# eIQ® Neutron NPU for MCUs Lab
# Part 1: Mobilenet

# Contents

# 1    Lab Overview

This document will cover how to convert models using NXP's eIQ Toolkit and will also highlight the performance improvements that can be achieved with the eIQ Neutron NPU on MCX devices.

# 2    Software and Hardware Installation

This section will cover the hardware and software needed for this lab.

## 2.1 Hardware

The <u>FRDM-MCXN947</u> development board is used in this lab

## 2.2 NXP Software Installation

1.  Install <u>MCUXpresso IDE v11.9.1</u> or later
2.  Install <u>eIQ Toolkit 1.11.4</u> - During installation it is not required to install the optional components offered by the install wizard. Everything needed for this lab is already installed by default.
3.  A quantized Mobilenet model can be found attached to this post or in the i.MX RT1170 MCUXpresso SDK v2.15.0 at **\boards\evkmimxrt1170\eiq_examples\tflm_label_image\doc**
4.  Download <u>MCUXpresso SDK 2.14.0 for FRDM-MCXN947</u>
    It includes the eIQ software platform and demos. It must be 2.14.0 to ensure compatibility with eIQ Toolkit v1.11.4
    a)  On the SDK builder page, make sure to select the "**eIQ**" middleware and that you're downloading version **2.14.0**



    b)  Then click on the **Download SDK** button and accept the license agreement to download the zip file.

# 3    Label Image Example

This section will use the eIQ Label Image example found in the MCX N MCUXpresso SDK to showcase how the eIQ Neutron NPU can significantly decrease inference times for quantized models.
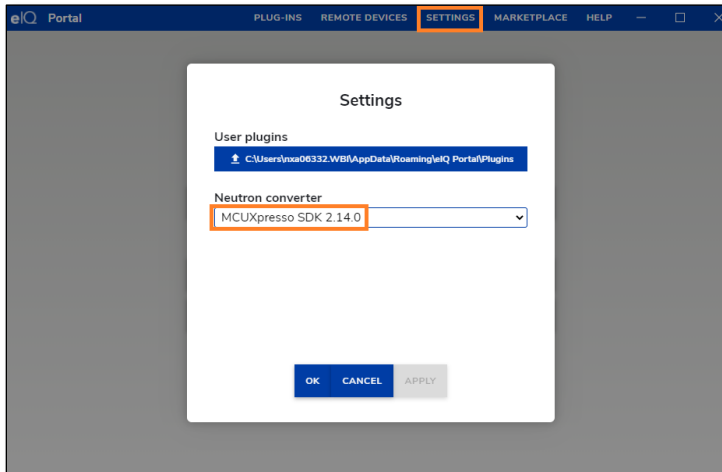
## 3.1 Convert Models

Use eIQ Toolkit to convert a pre-existing Mobilenet model into a Neutron optimized model.
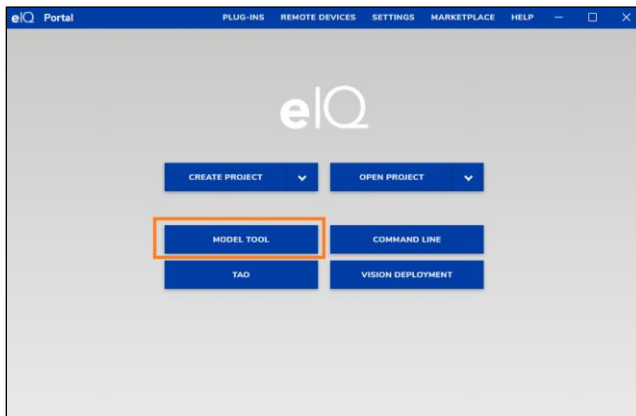
1. Download and unzip the Mobilenet v1 image classification TFLite model attached to this Community post or use this direct link.
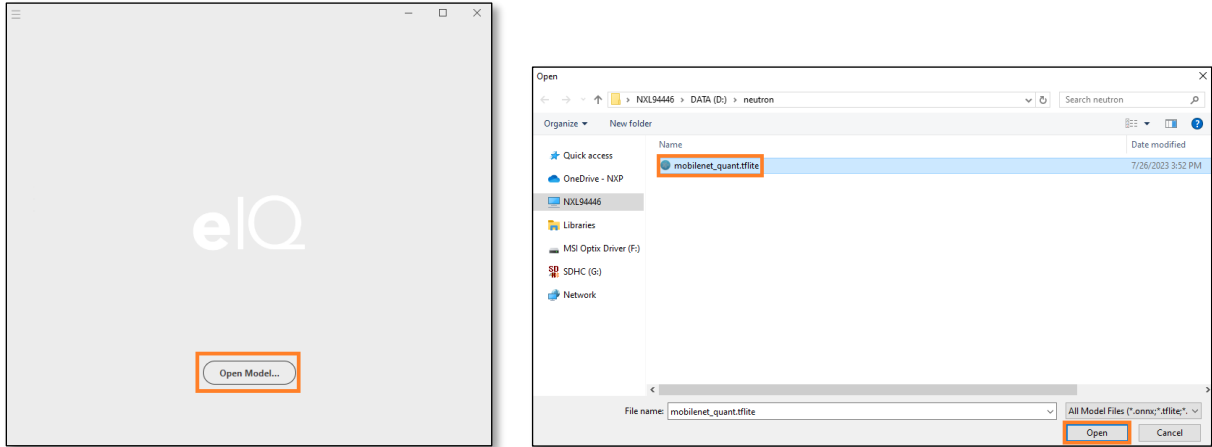2. Open eIQ Portal which is the GUI interface to eIQ Toolkit



