

# Evasion Attacks Against Machine Learning Models (Non-traditional Attacks)

# Recall: Adversarial Examples

- FGSM

$$\boldsymbol{\eta} = \epsilon \text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$

- Optimization based attack

$$\begin{aligned} \min d(x, x') + g(x') \\ \text{s.t. } x' \text{ is "valid"} \end{aligned}$$

- DeepFool

- Greedy algorithm to move the instance towards the nearest boundary

- JSMA (Jacobian-based Saliency Map Approach)

- Compute the saliency map for an  $X$  regarding to target  $y^*$ ; modify the max pixel each time

- BIM (Basic Iterative Method)

- Apply FGSM multiple times with small step size

# Recall: Interesting topic: how to analyze transferability?

Lower bound of adversarial transferability:

**Lemma 1.** *Let  $f, g : \mathcal{X} \rightarrow \mathcal{Y}$  be classifiers,  $\delta, \rho, \epsilon \in (0, 1)$  be constants, and  $\mathcal{A}(\cdot)$  be an attack strategy. Suppose that  $\mathcal{A}(\cdot)$  is  $\rho$ -covert and  $f, g$  have risk at most  $\epsilon$ . Then  $\Pr(f(\mathcal{A}(x)) \neq g(\mathcal{A}(x))) \leq 2\epsilon + \rho$  for a random instance  $x \sim P_X$ .*

**Theorem 3.** *Let  $f, g : \mathcal{X} \rightarrow \mathcal{Y}$  be classifiers ( $\mathcal{Y} \in \{-1, 1\}$ ),  $\delta, \rho, \epsilon \in (0, 1)$  be constants, and  $\mathcal{A}(\cdot)$  an attack strategy. Suppose that  $\mathcal{A}(\cdot)$  is  $\rho$ -covert and  $f, g$  have risk at most  $\epsilon$ . Given random instance  $x \in \mathcal{X}$ , if  $\mathcal{A}(\cdot)$  is  $(\delta, g)$ -effective, then it is also  $(\delta + 4\epsilon + \rho, f)$ -effective.*

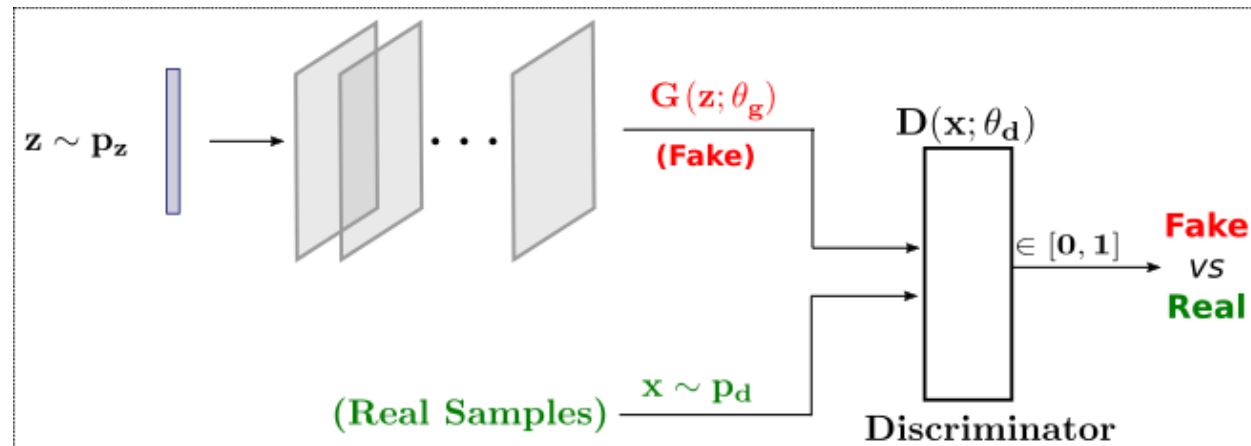
Lower bound of adversarial transferability?

# Generating Adversarial Examples with Adversarial Networks

- How can we generate more realistic adversarial examples?
- How can we generate diverse adversarial examples?
- How to perform blackbox attack efficiently?

# Generating Adversarial Examples with Adversarial Networks

- Generative adversarial networks (GANs)

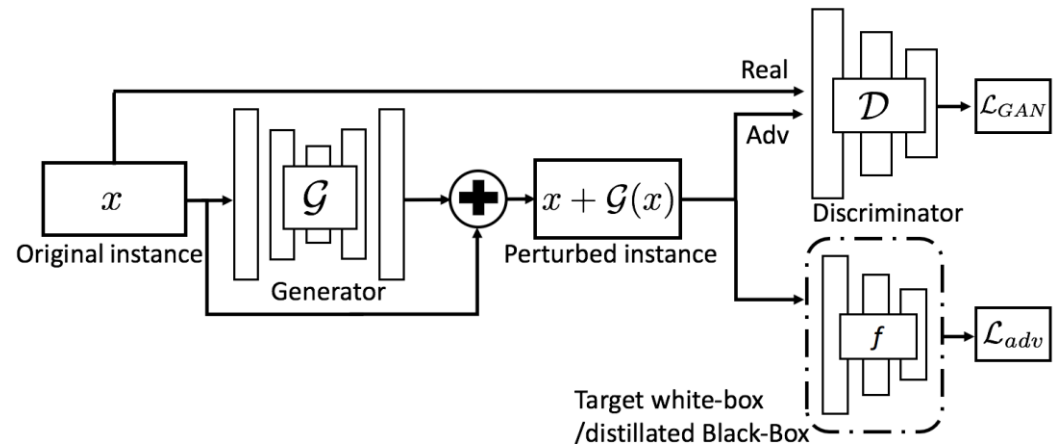


- Generate more realistic instances
- Approximate certain distribution
- Efficient once the generator is trained

