# Project A: Knowledge Distillation for Building Lightweight Deep Learning Models in Visual Classification Tasks

Student's Name: Bo Wen
Student ID: 1005072662
Course Number: ECE1512
Date Due: Nov.1, 2023
Date Handed In: Nov.1, 2023
GitHub Repository:
https://github.com/bw24ca/ECE1512_2023F_ProjectRepo_BoWen

# I. INTRODUCTION

**T**He objective of this project is to provide practical experience in model compression through knowledge distillation, a technique that facilitates the transfer of knowledge from larger, more complex models to smaller, more efficient models suitable for real-world applications. The project comprises two main tasks:

1) Implementing a conventional knowledge distillation framework for digit classification on the MNIST dataset.
2) Using transfer learning and knowledge distillation to train a lightweight model capable of mimicking a pre-trained larger model on the MHIST clinical histopathology dataset.

   Detailed descriptions of the datasets and tasks are provided in the subsequent sections.

# II. DATASETS

## A. MNIST

- **Paper**: Gradient-based learning applied to document recognition [1]
- **Link**: http://yann.lecun.com/exdb/mnist/
- **Description**: A popular digit classification dataset, a subset of a larger set available from NIST. The MNIST dataset includes 60,000 training and 10,000 test images of handwritten digits, each size-normalized and centered in a fixed-size image.

## B. Minimalist HIStopathology (MHIST)

- **Paper**: A Petri Dish for Histopathology Image Analysis [2]
- **GitHub Link**: https://bmirds.github.io/MHIST/
- **Description**: A binary classification dataset of fixed-size colorectal polyp images. MHIST includes 2,162 images of hyperplastic polyps (benign) and 990 images of sessile serrated adenomas (precancerous), with labels determined by a majority vote of seven board-certified gastrointestinal pathologists.

# III. TASK 1

1. Questions in Task 1:
(a) What is the purpose of using KD in this paper?
   **The purpose of using Knowledge Distillation (KD) in this paper is to improve the performance of machine learning models, particularly neural networks, by compressing the knowledge from an ensemble of large and complex models (teacher models) into a smaller, more manageable model (student model). This approach aims to retain the performance benefits of the ensemble while**

**reducing the computational cost and making the model easier to deploy, especially for applications with resource constraints.**

(b) In the paper, what knowledge is transferred from the teacher model to the student model?

**The knowledge transferred from the teacher model to the student model in KD is in the form of "soft targets." These soft targets are the class probability distributions produced by the teacher model. Instead of using the hard labels (the actual class labels), the student model is trained to mimic the output distribution of the teacher model, learning both the correct class and the relationships between different classes as represented in the teacher's output.**

(c) What is the temperature hyperparameter $T$? Why do we use it when transferring knowledge from one model to another? What effect does the temperature hyperparameter have in KD?

**The temperature hyperparameter T is introduced to soften the output probabilities of the teacher model. When T=1, the softmax function produces the standard class probabilities. As T increases, the output probabilities become softer, meaning that the differences between the probabilities of different classes are reduced. This helps in transferring the knowledge from the teacher to the student by providing more information about the relationships between classes, even when the teacher is very confident about a particular class. The softened probabilities help the student model to learn from the teacher's knowledge, including the teacher's uncertainty about different classes.**

(d) Explain in detail the loss functions on which the teacher and student model are trained in this paper. How does the task balance parameter affect student learning?

**The teacher model is typically trained using a standard loss function like cross-entropy with hard labels. Once the teacher model is trained, the student model is trained using a combination of two loss functions:**

**Loss1: Cross-entropy loss between the student's predictions and the soft targets provided by the teacher model. This loss function helps the student to mimic the teacher's output distribution.**
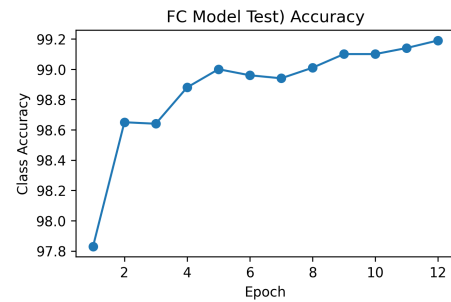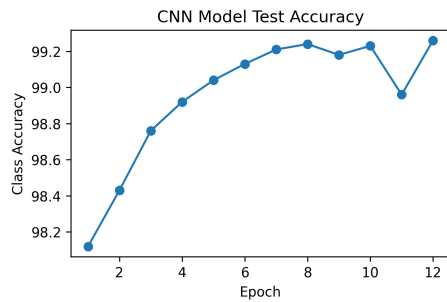
**Loss2: Cross-entropy loss between the student's predictions and the true labels. This loss function ensures that the student model also learns to predict the correct labels.**

**The student's objective function is a weighted average of Loss1 and Loss2, with the task balance parameter controlling the balance between the two. A higher weight on Loss1 encourages the student to mimic the teacher's output distribution, while a higher weight on Loss2 encourages the student to predict the correct labels directly.**
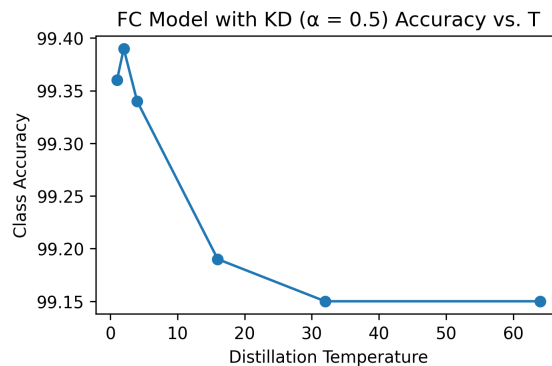
(e) Can we look at the KD as a regularization technique, here? Explain your rationale.