3. After it opens, go to the menu bar and click on Settings. In the dialog box that comes up, make sure the Neutron Converter SDK matches the SDK version that you are using. In this case it should be set to **MCUXPresso SDK 2.14.0** by default and require no change. It is very important that the Neutron Converter tool version is targeted to the eIQ Neutron libraries used in the targeted SDK. Click on OK to save the setting and go back to the main screen.
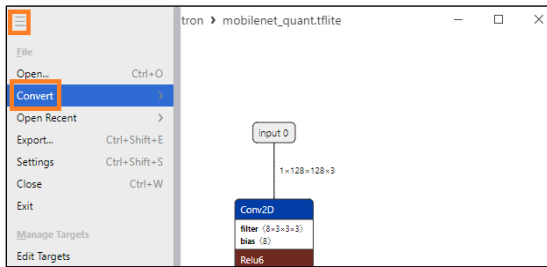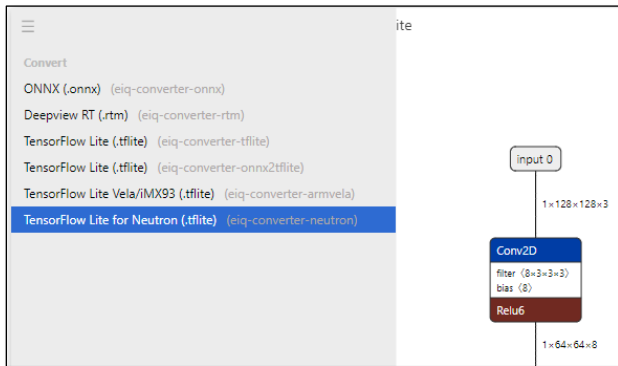


4. Click on Model Tool

5. Click on Open Model and select the **mobilenet_quant.tflite** model that was unzipped from the file downloaded earlier.



6. After the model is opened you'll see the various layers. Then click on the upper left corner to find the menu and select **Convert**
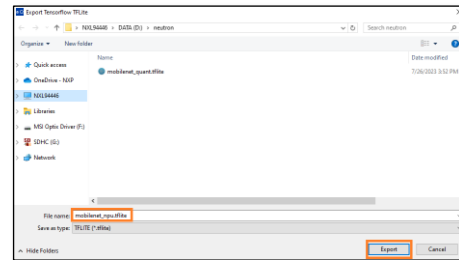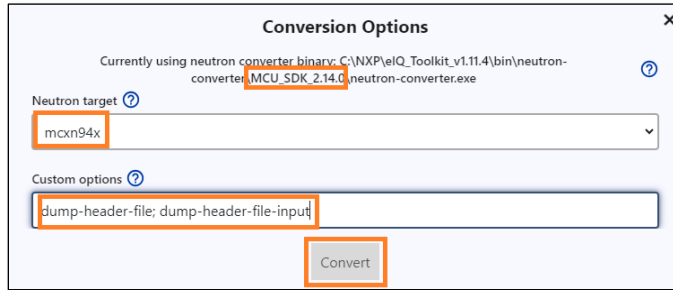


