

TensorFlow Lite - based TinyML implementation in X-HEEP

Special Project per il corso di Architetture dei Sistemi di Elaborazione

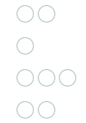
Eduard Antonovic Occhipinti

2024-10-04

Politecnico di Torino

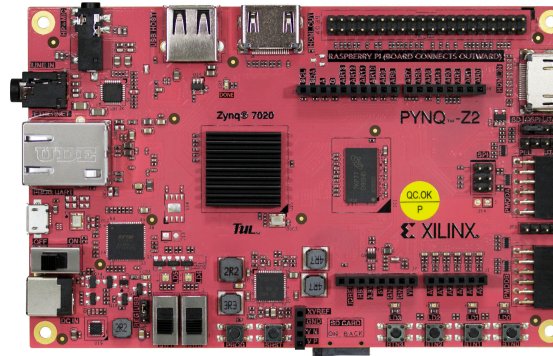
Outline

Introduzione	1	Quantizzazione	19
TensorFlow	3	Pruning	20
TensorFlow Model Optimization Toolkit	5	FastestDet	22
Quantization	6	Retraining	24
Weights Clustering	8	Manipolazione modello originale	26
Pruning	9	Conversione a <code>header file</code>	28
Collaborative Optimizations	10	Utilizzo	31
LeNet-5	11		
Modello sequenziale	13		
Modello funzionale	14		
Modello sottoclassato	15		
YOLOv8	16		



Introduzione

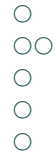
- Il progetto si presuppone di studiare le tecniche disponibili attualmente per la compressione di modelli pre-esistenti in maniera da renderli adatti all'esecuzione su dispositivi embedded.
 - In particolare l'esecuzione su una board RISC-V: la PYNQ



- Il progetto ha fornito inoltre un'opportunità per imparare a utilizzare una libreria di machine learning (TensorFlow) e mettere mano su tecniche più avanzate di manipolazione di modelli.

Outline

Introduzione	1	Quantizzazione	19
TensorFlow	3	Pruning	20
TensorFlow Model Optimization Toolkit	5	FastestDet	22
Quantization	6	Retraining	24
Weights Clustering	8	Manipolazione modello originale	26
Pruning	9	Conversione a <code>header file</code>	28
Collaborative Optimizations	10	Utilizzo	31
LeNet-5	11		
Modello sequenziale	13		
Modello funzionale	14		
Modello sottoclassato	15		
YOLOv8	16		



TensorFlow

TensorFlow è una libreria open-source per il machine learning sviluppata da Google.

- Fornisce un'interfaccia per la costruzione di modelli ed il training
- Supporta varie tecniche di compressione per ridurre la dimensione dei modelli

Keras

Una **API** di alto livello per la costruzione di modelli in TensorFlow.

TensorFlow Lite

Uno stack software per eseguire modelli TensorFlow su dispositivi embedded (in particolare col sottoinsieme di librerie denominate TensorFlow Lite for Microcontrollers).

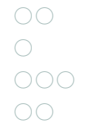


TensorFlow Model Optimization Toolkit

Il *TensorFlow Model Optimization Toolkit* è una suite di strumenti utili per la compressione di modelli TensorFlow.

È stato utilizzato in maniera estensiva per il progetto. Prima di arrivare al risultato finale ho deciso di esplorare tutte le tecniche offerte dal toolkit:

- **Quantization**
- **Pruning**
- **Weights Clustering**
- **Collaborative Optimizations**



Quantization

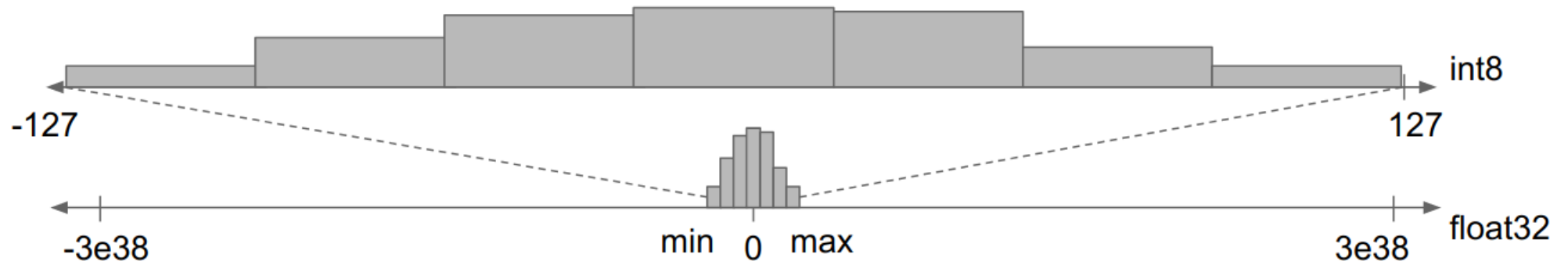


Figura 2: Conversione «lossy» da FP32 a INT8



Quantization

Quantization Aware Training

- Introduce l'errore di quantizzazione come rumore durante il training e come parte della loss, che l'algoritmo di ottimizzazione cerca di minimizzare
- Di conseguenza, il modello apprende parametri più robusti alla quantizzazione
- Sempre opportuno applicarlo ove possibile

