

# Game of Deception & Truth



Mentor: Jerry Liao

Member: Shih-Hsuan Lin, Aaron Kao, Johnson Chang, Irene Chen



# Outline

## Background

- Adversarial Attack
- Adversarial Defense

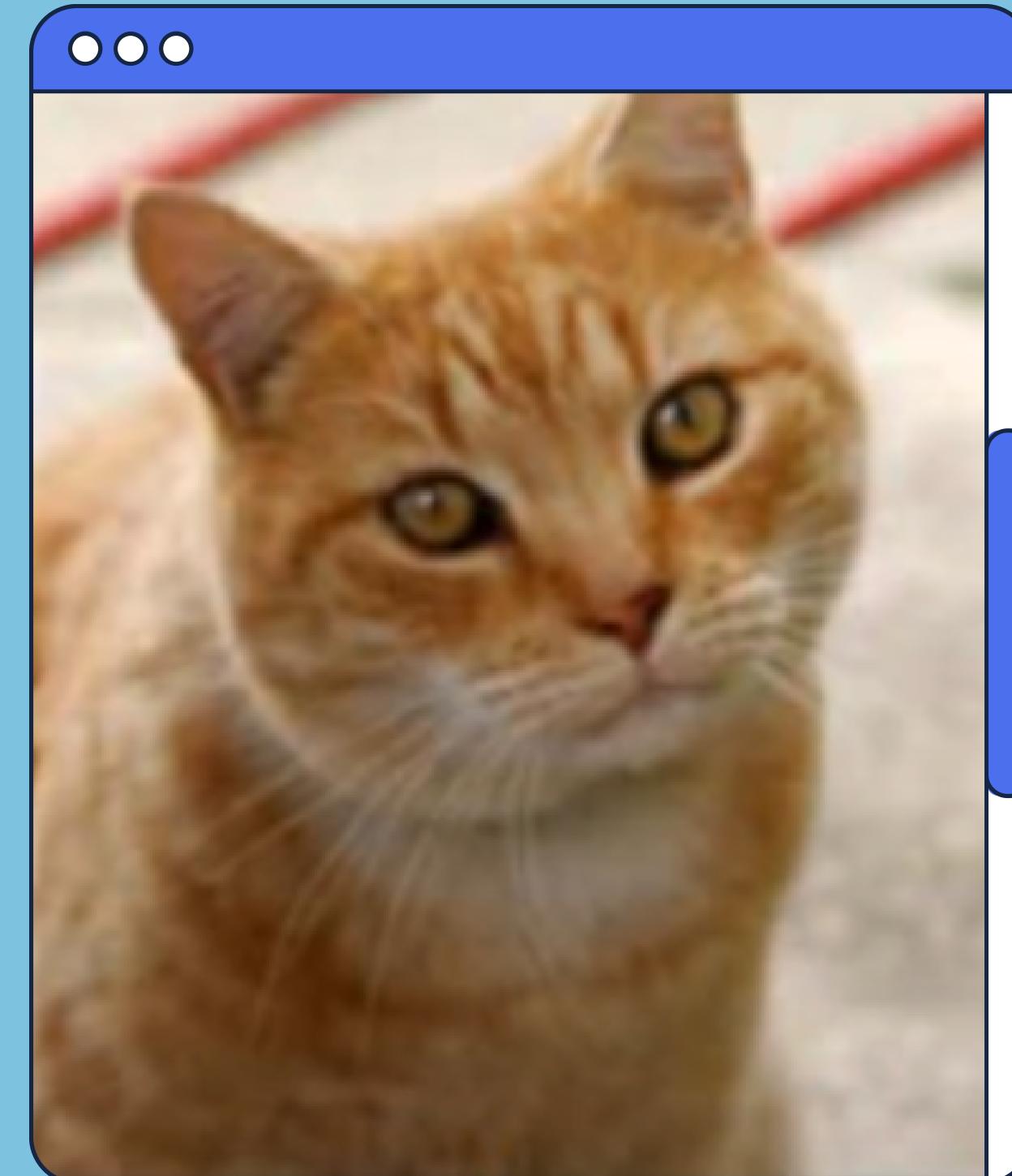
## Methodologies and Results

- Fast Gradient Sign Method (FGSM)
- DeepFool
- Adversarial Perturbation Elimination with GAN (APE-GAN)
- Diffusion model

## Summary & Future Works

# What Is Adversarial Attack?

Original Image

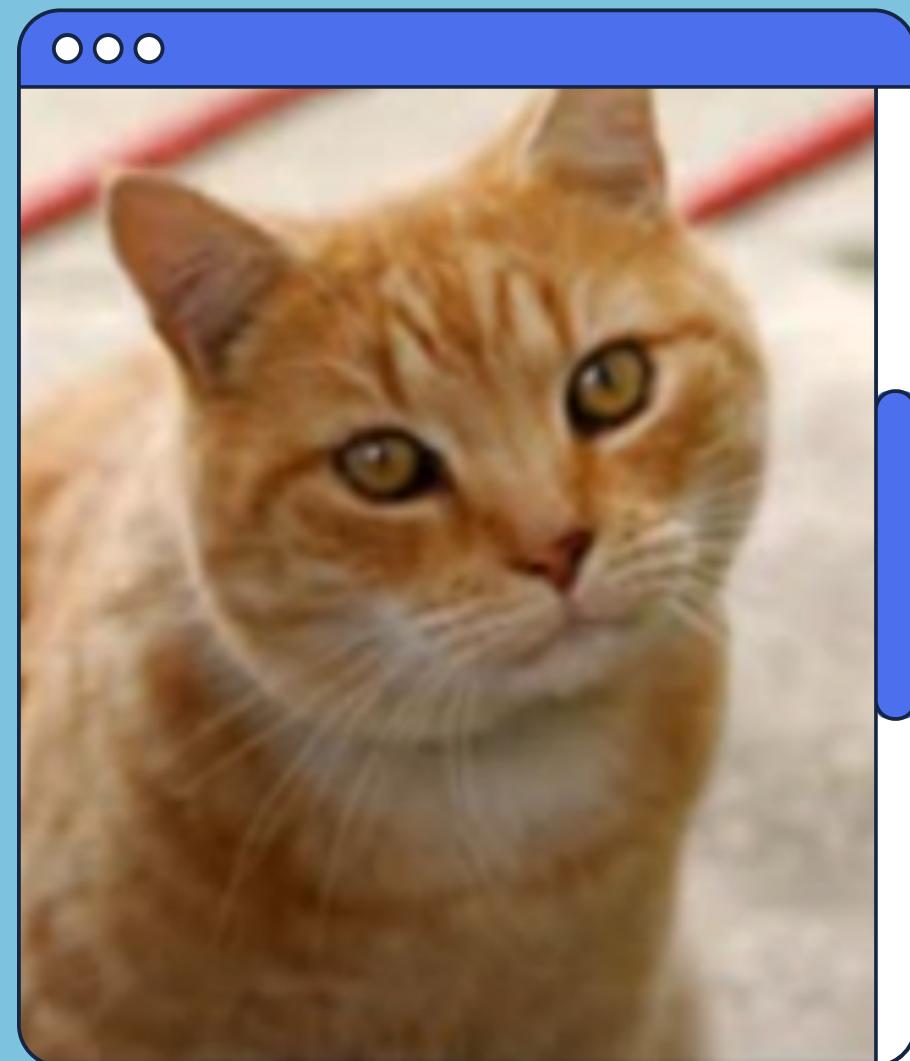


Tiger Cat. 0.64

Reference: Hung-Yi Lee, Machine Learning Slides

# What Is Adversarial Attack?

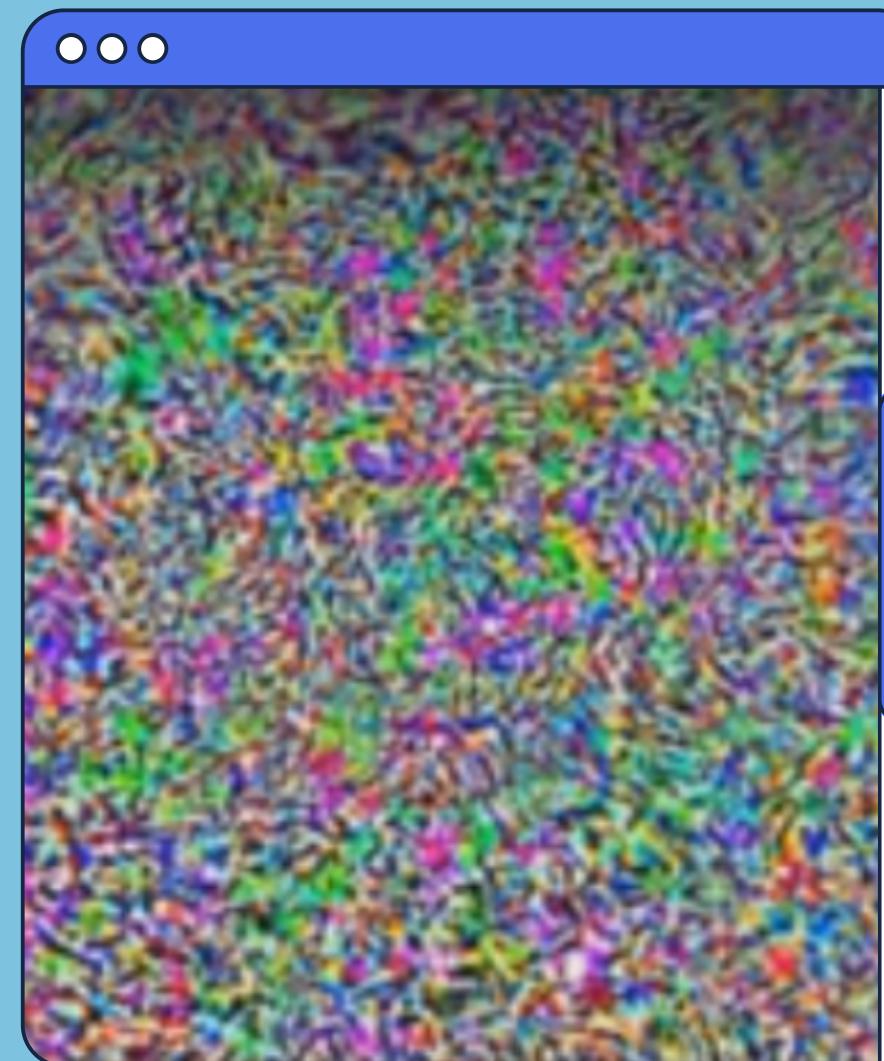
Original Image



Tiger Cat

0.64

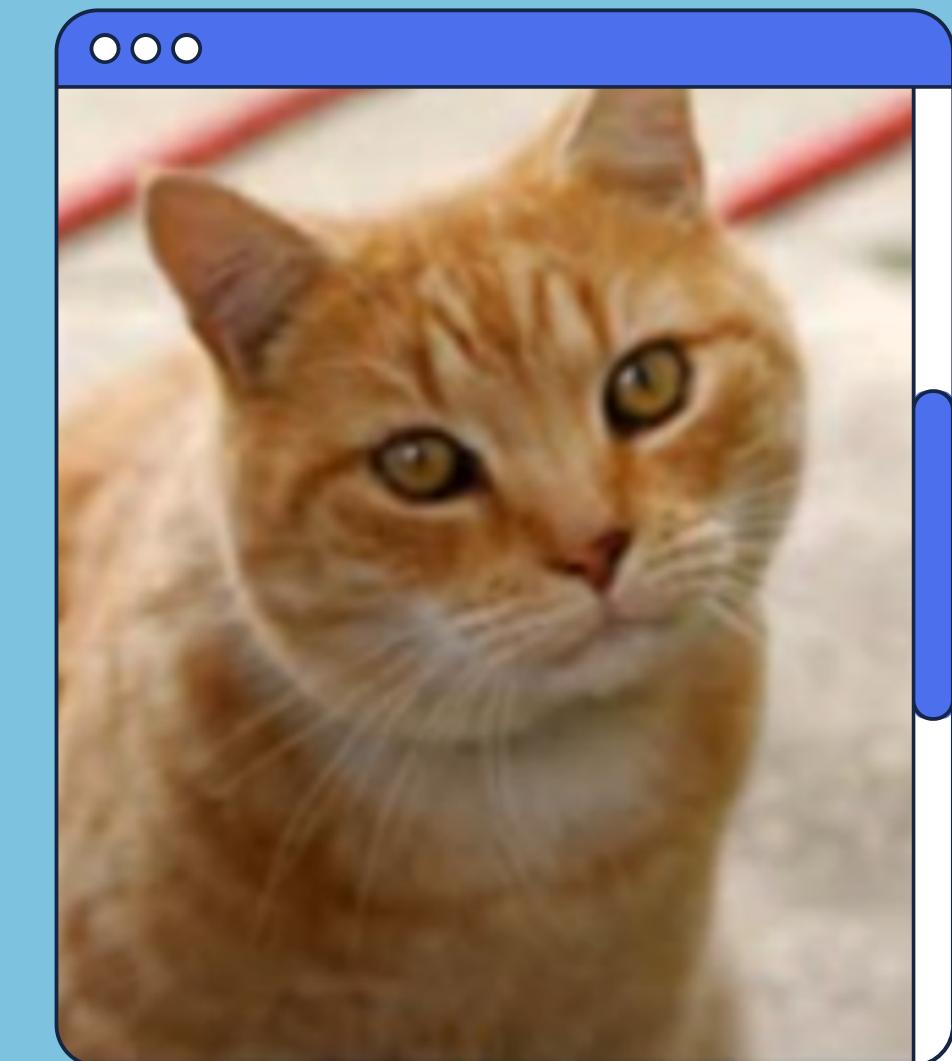
Noise



Star Fish

1.00

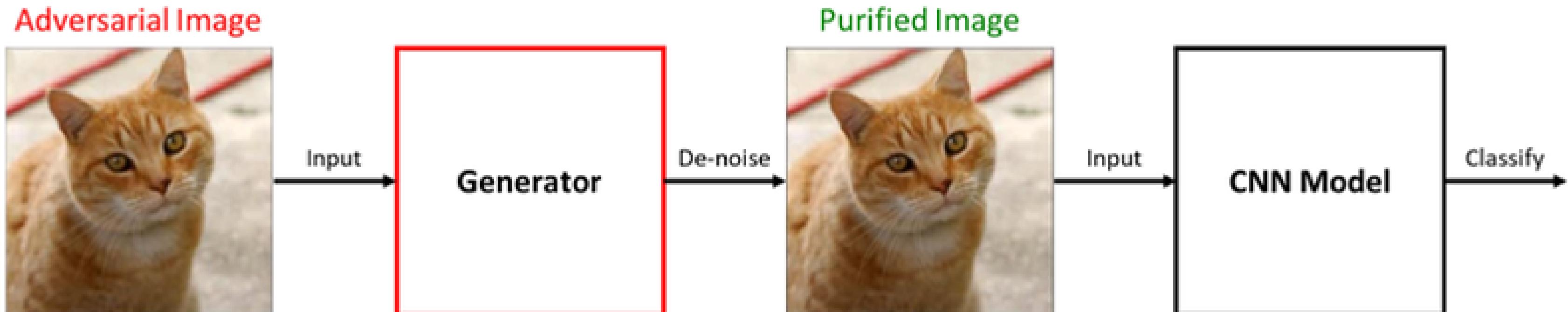
Adversarial Image



# How to Defend?

## Generative Model

- Adversarial Perturbation Elimination with GAN (APE-GAN)
- Diffusion Model

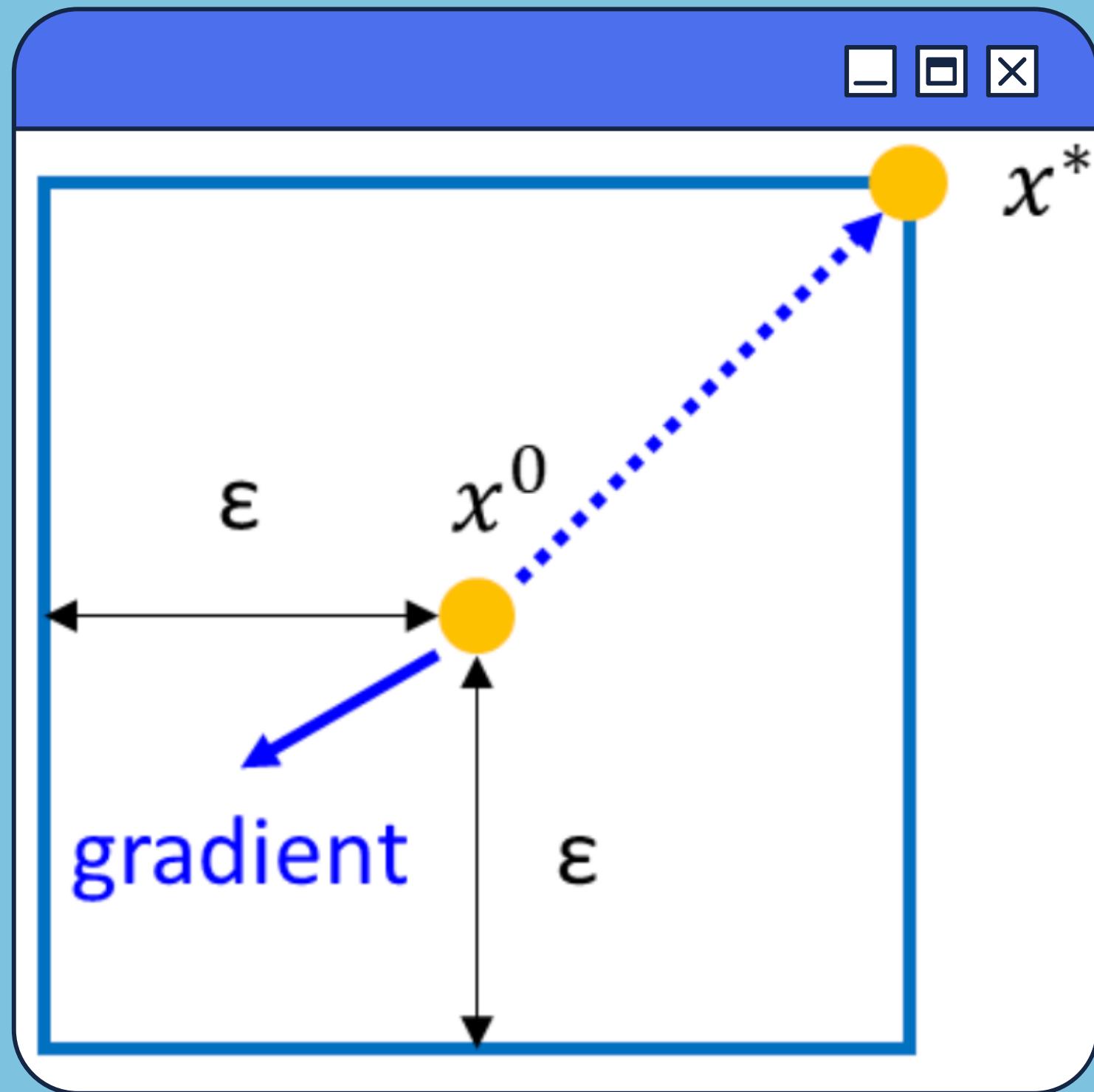




# Fast Gradient Sign Method (FGSM)



# Fast Gradient Sign Method (FGSM)

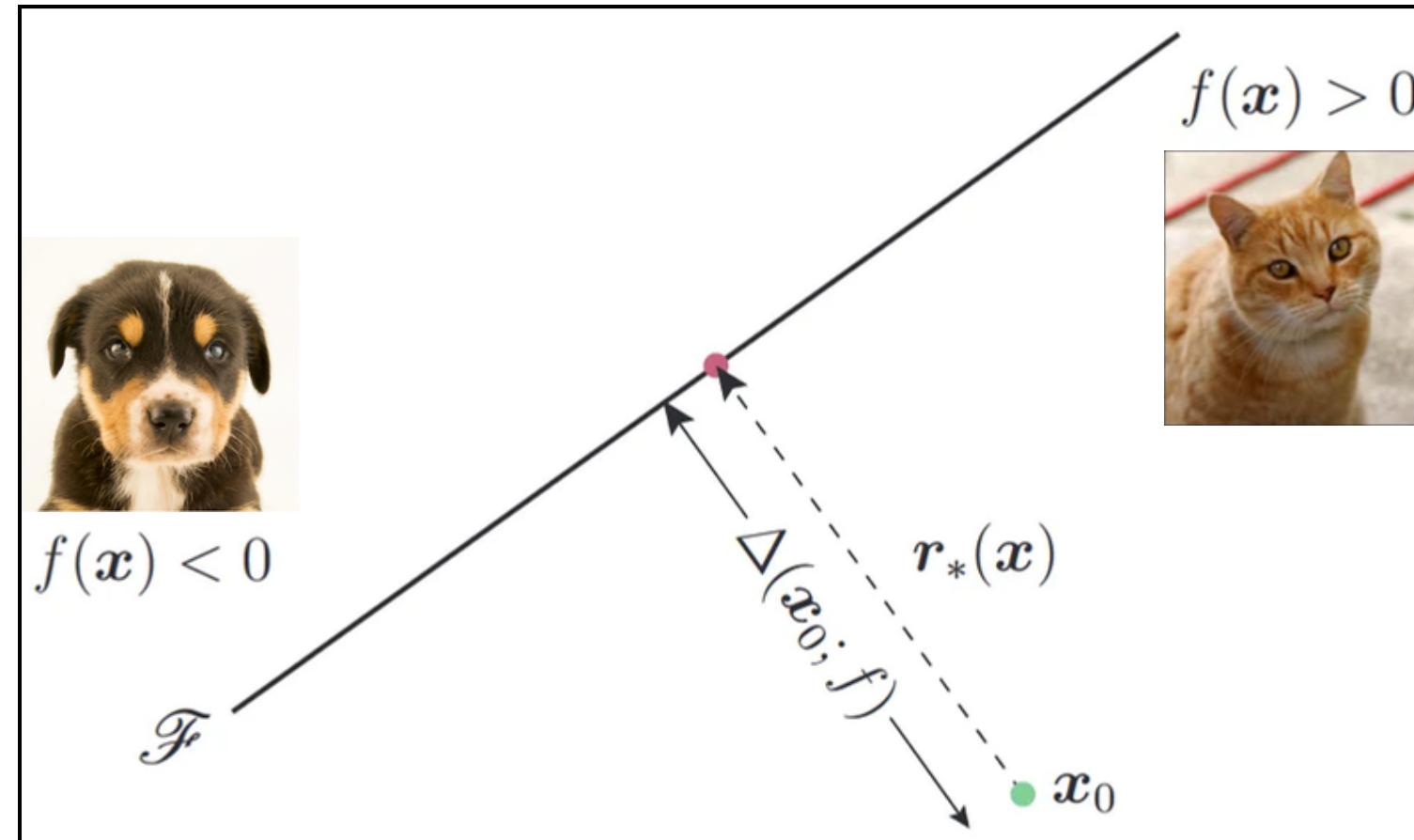


PLAY ▶



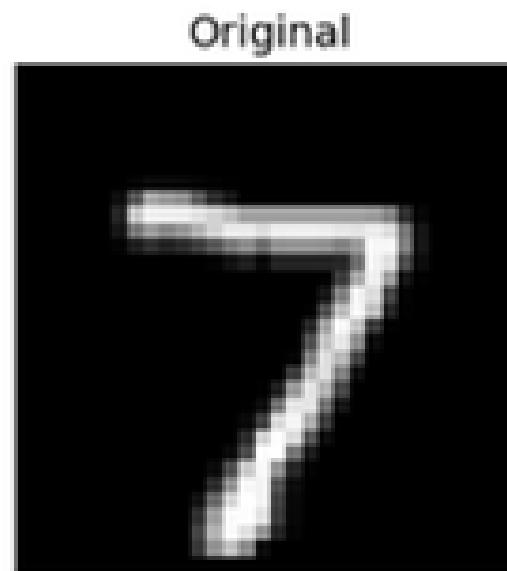
# DeepFool

- Use Binary Classifier as an example
- Obtain minimal perturbation (distance) → Iterate until label change