## Questions:

1. Can we generate more realistic adversarial examples?
2. Can we generate adversarial examples more efficiently?

# Generating Adversarial Examples with Adversarial Networks



Black-box can be performed here via distillation

$$\mathcal{L} = \mathcal{L}_{adv}^f + \alpha \mathcal{L}_{GAN} + \beta \mathcal{L}_{hinge}$$

$$\mathcal{L}_{GAN} = \mathbb{E}_{x \sim \mathcal{P}_{data}(x)} \log \mathcal{D}(x) + \mathbb{E}_{x \sim \mathcal{P}_{data}(x)} \log(1 - \mathcal{D}(x + \mathcal{G}(x)))$$

$$\mathcal{L}_{adv}^f = \mathbb{E}_x \ell_f(x + \mathcal{G}(x), t)$$

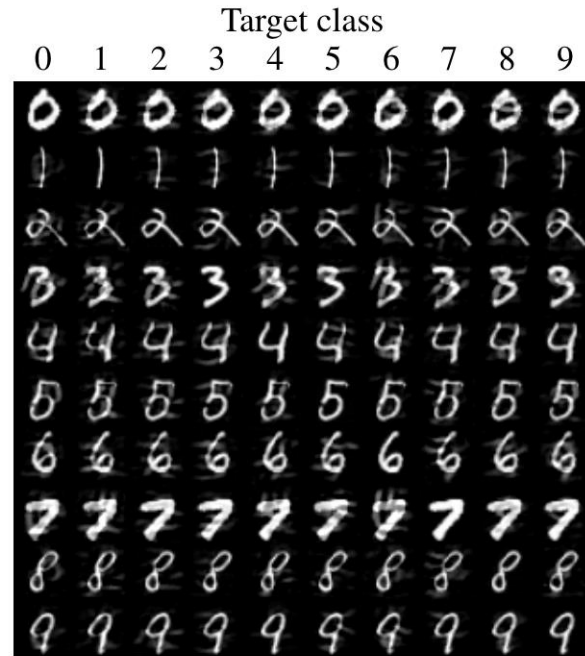
$$\mathcal{L}_{hinge} = \mathbb{E}_x \max(0, \|\mathcal{G}(x)\|_2 - c)$$

# Generating Adversarial Examples with Adversarial Networks

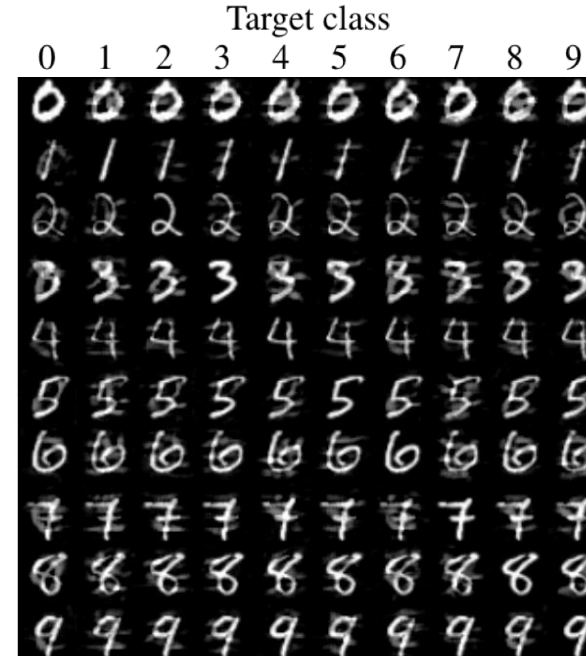
- Advantages

	FGSM	Opt.	Trans.	AdvGAN
Run time	0.06s	>3h	-	<0.01s
Targeted Attack	✓	✓	Ens.	✓
Black-box Attack			✓	✓

# Generating Adversarial Examples with Adversarial Networks



Semi-white box attack on MNIST

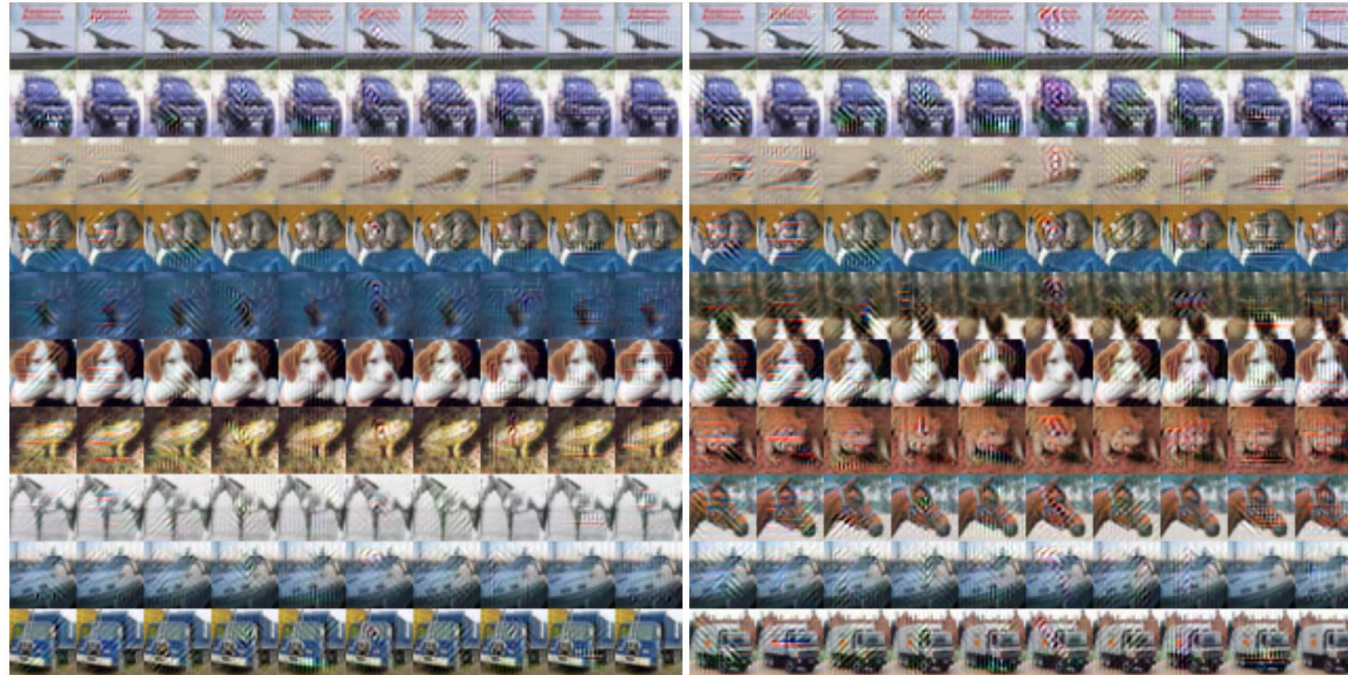


Black-box attack on MNIST

The perturbed images are very close to the original ones. The original images lie on the diagonal.



