

Leveraging the power of C++ for efficient machine learning on embedded devices

Adrian Stanciu

adrian.stanciu.pub@gmail.com

NDC TechTown, 2023

About me

- ▶ I am a software engineer from Romania
- ▶ I have a bachelor's degree in computer science from Politehnica University of Bucharest
- ▶ I have a master's degree in computer and network security from Politehnica University of Bucharest
- ▶ I have 12 years of professional experience in Linux system programming
- ▶ My main programming languages are C and C++
- ▶ I am interested in algorithms, data structures and operating systems
- ▶ I work at Bitdefender, where I develop security software solutions that run on routers

Disclaimer

- ▶ This presentation is based on a personal project I worked on in my spare time and it is not directly related to the work I do at Bitdefender
- ▶ The outcome of this project is just a proof of concept
- ▶ All mistakes are mine

Agenda

- ▶ Motivation
- ▶ Image classification
- ▶ Hand gesture recognition
- ▶ Summary

Motivation

Machine Learning (ML)

- ▶ Subfield of Artificial Intelligence (AI)
- ▶ Enables computers to learn from data and then use that knowledge to make predictions
- ▶ Applications:
 - ▶ Computer vision
 - ▶ Medicine
 - ▶ Search engines

Embedded devices

- ▶ Computing devices designed to perform specific tasks within larger systems
- ▶ Applications:
 - ▶ Consumer electronics (mobile phones, smart TVs)
 - ▶ Home automation (thermostats, lighting control systems)
 - ▶ Medical equipment (pacemakers, insulin pumps)
- ▶ Characteristics:
 - ▶ Limited hardware resources
 - ▶ Low power consumption
 - ▶ May have real-time performance constraints

Machine learning on embedded devices

- ▶ Alternative to cloud-based machine learning
- ▶ Advantages:
 - ▶ Real-time processing
 - ▶ Low latency
 - ▶ Reduced bandwidth usage
 - ▶ Offline operation
 - ▶ Improved privacy
- ▶ Disadvantages:
 - ▶ Compatibility with various hardware and software platforms
 - ▶ Slower updates

Using C++ for machine learning on embedded devices

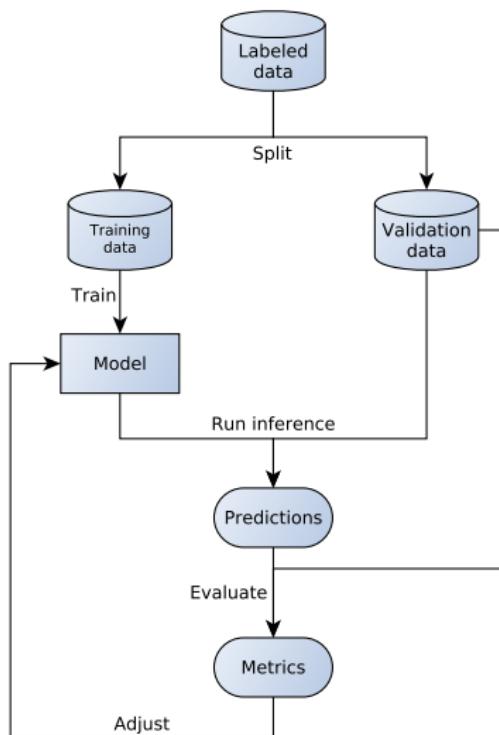
C++ is widely used in embedded systems:

- ▶ Was designed with efficiency in mind
- ▶ Provides both low-level access to hardware resources and high-level abstractions
- ▶ Is supported on most hardware and software platforms

Image classification

Supervised learning

Machine learning paradigm in which an algorithm learns from labeled training data to make predictions

