

# Universal Adversarial Perturbations Against Semantic Image Segmentation

Samim Zahoor<sup>1</sup>

<sup>1</sup>Department of Information Technology  
National Institute of Technology Srinagar

B.Tech Seminar, 2017

# Outline

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

## 1 Adversarial Perturbation

- What are Adversarial Perturbations?
- What are Universal Adversarial Perturbations?
- How to generate them?

## 2 Key Points and Insights

# Outline

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?  
How to generate  
them?

Key Points  
and Insights

## 1 Adversarial Perturbation

- What are Adversarial Perturbations?
- What are Universal Adversarial Perturbations?
- How to generate them?

## 2 Key Points and Insights

# What are Adversarial Perturbations?

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

- Imperceptible changes to the **inputs** to **deep network classifiers** that cause them to mis-predict labels.

$$\tilde{x} = x + \eta \quad (1)$$

# What are Adversarial Perturbations?

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?  
How to generate  
them?

Key Points  
and Insights

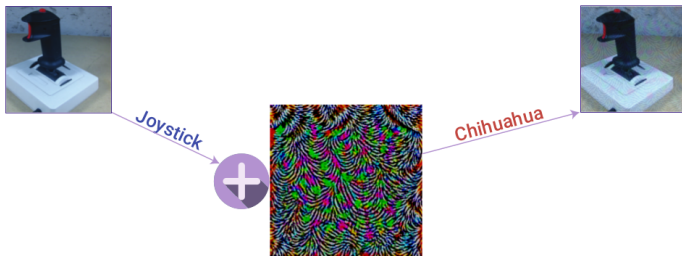


Figure: Adversarial Perturbation

# Outline

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

## 1 Adversarial Perturbation

- What are Adversarial Perturbations?
- What are Universal Adversarial Perturbations?
- How to generate them?

## 2 Key Points and Insights

# Universal Adversarial Perturbations

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

- A single perturbation ( $\eta$ ) that causes vast majority of the images to be mis-classified.
- Image agnostic perturbation.

# What are Adversarial Perturbations?

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

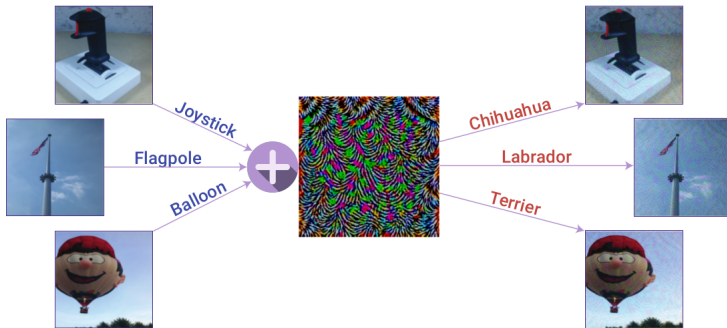


Figure: Universal Adversarial Perturbation



# Outline

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

## 1 Adversarial Perturbation

- What are Adversarial Perturbations?
- What are Universal Adversarial Perturbations?
- How to generate them?

## 2 Key Points and Insights

# Generating Adversarial Perturbations

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

- Fast Gradient Sign Method (FGSM) by Goodfellow et al based on their hypothesis that neural networks are too linear to resist linear adversarial perturbations.

# Generating Adversarial Perturbations

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

- Fast Gradient Sign Method (FGSM) by Goodfellow et al based on their hypothesis that neural networks are too linear to resist linear adversarial perturbations.
- Defined by

$$\eta = \epsilon \text{sign}(\nabla_x J(\Theta, \mathbf{x}, y)). \quad (2)$$

$J$  is the cost function,  $x$  is the input,  $y$  the label and  $\epsilon$  is set to 0.25 for the max norm constraint.

