

File Edit View Insert Cell Kernel Widgets Help

Trusted

Python 3 (ipykernel) O



## The Ultimate Compression Pipeline ~ Ajay Maheshwari

In [25]: ! pip install -q tensorflow-model-optimization

```
In [26]: import tensorflow as tf
import tf_keras as keras

import numpy as np
import tempfile
import zipfile
import os
```

```
In [27]: def get_gzipped_model_size(model):
    with tempfile.NamedTemporaryFile(suffix=".h5") as temp_file:
        model.save(temp_file.name)

    _, zipped_file = tempfile.mkstemp('.zip')
    with zipfile.ZipFile(zipped_file, 'w', compression=zipfile.ZIP_DEFLATED) as f:
        f.write(temp_file.name)

    # print(f"Zipped model is saved at: {zipped_file}")

    x = os.path.getsize(zipped_file)

    os.remove(zipped_file)
    # print(f"Temporary zip file removed: {zipped_file}")

    return x / 1000

def print_model_weights_sparsity(model):
    for layer in model.layers:
        if isinstance(layer, keras.layers.Wrapper):
            weights = layer.trainable_weights
        else:
            weights = layer.weights
        for weight in weights:
            if "kernel" not in weight.name or "centroid" in weight.name:
                continue
            weight_size = weight.numpy().size
            zero_num = np.count_nonzero(weight == 0)
            print(
                f"{weight.name}: {zero_num/weight_size:.2%} sparsity ",
                f"({zero_num}/{weight_size})",
            )

def get_gzipped_model_size2(file):
    _, zipped_file = tempfile.mkstemp('.zip')
    with zipfile.ZipFile(zipped_file, 'w', compression=zipfile.ZIP_DEFLATED) as f:
        f.write(file)

    return os.path.getsize(zipped_file)/1000

def eval_model(interpreter):
    input_index = interpreter.get_input_details()[0]["index"]
    output_index = interpreter.get_output_details()[0]["index"]

    prediction_digits = []
    for i, test_image in enumerate(test_images):
        test_image = np.expand_dims(test_image, axis=0).astype(np.float32)
        interpreter.set_tensor(input_index, test_image)

        interpreter.invoke()

        output = interpreter.tensor(output_index)
        digit = np.argmax(output[0])
        prediction_digits.append(digit)

    prediction_digits = np.array(prediction_digits)
    accuracy = (prediction_digits == test_labels).mean()
    return accuracy

model_acc = []
model_sz = []
```

## Creating a Base Model - 1.0

```
In [28]: # 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

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=tf.nn.relu),
keras.layers.MaxPooling2D(pool_size=(2, 2)),
keras.layers.Flatten(),
keras.layers.Dense(10)
])

opt = keras.optimizers.Adam(learning_rate=1e-3)

# Train the digit classification model
model.compile(optimizer=opt,
              loss=keras.losses.SparseCategoricalCrossentropy(from_logits=True),
              metrics=['accuracy'])

model.fit(
    train_images,
    train_labels,
    validation_split=0.1,
    epochs=10
)

WARNING:absl:At this time, the v2.11+ optimizer `tf.keras.optimizers.Adam` runs slowly on M1/M2 Macs, please use the legacy TF-Keras optimizer instead, located at `tf.keras.optimizers.legacy.Adam`.

Epoch 1/10
1688/1688 [=====] - 5s 3ms/step - loss: 0.3053 - accuracy: 0.9149 - val_loss: 0.1260 - val_accuracy: 0.9688
Epoch 2/10
1688/1688 [=====] - 5s 3ms/step - loss: 0.1223 - accuracy: 0.9652 - val_loss: 0.0858 - val_accuracy: 0.9778
Epoch 3/10
1688/1688 [=====] - 5s 3ms/step - loss: 0.0893 - accuracy: 0.9745 - val_loss: 0.0754 - val_accuracy: 0.9782
Epoch 4/10
1688/1688 [=====] - 5s 3ms/step - loss: 0.0732 - accuracy: 0.9779 - val_loss: 0.0678 - val_accuracy: 0.9820
Epoch 5/10
1688/1688 [=====] - 5s 3ms/step - loss: 0.0632 - accuracy: 0.9814 - val_loss: 0.0652 - val_accuracy: 0.9820
Epoch 6/10
1688/1688 [=====] - 5s 3ms/step - loss: 0.0553 - accuracy: 0.9830 - val_loss: 0.0624 - val_accuracy: 0.9833
Epoch 7/10
1688/1688 [=====] - 5s 3ms/step - loss: 0.0494 - accuracy: 0.9843 - val_loss: 0.0662 - val_accuracy: 0.9810
Epoch 8/10
1688/1688 [=====] - 4s 3ms/step - loss: 0.0453 - accuracy: 0.9861 - val_loss: 0.0562 - val_accuracy: 0.9837
Epoch 9/10
1688/1688 [=====] - 4s 3ms/step - loss: 0.0407 - accuracy: 0.9880 - val_loss: 0.0625 - val_accuracy: 0.9825
Epoch 10/10
1688/1688 [=====] - 4s 3ms/step - loss: 0.0373 - accuracy: 0.9883 - val_loss: 0.0579 - val_accuracy: 0.9845

