

## EECS 445 Project 2

Shiyu Liu

### Question 1

(a) i. The mean of RGB color channels are [118.263 118.004 118.594] respectively and standard deviation of RGB color channels are [66.307 69.039 75.332] respectively.

ii. We extract the mean and standard deviation from the training set to ensure a consistency across all data partitions which makes it easier to compare the performances. It also helps the model to better generalize from the training data to unseen data.

(b)

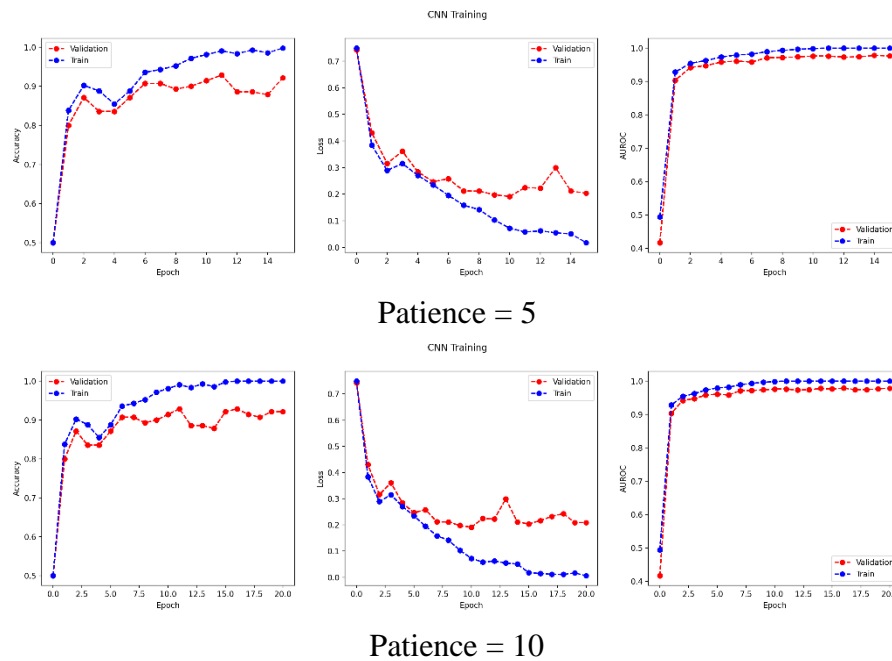


### Question 2

(a) Convolutional Layer 1:  $(3 \times 5 \times 5 + 1) \times 16 = 1216$   
Convolutional Layer 2:  $(16 \times 5 \times 5 + 1) \times 64 = 25664$   
Convolutional Layer 3:  $(64 \times 5 \times 5 + 1) \times 8 = 12808$   
Fully Connected Layer 1:  $(32 + 1) \times 2 = 66$   
Total learnable parameters:  $1216 + 25664 + 12808 + 66 = 39754$

(f) i. One possible reason for this behavior is overfitting. As the model gets more complex, it might capture noise that reduces validation performance. Another possible reason might be the early stopping criteria. If the patience parameter for early stopping is too small, the training may terminate prematurely based on minor fluctuations in the validation loss.

ii.



The model stops at epoch 20. Patience value of 10 works better because it results in a more stabilized validation loss compared to a patience value of 5. In a dataset with relatively large noise and more outliers, it might be better to increase patience value because it allows the model to recover from sudden changes and allow the model to continuously improve.

iii. The new size of the input to the fully connected layer is  $64 \times 2 \times 2 = 256$ .

	Epoch	Training AUROC	Validation AUROC
8 filters	10	0.9987	0.9761
64 filters	7	0.9993	0.9698

The Validation AUROC decreases as we increase the number of filters from 8 to 64. The mismatch between dataset and model complexity might account for this change: when dataset is relatively simple and can be well represented with a smaller number of filters, increasing the number of filters may lead to an overfitted model that struggles to capture relevant features which causes a decrease in the Validation AUROC.