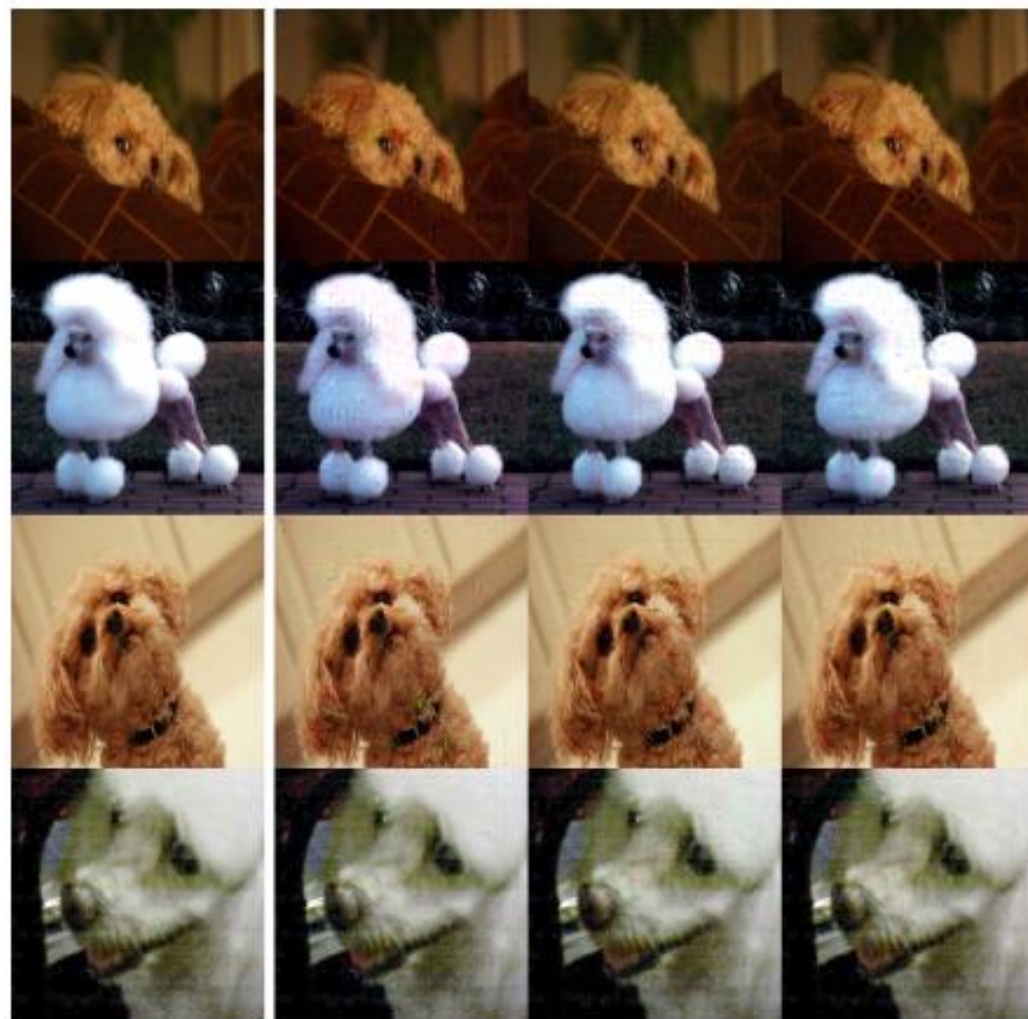
# Generating Adversarial Examples with Adversarial Networks



(a) Semi-whitebox setting

(b) Black-box setting

The perturbed images are very close to the original ones. The original images lie on the diagonal.



Poodle

Ambulance

Basketball

Electric guitar



(a) Strawberry



(b) Toy poodle



(c) Buckeye



(d) Toy poodle

# Attack Effectiveness Under Defenses

Data	Model	Defense	FGSM	Opt.	AdvGAN
MNIST	A	Adv.	4.3%	4.6%	<b>8.0%</b>
		Ensemble	1.6%	4.2%	<b>6.3%</b>
		Iter.Adv.	4.4%	2.96%	<b>5.6%</b>
	B	Adv.	6.0%	4.5%	<b>7.2%</b>
		Ensemble	2.7%	3.18%	<b>5.8%</b>
		Iter.Adv.	<b>9.0%</b>	3.0%	6.6%
	C	Adv.	2.7%	2.95%	<b>18.7%</b>
		Ensemble	1.6%	2.2%	<b>13.5%</b>
		Iter.Adv.	1.6%	1.9%	<b>12.6%</b>
CIFAR	ResNet	Adv.	13.10%	11.9%	<b>16.03%</b>
		Ensemble.	10.00%	10.3%	<b>14.32%</b>
		Iter.Adv	22.8%	21.4%	<b>29.47%</b>
	Wide ResNet	Adv.	5.04%	7.61%	<b>14.26%</b>
		Ensemble	4.65%	8.43%	<b>13.94 %</b>
		Iter.Adv.	14.9%	13.90%	<b>20.75%</b>

Attack success rate of adversarial examples generated by AdvGAN in semi-whitebox setting under defenses on MNIST and CIFAR-10

