

# Understanding Deep Neural Networks using Adversarial Attacks

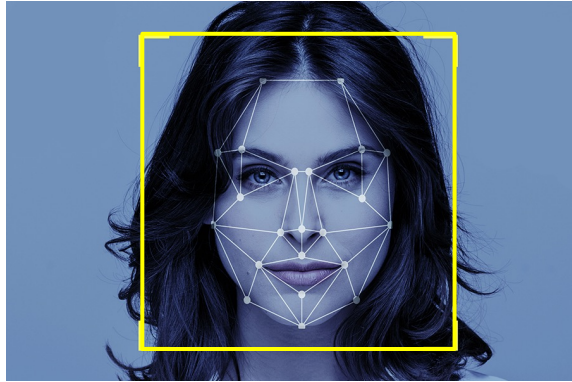
**Krishna Kanth Nakka**

August 15, 2022

**CVLab, EPFL**

# Motivation

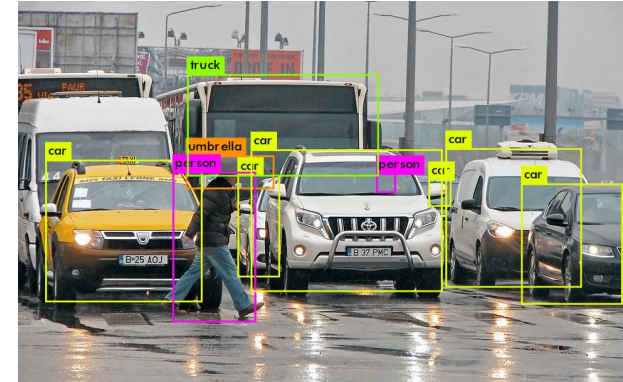
Understanding the behavior of DNNs in safety and security-critical applications is paramount



Biometric recognition



Scene segmentation



Object localization



Health-care applications



Self-driving cars

# Adversarial Examples (AE)

DNNs are sensitive to imperceptible perturbations



$x$

“panda”

57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$

=



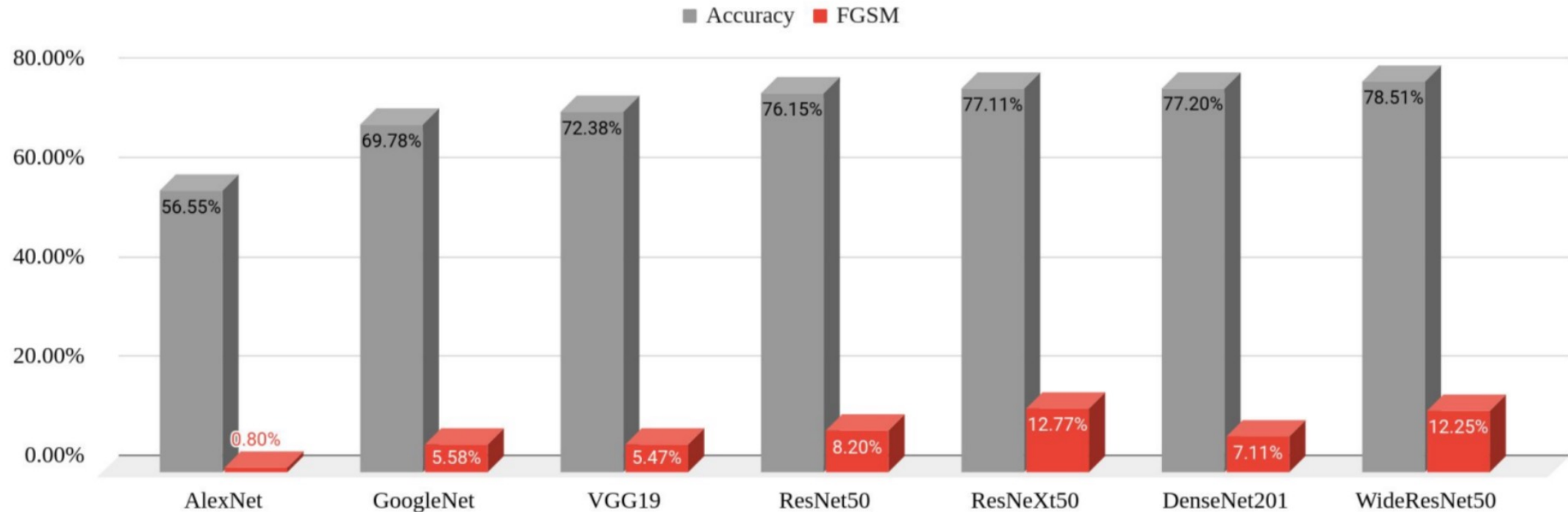
$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$   
“gibbon”