```

Out[28]: <tf\_keras.src.callbacks.History at 0x2a220e6d0>

### Evaluate the baseline model and save it for later usage

```

In [29]: _, baseline_model_accuracy = model.evaluate(
    test_images, test_labels, verbose=0)

print('Baseline test accuracy:', baseline_model_accuracy)

sz = get_gzipped_model_size(model)
print("Base model size: ", sz, ' KB' )

model_acc.append(baseline_model_accuracy)
model_sz.append(sz)

Baseline test accuracy: 0.9805999994277954
Base model size: 235.396 KB

/Users/ajaymaheshwari/anaconda3/lib/python3.11/site-packages/tf_keras/src/engine/training.py:3098: UserWarning: You are saving your model as an HDF5 file via `model.save()`. This file format is considered legacy. We recommend using instead the native TF-Keras format, e.g. `model.save('my_model.keras')`.
      saving_api.save_model()

```

----- Checkpoint Point 1 -----

### Pruning and then fine-tuning the model - 2.0

```

In [30]: import tensorflow_model_optimization as tfmot

pruning_params = {
    'pruning_schedule': tfmot.sparsity.keras.ConstantSparsity(0.5, begin_step=0, frequency=100)
}

callbacks = [
    tfmot.sparsity.keras.UpdatePruningStep()
]

pruned_model = model

prune_low_magnitude = tfmot.sparsity.keras.prune_low_magnitude

```

```

pruned_model = prune_low_magnitude(pruned_model, **pruning_params)

# learning rate for fine-tuning
opt = keras.optimizers.Adam(learning_rate=1e-5)

pruned_model.compile(
    loss=keras.losses.SparseCategoricalCrossentropy(from_logits=True),
    optimizer=opt,
    metrics=['accuracy'])

# Fine-tune model
pruned_model.fit(
    train_images,
    train_labels,
    epochs=3,
    validation_split=0.1,
    callbacks=callbacks)

stripped_pruned_model = tfmot.sparsity.keras.strip_pruning(pruned_model)

print_model_weights_sparsity(stripped_pruned_model)

pruned_model = stripped_pruned_model

```

WARNING:absl:At this time, the v2.11+ optimizer `tf.keras.optimizers.Adam` runs slowly on M1/M2 Macs, please use the legacy TF-Keras optimizer instead, located at `tf.keras.optimizers.legacy.Adam`.

```

Epoch 1/3
1688/1688 [=====] - 6s 3ms/step - loss: 0.1220 - accuracy: 0.9582 - val_loss: 0.0971 - val_accuracy: 0.9698
Epoch 2/3
1688/1688 [=====] - 5s 3ms/step - loss: 0.0723 - accuracy: 0.9764 - val_loss: 0.0726 - val_accuracy: 0.9785
Epoch 3/3
1688/1688 [=====] - 6s 4ms/step - loss: 0.0566 - accuracy: 0.9826 - val_loss: 0.0653 - val_accuracy: 0.9812
conv2d_3/kernel:0: 50.00% sparsity (54/108)
dense_3/kernel:0: 50.00% sparsity (10140/20280)