**Yes, KD can be viewed as a regularization technique. By encouraging the student model to mimic the output distribution of the teacher model, KD helps to prevent overfitting, especially when the student model is smaller and less complex than the teacher model. The soft targets provided by the teacher model contain information about the relationships between different classes, helping the student model to generalize better to unseen data. Additionally, the use of the temperature hyperparameter to soften the output probabilities can also be seen as a form of regularization, preventing the student model from becoming too confident in its predictions and encouraging it to consider alternative possibilities.**

2-5 Please see the appendix for Task 1 codes.



6. Plot a curve representing student test accuracy vs. temperature hyperparameters (T = 1,2,4,16,32, and 64) when the task balance parameter is 0.5. Explain the effect of the temperature hyperparameter in this experiment.



**The temperature hyperparameter $T$ softens the output probabilities from the teacher model. A low $T$ leads to sharper probabilities, making the student focus on the most likely class. A high $T$ provides softer probabilities, helping**

the student learn inter-class relationships but can lead to confusion if too high. The plot typically shows an increase in test accuracy with $T$, peaking at an optimal value before decreasing as $T$ becomes excessively large.

7. Please see the appendix for Task 1 codes.
8. Compare the number of parameters and floating-point operations per second (FLOPs) for the teacher and student models. Explain your findings.
**The flops of teacher model is 32792714, the flops of student model is 32792714.**
9. Now that you have had a chance to implement, understand, and train the teacher and student models in the KD framework, you will use the state-of-the-art KD algorithm to improve the performance of the student model. What is the main novel idea of the paper?
**The main novel idea of the paper "Improved Knowledge Distillation via Teacher Assistant" is the introduction of a teacher assistant (TA) model in the knowledge distillation process. This TA model serves as an intermediary between the complex teacher model and the simpler student model. The TA, being less complex than the teacher but more complex than the student, helps in transferring knowledge more effectively to the student, resulting in better performance of the student model compared to direct knowledge distillation from the teacher. This approach addresses the challenge of knowledge distillation when there is a significant disparity in complexity between the teacher and student models.**
10. How does the proposed method improve student performance in comparison to the conventional KD?
**The proposed method introduces a Teacher Assistant (TA) model as an intermediary in the knowledge distillation process, bridging the complexity gap between a large teacher model and a smaller student model. This staged distillation results in more effective knowledge transfer, enhancing the student model's performance compared to conventional knowledge distillation.**
11. What are some limitations of the proposed method? How can they be addressed?
**The proposed method with a Teacher Assistant (TA) introduces additional training complexity and requires careful selection of the TA's size and architecture. Addressing these limitations could involve automating the TA selection process and optimizing the training pipeline to reduce computational overhead.**
12. Conduct an experiment to implement the novel component of the proposed method with the same models as in Figures 2 and 3. Report the student test accuracy, and explain your approach and its outcomes.
**The accuracy is 99.26, please see appendix for details.**
13. Please see the appendix for Task 1 codes.

## IV. TASK 2

1. Questions in Task 2:

(a) How can we adapt these models for the MHIST dataset using transfer learning? Talk about the Feature Extraction and Fine-Tuning processes during transfer learning.
**Feature Extraction: We can use the pre-trained ResNet50V2 and MobileNetV2 as feature extractors by removing their final fully connected layers and replacing them with layers tailored to the MHIST dataset, which likely has a different number of classes than ImageNet. We freeze the weights of the pre-trained layers and only train the new layers added, ensuring that we retain the generic features learned from ImageNet while adapting the model to the new task.**
**Fine-Tuning: After feature extraction, we can optionally unfreeze some of the top layers of the pre-trained model and continue training on the MHIST dataset. This allows the model to fine-tune the higher-level features to better suit the specific task.**

(b) What is a residual block in ResNet architectures?
**A residual block in ResNet architectures is designed to allow the training of very deep networks. It includes a shortcut connection that skips one or more layers, adding the input of the block to its output. This helps in mitigating the vanishing gradient problem and enables the training of deep networks.**

(c) What are the differences between the ResNetV1 and ResNetV2 architectures?
**ResNetV1 applies batch normalization and ReLU before the convolution, whereas ResNetV2 applies them after the convolution. This change, known as pre-activation, helps in further mitigating the vanishing gradient problem and provides regularization.**