# Attack Effectiveness Under Defenses

## Black-Box Leaderboard (Original Challenge)

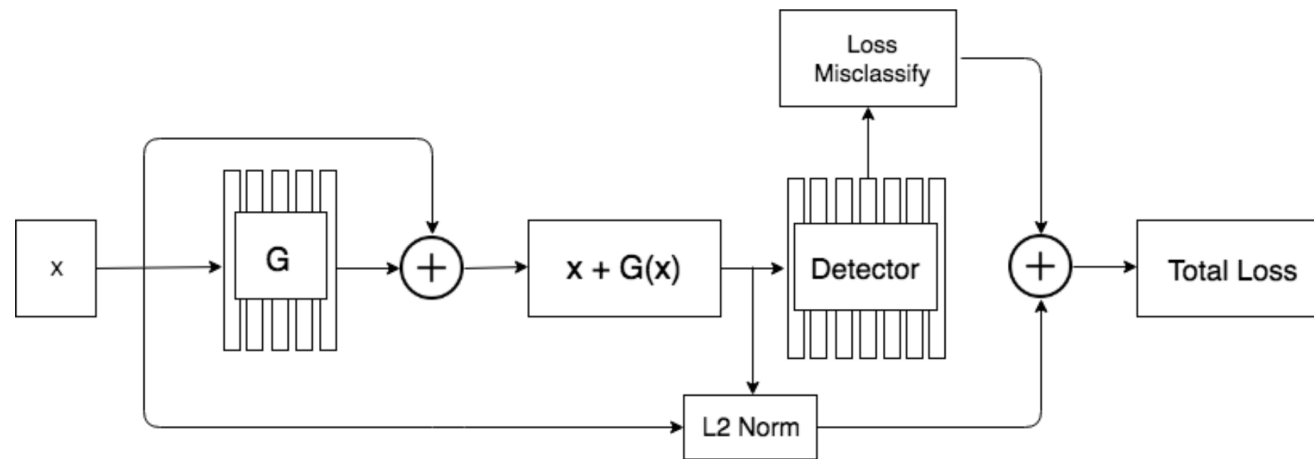
Attack	Submitted by	Accuracy	Submission Date
AdvGAN from <a href="#">"Generating Adversarial Examples with Adversarial Networks"</a>	AdvGAN	92.76%	Sep 25, 2017
PGD against three independently and adversarially trained copies of the network	<a href="#">Florian Tramèr</a>	93.54%	Jul 5, 2017
FGSM on the CW loss for model B from <a href="#">"Ensemble Adversarial Training [...]"</a>	<a href="#">Florian Tramèr</a>	94.36%	Jun 29, 2017
FGSM on the CW loss for the naturally trained public network	(initial entry)	96.08%	Jun 28, 2017
PGD on the cross-entropy loss for the naturally trained public network	(initial entry)	96.81%	Jun 28, 2017
Attack using Gaussian Filter for selected pixels on the adversarially trained public network	Anonymous	97.33%	Aug 27, 2017
FGSM on the cross-entropy loss for the adversarially trained public network	(initial entry)	97.66%	Jun 28, 2017
PGD on the cross-entropy loss for the adversarially trained public network	(initial entry)	97.79%	Jun 28, 2017

# Takeaways

- Adversarial examples and generative adversarial networks are different
- We can integrate them together to work better
- Generative models can indeed synthesize new types of adversarial examples
- Adversarial retraining based defense is not enough

# Similar work

- Adversarial Attacks on Face Detectors using Neural Net based Constrained Optimization

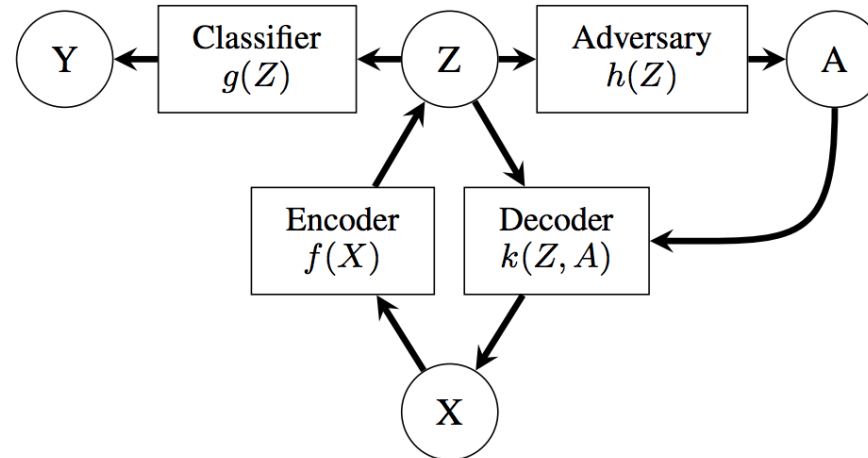


$$L_G(x, x') = \|x - x'\|_2^2 + \lambda \sum_{i=1}^N (Z(x'_i)_{\text{background}} - Z(x'_i)_{\text{face}})^+$$

Difference: attacking detector, face detection task

# Similar work

- Learning Adversarially Fair and Transferable Representations
  - Advocate representation learning as the key to mitigating unfair prediction
  - Propose and explore adversarial representation learning as a natural method of ensuring third parties act fairly



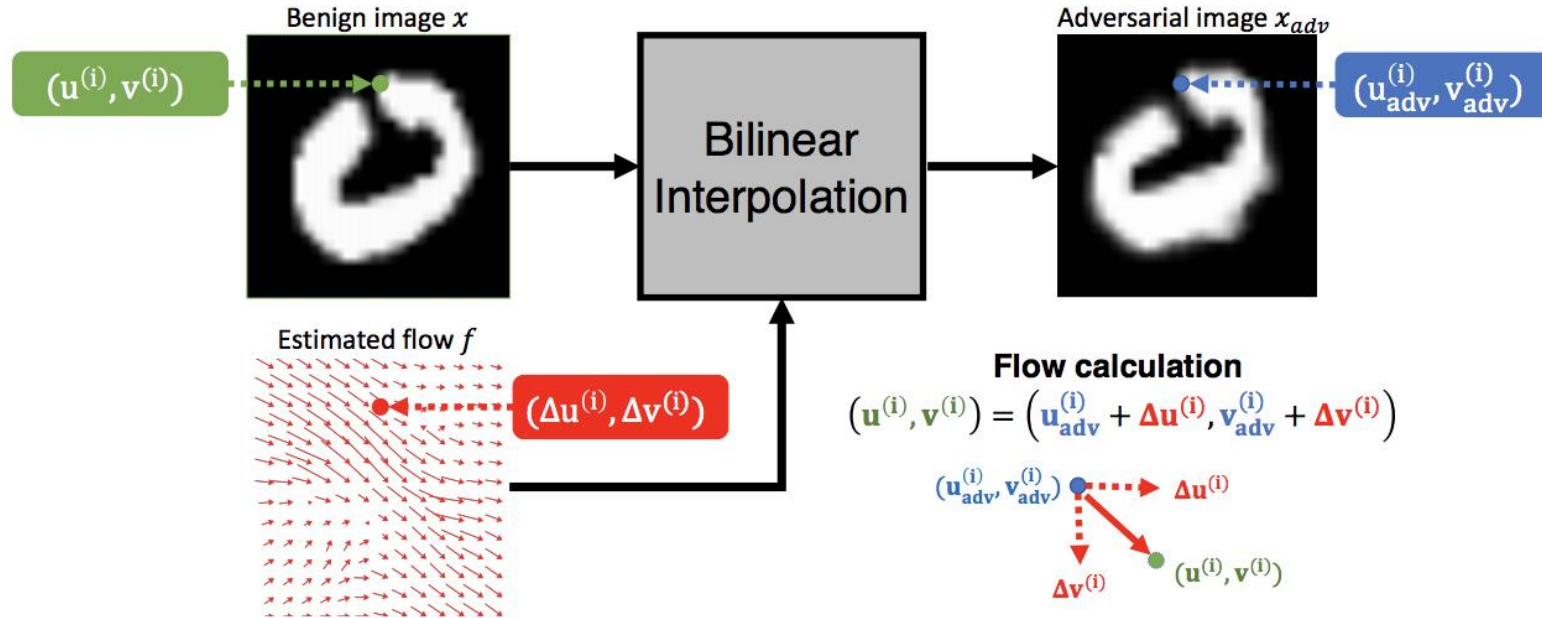
Difference: explore the fairness of machine learning from adversarial learning aspect; nice definition of fairness and theoretic analysis