```

```
In [31]: pruned_model.compile(
    loss=keras.losses.SparseCategoricalCrossentropy(from_logits=True),
    optimizer=opt,
    metrics=['accuracy'])
```

```
In [ ]:
```

```
In [32]: _, pruned_model_accuracy = pruned_model.evaluate(
    test_images, test_labels, verbose=0)

print('Pruned Model test accuracy:', pruned_model_accuracy)

sz = get_gzipped_model_size(pruned_model)
print("Stripped model size: ", sz, ' KB')

model_acc.append(pruned_model_accuracy)
model_sz.append(sz)
```

```
Pruned Model test accuracy: 0.9731000065803528
Stripped model size: 204.383 KB
```

```
----- Checkpoint Point 2 -----
```

## Knowledge Distillation - 3.0

```
In [33]: from keras import layers
from keras import ops
import numpy as np

class Distiller(keras.Model):
    def __init__(self, student, teacher):
        super().__init__()
        self.teacher = teacher
        self.student = student

    def compile(
        self,
        optimizer,
        metrics,
        student_loss_fn,
        distillation_loss_fn,
        alpha=0.1,
        temperature=3,
    ):
        super().compile(optimizer=optimizer, metrics=metrics)
        self.student_loss_fn = student_loss_fn
        self.distillation_loss_fn = distillation_loss_fn
        self.alpha = alpha
        self.temperature = temperature

    def compute_loss(
        self,
        x=None, y=None, y_pred=None, sample_weight=None, allow_empty=False
    ):
        teacher_pred = self.teacher(x, training=False)
        student_loss = self.student_loss_fn(y, y_pred)
```

```

student_loss = self.student_loss_fn(y_true, y_pred)

distillation_loss = self.distillation_loss_fn(
    ops.softmax(teacher_pred / self.temperature, axis=1),
    ops.softmax(y_pred / self.temperature, axis=1),
) * (self.temperature**2)

loss = self.alpha * student_loss + (1 - self.alpha) * distillation_loss
return loss

def call(self, x):
    return self.student(x)

```

```

In [34]: teacher = keras.Sequential(
    [
        keras.layers.InputLayer(input_shape=(28, 28)),
        keras.layers.Reshape(target_shape=(28, 28, 1)),
        keras.layers.Conv2D(filters=16, kernel_size=(3, 3), activation=tf.nn.relu),
        keras.layers.MaxPooling2D(pool_size=(2, 2)),
        keras.layers.Conv2D(filters=24, kernel_size=(3, 3), activation=tf.nn.relu),
        keras.layers.MaxPooling2D(pool_size=(2, 2)),
        keras.layers.Flatten(),
        keras.layers.Dense(units=128, activation=tf.nn.relu),
        keras.layers.Dense(10)
    ],
    name="teacher",
)

# base one
# teacher = 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=tf.nn.relu),
#         keras.layers.MaxPooling2D(pool_size=(2, 2)),
#         keras.layers.Flatten(),
#         keras.layers.Dense(10)
#     ],
#     name="teacher",
# )

teacher.compile(
    optimizer=keras.optimizers.Adam(),
    loss=keras.losses.SparseCategoricalCrossentropy(from_logits=True),
    metrics=[keras.metrics.SparseCategoricalAccuracy()],
)

```

WARNING:absl:At this time, the v2.11+ optimizer `tf.keras.optimizers.Adam` runs slowly on M1/M2 Macs, please use the legacy TF-Keras optimizer instead, located at `tf.keras.optimizers.legacy.Adam`.

```

In [35]: teacher.fit(train_images, train_labels, epochs=12)
teacher.evaluate(test_images, test_labels)

Epoch 1/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.1566 - sparse_categorical_accuracy: 0.9524
Epoch 2/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.0517 - sparse_categorical_accuracy: 0.9838
Epoch 3/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.0351 - sparse_categorical_accuracy: 0.9889
Epoch 4/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.0279 - sparse_categorical_accuracy: 0.9912
Epoch 5/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.0202 - sparse_categorical_accuracy: 0.9936
Epoch 6/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.0166 - sparse_categorical_accuracy: 0.9945
Epoch 7/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.0133 - sparse_categorical_accuracy: 0.9956
Epoch 8/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.0110 - sparse_categorical_accuracy: 0.9962
Epoch 9/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.0089 - sparse_categorical_accuracy: 0.9970
Epoch 10/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.0083 - sparse_categorical_accuracy: 0.9971
Epoch 11/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.0060 - sparse_categorical_accuracy: 0.9980
Epoch 12/12
1875/1875 [=====] - 8s 4ms/step - loss: 0.0071 - sparse_categorical_accuracy: 0.9976
313/313 [=====] - 1s 2ms/step - loss: 0.0358 - sparse_categorical_accuracy: 0.9914