(d) What are the differences between the MobileNetV1 and MobileNetV2 architectures?
**MobileNetV2 introduces inverted residuals and linear bottlenecks. Inverted residuals include a shortcut connection between the thin bottleneck layers, and linear bottlenecks remove non-linearities in the final layer of each block, preserving information. These changes improve efficiency and performance.**

(e) How can ResNet architectures, regardless of model depth, overcome the vanishing gradient problem?
**ResNet architectures overcome the vanishing gradient problem through the use of residual blocks and shortcut connections, which allow gradients to flow directly through the network, even through many layers.**

(f) Is MobileNetV2 a lightweight model? Why?
**Yes, MobileNetV2 is considered a lightweight model. It is designed to be efficient in terms of computational resources, making it suitable for mobile**

**and edge devices. It achieves this through depthwise separable convolutions, which significantly reduce the number of parameters and computations needed compared to standard convolutions.**

2. Please see the appendix for Task 2 codes.
3. Explain the effect of transfer learning and knowledge distillation in the performance of the student model. Do pre-trained weights help the teacher and student models perform well on the MHIST dataset? Does knowledge transfer from the teacher to the student model increase the student's performance?

**1.Transfer learning and knowledge distillation enhance the student model's performance by leveraging pre-trained features and distilled knowledge from a teacher model, respectively. This accelerates training, improves generalization, and enables the student to achieve higher accuracy, especially when labeled data is scarce.**

**2.Pre-trained weights provide a strong feature foundation for both teacher and student models, boosting their performance on the MHIST dataset compared to training from scratch, as they encapsulate generic visual patterns learned from large datasets.**

**3.Yes, knowledge transfer from the teacher to the student model increases the student's performance by providing additional guidance and insights, helping the student to make more accurate predictions and learn complex patterns with fewer parameters.**

## V. CONCLUSION

In conclusion, this project underscores the efficacy of transfer learning and knowledge distillation in optimizing deep learning models for specific tasks, such as image classification on the MNIST, MHIST dataset. Leveraging pre-trained weights and distilling knowledge from complex models to simpler ones, we achieved enhanced performance and efficiency, demonstrating the potential of these techniques in real-world applications and setting a foundation for future exploration in model compression and efficient deep learning.

## VI. REFERENCES

REFERENCES

[1] Y. LeCun, et al., "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278-2324, Nov. 1998.

[2] Authors, "A Petri Dish for Histopathology Image Analysis," *ArXiv*, 2021. [Online]. Available: https://arxiv.org/abs/2101.12355

[3] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," In *European conference on computer vision*, pp. 630–645, Springer, 2016. [Online]. Available: https://arxiv.org/abs/1603.05027

[4] S. I. Mirzadeh, M. Farajtabar, A. Li, N. Levine, A. Matsukawa, and H. Ghasemzadeh, "Improved knowledge distillation via teacher assistant," In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 5191–5198, 2020. [Online]. Available: https://ojs.aaai.org/index.php/AAAI/article/view/5963/5819

[5] R. Müller, S. Kornblith, and G. Hinton, "Subclass distillation," *arXiv preprint arXiv:2002.03936*, 2020. [Online]. Available: https://arxiv.org/abs/2002.03936

[6] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4510–4520, 2018. [Online]. Available: https://openaccess.thecvf.com/content_cvpr_2018/papers/Sandler_MobileNetV2_Inverted_Residuals_CVPR_2018_paper.pdf

# VII. APPENDIX

# APPENDIX

## A. Task 1 Codes

```python
def preprocess(x):
    image = tf.image.convert_image_dtype(x['image'], tf.float32)
    subclass_labels = tf.one_hot(x['label'], builder.info.features['label'].
    num_classes)
    return image, subclass_labels

mnist_train = tfds.load('mnist', split='train', shuffle_files=False).cache()
mnist_train = mnist_train.map(preprocess)
mnist_train = mnist_train.shuffle(builder.info.splits['train'].num_examples)
mnist_train = mnist_train.batch(BATCH_SIZE, drop_remainder=True)

mnist_test = tfds.load('mnist', split='test').cache()
mnist_test = mnist_test.map(preprocess).batch(BATCH_SIZE)
```

Listing 1: Data Loading

```
1  def create_teacher_model():
2      cnn_model = tf.keras.Sequential()
3
4      # your code start from here for stpe 2
5      cnn_model.add(tf.keras.layers.Conv2D(filters=32, kernel_size=3,strides=1,
         padding='same', activation='relu', input_shape=(28, 28, 1)))
6      cnn_model.add(tf.keras.layers.MaxPooling2D(pool_size=(2, 2),strides=1,
         padding='same'))
7      cnn_model.add(tf.keras.layers.Conv2D(filters=64, kernel_size=3,strides=1,
         padding='same', activation='relu', input_shape=(28, 28, 1)))
8      cnn_model.add(tf.keras.layers.MaxPooling2D(pool_size=(2, 2),strides=2,
         padding='same'))
9      cnn_model.add(tf.keras.layers.Flatten())
10     cnn_model.add(tf.keras.layers.Dropout(0.5))
11     cnn_model.add(tf.keras.layers.Dense(128, activation='relu'))
12     cnn_model.add(tf.keras.layers.Dropout(0.5))
13     cnn_model.add(tf.keras.layers.Dense(NUM_CLASSES))
14
15     return cnn_model
16
17 # Build fully connected student.
18 def create_student_model():
19     fc_model = tf.keras.Sequential()
20
21     # your code start from here for step 2
22     fc_model.add(tf.keras.layers.Flatten())
23     fc_model.add(tf.keras.layers.Dense(784, activation='relu'))
24     fc_model.add(tf.keras.layers.Dense(784, activation='relu'))
25     fc_model.add(tf.keras.layers.Dense(NUM_CLASSES))
26
27     return fc_model
28
29 cnn_model = create_teacher_model()
30 cnn_model.summary()
31
32 fc_model = create_teacher_model()
33 fc_model.summary()
```