99.3 % confidence

Key properties of adversarial examples

- Small perturbation
- High confidence
- Transferability

# DNNs performance drops significantly with single step FGSM attack



Perturbation norm set to 8.  
Results are reported on ImageNet validation set

# Implications of adversarial attacks against autonomous vehicles



STOP sign

+



Perturbation

=



Max speed 100

# A unifying perspective of thesis

1. Understand the underlying working mechanisms of adversarial attacks on DNNs

2. Design adversarial attacks to both fool and explain the DNNs

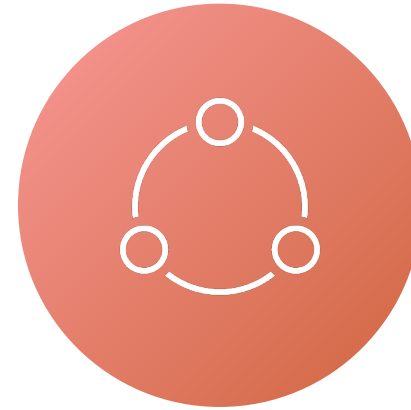
# Focus areas



Interpretable  
models



Adversarial  
Defense



Adversarial  
Attacks

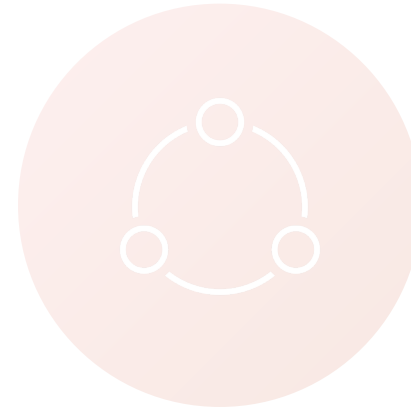
# Focus areas



Interpretable  
models



Adversarial  
Defense



Adversarial  
Attacks

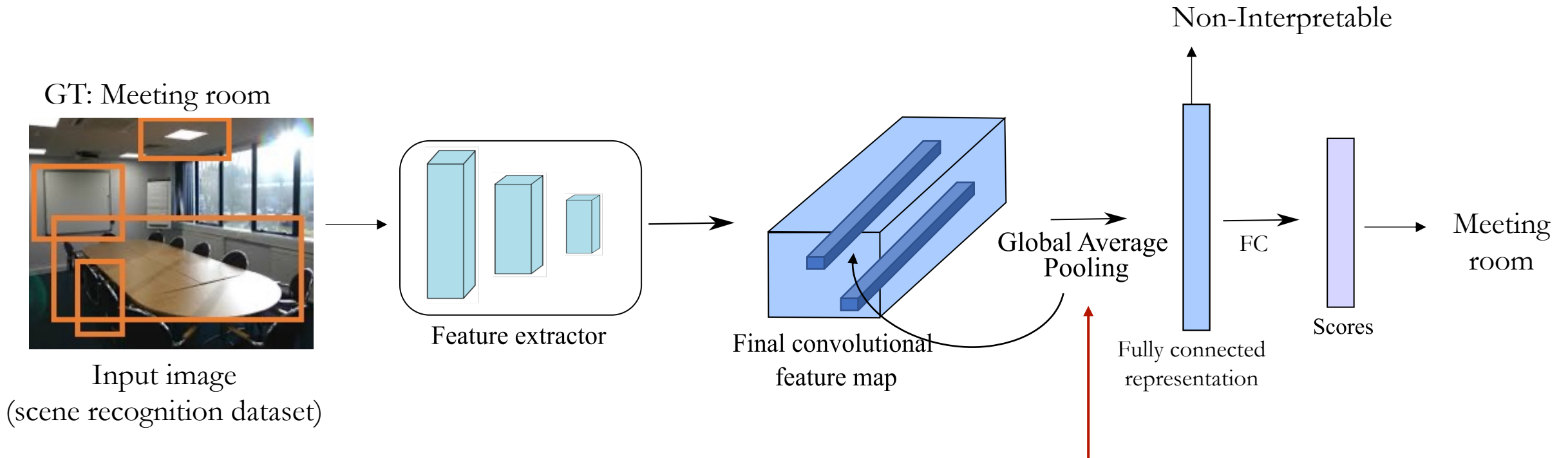


In order to build trust in safety-critical systems, we need to build transparent models that **have the ability to explain why they predict what they predict**



