, digital signal processors (DSP), application specific integrated circuit (ASIC), or field-programmable gate arrays (FPGA)
- is designed to perform a dedicated function under **tight constraints**, either as an independent system or as a part of a large system
- stores programming instructions in **read-only memory** or flash memory chips

# Embedded Systems: Hardware

| Board | MCU / ASIC | Clock | Memory | Sensors | Radio |
|---|---|---|---|---|---|
| Himax<br>WE-I Plus EVB | HX6537-A<br>32-bit EM9D DSP | 400 MHz | 2MB flash<br>2MB RAM | Accelerometer, Mic, Camera | None |
| Arduino<br>Nano 33 BLE Sense | 32-bit<br>nRF52840 | 64 MHz | 1MB flash<br>256kB RAM | Mic, IMU, Temp, Humidity, Gesture, Pressure, Proximity, Brightness, Color | BLE |
| SparkFun<br>Edge 2 | 32-bit<br>ArtemisV1 | 48 MHz | 1MB flash<br>384kB RAM | Accelerometer, Mic, Camera | BLE |
| Espressif<br>EYE | 32-bit<br>ESP32-D0WD | 240 MHz | 4MB flash<br>520kB RAM | Mic, Camera | WiFi, BLE |

# Embedded Systems: Hardware



Power

Pipeline Complexity

Applications Processors

Real-time Processors

Microcontrollers

A5, A9, R4, R5, R7, M3, M4, M0+, M0

Low Cost

Low Power

# Embedded Systems: Software

Mobile OS

Embedded Sys.

# Embedded Systems: Software



Software

TF Micro Application

Arduino

mbed OS

Nano 33 BLE Sense
**Hardware**

Windows
iOS
mac OS
Android
Linux

Mobile OS

RTOS
arm MBED OS

Embedded Sys.

# Microprocessor vs Microcontroller

| Microprocessor | Micro Controller |
|---|---|
| Microprocessor is heart of Computer system. | Micro Controller is a heart of embedded system. |
| It is just a processor. Memory and I/O components have to be connected externally | Micro controller has external processor along with internal memory and i/O components |
| Since memory and I/O has to be connected externally, the circuit becomes large. | Since memory and I/O are present internally, the circuit is small. |
| Cannot be used in compact systems and hence inefficient | Can be used in compact systems and hence it is an efficient technique |
| Cost of the entire system increases | Cost of the entire system is low |
| Due to external components, the entire power consumption is high. Hence it is not suitable to used with devices running on stored power like batteries. | Since external components are low, total power consumption is less and can be used with devices running on stored power like batteries. |
| Most of the microprocessors do not have power saving features. | Most of the micro controllers have power saving modes like idle mode and power saving mode. This helps to reduce power consumption even further. |
| Since memory and I/O components are all external, each instruction will need external operation, hence it is relatively slower. | Since components are internal, most of the operations are internal instruction, hence speed is fast. |
| Microprocessor have less number of registers, hence more operations are memory based. | Micro controller have more number of registers, hence the programs are easier to write. |
| Microprocessors are based on von Neumann model/architecture where program and data are stored in same memory module | Micro controllers are based on Harvard architecture where program memory and Data memory are separate |
| Mainly used in personal computers | Used mainly in washing machine, MP3 players |

# Tensorflow – Tensorflow Lite – Tensorflow Lite Micro

Library = 400MB

# **Tensorflow –** Tensorflow Lite – Tensorflow Lite Micro



Virtual Memory



Multi-Threading

# **Tensorflow – Tensorflow Lite – Tensorflow Lite Micro**

| | Microprocessor | > | Microcontroller |
|---|---|---|---|
| Platform | | | |
| Compute | 1GHz–4GHz | ~10X | 1MHz–400MHz |
| Memory | 512MB–64GB | ~10000X | 2KB–512KB |
| Storage | 64GB–4TB | ~100000X | 32KB–2MB |
| Power | 30W–300W | ~1000X | 150µW–23.5mW |