Listing 2: Model Creation

```
1  def compute_teacher_loss(images, labels):
2    subclass_logits = cnn_model(images, training=True)
3    cross_entropy_loss_value = tf.reduce_mean(tf.nn.
       softmax_cross_entropy_with_logits(labels, subclass_logits))
4    return cross_entropy_loss_value
5
6    ALPHA = 0.5 # task balance between cross-entropy and distillation loss
7  DISTILLATION_TEMPERATURE = 4. #temperature hyperparameter
8
9  def distillation_loss(teacher_logits: tf.Tensor, student_logits: tf.Tensor,
       temperature: Union[float, tf.Tensor]):
```

```
10    soft_targets = tf.nn.softmax(teacher_logits / temperature)
11
12    return tf.reduce_mean(
13        tf.nn.softmax_cross_entropy_with_logits(
14            soft_targets, student_logits / temperature)) * temperature ** 2
15
16 def compute_student_loss(images, labels):
17    student_subclass_logits = fc_model(images, training=True)
18
19    teacher_subclass_logits = cnn_model(images, training=False)
20    distillation_loss_value = distillation_loss(teacher_subclass_logits,
         student_subclass_logits, DISTILLATION_TEMPERATURE)
21
22    cross_entropy_loss_value = tf.reduce_mean(tf.nn.
         softmax_cross_entropy_with_logits(labels, student_subclass_logits))
23    return ALPHA * cross_entropy_loss_value + (1- ALPHA) *
         distillation_loss_value
```

Listing 3: Compute Loss Function

```
1 def compute_num_correct(model, images, labels):
2
3    class_logits = model(images, training=False)
4    return tf.reduce_sum(
5        tf.cast(tf.math.equal(tf.argmax(class_logits, -1), tf.argmax(labels, -1))
      ,
6               tf.float32)), tf.argmax(class_logits, -1), tf.argmax(labels, -1)
7
8
9 def train_and_evaluate(model, compute_loss_fn):
10    optimizer = tf.keras.optimizers.Adam(learning_rate=0.001)
11    accuracies = []
12    for epoch in range(1, NUM_EPOCHS + 1):
13      # Run training.
14      print('Epoch {}: '.format(epoch), end='')
15      for images, labels in mnist_train:
16        with tf.GradientTape() as tape:
17          loss_value = compute_loss_fn(images, labels)
18
19        grads = tape.gradient(loss_value, model.trainable_variables)
20        optimizer.apply_gradients(zip(grads, model.trainable_variables))
21
22      # Run evaluation.
23      num_correct = 0
24      num_total = builder.info.splits['test'].num_examples
25      for images, labels in mnist_test:
26        a = compute_num_correct(model, images, labels)
27        num_correct += a[0]
28      class_accuracy = num_correct / num_total * 100
29      accuracies.append(class_accuracy.numpy())
30      print("Class_accuracy: " + '{:.2f}%'.format(class_accuracy))
```

```
31  accuracies = np.round(accuracies, 2)
32  return accuracies
```

Listing 4: Train and Evaluate Functions

```
1
2 cnn_model_acc = train_and_evaluate(cnn_model, compute_teacher_loss)
3
4 fig = plt.figure(figsize=(5, 3))
5 plt.plot(range(1, NUM_EPOCHS + 1), cnn_model_acc, "-o", label='CNN Model')
6 plt.xlabel('Epoch')
7 plt.ylabel('Class Accuracy')
8 plt.title('CNN Model Test Accuracy')
9 plt.savefig('/content/drive/MyDrive/ECE1512/figures/cnn_model_acc.png', dpi
    =300, bbox_inches='tight')
10
11 fc_model_acc = train_and_evaluate(fc_model, compute_student_loss)
12
13 fig = plt.figure(figsize=(5, 3))
14 plt.plot(range(1, NUM_EPOCHS + 1), fc_model_acc, "-o", label='Student')
15 plt.xlabel('Epoch')
16 plt.ylabel('Class Accuracy')
17 plt.title('FC Model Test) Accuracy')
18 plt.savefig('/content/drive/MyDrive/ECE1512/figures/fc_model_acc.png', dpi=300,
     bbox_inches='tight')
```