Our work focuses on **Bag-of-visual words (BoW) pooling** architectures to understand the decisions of DNNs

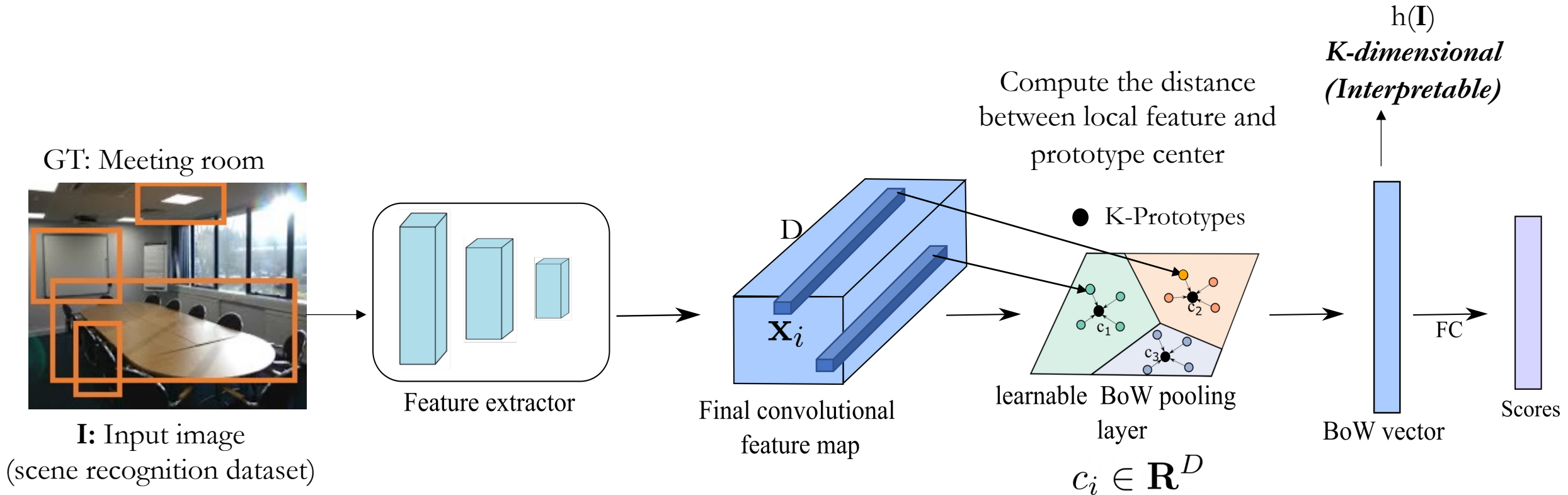
# Standard DNN architectures are difficult to understand how they reach to their decisions



presence of table, board, chairs, light bulb indicates **meeting room** class

Instead of GAP, replace with **learnable BoW pooling** layer to learn interpretable representation