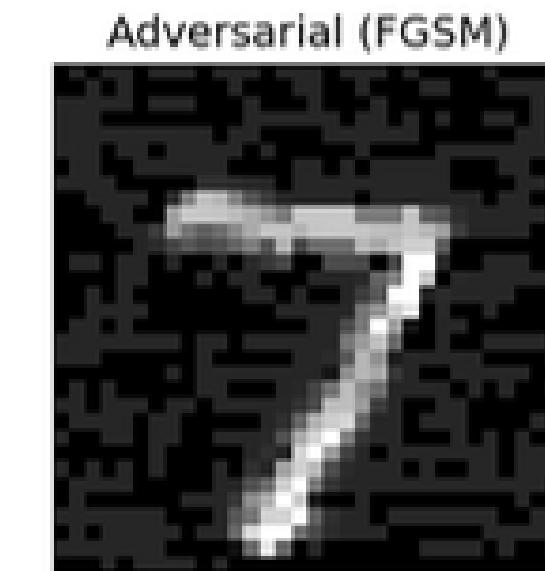
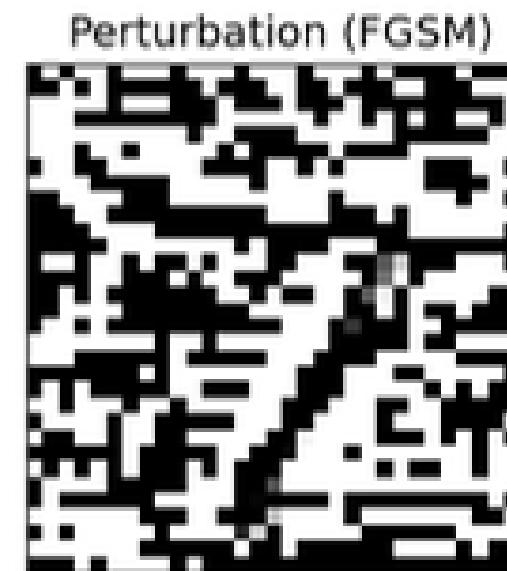


# FGSM vs DeepFool

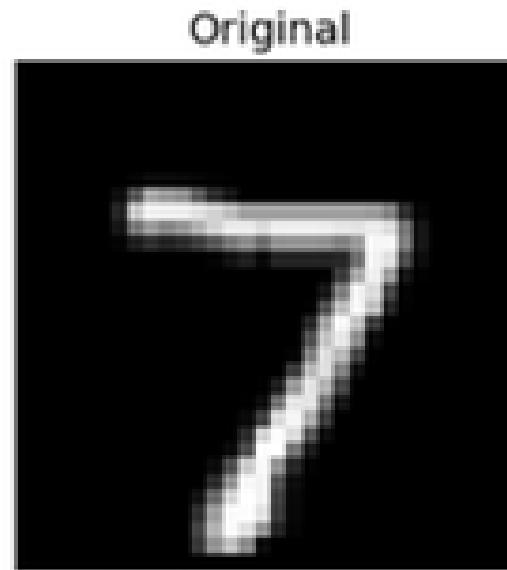
- DeepFool add less perturbation and more misclassification



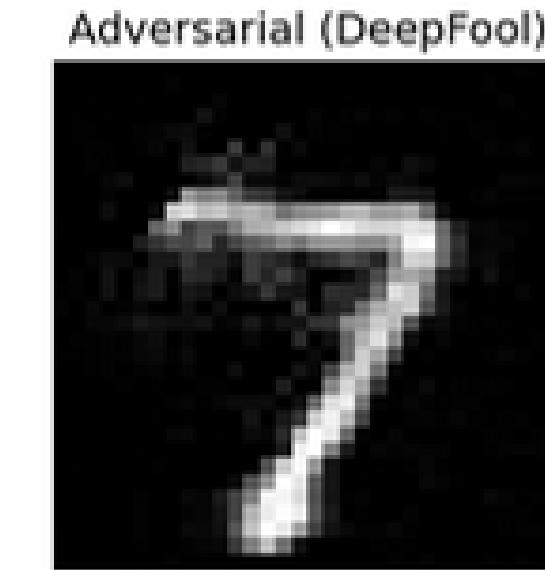
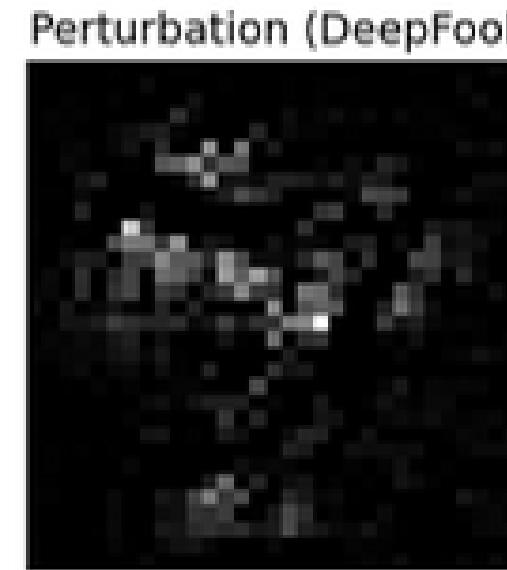
True label: 7  
Pred label: 7



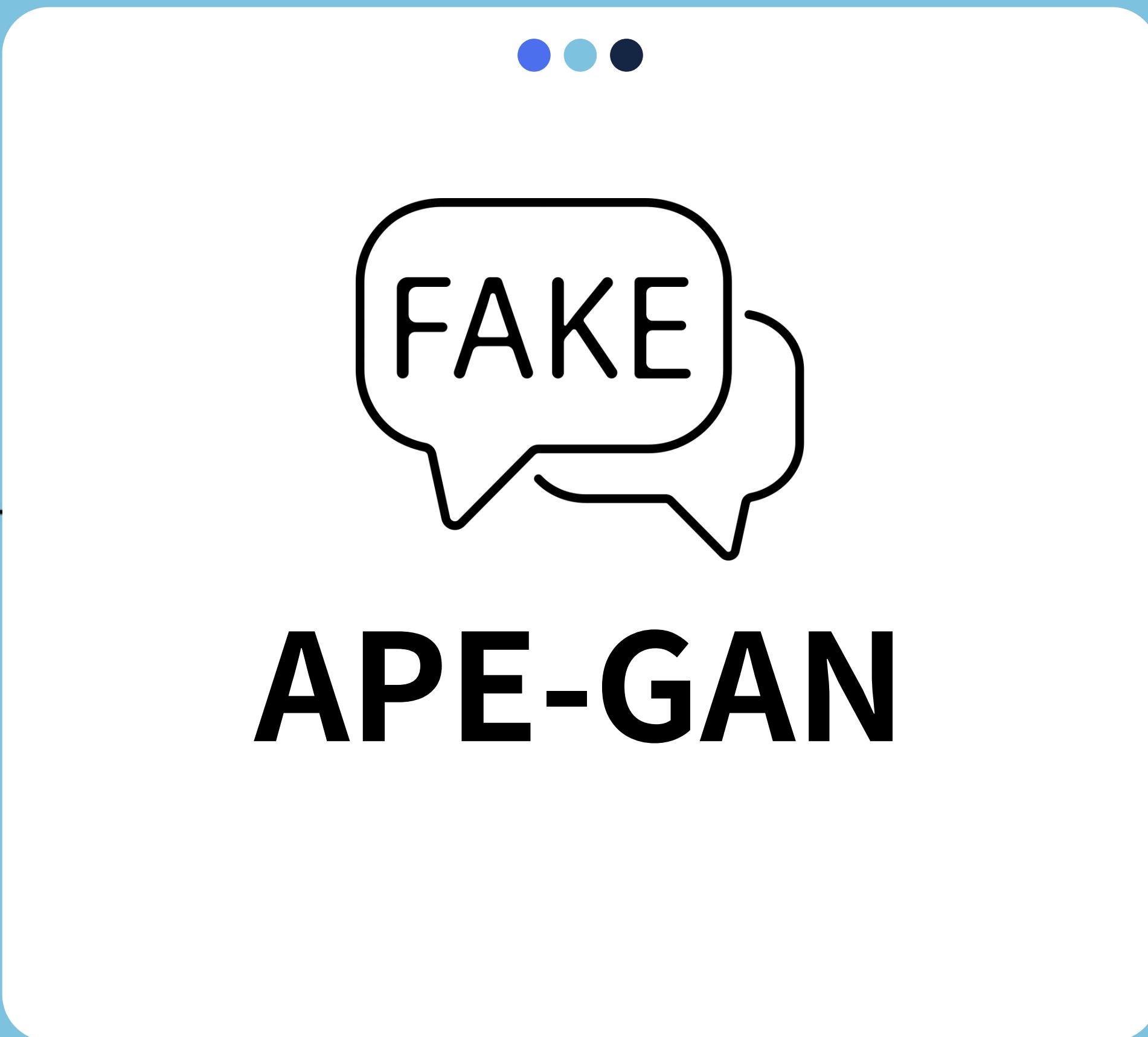
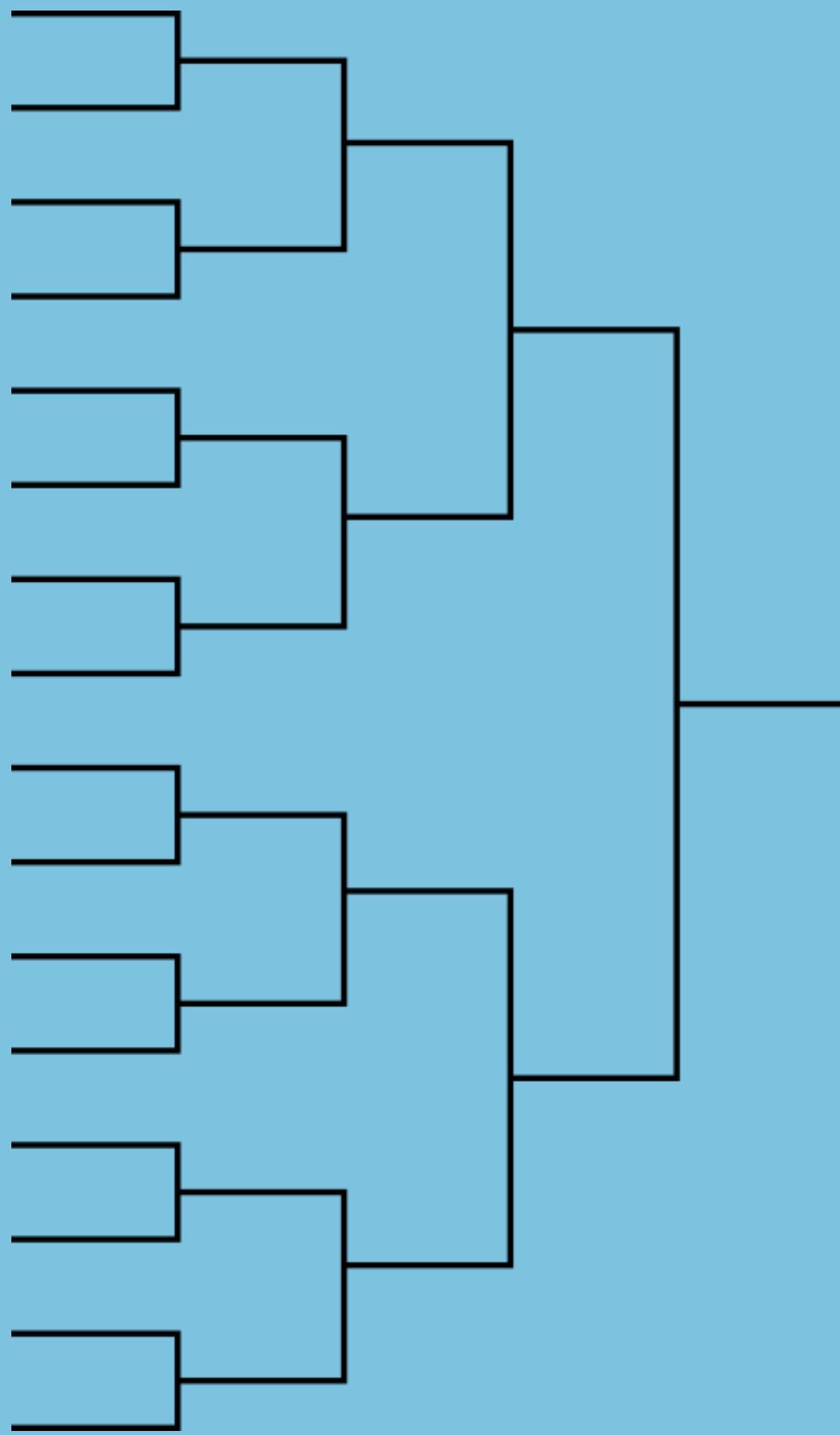
Pred label: 7



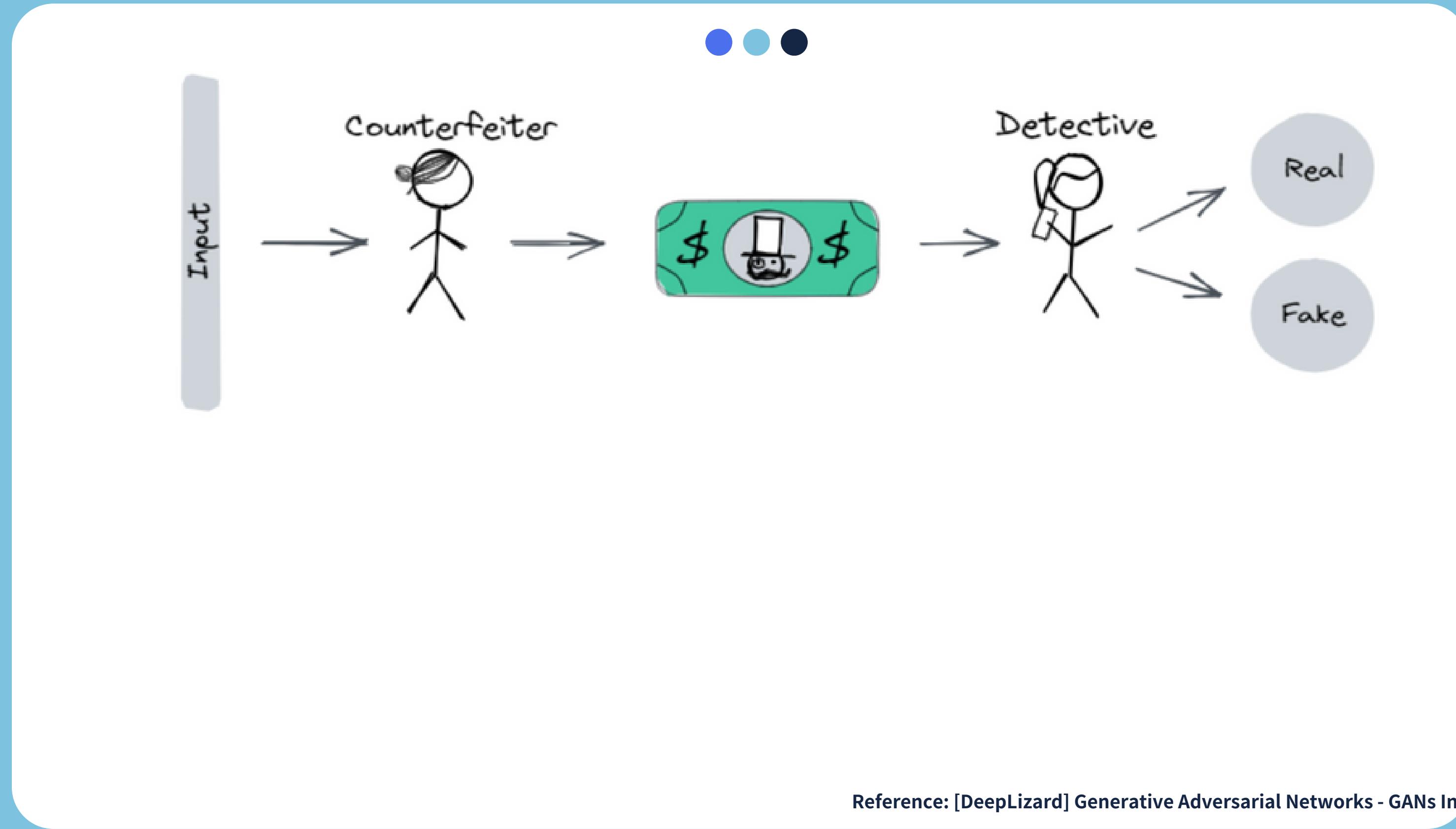
True label: 7  
Pred label: 7



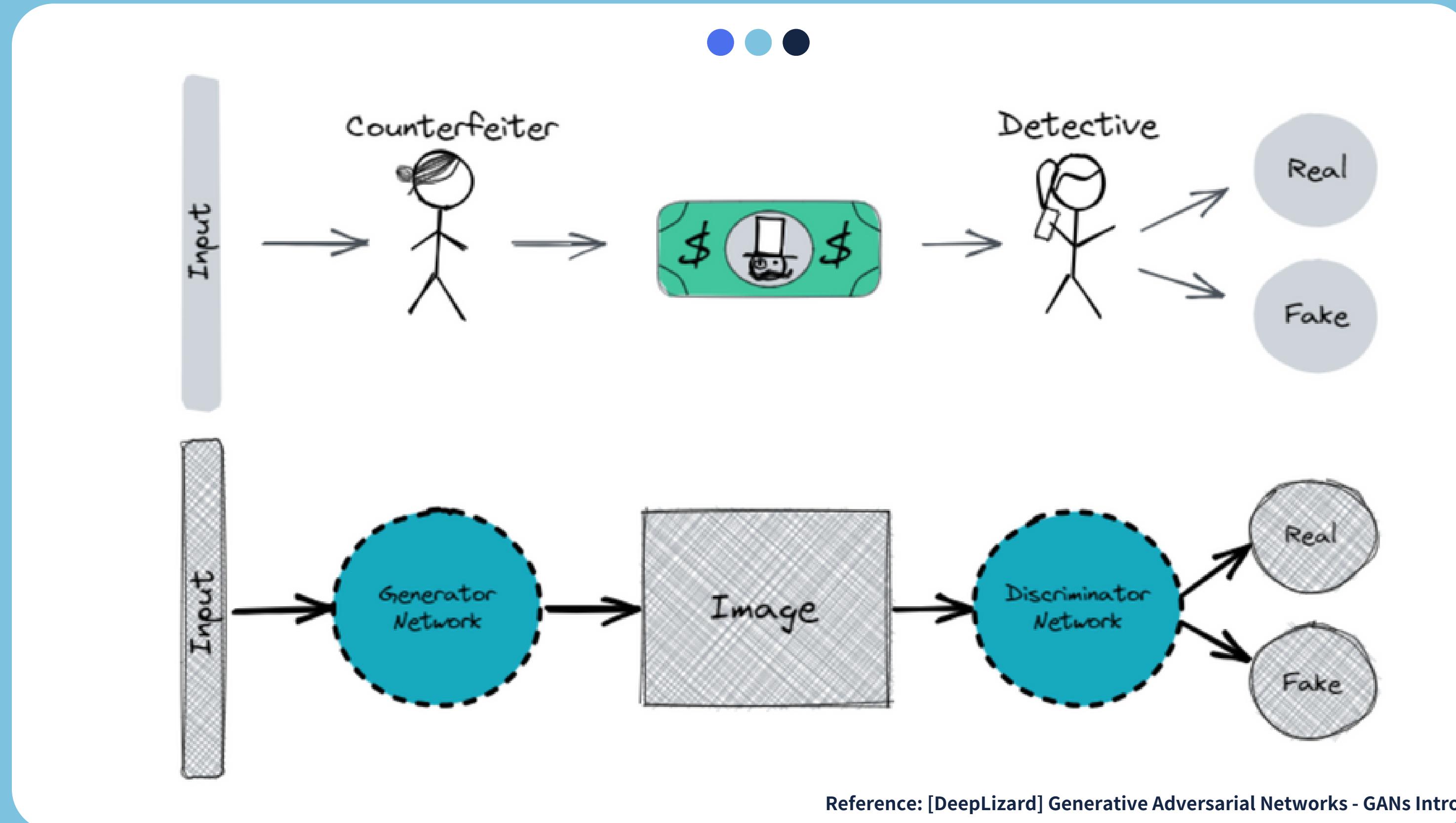
Pred label: 9



# Generative Adversarial Network (GAN)



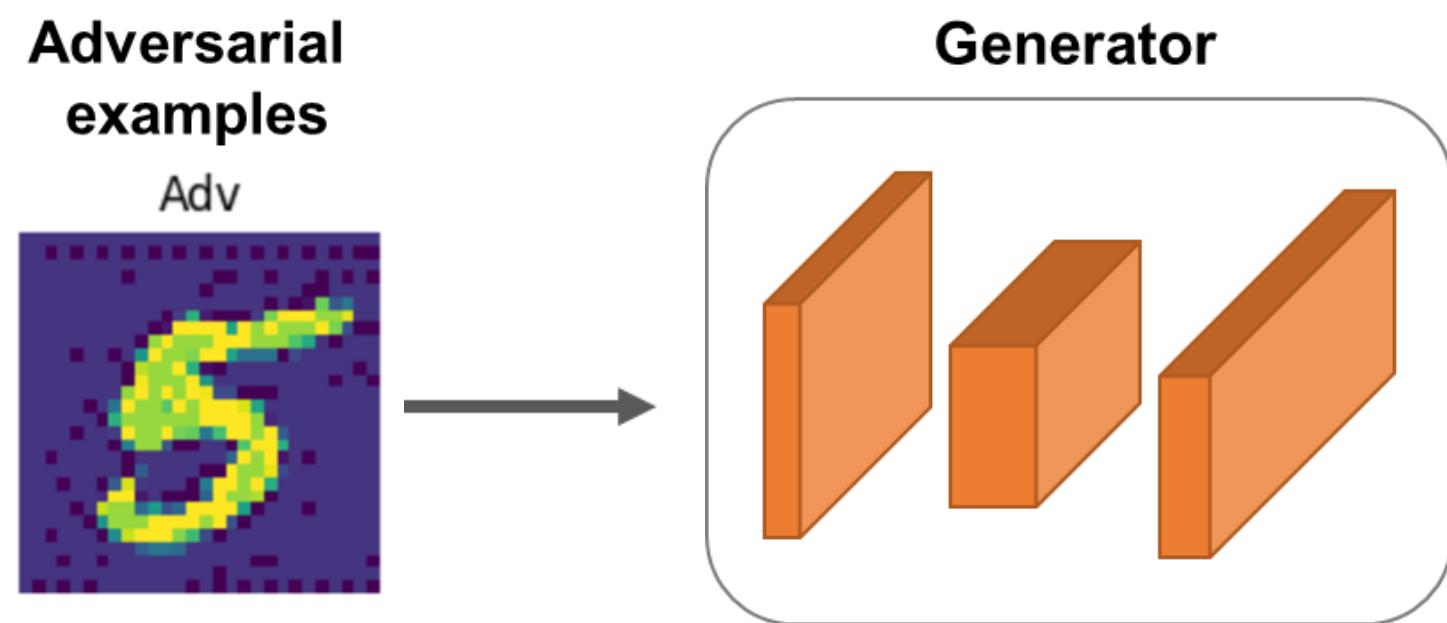
# Generative Adversarial Network (GAN)



