

---

# **Applications of Adversarial Attacks on Matching-based Algorithms**

**CS588**

**Speaker: Woo Jae Kim**

# Class Objectives

---

- **Up-to-date matching-based algorithms**
- **Adversarial attacks**
- **Combining adversarial attacks with matching-based algorithms**

---

# Matching-based Algorithms

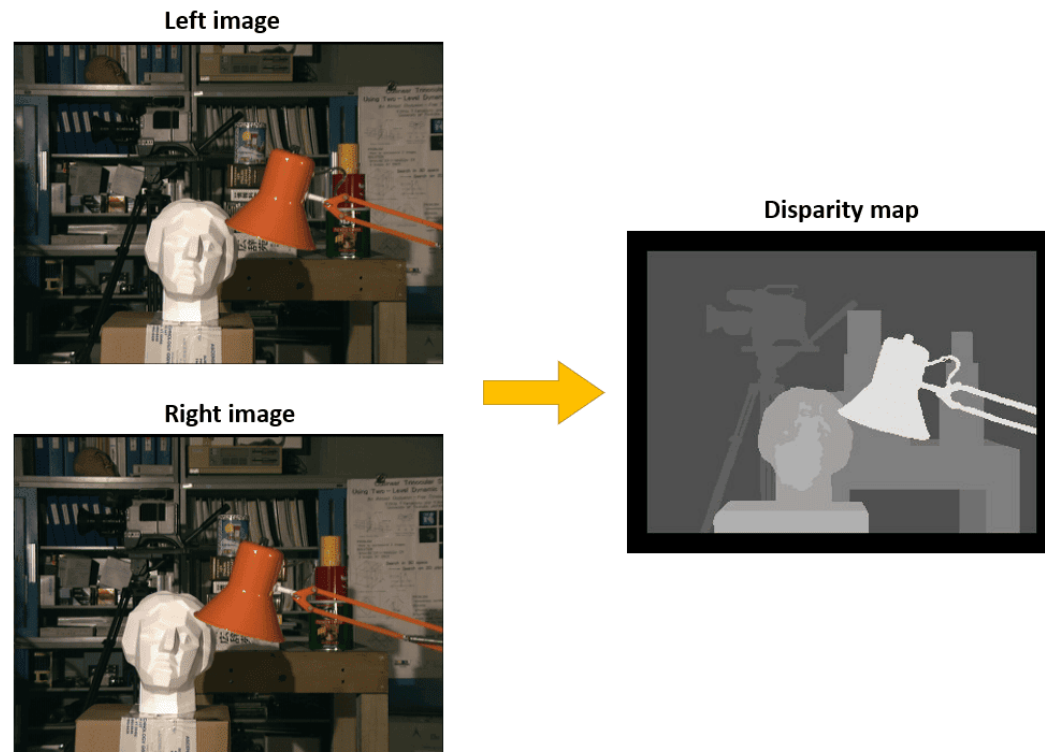
# Matching-based Algorithms

- What is matching-based algorithm?
- Finding correspondence or similarity between multiple sets of data
- Can be used to estimate 3D geometric structure of the data



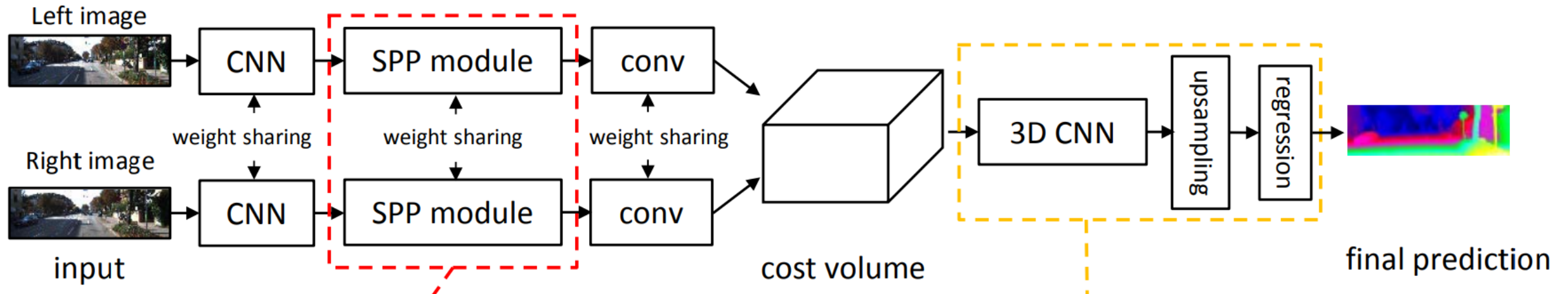
# Stereo Matching

- Uses two images taken from slightly different viewpoints
- Goal: Infer disparity or depth maps given the image pair



