

# COS30082 Applied Machine Learning

## Assignment - Image Classification

Student Name: Tran Hung Quoc Tuan

Student ID: 105000908

### 1. Methodology

This section details the process of data preparation, the architectural design of the models, and the strategies implemented to mitigate overfitting.

#### 1.1 Data Preparation

The provided dataset consisted of a single train folder containing 10,000 images organized into 10 distinct classes. As per the assignment requirement to create our own evaluation sets, the split-folders library was employed to partition this data into three separate, non-overlapping directories:

- **Training Set (70%):** Used for training the models.
- **Validation Set (15%):** Used for hyperparameter tuning and to monitor EarlyStopping.
- **Testing Set (15%):** A completely unseen set reserved for the final model evaluation.

```
import tensorflow as tf
from tensorflow.keras.preprocessing.image import ImageDataGenerator
from tensorflow.keras.callbacks import EarlyStopping
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix
import numpy as np
import os
import splitfolders

#1. Define Kaggle Paths
input_folder = '/kaggle/input/image-classification-train-zip/train'
output_folder = '/kaggle/working/data_split/'

#2. Split Data (Only needs to run once per session)
print("Splitting source data into train, val, and test sets...")
try:
    # This reads from /kaggle/input and writes to /kaggle/working/
    splitfolders.ratio(input_folder,
                      output=output_folder,
                      seed=42,
                      ratio=(.7, .15, .15)) # 70% train, 15% val, 15% test
    print(f'Data successfully split into {output_folder}')
except Exception as e:
    print(f'Could not split folders (maybe already split?): {e}')
```

Two separate ImageDataGenerator instances were created. For the **Training Set**, data augmentation (including random rotation, width/height shifts, zoom, and horizontal flips)

was applied to increase data variance and combat overfitting. For the **Validation** and **Testing Sets**, only pixel normalization (rescaling to 1./255) was performed to ensure a consistent and realistic evaluation. All images were resized to (128, 128) pixels.

```
#3. Define Image Parameters
IMG_SIZE = (128, 128)
BATCH_SIZE = 32
seed = 42

#4. Create SEPARATE Generators
# Generator for TRAINING data (WITH augmentation)
train_datagen = ImageDataGenerator(
    rescale=1./255,
    rotation_range=20,
    width_shift_range=0.1,
    height_shift_range=0.1,
    zoom_range=0.1,
    horizontal_flip=True
)

# Generator for VALIDATION and TEST data (RESCALE ONLY)
val_test_datagen = ImageDataGenerator(rescale=1./255)

# 5. Point Generators to the NEW Folders in /kaggle/working/
train_gen = train_datagen.flow_from_directory(
    os.path.join(output_folder, 'train'),
    target_size=IMG_SIZE,
    batch_size=BATCH_SIZE,
    class_mode='categorical',
    shuffle=True,
    seed=seed
)

val_gen = val_test_datagen.flow_from_directory(
    os.path.join(output_folder, 'val'),
    target_size=IMG_SIZE,
    batch_size=BATCH_SIZE,
    class_mode='categorical',
    shuffle=False, # No shuffle for validation
    seed=seed
)

# This is your UNSEEN test set for final evaluation
test_gen = val_test_datagen.flow_from_directory(
    os.path.join(output_folder, 'test'),
    target_size=IMG_SIZE,
    batch_size=BATCH_SIZE,
    class_mode='categorical',
    shuffle=False,
    seed=seed
)
```

```
Splitting source data into train, val, and test sets...
Copying files: 10008 files [00:26, 371.06 files/s]
Data successfully split into '/kaggle/working/data_split/'
Found 7005 images belonging to 10 classes.
Found 1500 images belonging to 10 classes.
Found 1503 images belonging to 10 classes.

Generators created.
Classes detected: {'Amphibia': 0, 'Animalia': 1, 'Arachnida': 2, 'Aves': 3, 'Fungi': 4, 'Insecta': 5, 'Mammalia': 6, 'Mollusca': 7, 'Plantae': 8, 'Reptilia': 9}
```

## 1.2 Model Architectures

To satisfy the assignment's requirement for model comparison, two distinct architectures were implemented:

1. **Model A: Baseline CNN (Sequential Model)** This model was a simple Convolutional Neural Network built from scratch. The architecture consisted of three Conv2D layers with increasing filter sizes (32, 64, 128), each followed by a MaxPooling2D layer. The feature maps were then flattened and passed through a Dense layer (128 units) with Dropout(0.5) for regularization, and finally to a 10-unit softmax output layer for classification.

```
#2. Baseline CNN Model (Run 1)
print("\nStarting Run 1: Baseline CNN Model")

model_a = tf.keras.Sequential([
    tf.keras.layers.Conv2D(32, (3,3), activation='relu', input_shape=(128,128,3)),
    tf.keras.layers.MaxPooling2D(2,2),
    tf.keras.layers.Conv2D(64, (3,3), activation='relu'),
    tf.keras.layers.MaxPooling2D(2,2),
    tf.keras.layers.Conv2D(128, (3,3), activation='relu'),
    tf.keras.layers.MaxPooling2D(2,2),
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(128, activation='relu'),
    tf.keras.layers.Dropout(0.5), # Anti-overfitting
    tf.keras.layers.Dense(10, activation='softmax') # 10 classes
])

model_a.compile(optimizer=tf.keras.optimizers.Adam(1e-3),
                loss='categorical_crossentropy', # Correct loss for multi-class
                metrics=['accuracy']) # 'accuracy' is Top-1 Accuracy

model_a.summary()

# This callback stops training when validation loss stops improving
callback = EarlyStopping(
    monitor='val_loss',
    patience=3,
    restore_best_weights=True,
    verbose=1
)

history_a = model_a.fit(train_gen,
                       validation_data=val_gen,
                       epochs=15,
                       callbacks=[callback])
```

```
Epoch 1/15
219/219 ----- 164s 727ms/step - accuracy: 0.1260 - loss: 2.3128 - val_accuracy: 0.1713 - val_loss: 2.2607
Epoch 2/15
219/219 ----- 155s 705ms/step - accuracy: 0.1550 - loss: 2.2437 - val_accuracy: 0.2173 - val_loss: 2.1532
Epoch 3/15
219/219 ----- 157s 716ms/step - accuracy: 0.2134 - loss: 2.1689 - val_accuracy: 0.2567 - val_loss: 2.0827
Epoch 4/15
219/219 ----- 168s 767ms/step - accuracy: 0.2377 - loss: 2.1158 - val_accuracy: 0.2853 - val_loss: 2.0340
Epoch 5/15
219/219 ----- 155s 704ms/step - accuracy: 0.2658 - loss: 2.0583 - val_accuracy: 0.2813 - val_loss: 2.0321
Epoch 6/15
219/219 ----- 158s 719ms/step - accuracy: 0.2838 - loss: 2.0315 - val_accuracy: 0.3013 - val_loss: 1.9732
Epoch 7/15
219/219 ----- 156s 713ms/step - accuracy: 0.2937 - loss: 1.9763 - val_accuracy: 0.2680 - val_loss: 2.0304
Epoch 8/15
219/219 ----- 169s 769ms/step - accuracy: 0.2930 - loss: 1.9987 - val_accuracy: 0.3007 - val_loss: 1.9701
Epoch 9/15
219/219 ----- 155s 707ms/step - accuracy: 0.2902 - loss: 1.9794 - val_accuracy: 0.3160 - val_loss: 1.9431
Epoch 10/15
219/219 ----- 207s 731ms/step - accuracy: 0.2939 - loss: 1.9699 - val_accuracy: 0.3400 - val_loss: 1.9176
Epoch 11/15
219/219 ----- 154s 703ms/step - accuracy: 0.3138 - loss: 1.9593 - val_accuracy: 0.3420 - val_loss: 1.9110
Epoch 12/15
219/219 ----- 171s 779ms/step - accuracy: 0.3103 - loss: 1.9334 - val_accuracy: 0.3140 - val_loss: 1.9484
Epoch 13/15
219/219 ----- 154s 701ms/step - accuracy: 0.3204 - loss: 1.9207 - val_accuracy: 0.3480 - val_loss: 1.8927
Epoch 14/15
219/219 ----- 154s 700ms/step - accuracy: 0.3307 - loss: 1.9102 - val_accuracy: 0.3407 - val_loss: 1.8981
Epoch 15/15
219/219 ----- 160s 728ms/step - accuracy: 0.3202 - loss: 1.9024 - val_accuracy: 0.3267 - val_loss: 1.9120
Restoring model weights from the end of the best epoch: 13.
```