Post-training Quantization

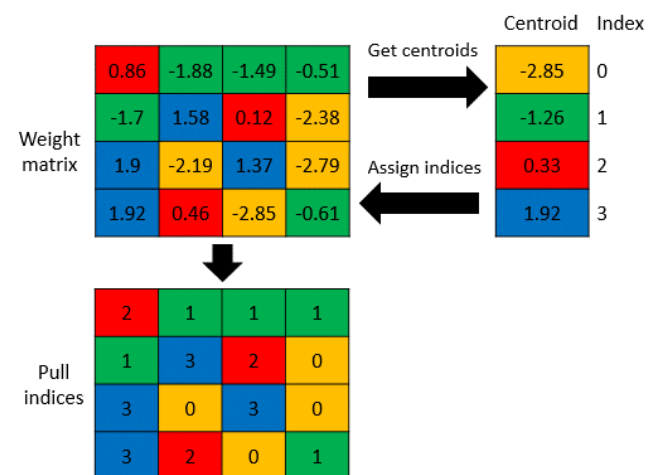
- Classica quantizzazione che comprime i parametri del modello finale per ridurre la dimensione
- Classico caso è la quantizzazione da FP32 a INT8 ma un ampio range di possibilità esiste, come FP16, INT4 o addirittura quantizzazione a 1 bit! ^[1]

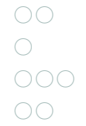
^[1]S. Ma *et al.*, «The Era of 1-bit LLMs: All Large Language Models are in 1.58 Bits». [Online]. Disponibile su: <https://arxiv.org/abs/2402.17764>



Weights Clustering

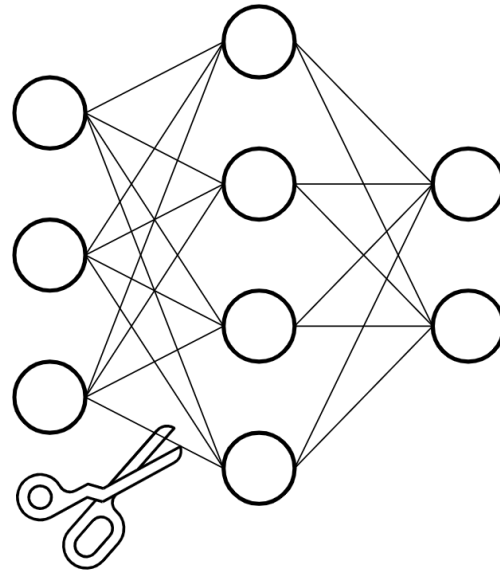
Consiste nel raggruppare i pesi in cluster ed approssimarli al centroide del cluster in maniera tale da diminuire il numero di valori distinti e ridurre la dimensione del modello.



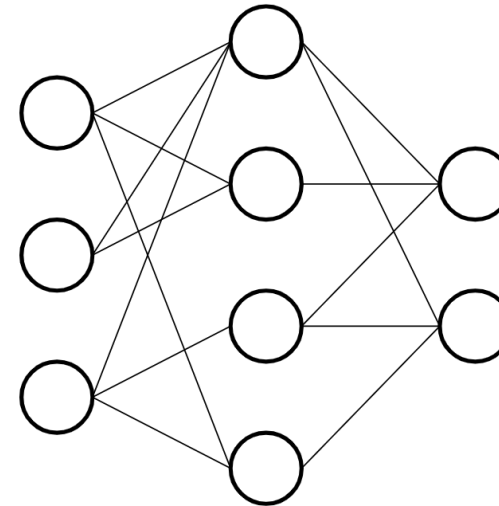


Pruning

Consiste nel rimuovere i pesi che sono stati identificati come non importanti.