# Tensorflow – **Tensorflow Lite** – Tensorflow Lite Micro



TFLite Conversion Process

**Tensorflow** – **Tensorflow Lite** – **Tensorflow Lite Micro**

Library = 1MB

# Tensorflow – Tensorflow Lite – **Tensorflow Lite Micro**

Library = 16KB

# Background: Memory Hierarchy



"SRAM"

"DRAM"

"SSD"

"HDD"

Image via teachbook.com.au

# Background: Memory in MCUs



Laptop

MacBook M1

Cache (24MB)

DRAM (16GB)  SSD (512 GB)

Volatile Memory:
Data is lost when the computer is turned off

Non-volatile Memory:
Data persists even when the computer is turned off

# Background: Memory in MCUs



**Laptop**

MacBook M1

Cache (24MB)

DRAM (16GB) · SSD (512 GB)

**On-Chip**

Cortex-M

Cache (4KB)

SRAM (256KB) · eFlash (1 MB)

Volatile Memory:
Data is lost when the computer is turned off

Non-volatile Memory:
Data persists even when the computer is turned off
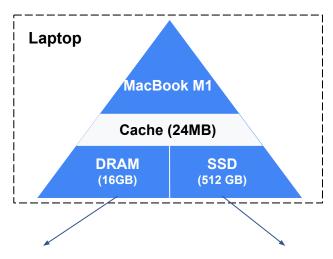
# ► **Interpreters** vs Compilers

Python: **Interpreted** language



```
Python 3.8.5 Shell
File   Edit   Shell   Debug   Options   Window   Help

Type "help", "copyright", "credits" or "l
>>> print('Hello, World!')
Hello, World!
>>> |
```

# **Interpreters** vs Compilers



Interpreter Memory Management

# Interpreters vs **Compilers**

Python: **Interpreted** language

C++: **Compiled** language



```
Python 3.8.5 Shell
File  Edit  Shell  Debug  Options  Window  Help

Type "help", "copyright", "credits" or "l
>>> print('Hello, World!')
Hello, World!
>>>
```

```cpp
helloWorld.cpp                          helloWorld.cpp
1  #include <iostream>
2
3  int main(){
4      std::cout << "Hello, World!" << std::endl;
5      return 0;
6  }
7
```

# **Roadmap**

- High-Level Tensorflow Lite Micro Sales Pitch

- Review of Important Concepts

- **Under the hood of Tensorflow Lite Micro**

- Tensorflow Lite Micro Demo

# Memory Allocation

# Memory Allocation

## Why care about Memory?

# Memory Allocation

## Why care about Memory?

Memory Limits

On-Chip

Cortex-M

Cache (4KB)

SRAM
(256KB)

eFlash
(1 MB)

# Memory Allocation

## Why care about Memory?

Memory Limits

Long-Running Applications



On-Chip

Cortex-M

Cache (4KB)

SRAM (256KB)  |  eFlash (1 MB)

Heap (Working Memory)

Fragmentation

# Memory Allocation

## Why care about Memory?

| Memory Limits | Long-Running Applications | Lack of OS Support |
|---|---|---|

**On-Chip**

Cortex-M

Cache (4KB)

SRAM (256KB) | eFlash (1 MB)

Fragmentation

~~malloc()~~
~~Virtual Memory~~

# How **TFL Micro** solves these challenges

1.  Ask developers to **supply a contiguous area of memory** to the interpreter, and in return the framework avoids any other memory allocations

```cpp
constexpr int kTensorArenaSize = 2000;
uint8_t tensor_arena[kTensorArenaSize];

...

static tflite::MicroInterpreter static_interpreter(model, resolver,
  tensor_arena, kTensorArenaSize, error_reporting);
```

# How **TFL Micro** solves these challenges

1. Ask developers to **supply a contiguous area of memory** to the interpreter, and in return the framework avoids any other memory allocations

2. Framework **guarantees that it won't allocate from this "arena" after initialization**, so long-running applications won't fail due to fragmentation

# How **TFL Micro** solves these challenges

1. Ask developers to **supply a contiguous area of memory** to the interpreter, and in return the framework avoids any other memory allocations