7. It may need to search for plugins. Let it.

8. After the plugins have been found, go back to the menu options and select Convert and then select **TensorFlow Lite for Neutron (.tflite)**

9. In the dialog box that comes up, select **mcxn94x** as the target and then in the custom options field put: **dump-header-file; dump-header-file-input**
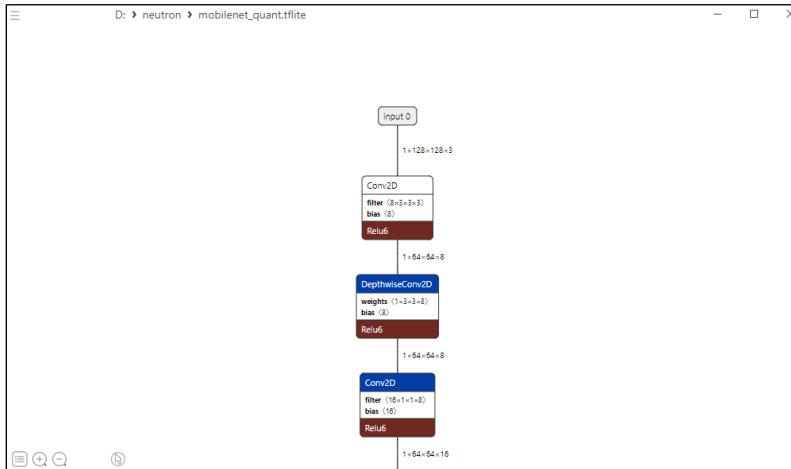
These two options will generate a C array of both the converted model and the input model, which we'll use to compare the performance of them. In typical situations you would only need the dump-header-file option though. Also make sure the Neutron Converter tool for MCU_SDK_2.14.0 is being used. Click on **Convert** and save the converted TFLite file on your hard drive as **mobilenet_npu.tflite**
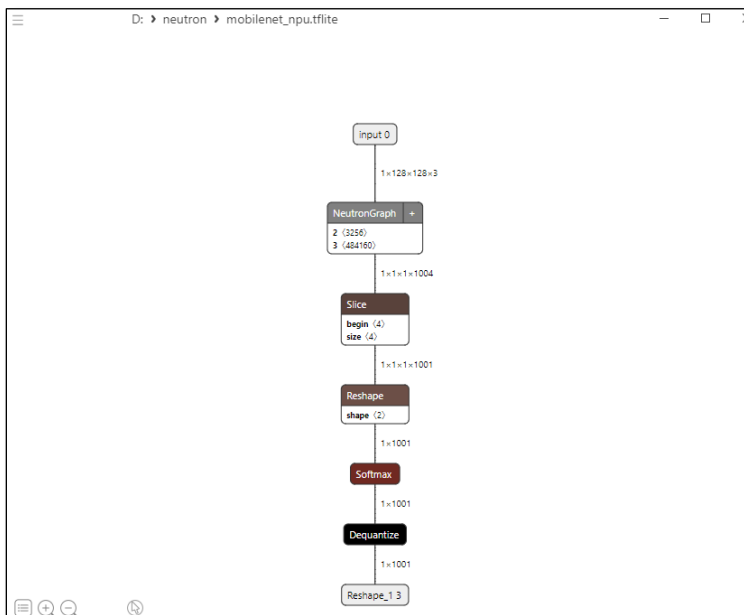
## 3.2 View Models

After conversion, you'll see the newly converted model optimized for the eIQ Neutron NPU pop up with the Model Tool. Take a moment to look at the original model compared to the new converted model.

1. The original TFLite file: **mobilenet_quant.tflite**



2. The Neutron converted file: **mobilenet_npu.tflite**



3. You can see how almost all the operators in the original model were replaced with a **NeutronGraph** operator. Those **NeutronGraph** operaters are what will be executed on the eIQ Neutron NPU when this model is ran on the MCX N94x. Any layers not converted to a NeutronGraph operator will instead be ran on the Cortex-M33 core.

