

CSC345/M45 Big Data and Machine Learning

Coursework: Object Recognition

Author: Linda Mafunu
Student number: 2216686
Date: 13/12/2023

Table of Contents

Introduction.....	2
Method.....	2
Task1: Download the <i>CIFAR100</i> dataset.	2
Task2: Computing Features and Visualizing Images	2
Task3: Learning Algorithms	2
Results.	3
Task 4: Evaluate the Models	3
Conclusion.	4
Task5: Benchmark and Discussion	4
References.	5
Appendices.	6

Introduction.

Classification is a type of supervised learning whereby the data is labelled with predefined classes. The goal was to use Machine Learning (ML) methodologies to perform image multi-classification. The CIFAR100 dataset was used to classify in the challenge. The dimensions of each image used are (32,32,3), which correspond to the height, width, and colour channel (RGB). The dataset was used to train and assess the performance of the models. According to Ray (2015), using multiple machine learning algorithms is a good way to improve ML accuracy because some algorithms are better suited for certain datasets than others. The classification methods used for this assignment are Support Vector Classifier (SVC), Neural Network (NN), and Convolutional Neural Network (CNN). CNN According to the experimental results shown below, CNN is a better classifier than the other two approaches.

Method.

Task1: Download the *CIFAR100* dataset.

To accomplish this, imported the tensor flow package, which is utilised in model development. The downloaded dataset is divided into four NumPy array variables, which is why the NumPy package was imported to store our datasets. The dataset consists of 50000 training datasets and 10,000 testing datasets. Another NumPy array is declared to store the image target class names.

Task2: Computing Features and Visualizing Images

Slice data:

The training dataset was sliced from 50000 to 500, and the testing dataset was sliced from 10000 to 100. Because the hog functions accept image input, as (32,32,3) the data was transposed to be shown as (samples, height, width, color_channels) rather than (height, width, color_channels, samples)

Feature Extraction

The extracted features are a simplified representation of the image that incorporates only the most significant information about the image. Imported the skimage.feature package to use the *hog function*. As stated by Singh (2019), the *hog function* computes the histogram of oriented gradients. Had to loop through the entire dataset and extract features for each image when using the Hog function, which took a long time for 50000 dataset it took 30 minutes. Instead of using the Hog function, decided to flatten and then normalise the dataset, as implemented in lab4. Before normalising the pixels by dividing them by 255, the training and testing datasets were reshaped to (50000,3072) and (10000,3072), respectively. This was done to convert the dataset from 4D to 2D and create a single long continuous linear vector that could be fed into neural network models. Ray (2015) claims that this is one of the methods used to improve model performance. However, it did not seem to be useful for SVC and Linear SVC since the dataset size was too enormous, requiring a long time to fit the dataset into the object, despite LSVC's documentation promising support for larger datasets.

Task3: Learning Algorithms

Build ML Models

First, defined the ML. model by selecting the type of model needed, followed by determining the network topology, in this case a sequential model. By doing so, the model's layers are configured, each with numerous nodes and activation functions, and the levels are integrated into a cohesive model. The second stage is model compilation, which involves running a method to compile the model with the supplied configuration, which will build the necessary data structures for efficient model usage. To archive this, select an optimization loss function (Sparse Categorical Cross entropy or binary cross entropy). This involves the selection of an

optimisation strategy (Adam) as well as metrics (accuracy) to monitor during the model training process. Following that, the model is evaluated using the validation dataset, which is 20% of the training dataset. This is done to ensure that the model fulfils its intended purpose and to assess the loss and metrics. Finally, predict the test picture dataset class labels from our dataset's 20 target classes to determine if the model is classifying images correctly.

Neural Network

A fully connected neural network is composed of many densely connected layers that connect every neuron in one layer to every neuron in the next using input from the hog function or flattened data and the ReLU action, which aids in learning non-linear correlations between components and makes the network more resilient to detecting diverse patterns.

Create SVC and Linear SVC

To modify image characteristics to be on a same scale, used the Standard Scaler class to standardise the dataset. The dataset was then modified to match the destination system (SVC). Created a SCV object and fitted the training dataset and its labels into it. SVC is typically used to work with small complex datasets. Finally, the object was evaluated once the classes of the testing dataset were predicted.

Build CNN.

CNN is a deep learning method that assumes the inputs are images and allows certain features to be encoded into the model architecture for image classification. Each layer of a CNN uses a differentiable function (ReLU) to transform one volume of activations to another. CNNs are effective because they can do automatic feature extraction at a large scale. SoftMax is used at the output to interpret raw classifier scores as probabilities.

Results.

Task 4: Evaluate the Models

Evaluation measure performance of independent blind test data. This was done using the confusion matrix, accuracy score, and classification report. In machine learning according to Bharathi (2021), a confusion matrix is a table used to assess the effectiveness of a classification model. Imported pandas and seaborn packages to display the results in the Appendices section as a heat map. In multilabel classification, the accuracy score computes the subset accuracy. The precision, recall, and f1-score for each target class are displayed in the classification report. Precision indicates how well the model predicts a certain class, recall indicates how many times the model detected a specific class, and the f1-score weighted average of precision and recall.

Refer to Appendices for screenshots of results.

	SVC (500 dataset)	NN	CN N1	CNN2	CNN3	2 nd CNN1	2 nd CNN2	2 nd CNN3
Accuracy Score %	12	10.34	5	19.09	25.18	44.98	49.2	53.56
Time lapse for training for 50000 image datasets	> 1hour	9.45s	1m 52.7s	2m 2.9s	4m 54.7s	1m 51s	3m 45.2s	5m 56s
Classification Bias towards target name	Medium sized mammals	Fruit and vegetables	Large outdoor scenes	Insects, large carnivores	Insects	People	Household furniture	Large outdoor scenes

Lowest Classification towards target name	Most target classes except people, large man-made outdoor things, and household electrical devices	fish, flowers, household furniture, large carnivores, non-insect invertebrates, people, reptiles, small mammals	Every target class	vehicles2	Small mammals	Small mammals	reptiles	non-insect invertebrates
Correct Classification	-	-	-	Some target classes	some target classes	Most target classes	Most targets classes	Most targets

Improvement Method

Image Data Generator, according to PyImageSerach (2019), can be used to improve model accuracy since it allows for image augmentation to reduce the possibility of overfitting data by loading the image dataset into memory and generating batches of augmented data. This method allows us to artificially increase the dataset's size and diversity, which improves model generalisation and the CNN's ability to learn more robust features.