```

```
Out[35]: [0.03579017519950867, 0.9914000034332275]
```

```

In [36]: distiller = Distiller(student=pruned_model, teacher=teacher)
distiller.compile(
    optimizer=keras.optimizers.Adam(),
    metrics=[keras.metrics.SparseCategoricalAccuracy()],
    student_loss_fn=keras.losses.SparseCategoricalCrossentropy(from_logits=True),
    distillation_loss_fn=keras.losses.KLDivergence(),
    alpha=0.1,
    temperature=10,
)

```

WARNING:absl:At this time, the v2.11+ optimizer `tf.keras.optimizers.Adam` runs slowly on M1/M2 Macs, please use the legacy TF-Keras optimizer instead, located at `tf.keras.optimizers.legacy.Adam`.

```

In [37]: # Distill teacher se student
distiller.fit(train_images, train_labels, epochs=5)
distiller.evaluate(test_images, test_labels)

Epoch 1/5
1875/1875 [=====] - 11s 6ms/step - sparse_categorical_accuracy: 0.9734

```

```
Epoch 2/5
1875/1875 [=====] - 10s 5ms/step - sparse_categorical_accuracy: 0.9712
Epoch 3/5
1875/1875 [=====] - 10s 6ms/step - sparse_categorical_accuracy: 0.9683
Epoch 4/5
1875/1875 [=====] - 11s 6ms/step - sparse_categorical_accuracy: 0.9671
Epoch 5/5
1875/1875 [=====] - 11s 6ms/step - sparse_categorical_accuracy: 0.9674
313/313 [=====] - 0s 881us/step - sparse_categorical_accuracy: 0.9695
```

```
Out[37]: 0.9695000052452087
```

```
In [38]: _, acc = distiller.student.evaluate(
    test_images, test_labels, verbose=0)

print('Distilled Model test accuracy:', acc)

sz = get_gzipped_model_size(distiller.student)
print("Distilled model size: ", sz, ' KB')

model_acc.append(acc)
model_sz.append(sz)
```

```
Distilled Model test accuracy: 0.9695000052452087
Distilled model size: 231.089 KB
```

```
----- Checkpoint Point 3 -----
```