4. Take a look at the file size of each of those header files and you can see that, in general, the NPU converted file will take up less flash space. Note that this might be counter-acted by the slightly increased size required for using the eIQ Neutron libraries.

5. During the conversion process the **dump-header-file** argument generated a header file for the NPU optimized model that can be used in the eIQ MCXUpresso SDK projects.

6. The **dump-header-file-input** argument generated a header file for the original non-converted model. This will be used so the inference time of the original model that only runs on the Cortex-M33 core can be compared to the NPU converted model that makes use of the eIQ Neutron NPU.

7. So let's run these models to see the performance improvements.
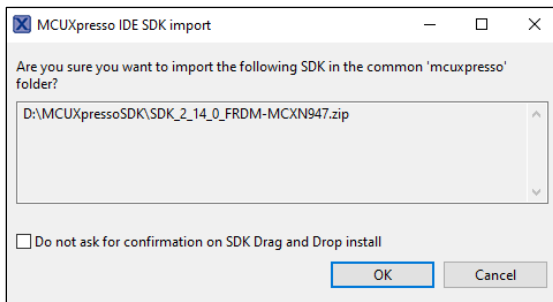
## 3.3 Run Models

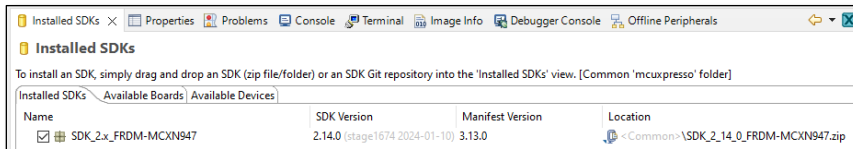Now let's use the MCUXpresso SDK eIQ Label Images example to run the models and see how long the inference time is.

8. Open MCUXpresso IDE and select a workspace location in an empty directory

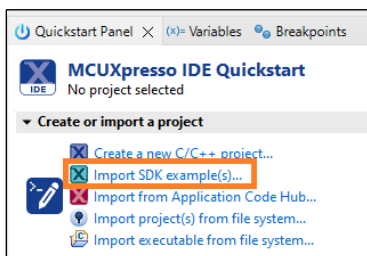9. Close the Welcome Screen tab by clicking on the X in that tab



10. Drag-and-drop the 2.14 FRDM-MCXN947 SDK zip file into the Installed SDKs window, located on a tab at the bottom of the screen named "Installed SDKs". You will get the following pop-up, so hit OK.
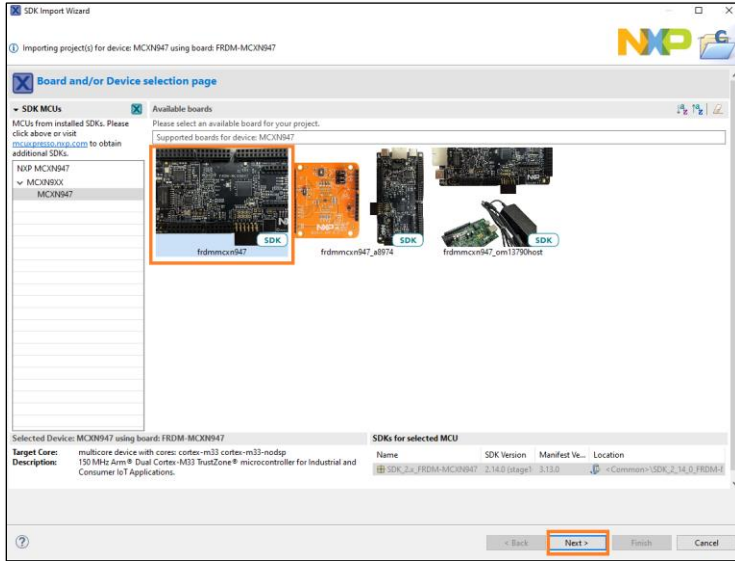