Conclusion.

Task5: Benchmark and Discussion

There is enough evidence to conclude that SVC is not an appropriate classification methodology for large datasets, which is why the data was sliced. However, using smaller datasets is not a good model training methodology because the model may have biased outliers due to noticing patterns that do not exist, resulting in overfitting. More training data, according to Ray (2015), allows the "data to tell for itself" rather than relying on assumptions and weak correlations. According to the SVC documentation, linear SVC is supposed to work for large datasets, but fitting the 50000-image dataset took more than an hour, which is inefficient.

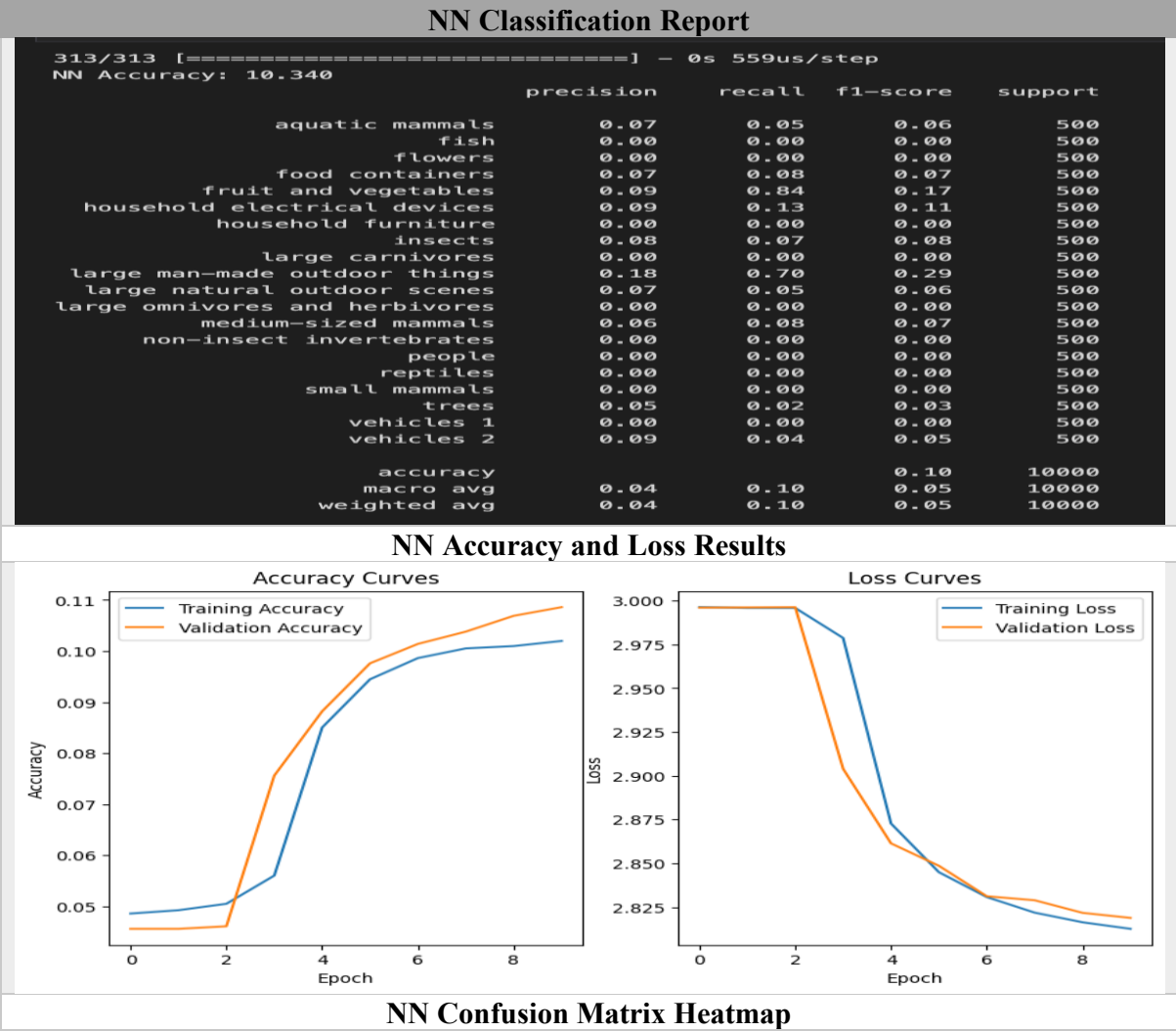
To try to solve the problem raised by SVC objects, NN was used to classify since it can handle large datasets. Even though using simple NN reduced model training time and allowed us to use the entire dataset, the model's accuracy score was very low. One significant reason for this poor performance is that NN does not automatically extract features like CNN when training models, resulting in large losses when training model. The image features had to be extracted manually first using the hog function, but the function was taking too long, so the method from lab4 was used. However, none of the images were correctly classified in the end, indicating that the model was not performing its intended function.