## Weight Clustering - 4.0

```
In [39]: def print_model_weight_clusters(model):
    for layer in model.layers:
        if isinstance(layer, keras.layers.Wrapper):
            weights = layer.trainable_weights
        else:
            weights = layer.weights
        for weight in weights:
            # ignore auxiliary quantization weights
            if "quantize_layer" in weight.name:
                continue
            if "kernel" in weight.name:
                unique_count = len(np.unique(weight))
                print(
                    f'{layer.name}/{weight.name}: {unique_count} clusters'
                )
```

```
In [40]: import tensorflow_model_optimization as tfmot
from tensorflow_model_optimization.python.core.clustering.keras.experimental import (
    cluster,
)

cluster_weights = tfmot.clustering.keras.cluster_weights
CentroidInitialization = tfmot.clustering.keras.CentroidInitialization

cluster_weights = cluster.cluster_weights

clustering_params = {
    'number_of_clusters': 8,
    'cluster_centroids_init': CentroidInitialization.KMEANS_PLUS_PLUS,
    'preserve_sparsity': True
}

sparsity_clustered_model = cluster_weights(distiller.student, **clustering_params)

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

print('Train sparsity preserving clustering model:')
sparsity_clustered_model.fit(train_images, train_labels, epochs=3, validation_split=0.1)
```

```
Train sparsity preserving clustering model:
Epoch 1/3
1688/1688 [=====] - 7s 4ms/step - loss: 0.0857 - accuracy: 0.9750 - val_loss: 0.0626 - val_accuracy: 0.9833
Epoch 2/3
1688/1688 [=====] - 6s 4ms/step - loss: 0.0625 - accuracy: 0.9805 - val_loss: 0.0705 - val_accuracy: 0.9773
Epoch 3/3
1688/1688 [=====] - 7s 4ms/step - loss: 0.0598 - accuracy: 0.9814 - val_loss: 0.0566 - val_accuracy: 0.9825
```

```
Out[40]: <tf_keras.src.callbacks.History at 0x1687c6990>
```

```
In [41]: stripped_clustered_model = tfmot.clustering.keras.strip_clustering(sparsity_clustered_model)

print("Model sparsity:\n")
print_model_weights_sparsity(stripped_clustered_model)

print("\nModel clusters:\n")
print_model_weight_clusters(stripped_clustered_model)

Model sparsity:

kernel:0: 32.41% sparsity (35/108)
kernel:0: 71.27% sparsity (14454/20280)

Model clusters:
```

```

conv2d_3/kernel:0: 8 clusters
dense_3/kernel:0: 8 clusters

In [42]: stripped_clustered_model.compile(optimizer=opt,
                                         loss=keras.losses.SparseCategoricalCrossentropy(from_logits=True),
                                         metrics=['accuracy'])

In [43]: _, stripped_clustered_model_accuracy = stripped_clustered_model.evaluate(
    test_images, test_labels, verbose=0)

print('Clustered Model test accuracy:', stripped_clustered_model_accuracy)

sz = get_gzipped_model_size(stripped_clustered_model)
print("Clustered model size: ", sz, ' KB')

model_acc.append(stripped_clustered_model_accuracy)
model_sz.append(sz)

Clustered Model test accuracy: 0.9782999753952026
Clustered model size: 165.816 KB

