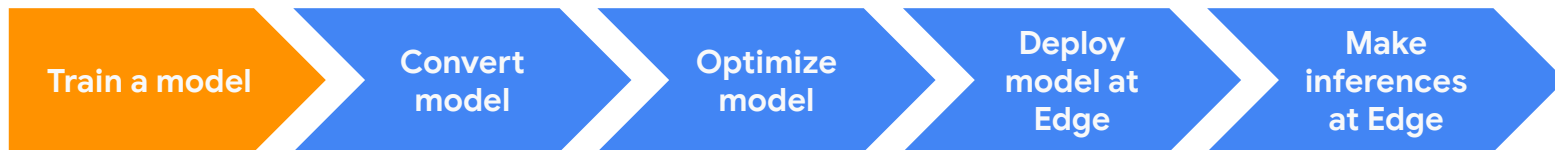
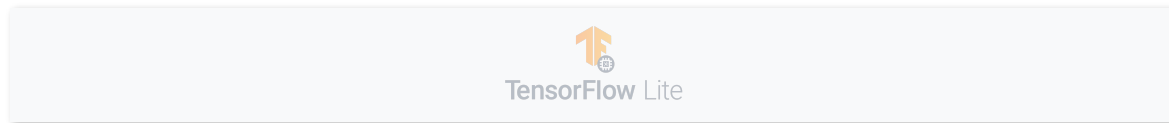
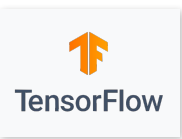


# Quantization-aware training

Optimizing for TFLite from the beginning...



Laurence Moroney, Google





**Train a model**

**Convert  
model**

**Optimize  
model**

**Deploy  
model at  
Edge**

**Make  
inferences  
at Edge**



Using QAT

```
# Load MNIST dataset
mnist = keras.datasets.mnist

(train_images, train_labels), (test_images, test_labels) =
    mnist.load_data()

# Normalize the input image so that each pixel value is between 0 to 1.
train_images = train_images / 255.0
test_images = test_images / 255.0

# Define the model architecture.
model = keras.Sequential([
    keras.layers.InputLayer(input_shape=(28, 28)),
    keras.layers.Reshape(target_shape=(28, 28, 1)),
    keras.layers.Conv2D(filters=12, kernel_size=(3, 3), activation='relu'),
    keras.layers.MaxPooling2D(pool_size=(2, 2)),
    keras.layers.Flatten(),
    keras.layers.Dense(10)
])
```

loss: 0.2724 - accuracy: 0.9244 -  
val\_loss: 0.1085 - val\_accuracy: 0.9695

```
import tensorflow_model_optimization as tfmot

quantize_model = tfmot.quantization.keras.quantize_model

# q_aware stands for for quantization aware.
q_aware_model = quantize_model(model)

# `quantize_model` requires a recompile.
q_aware_model.compile(optimizer='adam',
                      loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),
                      metrics=['accuracy'])
```

```
import tensorflow_model_optimization as tfmot
```

```
quantize_model = tfmot.quantization.keras.quantize_model
```

```
# q_aware stands for quantization aware.
```

```
q_aware_model = quantize_model(model)
```

```
# `quantize_model` requires a recompile.
```

```
q_aware_model.compile(optimizer='adam',  
                      loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),  
                      metrics=['accuracy'])
```



```
import tensorflow_model_optimization as tfmot
```

```
quantize_model = tfmot.quantization.keras.quantize_model
```

```
# q_aware stands for for quantization aware.
```

```
q_aware_model = quantize_model(model)
```

```
# `quantize_model` requires a recompile.
```

```
q_aware_model.compile(optimizer='adam',  
                      loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),  
                      metrics=['accuracy'])
```

```
import tensorflow_model_optimization as tfmot
```

```
quantize_model = tfmot.quantization.keras.quantize_model
```

```
# q_aware stands for for quantization aware.
```

```
q_aware_model = quantize_model(model)
```

```
# `quantize_model` requires a recompile.
```

```
q_aware_model.compile(optimizer='adam',  
                      loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),  
                      metrics=['accuracy'])
```

```
import tensorflow_model_optimization as tfmot
```

```
quantize_model = tfmot.quantization.keras.quantize_model
```

```
# q_aware stands for quantization aware.
```

```
q_aware_model = quantize_model(model)
```

```
# `quantize_model` requires a recompile.
```

```
q_aware_model.compile(optimizer='adam',  
                      loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),  
                      metrics=['accuracy'])
```

loss: 0.2724 - accuracy: 0.9244 -  
val\_loss: 0.1085 - val\_accuracy: 0.9695

loss: 0.2724 - accuracy: 0.9244 -  
val\_loss: 0.1085 - val\_accuracy: 0.9695

loss: 0.1315 - accuracy: 0.9589 -  
val\_loss: 0.1360 - val\_accuracy: 0.9600



Your Turn