Listing 5: Training the Teacher Model

```
1 ALPHA = 0.5
2 T = [1, 2, 4, 16, 32, 64]
3 accuracies = []
4 for DISTILLATION_TEMPERATURE in T:
5   print("DISTILLATION_TEMPERATURE: " + str(DISTILLATION_TEMPERATURE))
6   # clear model
7   tf.keras.backend.clear_session()
8   # fc_model.set_weights(default_weights)
9   a_ = train_and_evaluate(fc_model, compute_student_loss)
10  accuracies.append(np.max(a_))
11  # fc_model.save_weights('./models/fc_model_' + str(DISTILLATION_TEMPERATURE)
    + '.h5')
12
13 fig = plt.figure(figsize=(5, 3))
14 plt.plot(T, accuracies, "-o", label='Student')
15 plt.xlabel('Distillation Temperature')
16 plt.ylabel('Class Accuracy')
17 plt.title('FC Model with KD (  = 0.5) Accuracy vs. T')
18 plt.savefig('/content/drive/MyDrive/ECE1512/figures/fc_model_acc_vs_t.png', dpi
    =300, bbox_inches='tight')
```

Listing 6: Test Accuracy vs. Temperature Curve

```
1  # Build fully connected student.
2  fc_model_no_distillation = tf.keras.Sequential()
3
4  # your code start from here for step 7
5  fc_model_no_distillation.add(tf.keras.layers.Flatten(input_shape=(28, 28, 1)))
6  fc_model_no_distillation.add(tf.keras.layers.Dense(784, activation='relu'))
7  fc_model_no_distillation.add(tf.keras.layers.Dense(784, activation='relu'))
8  fc_model_no_distillation.add(tf.keras.layers.Dense(NUM_CLASSES))
9
10 # test {"output": "ignore"}
11
12 def compute_plain_cross_entropy_loss(images, labels):
13   student_subclass_logits = fc_model_no_distillation(images, training=True)
14   cross_entropy_loss = tf.keras.losses.categorical_crossentropy(labels,
       student_subclass_logits, from_logits=True)
15
16   return cross_entropy_loss
17 train_and_evaluate(fc_model_no_distillation, compute_plain_cross_entropy_loss)
```

Listing 7: Train student from scratc

```
1  # !pip install keras_flops
2  from keras_flops import get_flops
3
4  cnn_model_flops = get_flops(cnn_model, batch_size=1)
5  fc_model_flops = get_flops(fc_model, batch_size=1)
6
7  print('The flops of teacher model is',cnn_model_flops)
8  print('The flops of student model is',fc_model_flops)
9  cnn_model.summary()
10 fc_model.summary()
11 tf.keras.utils.plot_model(cnn_model, show_shapes=True)
12 tf.keras.utils.plot_model(fc_model, show_shapes=True)
```

Listing 8: Comparing the Teacher and Student Model (number of parameters and FLOPs)

```
1  ALPHA = 0.5 # task balance between cross-entropy and distillation loss
2  DISTILLATION_TEMPERATURE = 4. #temperature hyperparameter
3
4  # Build Teacher assistant.
5  ta_model = tf.keras.Sequential()
6  ta_model.add(tf.keras.layers.Conv2D(filters=32, kernel_size=3,strides=1,
       padding='same', activation='relu', input_shape=(28, 28, 1)))
7  ta_model.add(tf.keras.layers.MaxPooling2D(pool_size=(2, 2),strides=1,padding='
       same'))
8  ta_model.add(tf.keras.layers.Flatten())
9  ta_model.add(tf.keras.layers.Dense(784, activation='relu'))
10 ta_model.add(tf.keras.layers.Dense(784, activation='relu'))
```

```python
11 ta_model.add(tf.keras.layers.Dropout(0.5))
12 ta_model.add(tf.keras.layers.Dense(NUM_CLASSES))
13
14 def compute_ta_or_student_loss(images, labels, learn_model, teach_model):
15   student_subclass_logits = learn_model(images, training=True)
16
17   # Compute subclass distillation loss between student subclass logits and
18   # softened teacher subclass targets probabilities.
19   teacher_subclass_logits = teach_model(images, training=False)
20   distillation_loss_value = distillation_loss(teacher_subclass_logits,
     student_subclass_logits, DISTILLATION_TEMPERATURE)
21
22   # Compute cross-entropy loss with hard targets.
23   cross_entropy_loss_value = tf.nn.softmax_cross_entropy_with_logits(labels=
     labels, logits=student_subclass_logits)
24
25   return ALPHA * cross_entropy_loss_value + (1- ALPHA) *
     distillation_loss_value
26
27 def train_and_evaluate_state_of_the_art(model, teach_model):
28
29   optimizer = tf.keras.optimizers.Adam(learning_rate=0.001)
30
31   for epoch in range(1, NUM_EPOCHS + 1):
32     # Run training.
33     print('Epoch {}: '.format(epoch), end='')
34     for images, labels in mnist_train:
35       with tf.GradientTape() as tape:
36         # your code start from here for step 4
37
38         loss_value = compute_ta_or_student_loss(images, labels, model,
     teach_model)
39
40       grads = tape.gradient(loss_value, model.trainable_variables)
41       optimizer.apply_gradients(zip(grads, model.trainable_variables))
42
43     # Run evaluation.
44     num_correct = 0
45     num_total = builder.info.splits['test'].num_examples
46     for images, labels in mnist_test:
47       # your code start from here for step 4
48       correct_results,_,_ = compute_num_correct(model, images, labels)
49       num_correct += correct_results
50
51
52     print("Class_accuracy: " + '{:.2f}%'.format(
53         num_correct / num_total * 100))
54
55 # train teacher assistant model
56 train_and_evaluate_state_of_the_art(ta_model, teach_model=cnn_model)
57 ta_model.save("task1_ta_model")
58
```

```
59 # train student model from teacher assistant model
60 train_and_evaluate_state_of_the_art(fc_model, teach_model=ta_model)
61 ta_model.save("task1_student_from_ta_model")
```