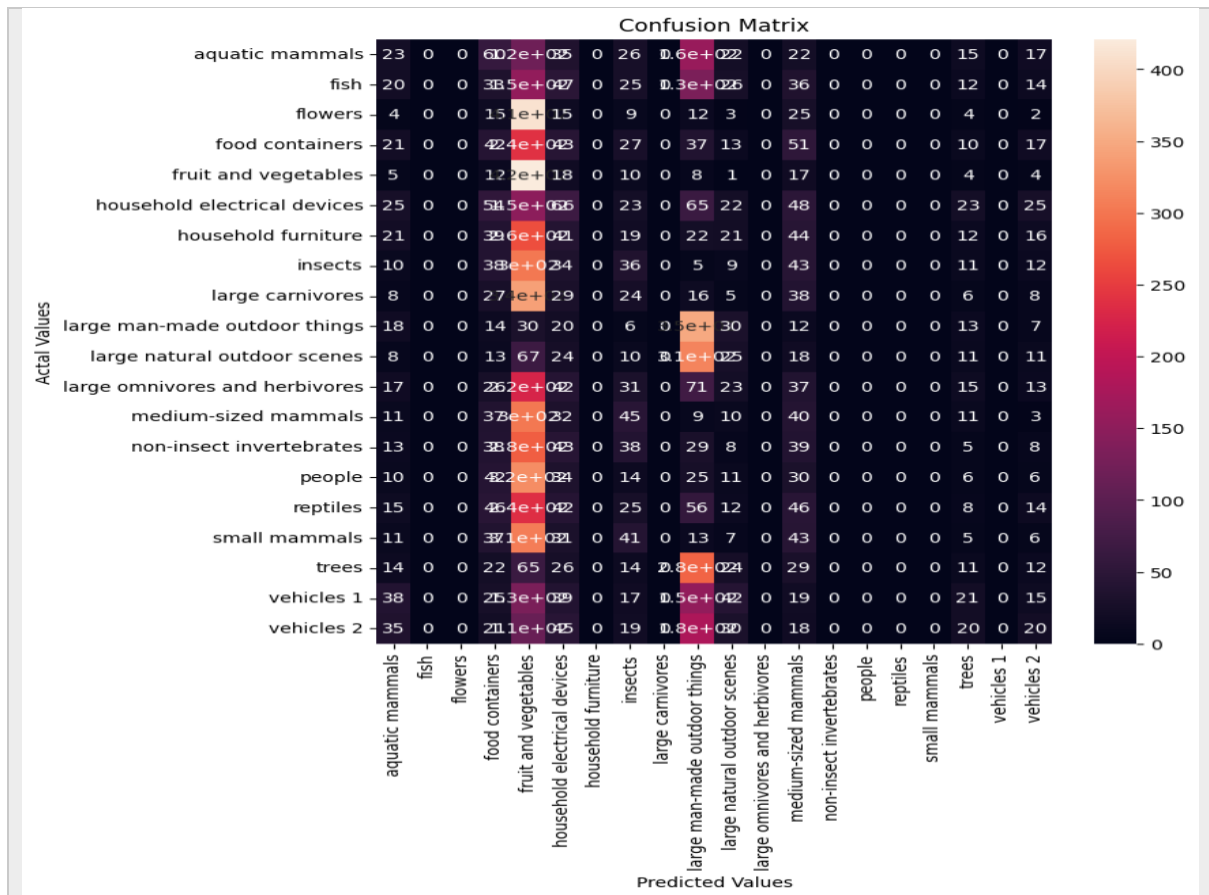
Another classification method, deep learning with CNN, was used to improve the classification. Deep learning automatically learns features from data, eliminating the need for manual extraction of image features. The results show that the more layers you add to the CNN, the better the model performs. The results show that more images are being classified correctly, something that the SVC and NN models did not do. Used the Image Data Processed methodology proposed by Ray (2015) and GeeksforGeeks (2023) to further improve our model performance. Image augmentation is one of the methods to be used to improve model accuracy. The model's performance increased from 25% to 53%, and the training and validation accuracy increased after 2 epochs for the 2nd CNN in the 1st CNN it increased after 6 epochs. In conclusion, CNN is the most effective ML methodology for multi-classification of images for large datasets; however, there is a need to investigate Transfer learning methodology to see how it improves model precision in the future.

References.

- Ray, S. (2015, December 28). *8 Proven Ways for boosting the “Accuracy” of a Machine Learning Model*. Analytics Vidhya. <https://www.analyticsvidhya.com/blog/2015/12/improve-machine-learning-results/>
- `sklearn.svm.LinearSVC`. (n.d.). Scikit-Learn. <https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html#sklearn.svm.LinearSVC>
- Bharathi. (2021, June 24). *Confusion Matrix for Multi-Class Classification*. Analytics Vidhya. <https://www.analyticsvidhya.com/blog/2021/06/confusion-matrix-for-multi-class-classification/>
- Keras ImageDataGenerator and Data Augmentation*. (2019, July 8). PyImageSearch. <https://pyimagesearch.com/2019/07/08/keras-imagedatagenerator-and-data-augmentation/>
- `keras.fit()` and `keras.fit_generator()`. (2019, June 12). GeeksforGeeks. https://www.geeksforgeeks.org/keras-fit-and-keras-fit_generator/
- Multiclass image classification using Transfer learning*. (2021, October 15). GeeksforGeeks. <https://www.geeksforgeeks.org/multiclass-image-classification-using-transfer-learning/>
- Singh, A. (2019, September 4). *Feature Descriptor / Hog Descriptor Tutorial*. Analytics Vidhya. <https://www.analyticsvidhya.com/blog/2019/09/feature-engineering-images-introduction-hog-feature-descriptor/>
- `sklearn.svm.LinearSVC`. (n.d.). Scikit-Learn. <https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html#sklearn.svm.LinearSVC>
- `tf.keras.preprocessing.image.ImageDataGenerator`. (n.d.). TensorFlow. https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/image/ImageDataGenerator
- `tf.keras.utils.to_categorical` | TensorFlow Core v2.5.0. (n.d.). TensorFlow. https://www.tensorflow.org/api_docs/python/tf/keras/utils/to_categorical
- `keras.fit()` and `keras.fit_generator()`. (2019, June 12). GeeksforGeeks. https://www.geeksforgeeks.org/keras-fit-and-keras-fit_generator/