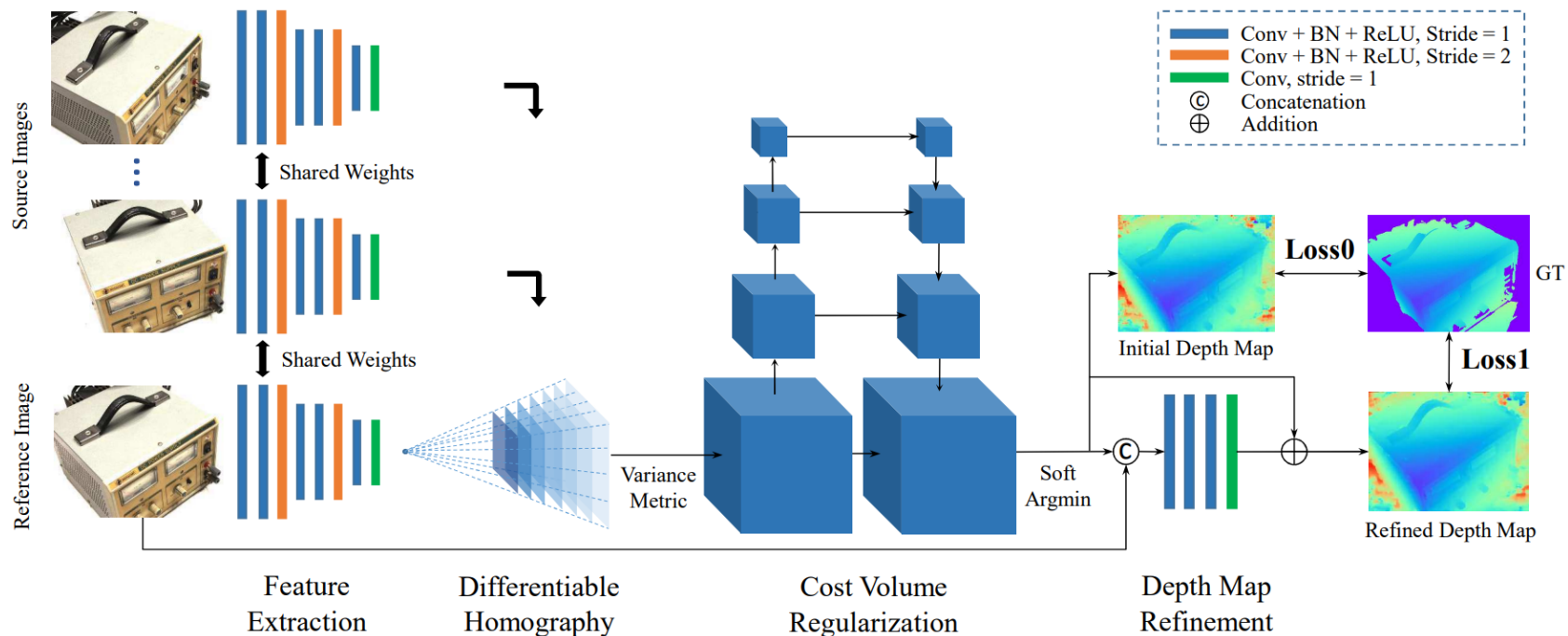
# Stereo Matching

- Current methods use CNN to extract features from each image
- It then constructs a cost volume, from which the final depth is predicted



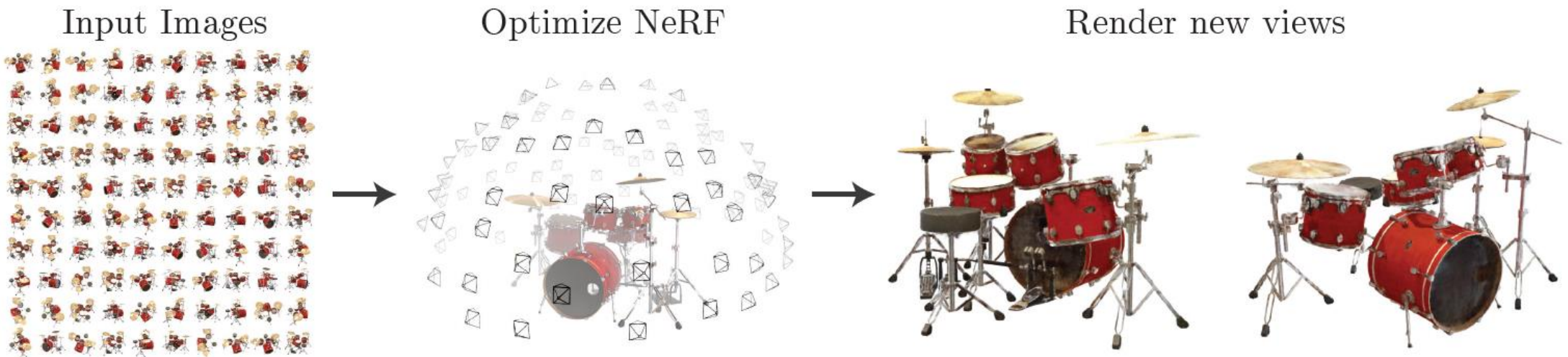
# Multi-view Stereo

- Multi-view stereo uses more than two images captured multiple viewpoints
- Can potentially provide more accurate and detailed 3D reconstructions
- Goal: Learn correspondence among images to infer depth of objects



