Lab04 (DLonEdge) Screenshots

1. Introduction
   Edge analytics with real-time processing capabilities is challenging but important and inevitable due to privacy/security concerns. However, edge devices like RaspberryPi are constrained with limited hardware resources, which at times are not sufficient to run complex deep learning models. These models require lot of computational resource and memory due to their size and complex architecture. Therefore, in such scenarios, we optimize the model such that it can run efficiently with reduced inference time critical for real-time analytics. Optimization can be achieved by combination of techniques like quantization and converting trained model into architecture specific lite model.

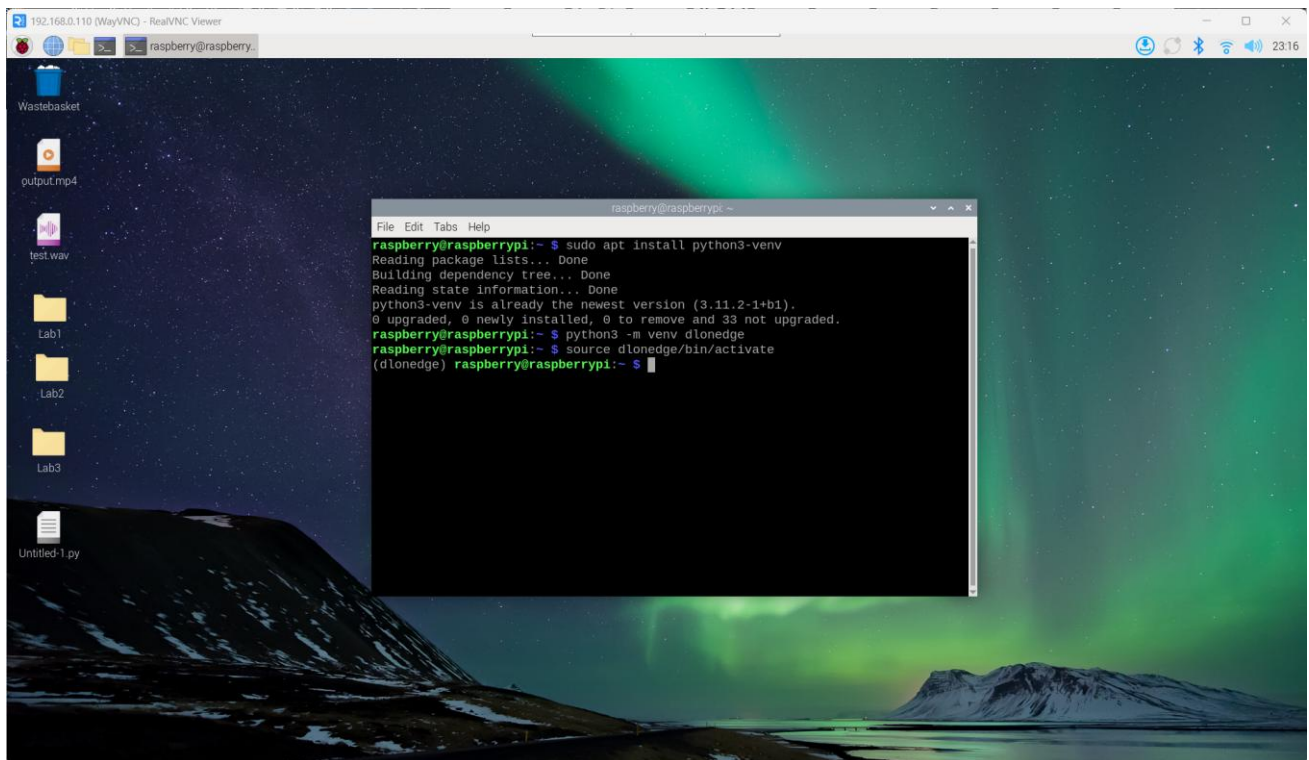2. Running Deep Learning Model On RaspberryPi