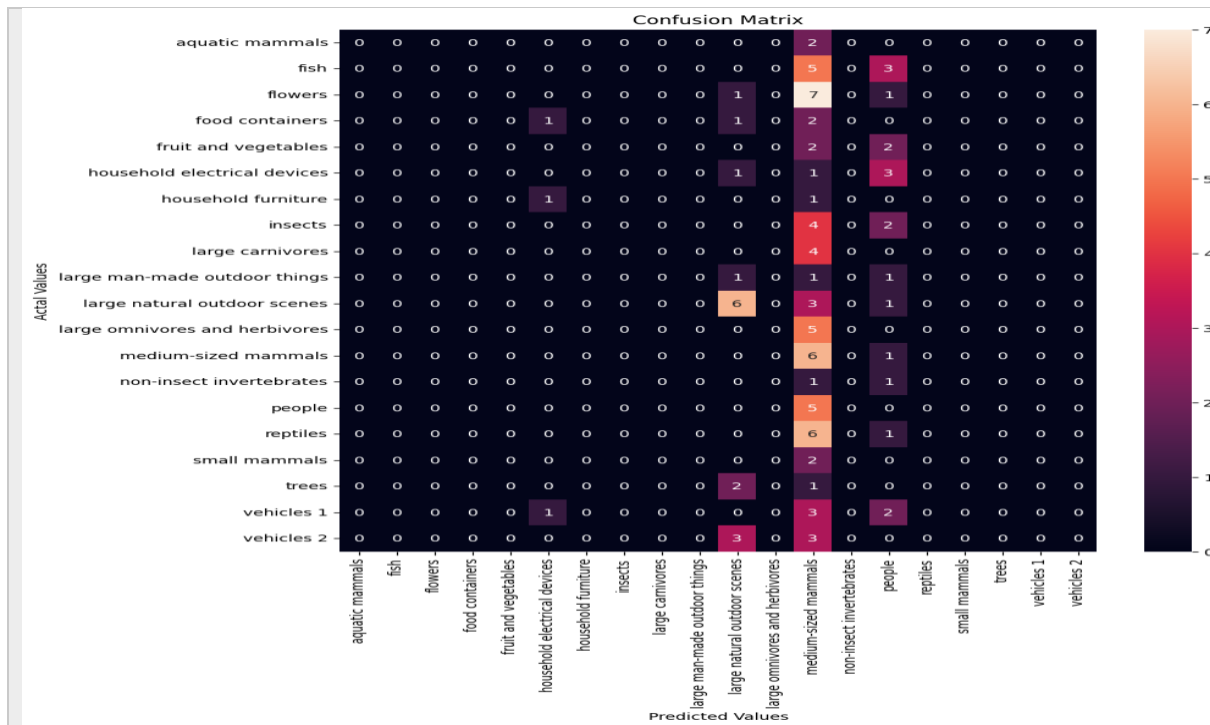
Appendices.