2. **Model B: Transfer Learning (MobileNetV2)** This model leveraged transfer learning, using the MobileNetV2 architecture pre-trained on the ImageNet dataset. The process was conducted in two phases:

**Phase 1 (Feature Extraction):** The base MobileNetV2 model was frozen (base.trainable = False). A new classifier head—consisting of GlobalAveragePooling2D, Dropout(0.3), and a 10-unit softmax layer—was trained on top.

```
#3.Transfer Learning Model (Run 2)
print("\nStarting Run 2: MobileNetV2 (Feature Extraction)")

# Load the base model without its classifier
base = tf.keras.applications.MobileNetV2(
    include_top=False,
    weights='imagenet',
    input_shape=(128,128,3)
)

# Freeze the base layers
base.trainable = False

model_b = tf.keras.Sequential([
    base,
    tf.keras.layers.GlobalAveragePooling2D(),
    tf.keras.layers.Dropout(0.3), # Anti-overfitting
    tf.keras.layers.Dense(10, activation='softmax')
])

model_b.compile(optimizer=tf.keras.optimizers.Adam(1e-4), # Slower LR
    loss='categorical_crossentropy',
    metrics=['accuracy'])

model_b.summary()

history_b = model_b.fit(train_gen,
    validation_data=val_gen,
    epochs=10,
    callbacks=[callback])
```

940c464/940c464 0s 0us/step  
Model: "sequential\_1"

Layer (type)	Output Shape	Param #
mobilenetv2_1.00_128 (Functional)	(None, 4, 4, 1280)	2,257,984
global_average_pooling2d (GlobalAveragePooling2D)	(None, 1280)	0
dropout_1 (Dropout)	(None, 1280)	0
dense_2 (Dense)	(None, 10)	12,810

Total params: 2,270,794 (8.66 MB)  
Trainable params: 12,810 (50.04 KB)  
Non-trainable params: 2,257,984 (8.61 MB)

Epoch 1/10  
219/219 134s 578ms/step - accuracy: 0.1176 - loss: 2.8567 - val\_accuracy: 0.3020 - val\_loss: 2.0119  
Epoch 2/10  
219/219 120s 549ms/step - accuracy: 0.2742 - loss: 2.2084 - val\_accuracy: 0.4473 - val\_loss: 1.6771  
Epoch 3/10  
219/219 142s 547ms/step - accuracy: 0.3994 - loss: 1.8248 - val\_accuracy: 0.5220 - val\_loss: 1.5011  
Epoch 4/10  
219/219 120s 547ms/step - accuracy: 0.4245 - loss: 1.7532 - val\_accuracy: 0.5493 - val\_loss: 1.3996  
Epoch 5/10  
219/219 120s 547ms/step - accuracy: 0.4679 - loss: 1.6507 - val\_accuracy: 0.5660 - val\_loss: 1.3348  
Epoch 6/10  
219/219 120s 548ms/step - accuracy: 0.4888 - loss: 1.5384 - val\_accuracy: 0.5773 - val\_loss: 1.2904  
Epoch 7/10  
219/219 120s 549ms/step - accuracy: 0.5169 - loss: 1.4530 - val\_accuracy: 0.5820 - val\_loss: 1.2609  
Epoch 8/10  
219/219 120s 546ms/step - accuracy: 0.5172 - loss: 1.4399 - val\_accuracy: 0.6000 - val\_loss: 1.2288  
Epoch 9/10  
219/219 125s 572ms/step - accuracy: 0.5367 - loss: 1.3903 - val\_accuracy: 0.6040 - val\_loss: 1.2079  
Epoch 10/10  
219/219 120s 548ms/step - accuracy: 0.5354 - loss: 1.3882 - val\_accuracy: 0.6107 - val\_loss: 1.1913  
Restoring model weights from the end of the best epoch: 10.