Before pruning

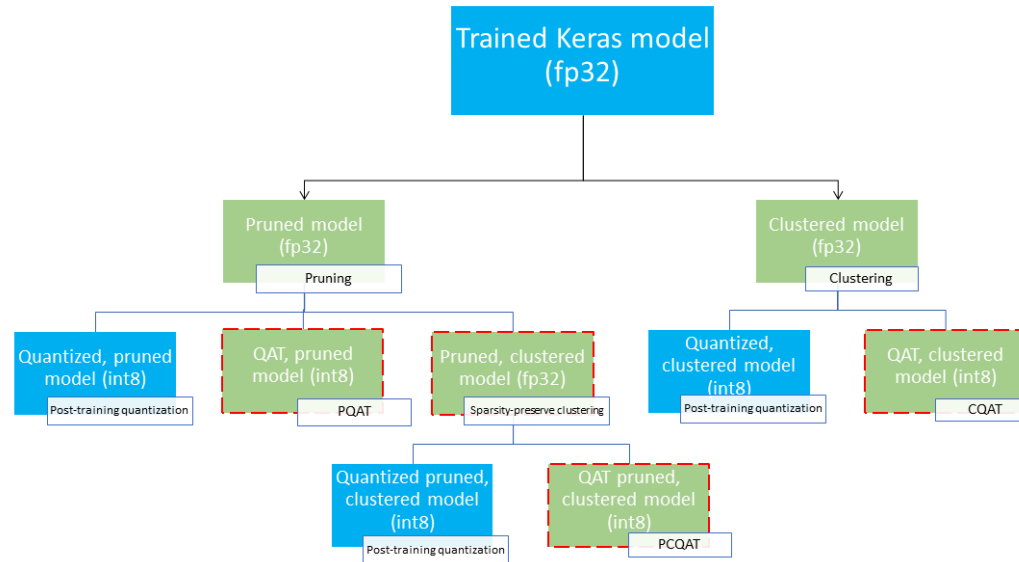


After pruning



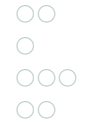
Collaborative Optimizations

Consiste nel combinare le varie tecniche di ottimizzazione in maniera tale che non interferiscano le une con le altre.



Outline

Introduzione	1	Quantizzazione	19
TensorFlow	3	Pruning	20
TensorFlow Model Optimization Toolkit	5	FastestDet	22
Quantization	6	Retraining	24
Weights Clustering	8	Manipolazione modello originale	26
Pruning	9	Conversione a <code>header file</code>	28
Collaborative Optimizations	10	Utilizzo	31
LeNet-5	11		
Modello sequenziale	13		
Modello funzionale	14		
Modello sottoclassato	15		
YOLOv8	16		

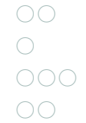


LeNet-5

Per iniziare a prendere mano con TensorFlow, è stato scelto di implementare un modello classico per architettura x86: **LeNet-5**

Per prendere la mano anche con i diversi tipi di modelli supportati da TensorFlow e rendersi meglio conto delle limitazioni e vantaggi di ciascuno di essi, questi è stato implementato in tre versioni:

- Un modello **sequenziale**
- Un modello **funzionale**
- Un modello **sottoclassato**



Modello sequenziale

È il modello più semplice da implementare, meno flessibile ma meglio supportato.

I modelli sequenziali permettono una più facile manipolazione dei layer e sono più adatti ad ottimizzazioni iterative.

- Modello scelto per studiare le varie ottimizzazioni disponibili

```
1 model = keras.models.Sequential(  
2     [  
3         keras.layers.Flatten(input_shape=(32, 32, 1)),  
4         keras.layers.Dense(128, activation="relu"),  
5         keras.layers.Dropout(0.2),  
6         keras.layers.Dense(10),  
7     ]  
8 )
```



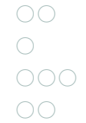


Modello funzionale

I modelli funzionali sono più flessibili e permettono di creare modelli più complessi, con possibilità di creare topology non lineari, multipli input e output, etc.

```
1 inputs = keras.Input(shape=(32, 32, 1))
2
3 x = layers.Conv2D(6, 5, activation="tanh")(inputs)
4 x = layers.AveragePooling2D(2)(x)
5 x = layers.Activation("sigmoid")(x)
6 x = layers.Conv2D(16, 5, activation="tanh")(x)
7 x = layers.AveragePooling2D(2)(x)
8 x = layers.Activation("sigmoid")(x)
9 x = layers.Conv2D(120, 5, activation="tanh")(x)
10 x = layers.Flatten()(x)
11 x = layers.Dense(84, activation="tanh")(x)
12
13 outputs = layers.Dense(10, activation="softmax")(x)
14 model = keras.Model(inputs=inputs, outputs=outputs)
```