(g) i.

	Training	Validation	Testing
Accuracy	0.981	0.9143	0.6161
AUROC	0.9987	0.9761	0.684

ii. There is some evidence of overfitting: the training accuracy is very close to 1 and much higher than validation accuracy.

iii. The performance of validation is much better than test. One possible hypothesis for this issue is lack of interpretability of CNN: the model might not emphasize the pixels that represent the Pantheon and Hofburg Imperial Palace but on other pixels that are irrelevant for the subject.

### Question 3

(a)

$$\alpha'_1 = \frac{1}{16} \sum_i \sum_j \frac{\partial y'_i}{\partial A_{ij}^{(1)}} = \frac{1}{16} (-1 + 1 - 2 - 1 - 1 + 1 + 1 + 1 + 2 + 2) = \frac{3}{16}$$

$$\alpha'_2 = \frac{1}{16} \sum_i \sum_j \frac{\partial y'_i}{\partial A_{ij}^{(2)}} = \frac{1}{16} (1 + 2 + 2 + 2 + 2 + 1 + 1 - 1 - 2 - 1) = \frac{7}{16}$$

$$\begin{aligned} \mathcal{L}_1 &= \text{ReLU} \left( \frac{3}{16} A^{(1)} + \frac{7}{16} A^{(2)} \right) \\ &= \text{ReLU} \left( \frac{3}{16} \begin{bmatrix} 1 & 1 & 2 & 1 \\ 1 & 2 & 1 & 0 \\ 0 & 1 & 0 & -1 \\ 0 & 1 & -2 & -2 \end{bmatrix} + \frac{7}{16} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 1 & 0 \\ -1 & -1 & 1 & 0 \end{bmatrix} \right) \\ &= \text{ReLU} \left( \begin{bmatrix} \frac{5}{8} & \frac{5}{8} & \frac{13}{16} & \frac{5}{8} \\ \frac{17}{16} & \frac{5}{4} & \frac{17}{16} & \frac{7}{8} \\ \frac{7}{8} & \frac{17}{16} & \frac{7}{16} & -\frac{3}{16} \\ -\frac{7}{16} & -\frac{1}{4} & -\frac{13}{16} & -\frac{7}{8} \end{bmatrix} \right) = \begin{bmatrix} \frac{5}{8} & \frac{5}{8} & \frac{13}{16} & \frac{5}{8} \\ \frac{17}{16} & \frac{5}{4} & \frac{17}{16} & \frac{7}{8} \\ \frac{7}{8} & \frac{17}{16} & \frac{7}{16} & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \end{aligned}$$

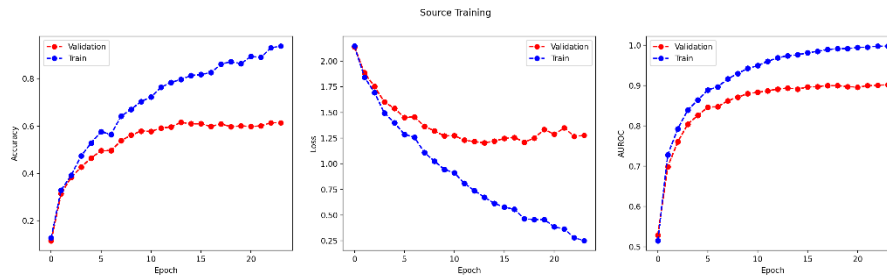
(b) The CNN appears to use features with light colors to identify the Hofburg Imperial Palace class. Many white and light blue objects such as the sky and clothes have high values in the Viridis color scale.

(c) The Grad-Cam visualizations is inline with my predictions since the pixels that have high intensity is irrelevant to the prediction of the Hofburg Imperial Palace class. The pixels that are relevant to the architectures doesn't have a high weight so the model doesn't perform well on the test set.