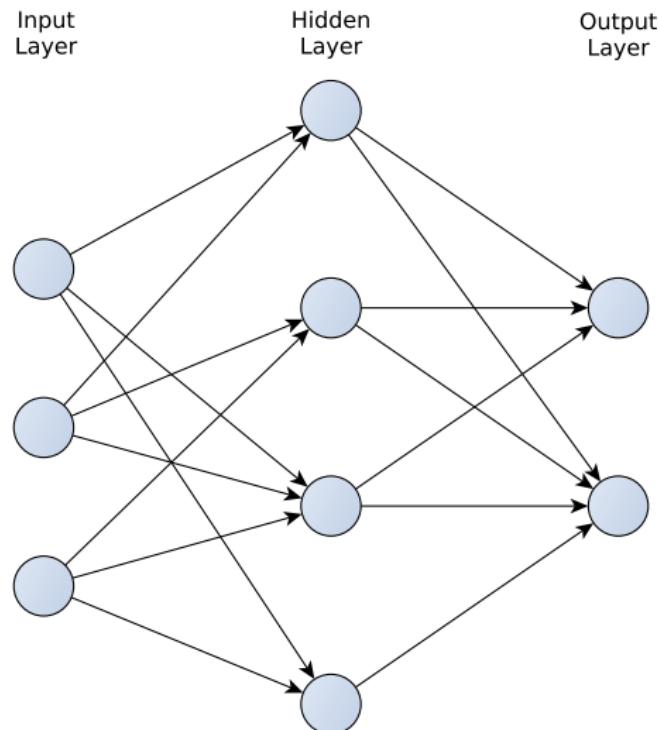


Supervised learning



Source: xkcd.com

Neural network (NN)

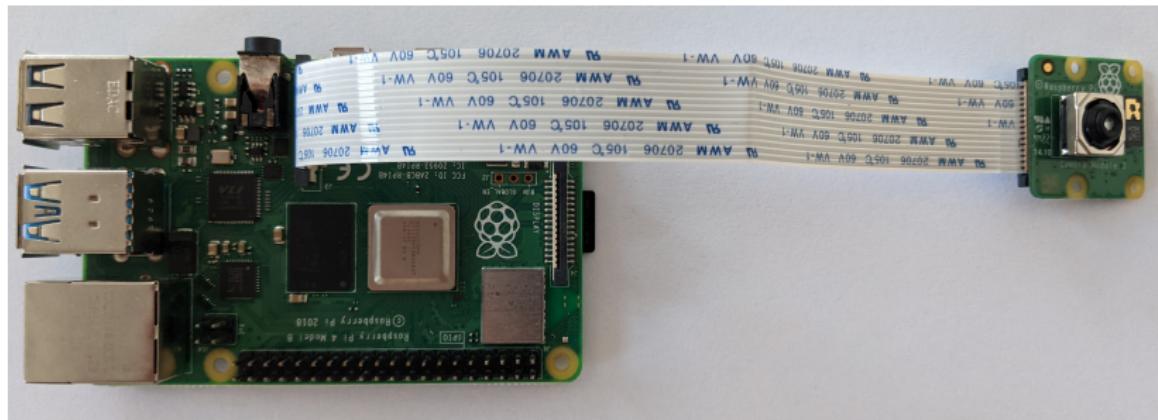


Convolutional neural network (CNN)

- ▶ Efficient in image classification
- ▶ A convolutional layer can apply filters to detect:
 - ▶ Edges
 - ▶ Shapes
 - ▶ Objects

Hardware setup

- ▶ Raspberry Pi 4 model B:
 - ▶ Quad core Cortex-A72 (ARM v8) 64-bit SoC @ 1.8GHz
 - ▶ 4GB RAM
- ▶ 8GB microSD card
- ▶ Raspberry Pi Camera module 3



Software dependencies

- ▶ TensorFlow Lite
- ▶ OpenCV

What is a tensor?

An array in which data (both input and output) is organized

MobileNet pre-trained model

- ▶ CNN architecture created by Google
- ▶ Accepts 224x224-pixel images, with 3 color channels per pixel (RGB)
- ▶ Trained on 1000 classes
- ▶ Labels are stored separate from the model:
 - ▶ `labels_mobilenet_quant_v1_224.txt` (11KB)
 - ▶ `mobilenet_v1_1.0_224_quant.tflite` (4.1MB)

Image classification algorithm

1. Load model and labels
2. Build interpreter
3. Allocate input and output tensors
4. Read image
5. Resize image
6. Copy resized image to input tensor
7. Run inference
8. Extract results from output tensor