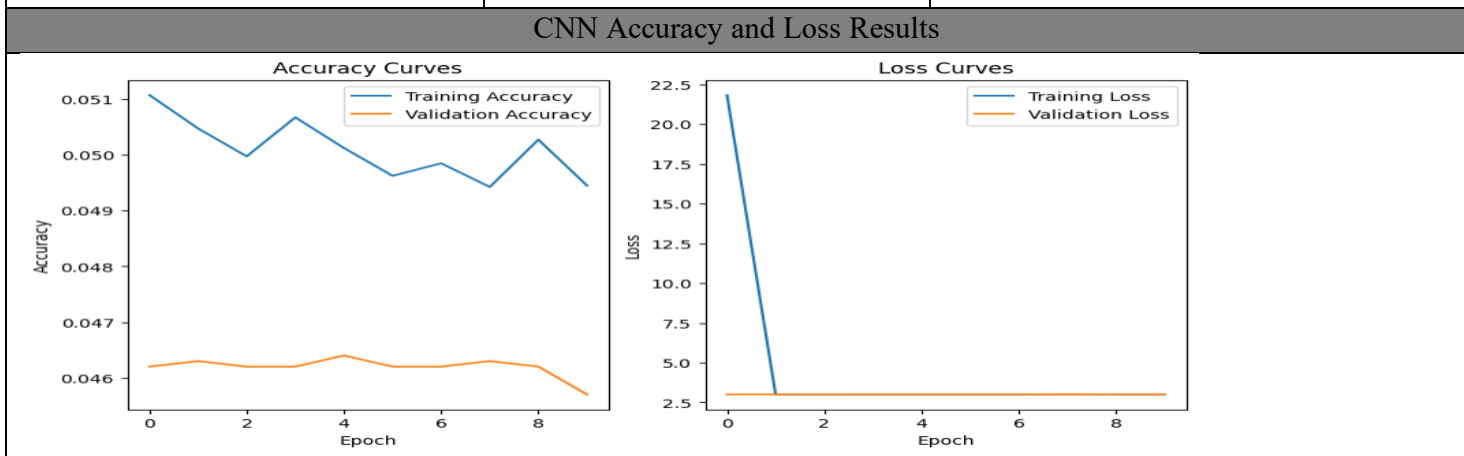


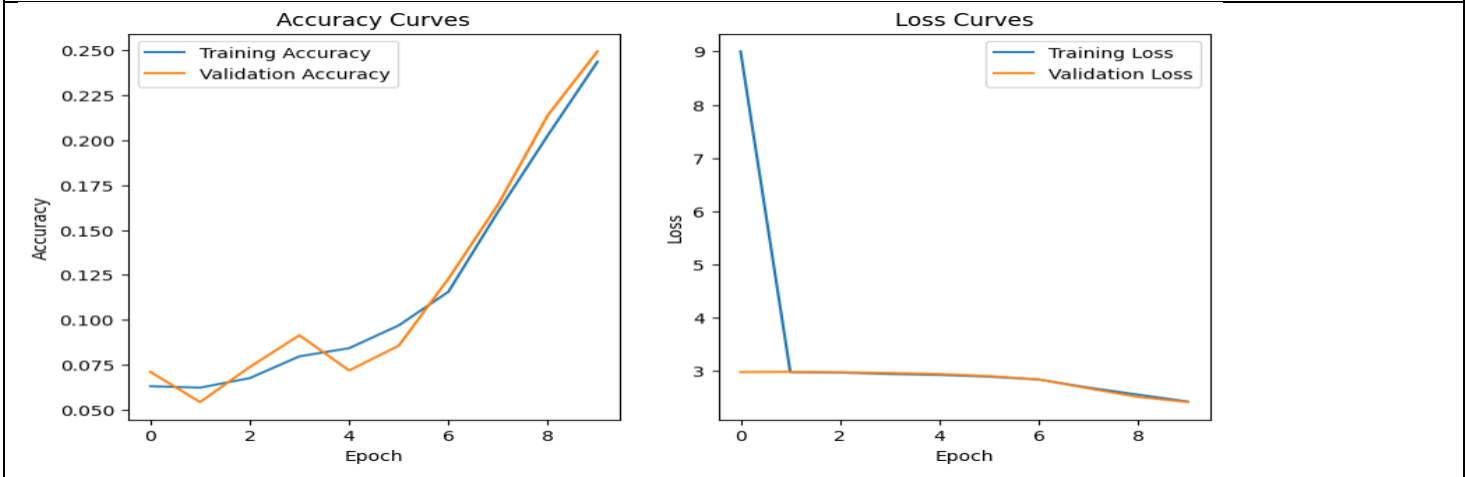
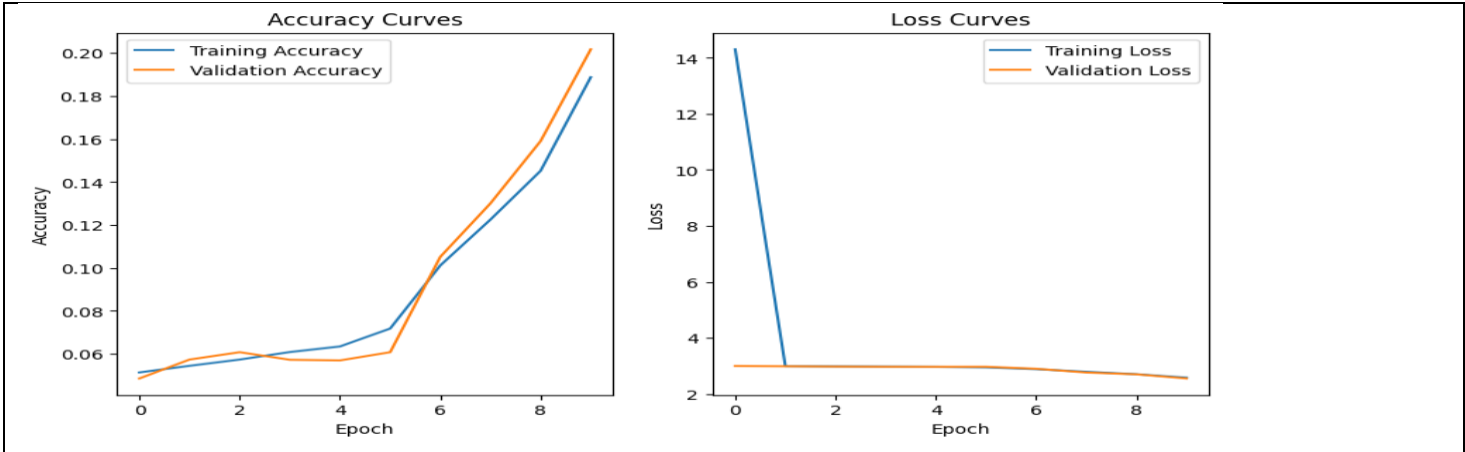
SVC Classification Report				
SVC Accuracy: 12.000				
	precision	recall	f1-score	support
aquatic mammals	0.00	0.00	0.00	2
fish	0.00	0.00	0.00	8
flowers	0.00	0.00	0.00	9
food containers	0.00	0.00	0.00	4
fruit and vegetables	0.00	0.00	0.00	4
household electrical devices	0.00	0.00	0.00	5
household furniture	0.00	0.00	0.00	2
insects	0.00	0.00	0.00	6
large carnivores	0.00	0.00	0.00	4
large man-made outdoor things	0.00	0.00	0.00	3
large natural outdoor scenes	0.40	0.60	0.48	10
large omnivores and herbivores	0.00	0.00	0.00	5
medium-sized mammals	0.09	0.86	0.17	7
non-insect invertebrates	0.00	0.00	0.00	2
people	0.00	0.00	0.00	5
reptiles	0.00	0.00	0.00	7
small mammals	0.00	0.00	0.00	2
trees	0.00	0.00	0.00	3
vehicles 1	0.00	0.00	0.00	6
vehicles 2	0.00	0.00	0.00	6
accuracy			0.12	100
macro avg	0.02	0.07	0.03	100
weighted avg	0.05	0.12	0.06	100

