# Post-training integer quantization

> Indented block

[View on TensorFlow.org](#)   [Run in Google Colab](#)   [View source on GitHub](#)

## Overview

[TensorFlow Lite](#) now supports converting all model values (weights and activations) to 8-bit intege
TensorFlow Lite's flat buffer format. This results in a 4x reduction in model size and a 3 to 4x perfo
In addition, this fully quantized model can be consumed by integer-only hardware accelerators.

In contrast to [post-training "on-the-fly" quantization](#)—which stores only the weights as 8-bit integer
weights *and* activations during model conversion.

In this tutorial, you'll train an MNIST model from scratch, check its accuracy in TensorFlow, and the
flatbuffer with full quantization. Finally, you'll check the accuracy of the converted model and comp

# Build an MNIST model

# Setup

```
1 import logging
2 logging.getLogger("tensorflow").setLevel(logging.DEBUG)
3
4 try:
5   # %tensorflow_version only exists in Colab.
6   import tensorflow.compat.v2 as tf
7 except Exception:
8   pass
9 tf.enable_v2_behavior()
10
11 from tensorflow import keras
12 import numpy as np
13 import pathlib
```

# Train and export the model

```
 1 # Load MNIST dataset
 2 mnist = keras.datasets.mnist
 3 (train_images, train_labels), (test_images, test_labels) = mnist.load_data()
 4
 5 # Normalize the input image so that each pixel value is between 0 to 1.
 6 train_images = train_images / 255.0
 7 test_images = test_images / 255.0
 8
 9 # Define the model architecture
10 model = keras.Sequential([
11   keras.layers.InputLayer(input_shape=(28, 28)),
12   keras.layers.Reshape(target_shape=(28, 28, 1)),
13   keras.layers.Conv2D(filters=12, kernel_size=(3, 3), activation=tf.nn.relu),
14   keras.layers.MaxPooling2D(pool_size=(2, 2)),
15   keras.layers.Flatten(),
16   keras.layers.Dense(10, activation=tf.nn.softmax)
17 ])
18
19 # Train the digit classification model
20 model.compile(optimizer='adam',
21               loss='sparse_categorical_crossentropy',
22               metrics=['accuracy'])
23 model.fit(
24   train_images,
25   train_labels,
26   epochs=1,
27   validation_data=(test_images, test_labels)
28 )
```

This training won't take long because you're training the model for just a single epoch, which trains

## Convert to a TensorFlow Lite model

Using the Python TFLiteConverter, you can now convert the trained model into a TensorFlow Lite m

Now load the model using the TFLiteConverter:

```
 1 converter = tf.lite.TFLiteConverter.from_keras_model(model)
 2 tflite_model = converter.convert()
```

Write it out to a .tflite file:

```
 1 tflite_models_dir = pathlib.Path("/tmp/mnist_tflite_models/")
 2 tflite_models_dir.mkdir(exist_ok=True, parents=True)
```

```
 1 tflite_model_file = tflite_models_dir/"mnist_model.tflite"
 2 tflite_model_file.write_bytes(tflite_model)
```

Now you have a trained MNIST model that's converted to a .tflite file, but it's still using 32-bit flo

So let's convert the model again, this time using quantization...

## Convert using quantization

First, first set the `optimizations` flag to optimize for size:

```
1 converter.optimizations = [tf.lite.Optimize.OPTIMIZE_FOR_SIZE]
```

Now, in order to create quantized values with an accurate dynamic range of activations, you need t

```
1 mnist_train, _ = tf.keras.datasets.mnist.load_data()
2 images = tf.cast(mnist_train[0], tf.float32) / 255.0
3 mnist_ds = tf.data.Dataset.from_tensor_slices((images)).batch(1)
4 def representative_data_gen():
5   for input_value in mnist_ds.take(100):
6     yield [input_value]
7
8 converter.representative_dataset = representative_data_gen
```

Finally, convert the model to TensorFlow Lite format:

```
1 tflite_model_quant = converter.convert()
2 tflite_model_quant_file = tflite_models_dir/"mnist_model_quant.tflite"
3 tflite_model_quant_file.write_bytes(tflite_model_quant)
```

Note how the resulting file is approximately `1/4` the size:

```
1 !ls -lh {tflite_models_dir}
```

Your model should now be fully quantized. However, if you convert a model that includes any oper
those ops are left in floating point. This allows for conversion to complete so you have a smaller a
be compatible with some ML accelerators that require full integer quantization. Also, by default, th
outputs, which also is not compatible with some accelerators.