# NetBoW: DNN with distance-based learnable pooling layer



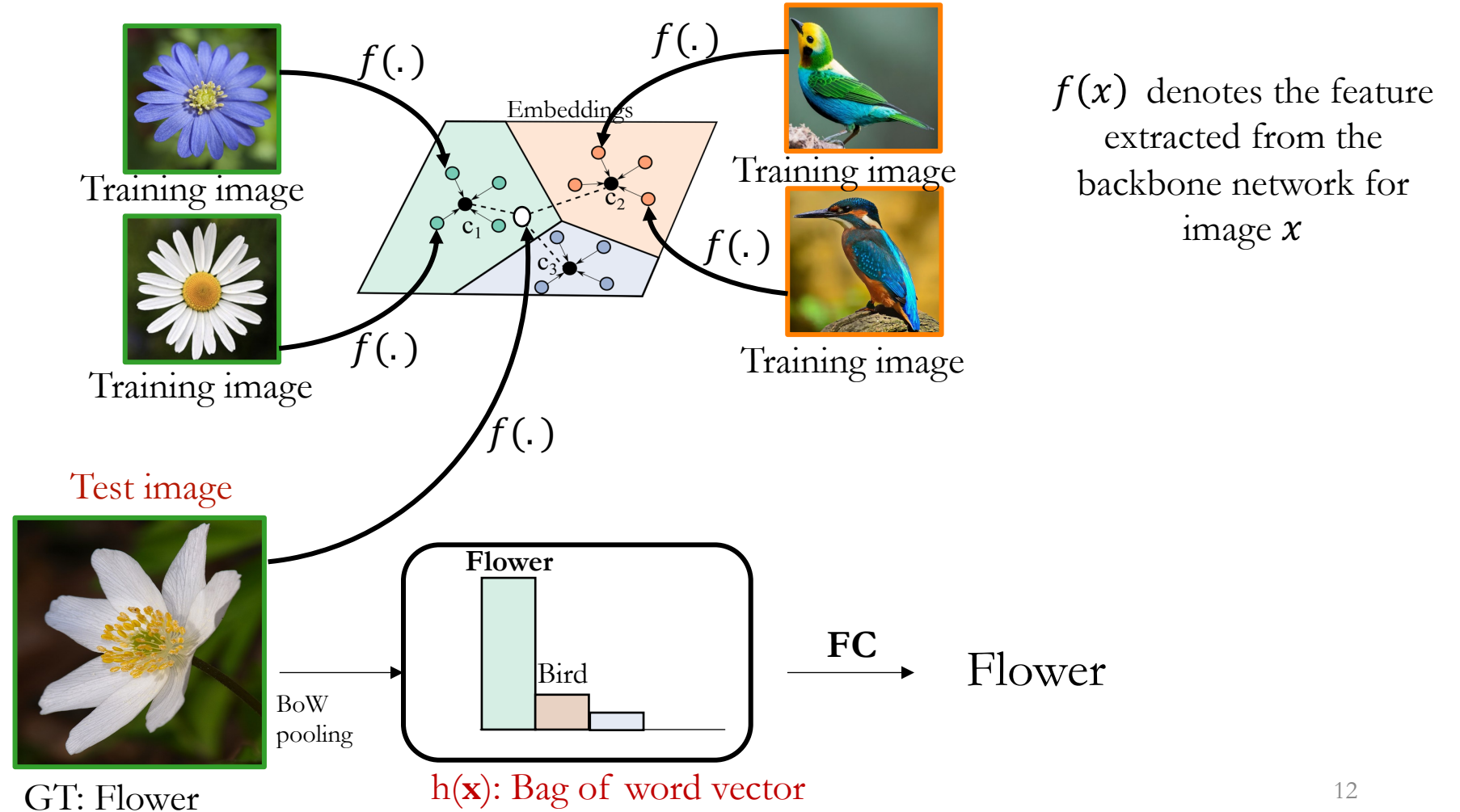
presence of table, board, chairs, light bulb indicates **meeting room** class

$$h_k(\mathbf{I}) = \sum_{i=1}^N a_k(\mathbf{x}_i)$$

$$a_k(\mathbf{x}_i) = \frac{e^{-\alpha \|\mathbf{x}_i - \mathbf{c}_k\|^2}}{\sum_{k'} e^{-\alpha \|\mathbf{x}_i - \mathbf{c}_{k'}\|^2}}$$

Network learns the 1. weights of backbone, 2. prototype centers ( $\mathbf{c}_1, \mathbf{c}_2, \mathbf{c}_3$ ), and 3. FC layer jointly