2. Framework **guarantees that it won't allocate from this "arena" after initialization**, so long-running applications won't fail due to fragmentation

3. Ensures clear budget for the memory used by ML, and that the **framework has no dependency on OS facilities needed by malloc or new**

`uint8_t tensor_arena[`**`kTensorArenaSize`**`]`

| Operator Variables | Interpreter State | Operator Inputs and Outputs |
|---|---|---|

# Loading the Model

```
model = tflite::GetModel(g_model);
```

**g_model** variable is an array of bytes that we exported from TensorFlow

Holds all of the information about the model, its operators, their connections, and the trained weights

## `g_model` FlatBuffer **Format**

**Metadata (version, quantization ranges, etc)**

| Name | Args | Input | Output | Weights |
|--------|------|-------|--------|---------|
| Conv2D | 3x3 | 0 | 1 | 2 |
| FC | - | 1 | 3 | 4 |
| Softmax | - | 3 | 5 | - |

**Weight Buffers**

| Index | Type | Values |
|-------|-------|-------------------|
| 2 | Float | 0.01, 7.45, 9.23, ... |
| 4 | Int8 | 34, 19, 243, ... |
| ... | ... | ... |

# Loading the Model

```
model = tflite::GetModel(g_model);
```

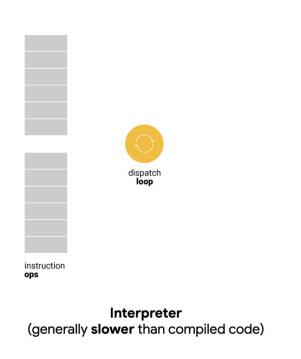**GetModel** call lets us access the model's data elements

This call doesn't do any data movement or copying, so it's very memory efficient, it just creates a thin wrapper to read the information in-place.
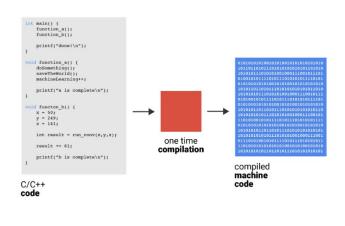
# Loading the Model

```
model = tflite::GetModel(g_model);
```

The resulting model is a class object that lets us access
the model weights and operations.

# Interpreter-based design for machine learning



```
int main() {
    function_a();
    function_b();

    printf("done!\n");
}

void function_a() {
    doSomething();
    saveTheWorld();
    machineLearning++;

    printf("a is complete\n");
}

void functon_b() {
    x = 50;
    y = 249;
    z = 141;

    int result = run_conv(x,y,z);

    result += 61;

    printf("b is complete\n");
}
```

dispatch
**loop**

instruction
**ops**

one time
**compilation**

C/C++
**code**

compiled
**machine
code**

**Interpreter**
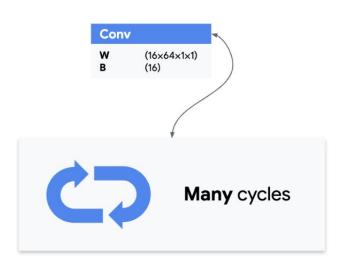(generally **slower** than compiled code)

**Compiler**
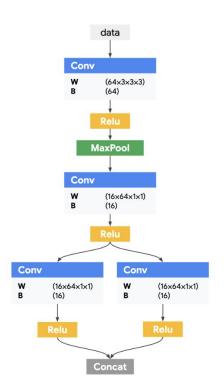(generally **faster** than interpreted code)

# Interpreter-based design for machine learning

- A single line of ML code can take up to **millions of operations**

- The bottleneck is **matrix multiplications** and other operations, so an interpreter is used.