Fig. 1. Screenshot of activating a virtual environment named "dlonedge" to avoid conflicts in libraries.
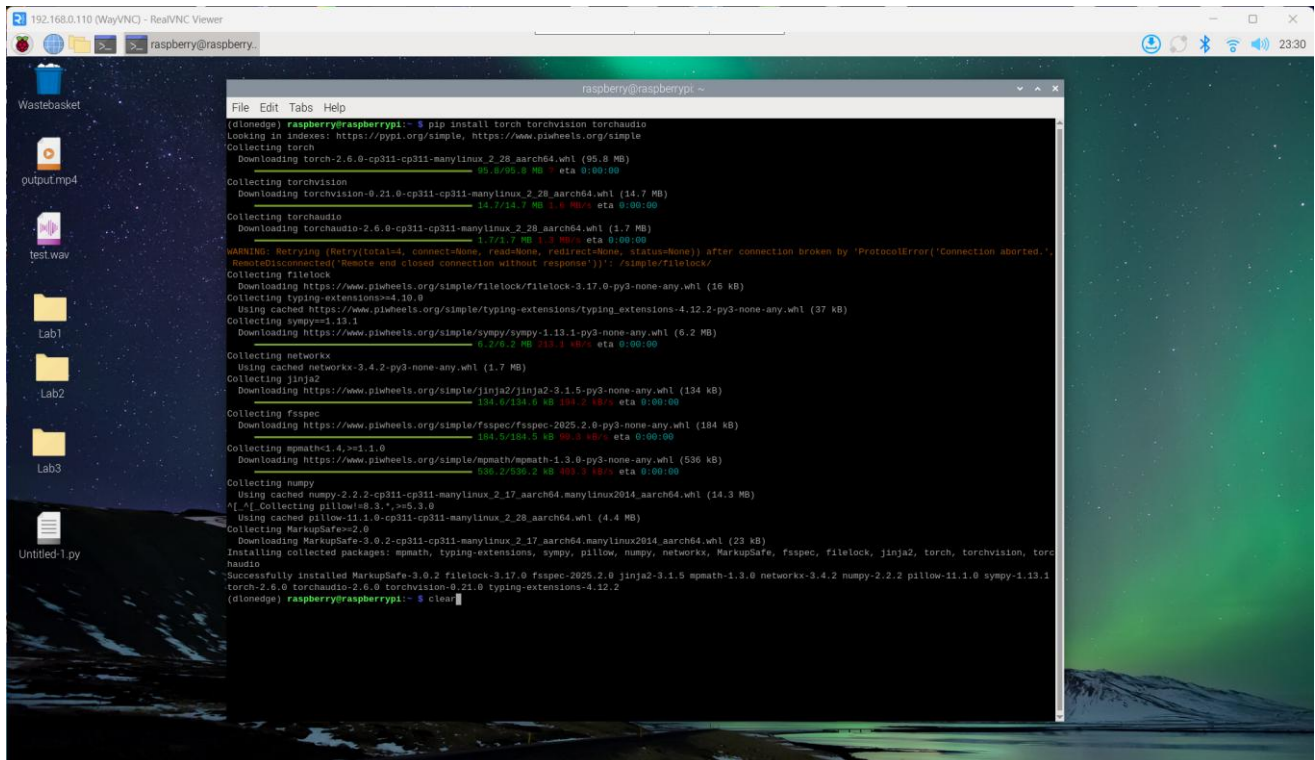
# Lab04 (DLonEdge) Screenshots