# Spatially Transformed Adversarial Examples

- Realistic attacks are possible with generative models
- What if we do not directly manipulate the value of pixels?
- What else can we modify? (2D, 3D)
- Potential topic: how to attack 3D point clouds?

# Spatially Transformed Adversarial Examples



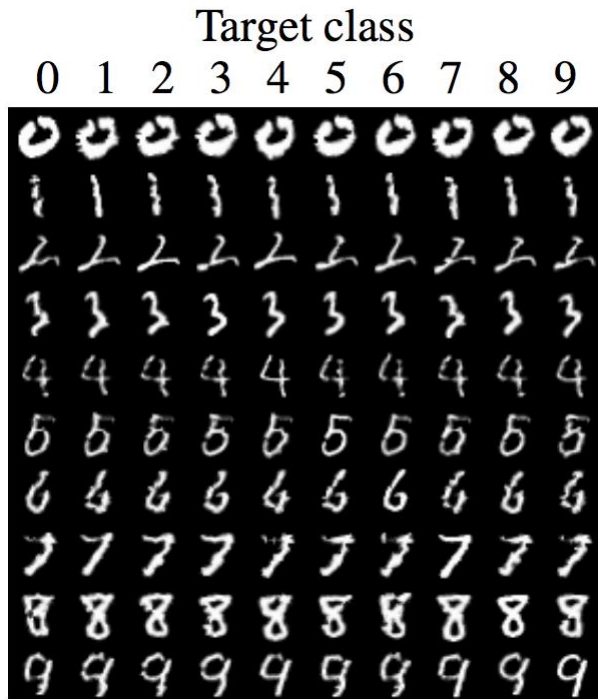
$$f^* = \operatorname{argmin}_f \mathcal{L}_{adv}(x, f) + \tau \mathcal{L}_{flow}(f),$$

$$\mathcal{L}_{adv}(x, f) = \max_{i \neq t} (\max g(\mathbf{x}_{adv})_i - g(\mathbf{x}_{adv})_t, \kappa)$$

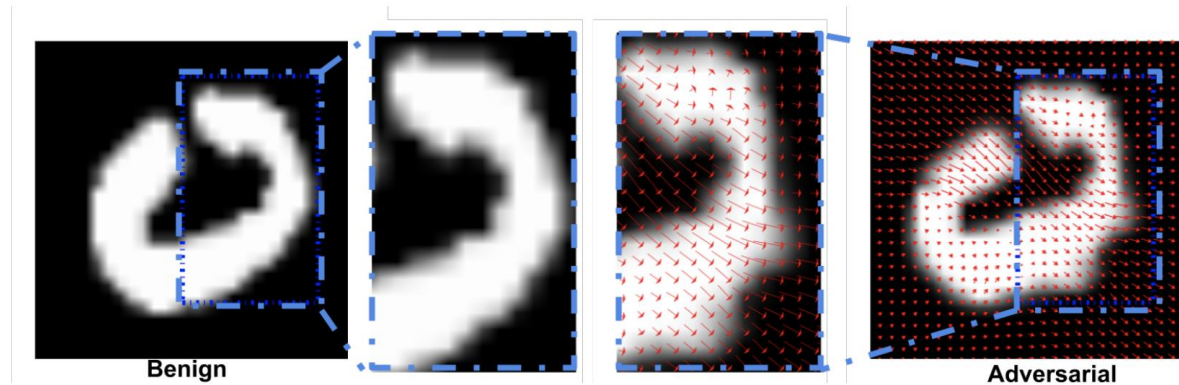
$$\mathcal{L}_{flow}(f) = \sum_p \sum_{q \in \mathcal{N}(p)} \sqrt{\|\Delta u^{(p)} - \Delta u^{(q)}\|_2^2 + \|\Delta v^{(p)} - \Delta v^{(q)}\|_2^2}.$$

$$\mathbf{x}_{adv}^{(i)} = \sum_{q \in \mathcal{N}(u^{(i)}, v^{(i)})} \mathbf{x}^{(q)} (1 - |u^{(i)} - u^{(q)}|) (1 - |v^{(i)} - v^{(q)}|)$$

# Examples generated by stAdv



Adversarial examples generated by stAdv on MNIST  
The ground truth images are shown in the diagonal



Flow visualization on MNIST. The digit "0" is misclassified as "2".