Conv

| W | (16×64×1×1) |
| B | (16) |

**Many** cycles

# ML models have low interpreter overhead



| Model | Total Cycles | Calculation Cycles | Interpreter Overhead |
|-------|--------------|--------------------|--------------------|
| Visual Wake Words (Ref) | 18,990.8K | 18,987.1K | < 0.1% |
| Google Hotword (Ref) | 36.4K | 34.9K | 4.1% |

# Example interpreter execution

```
if (op_type == CONV2D) {
  Convolution2d(conv_size, input, output, weights);
} else if (op_type == FULLY_CONNECTED) {
  FullyConnected(input, output, weights)

}
```

# NN Operator Support

TensorFlow

↓

1,400 operators

↓

Take up Flash Memory
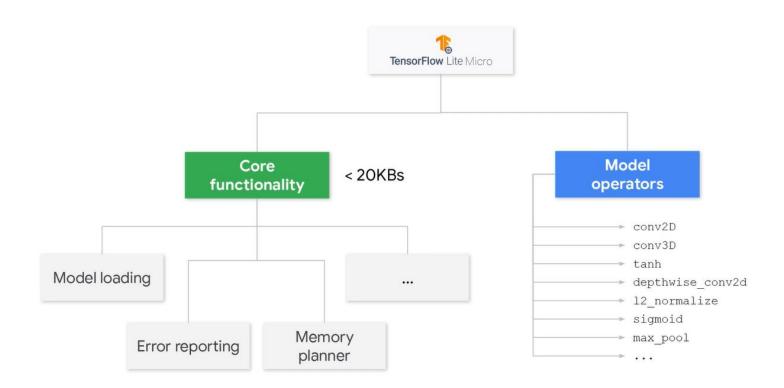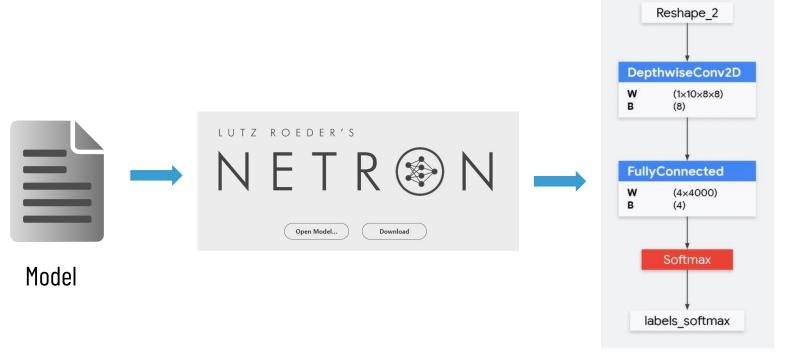
# NN Operator Support

# Choosing the Necessary Operators



Model



Visible Operators

# Choosing the Necessary Operators

Model

LUTZ ROEDER'S

NETRON

Open Model...    Download

Reshape_2

**DepthwiseConv2D**
W    (1×10×8×8)
B    (8)

**FullyConnected**
W    (4×4000)
B    (4)

Softmax

labels_softmax

Visible Operators

# OpResolver



```cpp
static tflite::MicroMutableOpResolver<4> micro_op_resolver(error_reporter);

if (micro_op_resolver.AddDepthwiseConv2D() != kTfLiteOk) {
  return;
}
if (micro_op_resolver.AddFullyConnected() != kTfLiteOk) {
  return;
}
if (micro_op_resolver.AddSoftmax() != kTfLiteOk) {
  return;
}
if (micro_op_resolver.AddReshape() != kTfLiteOk) {
  return;
}
```

# ▶ **Roadmap**

- High-Level Tensorflow Lite Micro Sales Pitch

- Review of Important Concepts

- Under the hood of Tensorflow Lite Micro

- Tensorflow Lite Micro Demo

# TFLM Workflow



specify ops you're using with **OpResolver** → load a model into an **Interpreter** → into model inputs, copy the **Input Data** → run the model via **Invoke()** → from model outputs, read the **Output Data** → **Action**