# 3D Matching (Novel-view Synthesis)

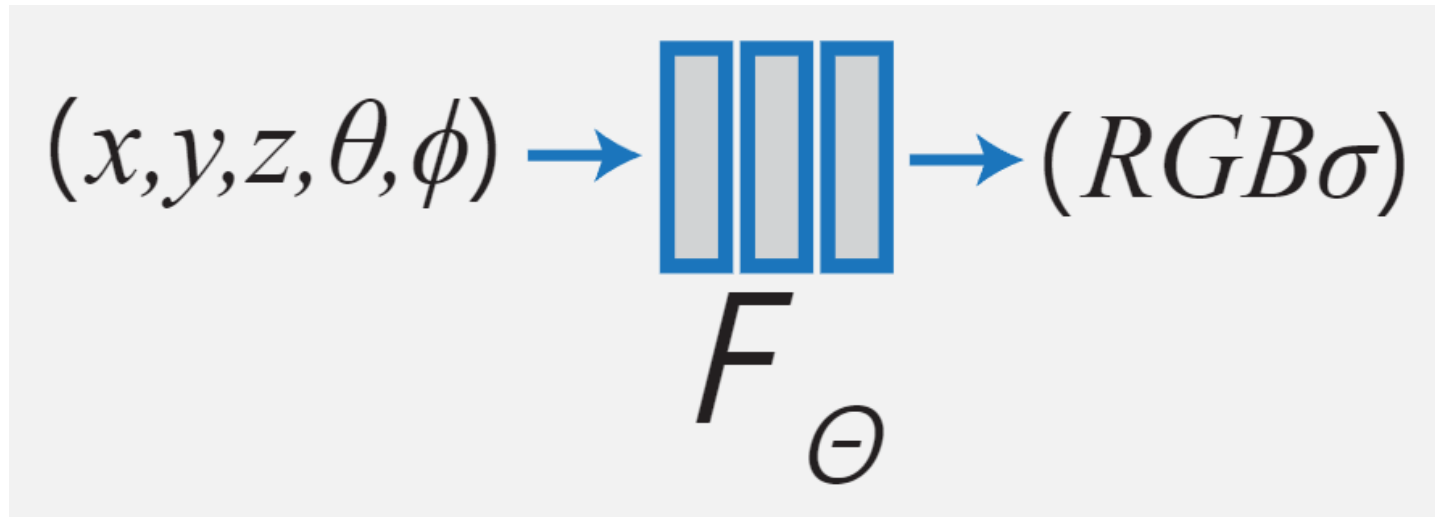
- More current methods use correspondence among multi-view images to model a continuous 3D representation of a scene
- This allows a more accurate and seamless prediction of the 3D geometry
- Became highly popular with introduction of Neural Radiance Fields (NeRF)





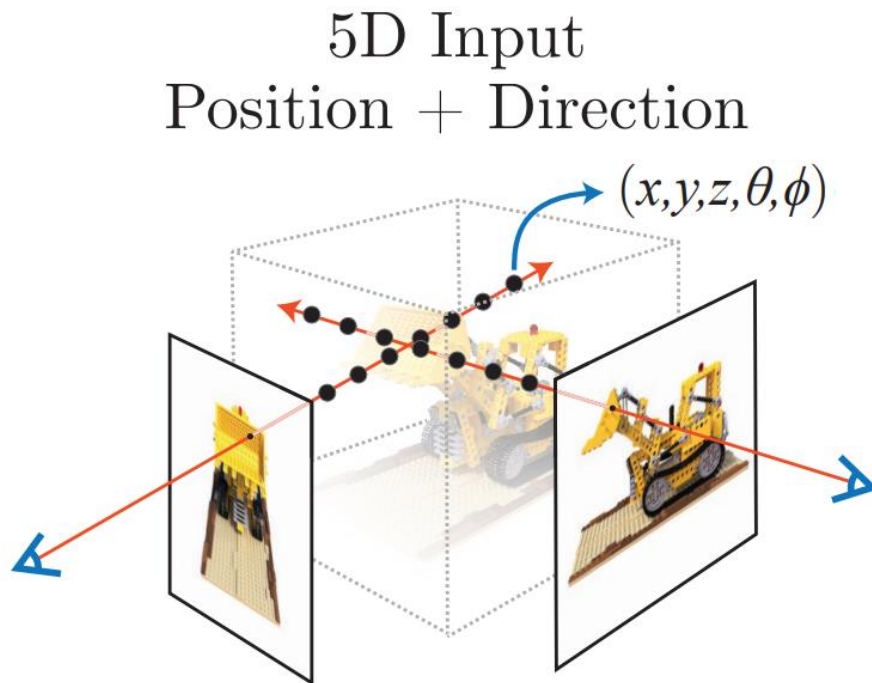
# Neural Radiance Fields

- NeRF uses an MLP network  $F_{\Theta}$  to represent a 3D scene
- Given the position  $(x, y, z)$  and direction  $(\theta, \phi)$  of a 3D point,
- MLP  $F_{\Theta}$  predicts the color  $RGB$  and density  $\sigma$