# Attack Effectiveness Under Defenses

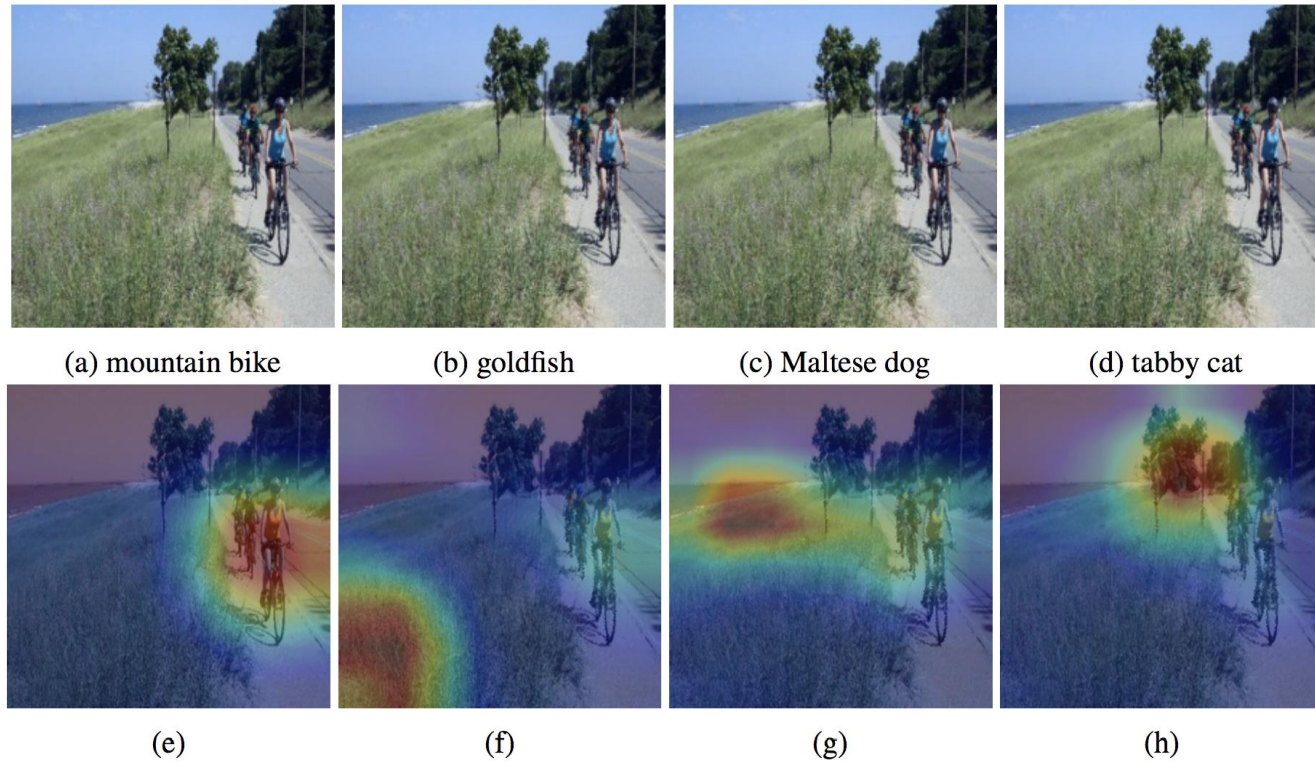
Model	Def.	FGSM	C&W.	stAdv
A	Adv.	4.3%	4.6%	<b>32.62%</b>
	Ens.	1.6%	4.2%	<b>48.07%</b>
	PGD	4.4%	2.96%	<b>48.38%</b>
B	Adv.	6.0%	4.5%	<b>50.17%</b>
	Ens.	2.7%	3.18%	<b>46.14%</b>
	PGD	9.0%	3.0%	<b>49.82%</b>
C	Adv.	3.22%	0.86%	<b>30.44%</b>
	Ens.	1.45%	0.98%	<b>28.82%</b>
	PGD	2.1%	0.98%	<b>28.13%</b>

Model	Def.	FGSM	C&W.	stAdv
ResNet32	Adv.	13.10%	11.9%	<b>43.36%</b>
	Ens.	10.00%	10.3%	<b>36.89%</b>
	PGD	22.8%	21.4%	<b>49.19%</b>
wide ResNet34	Adv.	5.04%	7.61%	<b>31.66%</b>
	Ens.	4.65%	8.43%	<b>29.56%</b>
	PGD	14.9%	13.90%	<b>31.6%</b>

Attack success rate of adversarial examples generated by stAdv against different models under standard defense on MNIST and CIFAR-10

# Attention of Networks



**CAM attention visualization for ImageNet inception\_v3 model. (a) the original image and (b)-(d) are stAdv adversarial examples targeting different classes. Row 2 shows the attention visualization for the corresponding images above.**

inception\_v3 model



(a) Benign

(b) FGSM

(c) C&W

(d) StAdv

Adversarial trained  
inception\_v3 model



(e) Benign

(f) FGSM

(g) C&W

(h) StAdv

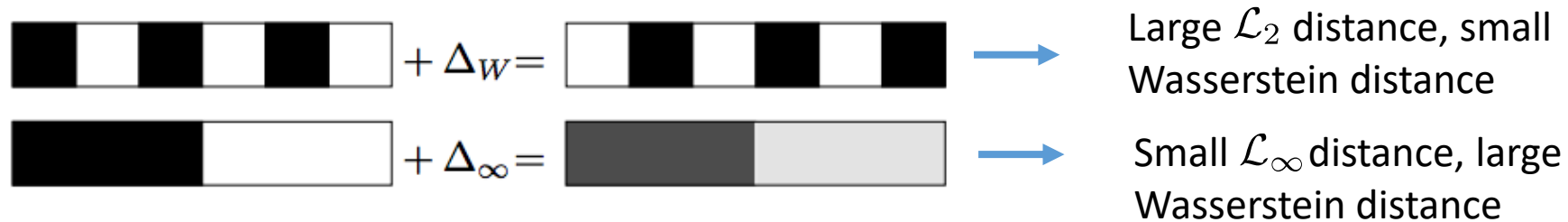
**CAM attention visualization for ImageNet inception\_v3 model. Column 1 shows the CAM map corresponding to the original image. Column 2-4 show the adversarial examples generated by different methods. (a) and (e)-(g) are labeled as the ground truth “cinema”, while (b)-(d) and (h) are labeled as the adversarial target “missile.”**

# Takeaways

- Instead of manipulating the pixel values, we can also move the position of pixels to generate adversarial examples for 2D images
- For 3D, you can add points, what else?
- It is impossible to tell/detect adversarial perturbation from network attention
- A lot of diverse adversarial examples can be explored