# Initialization

**Declare Variables**
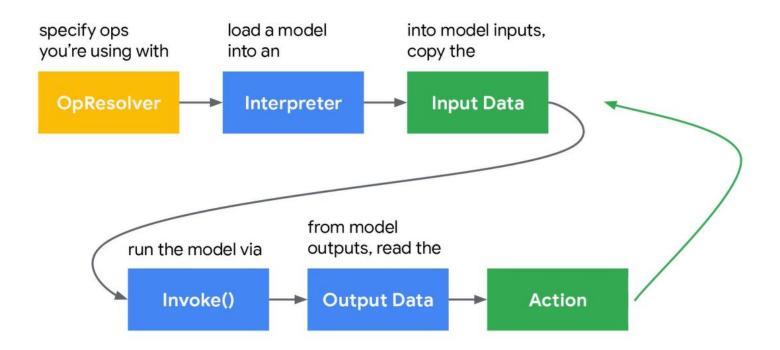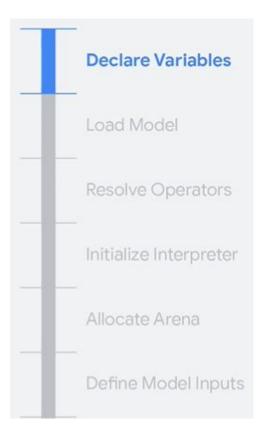
Load Model

Resolve Operators

Initialize Interpreter

Allocate Arena

Define Model Inputs

```cpp
// Globals, used for compatibility
// with Arduino-style sketches.

namespace {
    tflite::ErrorReporter* error_reporter = nullptr;
    const tflite::Model* model = nullptr;
    tflite::MicroInterpreter* interpreter = nullptr;
    TfLiteTensor* model_input = nullptr;
    FeatureProvider* feature_provider = nullptr;
    RecognizeCommands* recognizer = nullptr;
    int32_t previous_time = 0;

    // Create an area of memory to use for input,
    // output, and intermediate arrays.
    // The size of this will depend on the model
    // you're using, and may need to be
    // determined by experimentation.

    constexpr int kTensorArenaSize = 10 * 1024;
    uint8_t tensor_arena[kTensorArenaSize];
    int8_t feature_buffer[kFeatureElementCount];
    int8_t* model_input_buffer = nullptr;
}
```

# ▶ Initialization

**Declare Variables**

Load Model

Resolve Operators

Initialize Interpreter

Allocate Arena

Define Model Inputs

```cpp
// Globals, used for compatibility
// with Arduino-style sketches.

namespace {
    tflite::ErrorReporter* error_reporter = nullptr;
    const tflite::Model* model = nullptr;
    tflite::MicroInterpreter* interpreter = nullptr;
    TfLiteTensor* model_input = nullptr;
    FeatureProvider* feature_provider = nullptr;
    RecognizeCommands* recognizer = nullptr;
    int32_t previous_time = 0;

    // Create an area of memory to use for input,
    // output, and intermediate arrays.
    // The size of this will depend on the model
    // you're using, and may need to be
    // determined by experimentation.

    constexpr int kTensorArenaSize = 10 * 1024;
    uint8_t tensor_arena[kTensorArenaSize];
    int8_t feature_buffer[kFeatureElementCount];
    int8_t* model_input_buffer = nullptr;
}
```
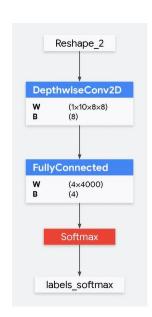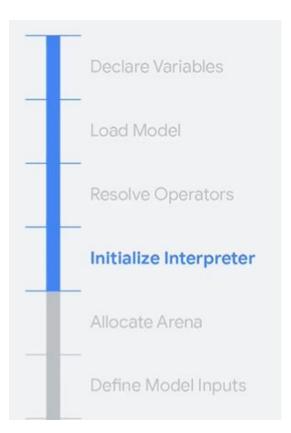
# ▶ Initialization