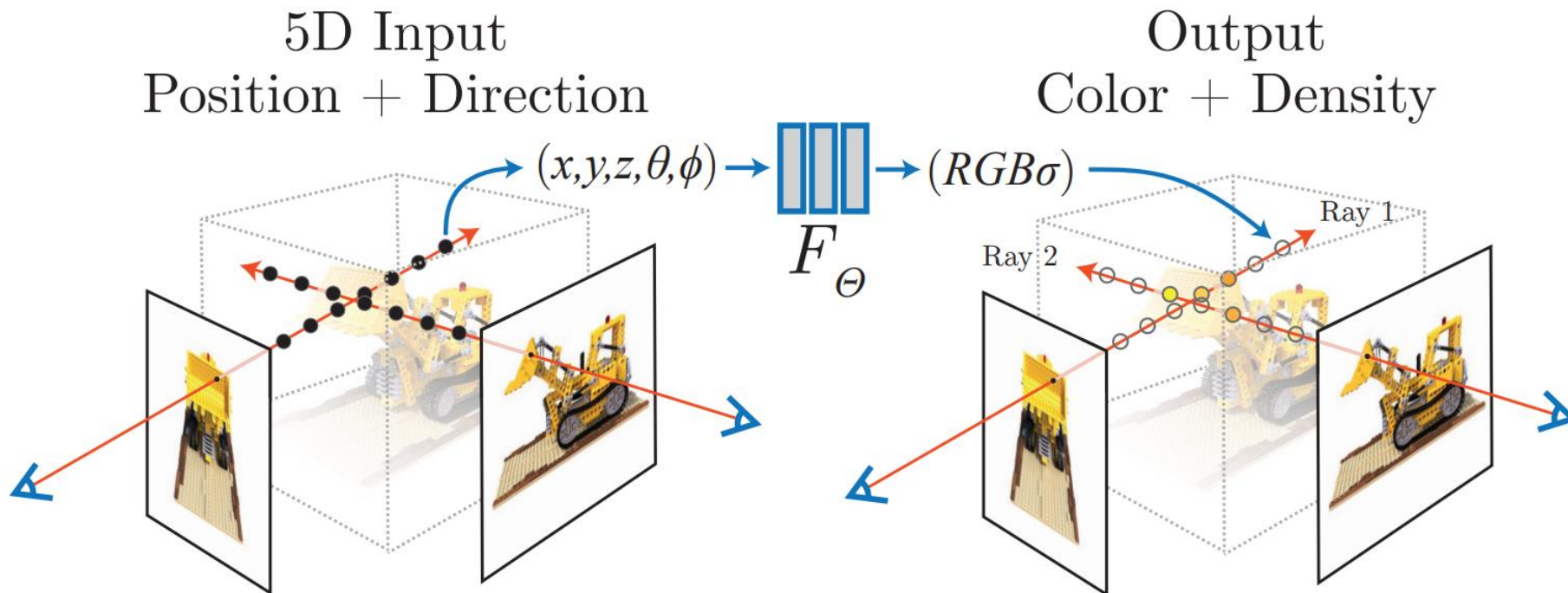
# Neural Radiance Fields: Input

- A set of images captured on a scene
- Randomly sample points along the rays from camera origin to each pixel of the image



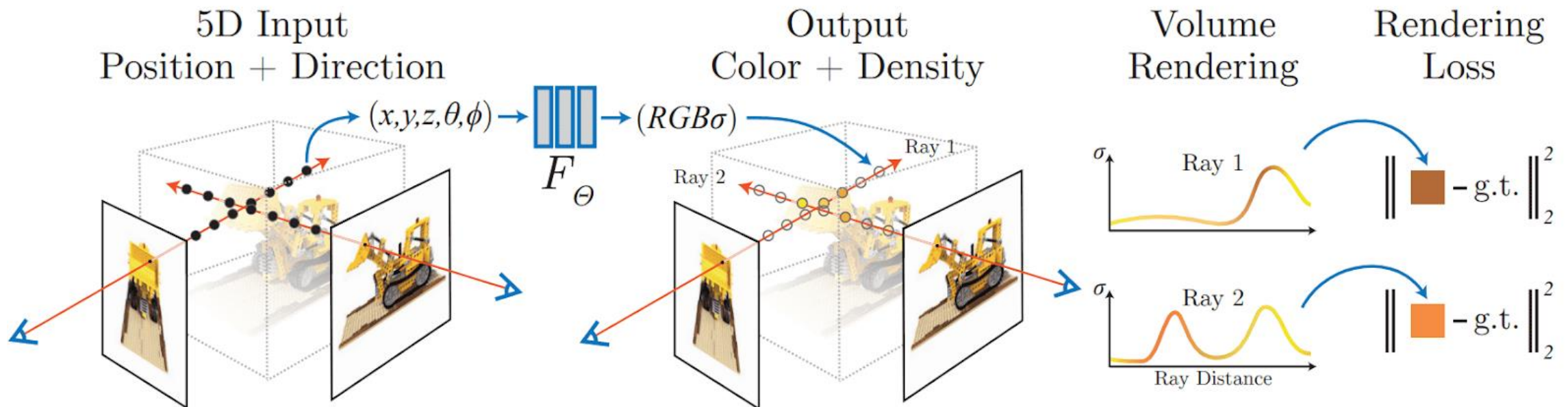
# Neural Radiance Fields: Output

- Pass the pos  $(x, y, z)$  and dir  $(\theta, \phi)$  of each point to the MLP network  $F_{\Theta}$
- MLP then predicts color RGB and density  $\sigma$  for each point



# Neural Radiance Fields: Volume Rendering

- Since NeRF models a continuous 3D volume, it can render images to an arbitrary view
- Volume rendering: Accumulation of colors and density of multiple points on a ray





# NeRF for Novel-View Synthesis

- NeRF shows excellent performance at modeling 3D scene & rendering arbitrary views