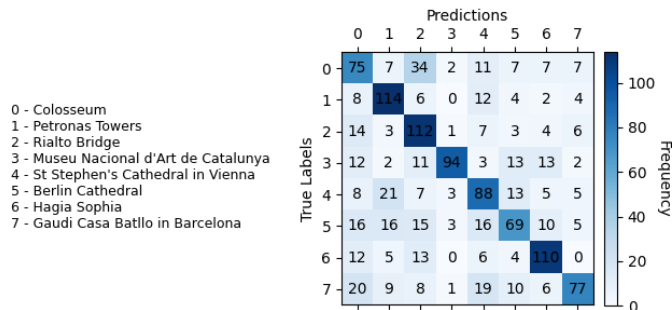
## Question 4

### 4.1 Transfer Learning

(c) The epoch number with the lowest validation loss is epoch 13.



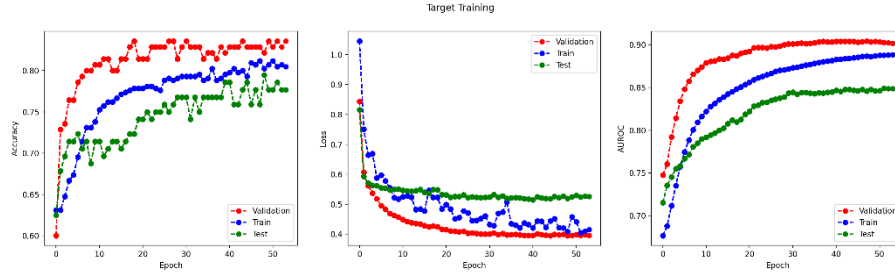
(d)



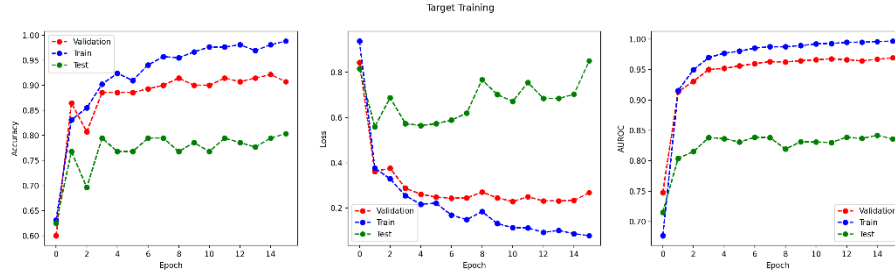
The most accurate landmark is 1 – Petronas Towers and the least accurate landmark is 5 - Berlin Cathedral. A possible reason for this result is because Petronas Towers has the lightest color compared with its surroundings whereas the Berlin Cathedral has the darkest color compared with its surroundings. If the landmark has light colors, the model can capture its shape and make more accurate predictions.

(f)

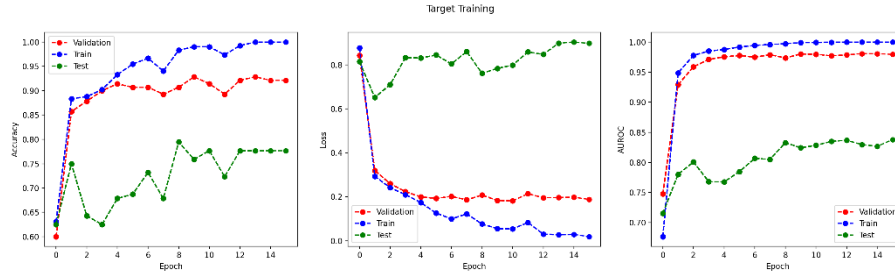
	AUROC		
	TRAIN	VAL	TEST
Freeze all CONV layers (Fine-tune FC layer)	0.8878	0.9027	0.8489
Freeze first two CONV layers (Fine-tune last CONV and FC layers)	0.9922	0.9663	0.8307
Freeze first CONV layer (Fine-tune last 2 conv. and fc layers)	0.9994	0.9798	0.8284
Freeze no layers (Fine-tune all layers)	1.0	0.9837	0.8138
No Pretraining or Transfer Learning (Section 2(g) performance)	0.9987	0.9761	0.684