Declare Variables

**Load Model**

Resolve Operators

Initialize Interpreter

Allocate Arena

Define Model Inputs

```cpp
model =
tflite::GetModel(g_model);
```

```cpp
34
35  const unsigned char g_model[] DATA_ALIGN_ATTRIBUTE = {
36    0x20, 0x00, 0x00, 0x00, 0x54, 0x46, 0x4c, 0x33, 0x00, 0x00, 0x00, 0x00,
37    0x00, 0x00, 0x00, 0x12, 0x00, 0x1c, 0x00, 0x04, 0x00, 0x08, 0x00, 0x0c, 0x00,
38    0x10, 0x00, 0x14, 0x00, 0x00, 0x00, 0x00, 0x18, 0x00, 0x12, 0x00, 0x00, 0x00,
39    0x03, 0x00, 0x00, 0x00, 0x94, 0x48, 0x00, 0x00, 0x34, 0x42, 0x00, 0x00,
40    0x1c, 0x42, 0x00, 0x00, 0x3c, 0x00, 0x00, 0x00, 0x04, 0x00, 0x00, 0x00,
41    0x01, 0x00, 0x00, 0x00, 0x0c, 0x00, 0x00, 0x00, 0x08, 0x00, 0x0c, 0x00,
42    0x04, 0x00, 0x08, 0x00, 0x08, 0x00, 0x00, 0x00, 0x08, 0x00, 0x00, 0x00,
43    0x0b, 0x00, 0x00, 0x00, 0x13, 0x00, 0x00, 0x00, 0x6d, 0x69, 0x6e, 0x5f,
44    0x72, 0x75, 0x6e, 0x74, 0x69, 0x6d, 0x65, 0x5f, 0x76, 0x65, 0x72, 0x73,
45    0x69, 0x6f, 0x6e, 0x00, 0x0c, 0x00, 0x00, 0x00, 0x04, 0x41, 0x00, 0x00,
46    0xb4, 0x41, 0x00, 0x00, 0x24, 0x03, 0x00, 0x00, 0xf4, 0x02, 0x00, 0x00,
47    0xec, 0x02, 0x00, 0x00, 0xe4, 0x02, 0x00, 0x00, 0xc4, 0x02, 0x00, 0x00,
48    0xbc, 0x02, 0x00, 0x00, 0x2c, 0x00, 0x00, 0x00, 0x24, 0x00, 0x00, 0x00,
49    0x1c, 0x00, 0x00, 0x00, 0x04, 0x00, 0x00, 0x00, 0x16, 0xbd, 0xff, 0xff,
50    0x04, 0x00, 0x00, 0x00, 0x05, 0x00, 0x00, 0x00, 0x31, 0x2e, 0x35, 0x2e,
51    0x30, 0x00, 0x00, 0x00, 0x94, 0xba, 0xff, 0xff, 0x98, 0xba, 0xff, 0xff,
52    0x32, 0xbd, 0xff, 0xff, 0x04, 0x00, 0x00, 0x00, 0x80, 0x02, 0x00, 0x00,
53    0xfa, 0xee, 0x28, 0xc4, 0xee, 0xfe, 0xcf, 0x0f, 0x1e, 0xf7, 0x1f, 0x06,