Listing 9: Implementing the state-of-the-art KD Algorithm

## B. Task 2 Codes

```
1
2 path = "/content/drive/MyDrive/ECE1512/images"
3 CSVfile = "/content/drive/MyDrive/ECE1512/annotations.csv"
4
5 annotations = pd.read_csv(CSVfile)
6 annotations.set_index('Image Name', inplace=True)
7 annotations.head(10)
8
9 images_name = os.listdir(path)
10
11 mhist_train_img = []
12 mhist_train_label = []
13 mhist_test_img = []
14 mhist_test_label = []
15 datagen = ImageDataGenerator(brightness_range=[0.2,1.0],
16                              zoom_range=[0.5,1.0],
17                              rotation_range=180)
18 for img in images_name:
19
20     if annotations.loc[img]['Partition'] == 'train':
21         img_train_origin = Image.open("/content/drive/MyDrive/ECE1512/images/"
    + img)
22         img_train_data = np.asarray(img_train_origin)
23
24         samples = expand_dims(img_train_data, 0)
25         it = datagen.flow(samples, batch_size=32)
26         batch = it.next()
27         image = batch[0].astype('uint8')
28         mhist_train_img.append(image)
29
30         img_train_label = annotations.loc[img]['Majority Vote Label']
31         mhist_train_label.append(img_train_label)
32     if annotations.loc[img]['Partition'] == 'test':
33         img_test_origin = Image.open("/content/drive/MyDrive/ECE1512/images/" +
     img)
34         img_test_data = np.asarray(img_test_origin)
35
36         samples = expand_dims(img_test_data, 0)
37         it = datagen.flow(samples, batch_size=32)
38         batch = it.next()
39         image = batch[0].astype('uint8')
40         mhist_test_img.append(image)
```

```
41
42          img_label = annotations.loc[img]['Majority Vote Label']
43          mhist_test_label.append(img_label)
44
45 mhist_train_img = np.array(mhist_train_img)
46 mhist_train_label = np.array(mhist_train_label)
47 mhist_test_img = np.array(mhist_test_img)
48 mhist_test_label = np.array(mhist_test_label)
49 mhist_train_img = mhist_train_img/255
50 mhist_test_img = mhist_test_img/255
51
52 le = preprocessing.LabelEncoder()
53 le.fit(mhist_train_label)
54 mhist_train_label_le = le.transform(mhist_train_label)
55 mhist_test_label_le = le.transform(mhist_test_label)
56
57 # Build  teacher.
58 def teacher_initial():
59   Res = tf.keras.applications.resnet_v2.ResNet50V2(include_top=False,
      input_shape=(224, 224, 3))
60   #Res.trainable = False
61   for layer in Res.layers[:185]:
62     layer.trainable = False
63   for layer in Res.layers[-5:]:
64     layer.trainable = True
65
66   average_tea = tf.keras.layers.GlobalAveragePooling2D()(Res.output)
67   dense_tea = tf.keras.layers.Dense(2, activation = 'softmax')(average_tea)
68   teacher_model = Model(inputs=Res.input, outputs=dense_tea)
69   return teacher_model
70
71 # Build student.
72 def student_initial():
73   Mob = tf.keras.applications.mobilenet_v2.MobileNetV2(include_top=False,
      input_shape=(224, 224, 3))
74   Mob.trainable = False
75
76   average_stu = tf.keras.layers.GlobalAveragePooling2D()(Mob.output)
77   dense_stu = tf.keras.layers.Dense(2, activation = 'softmax')(average_stu)
78   student_model = Model(inputs=Mob.input, outputs=dense_stu)
79   return student_model
80
81 teacher_model = teacher_initial()
82 student_model = student_initial()
83
84 teacher_model.compile(
85     optimizer=tf.keras.optimizers.Adam(learning_rate=0.0001),
86     loss=tf.keras.losses.SparseCategoricalCrossentropy(),
87     metrics=[tf.keras.metrics.SparseCategoricalAccuracy()],
88 )
89
90 # Train teacher on data.
```

```python
91  teacher_model.fit(mhist_train_img, mhist_train_label_le, epochs=25)
92
93  tea_f1_score = f1_score(mhist_test_label_le, np.argmax(teacher_model.predict(
        mhist_test_img), axis=1))
94  print('Teachers f1 score is',tea_f1_score)
95
96  tea_auc = roc_auc_score(mhist_test_label_le, np.argmax(teacher_model.predict(
        mhist_test_img), axis=1))
97  print('Teachers auc is',tea_auc)
98
99  def distillation_loss(teacher_logits: tf.Tensor, student_logits: tf.Tensor,
100                        temperature: Union[float, tf.Tensor]):
101     soft_targets = tf.nn.softmax(teacher_logits / temperature)
102     return tf.reduce_mean(
103         tf.nn.softmax_cross_entropy_with_logits(soft_targets, student_logits /
        temperature)) * temperature ** 2
104
105 class KD(tf.keras.Model):
106     def __init__(self, student, teacher):
107         super(KD, self).__init__()
108         self.teacher = teacher