# APE-GAN



- Build Generator



**Loss Function**

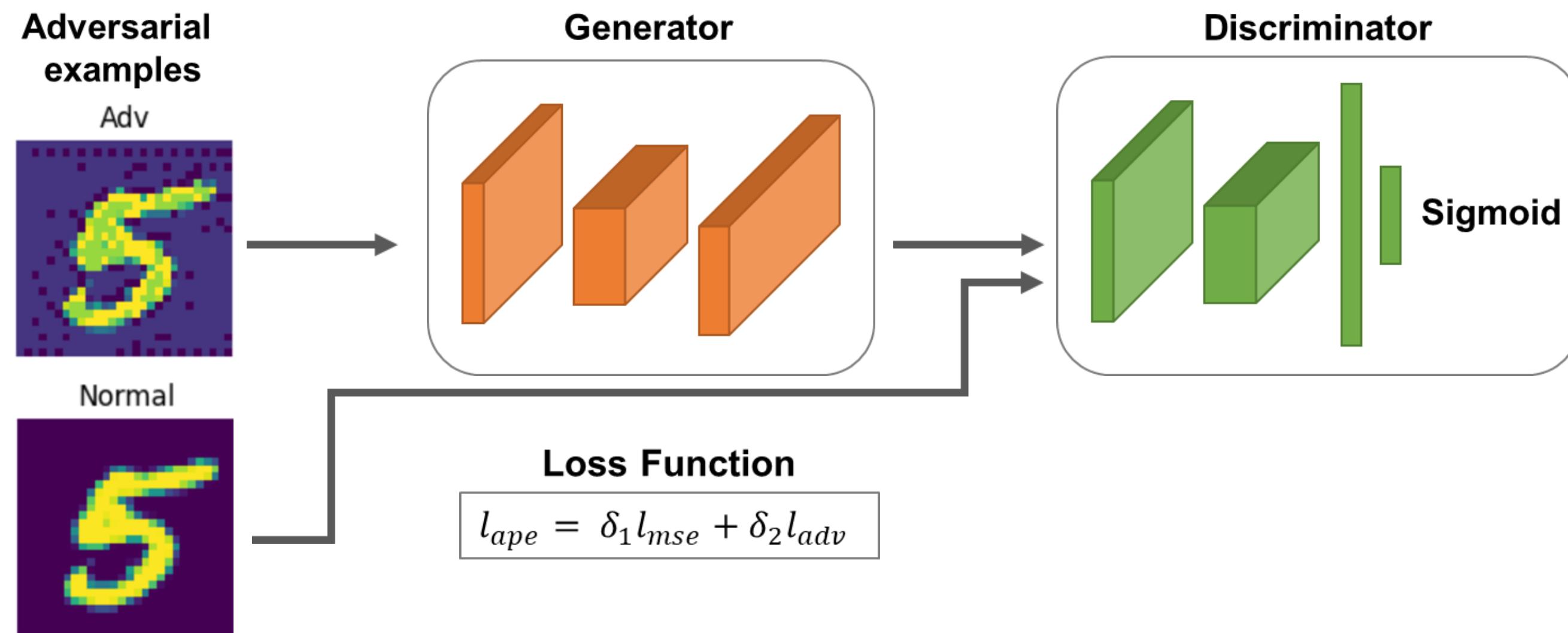
$$l_{ape} = \delta_1 l_{mse} + \delta_2 l_{adv}$$



# APE-GAN



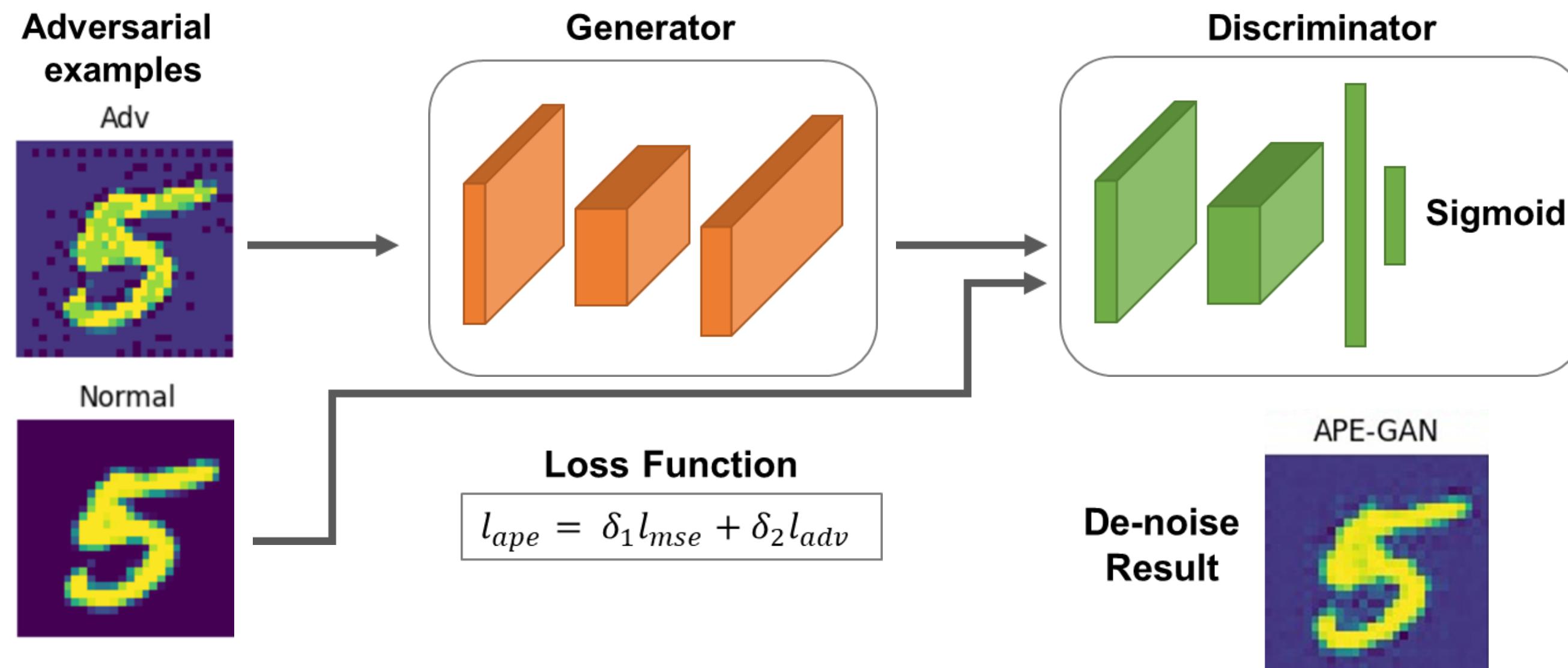
- Build Discriminator

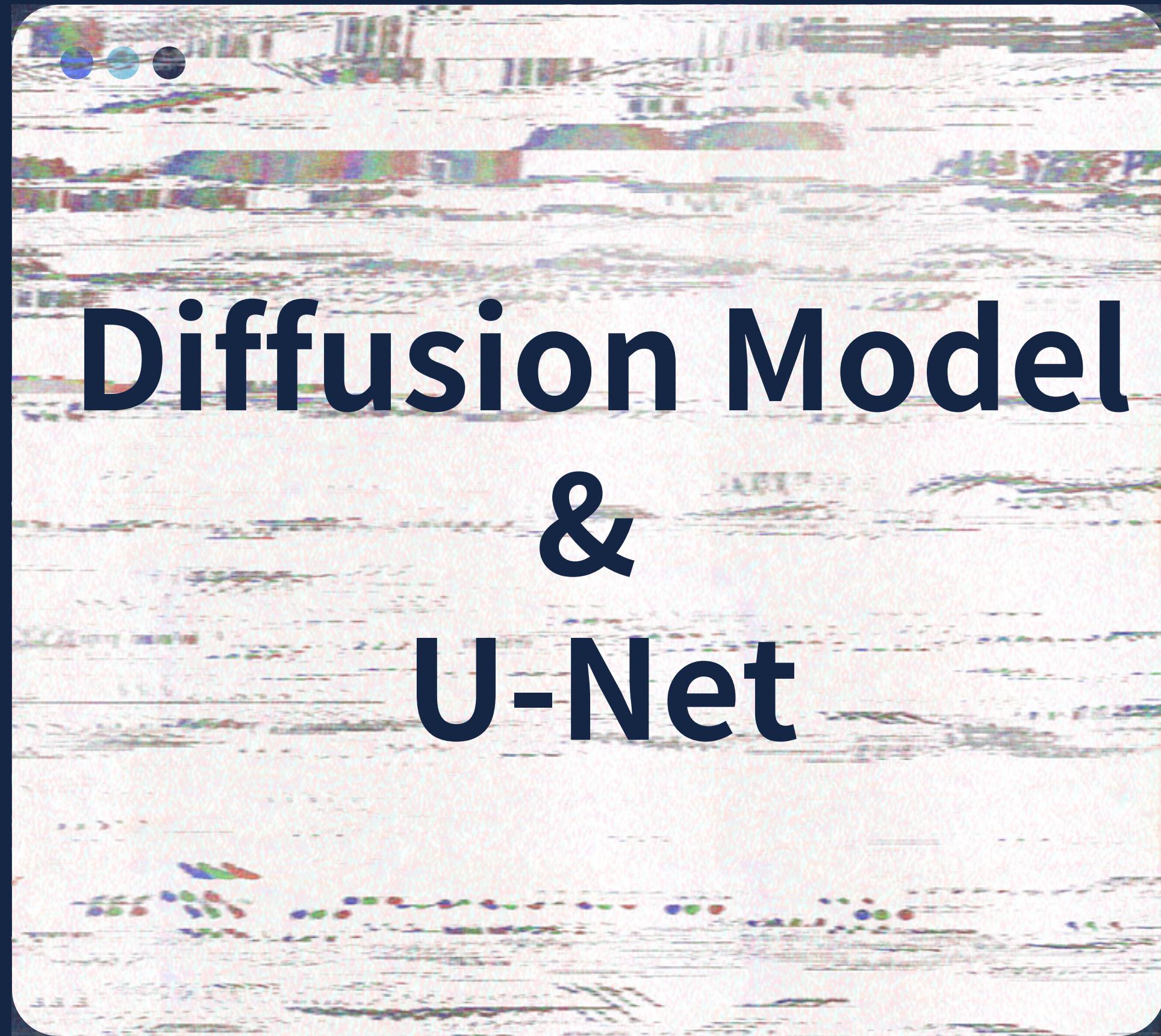


# APE-GAN



- Train generator twice + discriminator one time in each epoch





# Diffusion Model



## Forward Diffusion Process



Data

Noise



# Diffusion Model



**Reverse Denoise Process(Generative)**

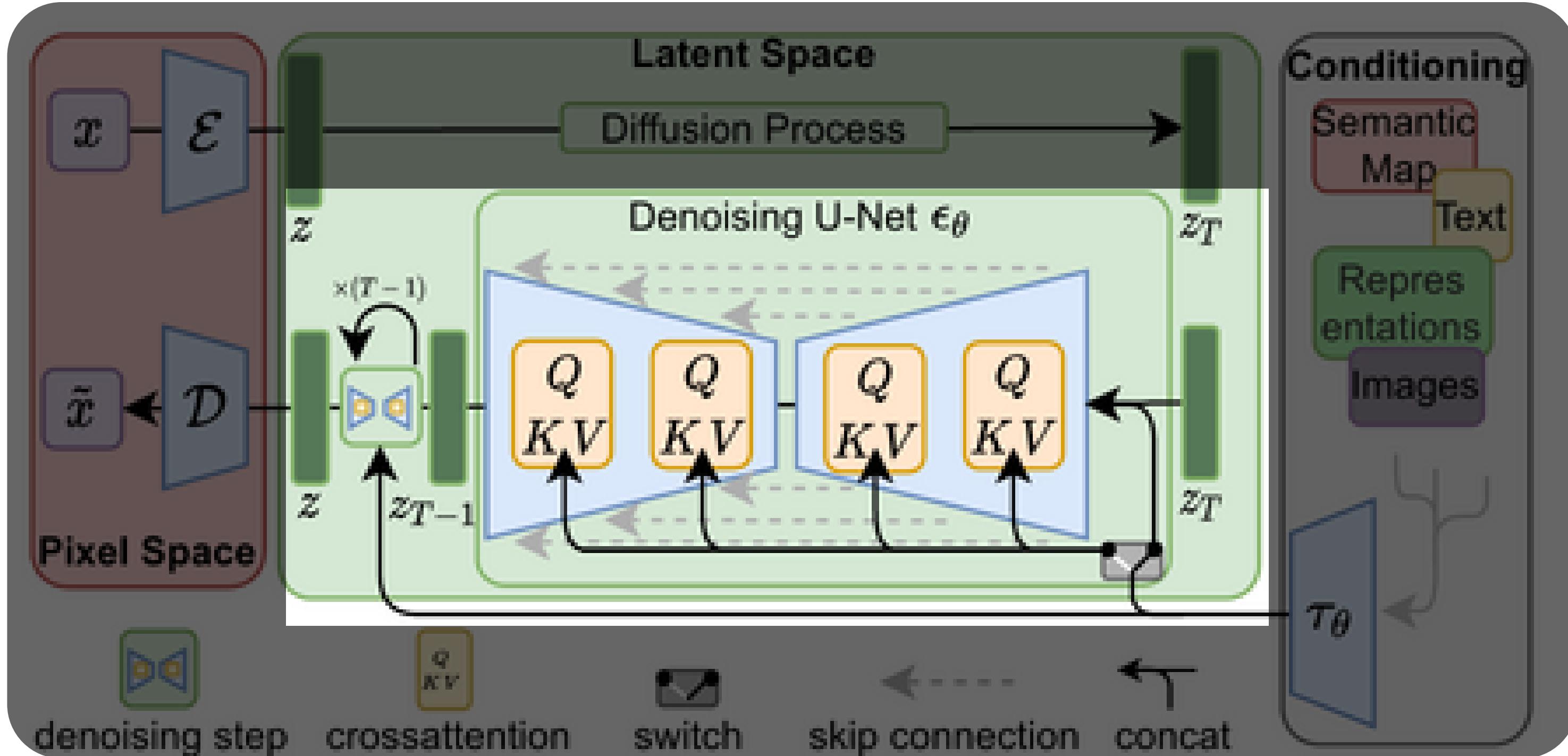


**Data**

**Noise**



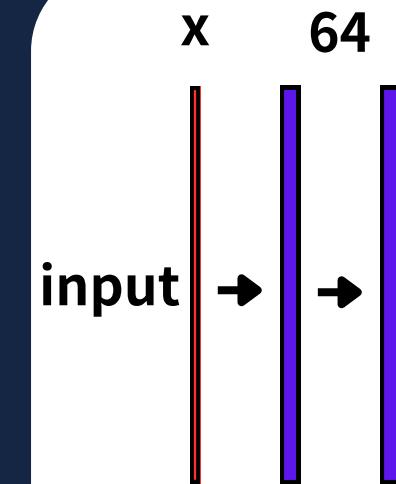
# Diffusion Architecture



Reference: <https://arxiv.org/abs/2112.10752>

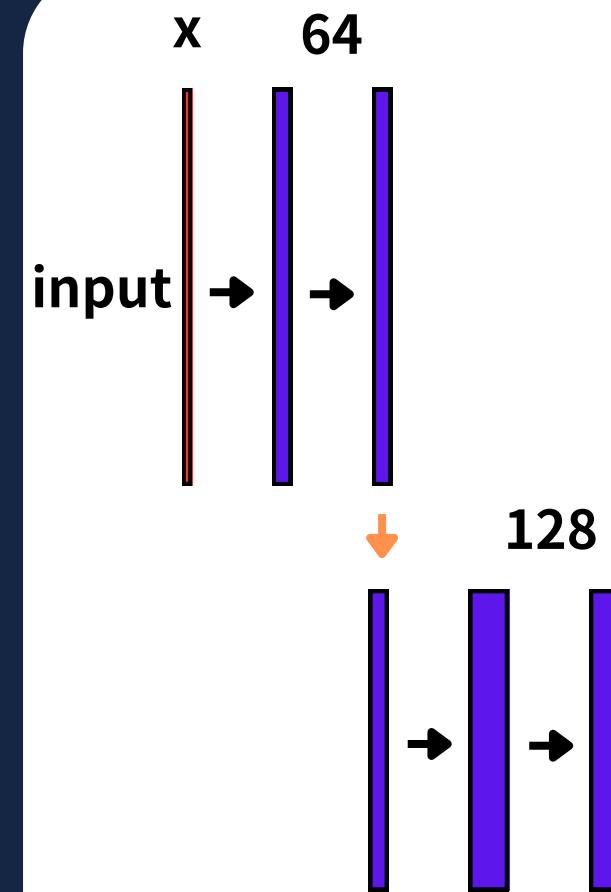


# U-Net Architecture



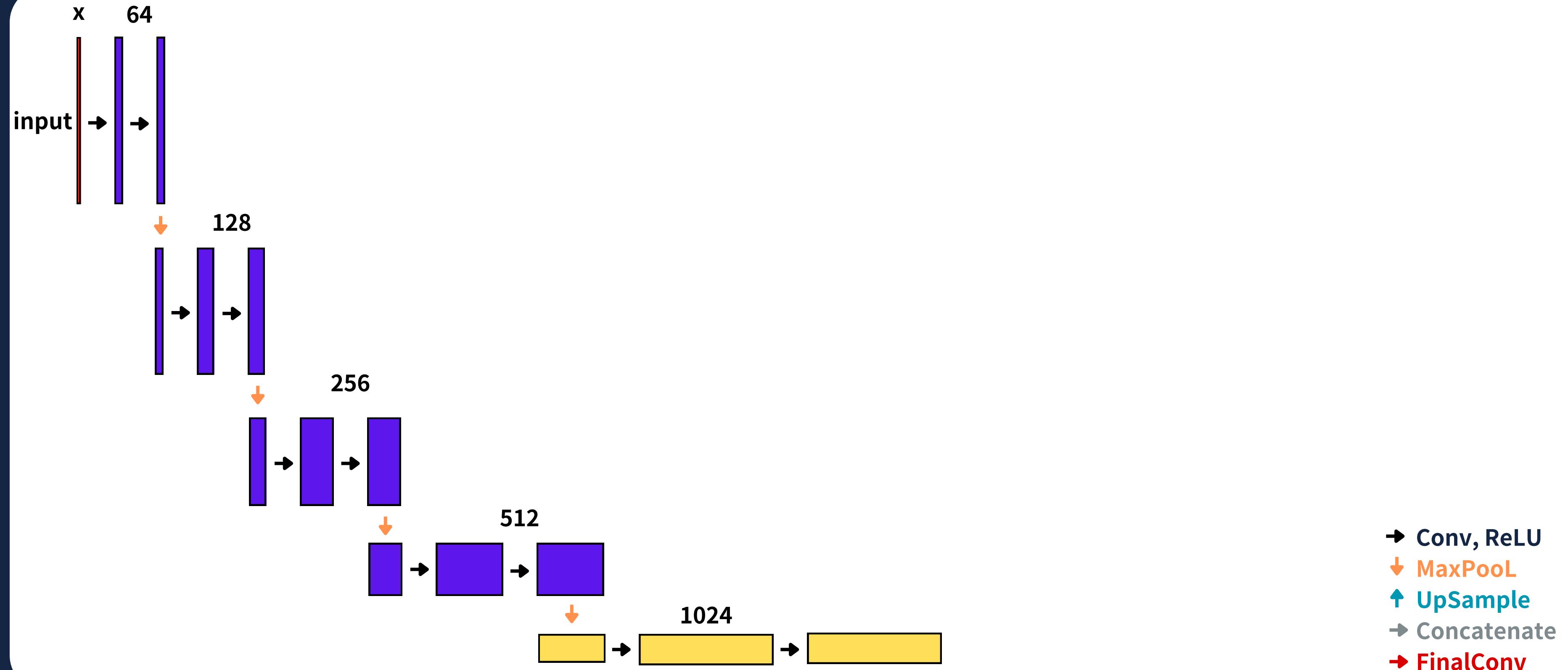
- Conv, ReLU
- ↓ MaxPool
- ↑ UpSample
- Concatenate
- FinalConv