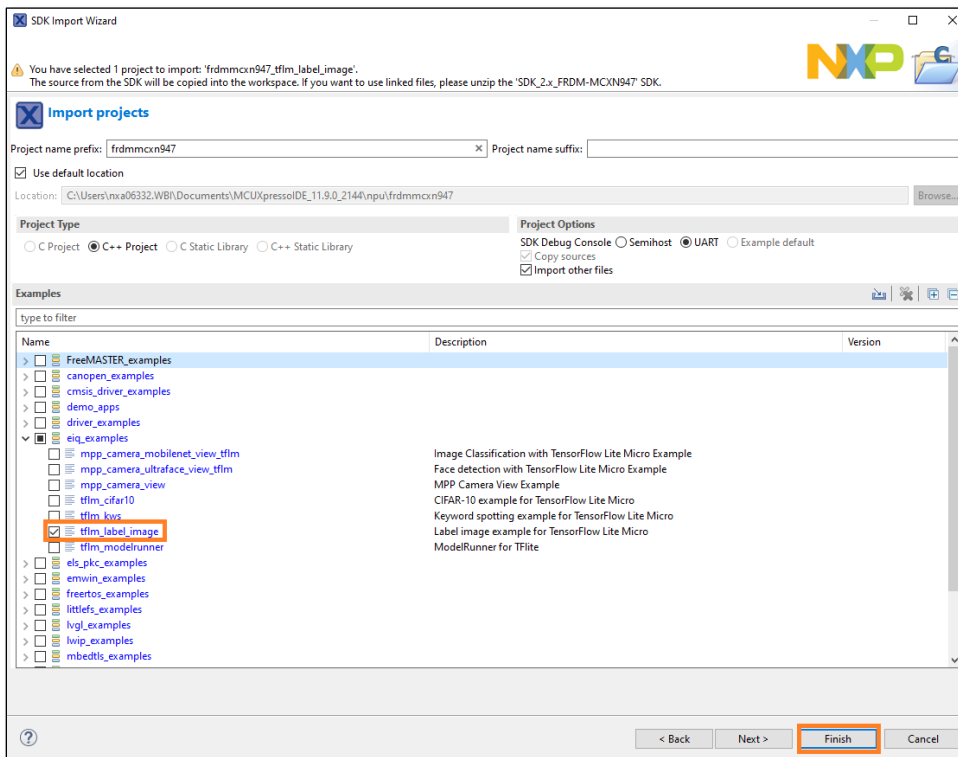11. Once imported, the Installed SDK panel will look something like this:



12. Next import the desired project. In the Quickstart Panel, select **Import SDK examples(s)...**
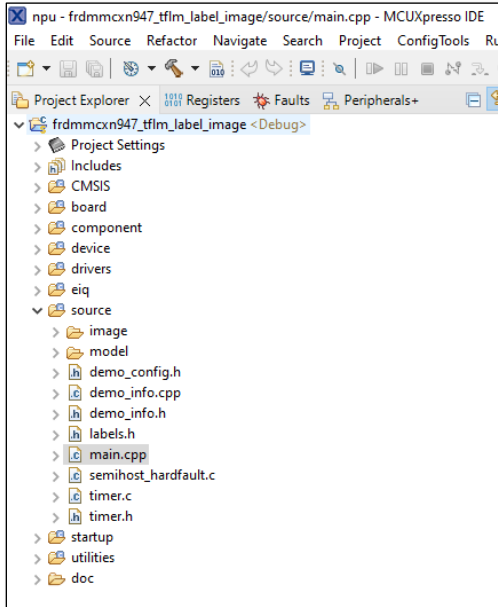
13. Select the **frdmmcxn947** board and click on **Next**