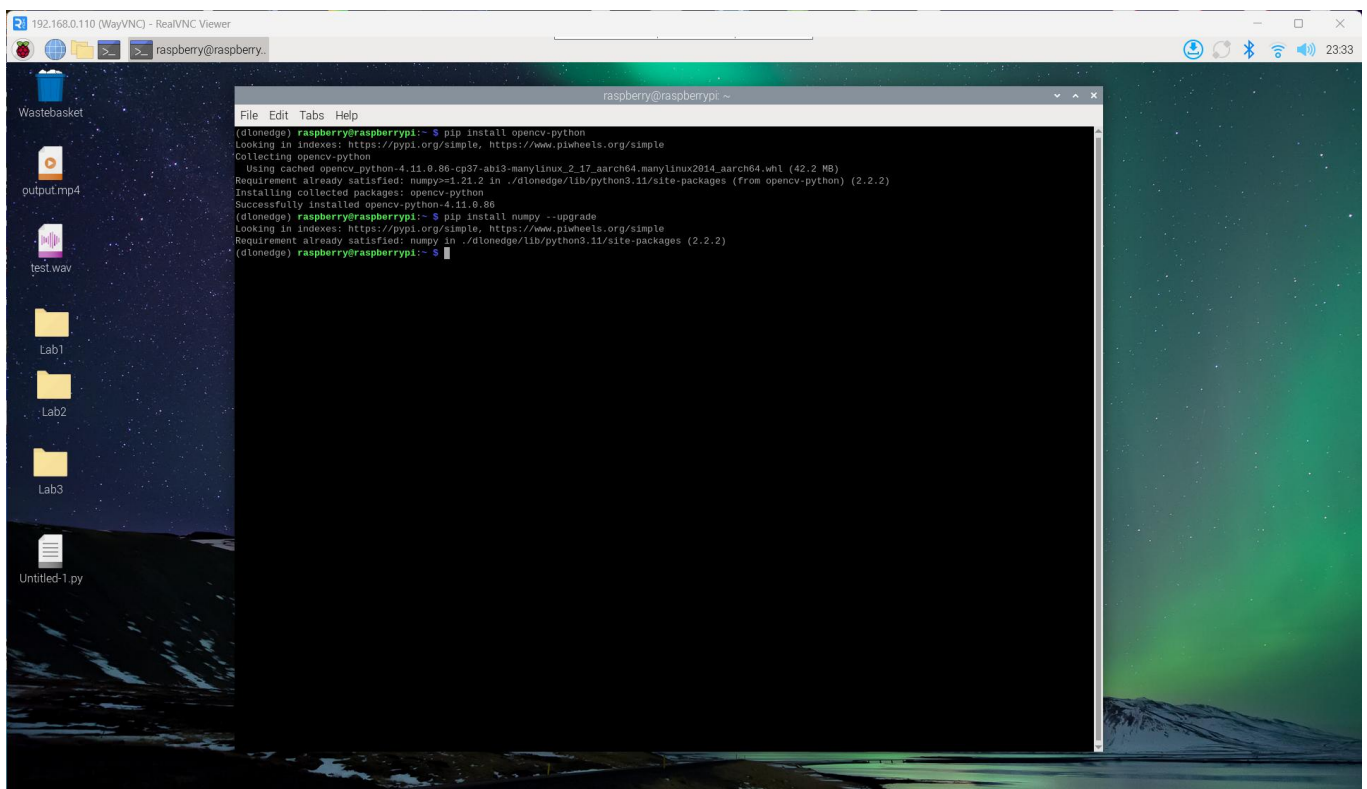Fig. 2. Screenshot of installing PyTorch (`torch, torchvision, torchaudio`)



Fig. 2. Screenshot of installing OpenCV (`opencv-python`) and upgrade `numpy`
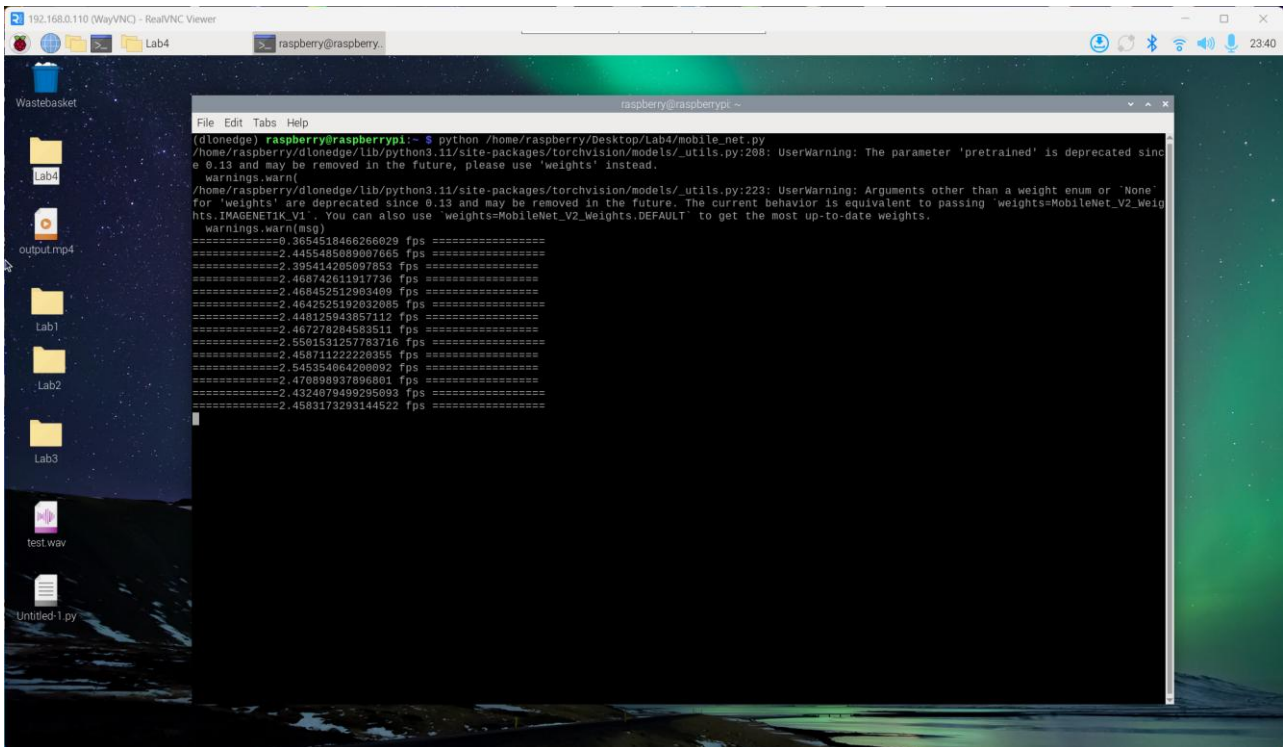
Lab04 (DLonEdge) Screenshots