Modello sottoclassato

Permette alta flessibilità e controllo, ma è più complesso da implementare e meno supportato.

```

1 class Lenet(keras.Model):
2     def __init__(self):
3         super().__init__()
4
5         self.conv1 = self._make_conv_layer(6, 5, input_shape=[32, 32, 1])
6         self.conv2 = self._make_conv_layer(16, 5)
7         self.conv3 = self._make_conv_layer(
8             120, 5, pooling=False, flatten=True
9         )
10        self.dense1 = layers.Dense(84, activation="tanh")
11        self.dense2 = layers.Dense(10, activation="softmax")
12
13    def call(self, inputs):
14        x = self.conv1(inputs)
15        x = self.conv2(x)
16        x = self.conv3(x)
17        x = self.dense1(x)
18        x = self.dense2(x)
19
20        return x

```

```

21 def _make_conv_layer(
22     self, filters, kernel_size, ish=None, pooling=True, flatten=False
23 ):
24     l = keras.Sequential()
25
26     if input_shape is not None:
27         l.add(
28             layers.Conv2D(
29                 filters, kernel_size, activation="tanh", ish=input_shape
30             )
31         )
32     else:
33         l.add(layers.Conv2D(filters, kernel_size, activation="tanh"))
34     if pooling:
35         l.add(layers.AveragePooling2D(2))
36         l.add(layers.Activation("sigmoid"))
37     if flatten:
38         l.add(layers.Flatten())
39
40     return l

```


Outline

Introduzione	1	Quantizzazione	19
TensorFlow	3	Pruning	20
TensorFlow Model Optimization Toolkit	5	FastestDet	22
Quantization	6	Retraining	24
Weights Clustering	8	Manipolazione modello originale	26
Pruning	9	Conversione a <code>header file</code>	28
Collaborative Optimizations	10	Utilizzo	31
LeNet-5	11		
Modello sequenziale	13		
Modello funzionale	14		
Modello sottoclassato	15		
YOLOv8	16		



YOLOv8

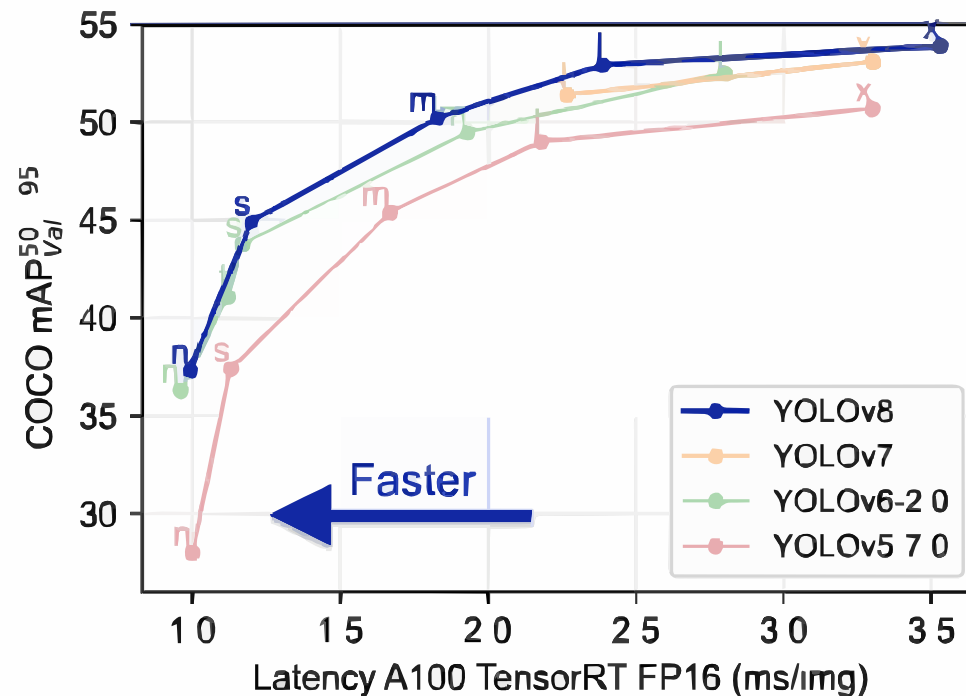
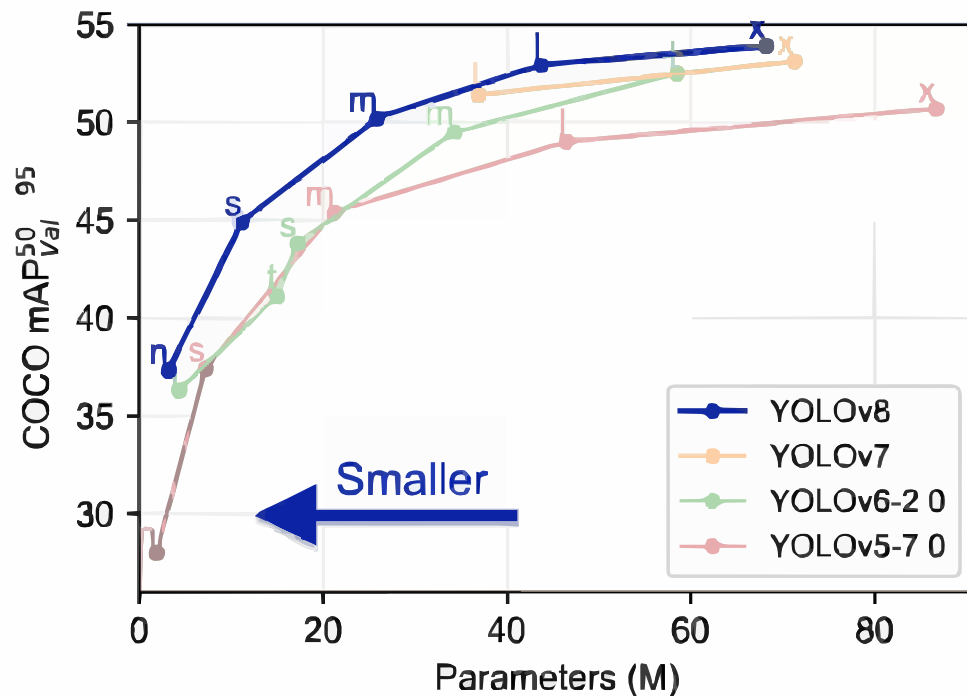
Un primo tentativo di applicare le tecniche studiate è stato fatto su un modello di object detection: YOLOv8.