Steps 4-8 can be repeated multiple times

1. Load model and labels

```
1 // const char *model_path;
2 // const char *labels_path;
3
4 std::unique_ptr<tflite::FlatBufferModel> model{
5     tflite::FlatBufferModel::BuildFromFile(model_path)
6 };
7
8 std::ifstream labels_ifs{labels_path};
9 std::string label;
10 std::vector<std::string> labels;
11 while (getline(labels_ifs, label)) {
12     labels.push_back(std::move(label));
13 }
```

2. Build interpreter

```
1 // std::unique_ptr<tflite::FlatBufferModel> model;
2 // int num_threads;
3
4 tflite::ops::builtin::BuiltinOpResolver resolver;
5 tflite::InterpreterBuilder builder{*model, resolver};
6
7 builder.SetNumThreads(num_threads);
8
9 std::unique_ptr<tflite::Interpreter> interpreter;
10 builder(&interpreter);
```

3. Allocate input and output tensors

```
1 // std::unique_ptr<tflite::Interpreter> interpreter;  
2  
3 interpreter->AllocateTensors();
```

4. Read image

```
1 // const char *image_path;  
2  
3 cv::Mat image{cv::imread(image_path)};
```

5. Resize image

```
1 // cv::Mat image;
2 // int required_image_width;
3 // int required_image_height;
4
5 cv::Mat resized_image;
6 cv::resize(
7     image,
8     resized_image,
9     cv::Size{required_image_width, required_image_height}
10 );
```

6. Copy resized image to input tensor

```
1 // std::unique_ptr<tflite::Interpreter> interpreter;
2 // cv::Mat resized_image;
3
4 std::memcpy(
5     interpreter->typed_input_tensor<uint8_t>(0),
6     resized_image.data,
7     resized_image.total() * resized_image.elemSize()
8 );
```

7. Run inference

```
1 // std::unique_ptr<tflite::Interpreter> interpreter;  
2  
3 interpreter->Invoke();
```

8. Extract results from output tensor

```
1 // std::unique_ptr<tflite::Interpreter> interpreter;
2 // std::vector<std::string> labels;
3
4 std::span<uint8_t> outputs{
5     interpreter->typed_output_tensor<uint8_t>(0),
6     labels.size()
7 };
8
9 std::vector<float> probabilities;
10 probabilities.reserve(labels.size());
11 for (auto output : outputs) {
12     probabilities.push_back(
13         1.0f * output / std::numeric_limits<uint8_t>::max()
14     );
15 }
```

Demo

```
$ libcamera-jpeg --width=640 --height=480 -o out.jpeg
```



Demo - Good results



0.96 | tennis ball

Demo - Bad results



0.30 | vase

0.25 | bell pepper
ambiguous results

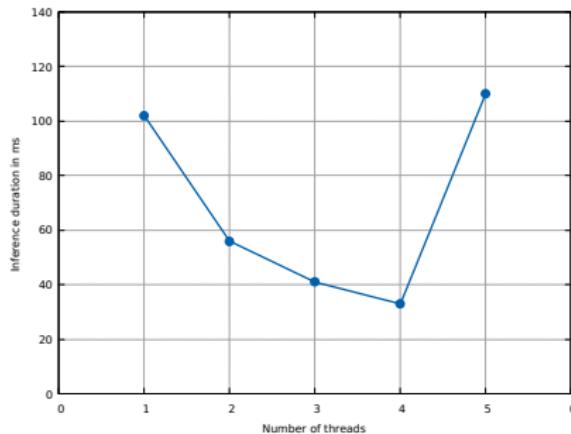
Demo - Bad results explanation

```
$ grep "tomato" labels_mobilenet_quant_v1_224.txt  
$
```