Fig. 3. Screenshot of doing model inference on local pre-trained MobileNetV2 model, with no optimization of model and could only achieve 2-3 FPS on Raspberry Pi 3B+
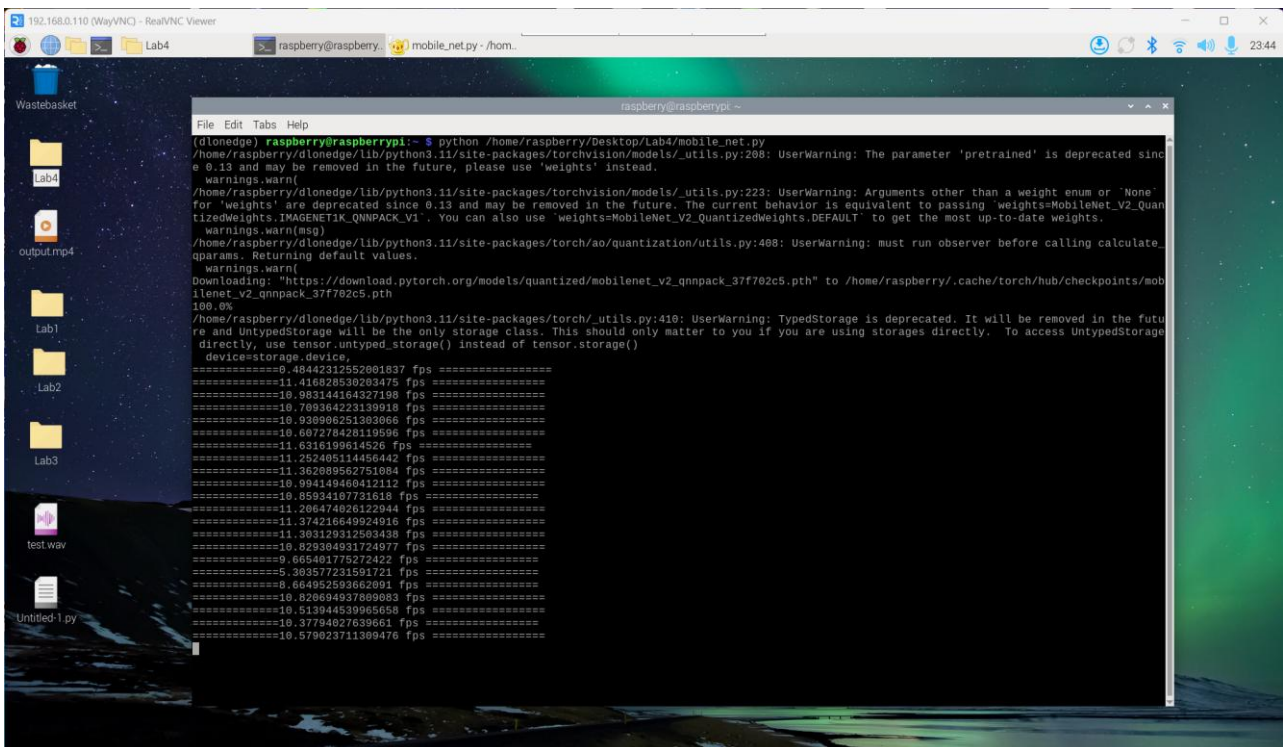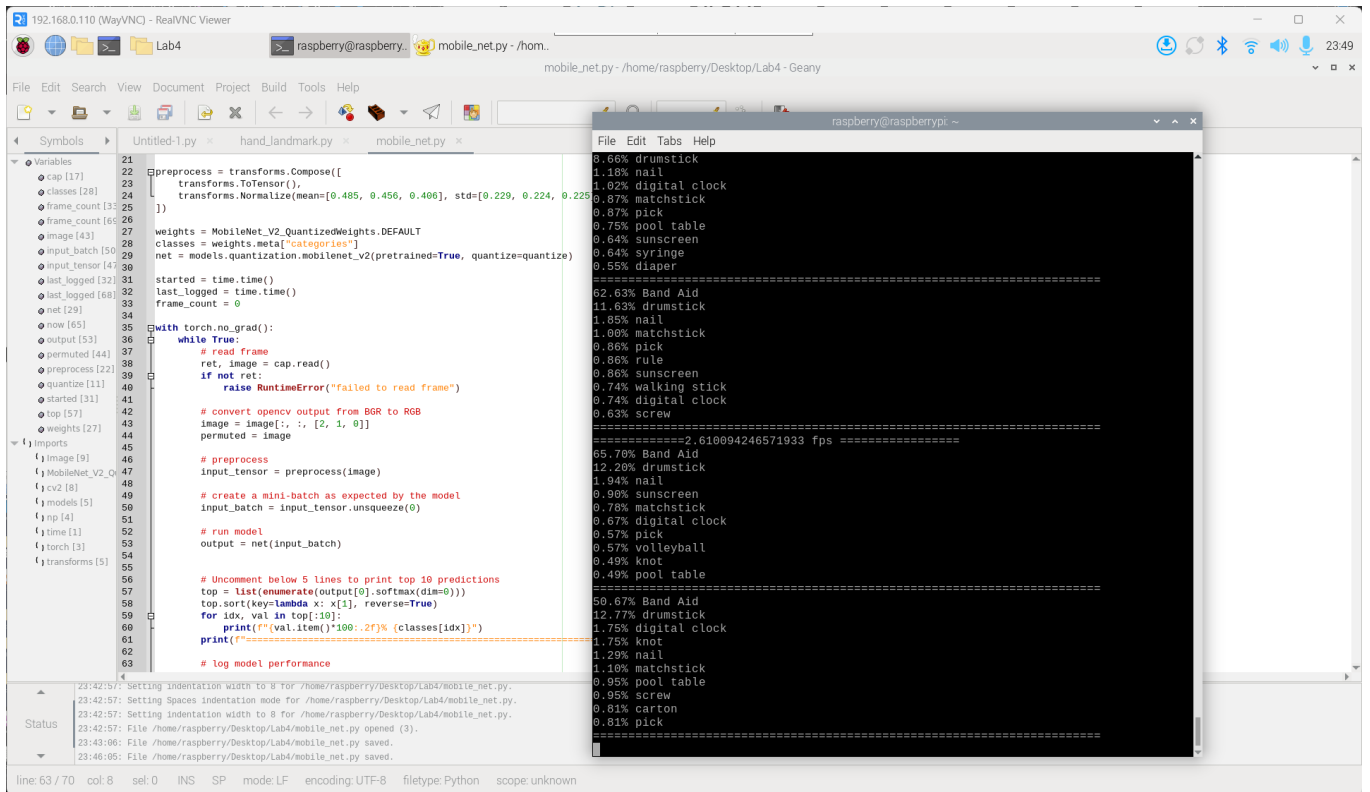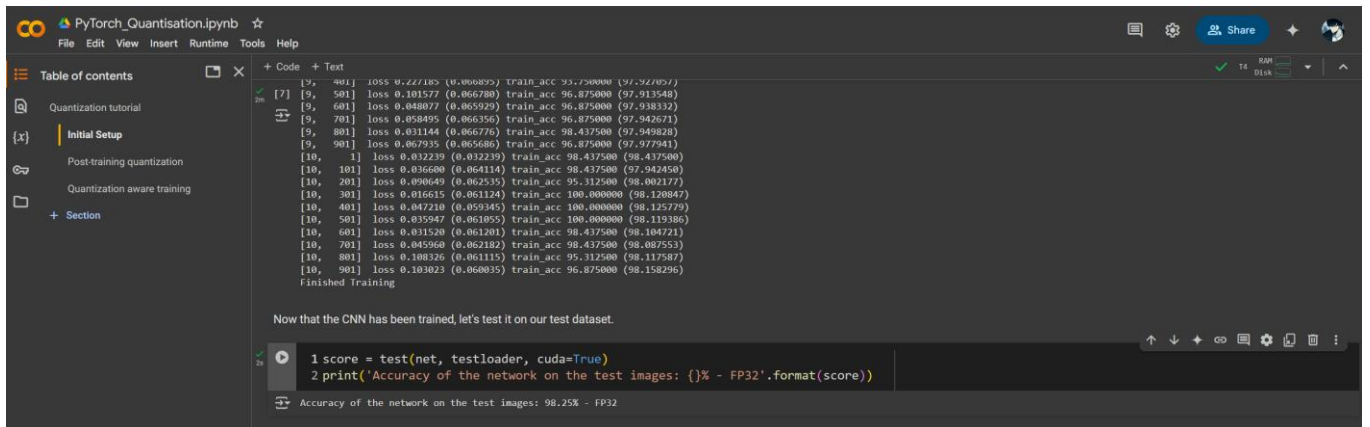