- Al momento dei miei esperimenti questi corrispondeva allo stato dell'arte in termini di performance e dimensione ma nel frattempo è stato superato con due nuove versioni ^[1] ^[2] (la ricerca in questo campo è molto attiva).
- Mentre scrivevo questa slide un'altra versione, YOLO11 è stata rilasciata.

^[1]C.-Y. Wang, I.-H. Yeh, e H.-Y. M. Liao, «YOLOv9: Learning What You Want to Learn Using Programmable Gradient Information». [Online]. Disponibile su: <https://arxiv.org/abs/2402.13616>

^[2]A. Wang *et al.*, «YOLOv10: Real-Time End-to-End Object Detection». [Online]. Disponibile su: <https://arxiv.org/abs/2405.14458>

In particolare con il checkpoint `yolo8n`, il più piccolo tra quelli disponibili come possiamo vedere nel seguente grafico.





Quantizzazione

Una prima prova è stata fatta cercando di quantizzare a **INT8** il modello **torch** (post-training)

Risultati

- Riduzione della dimensione del modello dal checkpoint da 6.2 MB (o 12.2 MB se consideriamo il modello **.onnx** usato come step intermedio per la conversione) a 3.34 MB
- Perdita di performance **mAP50** e **mAP50-95** molto bassa

```
1 from ultralytics import YOLO
2
3 results_quant = YOLO(baseline_quantized).val(data="coco128.yaml")
```



```
1 {
2 |   'metrics/precision(B)': 0.6401136562982888,
3 |   'metrics/recall(B)': 0.5371329744029286,
4 |   'metrics/mAP50(B)': 0.6049606058248067,
5 |   'metrics/mAP50-95(B)': 0.4455822678794395,
6 |   'fitness': 0.4615201016739762
7 }
```

ORIGINAL_FP32

```
1 {
2 |   'metrics/precision(B)': 0.6693744628674021,
3 |   'metrics/recall(B)': 0.5357106358126033,
4 |   'metrics/mAP50(B)': 0.6089456673483312,
5 |   'metrics/mAP50-95(B)': 0.4540750785376153,
6 |   'fitness': 0.4695621374186869
7 }
```

QUANTIZED_INT8



Pruning

Un tentativo di ridurre ulteriormente la dimensione finale del modello è stata fatta utilizzando il pruning di torch.

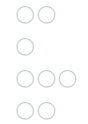
- Non essendo tutti i layer compatibili ho isolato solo i layer di tipo `Conv2d`.

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.utils.prune as prune
4 import YOLO
5
6 model = YOLO("yolov8n.pt")
7 PRUNING_AMOUNT = 0.1
8
9 for name, m in model.named_modules():
10     if isinstance(m, nn.Conv2d):
11         prune.ll_unstructured(m, name="weight", amount=PRUNING_AMOUNT)
12         prune.remove(m, "weight") # make permanent
```

- Purtroppo i risultati rendono il modello molto meno efficace dell'originale.