Performance

- ▶ Compilation duration: 35s
- ▶ Binary size: 67KB
- ▶ Running duration: 1s

Number of threads	Inference duration (ms)
1	102
2	56
3	41
4	33
5	110



- ▶ Memory consumption with 4 threads: 93MB

Comparision with Python

Comparision made with TensorFlow Lite's `label_image.py`

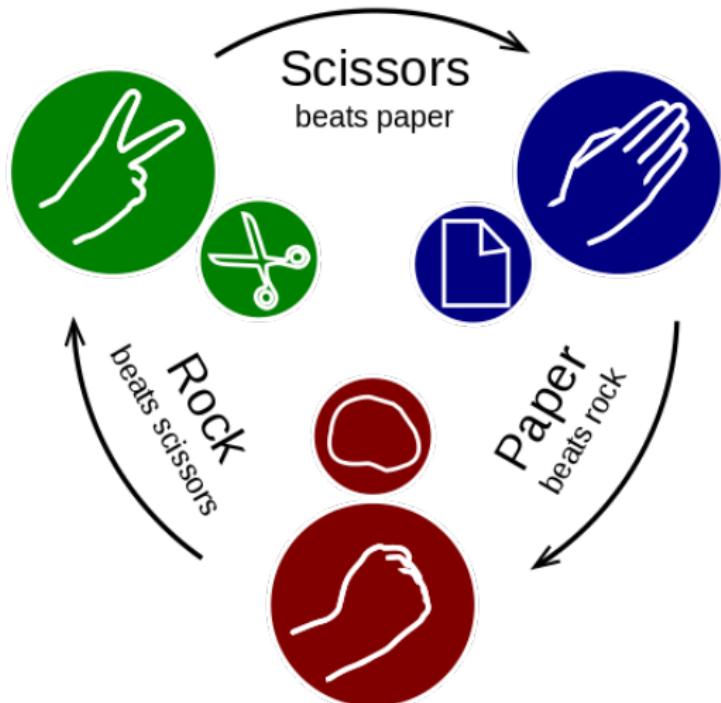
	C++	Python
Running duration (s)	1	7
Number of threads	Inference duration (ms)	
1	102	118
2	56	62
3	41	42
4	33	33
5	110	170

Number of threads	Inference duration (ms)	
	C++	Python
1	102	118
2	56	62
3	41	42
4	33	33
5	110	170

Memory consumption with 4 threads (MB)	C++	Python
93	320	

Hand gesture recognition

Rock-Paper-Scissors



Source: wikipedia.org

Data

- ▶ Google MediaPipe's Rock-Paper-Scissors dataset for hand gesture recognition
 - ▶ Contains 125 images for each class
 - ▶ Images have various sizes
 - ▶ Images have 3 color channels per pixel (RGB)
- ▶ Laurence Moroney's Rock-Paper-Scissors dataset (published by Sani Kamal on Kaggle)
 - ▶ Contains 964 images for each class split into training set (840) and testing set (124)
 - ▶ Images are 300x300 pixels
 - ▶ Images have 3 color channels per pixel (RGB)

Train a CNN-based Rock-Paper-Scissors model

- ▶ Transfer learning from a ResNet50 model
- ▶ Run in Python, on a laptop

Train on small dataset

- ▶ Random shuffle and split the images into:
 - ▶ 80% for training (300 images)
 - ▶ 20% for validation (75 images)

```
1 Epoch 1/10
2 9s - loss: 2.0545 - accuracy: 0.5700 - val_loss: 3.0788 - val_accuracy: 0.4800
3 Epoch 2/10
4 5s - loss: 1.4671 - accuracy: 0.7300 - val_loss: 1.5156 - val_accuracy: 0.7467
5 Epoch 3/10
6 5s - loss: 0.6304 - accuracy: 0.8567 - val_loss: 0.1994 - val_accuracy: 0.9467
7 Epoch 4/10
8 5s - loss: 0.3750 - accuracy: 0.9100 - val_loss: 0.3066 - val_accuracy: 0.9067
9 Epoch 5/10
10 5s - loss: 0.4212 - accuracy: 0.8933 - val_loss: 0.2647 - val_accuracy: 0.9333
11 Epoch 6/10
12 5s - loss: 0.3357 - accuracy: 0.8933 - val_loss: 0.1780 - val_accuracy: 0.9200
13 Epoch 7/10
14 5s - loss: 0.1435 - accuracy: 0.9533 - val_loss: 0.1747 - val_accuracy: 0.9733
15 Epoch 8/10
16 5s - loss: 0.2573 - accuracy: 0.9267 - val_loss: 0.1982 - val_accuracy: 0.9333
17 Epoch 9/10
18 5s - loss: 0.1846 - accuracy: 0.9367 - val_loss: 0.1352 - val_accuracy: 0.9867
19 Epoch 10/10
20 5s - loss: 0.1590 - accuracy: 0.9433 - val_loss: 0.1135 - val_accuracy: 0.9600
```

