

## **Training and Validation results**

```
Epoch 17/50
313/313 - 0s - 1ms/step - accuracy: 0.8430 - loss: 0.4423 - val_accuracy: 0.7887 - val_loss: 0.6714 - learning_rate: 1.2500e-04
Epoch 18/50
313/313 - 0s - 1ms/step - accuracy: 0.8451 - loss: 0.4289 - val_accuracy: 0.7787 - val_loss: 0.7229 - learning_rate: 1.2500e-04
Epoch 19/50
313/313 - 0s - 1ms/step - accuracy: 0.8493 - loss: 0.4208 - val_accuracy: 0.7844 - val_loss: 0.7117 - learning_rate: 1.2500e-04
Epoch 20/50
313/313 - 0s - 1ms/step - accuracy: 0.8523 - loss: 0.4084 - val_accuracy: 0.7870 - val_loss: 0.7041 - learning_rate: 1.2500e-04
Epoch 21/50
313/313 - 0s - 1ms/step - accuracy: 0.8592 - loss: 0.3936 - val_accuracy: 0.7829 - val_loss: 0.7329 - learning_rate: 6.2500e-05
Epoch 22/50
313/313 - 0s - 2ms/step - accuracy: 0.8602 - loss: 0.3872 - val_accuracy: 0.7895 - val_loss: 0.6901 - learning_rate: 6.2500e-05
Epoch 23/50
313/313 - 0s - 1ms/step - accuracy: 0.8610 - loss: 0.3836 - val_accuracy: 0.7881 - val_loss: 0.7037 - learning_rate: 6.2500e-05
Epoch 24/50
313/313 - 0s - 2ms/step - accuracy: 0.8655 - loss: 0.3764 - val_accuracy: 0.7888 - val_loss: 0.7033 - learning_rate: 3.1250e-05
Epoch 25/50
313/313 - 0s - 1ms/step - accuracy: 0.8659 - loss: 0.3719 - val_accuracy: 0.7920 - val_loss: 0.6873 - learning_rate: 3.1250e-05
Epoch 26/50
313/313 - 0s - 1ms/step - accuracy: 0.8669 - loss: 0.3699 - val_accuracy: 0.7871 - val_loss: 0.7118 - learning_rate: 3.1250e-05
Epoch 27/50
313/313 - 0s - 2ms/step - accuracy: 0.8681 - loss: 0.3660 - val_accuracy: 0.7915 - val_loss: 0.6998 - learning_rate: 1.5625e-05
```

Figure 1: Last 10 epochs ran for CNN model

Referring to figure 1, we can observe that the final training accuracy and loss reached 0.8681 and 0.3660, and the final validation accuracy and loss reached 0.7915 and 0.6998.

With the CNN reaching the smallest validation loss of 0.6714 and validation accuracy of 0.7887 in epoch 17.

```
Epoch 90/100
44/44 227s 5s/step - accuracy: 0.8660 - loss: 0.3820 - val_accuracy: 0.8194 - val_loss: 0.5286 - learning_rate: 0.0010
Epoch 91/100
44/44 228s 5s/step - accuracy: 0.8621 - loss: 0.3918 - val_accuracy: 0.8316 - val_loss: 0.4992 - learning_rate: 0.0010
Epoch 92/100
44/44 226s 5s/step - accuracy: 0.8638 - loss: 0.3934 - val_accuracy: 0.8212 - val_loss: 0.5183 - learning_rate: 0.0010
Epoch 93/100
44/44 227s 5s/step - accuracy: 0.8682 - loss: 0.3749 - val_accuracy: 0.8228 - val_loss: 0.5252 - learning_rate: 0.0010
Epoch 94/100
44/44 227s 5s/step - accuracy: 0.8653 - loss: 0.3821 - val_accuracy: 0.8284 - val_loss: 0.5037 - learning_rate: 0.0010
Epoch 95/100
44/44 227s 5s/step - accuracy: 0.8682 - loss: 0.3760 - val_accuracy: 0.8254 - val_loss: 0.5150 - learning_rate: 0.0010
Epoch 96/100
44/44 226s 5s/step - accuracy: 0.8732 - loss: 0.3694 - val_accuracy: 0.8158 - val_loss: 0.5352 - learning_rate: 0.0010
Epoch 97/100
44/44 227s 5s/step - accuracy: 0.8744 - loss: 0.3620 - val_accuracy: 0.8284 - val_loss: 0.5050 - learning_rate: 0.0010
Epoch 98/100
44/44 227s 5s/step - accuracy: 0.8709 - loss: 0.3633 - val_accuracy: 0.8236 - val_loss: 0.5143 - learning_rate: 0.0010
Epoch 99/100
44/44 229s 5s/step - accuracy: 0.8808 - loss: 0.3491 - val_accuracy: 0.8344 - val_loss: 0.4837 - learning_rate: 2.0000e-04
Epoch 100/100
44/44 227s 5s/step - accuracy: 0.8950 - loss: 0.3038 - val_accuracy: 0.8372 - val_loss: 0.4914 - learning_rate: 2.0000e-04
```

Figure 2: Last 10 epochs ran for ViT model

Referring to figure 2, we can observe that the final training accuracy and loss reached 0.8950 and 0.3038, and the final validation accuracy and loss reached 0.8372 and 0.4914.

With the ViT reaching the smallest validation loss of 0.4914 and validation accuracy of 0.8372 in epoch 100.

## **Performance comparison**

The ViT performed better than the CNN both in terms of a moderately lower validation loss score and a slightly higher validation accuracy score.

## **Training efficiency and learning dynamics**

The CNN was able to reach its peak and plateaued early on at epoch 17, and was stopped prematurely by the implemented early stopping procedure.

Whereas we can observe that even in its final few runs, the ViT displayed continued growth, and likely could have gained an even better validation accuracy and loss score, if allowed to run for more epochs. However, the ViT is limited by its much slower training time (average of 227s per epoch in ViT vs average of 300ms per epoch in CNN).

This is explained by the nature of CNN, utilising convolutional filters that while allowing the model to quickly learn local relationships between neighbouring pixels. This allows CNN to quickly reach its peak but struggle to see much gains in performance.

But in a ViT, utilising self attention and being able to extract global complex patterns (not just limited to local relationships) in images allows the ViT to outperform the CNN in the long run. But without the spatial assumption in ViT when it first performs patching, the ViT has to learn position embeddings of the image patches. This is observed in the training results where the ViT requires many epochs (16 epochs on ViT) to hit a validation accuracy of 0.70, while the CNN easily achieves that in 5 epochs.