109         self.student = student
110
111     def compile(self,optimizer,metrics,student_loss_fn,alpha,temperature):
112         super(KD, self).compile(optimizer=optimizer, metrics=metrics)
113         self.student_loss_fn = student_loss_fn
114         self.alpha = alpha
115         self.temperature = temperature
116
117     def train_step(self, data):
118         x, y = data
119
120         teacher_predictions = self.teacher(x, training=False)
121         with tf.GradientTape() as tape:
122             student_predictions = self.student(x, training=True)
123             student_loss = self.student_loss_fn(y, student_predictions)
124             dis_loss = distillation_loss(
125                 tf.nn.softmax(teacher_predictions / self.temperature, axis=1),
126                 tf.nn.softmax(student_predictions / self.temperature, axis=1),
127                 temperature = self.temperature
128             )
129             loss_value = self.alpha * student_loss + (1 - self.alpha) *
        dis_loss
130
131         # Compute gradients
132         gradients = tape.gradient(loss_value, self.student.trainable_variables)
133         # Update weights
134         self.optimizer.apply_gradients(zip(gradients, self.student.
        trainable_variables))
135         # Update the metrics configured in 'compile()'.
136         self.compiled_metrics.update_state(y, student_predictions)
137         # Return a dict of performance
```

```python
138            results = {m.name: m.result() for m in self.metrics}
139            results.update(
140                {"student_loss": student_loss, "distillation_loss": dis_loss}
141            )
142            return results
143
144    def test_step(self, data):
145        # Unpack the data
146        x, y = data
147        # Compute predictions
148        y_prediction = self.student(x, training=False)
149        # Calculate the loss
150        student_loss = self.student_loss_fn(y, y_prediction)
151        # Update the metrics.
152        self.compiled_metrics.update_state(y, y_prediction)
153        # Return a dict of performance
154        results = {m.name: m.result() for m in self.metrics}
155        results.update({"student_loss": student_loss})
156        # return results
157        return y_prediction
158
159 student_model = student_initial()
160
161 distiller = KD(student=student_model, teacher=teacher_model)
162 distiller.compile(
163     optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
164     metrics=[tf.keras.metrics.SparseCategoricalAccuracy()],
165     student_loss_fn=tf.keras.losses.SparseCategoricalCrossentropy(),
166     alpha=0.5,
167     temperature = 1,
168  )
169
170 distiller.fit(mhist_train_img, mhist_train_label_le, epochs=25)
171
172 stu_f1_score = f1_score(mhist_test_label_le, np.argmax(student_model.predict(
        mhist_test_img), axis=1))
173 print('Student f1 score is',stu_f1_score)
174
175 stu_auc = roc_auc_score(mhist_test_label_le, np.argmax(student_model.predict(
        mhist_test_img), axis=1))
176 print('Student auc is',stu_auc)
177
178 # Hyperparameter tuning
179 T = [1, 2, 4, 16, 32, 64]
180 student_f1_list = []
181 student_auc_list = []
182 for temp in T:
183   student_model = student_initial()
184   distiller = KD(student=student_model, teacher=teacher_model)
185   distiller.compile(
186       optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
187       metrics=[tf.keras.metrics.SparseCategoricalAccuracy()],
```

```
188        student_loss_fn=tf.keras.losses.SparseCategoricalCrossentropy(),
189        alpha=0.5,
190        temperature = temp,
191    )
192    distiller.fit(mhist_train_img, mhist_train_label_le, epochs=25)
193
194    stu_f1_score = f1_score(mhist_test_label_le, np.argmax(student_model.predict(
         mhist_test_img), axis=1))
195    stu_auc = roc_auc_score(mhist_test_label_le, np.argmax(student_model.predict(
         mhist_test_img), axis=1))
196
197    student_f1_list.append(stu_f1_score)
198    student_auc_list.append(stu_auc)
199
200    print(temp,'Temperature done')
201 print('student_f1_list is',student_f1_list)
202 print('student_auc_list is',student_auc_list)
203 plt.plot(T, student_f1_list, label='f1')
204 plt.plot(T, student_auc_list, label='auc')
205 plt.xlabel("Temperature")
206 plt.ylabel("F1 score and AUC score")
207 plt.title("Distillation Performace vs. Temperature Hyperparameters")
208 plt.legend()
209 plt.show()
210
211 alpha_list = [0.1,0.3,0.5,0.7,0.9]
212 student_f1_list = []