- ▶ Duration: 2m15s
- ▶ Model size: 23MB

Test on small dataset

Using the testing dataset of the big dataset (124 images per class)

Gesture	Correct predictions	Accuracy (%)
Rock	108	87.09
Paper	124	100
Scissors	8	6.45
All	240 out of 372	64.51

Train on big dataset

- ▶ Random shuffle and split the training set:
 - ▶ 80% for training (2016 images)
 - ▶ 20% for validation (504 images)

```
1 Epoch 1/10
2 37s - loss: 1.0362 - accuracy: 0.8284 - val_loss: 0.1478 - val_accuracy: 0.9385
3 Epoch 2/10
4 35s - loss: 0.1554 - accuracy: 0.9444 - val_loss: 0.0386 - val_accuracy: 0.9861
5 Epoch 3/10
6 35s - loss: 0.1092 - accuracy: 0.9608 - val_loss: 0.0256 - val_accuracy: 0.9901
7 Epoch 4/10
8 35s - loss: 0.0918 - accuracy: 0.9668 - val_loss: 0.0133 - val_accuracy: 0.9980
9 Epoch 5/10
10 35s - loss: 0.0589 - accuracy: 0.9797 - val_loss: 0.0325 - val_accuracy: 0.9861
11 Epoch 6/10
12 35s - loss: 0.0500 - accuracy: 0.9826 - val_loss: 0.0091 - val_accuracy: 0.9980
13 Epoch 7/10
14 35s - loss: 0.0361 - accuracy: 0.9896 - val_loss: 0.0070 - val_accuracy: 1.0000
15 Epoch 8/10
16 35s - loss: 0.0410 - accuracy: 0.9856 - val_loss: 0.0071 - val_accuracy: 1.0000
17 Epoch 9/10
18 35s - loss: 0.0310 - accuracy: 0.9886 - val_loss: 0.0074 - val_accuracy: 1.0000
19 Epoch 10/10
20 35s - loss: 0.0395 - accuracy: 0.9866 - val_loss: 0.0046 - val_accuracy: 1.0000
```

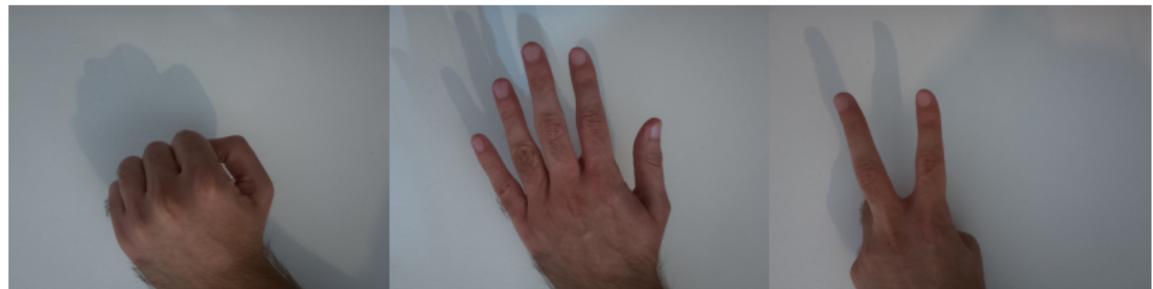
- ▶ Duration: 7m10s
- ▶ Model size: 23MB

Test on big dataset

Using its own testing dataset (124 images per class)

Gesture	Correct predictions	Accuracy (%)
Rock	124	100
Paper	122	98.38
Scissors	97	78.22
All	343 out of 372	92.20