---

# Adversarial Attacks

# Adversarial Attacks

- Many deep learning models are known to be vulnerable to adversarial attacks



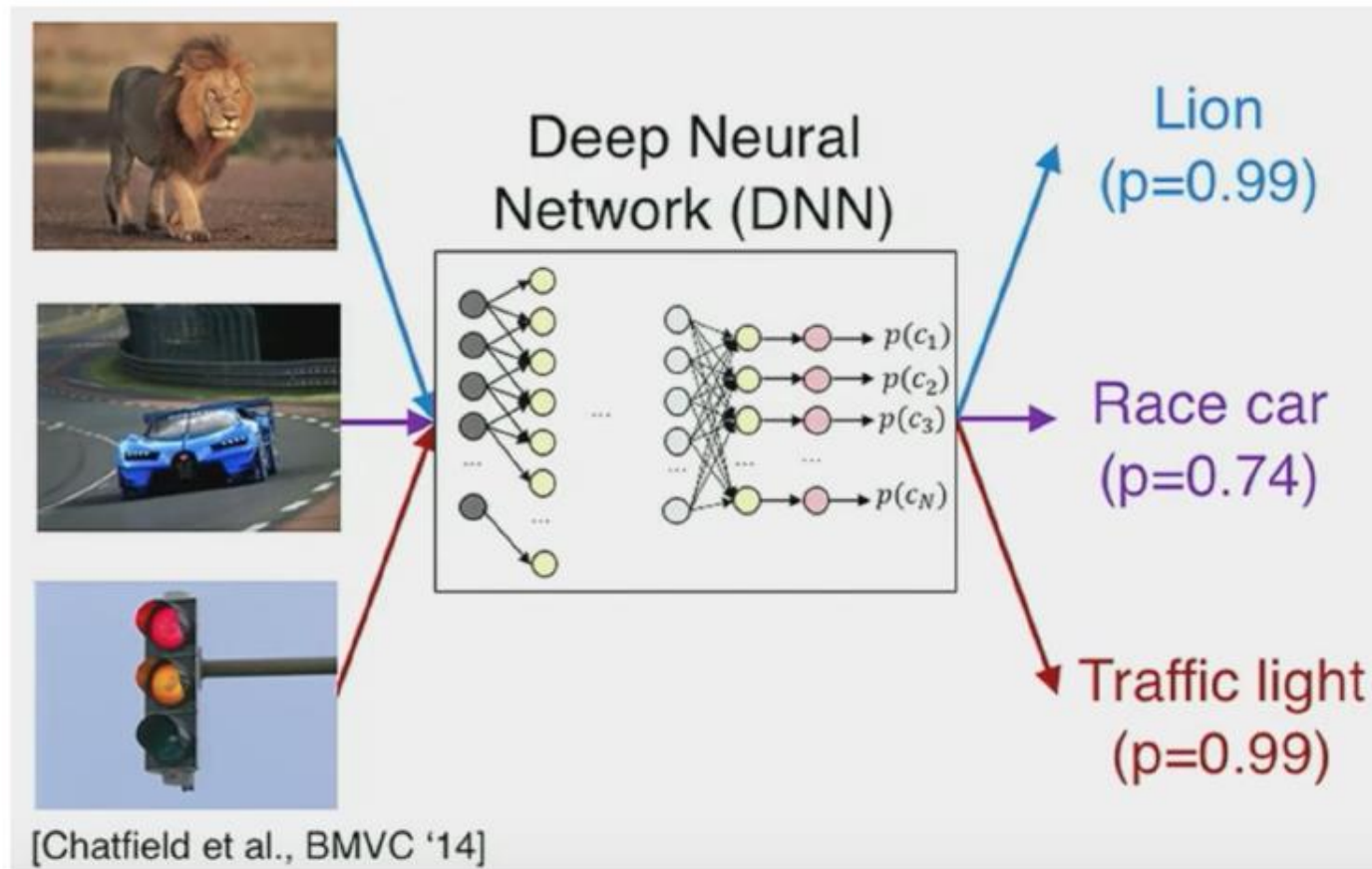


# Adversarial Attacks

- Many deep learning models are known to be vulnerable to adversarial attacks

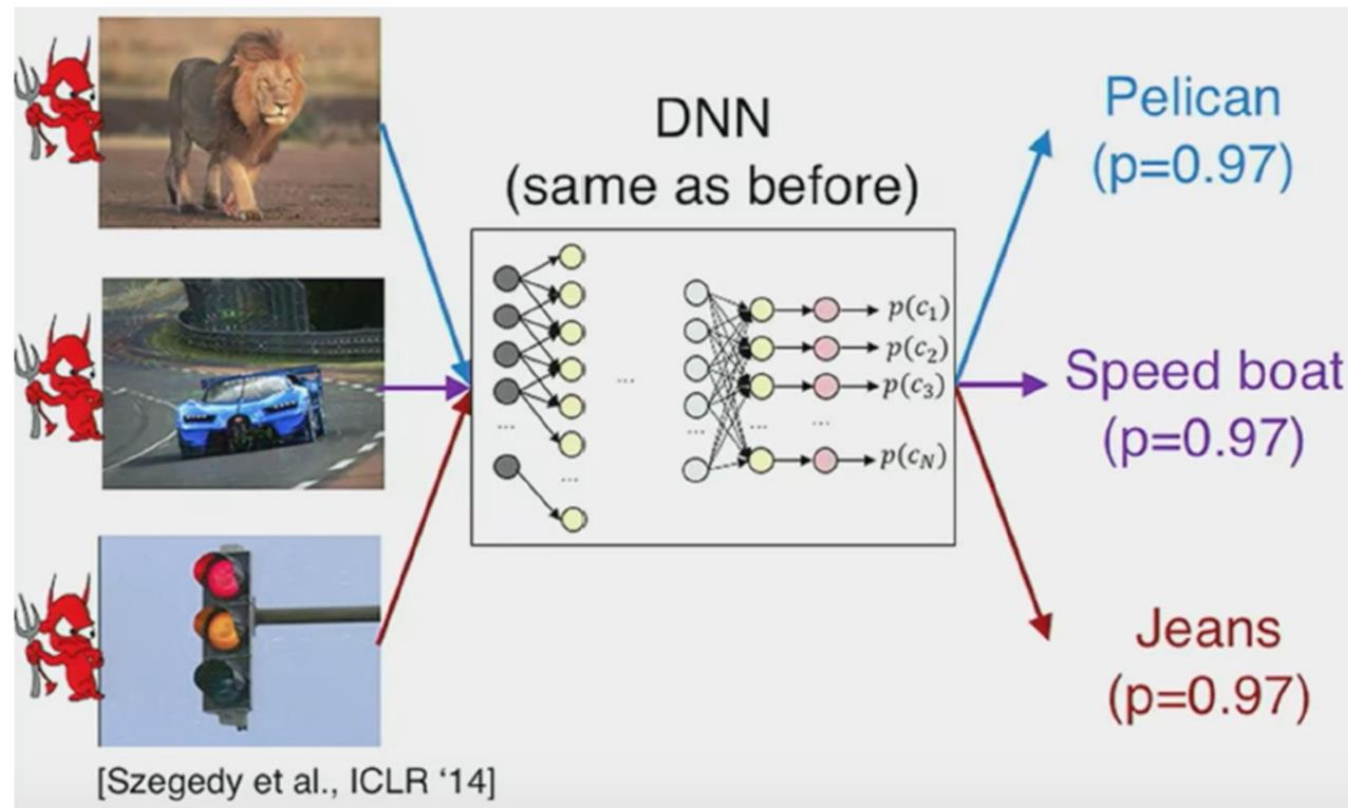


# Adversarial Attacks



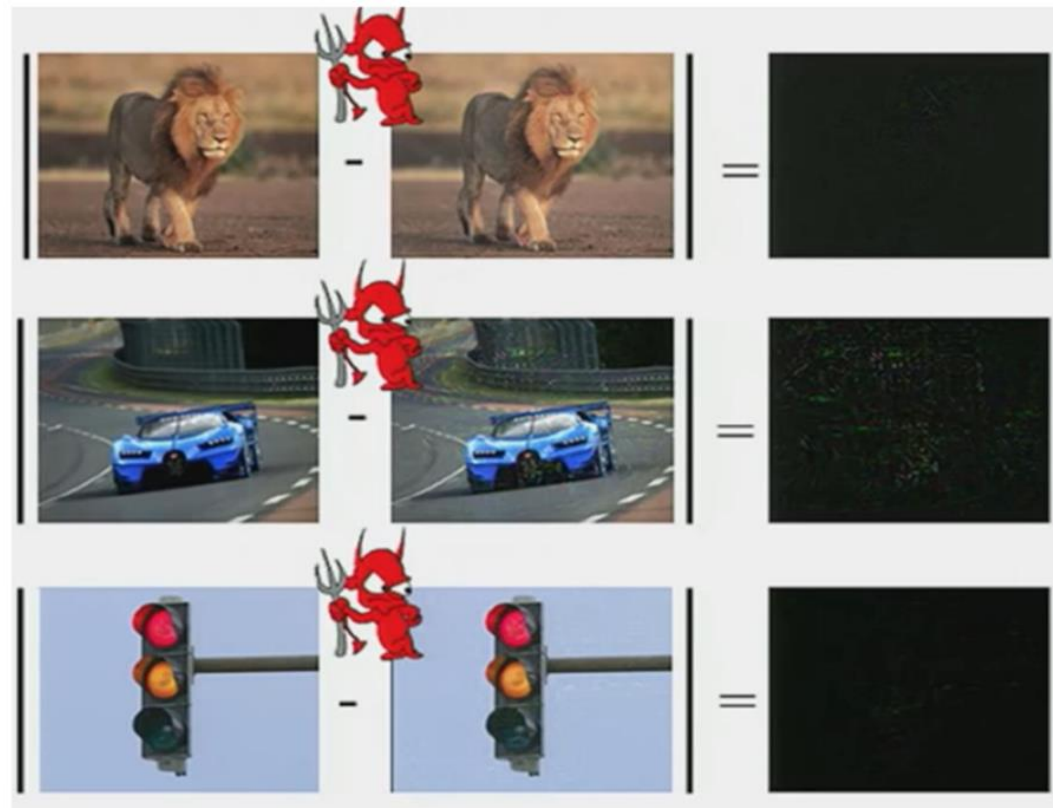
