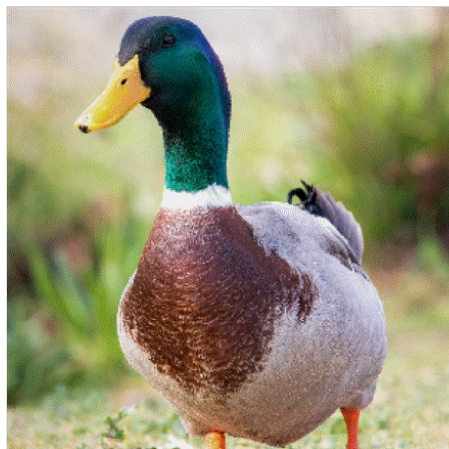# Adversarial

Eric Wong
9/8/2022

# Noise attack
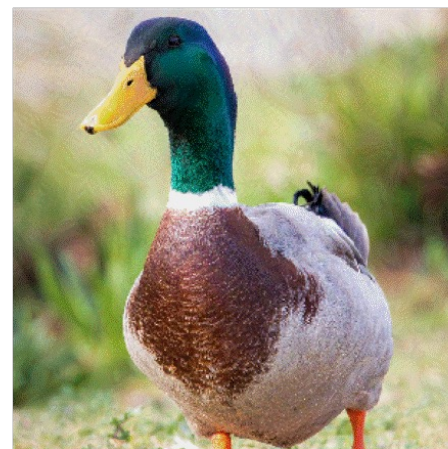


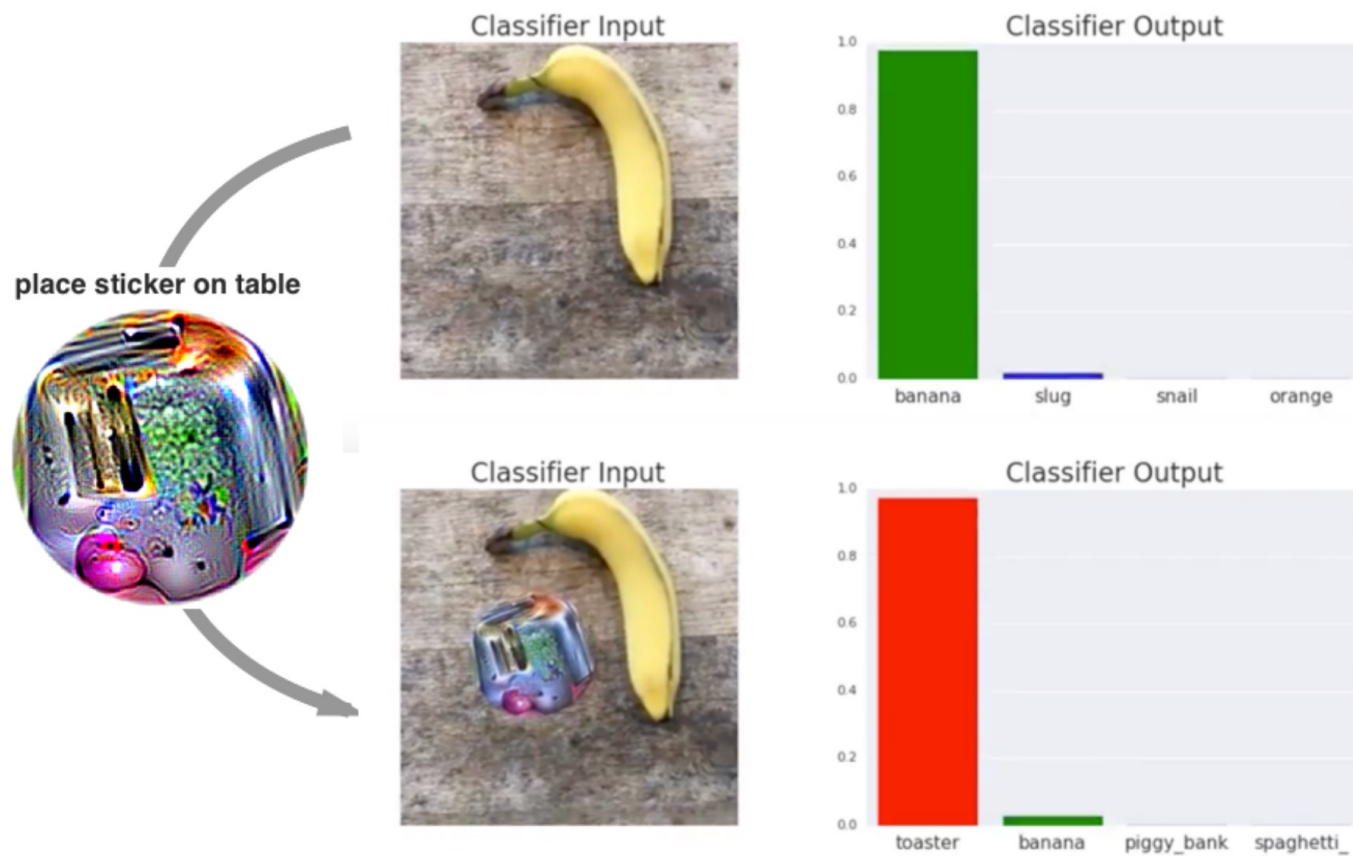"Duck"   +   =   "Hermit crab"