14. Under the **eiq_examples** category, select the **tflm_label_image** example. Then click on **Finish** to select that project.
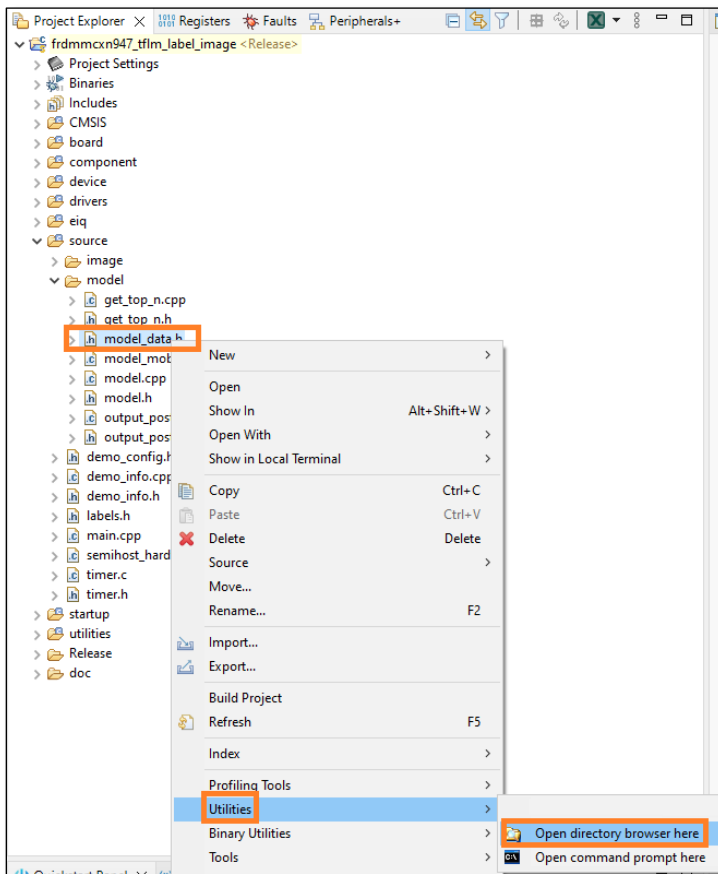
15. It should look like the following when done:



16. Now we need to import the models that were generated in the last section into this project.

17. Navigate down to the **source->model** folder and right click on **model_data.h** and then select **Utilities->Open directory browser here** to open the location of that file on your hard drive.

18. Now copy and paste the two .h header files that were generated in the previous section into this file location. It should look like the following when complete:



19. Now we need to slightly modify those two header files to add some information for the eIQ MCUXpresso SDK project that describe how much memory this model will require and to describe some of the normalization values that this model uses:

a) Open **model_data.h,** which contains the default model for this example, and find the following section of code and copy it.

b) Then in both **mobilenet_quant.h** and **mobilenet_ npu.h** copy that code above the array, overwriting the default #defines above the array. Make sure not to erase the commented lines at the top of the file as those comments will be used later. Note that the **MODEL_NAME** can be changed to "mobilenet".



20. After changing both files, now double click on **model.cpp** to open it.



21. Go to line 27 and change it to point to the non-NPU accelerated model in **mobilenet_quant.h**. It should look like the following after changed:

22. Next look at line 42 in that same **model.cpp** file to find where the model is loaded by the TFLM inference engine using the C array name **model_data**. Because the model array name in the new header file is the same as the original header file we replaced, **no change is needed here**. This is just for informational purposes only.

```
38 status_t MODEL_Init(void)
39 {
40     // Map the model into a usable data structure. This doesn't involve any
41     // copying or parsing, it's a very lightweight operation.
42     s_model = tflite::GetModel(model_data);
43     if (s_model->version() != TFLITE_SCHEMA_VERSION)
44     {
45         PRINTF("Model provided is schema version %d not equal "
46                 "to supported version %d.",
47                 s_model->version(), TFLITE_SCHEMA_VERSION);
48         return kStatus_Fail;
49     }
50
```

23. Compile the project by clicking on the Build button in the Quickstart Panel in the lower left hand corner. Make sure the correct project is listed as well.

   *Note: Because TFLM is now deployed as a library in MCUXpresso SDK 2.14, the default Debug no-optimization target can be used and still have the same inference time as the Release high-optimization target.*

   

24. Connect a USB C cable from your computer to the USB port on the FRDM-MCXN947 at J17

25. Open TeraTerm or other terminal program, and connect to the virtual COM port that debugger or UART-to-USB converter enumerated as. Use 115200 baud, 1 stop bit, no parity. There is a built-in serial terminal in MCUXpresso IDE that can be used as well:

26. Then in MCUXpresso IDE click on Debug



27. You'll see the following dialog box come up. It should list your debugger hardware. If there is a notice about a Firmware update available, that can be ignored for now. Click on **OK** to connect to the target.



28. You should see MCUXpresso IDE connect and download the program to your board in the Console tab.

29. You may see this dialog box come up the first time you run the program. Keep the default selection to run on the primary core and hit OK

30. Once complete, it will pause at the start of main(). Hit the Resume icon to run the program and look at the terminal tab.



