Defense Implemented: Input Transformation

Input Transformation Approach: Adversarial Logit Pairing

| Metric | Clean | No Defense Attack | w/ Defense Attack |
|---|---|---|---|
| Loss | 0.00138 | 3.72 | 0.00273 |
| Accuracy | 100% | 20% | 100% |

Example Misclassifications:

```
---------------------------------------------------------
No misclassified images for stage: Clean
Attack: fgsm_pgd_attack
Dataset: MNIST
Training Epochs: 10
Trained Clean Images: 54210
Test Images: 140
Accuracy: 1.00
Precision: 1.00
Recall: 1.00
F1-score: 1.00
ROC AUC score is not defined for a single class.
---------------------------------------------------------
```

```
---------------------------------------------------------
Number of misclassified images for No Defense Attack: 51
Attack: fgsm_pgd_attack
Dataset: MNIST
Training Epochs: 10
Adversarial Training Images: 27105
Test Images: 140
Accuracy: 0.20
Precision: 0.35
Recall: 0.31
F1-score: 0.33
ROC AUC Score: 1.00
---------------------------------------------------------
```

Misclassifications:
0 -> 6: 29
0 -> 2: 8
0 -> 7: 1
0 -> 5: 4
0 -> 8: 3
0 -> 4: 2
0 -> 9: 4

```
---------------------------------------------------------
```

No misclassified images for stage: w/ Defense Attack
Attack: fgsm_pgd_attack
Dataset: MNIST
Training Epochs: 10
Retrained Clean and Adversarial Images: 81315
Test Images: 140
Accuracy: 1.00
Precision: 1.00
Recall: 1.00
F1-score: 1.00
ROC AUC score is not defined for a single class.
---------------------------------------------------------