# Adversarial Attacks

- Adversarial example guides classifier to make wrong predictions



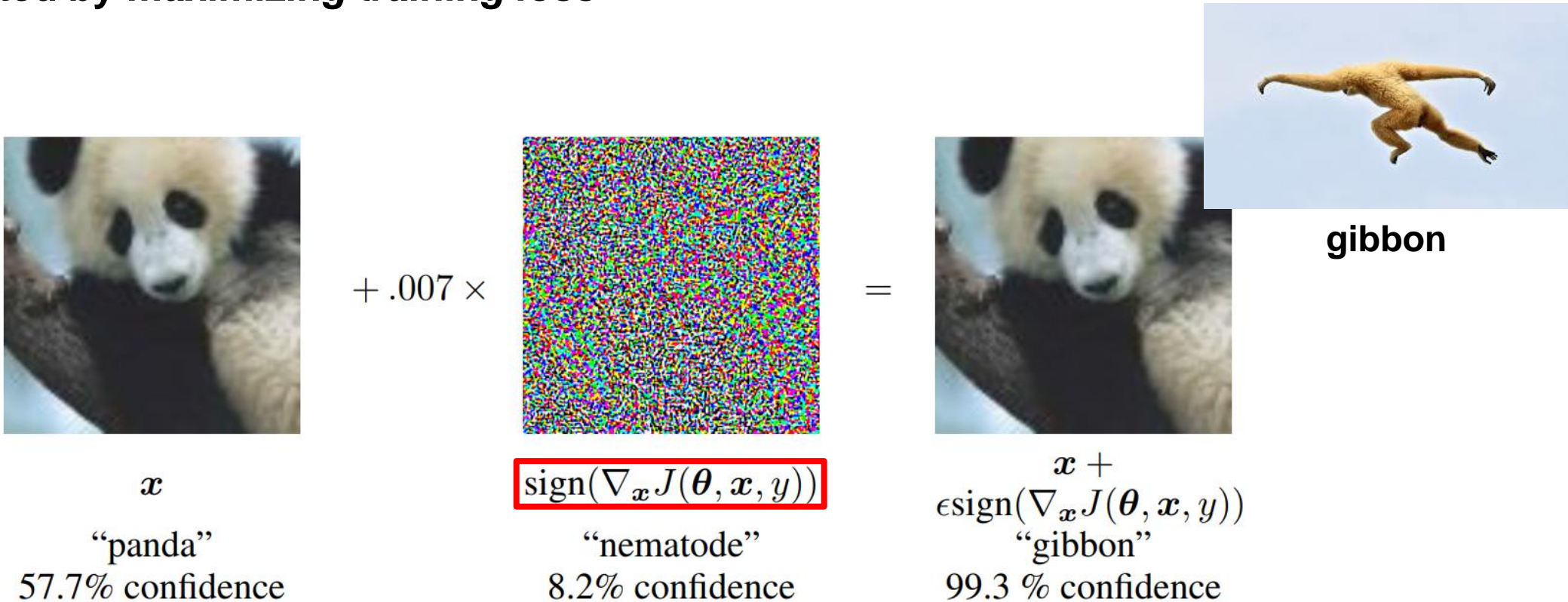
# Adversarial Attacks

- Difference between adversarial image and natural image is hardly noticeable



# How does Adversarial Attack Work?

- Perturbation maliciously designed to fool machine learning models
- Formulated by *maximizing training loss*