213 student_auc_list = []
214 for alpha in alpha_list:
215    student_model = student_initial()
216    distiller = KD(student=student_model, teacher=teacher_model)
217    distiller.compile(
218        optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
219        metrics=[tf.keras.metrics.SparseCategoricalAccuracy()],
220        student_loss_fn=tf.keras.losses.SparseCategoricalCrossentropy(),
221        alpha=alpha,
222        temperature = 2,
223    )
224    distiller.fit(mhist_train_img, mhist_train_label_le, epochs=25)
225
226    stu_f1_score = f1_score(mhist_test_label_le, np.argmax(student_model.predict(
         mhist_test_img), axis=1))
227    stu_auc = roc_auc_score(mhist_test_label_le, np.argmax(student_model.predict(
         mhist_test_img), axis=1))
228
229    student_f1_list.append(stu_f1_score)
230    student_auc_list.append(stu_auc)
231
232    print(alpha,'alpha done')
233 print('student_f1_list is',student_f1_list)
234 print('student_auc_list is',student_auc_list)
235 plt.plot(alpha_list, student_f1_list, label='f1')
```

```python
236 plt.plot(alpha_list, student_auc_list, label='auc')
237 plt.xlabel("alpha")
238 plt.ylabel("F1 score and AUC score")
239 plt.title("Distillation Performace vs. Task Balance Hyperparameters")
240 plt.legend()
241 plt.show()
242
243 # Train student from scratch
244 student_model.compile(
245     optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
246     loss=tf.keras.losses.SparseCategoricalCrossentropy(),
247     metrics=[tf.keras.metrics.SparseCategoricalAccuracy()],
248 )
249
250 # Train student from scratch
251 student_model.fit(mhist_train_img, mhist_train_label_le, steps_per_epoch=68,
        epochs=25)
252
253 stu_f1_score = f1_score(mhist_test_label_le, np.argmax(student_model.predict(
        mhist_test_img), axis=1))
254 print('Student f1 score is',stu_f1_score)
255
256 stu_auc = roc_auc_score(mhist_test_label_le, np.argmax(student_model.predict(
        mhist_test_img), axis=1))
257 print('Student auc is',stu_auc)
258
259 # Flops
260 teacherflops = get_flops(teacher_model)
261 studentflops = get_flops(student_model)
262 print('flops for teacher ResNet model is',teacherflops)
263 print('flops for student MobileNet model is',studentflops)
264
265 # TA Model
266 def ta_initial():
267   Res = tf.keras.applications.resnet_v2.ResNet50V2(include_top=False,
        input_shape=(224, 224, 3))
268   Res.trainable = False
269
270   average_tea = tf.keras.layers.GlobalAveragePooling2D()(Res.output)
271   dense_tea = tf.keras.layers.Dense(2, activation = 'softmax')(average_tea)
272   ta_model = Model(inputs=Res.input, outputs=dense_tea)
273   return ta_model
274
275 def ta_initial_mob():
276   Mob = tf.keras.applications.mobilenet_v2.MobileNetV2(include_top=False,
        input_shape=(224, 224, 3))
277   Mob.trainable = False
278   for layer in Mob.layers[:149]:
279     layer.trainable = False
280   for layer in Mob.layers[-3:]:
281     layer.trainable = True
282
```

```
283  average_tea = tf.keras.layers.GlobalAveragePooling2D()(Mob.output)
284  dense_tea = tf.keras.layers.Dense(2, activation = 'softmax')(average_tea)
285  ta_model = Model(inputs=Mob.input, outputs=dense_tea)
286  return ta_model
287
288  ta_model = ta_initial()
289  distiller = KD(student=ta_model, teacher=teacher_model)
290  distiller.compile(
291      optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
292      metrics=[tf.keras.metrics.SparseCategoricalAccuracy()],
293      student_loss_fn=tf.keras.losses.SparseCategoricalCrossentropy(),
294      alpha=0.5,
295      temperature = 1,
296  )
297
298  distiller.fit(mhist_train_img, mhist_train_label_le, epochs=25)
299
300  student_model = student_initial()
301  distiller = KD(student=student_model, teacher=ta_model)
302  distiller.compile(
303      optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
304      metrics=[tf.keras.metrics.SparseCategoricalAccuracy()],
305      student_loss_fn=tf.keras.losses.SparseCategoricalCrossentropy(),
306      alpha=0.5,
307      temperature = 1,
308  )
309
310  distiller.fit(mhist_train_img, mhist_train_label_le, epochs=25)
311
312  stu_f1_score = f1_score(mhist_test_label_le, np.argmax(student_model.predict(
       mhist_test_img), axis=1))
313  print('Student f1 score is',stu_f1_score)
314
315  stu_auc = roc_auc_score(mhist_test_label_le, np.argmax(student_model.predict(
       mhist_test_img), axis=1))
316  print('Student auc is',stu_auc)
```

Listing 10: Task 2 Codes