Fig. 4. Screenshot of doing model inference on pre-trained local MobileNetV2 model, with optimization of model (with quantization enabled) and could achieve 5-12 FPS on Raspberry Pi 3B+
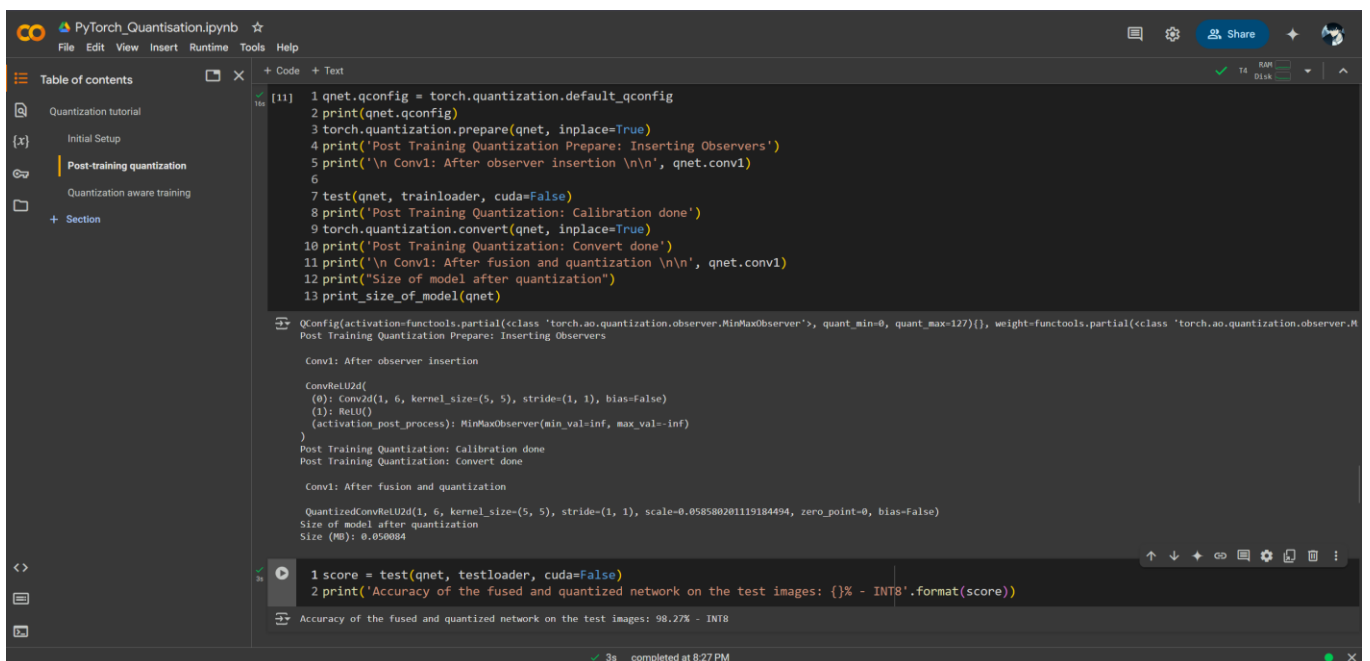
# Lab04 (DLonEdge) Screenshots