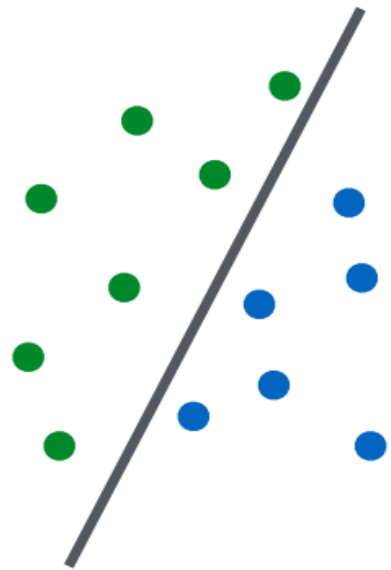
The diagram illustrates the process of an adversarial attack. It shows three stages:

- Original Image:** A panda image labeled  $x$  with a confidence of 57.7% for the class "panda".
- Perturbation:** A small, noisy image labeled  $\epsilon \text{sign}(\nabla_x J(\theta, x, y))$  with a confidence of 8.2% for the class "nematode".
- Adversarial Example:** The original image plus the perturbation, labeled  $x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$ , with a confidence of 99.3% for the class "gibbon".

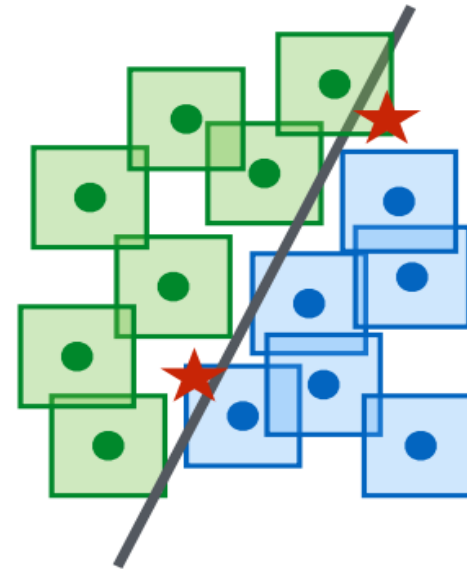
An inset image shows a gibbon in a jumping pose, labeled "gibbon".

# How does Adversarial Attack Work?

- Adversarial examples are designed to cross the decision boundary of models
- The degree of perturbation on the data should be minimal



(a) Trained binary classifier



(b) **Adversarial attack** crossing decision boundary

# Applicability of Adversarial Attacks

- Adversarial examples can exist on possibly any:
  - Deep neural network (MLP, CNN, ViT, ...)
  - Form of data (image, video, point cloud, mesh, ...)
  - Task (classification, localization, generation, ...)

## On the Robustness of Vision Transformers to Adversarial Examples

Kaleel Mahmood

Department of Computer Science and Engineering  
University of Connecticut, CT, 06269, USA [kaleel.mahmood@uconn.edu](mailto:kaleel.mahmood@uconn.edu)

Rigel Mahmood

Department of Computer Science and Engineering  
University of Connecticut, CT, 06269, USA

Marten Van Dijk

CWI, Amsterdam  
The Netherlands

## Generating 3D Adversarial Point Clouds

Chong Xiang

Shanghai Jiao Tong University  
Shanghai, China  
[xiangchong97@gmail.com](mailto:xiangchong97@gmail.com)

Charles R. Qi

Facebook AI Research  
California, USA  
[charlesq34@gmail.com](mailto:charlesq34@gmail.com)

Bo Li

University of Illinois at Urbana-Champaign  
Illinois, USA  
[1xbosky@gmail.com](mailto:1xbosky@gmail.com)

## On the Robustness of ChatGPT: An Adversarial and Out-of-distribution Perspective

Jindong Wang<sup>1,\*</sup>, Xixu Hu<sup>1,2,†</sup>, Wenxin Hou<sup>3,†</sup>, Hao Chen<sup>4</sup>, Runkai Zheng<sup>1,5,‡</sup>, Yidong Wang<sup>6</sup>, Linyi Yang<sup>7</sup>, Wei Ye<sup>6</sup>, Haojun Huang<sup>3</sup>, Xiubo Geng<sup>3</sup>, Binxing Jiao<sup>3</sup>, Yue Zhang<sup>7</sup>, Xing Xie<sup>1</sup>

<sup>1</sup>Microsoft Research, <sup>2</sup>City University of Hong Kong, <sup>3</sup>Microsoft STCA, <sup>4</sup>Carnegie Mellon University, <sup>5</sup>Chinese University of Hong Kong (Shenzhen), <sup>6</sup>Peking University, <sup>7</sup>Westlake University

<https://github.com/microsoft/robustlearn>

# Adversarial Attack on Matching-based Algorithms

---

- There also have been adversarial attacks on matching-based algorithm
  - NeRFool investigates the adversarial vulnerability of Generalizable NeRFs
- 