31. Looks like there's an error: **Didn't find op for builtin opcode 'CONV_2D'**



32. Stop the debugger by clicking on the red square



33. When changing models, the list of operators needs to be updated as well. To fix the error, go back to MCUXpresso IDE and open the **model_mobilenet_ops_npu.cpp** file.

34. Inside the **MODEL_GetOpsResolver** function is a list of operators. If you open the **mobilenet_quant.h** header file you'll also find a list of operators used by the model in the comment block at the top of the file.



```
mobilenet_quant.h  ×
  1  // Copyright 2022 NXP.
  2  // All rights reserved.
  3  // SPDX-License-Identifier: BSD-3-Clause
  4
  5  // Neutron Converter Version: 1.2.0+0X84d37e1f
  6  // Neutron Microcode Version: 0X84d37e1f
  7
  8  /*
  9  // Register operators for TFLite Micro
 10  microOpResolver.AddConv2D();
 11  microOpResolver.AddDepthwiseConv2D();
 12  microOpResolver.AddAveragePool2D();
 13  microOpResolver.AddReshape();
 14  microOpResolver.AddSoftmax();
 15  microOpResolver.AddDequantize();
 16  */
 17
 18  #ifdef __arm__
 19  #include <cmsis_compiler.h>
 20  #else
 21  #define __ALIGNED(x) __attribute__((aligned(x)))
 22  #endif
 23
```

35. Copy that list into the **MODEL_GetOpsResolver** function replacing the original list.

```
  8  #include "tensorflow/lite/micro/kernels/micro_ops.h"
  9  #include "tensorflow/lite/micro/micro_mutable_op_resolver.h"
 10  #include "tensorflow/lite/micro/kernels/neutron/neutron.h"
 11
 12  tflite::MicroOpResolver &MODEL_GetOpsResolver()
 13  {
 14      static tflite::MicroMutableOpResolver<5> s_microOpResolver;
 15
 16      microOpResolver.AddConv2D();
 17      microOpResolver.AddDepthwiseConv2D();
 18      microOpResolver.AddAveragePool2D();
 19      microOpResolver.AddReshape();
 20      microOpResolver.AddSoftmax();
 21      microOpResolver.AddDequantize();
 22
 23      return s_microOpResolver;
 24  }
```
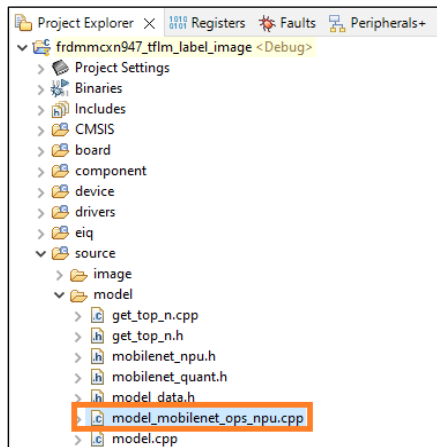
36. Also update the number of operators, the name of the operator variable (**s_microOpResolver** to **microOpResolver**), and the return variable name so they all use the **microOpResolver** name.

```
#include "tensorflow/lite/micro/kernels/micro_ops.h"
#include "tensorflow/lite/micro/micro_mutable_op_resolver.h"
#include "tensorflow/lite/micro/kernels/neutron/neutron.h"

tflite::MicroOpResolver &MODEL_GetOpsResolver()
{
    static tflite::MicroMutableOpResolver<6> microOpResolver;

    microOpResolver.AddConv2D();
    microOpResolver.AddDepthwiseConv2D();
    microOpResolver.AddAveragePool2D();
    microOpResolver.AddReshape();
    microOpResolver.AddSoftmax();
    microOpResolver.AddDequantize();

    return microOpResolver;
}
```

37. Recompile and reprogram the board using the previous steps.
38. You should now see the following on the serial terminal:

```
Installed SDKs  Properties  Problems  Console  Terminal  ×
COM4  ×
Label image example using a TensorFlow Lite Micro model.
Detection threshold: 23%
Model: mobilenet

Static data processing:
----------------------------------------
    Inference time: 336 ms
    Detected: stopwatch (87%)
----------------------------------------

Camera data processing:
Camera input is currently not supported on this device
```

39. Stop the debugger by clicking on the red square



40. Now run the program again, but this time with the Neutron NPU accelerated version of the model.

41. Re-open **model.cpp** and this time change line 27 to point to the Neutron NPU converted version of the model in the **mobilenet_ npu.h** file:



42. Re-open **model_mobilenet_ops_npu.cpp** and update the **MODEL_GetOpsResolver** function with the operators listed in the comment block of the **mobilenet_npu.h** file. Also make sure to update the array size.