Fig. 5. Screenshot of printing the top 10 predictions made by the local pre-trained MobileNetV2 model according to streaming data sent through the webcam

3. Quantization using Pytorch



Fig. 6. Training of the CNN on the training dataset

Fig. 7. Testing the CNN on the test dataset, returning accuracy score of 98.25



Fig. 8. Post-training quantization, with accuracy score of 98.27%

# Lab04 (DLonEdge) Screenshots



```
 8                                          activation=MovingAverageMinMaxObserver.with_args(reduce_range=True),
 9                                          weight=MovingAverageMinMaxObserver.with_args(dtype=torch.qint8, qscheme=torch.per_tensor_symmetric))
10 print(qnet.qconfig)
11 torch.quantization.prepare(qnet, inplace=True)
12 print('Post Training Quantization Prepare: Inserting Observers')
13 print('\n Conv1: After observer insertion \n\n', qnet.conv1)
14
15 test(qnet, trainloader, cuda=False)
16 print('Post Training Quantization: Calibration done')
17 torch.quantization.convert(qnet, inplace=True)
18 print('Post Training Quantization: Convert done')
19 print('\n Conv1: After fusion and quantization \n\n', qnet.conv1)
20 print("Size of model after quantization")
21 print_size_of_model(qnet)
22 score = test(qnet, testloader, cuda=False)
23 print('Accuracy of the fused and quantized network on the test images: {}% - INT8'.format(score))
```

```
QConfig(activation=functools.partial(<class 'torch.ao.quantization.observer.MovingAverageMinMaxObserver'>, reduce_range=True){}, weight=functools.partial(<class 'torch.ao.quantization.observ
Post Training Quantization Prepare: Inserting Observers

 Conv1: After observer insertion

ConvReLU2d(
  (0): Conv2d(1, 6, kernel_size=(5, 5), stride=(1, 1), bias=False)
  (1): ReLU()
  (activation_post_process): MovingAverageMinMaxObserver(min_val=inf, max_val=-inf)
)
/usr/local/lib/python3.11/dist-packages/torch/ao/quantization/observer.py:229: UserWarning: Please use quant_min and quant_max to specify the range for observers.          reduce_
  warnings.warn(
Post Training Quantization: Calibration done
Post Training Quantization: Convert done

 Conv1: After fusion and quantization

 QuantizedConvReLU2d(1, 6, kernel_size=(5, 5), stride=(1, 1), scale=0.05842381715774536, zero_point=0, bias=False)
Size of model after quantization
Size (MB): 0.050084
Accuracy of the fused and quantized network on the test images: 98.25% - INT8
```