# U-Net Architecture

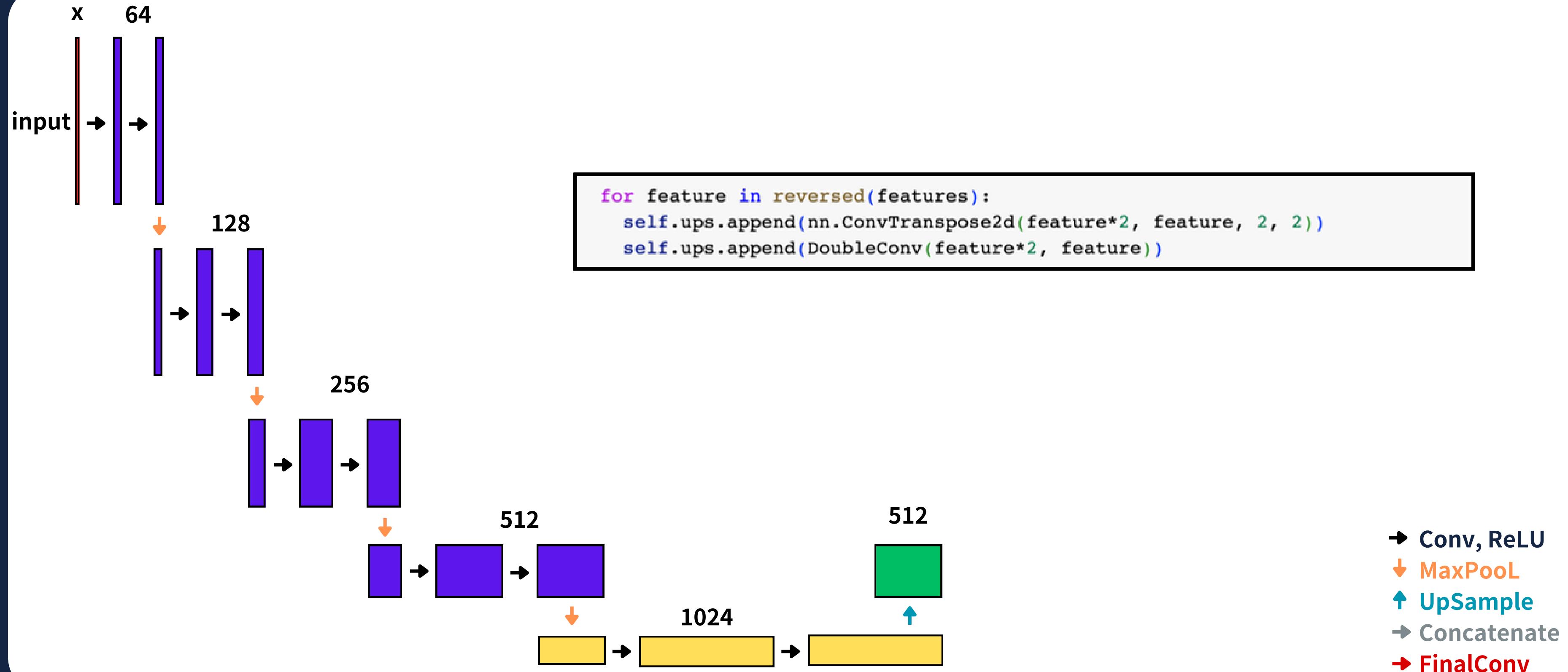


- Conv, ReLU
- ↓ MaxPool
- ↑ UpSample
- Concatenate
- FinalConv

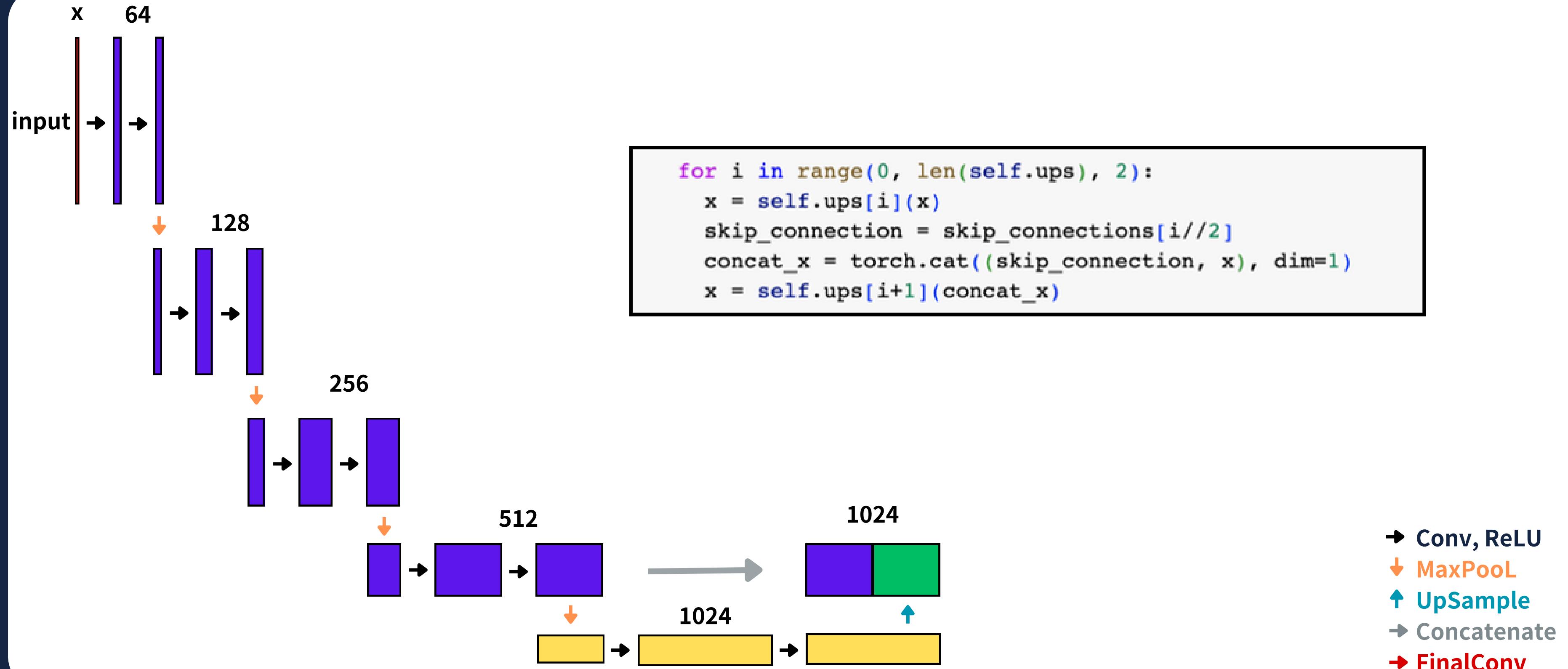
# U-Net Architecture



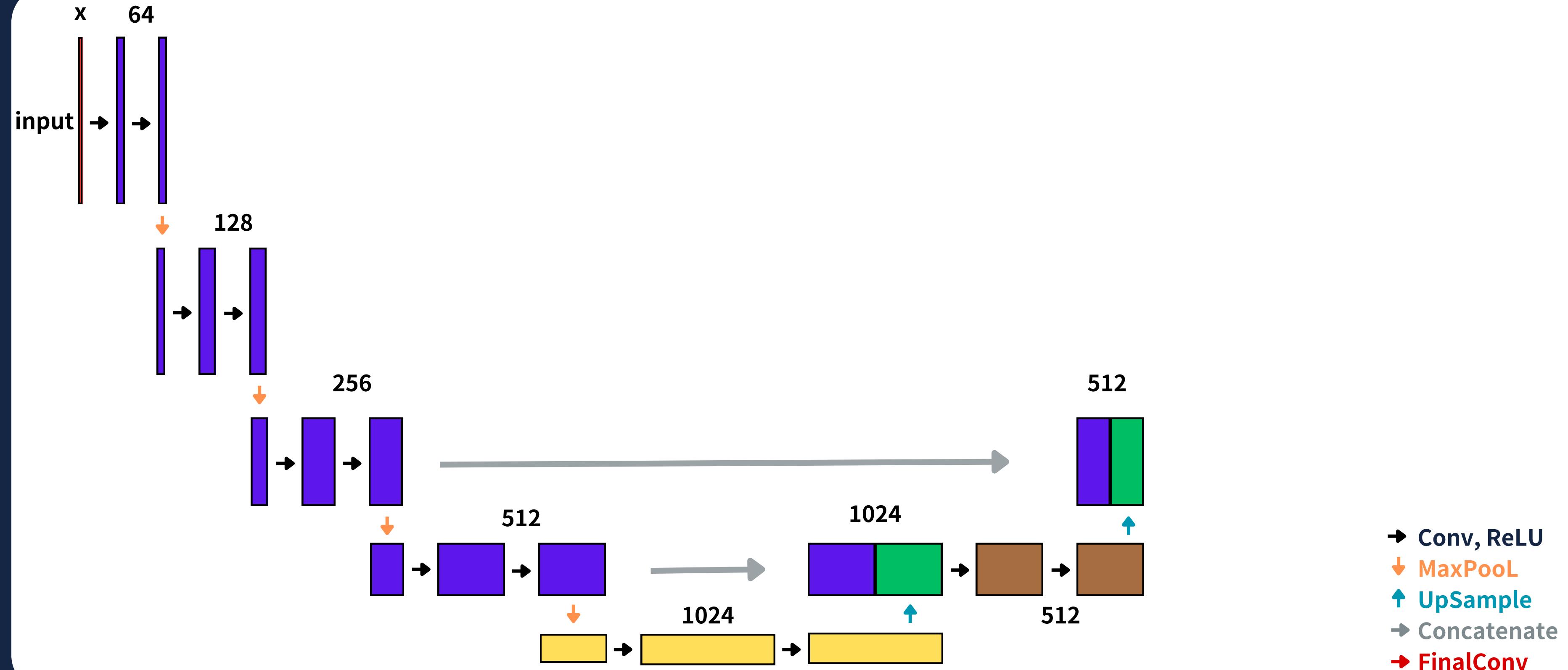
# U-Net Architecture



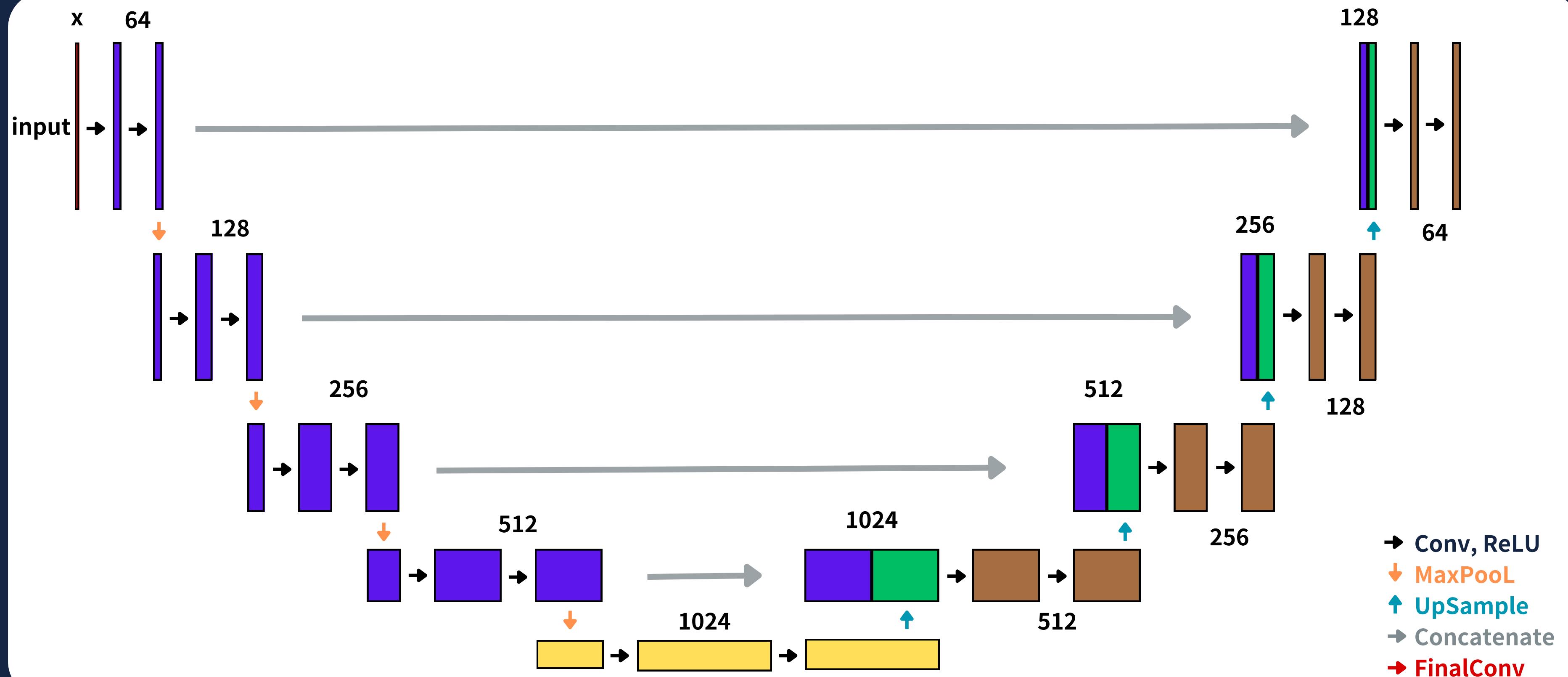
# U-Net Architecture



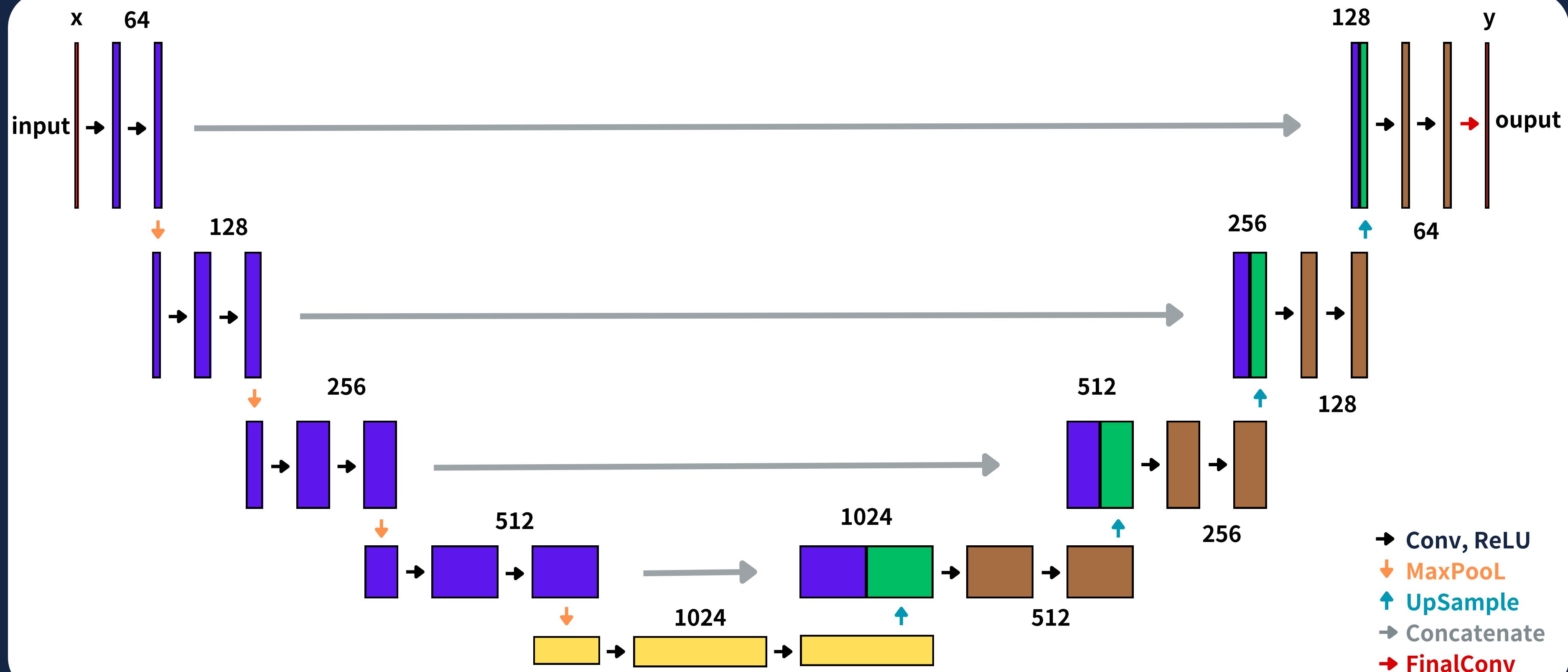
# U-Net Architecture



# U-Net Architecture



# U-Net Architecture

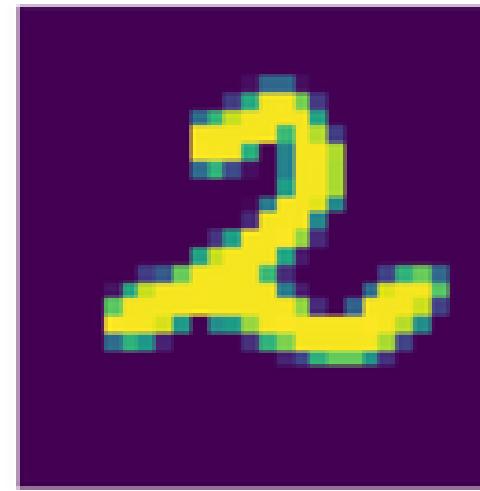




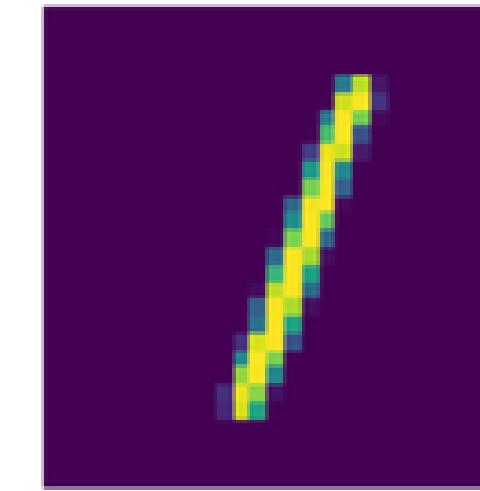
# Results of APE-GAN Defense



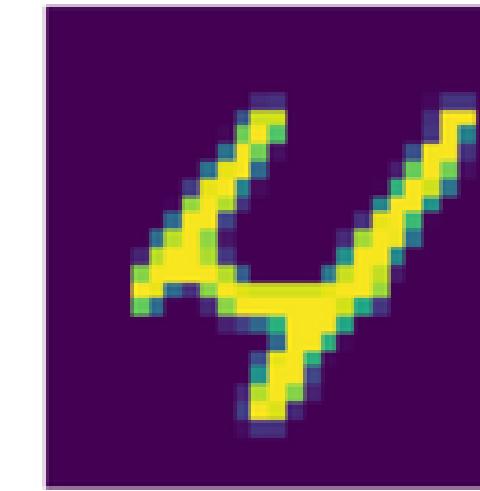
Normal



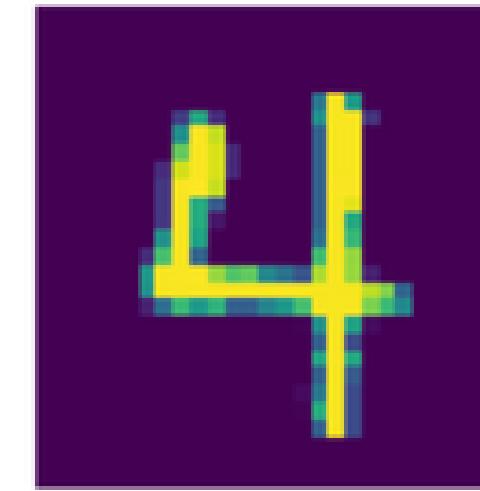
Normal



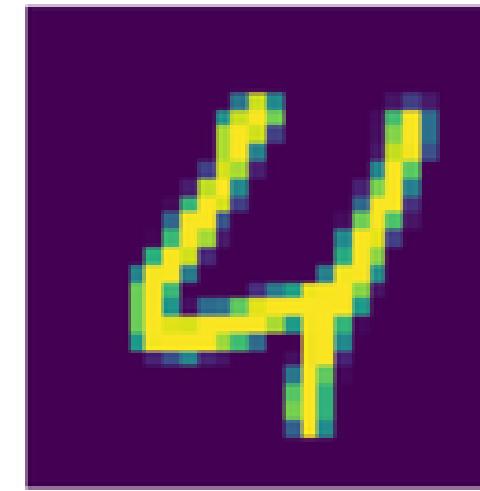
Normal



Normal

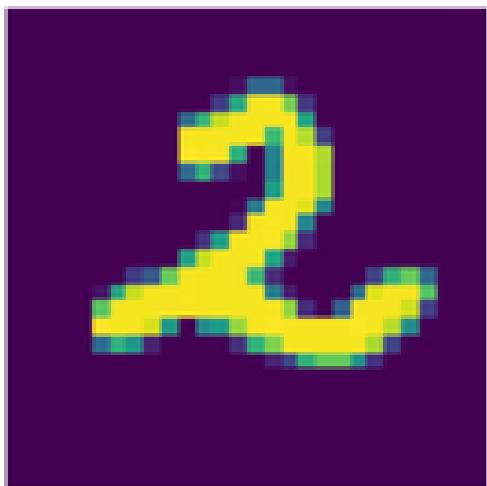


Normal

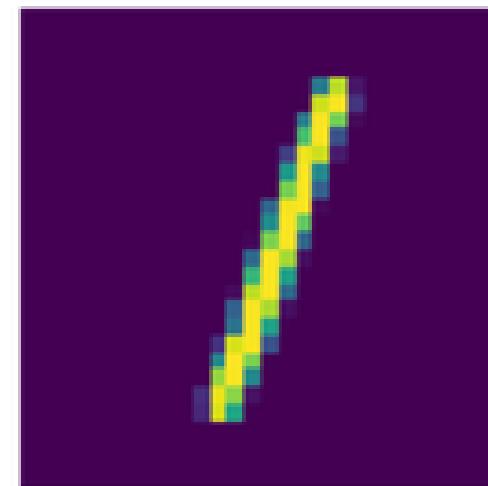




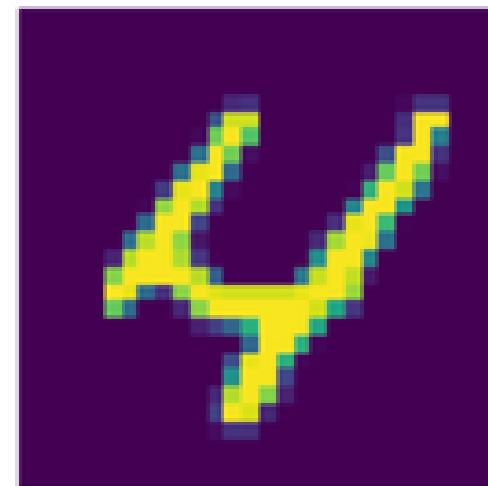
Normal



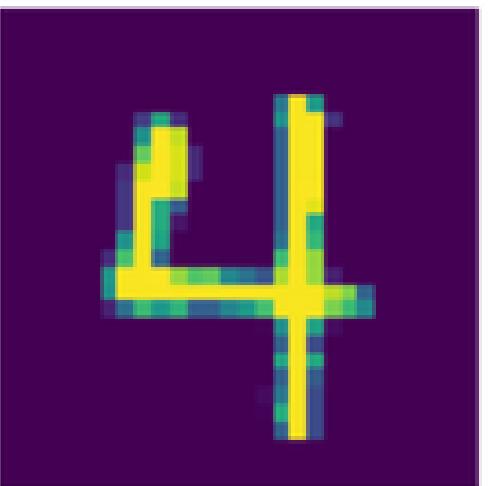
Normal



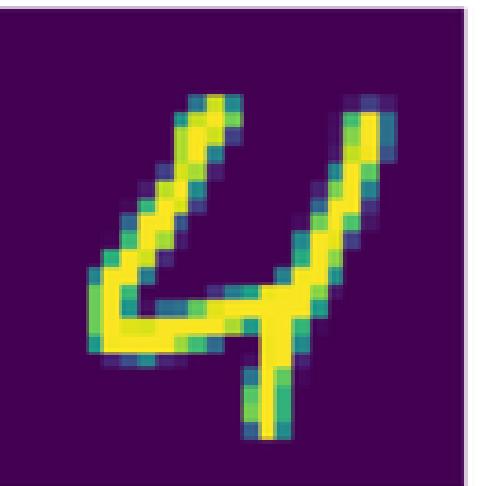
Normal



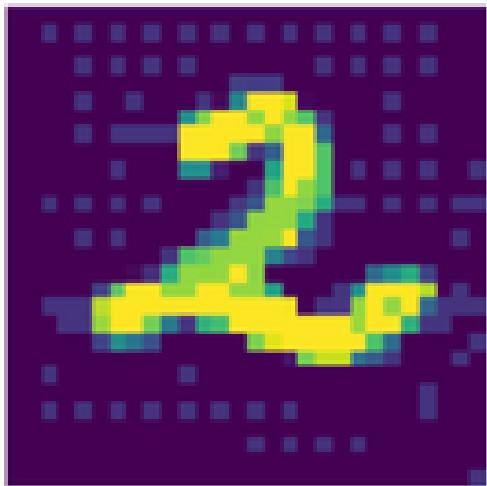
Normal



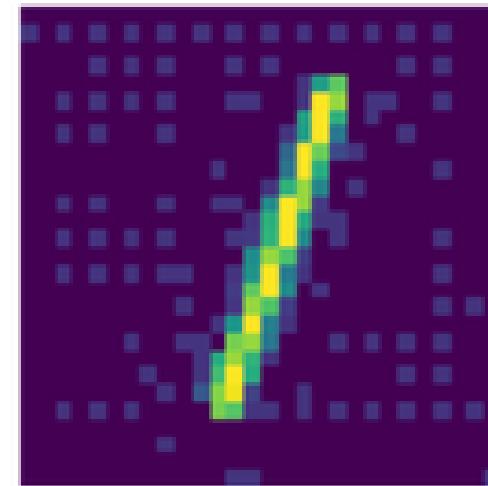
Normal



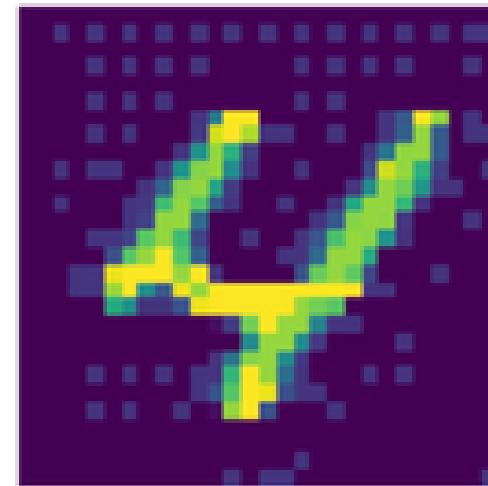
Adv



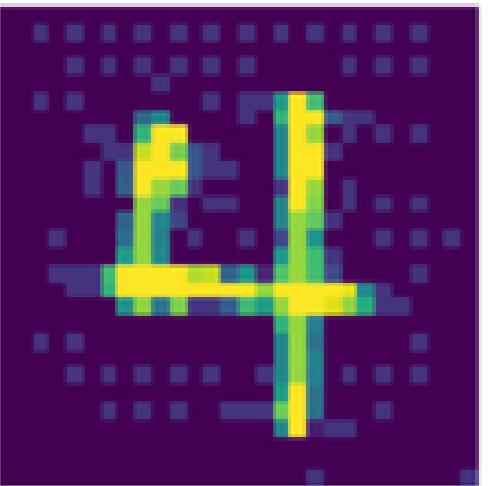
Adv



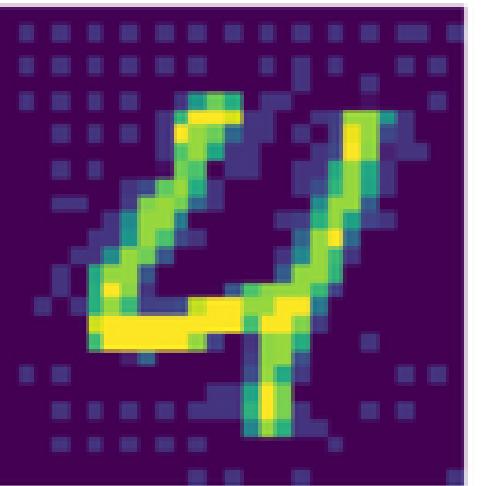
Adv



Adv

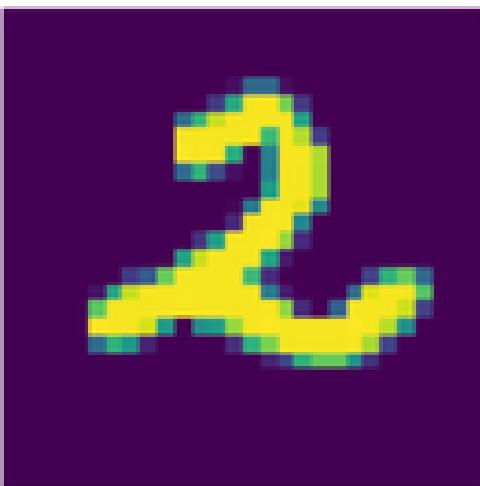


Adv

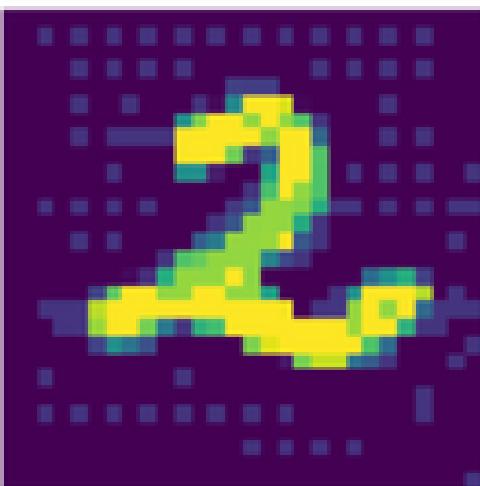




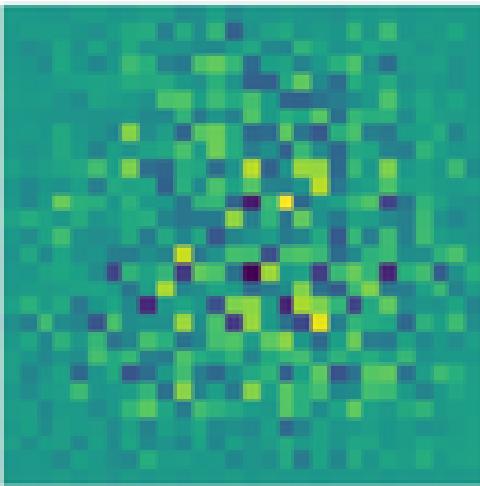
Normal



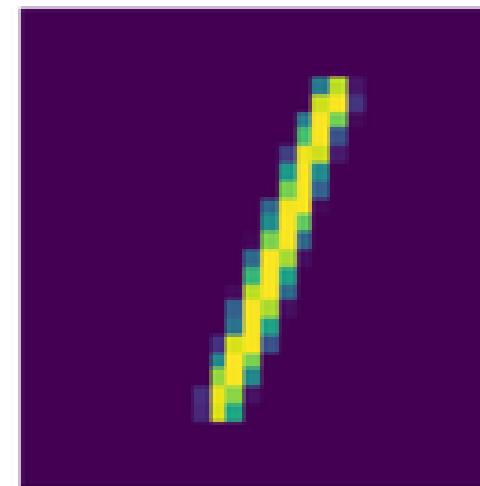
Adv



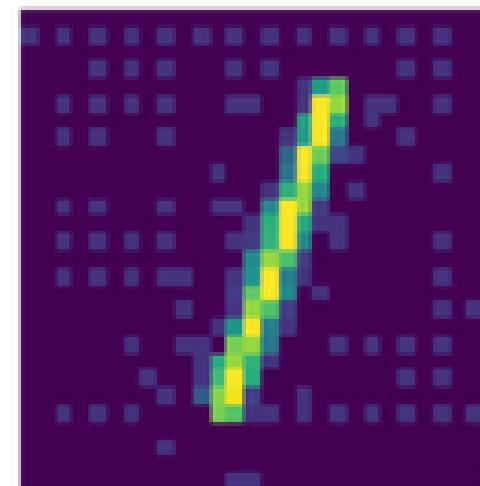