# Patch attack



place sticker on table

Classifier Input

Classifier Output

banana    slug    snail    orange

Classifier Input

Classifier Output

toaster    banana    piggy_bank    spaghetti_

Brown et al. 2017 "Adversarial Patch"

# 3D printed textures



classified as turtle     classified as rifle

classified as other

Athalye et al. 2017 "Synthesizing Robust Adversarial Examples"

# Glasses



(a)  (b)  (c)

Sharif et al. 2018 "A General Framework for Adversarial Examples with Objectives"
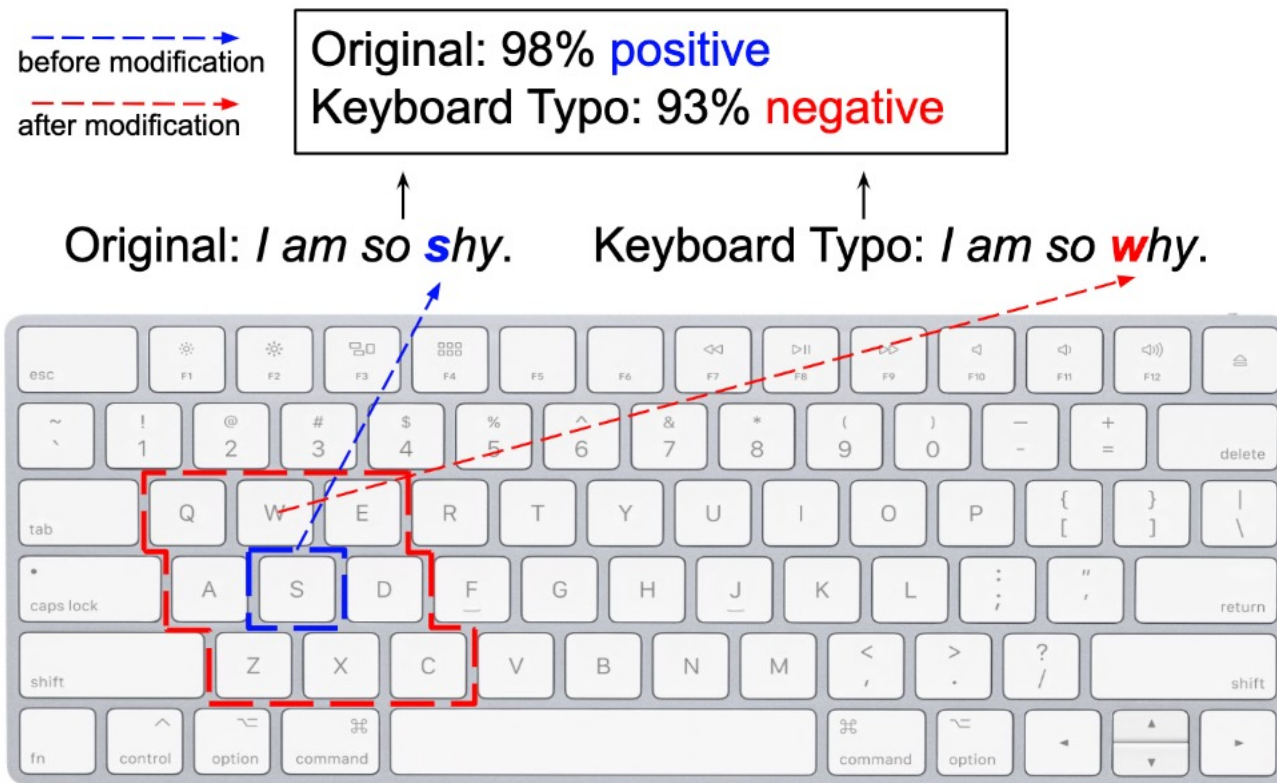
# Clothing