# Wasserstein Adversarial Examples via Projected Sinkhorn Iterations

- Another type of *spatial transformed* adversarial examples
- Beyond  $\mathcal{L}_p$  norm-bounded perturbation – Wasserstein distance
- Generate adversarial examples by projecting onto the Wasserstein ball based on Sinkhorn iteration






# Wasserstein Adversarial Examples via Projected Sinkhorn Iterations

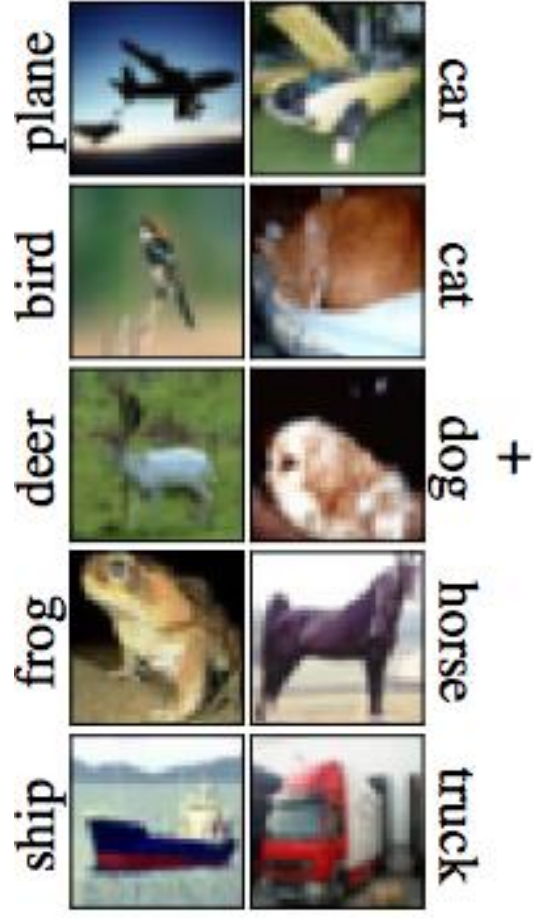
- Wasserstein distance: “earth mover’s distance”, the minimum cost of moving probability mass to change one distribution into another

$$x^{(t+1)} = \underset{\mathcal{B}(x, \epsilon)}{\text{proj}} \left( x^{(t)} + \arg \max_{\|v\| \leq \alpha} v^T \nabla \ell(x^{(t)}, y) \right)$$

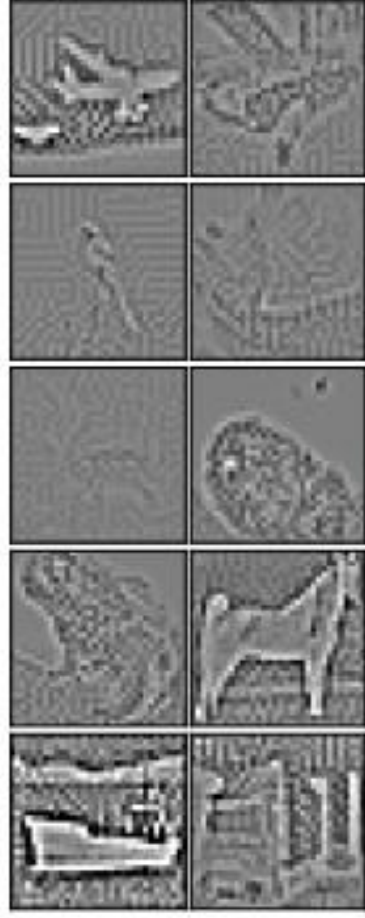
 Wasserstein ball

- Approximate the W-distance with entropy regularization on the transportation plan  $\Pi$  using Sinkhorn-Knopp matrix scaling.

$$\begin{aligned} & \underset{z \in \mathbb{R}_+^n, \Pi \in \mathbb{R}_+^{n \times n}}{\text{minimize}} && \frac{1}{2} \|w - z\|_2^2 + \frac{1}{\lambda} \sum_{ij} \Pi_{ij} \log(\Pi_{ij}) \\ & \text{subject to} && \Pi \mathbf{1} = x, \quad \Pi^T \mathbf{1} = z \\ & && \langle \Pi, C \rangle \leq \epsilon. \end{aligned}$$



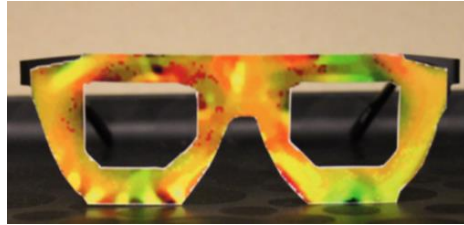
+



=



# Physical Attacks In Practice



Physical attack: Sharif et al., "Accessorize to a crime: real and stealthy attacks on state-of-the-art face recognition," CCS 2016



However, What We Can See Everyday...

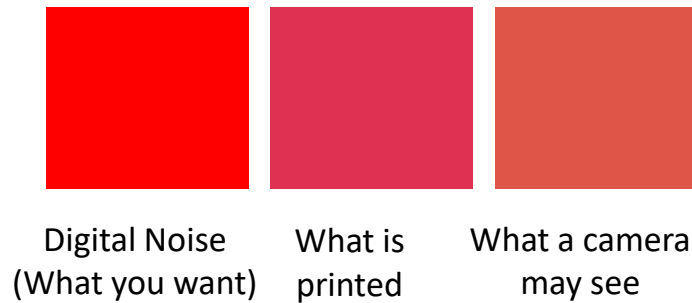


# The Physical World Is... Messy

Varying Physical Conditions (Angle, Distance, Lighting, ...) Physical Limits on Imperceptibility



Fabrication/Perception Error (Color Reproduction, etc.)