54    0x0d, 0xed, 0xe9, 0x83, 0x5c, 0xc9, 0x18, 0xe3, 0xf9, 0x14, 0x28, 0x2a,
55    0x09, 0xf2, 0x18, 0x34, 0x62, 0xea, 0xef, 0xd6, 0x36, 0xb7, 0x1e, 0xf7,
56    0x3b, 0x22, 0x28, 0x39, 0xc2, 0x9d, 0xf1, 0x07, 0x5e, 0x0b, 0x1e, 0x2c,
57    0x07, 0xdd, 0xfd, 0xc3, 0xd8, 0x4a, 0xf3, 0x28, 0xa7, 0x16, 0xd5, 0xf1,
58    0xc3, 0x05, 0xfd, 0x27, 0xcc, 0xba, 0x1e, 0xcb, 0xd7, 0x3d, 0xd4, 0x29,
59    0x00, 0xfd, 0x28, 0x44, 0xfb, 0xf2, 0xf3, 0xb6, 0x4f, 0xcf, 0x09, 0xf0,
60    0xfa, 0x45, 0x41, 0x49, 0x05, 0xc5, 0x17, 0x5d, 0x64, 0x00, 0xf8, 0xee,
61    0x48, 0x17, 0xf4, 0xe9, 0x2e, 0x4b, 0x2e, 0x3f, 0xdf, 0xee, 0xe4, 0x08,
62    0x38, 0xf1, 0x16, 0x13, 0x2f, 0x2a, 0xed, 0xc2, 0xbf, 0x36, 0xf4, 0x02,
63    0xcf, 0xaa, 0xd2, 0xfa, 0xac, 0x13, 0xf6, 0xe8, 0xb5, 0x68, 0x12, 0xb6,
64    0xce, 0x0e, 0xdf, 0x58, 0xe4, 0x49, 0x14, 0x15, 0x03, 0xed, 0xfa, 0xd4,
65    0x40, 0xa7, 0xf6, 0xca, 0xfb, 0x00, 0x4d, 0x5e, 0xe4, 0x55, 0x1d, 0x30,
66    0x45, 0xe2, 0xfc, 0x01, 0x48, 0x81, 0xe9, 0xf1, 0x1e, 0xfc, 0x21, 0x32,
67    0xed, 0x4b, 0xed, 0xfa, 0x2f, 0xd2, 0xfa, 0xfb, 0x4d, 0xa7, 0xed, 0xc7,
68    0x92, 0xdf, 0xe6, 0xdb, 0xf8, 0x1f, 0xd9, 0xfa, 0x91, 0xf5, 0xe5, 0xc5,
69    0x8c, 0x17, 0x0f, 0xb9, 0xd2, 0xc7, 0xfe, 0x68, 0xd3, 0x51, 0x2e, 0x49,
70    0x1f, 0xbd, 0x01, 0xeb, 0x31, 0x17, 0xf0, 0xef, 0xff, 0xb8, 0x5d, 0x62,
71    0x02, 0x0f, 0x1f, 0x78, 0x6a, 0xb0, 0xf9, 0xfe, 0x4f, 0xcc, 0xd3, 0xff,
72    0x0a, 0x96, 0x1e, 0x2c, 0xed, 0xbc, 0xf4, 0x0b, 0x42, 0xc8, 0xf1, 0xea,
73    0x6e, 0x58, 0xec, 0xc4, 0x99, 0xae, 0xdc, 0xd7, 0x12, 0x87, 0xd8, 0x06,
74    0xa2, 0xc2, 0xe6, 0xa2, 0x81, 0x24, 0xe9, 0xac, 0xce, 0xb6, 0x15, 0x6b,
75    0xba, 0x00, 0x19, 0x58, 0x29, 0xb6, 0xfe, 0x01, 0x25, 0x96, 0xd2, 0xec,
```