# Advantage: Interpretable BoW representation



# Interpretability by Design

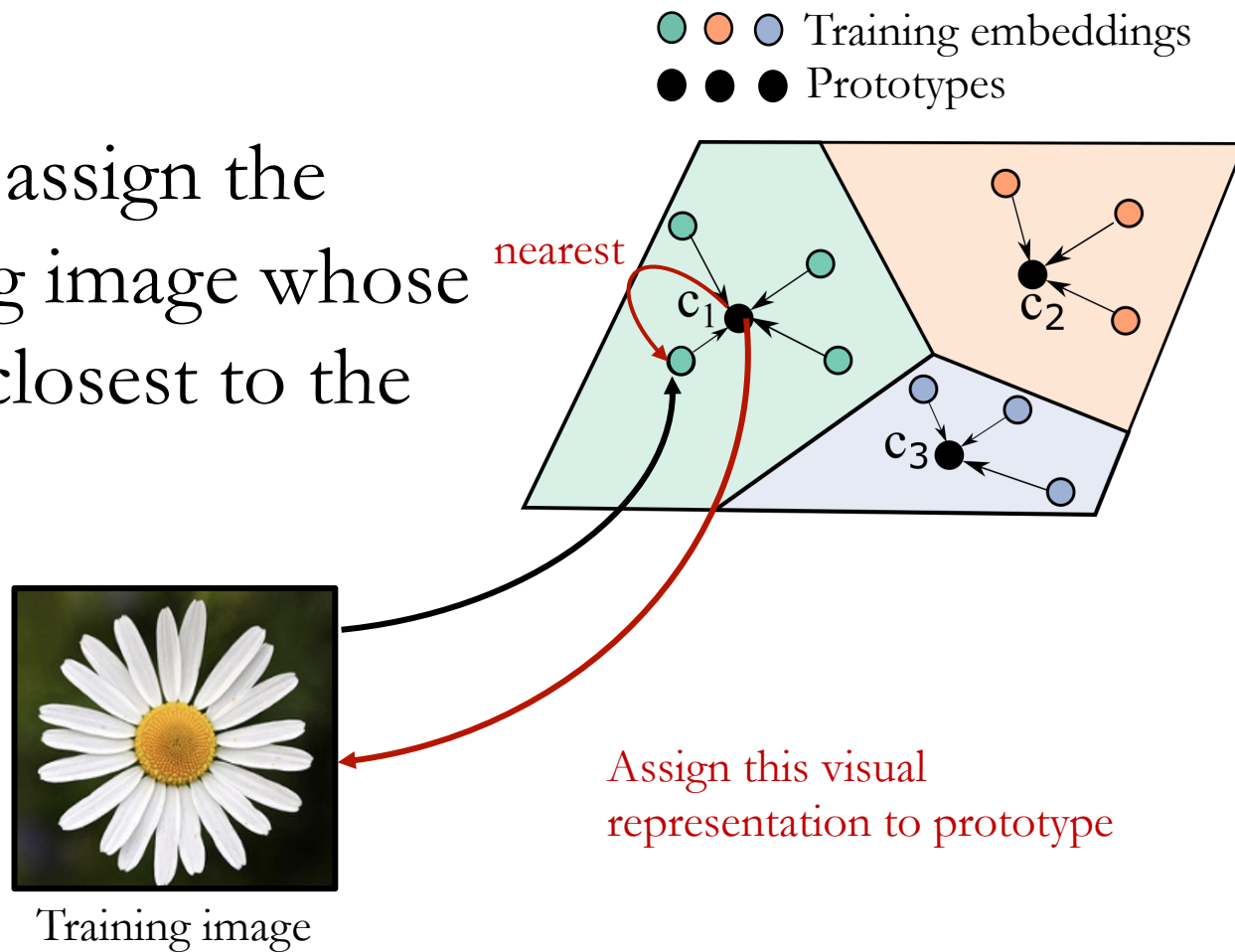
BoW Networks are interpretable since one can understand the reasons for particular output decision through the prototype activations and not in post-hoc manner



Diagnose the failure modes such as **adversarial examples, out-of-distribution examples** and analyze them more naturally and intrinsically

# What does prototype represent in input image space?

After training, assign the nearest training image whose embedding is closest to the prototype



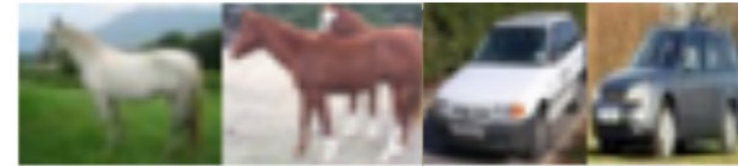
# NetBoW helps to understand the reasons for a label prediction through visual codebook