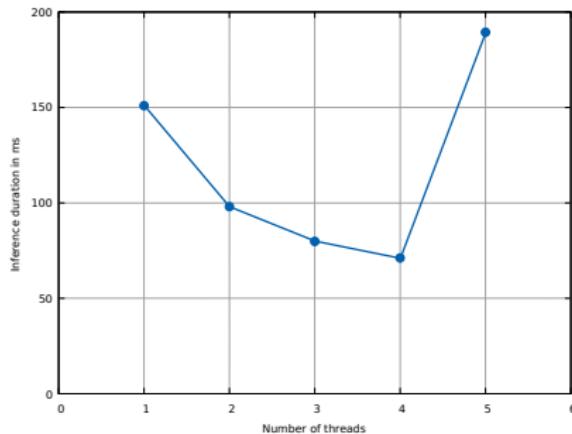
Demo



Performance

- ▶ Compilation duration: 35s
- ▶ Binary size: 67KB
- ▶ Running duration: 1s

Number of threads	Inference duration (ms)
1	151
2	98
3	80
4	71
5	189



- ▶ Memory consumption with 4 threads: 118MB

Live camera stream

- ▶ libcamera backend
- ▶ gstreamer middleware
- ▶ OpenCV frontend

Live camera stream

```
1 static constexpr const char *GstreamerPipeline{R"(  
2     libcamerasrc !  
3     video/x-raw,  
4     width=(int)640,  
5     height=(int)480,  
6     framerate=(fraction)10/1 !  
7     videoconvert !  
8     appsink  
9 )\"};  
10  
11 cv::VideoCapture camera{  
12     GstreamerPipeline,  
13     cv::CAP_GSTREAMER  
14 };  
15  
16 cv::Mat image;  
17 camera >> image;
```

Live inference - Algorithm

1. Create image classifier
2. Open camera
3. Do in a loop:
 - 3.1 Read image from camera
 - 3.2 Use classifier to run inference on image
 - 3.3 Show predictions

Live inference - Demo

Rock-Paper-Scissors game

On each round the AI:

1. Chooses Rock, Paper or Scissors uniformly at random
2. Infers the human's hand gesture
3. Computes the outcome

Rock-Paper-Scissors game - Demo

Summary

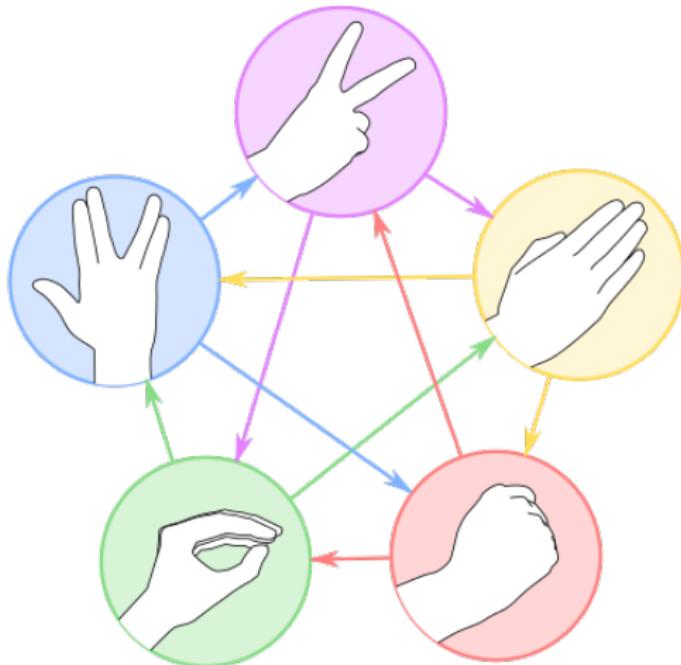
Recap

Task	Model	Inference
Image classification	pre-trained MobileNet	on-device in C++
Rock-Paper-Scissors hand gesture recognition	transfer learning from ResNet50	

Future work

- ▶ Reduce the size of the model and decrease inference duration by converting images to gray scale
- ▶ Use libcamera's API directly to reduce dependencies

Future work



Source: wikipedia.org

Conclusions

- ▶ Code isn't enough... data matters
- ▶ More diverse data leads to better models
- ▶ Building accurate models is an expert job
- ▶ Running on-device inference is here to stay

Repository

<https://github.com/adrian-stanciu/cpp-embedded-ml>

Resources

- ▶ <https://github.com/tensorflow/tensorflow/tree/master/tensorflow/lite/examples/python>
- ▶ <https://kaggle.com/datasets/sanikamal/rock-paper-scissors-dataset>
- ▶ <https://raspberrypi.com/products/camera-module-3>
- ▶ https://storage.googleapis.com/mediapipe-tasks/gesture_recognizer/rps_data_sample.zip
- ▶ <https://tensorflow.org/lite/examples>

Feedback is welcome

<https://forms.gle/BvSieuMRZbfKFpcq8>



Q&A

Thank you!