# Initialization

Declare Variables

Load Model

**Resolve Operators**

Initialize Interpreter

Allocate Arena

Define Model Inputs

```
static tflite::MicroMutableOpResolver<4>
micro_op_resolver(error_reporter);
if (micro_op_resolver.AddBuiltin(
    tflite::BuiltinOperator_DEPTHWISE_CONV_2D,
    tflite::ops::micro::Register_DEPTHWISE_CONV_2D()) != kTfLiteOk)
{
    return;
}
if (micro_op_resolver.AddBuiltin(
    tflite::BuiltinOperator_FULLY_CONNECTED,
    tflite::ops::micro::Register_FULLY_CONNECTED()) != kTfLiteOk)
{
    return;
}
if (micro_op_resolver.AddBuiltin(
    tflite::BuiltinOperator_SOFTMAX,
    tflite::ops::micro::Register_SOFTMAX()) != kTfLiteOk)
{
    return;
}
if (micro_op_resolver.AddBuiltin(
    tflite::BuiltinOperator_RESHAPE,
    tflite::ops::micro::Register_RESHAPE()) != kTfLiteOk)
{
    return;
}
```

Reshape_2

**DepthwiseConv2D**
W    (1×10×8×8)
B    (8)

**FullyConnected**
W    (4×4000)
B    (4)

Softmax

labels_softmax

58

# Initialization

Declare Variables

Load Model

Resolve Operators

**Initialize Interpreter**

Allocate Arena

Define Model Inputs

```cpp
// Build an interpreter to run the model with.

static tflite::MicroInterpreter static_interpreter(
    model, micro_op_resolver, tensor_arena,
    kTensorArenaSize, error_reporter);


interpreter = &static_interpreter;
```
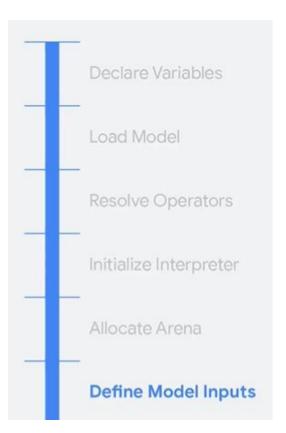
# ❯ Initialization

Declare Variables

Load Model

Resolve Operators

Initialize Interpreter

**Allocate Arena**

Define Model Inputs

```
// Allocate memory from the tensor_arena for
// the model's tensors.

TfLiteStatus allocate_status =
interpreter->AllocateTensors();
```

# Initialization

Declare Variables

Load Model

Resolve Operators

Initialize Interpreter

Allocate Arena

**Define Model Inputs**

```
model_input = interpreter->input(0);

model_input_buffer = model_input->data.int8;
```

# TFLM Executive Summary

- TFLM is built to fit with **embedded system constraints**

- Very **small** binary footprint

- **No** dynamic memory allocation

- **No** dependencies on standard C/C++ Libraries

- **No** operating system dependencies

# Roadmap

- High-Level Tensorflow Lite Micro Sales Pitch

- Review of Important Concepts

- Under the hood of Tensorflow Lite Micro

- Tensorflow Lite Micro Demo

# Extra Slides

# References

- David, Robert, et al. "Tensorflow lite micro: Embedded machine learning for tinyml systems." Proceedings of Machine Learning and Systems 3 (2021): 800-811.

- Chen, Tianqi, et al. "{TVM}: An automated {End-to-End} optimizing compiler for deep learning." 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18). 2018.

- UTensor (https://github.com/uTensor/uTensor)

- Microprocessor VS Microcontroller (https://www.linquip.com/blog/difference-between-microprocessor-and-microcontroller/)

- Embedded Systems (https://www.heavy.ai/technical-glossary/embedded-systems)

- TFLite Conversion (https://www.tensorflow.org/lite/models/convert)

- Edx TinyML Course (https://www.edx.org/professional-certificate/harvardx-tiny-machine-learning)

- Edx TinyML Course Material (https://github.com/tinyMLx/courseware/tree/master/edX#chapter-44-tensorflow-lite-micro)

# Suggested Readings

- Banbury, C., Reddi, V. J., Torelli, P., Holleman, J., Jeffries, N., Kiraly, C., ... & Xuesong, X. (2021). Mlperf tiny benchmark. arXiv preprint arXiv:2106.07597. 2106.07597.pdf (arxiv.org)

- David, R., Duke, J., Jain, A., Janapa Reddi, V., Jeffries, N., Li, J., ... & Rhodes, R. (2021). Tensorflow lite micro: Embedded machine learning for tinyml systems. Proceedings of Machine Learning and Systems, 3, 800-811. TensorFlow Lite Micro: Embedded Machine Learning on TinyML Systems (arxiv.org)

- Chen, Tianqi, et al. "{TVM}: An automated {End-to-End} optimizing compiler for deep learning." 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18). 2018. https://www.usenix.org/system/files/osdi18-chen.pdf

**Additional Links:**

- Memory Management in Tensorflow Lite Micro
  - Source: tflite-micro/memory_management.md at main · tensorflow/tflite-micro (github.com)
- MLPerf Tiny Deep Learning Benchmarks for Embedded Devices
  - Source: mlcommons/tiny: MLPerf™ Tiny is an ML benchmark suite for extremely low-power systems such as microcontrollers (github.com)