Background Modifications\* Image Courtesy, OpenAI



# An Optimization Approach To Creating Robust Physical Adversarial Examples

$$\operatorname{argmin}_{\delta} \lambda \|\delta\|_p + J(f_{\theta}(x + \delta), y^*)$$

Perturbation/Noise Matrix  $\rightarrow$   $\lambda \|\delta\|_p$   $\rightarrow$   $J(f_{\theta}(x + \delta), y^*)$   $\rightarrow$  Adversarial Target Label

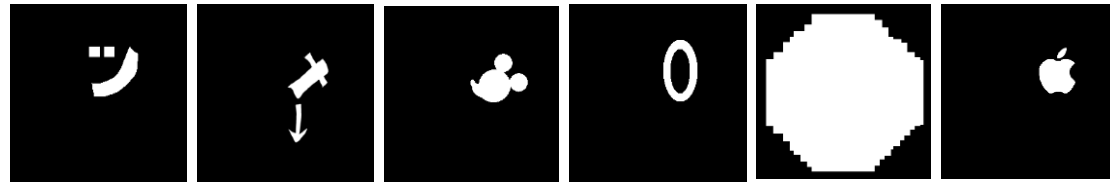
Lp norm (L-0, L-1, L-2, ...)      Loss Function

$$\operatorname{argmin}_{\delta} \lambda \|\delta\|_p + \frac{1}{k} \sum_{i=1}^k J(f_{\theta}(x_i + \delta), y^*)$$



# Optimizing Spatial Constraints (Handling Limits on Imperceptibility)

$$\operatorname{argmin}_{\delta} \lambda \|M_x \cdot \delta\|_p + \frac{1}{k} \sum_{i=1}^k J(f_{\theta}(x_i + M_x \cdot \delta), y^*)$$



Subtle Poster  
Camouflage Sticker

Mimic vandalism

“Hide in the human  
psyche”





Subtle Poster

## Lab Test Summary (Stationary)

Target Class: Speed Limit 45



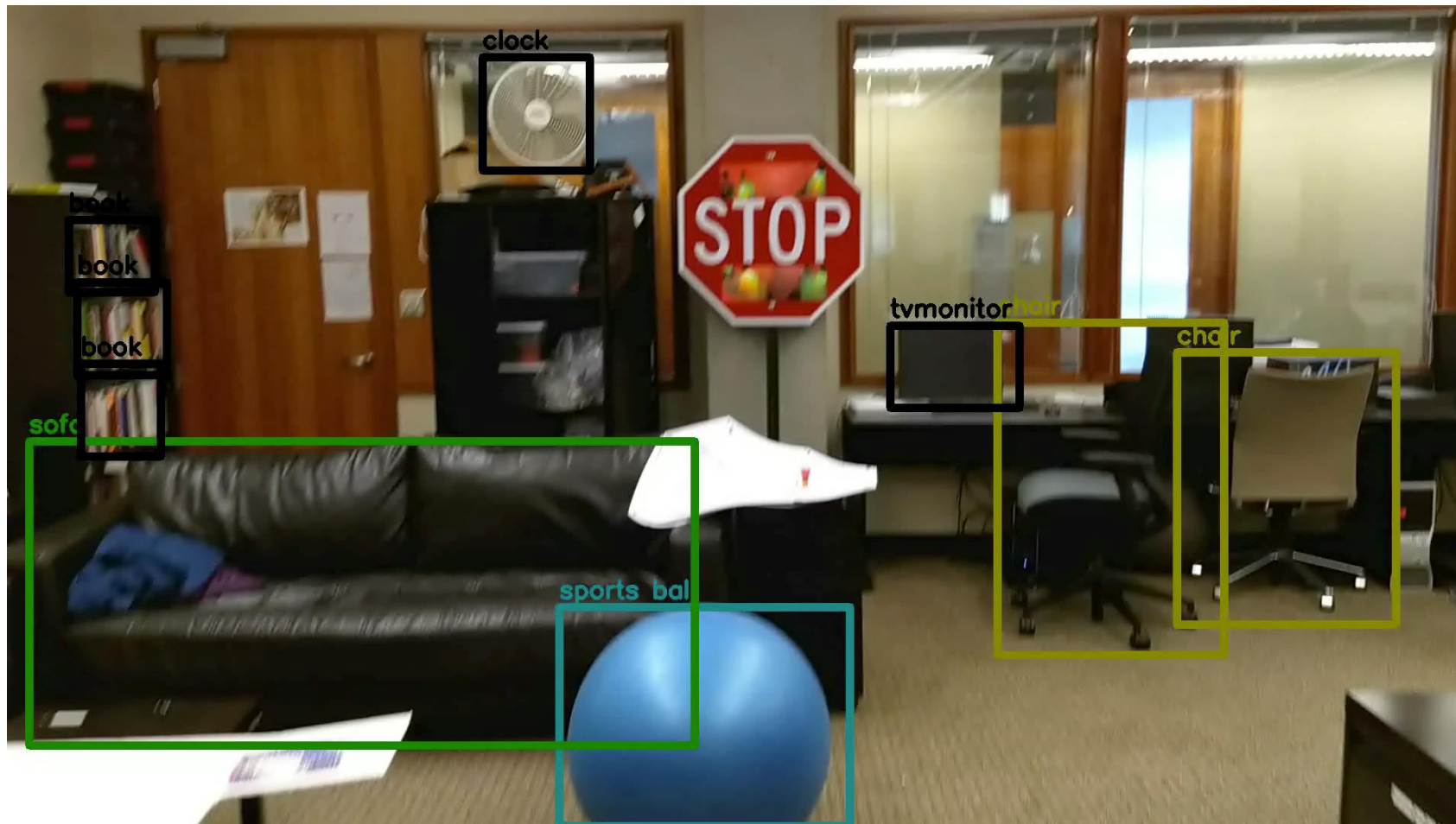
# Art Perturbation



# Subtle Perturbation



# Physical Attacks Against Detectors



# Physical Attacks Against Detectors



# Review format

- Summary
  - Goal
  - Contributions
  - Specific technique details/analysis
- Advantages
- Disadvantages
- Potential improvement and other thoughts

# Potential Final Project Topics

- Attacks against general machine learning models such as 3D reconstruction, BERT, and RL systems.
- Detection against attacks such as Deepfake.
- GWAS for AI
- Theoretically understanding of generative models from the game theoretic perspective
- Applications of GANs (GAN Zoo)
- Provable robustness for classifiers against different types of perturbation
- Differential private graphs, and robust graph neural networks
- Privacy analysis for generative models
- Certiably robust reinforcement learning
- Improve model robustness with unlabeled data via semi-supervised learning
- Robustness testing for different deep neural networks architectures
- Robust autoML
- Semantic Forensics
- Design an ensemble model which guarantees the diversity of the individual classifiers and therefore improve robustness