SVC Confusion Matrix Heatmap

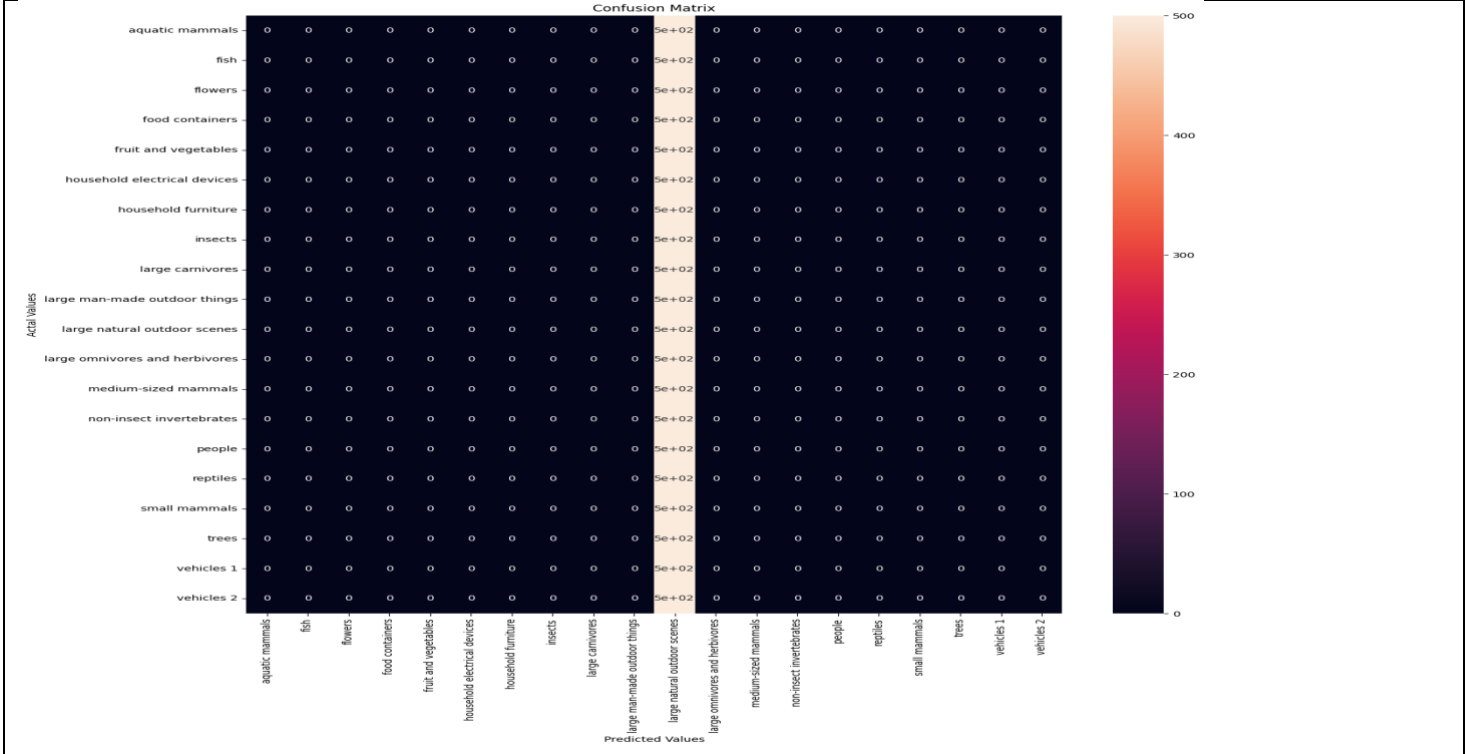


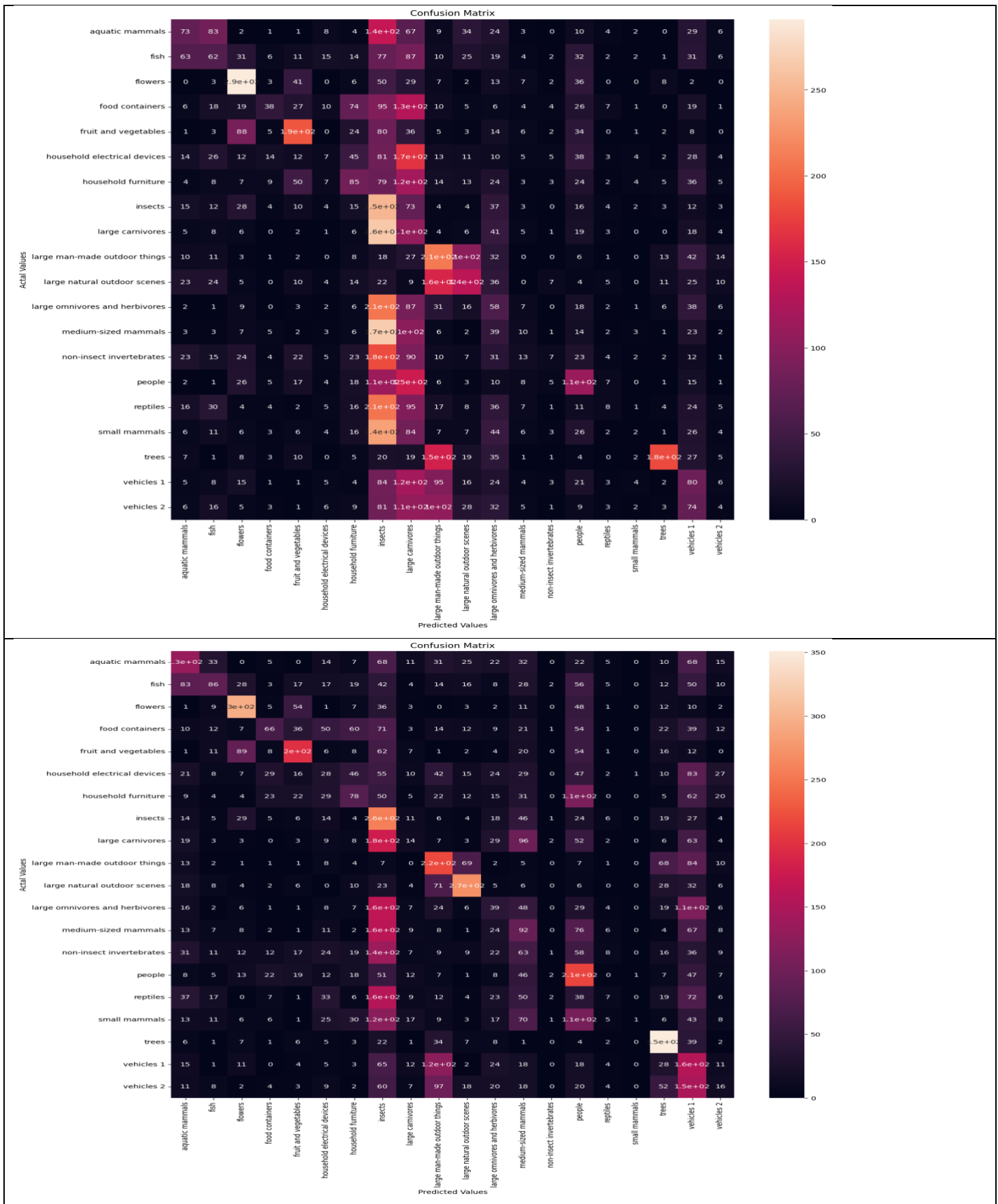
CNN Classification Report															
313/313 [=====] - 1s 3ms/step					313/313 [=====] - 2s 6ms/step					313/313 [=====] - 3s 9ms/step					
CNN Accuracy: 5.000					CNN Accuracy: 19.090					CNN Accuracy: 25.180					
	precision	recall	f1-score	support		precision	recall	f1-score	support		precision	recall	f1-score	support	
aquatic mammals	0.00	0.00	0.00	500	aquatic mammals	0.26	0.15	0.19	500	aquatic mammals	0.28	0.26	0.27	500	
fish	0.00	0.00	0.00	500	fish	0.18	0.12	0.15	500	fish	0.35	0.17	0.23	500	
flowers	0.00	0.00	0.00	500	flowers	0.49	0.58	0.53	500	flowers	0.55	0.59	0.57	500	
food containers	0.00	0.00	0.00	500	food containers	0.35	0.08	0.12	500	food containers	0.33	0.13	0.19	500	
fruit and vegetables	0.00	0.00	0.00	500	fruit and vegetables	0.45	0.38	0.41	500	fruit and vegetables	0.48	0.40	0.43	500	
household electrical devices	0.00	0.00	0.00	500	household electrical devices	0.08	0.01	0.02	500	household electrical devices	0.09	0.06	0.07	500	
household furniture	0.00	0.00	0.00	500	household furniture	0.21	0.17	0.19	500	household furniture	0.23	0.16	0.19	500	
insects	0.00	0.00	0.00	500	insects	0.10	0.50	0.16	500	insects	0.14	0.51	0.22	500	
large carnivores	0.00	0.00	0.00	500	large carnivores	0.06	0.21	0.10	500	large carnivores	0.09	0.03	0.04	500	
large man-made outdoor things	0.00	0.00	0.00	500	large man-made outdoor things	0.24	0.42	0.31	500	large man-made outdoor things	0.29	0.43	0.35	500	
large natural outdoor scenes	0.05	1.00	0.10	500	large natural outdoor scenes	0.30	0.27	0.29	500	large natural outdoor scenes	0.56	0.54	0.55	500	
large omnivores and herbivores	0.00	0.00	0.00	500	large omnivores and herbivores	0.10	0.12	0.11	500	large omnivores and herbivores	0.12	0.08	0.09	500	
medium-sized mammals	0.00	0.00	0.00	500	medium-sized mammals	0.10	0.02	0.03	500	medium-sized mammals	0.13	0.18	0.15	500	
non-insect invertebrates	0.00	0.00	0.00	500	non-insect invertebrates	0.15	0.01	0.03	500	non-insect invertebrates	0.08	0.00	0.00	500	
people	0.00	0.00	0.00	500	people	0.23	0.22	0.23	500	people	0.21	0.43	0.28	500	
reptiles	0.00	0.00	0.00	500	reptiles	0.13	0.02	0.03	500	reptiles	0.11	0.01	0.02	500	
small mammals	0.00	0.00	0.00	500	small mammals	0.06	0.00	0.01	500	small mammals	0.33	0.00	0.00	500	
trees	0.00	0.00	0.00	500	trees	0.74	0.36	0.49	500	trees	0.49	0.70	0.58	500	
vehicles 1	0.00	0.00	0.00	500	vehicles 1	0.14	0.16	0.15	500	vehicles 1	0.12	0.31	0.18	500	
vehicles 2	0.00	0.00	0.00	500	vehicles 2	0.05	0.01	0.01	500	vehicles 2	0.09	0.03	0.05	500	
accuracy			0.05	10000	accuracy			0.19	10000	accuracy			0.25	10000	
macro avg	0.00	0.05	0.00	10000	macro avg	0.22	0.19	0.18	10000	macro avg	0.25	0.25	0.22	10000	
weighted avg	0.00	0.05	0.00	10000	weighted avg	0.22	0.19	0.18	10000	weighted avg	0.25	0.25	0.22	10000	