## **NeRFool: Uncovering the Vulnerability of Generalizable Neural Radiance Fields against Adversarial Perturbations**

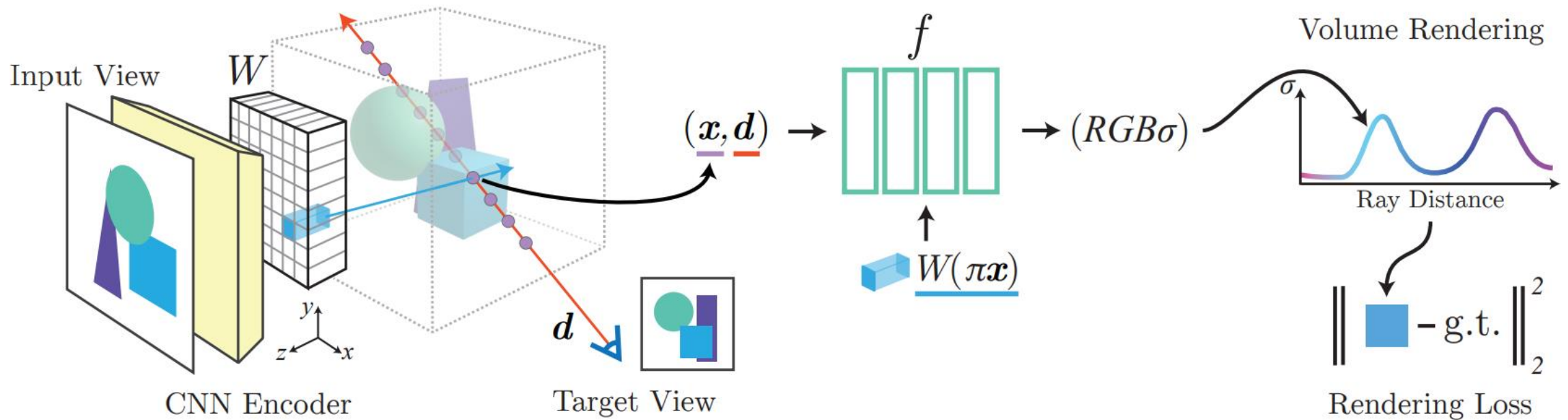
---

**Yonggan Fu<sup>1</sup> Ye Yuan<sup>1</sup> Souvik Kundu<sup>2</sup> Shang Wu<sup>1</sup> Shunyao Zhang<sup>3</sup> Yingyan (Celine) Lin<sup>1</sup>**



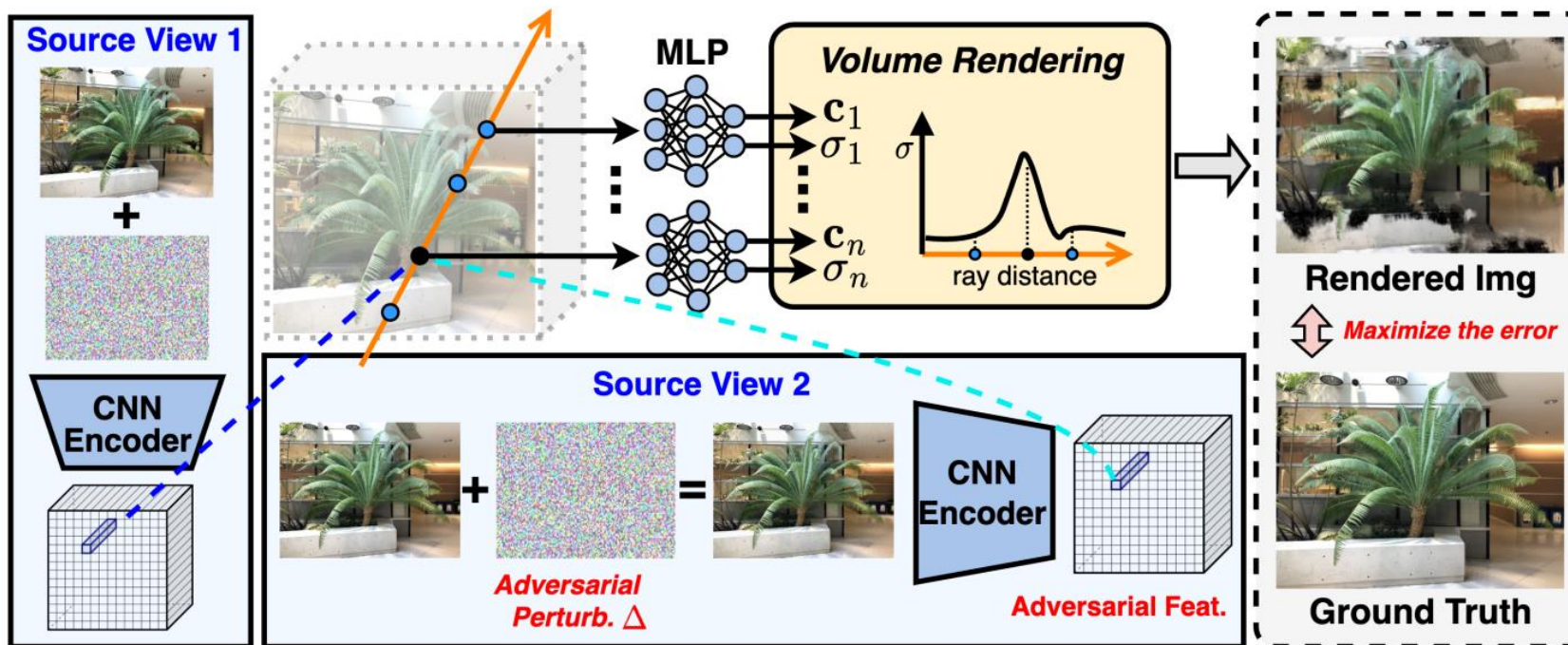
# Generalizable NeRFs

- Originally, a single NeRF is fitted to a single 3D scene or object
- Generalizable NeRFs generalize a single NeRF model to multiple scenes or objects
- They use a set of “support images” to condition the MLP network



# NeRFool

- The goal of NeRFool is to fool the target G-NeRF to render disrupted images
- To do so, it crafts “adversarial support images”





---

# Q&A