```
1 {
2 |   'metrics/precision(B)': 0.48268779643490495,
3 |   'metrics/recall(B)': 0.42881404068150936,
4 |   'metrics/mAP50(B)': 0.462184408704583,
5 |   'metrics/mAP50-95(B)': 0.3170316451700596,
6 |   'fitness': 0.33154692152351195
7 }
```

PRUNED

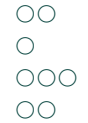


Altri appunti

- Il metodo consigliato da [ultralytics](https://github.com/ultralytics/ultralytics) per il pruning del modello al momento dei miei test non funzionava correttamente (<https://github.com/ultralytics/ultralytics/issues/3507>)
- Il metodo descritto da <https://github.com/VainF/Torch-Pruning> non era compatibile con tutti i layer del modello

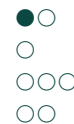
Outline

Introduzione	1	Quantizzazione	19
TensorFlow	3	Pruning	20
TensorFlow Model Optimization Toolkit	5	FastestDet	22
Quantization	6	Retraining	24
Weights Clustering	8	Manipolazione modello originale	26
Pruning	9	Conversione a <code>header file</code>	28
Collaborative Optimizations	10	Utilizzo	31
LeNet-5	11		
Modello sequenziale	13		
Modello funzionale	14		
Modello sottoclassato	15		
YOLOv8	16		



FastestDet

Dato gli scarsi risultati nel cercare di comprimere un modello già estremamente efficiente quale YOLOv8 senza comprometterne in maniera significativa le performance, ho esplorato altre opzioni finendo per trovare un modello con appena 250K parametri, FastestDet



Retraining

Una delle strategie adottate per ridurre la dimensione è stata quella di provare a re-trainare il modello con immagini più piccole e ad andare di conseguenza anche a modificare i vari layer a livello di codice.

DATASET:

```
- TRAIN: "/home/qiuqiu/Desktop/coco2017/train2017.txt"
- VAL: "/home/qiuqiu/Desktop/coco2017/val2017.txt"
```

```
+ TRAIN: "- "
```

```
+ VAL: "- "
```

```
  NAMES: "configs/coco.names"
```

MODEL:

```
  NC: 80
```

```
- INPUT_WIDTH: 352
```

```
- INPUT_HEIGHT: 352
```

```
+ INPUT_WIDTH: 128
```

```
+ INPUT_HEIGHT: 128
```

```
...
```

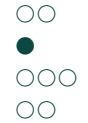
Diff



Retraining

Problemi nell'adattare il codice

- I `requirements.txt` non erano scritti in maniera corretta e causavano conflitti tra le dipendenze, questo ha comportato il dover adattare il codice per far sì che funzionasse con le nuove versioni di `torch`, linguaggio in cui il modello è implementato.
- Purtroppo anche dopo questo il training non è andato a buon fine in quanto il dataset originale (che come abbiamo visto in precedenza era hardcodato) non era in formato `coco` standard.



Manipolazione modello originale

Ho quindi deciso di riutilizzare il checkpoint in formato `.onnx` fornito dall'autore.

```
1 import tempfile
2 import onnx2tf
3
4 saved_model_fastest_det = tempfile.mkdtemp()
5
6 onnx2tf.convert(
7     ONNX_PATH,
8     output_integer_quantized_tflite=True,
9     output_h5=True,
10    output_folder_path=saved_model_fastest_det,
11 )
```



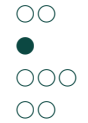
- La decisione di usare `.h5` come formato rende la conversione più robusta
- È possibile effettuare quantizzazione già in questo step di conversione

```
1 import numpy as np
2
3 def representative_dataset():
4     for _ in range(100):
5         data = np.random.rand(1, 352, 352, 3)
6         yield [data.astype(np.float32)]
```