**Phase 2 (Fine-Tuning):** The base model was unfrozen, and the top 30 layers were made trainable. The entire model was then re-compiled with a very low

learning rate (1e-5) to fine-tune the pre-trained weights specifically for our dataset.

```
#Transfer Learning Model (Run 3 - Fine-Tuned)

print("\nStarting Run 3: MobileNetV2 (Fine-Tuned Top Layers)")

# Unfreeze the base model
base.trainable = True

# Freeze all layers *except* the top 30
for layer in base.layers[:-30]:
    layer.trainable = False

# Re-compile the model with a very low learning rate
model_b.compile(optimizer=tf.keras.optimizers.Adam(1e-5),
                 loss='categorical_crossentropy',
                 metrics=['accuracy'])

model_b.summary() # Note the change in trainable params

# Continue training (fine-tuning)
history_finetune = model_b.fit(train_gen,
                               validation_data=val_gen,
                               epochs=10,
                               callbacks=[callback])
```

Starting Run 3: MobileNetV2 (Fine-Tuned Top Layers)  
Model: "sequential\_1"

Layer (type)	Output Shape	Param #
mobilenetv2_1.00_128 (Functional)	(None, 4, 4, 1280)	2,257,904
global_average_pooling2d (GlobalAveragePooling2D)	(None, 1280)	0
dropout_1 (Dropout)	(None, 1280)	0
dense_2 (Dense)	(None, 10)	12,810

Total params: 2,270,794 (8.66 MB)

Trainable params: 1,539,210 (5.87 MB)

Non-trainable params: 731,584 (2.79 MB)

Epoch 1/10  
219/219 — 152s 616ms/step - accuracy: 0.4718 - loss: 1.5972 - val\_accuracy: 0.6093 - val\_loss: 1.1724  
Epoch 2/10  
219/219 — 131s 597ms/step - accuracy: 0.5143 - loss: 1.4599 - val\_accuracy: 0.6180 - val\_loss: 1.1629  
Epoch 3/10  
219/219 — 136s 619ms/step - accuracy: 0.5458 - loss: 1.3339 - val\_accuracy: 0.6213 - val\_loss: 1.1477  
Epoch 4/10  
219/219 — 136s 619ms/step - accuracy: 0.5709 - loss: 1.2710 - val\_accuracy: 0.6233 - val\_loss: 1.1377  
Epoch 5/10  
219/219 — 132s 603ms/step - accuracy: 0.5676 - loss: 1.2749 - val\_accuracy: 0.6367 - val\_loss: 1.1205  
Epoch 6/10  
219/219 — 131s 599ms/step - accuracy: 0.5972 - loss: 1.2155 - val\_accuracy: 0.6400 - val\_loss: 1.1105  
Epoch 7/10  
219/219 — 131s 599ms/step - accuracy: 0.6071 - loss: 1.1678 - val\_accuracy: 0.6433 - val\_loss: 1.1025  
Epoch 8/10  
219/219 — 133s 605ms/step - accuracy: 0.6302 - loss: 1.1346 - val\_accuracy: 0.6453 - val\_loss: 1.0953  
Epoch 9/10  
219/219 — 131s 598ms/step - accuracy: 0.6301 - loss: 1.0998 - val\_accuracy: 0.6513 - val\_loss: 1.0929  
Epoch 10/10  
219/219 — 131s 598ms/step - accuracy: 0.6367 - loss: 1.0764 - val\_accuracy: 0.6547 - val\_loss: 1.0866  
Restoring model weights from the end of the best epoch: 10.

## 1.3 Overfitting Strategies

As the assignment explicitly warned of overfitting risks, the following strategies were employed:

- **Data Augmentation:** Applied to the training set to create new, varied image samples and prevent the model from memorizing the training data.
- **Dropout:** Implemented in both Model A (0.5) and Model B (0.3) to randomly deactivate neurons during training, forcing the network to learn more robust features.
- **Early Stopping:** Monitored the val\_loss with a patience=3. This callback automatically stopped the training process when the model's performance on the validation set ceased to improve, and restored the weights from the best-performing epoch.