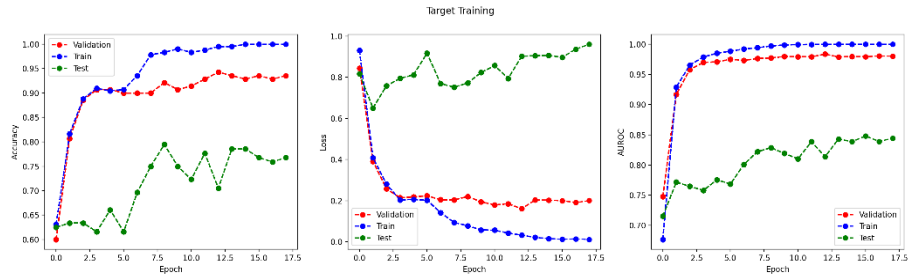
Freeze all CONV layers (Fine-tune FC layer)



Freeze first two CONV layers (Fine-tune last CONV and FC layers)



Freeze first CONV layer (Fine-tune last 2 conv. and fc layers)



Freeze no layers (Fine-tune all layers)

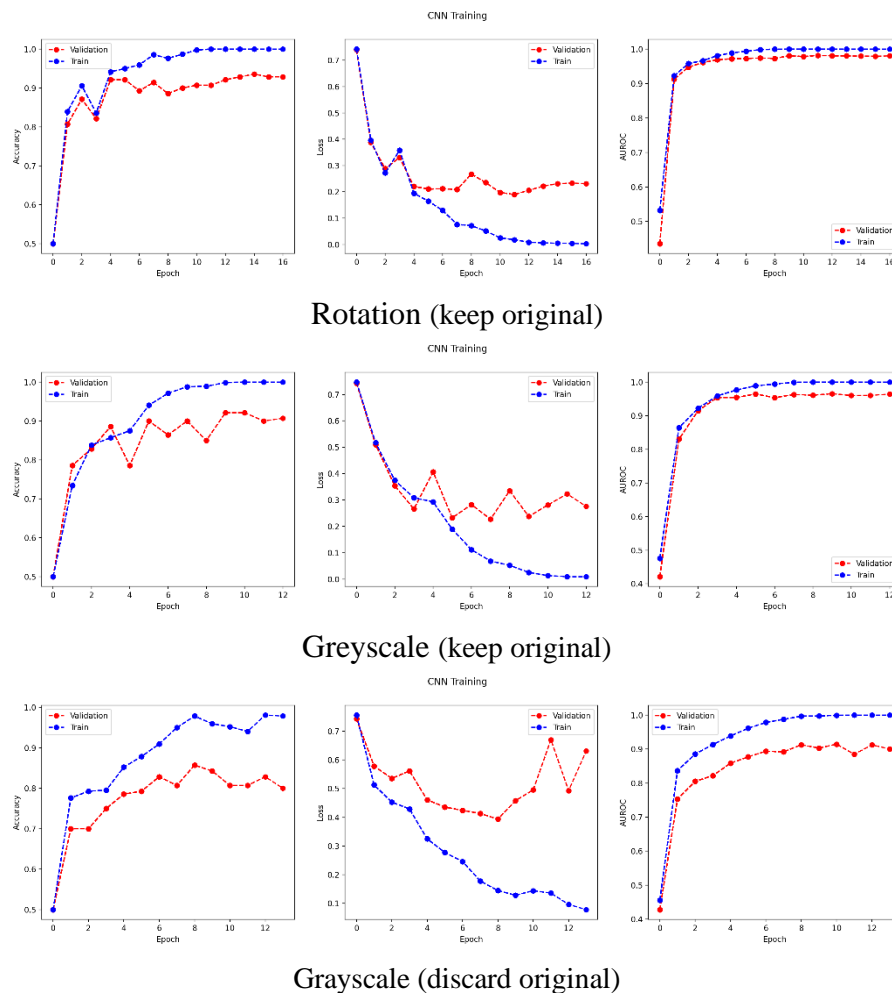
From the table above, we can see that while some models have lower training and validation AUROC than the model from 2g, all of them has higher testing AUROC than the original model (the model that freezes all layers have the best performance).