APE-GAN



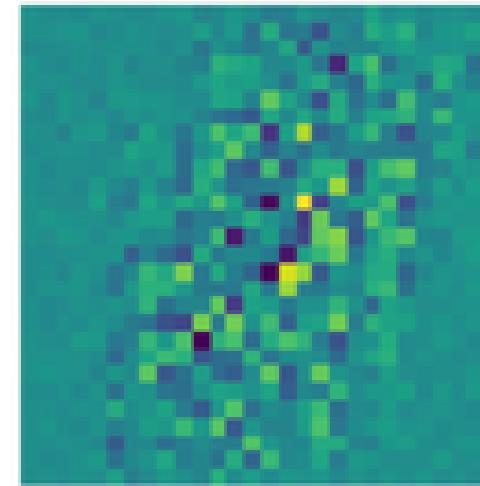
Normal



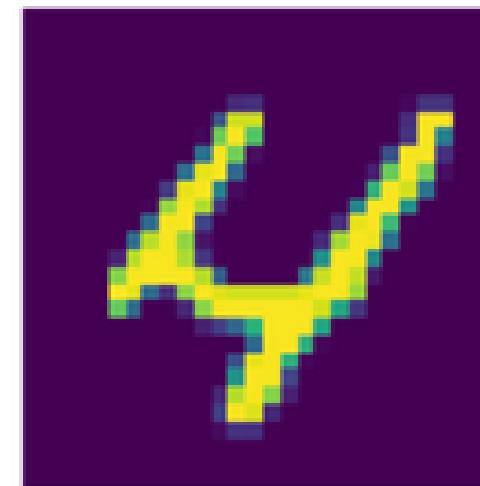
Adv



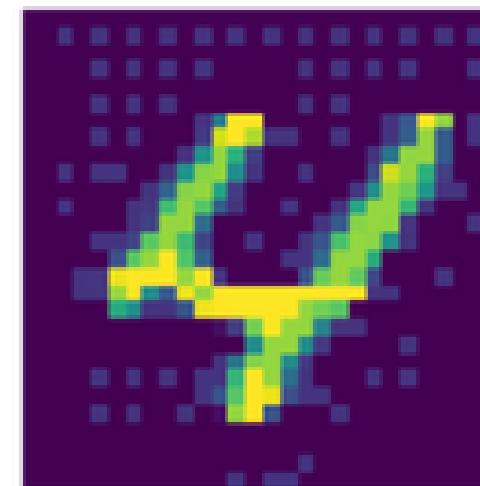
APE-GAN



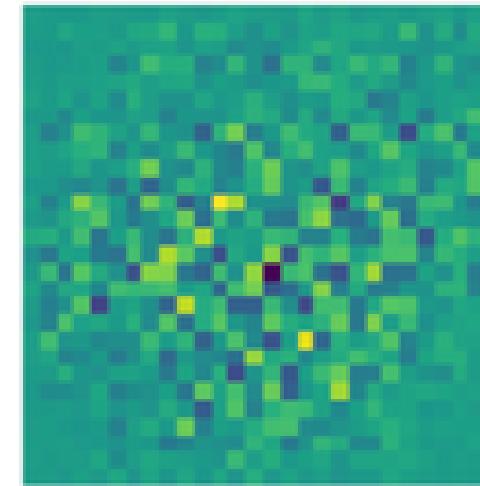
Normal



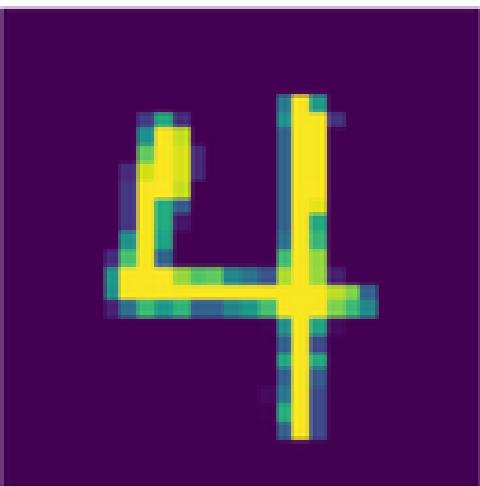
Adv



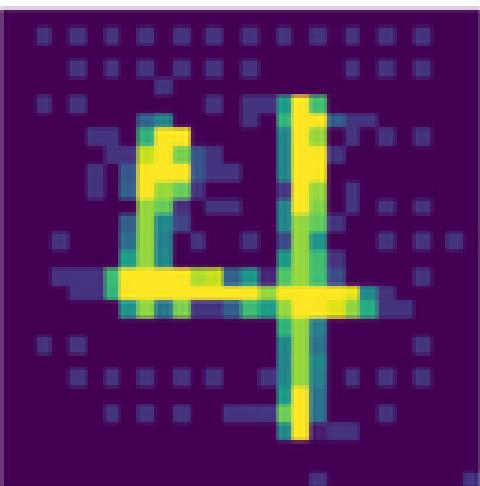
APE-GAN



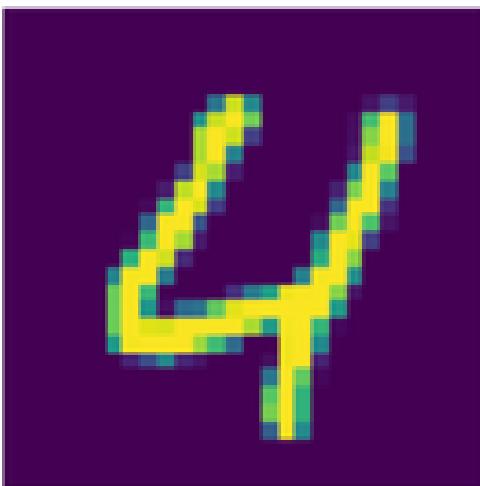
Normal



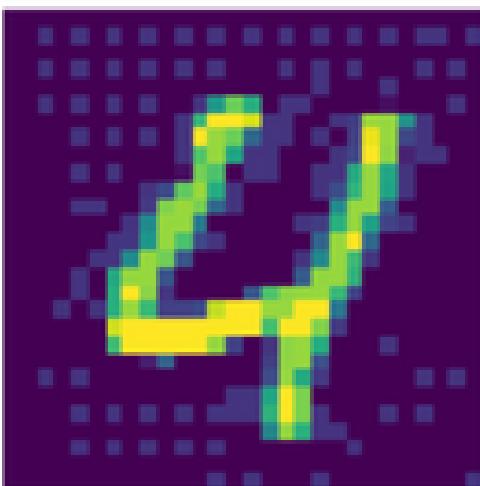
Adv



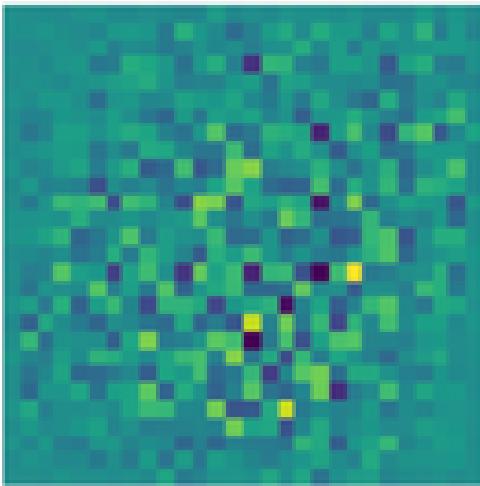
Normal



Adv

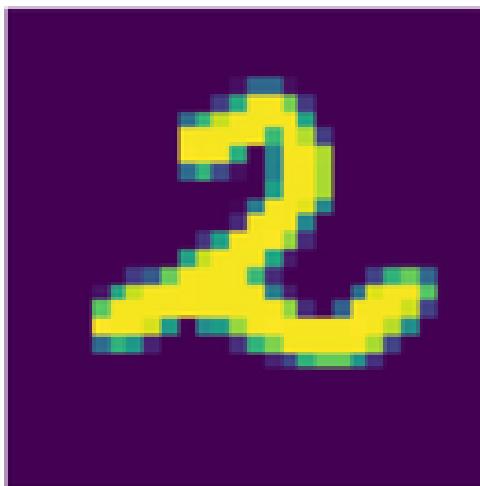


APE-GAN

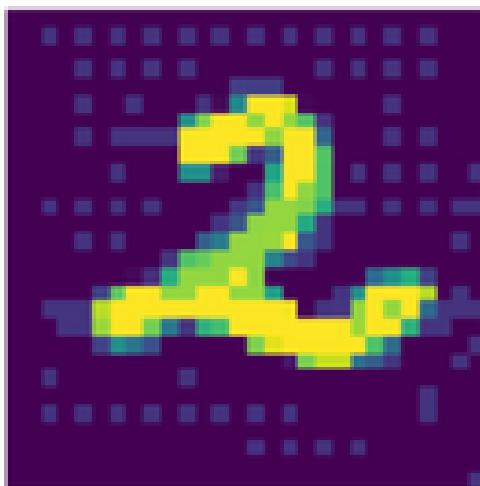




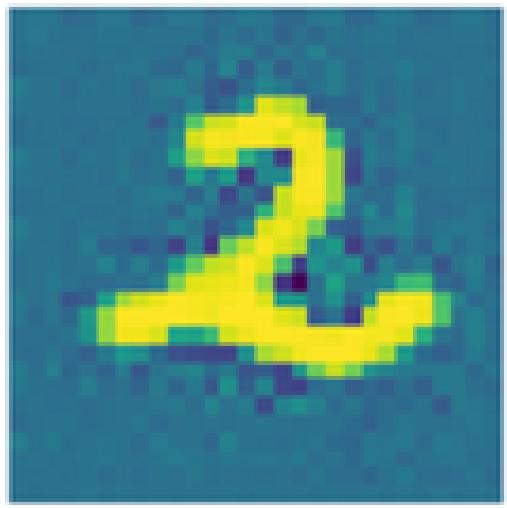
Normal



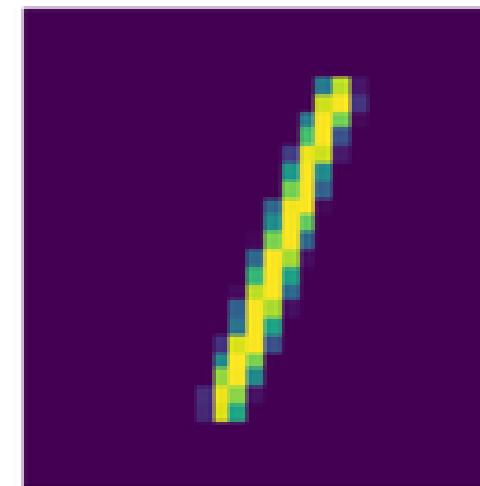
Adv



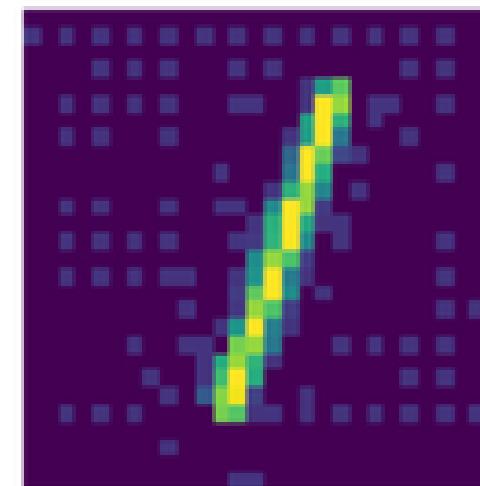
APE-GAN



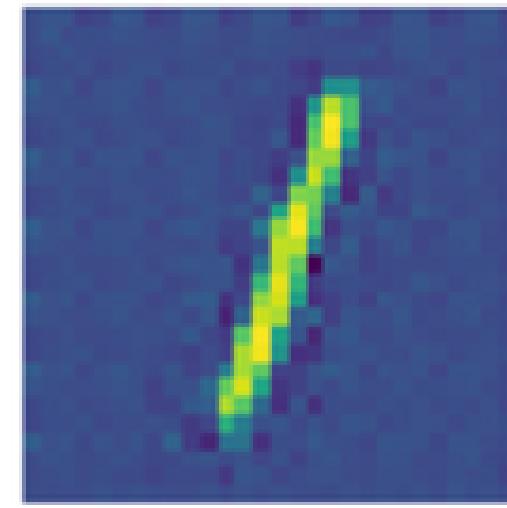
Normal



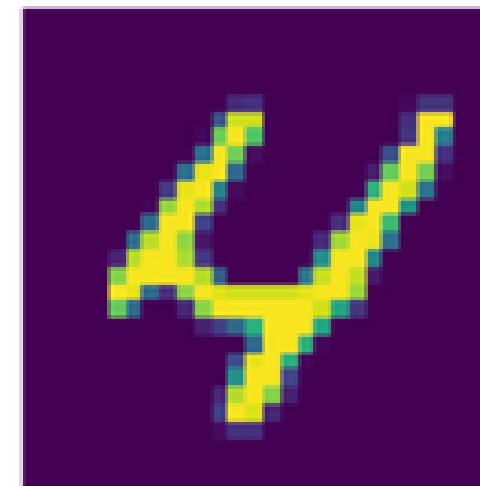
Adv



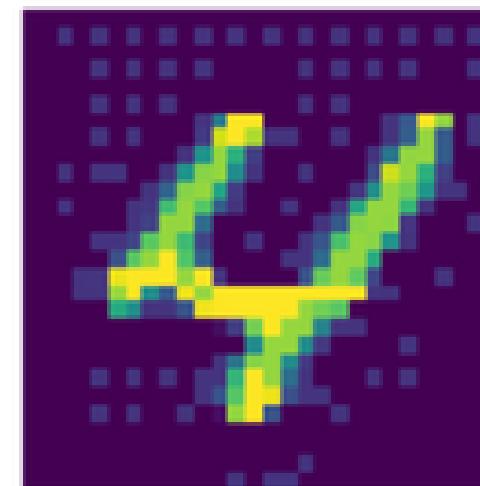
APE-GAN



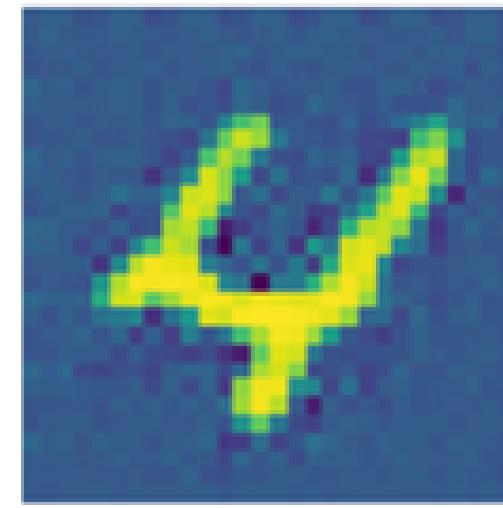
Normal



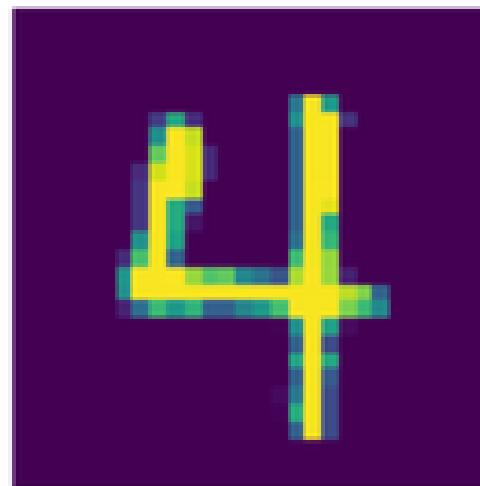
Adv



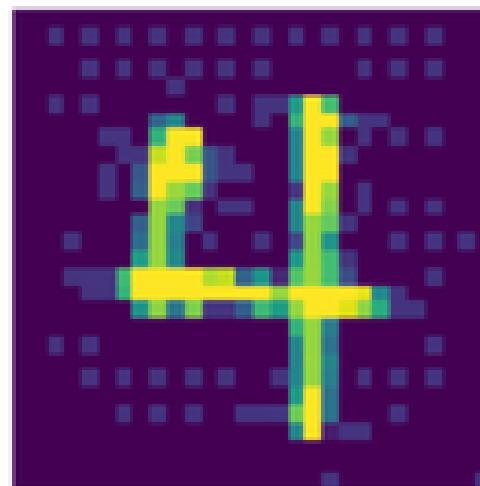
APE-GAN



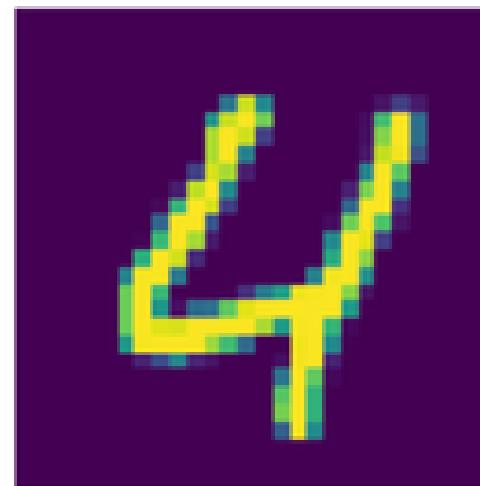
Normal



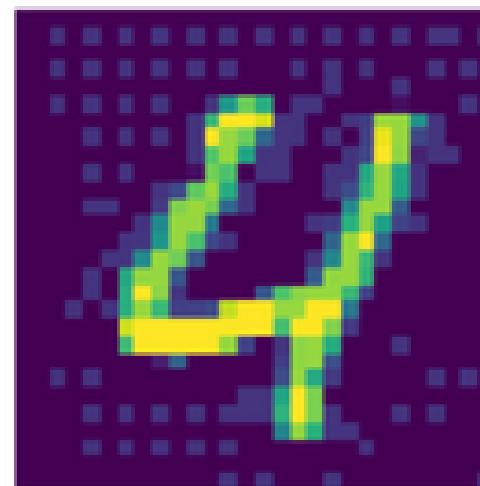
Adv



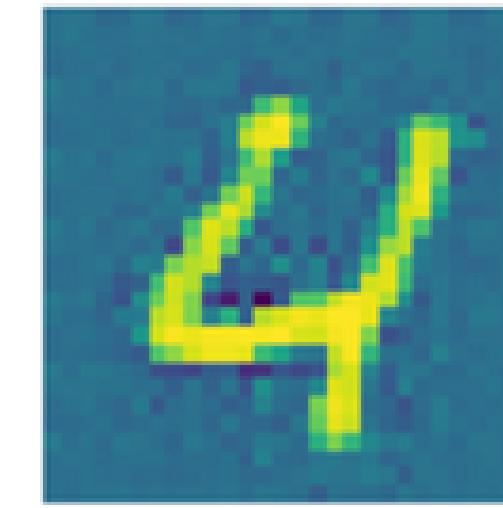
Normal



Adv



APE-GAN

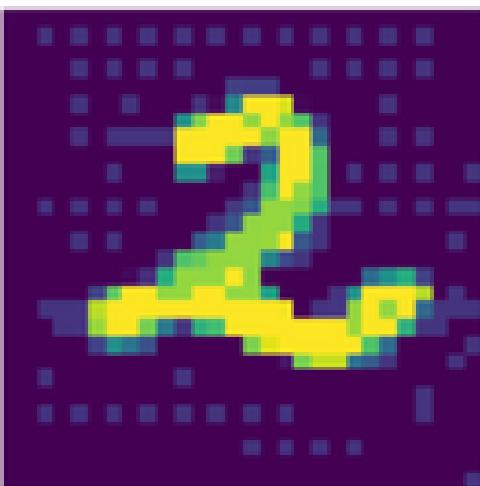




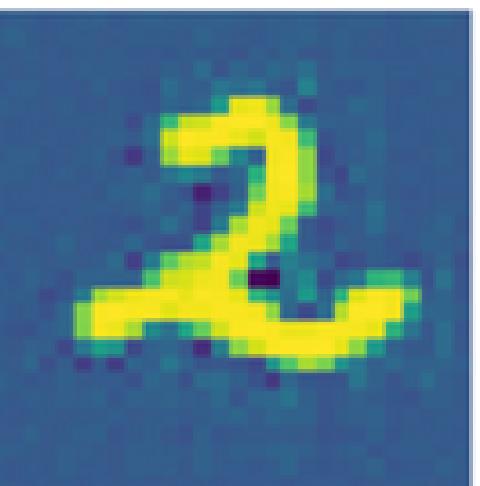
Normal



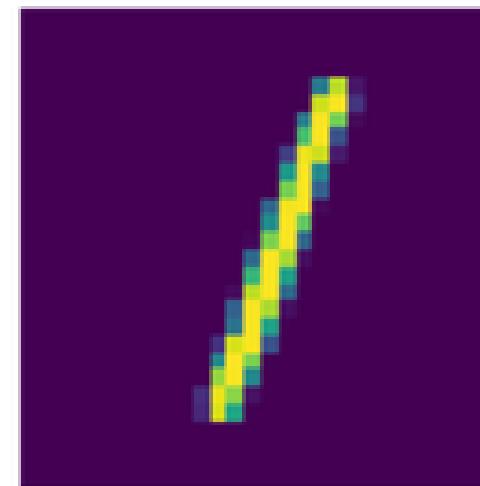
Adv



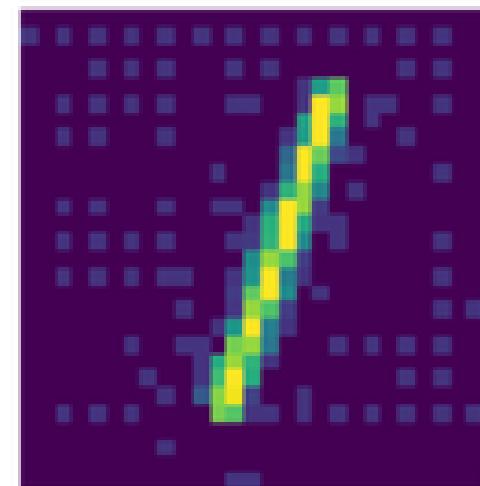
APE-GAN



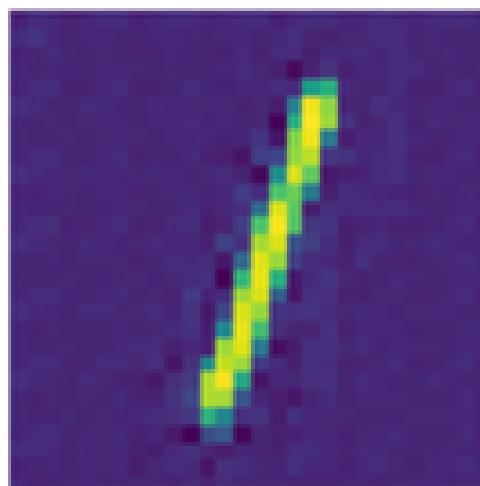
Normal