Confusion Matrix HeatMap





2nd CNN Classification Report

ON2 Accuracy: 44.980

	precision	recall	f1-score	support
aquatic mammals	0.06	0.08	0.07	500
fish	0.05	0.04	0.05	500
flowers	0.04	0.05	0.04	500
food containers	0.03	0.03	0.03	500
fruit and vegetables	0.06	0.05	0.05	500
household electrical devices	0.06	0.04	0.05	500
household furniture	0.03	0.02	0.02	500
insects	0.06	0.07	0.06	500
large carnivores	0.06	0.07	0.06	500
large man-made outdoor things	0.05	0.06	0.06	500
large natural outdoor scenes	0.06	0.04	0.05	500
large omnivores and herbivores	0.04	0.04	0.04	500
medium-sized mammals	0.05	0.03	0.04	500
non-insect invertebrates	0.05	0.04	0.05	500
people	0.05	0.08	0.06	500
reptiles	0.06	0.06	0.06	500
small mammals	0.03	0.02	0.02	500
trees	0.04	0.04	0.04	500
vehicles 1	0.03	0.02	0.02	500
vehicles 2	0.05	0.07	0.06	500
accuracy			0.05	10000
macro avg	0.05	0.05	0.05	10000
weighted avg	0.05	0.05	0.05	10000

ON2 Accuracy: 49.300

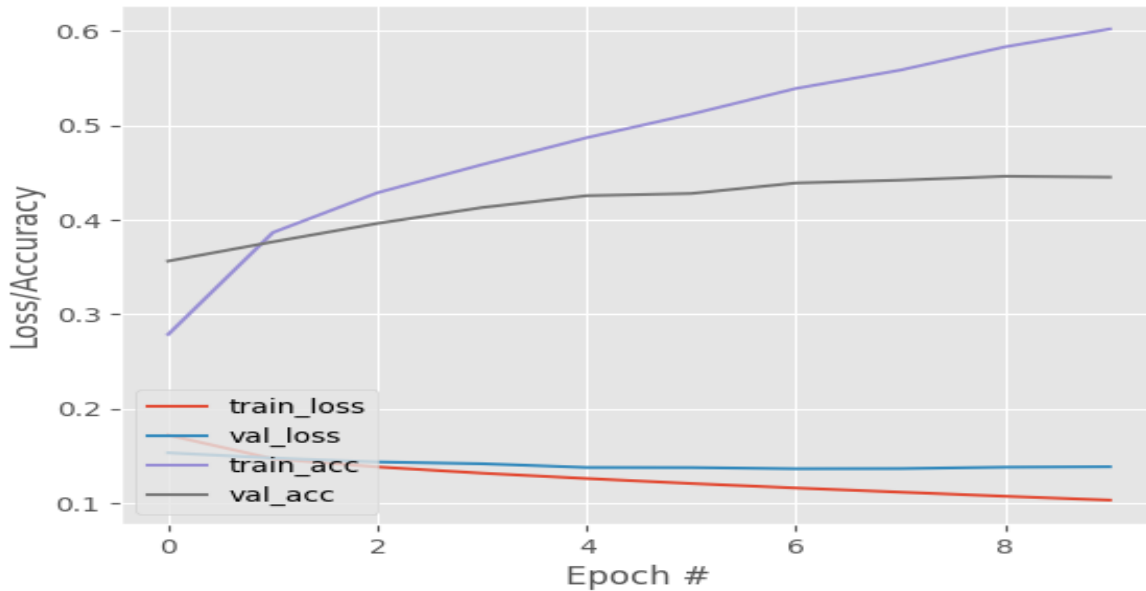
	precision	recall	f1-score	support
aquatic mammals	0.05	0.04	0.05	500
fish	0.06	0.04	0.05	500
flowers	0.05	0.07	0.06	500
food containers	0.05	0.07	0.06	500
fruit and vegetables	0.05	0.05	0.05	500
household electrical devices	0.06	0.07	0.06	500
household furniture	0.05	0.08	0.06	500
insects	0.06	0.04	0.05	500
large carnivores	0.06	0.08	0.07	500
large man-made outdoor things	0.04	0.05	0.05	500
large natural outdoor scenes	0.06	0.05	0.05	500
large omnivores and herbivores	0.07	0.06	0.06	500
medium-sized mammals	0.04	0.02	0.03	500
non-insect invertebrates	0.04	0.02	0.03	500
people	0.05	0.07	0.06	500
reptiles	0.04	0.02	0.03	500
small mammals	0.05	0.03	0.04	500
trees	0.04	0.05	0.04	500
vehicles 1	0.05	0.04	0.05	500
vehicles 2	0.05	0.06	0.05	500
accuracy			0.05	10000
macro avg	0.05	0.05	0.05	10000
weighted avg	0.05	0.05	0.05	10000