In addition, we can significantly improve on the accuracy simply by using a different quantization configuration. We repeat the same exercise with the recommended configuration for quantizing for arm64 architecture (qnnpack). This configuration does the following: Quantizes weights on a per-channel basis. It uses a histogram observer that collects a histogram of activations and then picks quantization parameters in an optimal manner.

```
1 qnet = Net(q=True)
2 load_model(qnet, net)
3 fuse_modules(qnet)
```

```
1 qnet.qconfig = torch.quantization.get_default_qconfig('qnnpack')
2 print(qnet.qconfig)
3
4 torch.quantization.prepare(qnet, inplace=True)
5 test(qnet, trainloader, cuda=False)
6 torch.quantization.convert(qnet, inplace=True)
7 print("Size of model after quantization")
8 print_size_of_model(qnet)
```
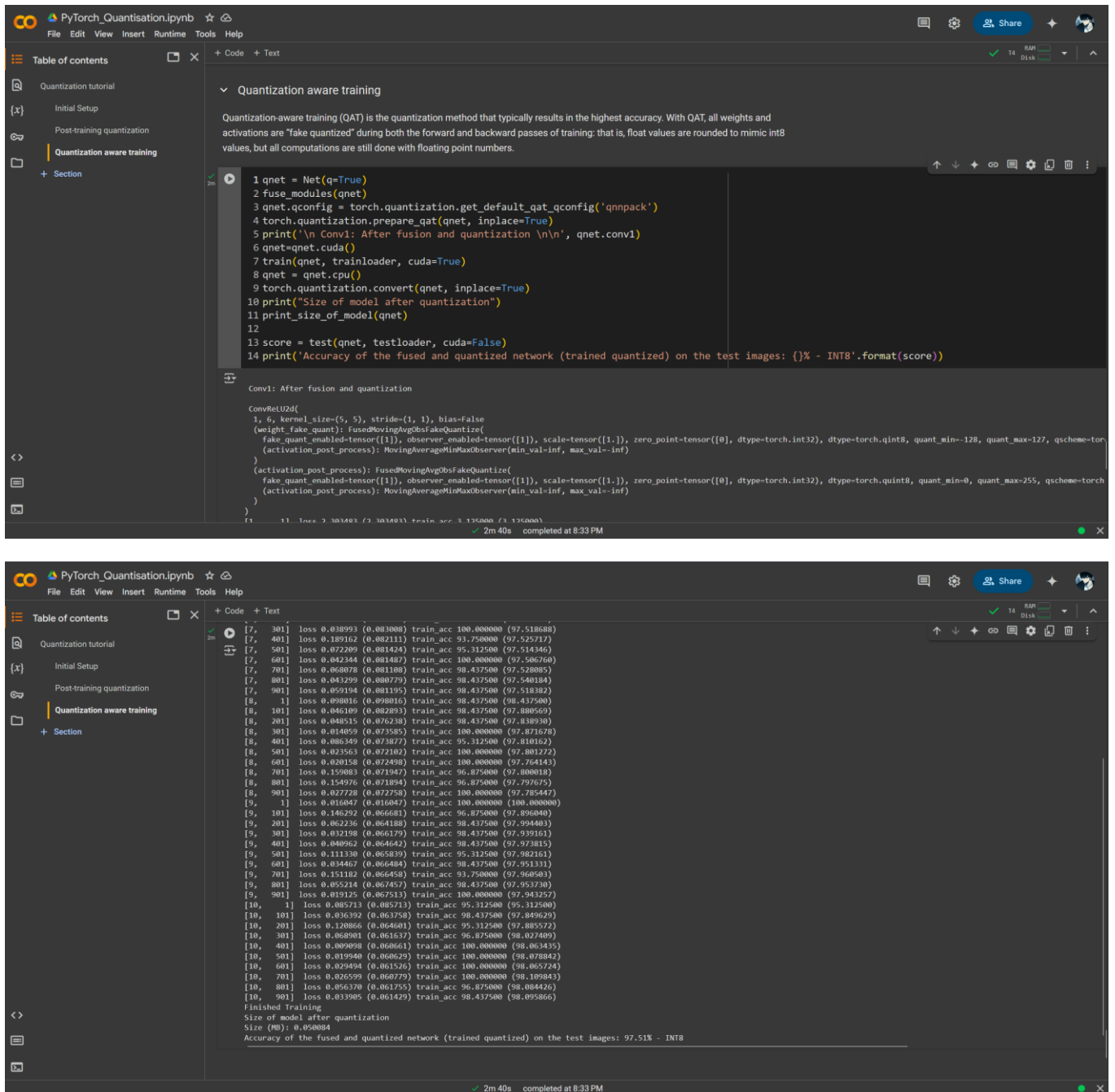
```
QConfig(activation=functools.partial(<class 'torch.ao.quantization.observer.HistogramObserver'>, reduce_range=False){}, weight=functools.partial(<class 'torch.ao.quantization.observer.MinMax
Size of model after quantization
Size (MB): 0.050084
```

```
1 score = test(qnet, testloader, cuda=False)
2 print('Accuracy of the fused and quantized network on the test images: {}% - INT8'.format(score))
```

```
Accuracy of the fused and quantized network on the test images: 98.14% - INT8
```

Fig 9. Result of a custom quantization configuration with accuracy of 98.25%

# Lab04 (DLonEdge) Screenshots





Fig 10. Quantization Aware Training (QAT), with accuracy result of 97.51%