```

----- Checkpoint Point 4 -----

## Quantization - 5.0

```

In [44]: quant_aware_annotation_model = tfmot.quantization.keras.quantize_annotation_model(
    stripped_clustered_model)
quant_model = tfmot.quantization.keras.quantize_apply(
    quant_aware_annotation_model,
    tfmot.experimental.combine.Default8BitClusterPreserveQuantizeScheme(preserve_sparsity=True))

quant_model.compile(optimizer='adam',
                     loss=keras.losses.SparseCategoricalCrossentropy(from_logits=True),
                     metrics=['accuracy'])
print('Training after quantization model:')
quant_model.fit(train_images, train_labels, batch_size=128, epochs=3, validation_split=0.1)

Training after quantization model:
Epoch 1/3
WARNING:tensorflow:Gradients do not exist for variables ['conv2d_3/kernel:0', 'dense_3/kernel:0'] when minimizing t
he loss. If you're using `model.compile()`, did you forget to provide a `loss` argument?
WARNING:tensorflow:Gradients do not exist for variables ['conv2d_3/kernel:0', 'dense_3/kernel:0'] when minimizing t
he loss. If you're using `model.compile()`, did you forget to provide a `loss` argument?
WARNING:tensorflow:Gradients do not exist for variables ['conv2d_3/kernel:0', 'dense_3/kernel:0'] when minimizing t
he loss. If you're using `model.compile()`, did you forget to provide a `loss` argument?
WARNING:tensorflow:Gradients do not exist for variables ['conv2d_3/kernel:0', 'dense_3/kernel:0'] when minimizing t
he loss. If you're using `model.compile()`, did you forget to provide a `loss` argument?

422/422 [=====] - 4s 8ms/step - loss: 0.0474 - accuracy: 0.9856 - val_loss: 0.0547 - val_a
ccuracy: 0.9852
Epoch 2/3
422/422 [=====] - 3s 8ms/step - loss: 0.0441 - accuracy: 0.9866 - val_loss: 0.0539 - val_a
ccuracy: 0.9853
Epoch 3/3
422/422 [=====] - 4s 8ms/step - loss: 0.0415 - accuracy: 0.9872 - val_loss: 0.0530 - val_a
ccuracy: 0.9865

Out[44]: <tf_keras.src.callbacks.History at 0x16d163110>

In [45]: print("Final Model clusters:")
print_model_weight_clusters(quant_model)
print("\nFinal Model sparsity:")
print_model_weights_sparsity(quant_model)

Final Model clusters:
quant_conv2d_3/conv2d_3/kernel:0: 8 clusters
quant_dense_3/dense_3/kernel:0: 8 clusters

Final Model sparsity:
conv2d_3/kernel:0: 32.41% sparsity (35/108)
dense_3/kernel:0: 71.48% sparsity (14497/20280)

In [46]: converter = tf.lite.TFLiteConverter.from_keras_model(quant_model)
converter.optimizations = [tf.lite.Optimize.DEFAULT]
final_tflite_model = converter.convert()
final_model_file = 'final_model.tflite'
# Save the model.
with open(final_model_file, 'wb') as f:
    f.write(final_tflite_model)

sz = get_gzipped_model_size2(final_model_file)
print("Final model size: ", sz, ' KB')

model_sz.append(sz)

INFO:tensorflow:Assets written to: /var/folders/px/z8lb6znd6q95tq6vlznyb0s40000gn/T/tmpdptprpy2s/assets
INFO:tensorflow:Assets written to: /var/folders/px/z8lb6znd6q95tq6vlznyb0s40000gn/T/tmpdptprpy2s/assets
Final model size: 6.547 KB
/Users/ajaymaheshwari/anaconda3/lib/python3.11/site-packages/tensorflow/lite/python/convert.py:964: UserWarning: St
atistics for quantized inputs were expected, but not specified; continuing anyway.
    warnings.warn(
W0000 00:00:1717013837.495217 269976 tf_tfl_flatbuffer_helpers.cc:390] Ignored output_format.
W0000 00:00:1717013837.495757 269976 tf_tfl_flatbuffer_helpers.cc:393] Ignored drop_control_dependency.
2024-05-30 01:47:17.496489: I tensorflow/cc/saved_model/reader.cc:831 Reading SavedModel from: /var/folders/px/z8lb

```

```
6znd6q95tq6vlznyb0s40000gn/T/tmpdtprry2s
2024-05-30 01:47:17.497870: I tensorflow/cc/saved_model/reader.cc:51] Reading meta graph with tags { serve }
2024-05-30 01:47:17.497875: I tensorflow/cc/saved_model/reader.cc:146] Reading SavedModel debug info (if present) f
rom: /var/folders/px/z8lb6znd6q95tq6vlznyb0s40000gn/T/tmpdtprry2s
2024-05-30 01:47:17.508730: I tensorflow/cc/saved_model/loader.cc:234] Restoring SavedModel bundle.
2024-05-30 01:47:17.545231: I tensorflow/cc/saved_model/loader.cc:218] Running initialization op on SavedModel bund
le at path: /var/folders/px/z8lb6znd6q95tq6vlznyb0s40000gn/T/tmpdtprry2s
2024-05-30 01:47:17.554274: I tensorflow/cc/saved_model/loader.cc:317] SavedModel load for tags { serve }; Status:
success: OK. Took 57787 microseconds.
```

```
In [47]: interpreter = tf.lite.Interpreter(final_model_file)
interpreter.allocate_tensors()

final_test_accuracy = eval_model(interpreter)

print('Final test accuracy:', final_test_accuracy)

model_acc.append(final_test_accuracy)

Final test accuracy: 0.9803
```

----- Checkpoint Point 5 -----

```
In [48]: for i in range(len(model_acc)):
    print(f"Accuracy = {round(model_acc[i]*100,2)} with size = {model_sz[i]} KB ")

Accuracy = 98.06 with size = 235.396 KB
Accuracy = 97.31 with size = 204.383 KB
Accuracy = 96.95 with size = 231.089 KB
Accuracy = 97.83 with size = 165.816 KB
Accuracy = 98.03 with size = 6.547 KB
```

----- Final Comparison Summary -----