```
+----------+---------+--------------------+--------------------+
| Metric   | Clean   | No Defense Attack   | w/ Defense Attack  |
+==========+=========+====================+====================+
| Loss     | 0.031   | 4.38               | 0.031              |
+----------+---------+--------------------+--------------------+
| Accuracy | 98%     | 12%                | 98%                |
+----------+---------+--------------------+--------------------+
```

Example Misclassifications:


------------------------------------------------------------
Number of misclassified images for Clean: 1
Attack: fgsm_cw_attack
Dataset: MNIST
Training Epochs: 10
Trained Clean Images: 64
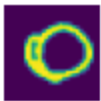Test Images: 140
Accuracy: 0.98
Precision: 0.99
Recall: 0.99
F1-score: 0.99
ROC AUC Score: 1.00
------------------------------------------------------------

Misclassifications:
0 -> 6: 1



------------------------------------------------------------
Number of misclassified images for No Defense Attack: 56
Attack: fgsm_cw_attack
Dataset: MNIST
Training Epochs: 10
Adversarial Training Images: 60
Test Images: 140
Accuracy: 0.12
Precision: 0.23
Recall: 0.20
F1-score: 0.22
ROC AUC Score: 1.00
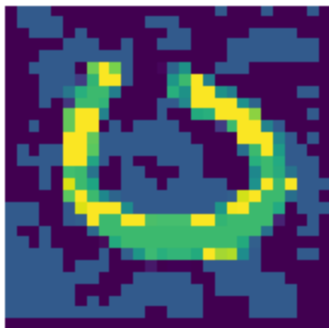------------------------------------------------------------

Misclassifications:
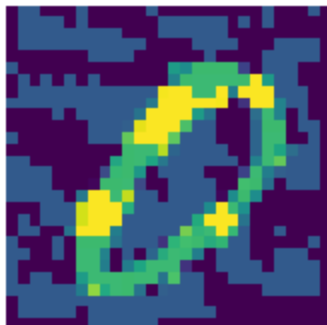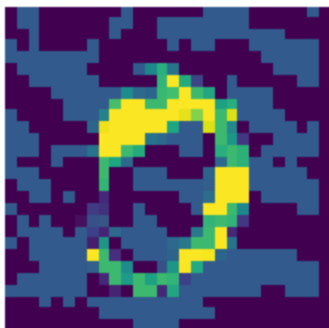0 -> 9: 9
0 -> 7: 7
0 -> 8: 3
0 -> 6: 21
0 -> 5: 14
0 -> 4: 2

0 -> 5

0 -> 9

0 -> 6

0 -> 9

0 -> 5

0 -> 7

0 -> 5

0 -> 7

0 -> 5

0 -> 5

0 -> 6



0 -> 7



0 -> 5



0 -> 6



0 -> 6



0 -> 7



0 -> 6



0 -> 9



0 -> 9



0 -> 6

0 -> 8

0 -> 5



0 -> 9

0 -> 6



0 -> 9

0 -> 6



0 -> 5

0 -> 5



0 -> 7

0 -> 6

0 -> 6



0 -> 6



0 -> 4



0 -> 9



0 -> 6



0 -> 6



0 -> 8



0 -> 6



0 -> 6



0 -> 5

0 -> 6                                              0 -> 6



0 -> 6                                              0 -> 6



0 -> 4                                              0 -> 5



```
-----------------------------------------------------------
Number of misclassified images for w/ Defense Attack: 1
Attack: fgsm_cw_attack
Dataset: MNIST
Training Epochs: 10
Retrained Clean and Adversarial Images: 124
Test Images: 140
Accuracy: 0.98
Precision: 0.99
Recall: 0.99
F1-score: 0.99
ROC AUC Score: 1.00
-----------------------------------------------------------

Misclassifications:
0 -> 6: 1
```
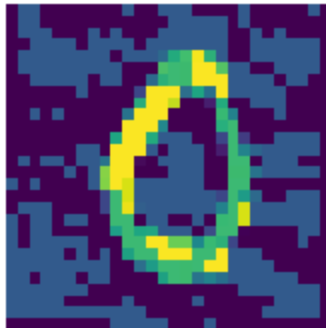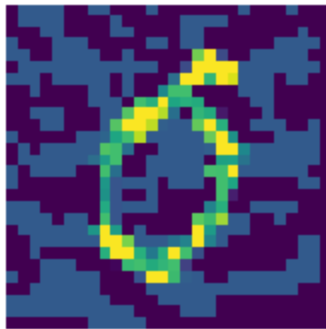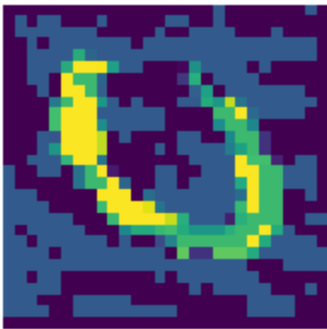
0 -> 6