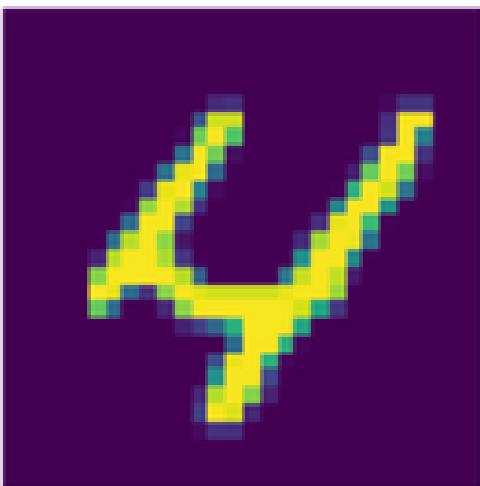
Adv



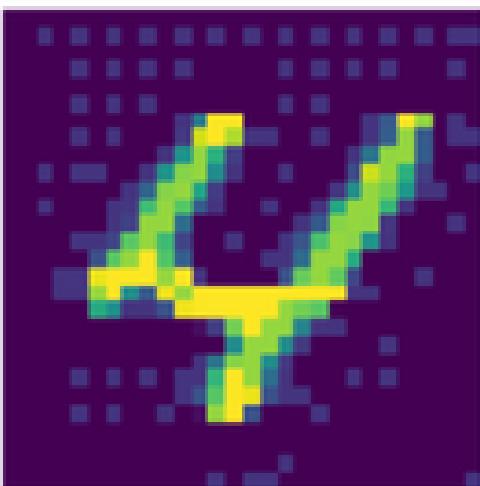
APE-GAN



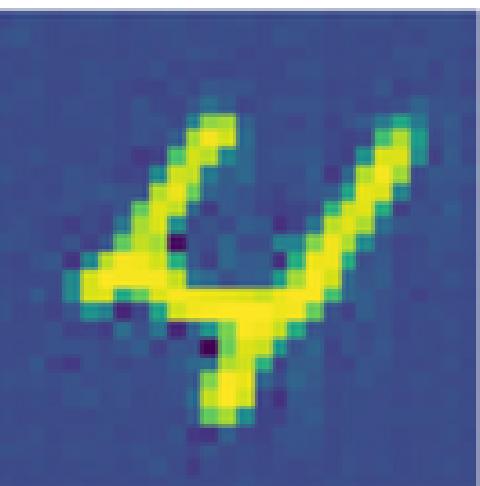
Normal



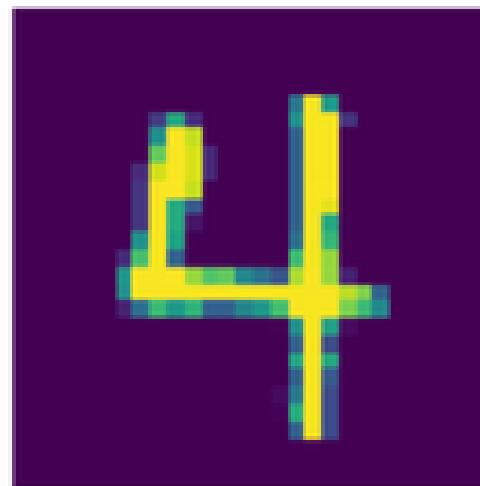
Adv



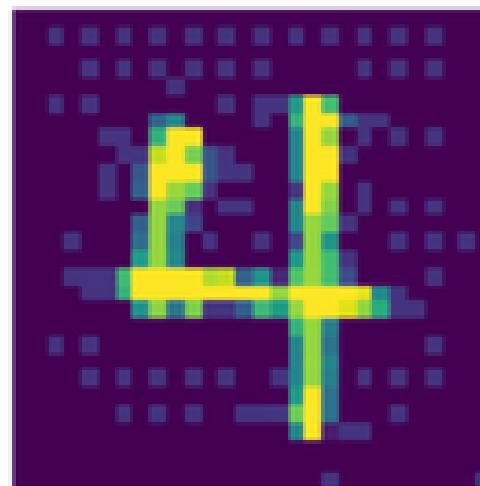
APE-GAN



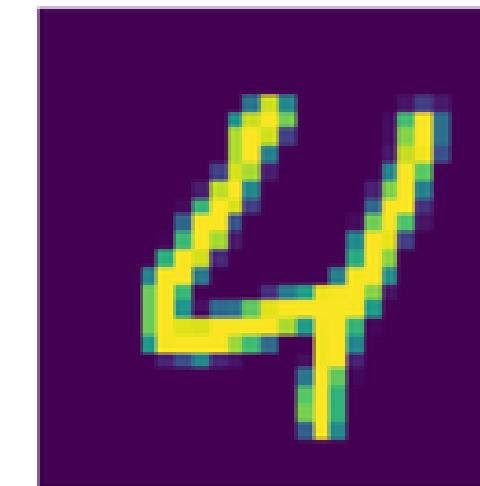
Normal



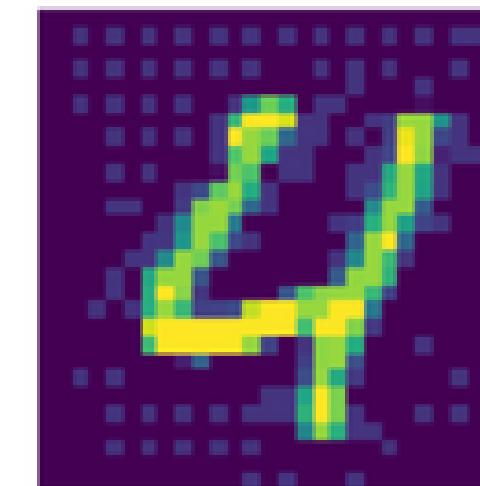
Adv



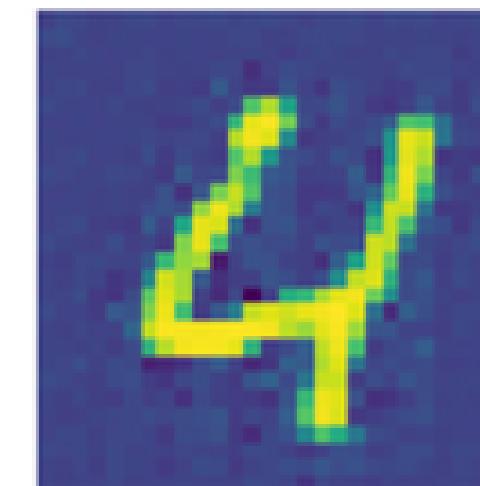
Normal



Adv

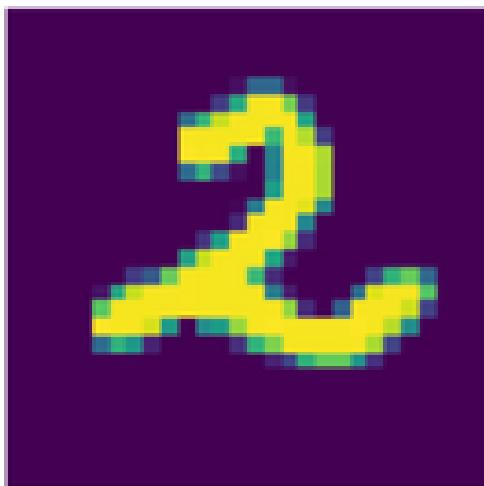


APE-GAN

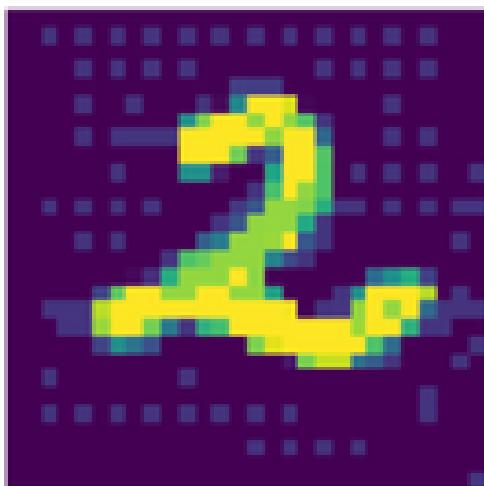




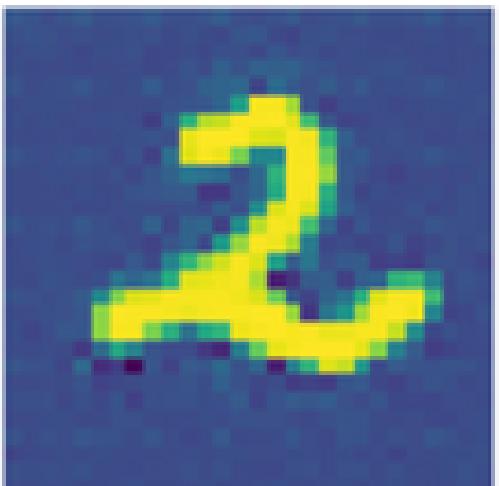
Normal



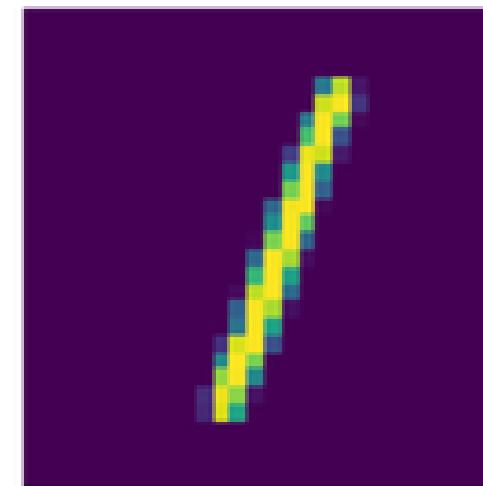
Adv



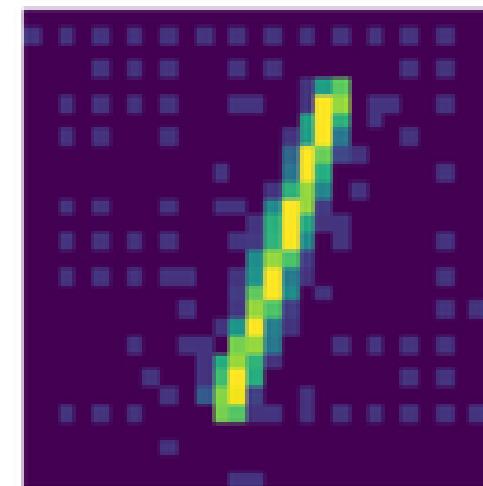
APE-GAN



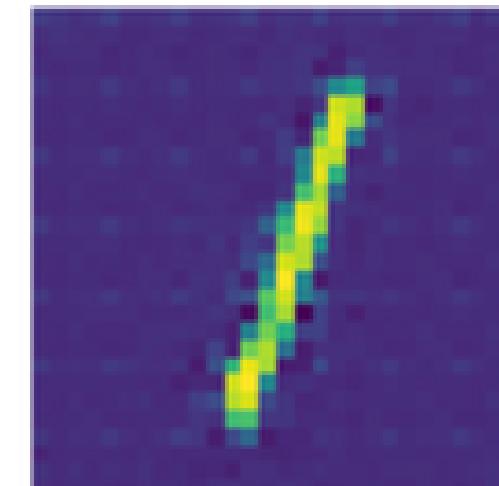
Normal



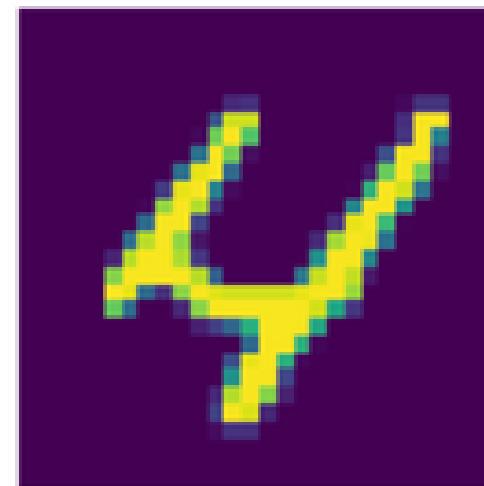
Adv



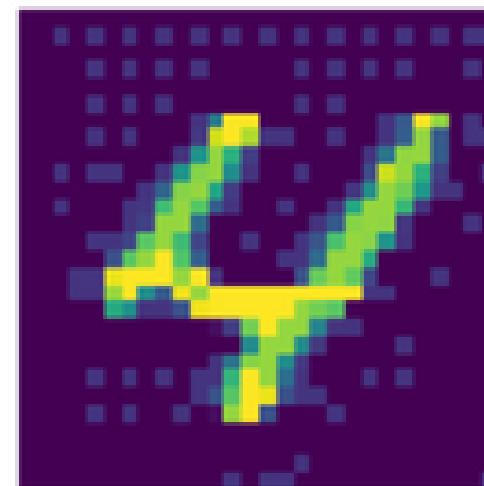
APE-GAN



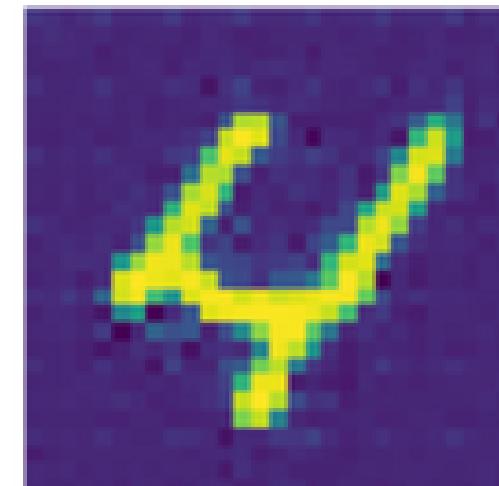
Normal



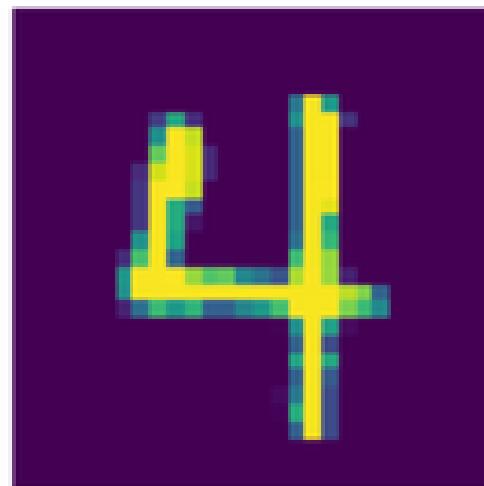
Adv



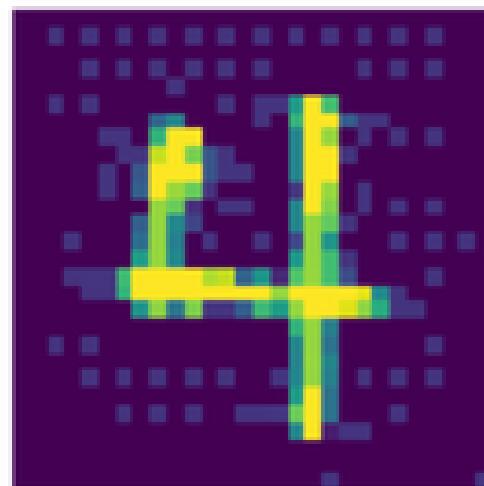
APE-GAN



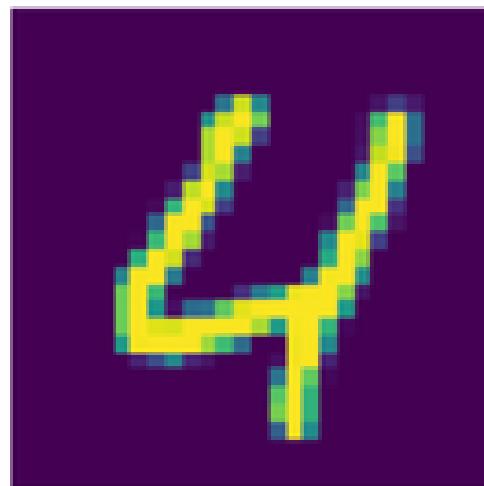
Normal



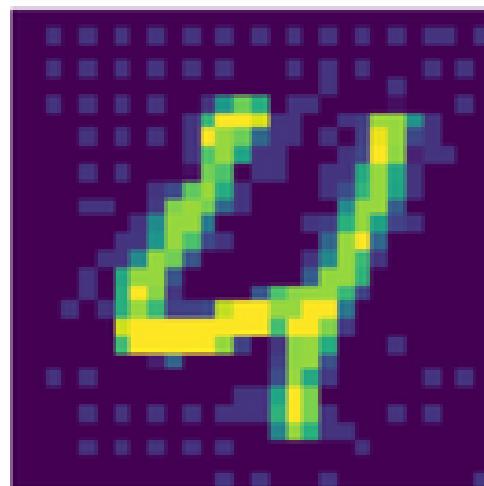
Adv



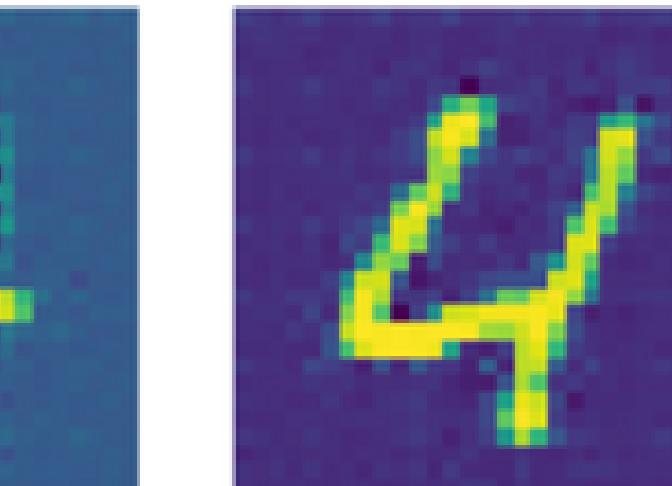
Normal



Adv

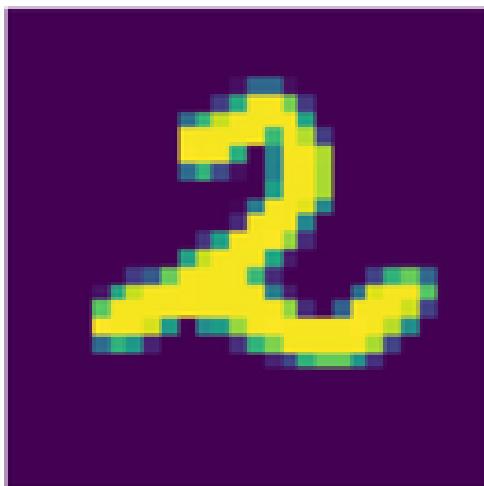


APE-GAN

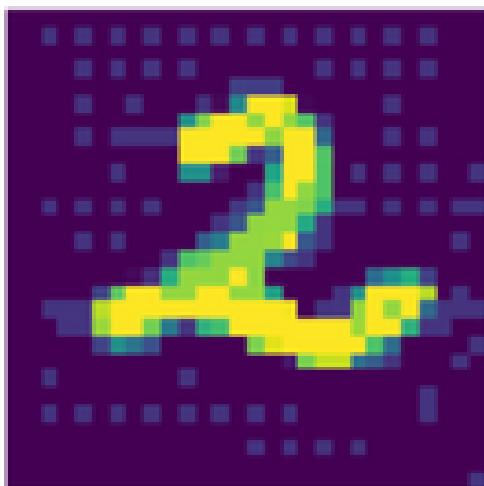




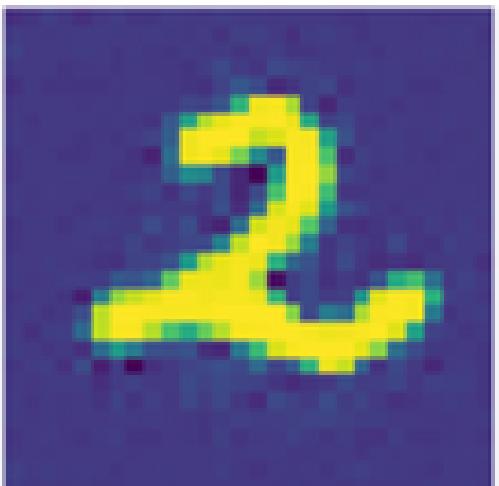
Normal



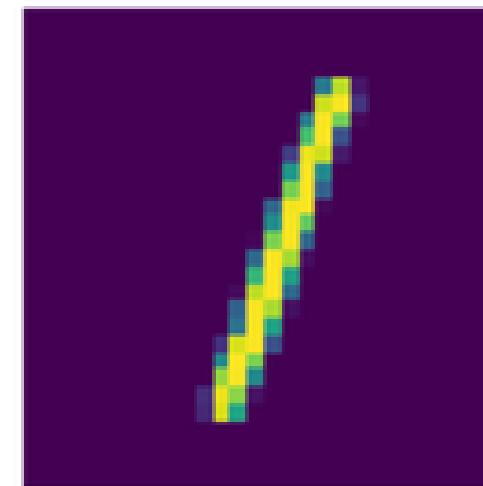
Adv



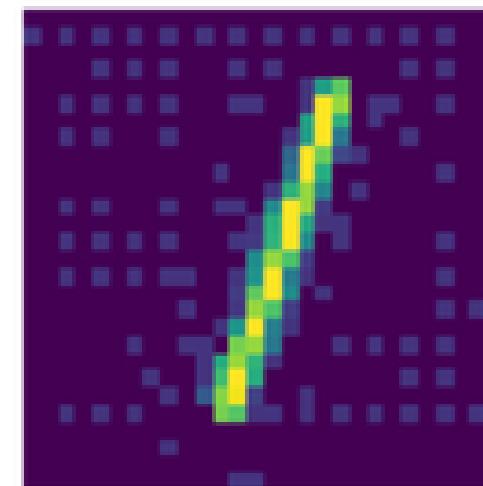
APE-GAN



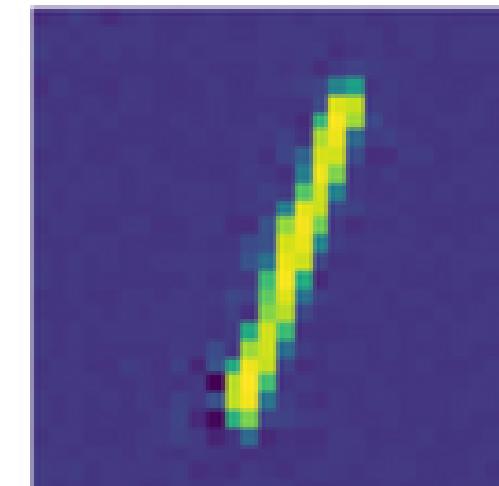
Normal



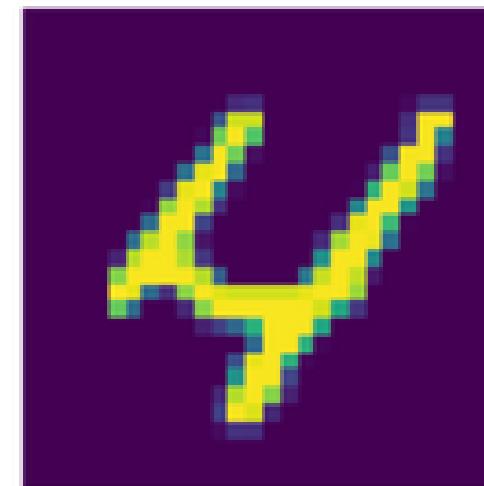
Adv