# Generating Adversarial Perturbations

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

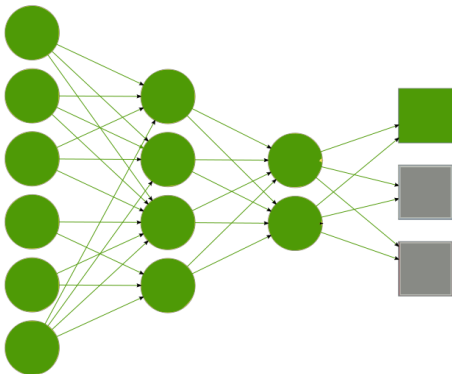
What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights



Forward Pass



# Generating Adversarial Perturbations

Universal  
Adversarial  
Perturbations

Samim Zahoor

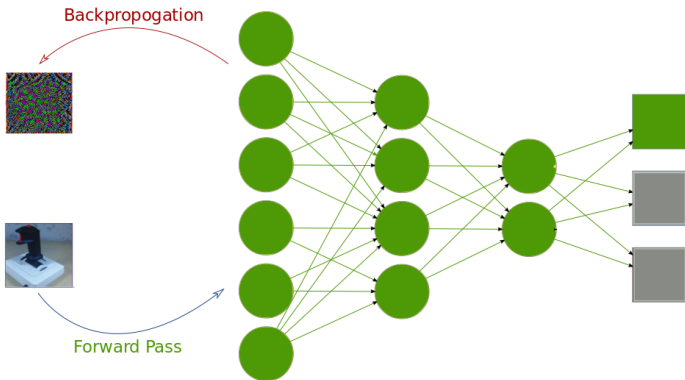
Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights



# Generating Universal Adversarial Perturbation

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

- A Slight Modification of the FGSM proposed by Dezfooli et al. produces image agnostic perturbations.
- They called it **DeepFool**.

# Deepfool

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

- For the first image, DeepFool identifies a standard image-dependent perturbation. For subsequent images, it is checked whether adding the previous adversarial perturbation already fools the classifier; if yes the algorithm continues with the next image, otherwise it updates the perturbation using DeepFool such that also the current image becomes adversarial. The algorithm stops once the perturbation is adversarial on a large fraction of the train set.

# Adversarial Perturbations for Semantic Image Segmentation

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

- Semantic image segmentation denotes a dense prediction task that addresses the what is where in an image? question by assigning a class label to each pixel of the image.



# Adversarial Perturbations for Semantic Image Segmentation

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?  
How to generate  
them?

Key Points  
and Insights

- Semantic image segmentation denotes a dense prediction task that addresses the what is where in an image? question by assigning a class label to each pixel of the image.



**Figure:** Dynamic target segmentation for hiding pedestrians.

- Recently, deep learning based approaches have become the dominant and best performing class of methods for this task and thus susceptible to adversarial examples.

# Generating Universal Adversarial Perturbations

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

- In the context of semantic image segmentation, the universal adversarial perturbations  $\Xi$  are generated on a set of  $m$  training examples.

# Generating Universal Adversarial Perturbations

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?

How to generate  
them?

Key Points  
and Insights

- In the context of semantic image segmentation, the universal adversarial perturbations  $\Xi$  are generated on a set of  $m$  training examples.
- $\Xi$  is defined by

$$\Xi^{(0)} = 0, \quad (3)$$

$$\Xi^{(n)} = \text{Clip}_{\epsilon} \{ \Xi^{(n)} - \alpha \text{sign}(\nabla(\Xi)) \}. \quad (4)$$

With  $\nabla(\Xi)$  being the loss gradient averaged over the entire training data for the semantic image segmentation loss function.

# Key Points and Insights

Universal  
Adversarial  
Perturbations

Samim Zahoor

Adversarial  
Perturbation

What are  
Adversarial  
Perturbations?

What are  
Universal  
Adversarial  
Perturbations?  
How to generate  
them?

Key Points  
and Insights

- Adversarial examples **Generalize** well across different architectures and data.
- Adversarial examples can be explained as a property of high-dimensional dot products. They are a result of models being too **linear, rather than too nonlinear**.
- Linear models lack the capacity to **resist** adversarial perturbation.
- Outlook
  - Future work to address how machine learning can become more robust against (adversarial) perturbations. (Adversarial Training)
  - How adversarial attacks can be detected.