43. Build and program the program as before.

44. This time you should see the following on the terminal. That's over a 28x improvement in inference time, with the same confidence percentage on this static image.
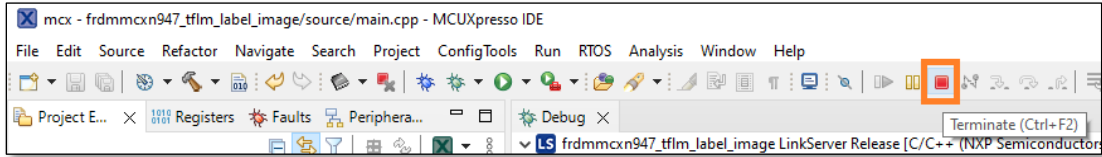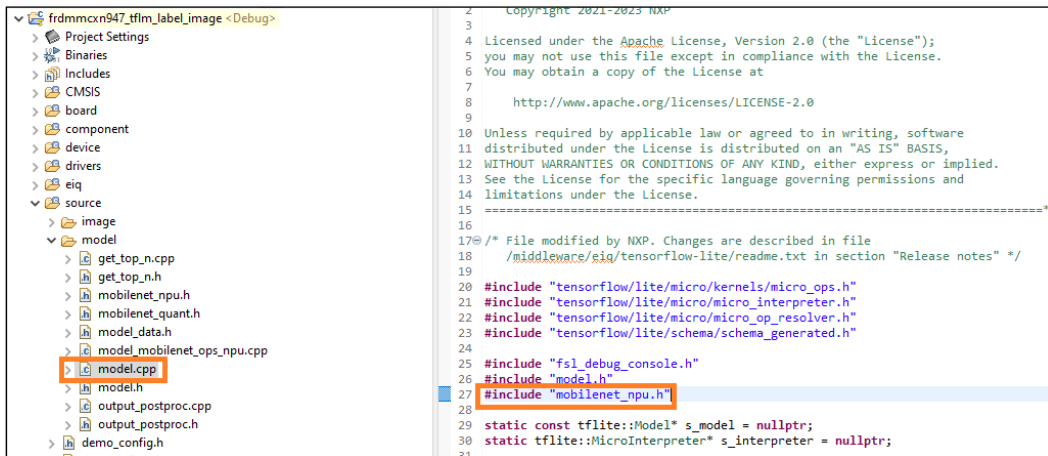


```
📁 Installed SDKs  ▣ Properties  👤 Problems  🖥 Console  💻 Terminal ✕
🖥 COM4 ✕
Label image example using a TensorFlow Lite Micro model.
Detection threshold: 23%
Model: mobilenet

Static data processing:
-----------------------------------------
    Inference time: 12 ms
    Detected: stopwatch (86%)
-----------------------------------------

Camera data processing:
Camera input is currently not supported on this device
```
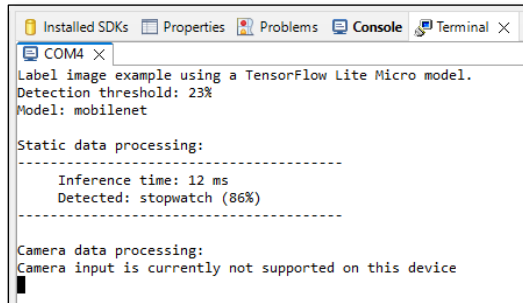
45. The decrease in inference time is very model dependent depending on how well that specific model could be optimized for the NPU.

# 4   Conclusion

This lab demonstrated how the eIQ Neutron NPU on MCX N devices can significantly decrease inference time on quantized models. These same steps can be used to benchmark other quantized models to see the performance improvements that the eIQ Neutron NPU can have.