Comparing with CNN in section 2, we can see that transfer learning is definitely helpful, since all of the test AUROCs from transfer learning models are better than the original model. The source task is also helpful because it helps us determine the epoch with the lowest validation loss in the source training data so we can better generalize the target model and reduces potential bias. I hypothesize it is helpful because the source task is like bootstrapping: it finds out the best result using the population and then use it to generalize one portion of it. It can achieve better performance because it leverages more knowledge gained from the source than just from the limited target data. When all layers

are frozen, the training curve takes very long time to converge and the training loss is lower than both validation and testing loss, whereas for other models, the convergence rate is faster and the validation losses are lower than the training losses. As more layers are frozen, the training and validation AUROC decreases whereas the testing AUROC increases. There are a few possible explanations for this: The original model is overfitted to the target. As more layers are frozen, less weight and bias are influenced by training, but this causes a higher testing result, indicating that the original model fits too closely to the training data. Another reason is the lack of relevance for the task. By freezing more layers, we restrict the model from learning task-specific features, causing a decrease in AUROC for both training and validation. However, the testing AUROC still improves because the model is trained more can still provide a valuable information for classification.

## 4.2 Data Augmentation

(b) iii.



	AUROC		
	TRAIN	VAL	TEST
Rotation (keep original)	1.0	0.9816	0.7423
Grayscale (keep original)	0.9993	0.9631	0.7682
Grayscale (discard original)	0.997	0.9129	0.8058
No augmentation (Section 2(g) performance)	0.9987	0.9761	0.684

All of these augmentations provide better performance than the original model. Grayscale (discard original) has the highest test AUROC comparing with other augmentations, but also have the lowest training and validation AUROC.

- (c) Yes, the validation and training plots change as I apply these augmentation techniques. One possible reason for this might be that these augmentations increase data variability. When modified data is introduced to the original set, the model may have stronger interpretability when dealing with a specific task because it has strong generalization power through training with the new data. Another possible reason is that the augmentation has regularization effect that prevents overfitting: The larger variety of data prevents the model to overfit to a specific form and also results in more changes in the training curve.

#### Question 5 Challenge

In my model, I used the same criterion and optimizer as Appendix B. I tried to use weight decay, but the result is better without it. I also tried patience value between 5 and 10 and finally decided on 7, a value in-between. Patience of 5 converges too quickly, leading to a lot of variations in the learning curves and patience of 10 takes too many epochs to converge (about 70). For the model architectures, I added more convolutional layers with larger input and output sizes, added more pooling layers and reduces the sizes of kernel, stride and padding. I gained inspiration from the architecture that I learned from another course I'm currently taking (EECS 442). Large input and output dimensions provide a richer source of information which allow for more complex features to be learned. Smaller kernel sizes also allow for more intricate feature extraction at different levels of abstraction. Combining these would potentially increase predictive power. For transfer learning, I tried to freeze different numbers of layers but they didn't produce better result than the original model, so I decide not to use it. For Data Augmentation, I tried all three kinds of augmentation we did in 4.2 and discovered that Grayscale (discard original) produces the best result, So I used it in my model. For model evaluation, I used lowest validation loss to find the best checkpoint and used test AUROC to determine which model is the best because it is a comprehensive metric that generalizes all aspects of the model's performance.