Wu et al. 2019 "Making an Invisibility Cloak: Real World Adversarial Attacks on Object Detectors"

# Camera stickers



Li et al. 2019 "Adversarial camera stickers: A physical camera based attack on deep learning sytems"

# Typos



before modification

after modification

Original: 98% positive
Keyboard Typo: 93% negative

Original: *I am so **shy**.*    Keyboard Typo: *I am so **why**.*

Sun et al. 2020 "Adv-BERT: BERT is not robust on misspellings! Generating nature adversarial samples on BERT"

# Sentences

| Label | Sentence |
|:---:|:---|
| P | I am currently trying to give this company another chance. I have had the same scheduling experience as others have written about. Wrote to them today |
| N | I am currently trying to give this company another *review*. I have had the same *dental experience about others or written with a name. Thanks* to them today |

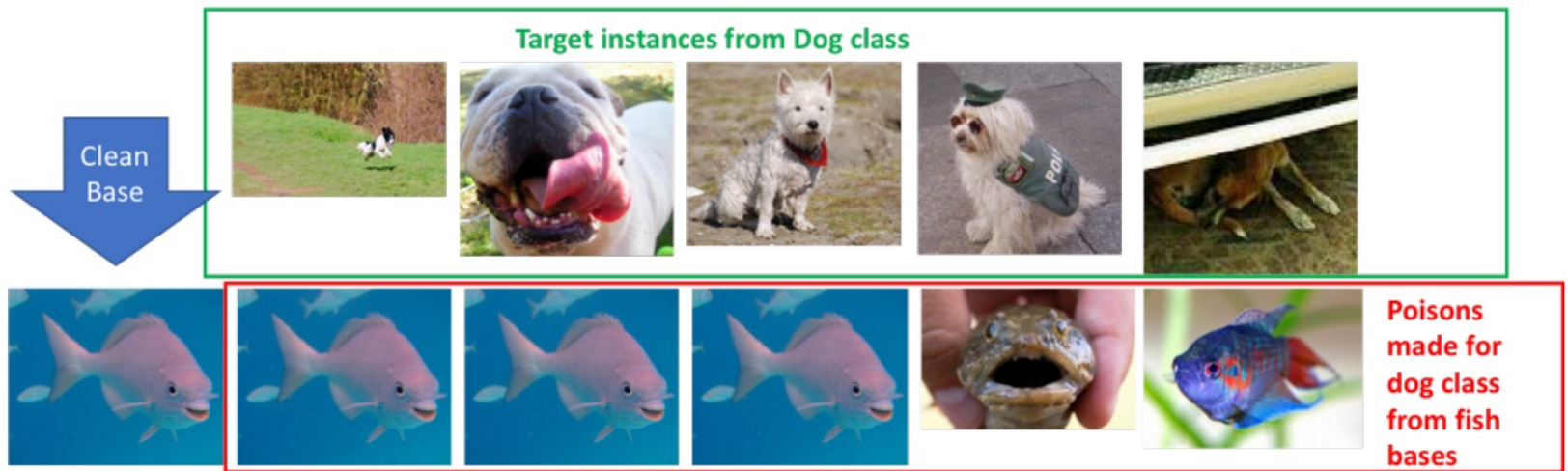Hsieh et al. 2019 "Natural Adversarial Sentence Generation with Graident based Perturbation"

# Speech recognition



Abdullah et al. 2020

# Data Poisoning

Eric Wong
9/22/2022

# "One-shot" poison



Target instances from Fish class

Clean Base

Poison instances made for fish class from dog base instances

Shafahi et al. 2020 "Poison Frogs! Targeted Clean-Label Poisoning Attacks on Neural Networks"

# "One-shot" poison



Shafahi et al. 2020 "Poison Frogs! Targeted Clean-Label Poisoning Attacks on Neural Networks"

# Meta Poisoning



Huang et al. 2020 "MetaPoison: Practical General-purpose Clean-label Data Poisoning"

# Backdoor triggers



Clean  Yellow Square  Bomb  Flower

Gu et al. 2017 "BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain"

# Hidden backdoors



Generating poison · Training · Testing

Saha et al. 2019 "Hidden Trigger Backdoor Attacks"

# Real & robust backdoors



**Generating Backdoor Face Images**

**Injecting Backdoor Face Images & Performing Transformations**

**Triggering Physical Backdoor Attacks**

# Distributed backdoor