```
1 import tensorflow as tf
2
3 converter = tf.lite.TFLiteConverter.from_saved_model(
4     saved_model_fastest_det
5 )
6
7 converter.optimizations = [tf.lite.Optimize.DEFAULT]
8 converter.representative_dataset = representative_dataset
9 converter.target_spec.supported_ops = [
10     tf.lite.OpsSet.TFLITE_BUILTINS_INT8
11 ]
12 converter.inference_input_type = tf.uint8
13 converter.inference_output_type = tf.uint8
```





Risultati

Quantizzando abbiamo ridotto la dimensione del modello da 1.2 MB a 0.42 MB, riuscendo con successo a scendere sotto il nostro target di 0.5 MB.

- Il modello riesce ad operare anche su immagini scalate, senza bisogno di re-training

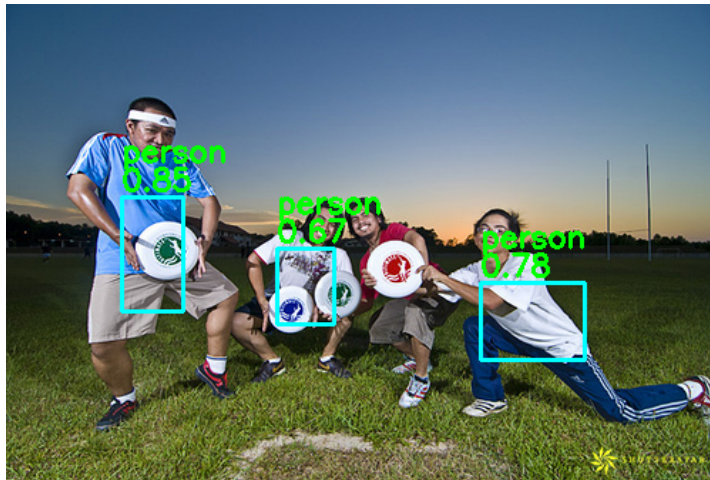
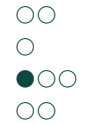


Figura 7: Inference su immagine originale



Figura 8: Inference su immagine 128x86

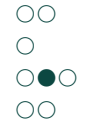


Conversione a header file

Andiamo adesso a convertire il modello in un **header file** in maniera tale da poterlo usare direttamente in un progetto **C/C++** per effettuare inference. Per fare ciò ho scritto due piccole funzioni, un' alternativa poteva essere quella di usare `xxd -i [file]` ma fare la conversione in maniera programmatica permette di avere più controllo sul risultato finale.

```
1 def c_style_hexdump(input, output, name):
2     with open(input, "rb") as f:
3         file = f.read()
4
5     file = bytearray(file)
6     _bytes = [f"0x{x:02x}" for x in file]
7     file = ",".join(_bytes)
8
9     with open(output, "w") as f:
10         f.write("#pragma once\n")
11         f.write("#include <stdalign.h>\n")
12         f.write(f"alignas(16) const unsigned char {name}[] = {{{file}}};\n")
13
14     return len(_bytes)
```

```
1 def build_header(output, names_with_sizes):
2     with open(output, "w") as f:
3         f.write('#pragma once\n#ifdef __cplusplus\nextern "C"\n{\n#endif\n')
4         f.write("#include <stdalign.h>\n\n")
5
6         for name, size in names_with_sizes:
7             f.write(
8                 f"alignas(16) extern const unsigned char {name}[{size}];\n"
9             )
10
11         f.write("\n#endif __cplusplus\n}\n#endif\n")
```



Conversione a header file

```
1 MODEL = "models/fastest_det_rom.c"
2 INPUT = "models/fastest_det_input.c"
3 TEST_IMAGE = "fastest_det/data/3.jpg"
4 HEADER = "models/fastest_det.h"
```



```
1 model_size = c_style_hexdump(
2     fastest_det_quant_file, MODEL, "model_data"
3 )
```

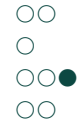


Codice 1: Generiamo il dump del modello, noteremo che la dimensione del file C generato è di ben 2.1 MB, ciò è causato da byte a zero di padding, che non influiscono sulla dimensione del modello caricato in memoria

```
1 from PIL import Image
2
3 img = Image.open(TEST_IMAGE)
4 img = img.convert("RGB")
5
6 img_array = np.array(img)
7 img_array = img_array.astype(np.uint8)
8
9 img.thumbnail((128, 128))
10 img_array_compressed = np.array(img)
11 img_array_compressed = img_array_compressed.astype(np.uint8)
12
13 _, imagebin = tempfile.mkstemp(".bin")
14
15 with open(imagebin, "wb") as f:
16     f.write(bytearray(img_array_compressed)) # type: ignore
17
18 input_size = c_style_hexdump(imagebin, INPUT, "input_data")
```



Codice 2: Comprimiamo l'immagine e la salviamo come file C



Conversione a header file

Generiamo infine l'header file che verrà incluso nel codice sorgente con `build_header(HEADER, [("input_data", input_size), ("model_data", model_size)])`.