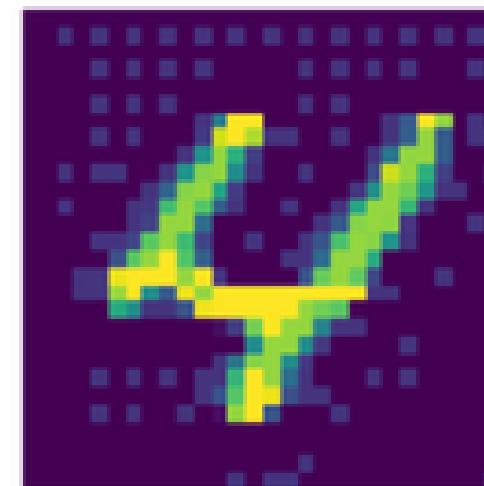
APE-GAN



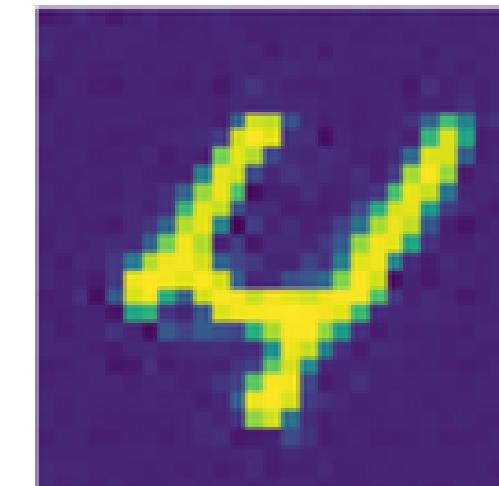
Normal



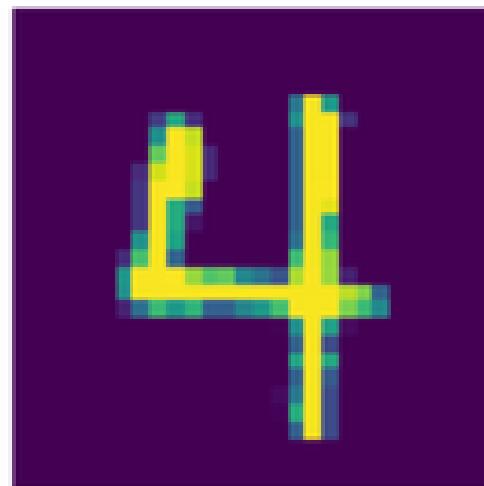
Adv



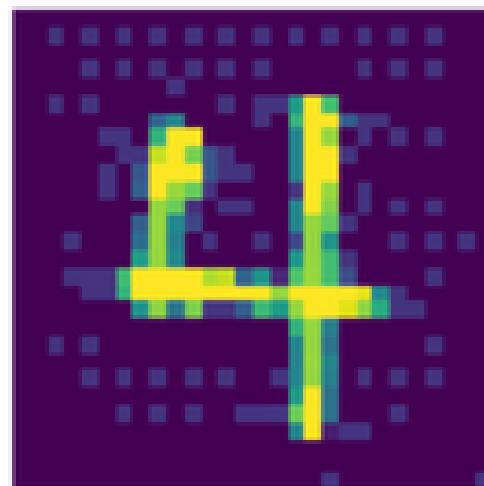
APE-GAN



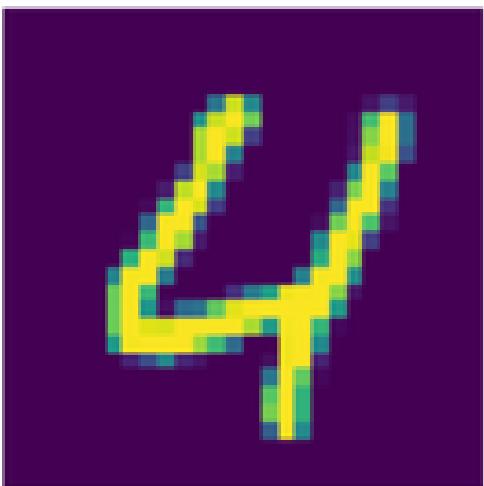
Normal



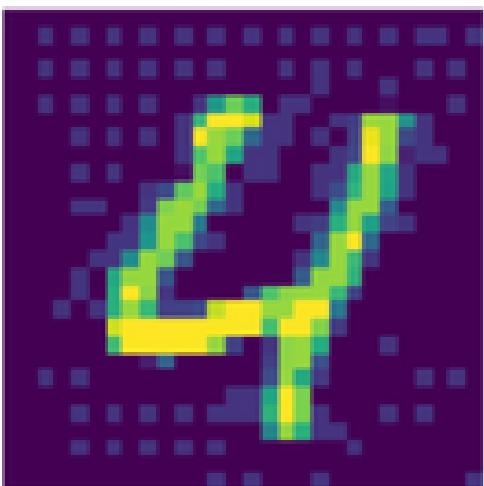
Adv



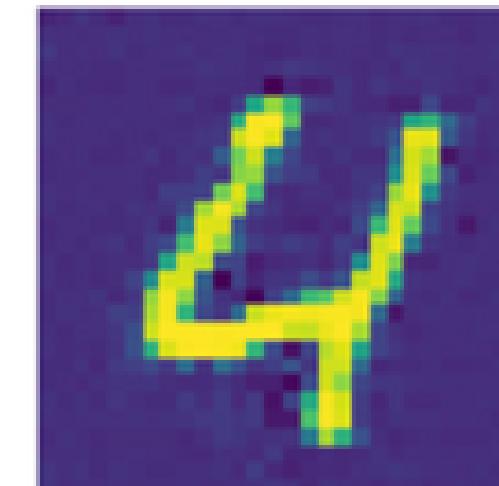
Normal



Adv

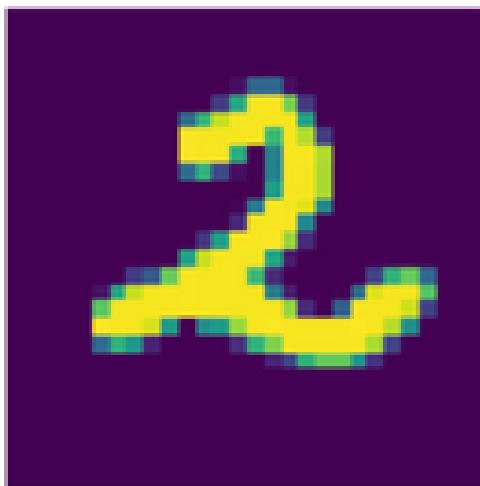


APE-GAN

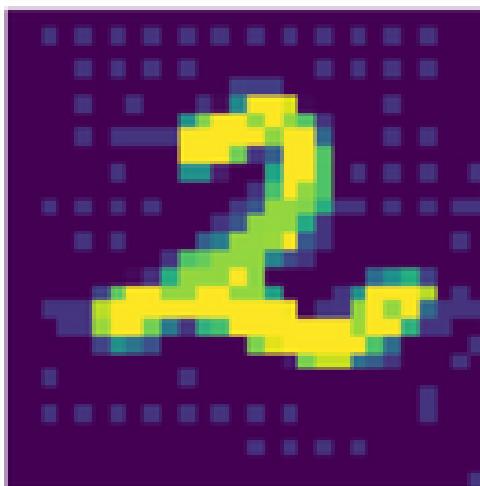




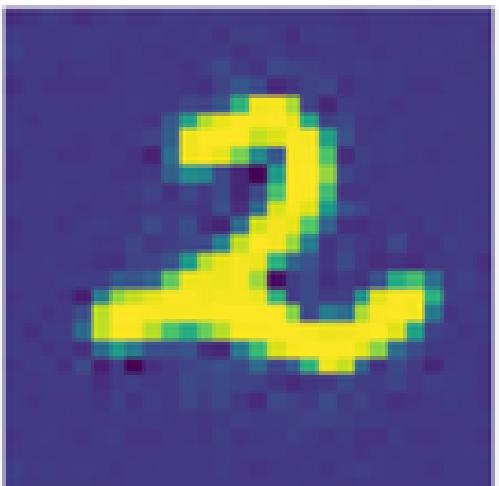
Normal



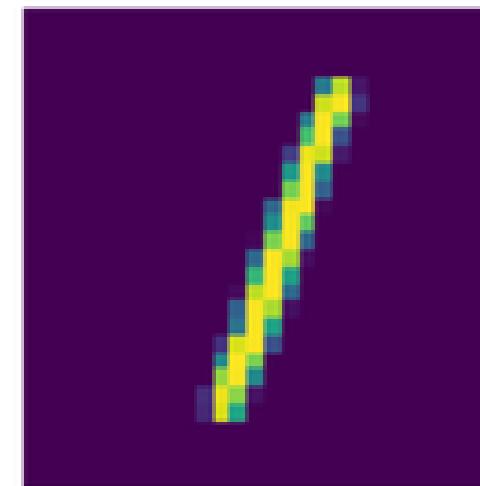
Adv



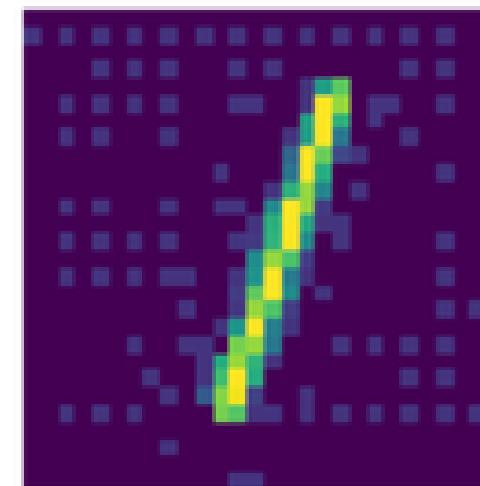
APE-GAN



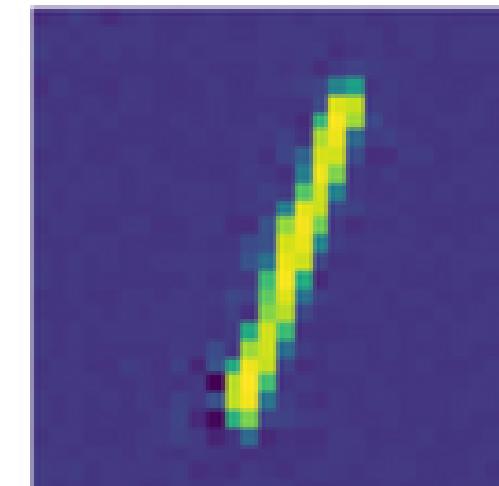
Normal



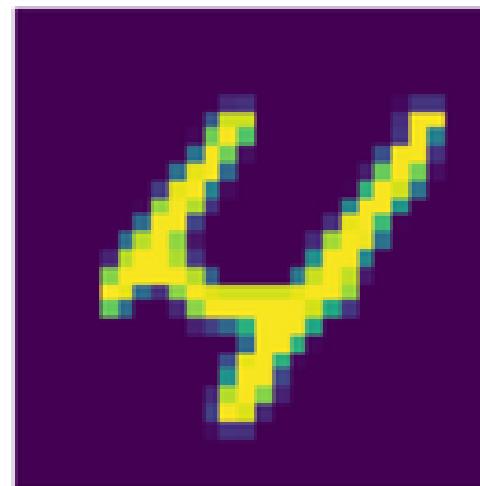
Adv



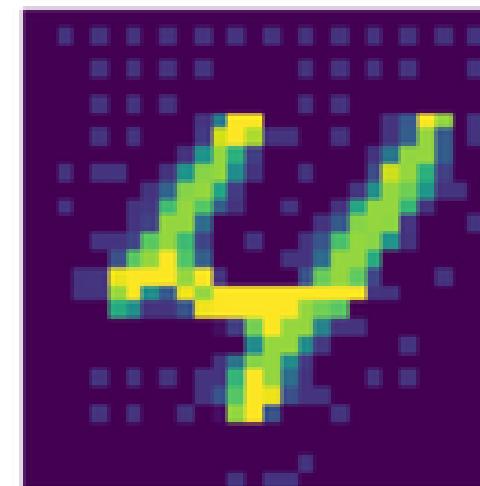
APE-GAN



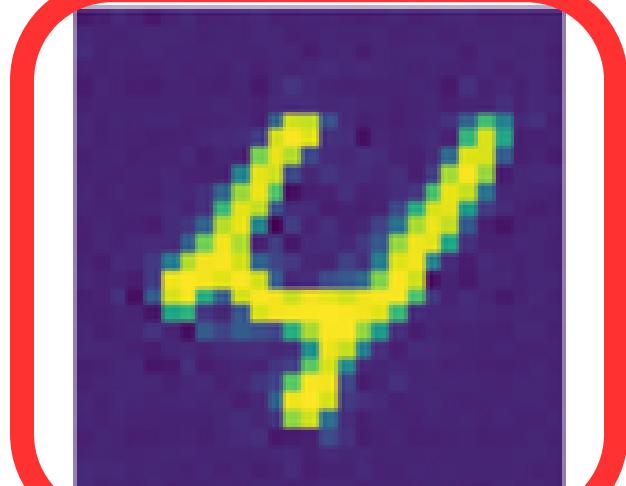
Normal



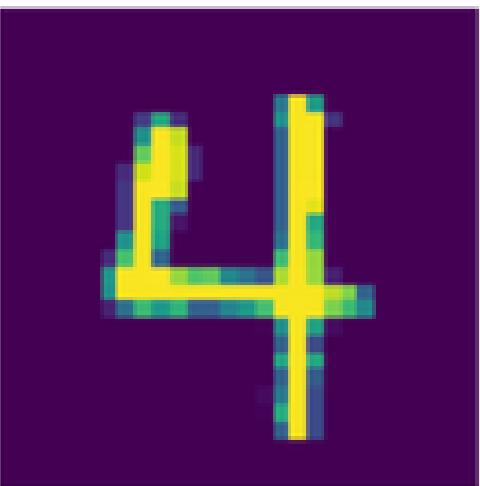
Adv



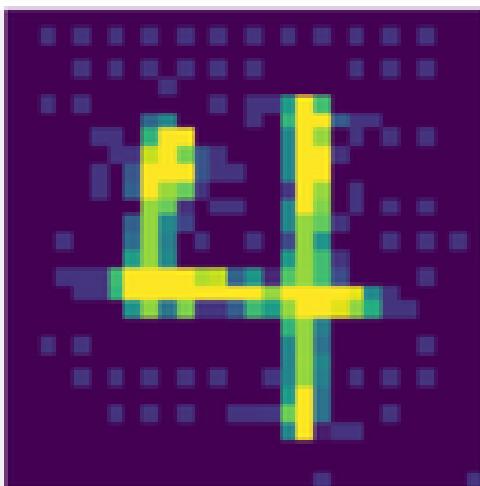
APE-GAN



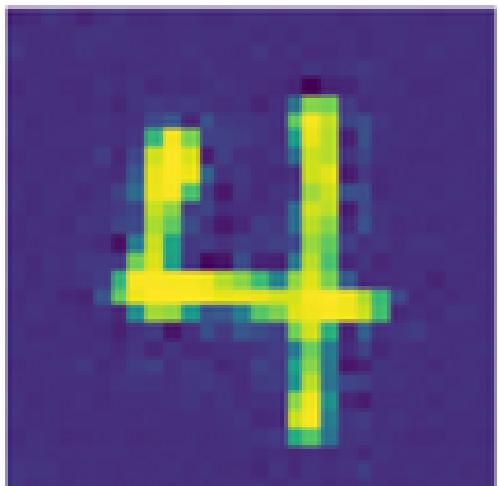
Normal



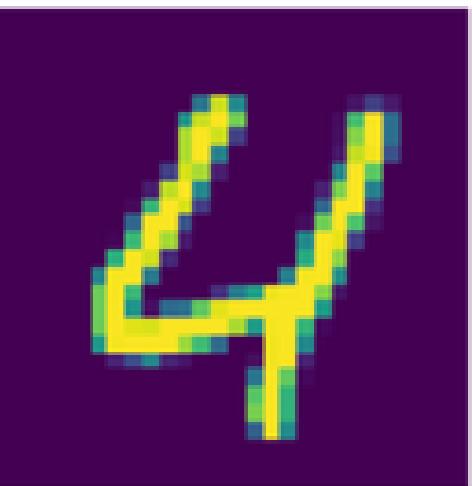
Adv



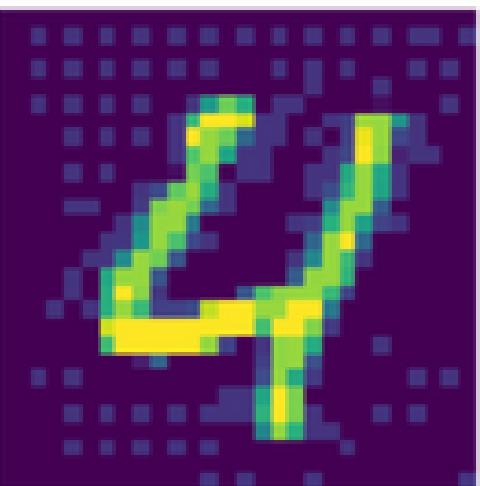
APE-GAN



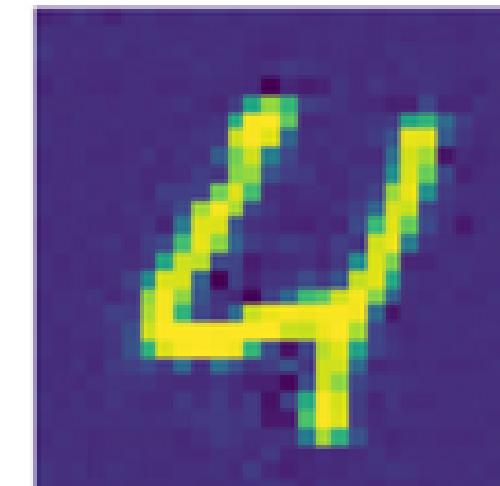
Normal



Adv



APE-GAN





# Results of U-Net Defense



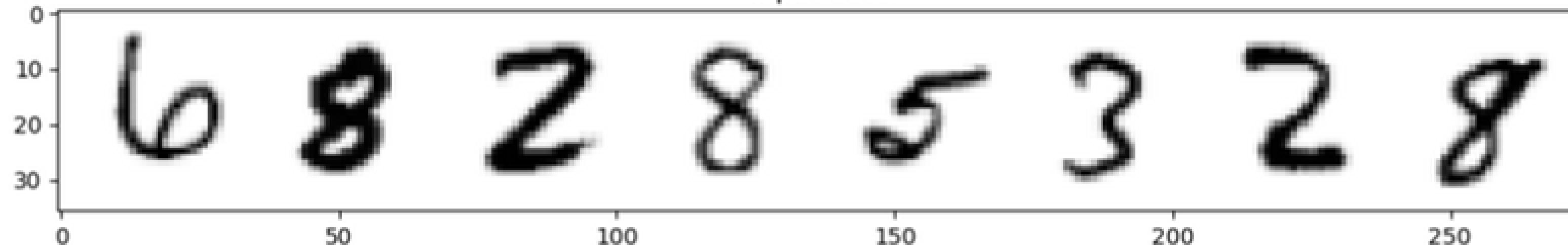
# FGSM Training U-Net



# FGSM Training U-Net



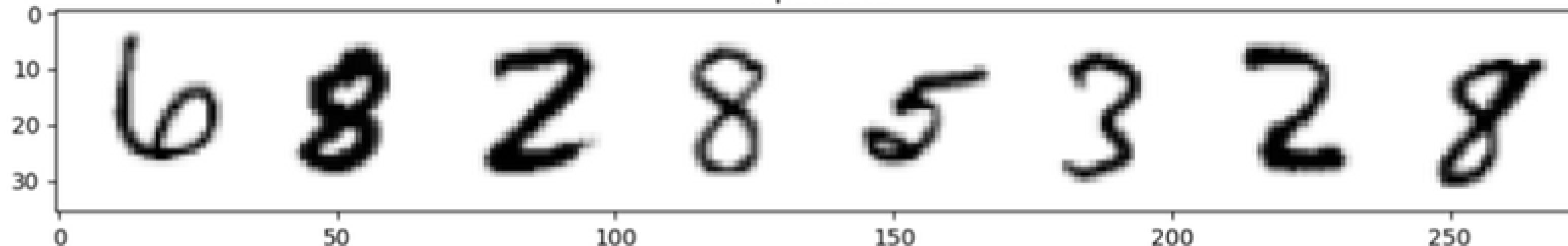
Input data



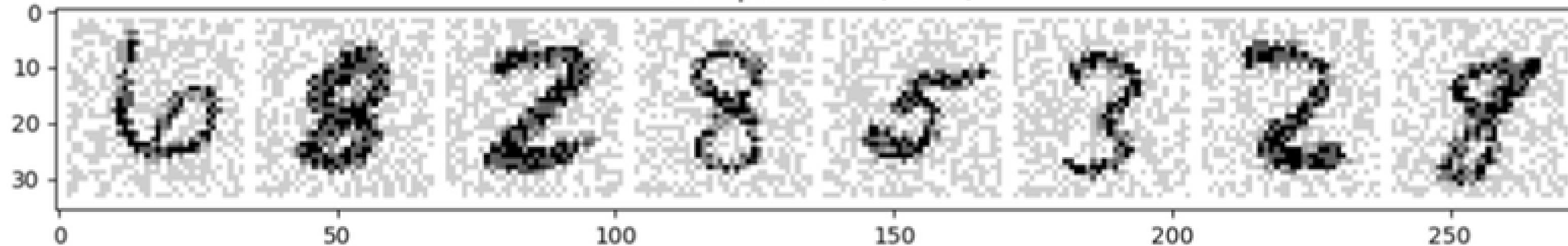
# FGSM Training U-Net



Input data



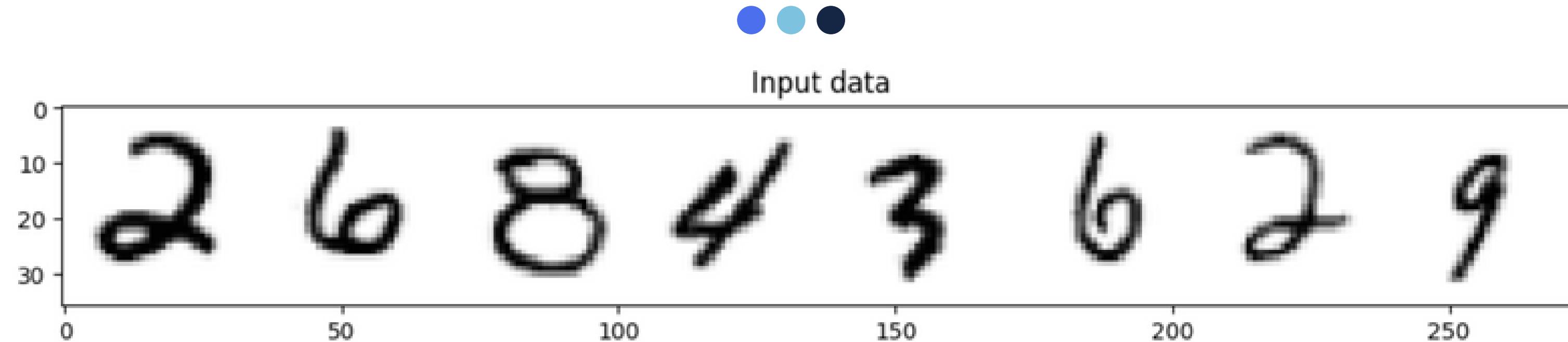
Corrupted data (FGSM)