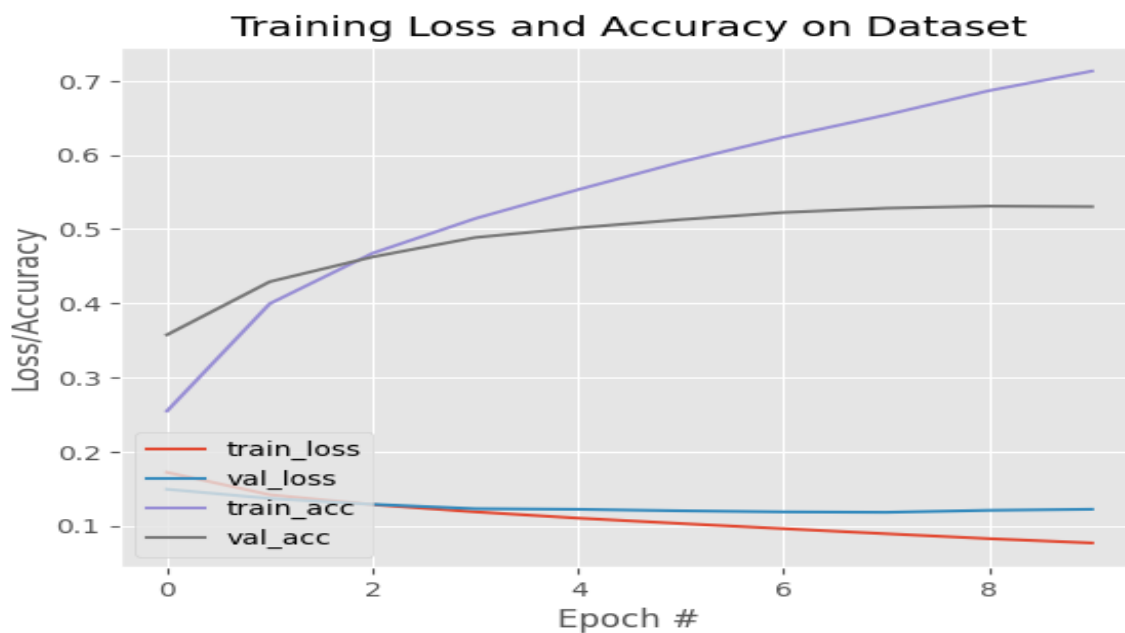
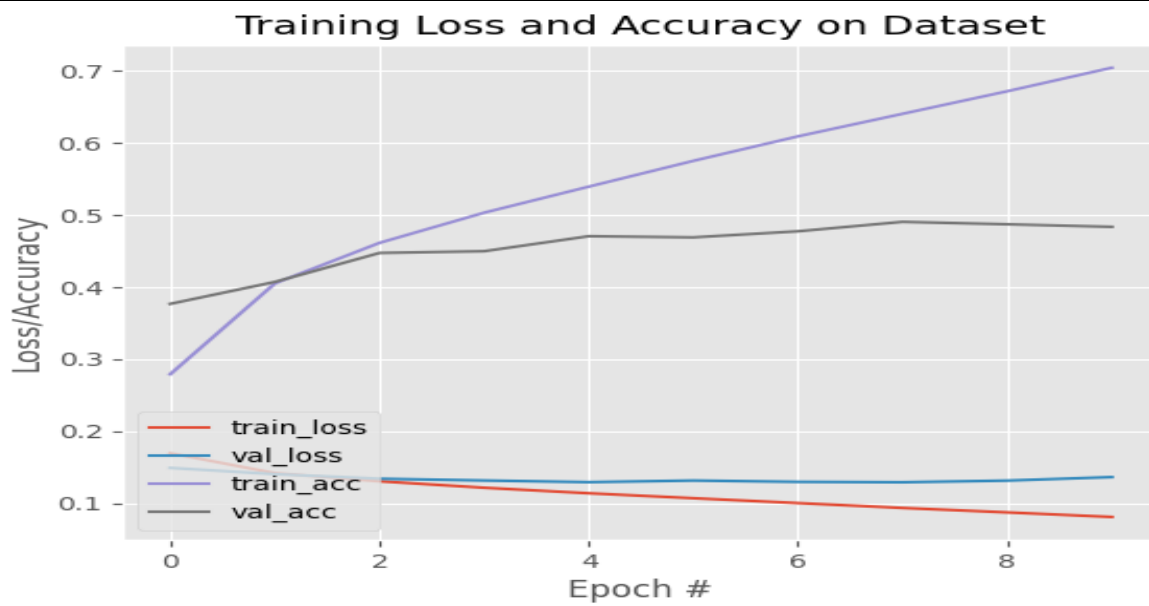
ON2 Accuracy: 53.560

	precision	recall	f1-score	support
aquatic mammals	0.05	0.05	0.05	500
fish	0.07	0.06	0.06	500
flowers	0.05	0.04	0.05	500
food containers	0.04	0.03	0.03	500
fruit and vegetables	0.05	0.07	0.06	500
household electrical devices	0.06	0.05	0.05	500
household furniture	0.05	0.05	0.05	500
insects	0.04	0.05	0.05	500
large carnivores	0.04	0.04	0.04	500
large man-made outdoor things	0.06	0.10	0.07	500
large natural outdoor scenes	0.04	0.04	0.04	500
large omnivores and herbivores	0.06	0.05	0.05	500
medium-sized mammals	0.06	0.08	0.07	500
non-insect invertebrates	0.06	0.04	0.05	500
people	0.04	0.05	0.04	500
reptiles	0.05	0.04	0.04	500
small mammals	0.05	0.05	0.05	500
trees	0.06	0.06	0.06	500
vehicles 1	0.04	0.04	0.04	500
vehicles 2	0.05	0.05	0.05	500
accuracy			0.05	10000
macro avg	0.05	0.05	0.05	10000
weighted avg	0.05	0.05	0.05	10000

2nd CNN Accuracy and Loss Results

Training Loss and Accuracy on Dataset





2nd CNN Confusion Matrix HeatMap