## 2. Results and Discussion

This section presents the final evaluation of the models on the unseen test set, followed by a comparative analysis of their performance.

```
#4. Evaluation Section
print("\nEvaluating Models on UNSEEN Test Set...")

def evaluate_model(model, name, generator):
    print(f"\n===== {name} =====")

    # 1. Get Top-1 Accuracy
    # .evaluate() returns [loss, accuracy]
    loss, top1_acc = model.evaluate(generator)
    print(f"Overall Top-1 Accuracy: {top1_acc:.4f}")

    # 2. Calculate Average Accuracy per Class

    # Get true labels from the generator
    y_true = generator.classes

    # Get predicted probabilities
    y_pred_probs = model.predict(generator)

    # Convert probabilities to class labels
    y_pred = np.argmax(y_pred_probs, axis=-1)

    # Generate the confusion matrix
    conf_mat = confusion_matrix(y_true, y_pred)

    # Calculate per-class accuracy
    # (Diagonal elements) / (Sum of elements in that row)
    class_acc = conf_mat.diagonal() / conf_mat.sum(axis=-1)

    # Calculate the average of per-class accuracies
    avg_acc = np.mean(class_acc)

    print(f"Average Accuracy per Class: {avg_acc:.4f}")

    # Print per-class accuracies to see which ones fail
    print("\nPer-Class Accuracy:")
    class_names = list(generator.class_indices.keys())
    for i, acc in enumerate(class_acc):
        print(f"({class_names[i]:<10}): {acc:.4f}")

    return top1_acc, avg_acc

# Evaluate both models on the TEST_GEN
acc_a, avg_a = evaluate_model(model_a, "Run 1: Baseline CNN", test_gen)
acc_b, avg_b = evaluate_model(model_b, "Run 3: MobileNetV2 Fine-Tuned", test_gen)
```

## 2.1 Evaluation Metrics

The models were evaluated using the two required metrics: **Top-1 Accuracy** and **Average Accuracy per Class**. Top-1 Accuracy measures the standard classification accuracy, while Average Accuracy per Class provides a crucial measure of model fairness and balance by averaging the accuracy of each individual class.

The final performance of both models on the test set is presented below:

```
Evaluating Models on UNSEEN Test Set...

===== Run 1: Baseline CNN =====
47/47 ----- 22s 473ms/step - accuracy: 0.2659 - loss: 1.9694
Overall Top-1 Accuracy: 0.3407
47/47 ----- 15s 314ms/step
Average Accuracy per Class: 0.3405

Per-Class Accuracy:
Amphibia : 0.1333
Animalia : 0.1333
Arachnida : 0.4803
Aves : 0.4533
Fungi : 0.3510
Insecta : 0.2800
Mammalia : 0.3733
Mollusca : 0.2067
Plantae : 0.5533
Reptilia : 0.4400

===== Run 3: MobileNetV2 Fine-Tuned =====
47/47 ----- 16s 342ms/step - accuracy: 0.6191 - loss: 1.1933
Overall Top-1 Accuracy: 0.6494
47/47 ----- 20s 387ms/step
Average Accuracy per Class: 0.6494

Per-Class Accuracy:
Amphibia : 0.4933
Animalia : 0.6533
Arachnida : 0.5921
Aves : 0.8067
Fungi : 0.6755
Insecta : 0.7200
Mammalia : 0.7600
Mollusca : 0.6200
Plantae : 0.5733
Reptilia : 0.6000
```

## 2.2 Model Performance and Discussion

The results demonstrate a stark contrast in performance. Model B (MobileNetV2 Fine-Tuned) significantly outperformed the Baseline CNN, achieving nearly double the accuracy in both metrics.

**Analysis of Model A (Baseline CNN):** The performance of Model A was extremely poor (34.07%). Analysis of the training logs revealed that the model suffered from severe underfitting; both training and validation accuracy remained low. This indicates that the simple CNN architecture was not complex enough to capture the features of the dataset, especially when combined with the difficulty added by strong data augmentation. A look at its per-class accuracy confirms this: the model completely failed to learn some classes, scoring only 13% on 'Amphibia' and 'Animalia'.

**Analysis of Model B (MobileNetV2 Fine-Tuned):** Model B's performance of 64.94% was a substantial improvement. This demonstrates the power of transfer learning; by

leveraging the rich features learned from ImageNet, the model could effectively classify the new dataset.

The most compelling result is the perfect match between its Top-1 Accuracy (64.94%) and its Average Accuracy per Class (64.94%). This indicates that Model B is exceptionally well-balanced. Unlike Model A, it learned to classify all 10 classes with a relatively consistent (and much higher) level of performance. For example, its worst-performing class, 'Amphibia' (49.33%), was still significantly better than Model A's best-performing class.

Conclusion: The assignment's prompt correctly identified overfitting as a key challenge. However, our experiments showed that for a simple CNN, the combination of strong anti-overfitting techniques (like augmentation) and a complex dataset resulted in underfitting. This reinforced the necessity of using a more powerful, pre-trained architecture. The MobileNetV2 Fine-Tuned model (Model B) was proven to be the superior solution, yielding significantly higher and more balanced accuracy, thus fulfilling the performance objectives of the assignment