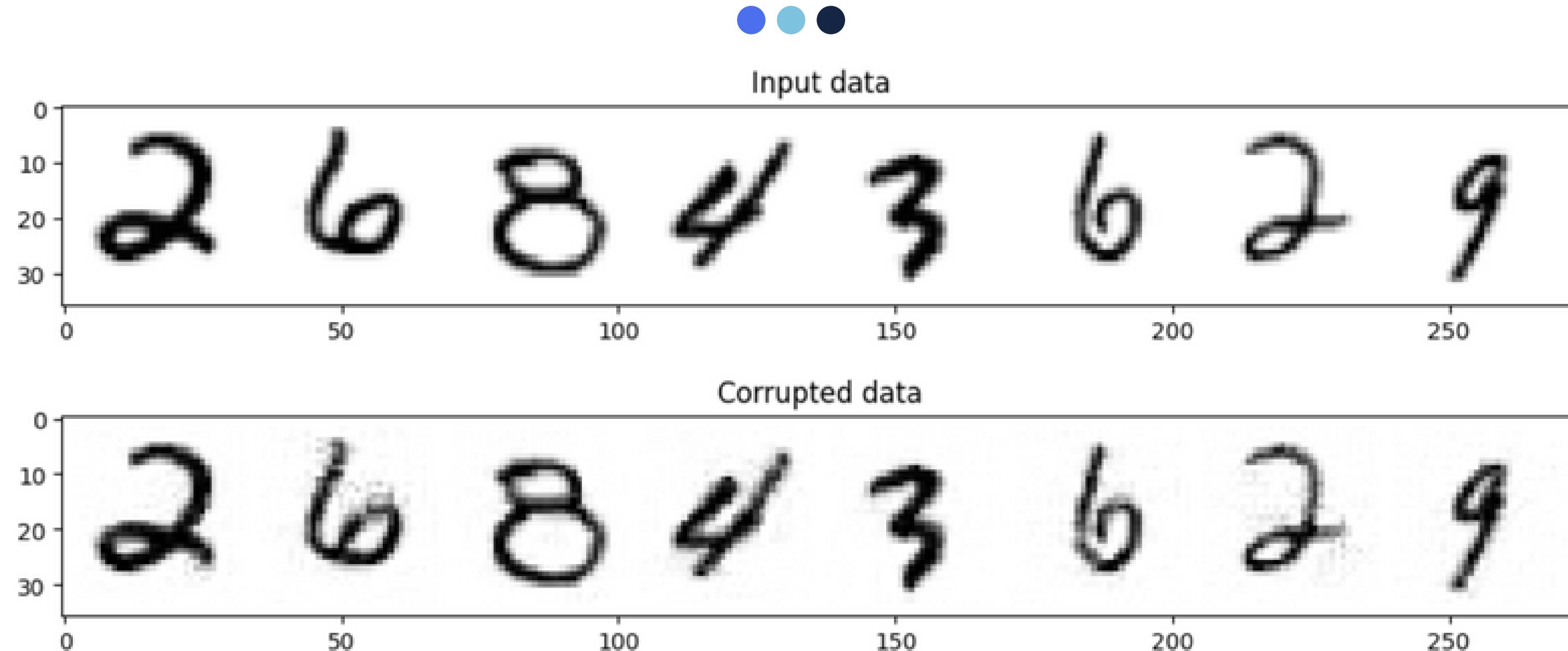
# DeepFool Attack



# DeepFool Attack 🔎



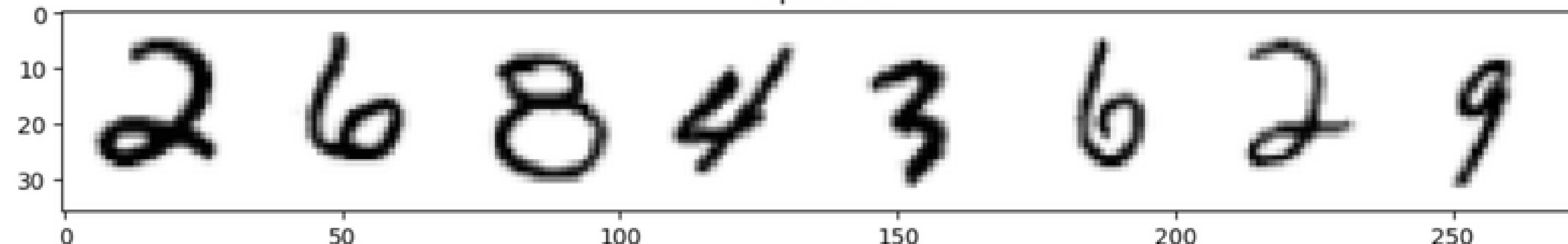
# DeepFool Attack 🔎



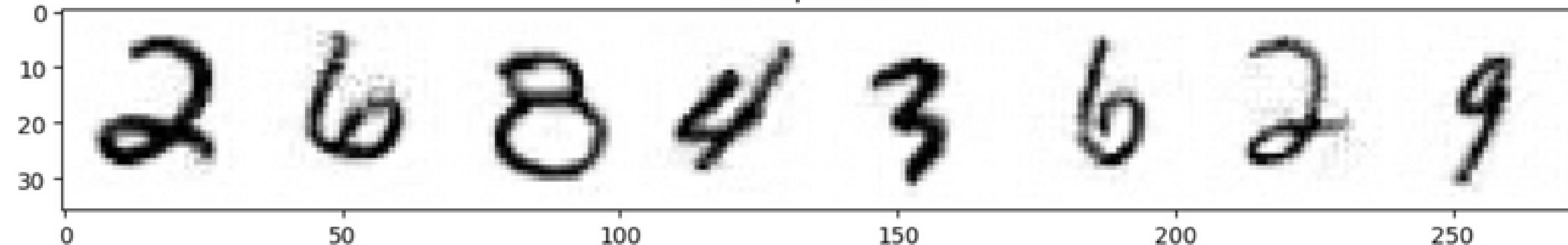
# U-Net Prediction



Input data



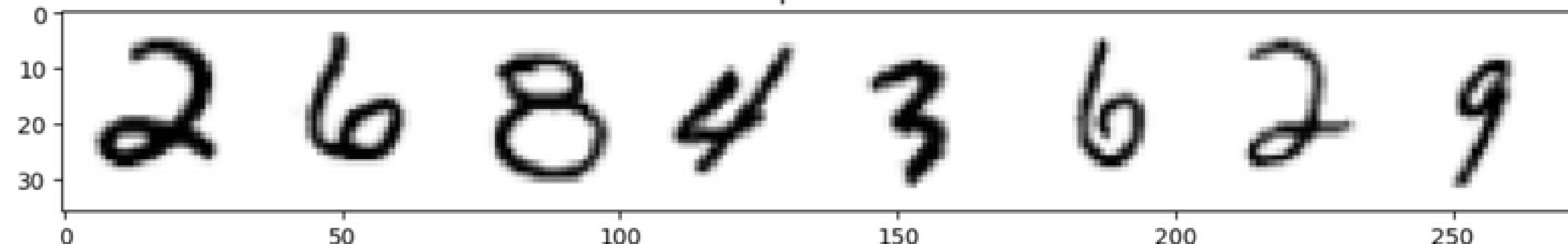
Corrupted data



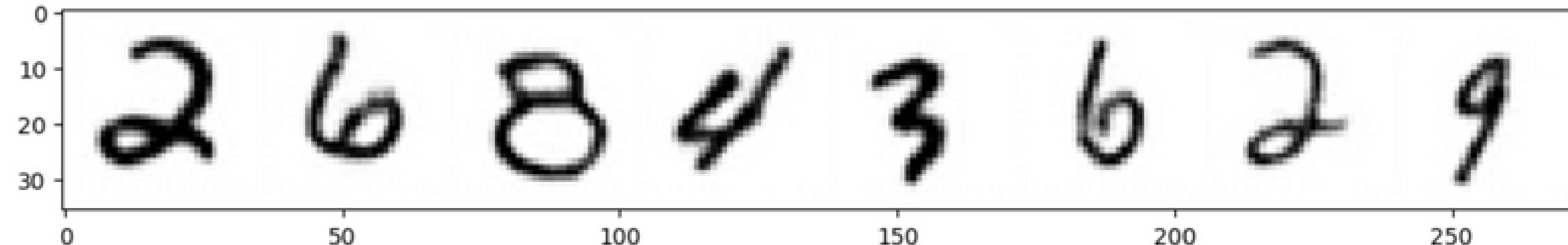
# U-Net Prediction



Input data



Network Predictions



# FGSM Training & Deepfool Attack

	Original Data	Deepfool Data	U-Net Data
accuracy	97.4%	1.9%	96.8%
	100% Accuracy	1.9% Accuracy	96.8% Accuracy

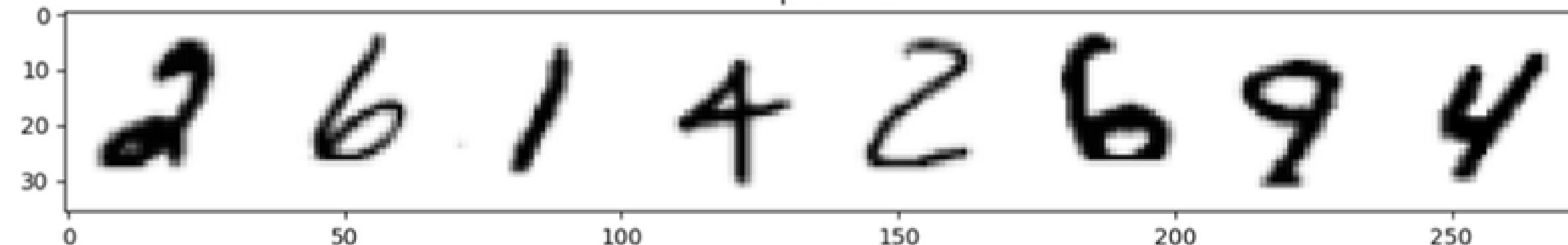
# DeepFool Training U-net



# DeepFool Training U-net



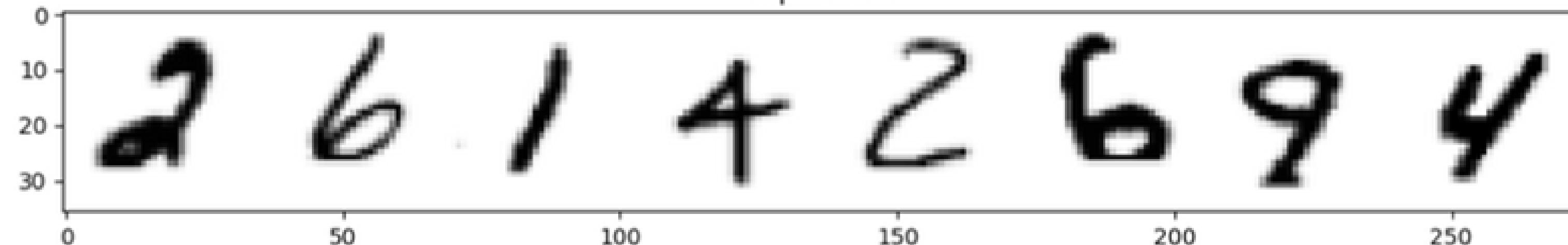
Input data



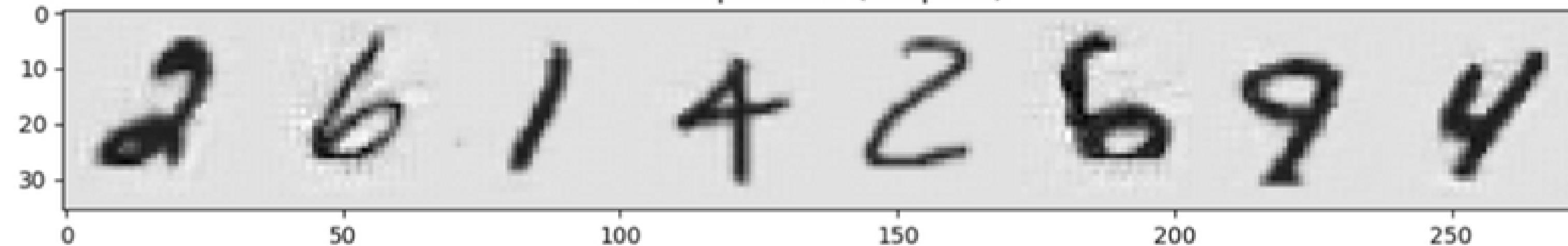
# DeepFool Training U-net



Input data



Corrupted data (DeepFool)



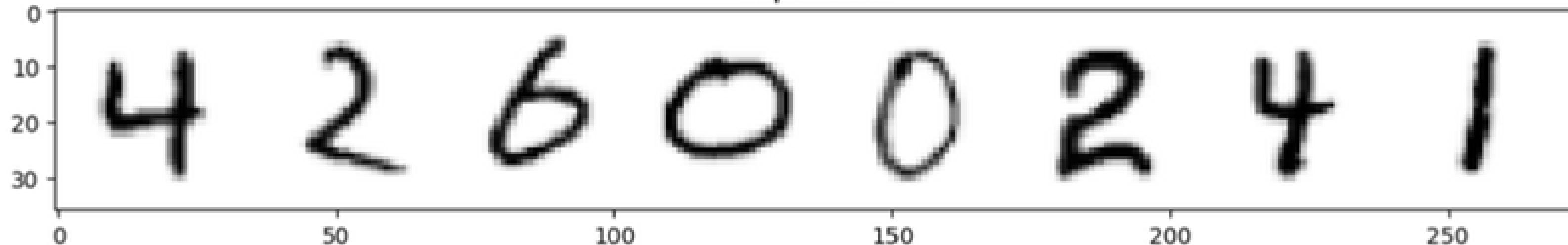
# FGSM Attack



# FGSM Attack



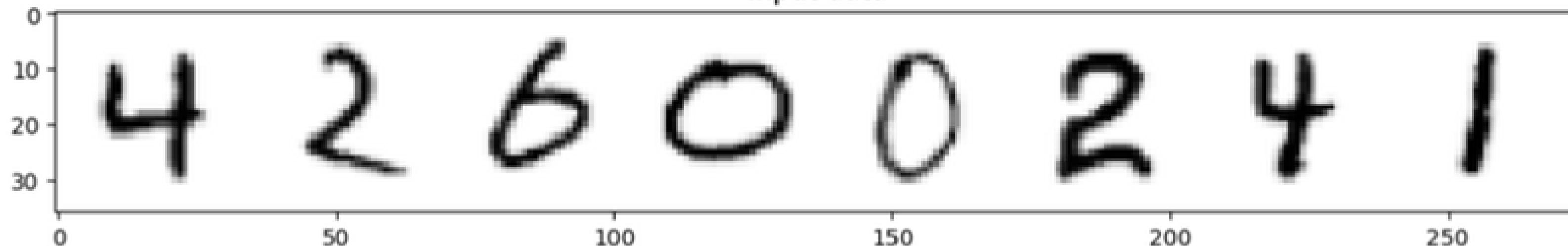
Input data



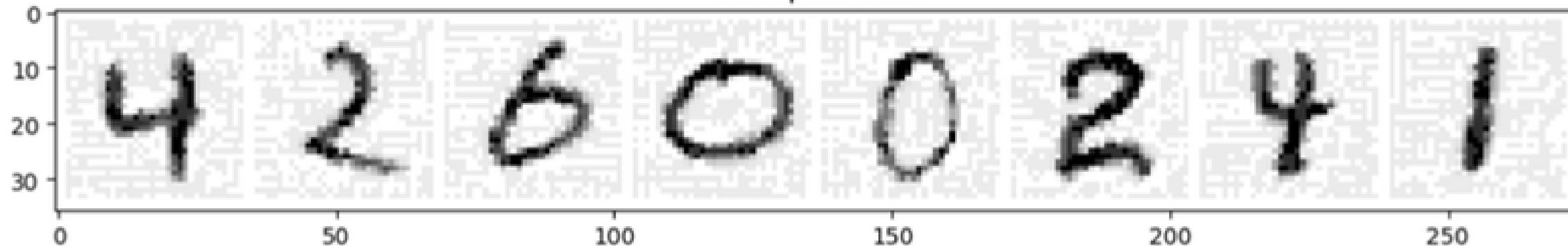
# FGSM Attack



Input data



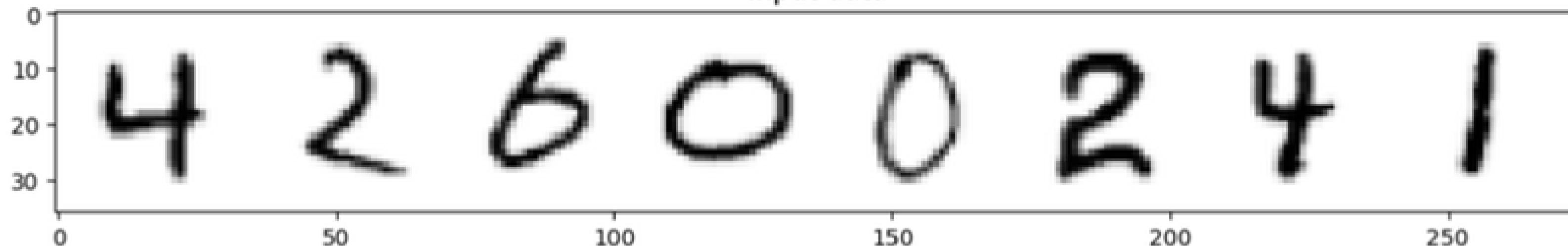
Corrupted data



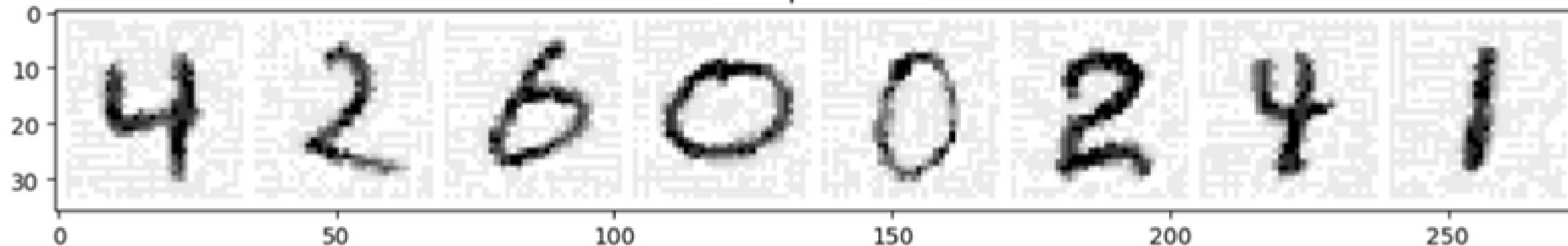
# U-Net Prediction



Input data



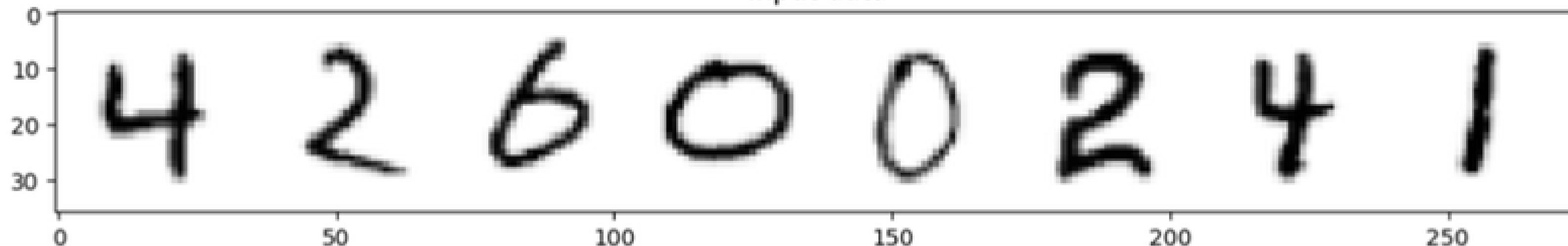
Corrupted data



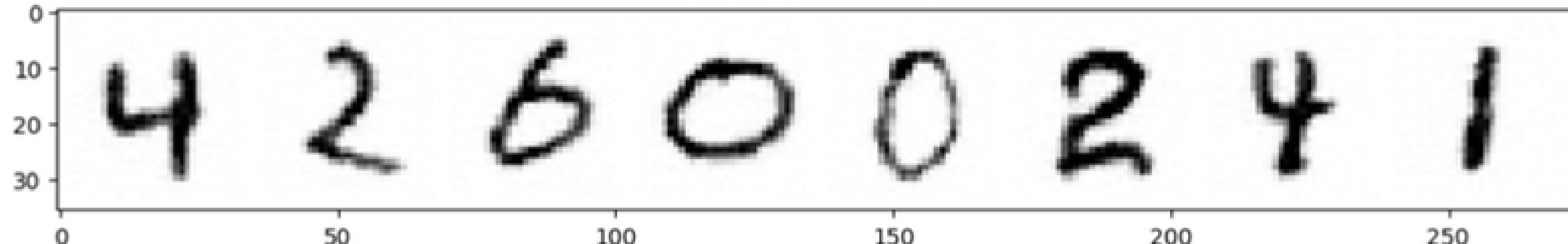
# U-Net Prediction



Input data



Network Predictions



# DeepFool Training & FGSM Attack

	Original Data	FGSM Data	U-net Data
accuracy	98.1%	22.4%	80.3%

# Attack Overview

	accuracy
Original	97.9%
Attacked by FGSM	22.5%
Attacked by Deepfool	1.9%

# U-Net Defense Overview

	training by FGSM	training by Deepfool
attack by FGSM	82.6%	80.3%
attack by Deepfool	96.8%	98.1%

# Summary

# Summary

1

Generate adversarial images by FGSM model and purify those images with APE-GAN & Diffusion model

# Summary

1

Generate adversarial images by FGSM model and purify those images with APE-GAN & Diffusion model

2

Try different model (Deepfool) to generate adversarial images

# Summary

1

Generate adversarial images by FGSM model and purify those images with APE-GAN & Diffusion model

2

Try different model (Deepfool) to generate adversarial images

3

We found that U-net is a strong model to purify image, better than APE-GAN

# Furniture Work

# Furture Work

1

Try FGSM & Deepfool to attack  
on different datasets, and find  
the purification result of  
different dataset

# Furture Work

1

Try FGSM & Deepfool to attack on different datasets, and find the purification result of different dataset

2

Find other models that can defend attack more effectively

# Furture Work

1

Try FGSM & Deepfool to attack on different datasets, and find the purification result of different dataset

2

Find other models that can defend attack more effectively

3

Try to defend different data type, such as video and text



# Thanks for Listing !