So to ensure that the converted model is fully quantized (make the converter throw an error if it en
and to use integers for the model's input and output, you need to convert the model again using th

```
1 converter.target_spec.supported_ops = [tf.lite.OpsSet.TFLITE_BUILTINS_INT8]
2 converter.inference_input_type = tf.uint8
3 converter.inference_output_type = tf.uint8
4
5 tflite_model_quant = converter.convert()
6 tflite_model_quant_file = tflite_models_dir/"mnist_model_quant_io.tflite"
7 tflite_model_quant_file.write_bytes(tflite_model_quant)
```

In this example, the resulting model size remains the same because all operations successfully qu
model now uses quantized input and output, making it compatible with more accelerators, such as

In the following sections, notice that we are now handling two TensorFlow Lite models: `tflite_mo` uses floating-point parameters, and `tflite_model_quant_file` is the same model converted with and output.

## Run the TensorFlow Lite models

Run the TensorFlow Lite model using the Python TensorFlow Lite Interpreter.

## Load the model into the interpreters

```
1 interpreter = tf.lite.Interpreter(model_path=str(tflite_model_file))
2 interpreter.allocate_tensors()
```

```
1 interpreter_quant = tf.lite.Interpreter(model_path=str(tflite_model_quant_file))
2 interpreter_quant.allocate_tensors()
3 input_index_quant = interpreter_quant.get_input_details()[0]["index"]
4 output_index_quant = interpreter_quant.get_output_details()[0]["index"]
```

## Test the models on one image

First test it on the float model:

```
1 test_image = np.expand_dims(test_images[0], axis=0).astype(np.float32)
2
3 input_index = interpreter.get_input_details()[0]["index"]
4 output_index = interpreter.get_output_details()[0]["index"]
5 interpreter.set_tensor(input_index, test_image)
6 interpreter.invoke()
7 predictions = interpreter.get_tensor(output_index)
```

```
1 import matplotlib.pylab as plt
2
3 plt.imshow(test_images[0])
4 template = "True:{true}, predicted:{predict}"
5 _ = plt.title(template.format(true= str(test_labels[0]),
6                               predict=str(np.argmax(predictions[0]))))
7 plt.grid(False)
```

Now test the quantized model (using the uint8 data):

```
1 input_index = interpreter_quant.get_input_details()[0]["index"]
2 output_index = interpreter_quant.get_output_details()[0]["index"]
3 interpreter_quant.set_tensor(input_index, test_image)
4 interpreter_quant.invoke()
5 predictions = interpreter_quant.get_tensor(output_index)
```

```
1 plt.imshow(test_images[0])
2 template = "True:{true}, predicted:{predict}"
3 _ = plt.title(template.format(true= str(test_labels[0]),
4                               predict=str(np.argmax(predictions[0]))))
5 plt.grid(False)
```

## Evaluate the models

```
 1 # A helper function to evaluate the TF Lite model using "test" dataset.
 2 def evaluate_model(interpreter):
 3   input_index = interpreter.get_input_details()[0]["index"]
 4   output_index = interpreter.get_output_details()[0]["index"]
 5
 6   # Run predictions on every image in the "test" dataset.
 7   prediction_digits = []
 8   for test_image in test_images:
 9     # Pre-processing: add batch dimension and convert to float32 to match with
10     # the model's input data format.
11     test_image = np.expand_dims(test_image, axis=0).astype(np.float32)
12     interpreter.set_tensor(input_index, test_image)
13
14     # Run inference.
15     interpreter.invoke()
16
17     # Post-processing: remove batch dimension and find the digit with highest
18     # probability.
19     output = interpreter.tensor(output_index)
20     digit = np.argmax(output()[0])
21     prediction_digits.append(digit)
22
23   # Compare prediction results with ground truth labels to calculate accuracy.
24   accurate_count = 0
25   for index in range(len(prediction_digits)):
26     if prediction_digits[index] == test_labels[index]:
27       accurate_count += 1
28   accuracy = accurate_count * 1.0 / len(prediction_digits)
29
30   return accuracy
```

```
1 print(evaluate_model(interpreter))
```

Repeat the evaluation on the fully quantized model using the uint8 data:

```
1 # NOTE: Colab runs on server CPUs, and TensorFlow Lite currently
2 # doesn't have super optimized server CPU kernels. So this part may be
3 # slower than the above float interpreter. But for mobile CPUs, considerable
4 # speedup can be observed.
5
6 print(evaluate_model(interpreter_quant))
```

In this example, you have fully quantized a model with almost no difference in the accuracy, compa