```
1 #pragma once
2 #include <stdalign.h>
3 alignas(16) const unsigned char input_data[] = {
4     0x19, 0x47, 0x6b, 0x1c, 0x4a, 0x6e, 0x1d,
5     ...
6 };
```

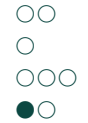
Codice 3: `models/fastest_det_input.c`

```
1 #pragma once
2 #include <stdalign.h>
3 alignas(16) const unsigned char model_data[] = {
4     0x1c, 0x00, 0x00, 0x00, 0x54, 0x46, 0x4c,
5     ...
6 };
```

Codice 4: `models/fastest_det_rom.c`

```
1 #pragma once
2
3 #ifdef __cplusplus
4     extern "C"
5     {
6 #endif
7
8     #include <stdalign.h>
9
10    alignas(16) extern const unsigned char input_data[33024];
11    alignas(16) extern const unsigned char model_data[420000];
12
13    #ifdef __cplusplus
14    }
15 #endif
```

Codice 5: `models/fastest_det.h`



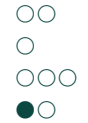
Utilizzo

Per utilizzare il modello, dobbiamo scrivere il codice **C++** relativo (che utilizzerà le librerie di **TensorFlow Lite (Micro)** in maniera appropriata).

```
1 #pragma message "hello_world_test.cc"
2 extern "C"
3 {
4 #include "fastest_det_test.h"
5 #include <math.h>
6 #include <stdio.h>
7 }
8
9 #include "models/fastest_det.h"
10 #include "tensorflow/lite/core/c/common.h"
11 #include "tensorflow/lite/micro/micro_interpreter.h"
12 #include "tensorflow/lite/micro/micro_log.h"
13 #include "tensorflow/lite/micro/micro_mutable_op_resolver.h"
14 #include "tensorflow/lite/micro/micro_profiler.h"
15 #include "tensorflow/lite/micro/recording_micro_interpreter.h"
16 #include "tensorflow/lite/micro/system_setup.h"
17 #include "tensorflow/lite/schema/schema_generated.h"
18
```



```
19 namespace
20 {
21 constexpr int kTensorArenaSize = 0x70000;
22 uint8_t tensor_arena[kTensorArenaSize];
23 const tflite::Model *model = nullptr;
24
25 using FastestDetOpResolver = tflite::MicroMutableOpResolver<11>;
26
27 TfLiteStatus RegisterOps(FastestDetOpResolver &op_resolver)
28 {
29     TF_LITE_ENSURE_STATUS(op_resolver.AddQuantize());
30     TF_LITE_ENSURE_STATUS(op_resolver.AddPad());
31     TF_LITE_ENSURE_STATUS(op_resolver.AddDepthwiseConv2D());
32     TF_LITE_ENSURE_STATUS(op_resolver.AddConv2D());
33     TF_LITE_ENSURE_STATUS(op_resolver.AddConcatenation());
34     TF_LITE_ENSURE_STATUS(op_resolver.AddTranspose());
35     TF_LITE_ENSURE_STATUS(op_resolver.AddReshape());
36     TF_LITE_ENSURE_STATUS(op_resolver.AddGather());
37 }
```

```

37     TF_LITE_ENSURE_STATUS(op_resolver.AddAveragePool2D());
38     TF_LITE_ENSURE_STATUS(op_resolver.AddLogistic());
39     TF_LITE_ENSURE_STATUS(op_resolver.AddSoftmax());
40     return kTfLiteOk;
41 }
42 } // namespace
43
44 TfLiteStatus load_model()
45 {
46     if (model != nullptr) { return kTfLiteOk; }
47     model = ::tflite::GetModel(model_data);
48     TFLITE_CHECK_EQ(model->version(), TFLITE_SCHEMA_VERSION);
49     return kTfLiteOk;
50 }
51
52 TfLiteStatus Infer(const char *data, size_t len, int8_t **out, size_t
*out_len)
53 {
54     if (model == nullptr) { return kTfLiteError; }
55     FastestDetOpResolver op_resolver;
56     TF_LITE_ENSURE_STATUS(RegisterOps(op_resolver));
57
58     tflite::MicroInterpreter interpreter(
59         model, op_resolver, tensor_arena, kTensorArenaSize);
60     TF_LITE_ENSURE_STATUS(interpreter.AllocateTensors());

```

```

61     TfLiteTensor *input = interpreter.input(0);
62     TFLITE_CHECK_NE(input, nullptr);
63     TfLiteTensor *output = interpreter.output(0);
64     TFLITE_CHECK_NE(output, nullptr);
65     memcpy(input->data.int8, data, len);
66     TF_LITE_ENSURE_STATUS(interpreter.Invoke());
67     *out = output->data.int8;
68     *out_len = output->bytes;
69
70     return kTfLiteOk;
71 }
72
73 extern "C" int init_tflite()
74 {
75     tflite::InitializeTarget();
76     TF_LITE_ENSURE_STATUS(load_model());
77     return kTfLiteOk;
78 }
79
80 extern "C" int
81 infer(const char *data, size_t len, int8_t **out, size_t *out_len)
82 {
83     return Infer(data, len, out, out_len);
84 }

```

Grazie dell'attenzione