# Synthesizing Robust Adversarial Examples

Anish Athalye\*, Logan Engstrom\*, Andrew Ilyas\*, Kevin Kwok

# Adversarial examples

# Adversarial examples

• Imperceptible perturbations to an input can change a neural network's prediction



88% tabby cat





**99% guacamole** 

# Adversarial examples

**Given:** Input image *x*, target label *y* 

**Optimize:** 

arg max  $\mathbf{X}'$ 

 $P\left(y \mid \mathbf{x}'\right)$ subject to  $d(\mathbf{x}, \mathbf{x}') < \epsilon$ 



#### step 00



# Do adversarial examples work in the physical world?

## Adversarial examples in the physical world



#### (a) Image from dataset

(b) Clean image

(c) Adv. image,  $\epsilon = 4$ 

(Kurakin et al. 2016)



# Before Foveation After Foveation

Foveation-based Mechanisms Alleviate Adversarial Examples (Luo et al. 2015)

## ... or not?



NO Need to Worry about Adversarial Examples in Object Detection in Autonomous Vehicles (Lu et al. 2017)



# Standard examples are fragile

#### Zoom: 1.000000x





# Are adversarial examples fundamentally fragile?

# Image processing pipeline



## optimize $P(y | \mathbf{x}')$ using gradient descent

# Physical world processing pipeline





these are randomized

Challenge: No direct control over model input

# **Attack: Expectation Over Transformation**



## optimize $\mathbb{E}_{t \sim T} \left[ P(y \mid t(\mathbf{x}')) \right]$ using gradient descent (sampling, chain rule, differentiating through t)

## is differentiable



#### these are randomized but the distribution T is known

# EOT produces robust examples



 $T = \{rescale from 1x to 5x\}$ 

# EOT produces robust physical-world examples





### T = {rescale + rotate + translate + skew}

# Can we make this work with 3D objects?

# Physical world 3D processing pipeline





PARAMETERS







is this differentiable?

MODEL

## PREDICTIONS







# Differentiable rendering



- Simplest renderer: linear transformation of texture

• For any pose, 3D rendering is differentiable with respect to texture

# EOT produces 3D adversarial objects









# ff, drop, lith, me agama ega U





![](_page_20_Picture_1.jpeg)

# EOT reliably produces 3D adversarial objects

	Inputs	Classification accuracy	Attacker success rate	<b>Distortion (I2)</b>
2D	Original	70%	N/A	0
	Adversarial	0.9%	96.4%	5.6 × 10-5
3D	Original	84%	N/A	0
	Adversarial	1.7%	84.0%	6.5 × 10 <sup>-5</sup>

# Implications

- Defenses based on randomized input transformations are insecure
- Adversarial examples / objects are a physical-world concern

**Poster (and live demo): 6:15 – 9:00pm @ Hall B #73**