MNIST



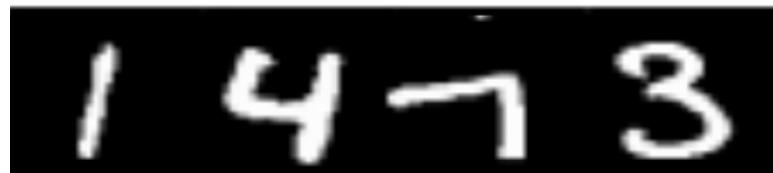
FMNIST



CIFAR10



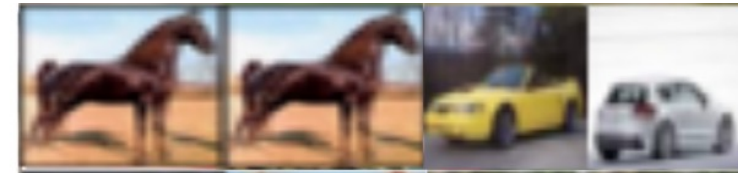
Highest activated prototype



MNIST

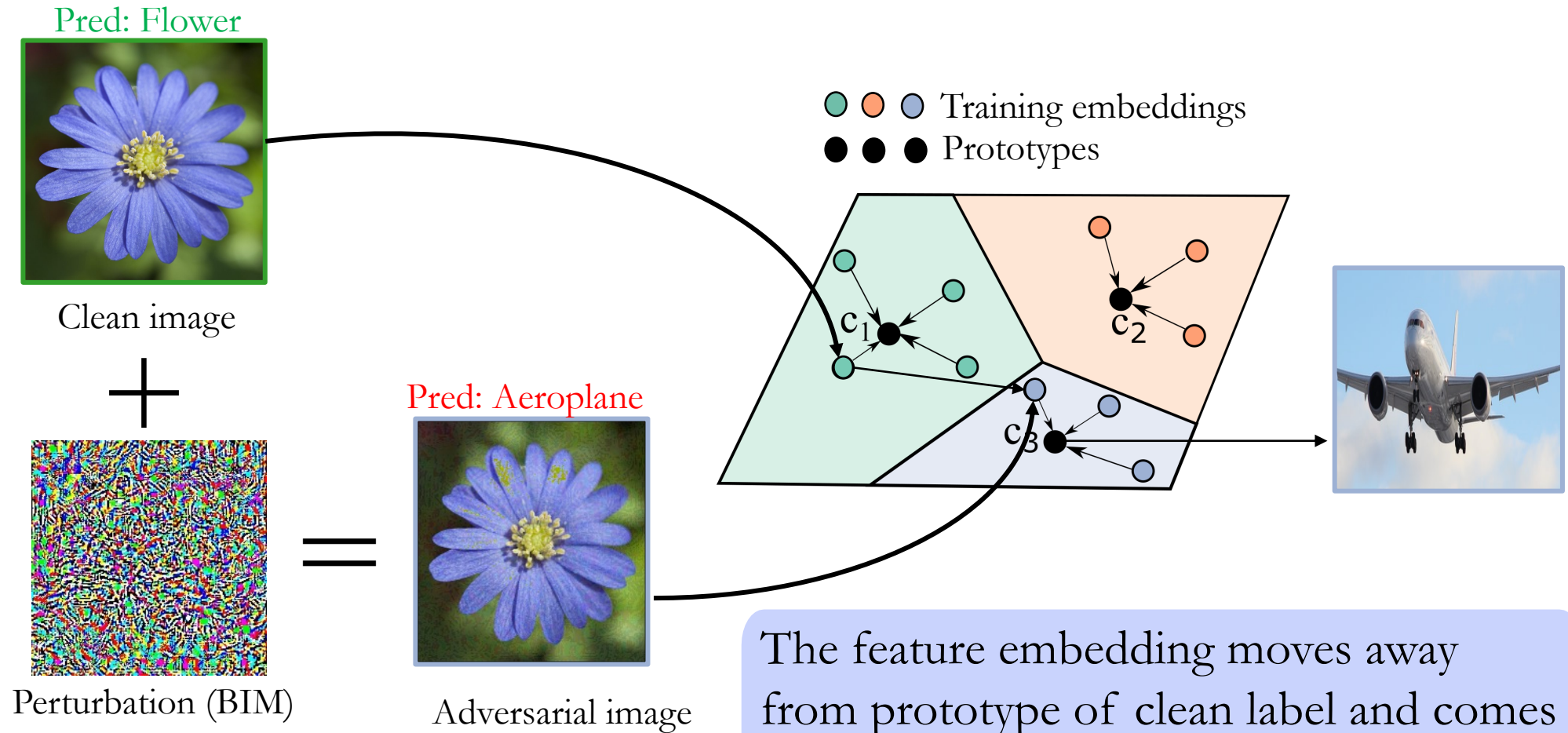


FMNIST



CIFAR10

# Can we understand the mechanism of adversarial examples through interpretable models in a better way?





# How can we use the interpretable BoW networks to detect the adversarial examples



## Key idea

Adversarial attacked image should activate the prototype of other class. Therefore, we detect the attacks by comparing the input image with the visual representation of activated prototype through an auxiliary **detector** network

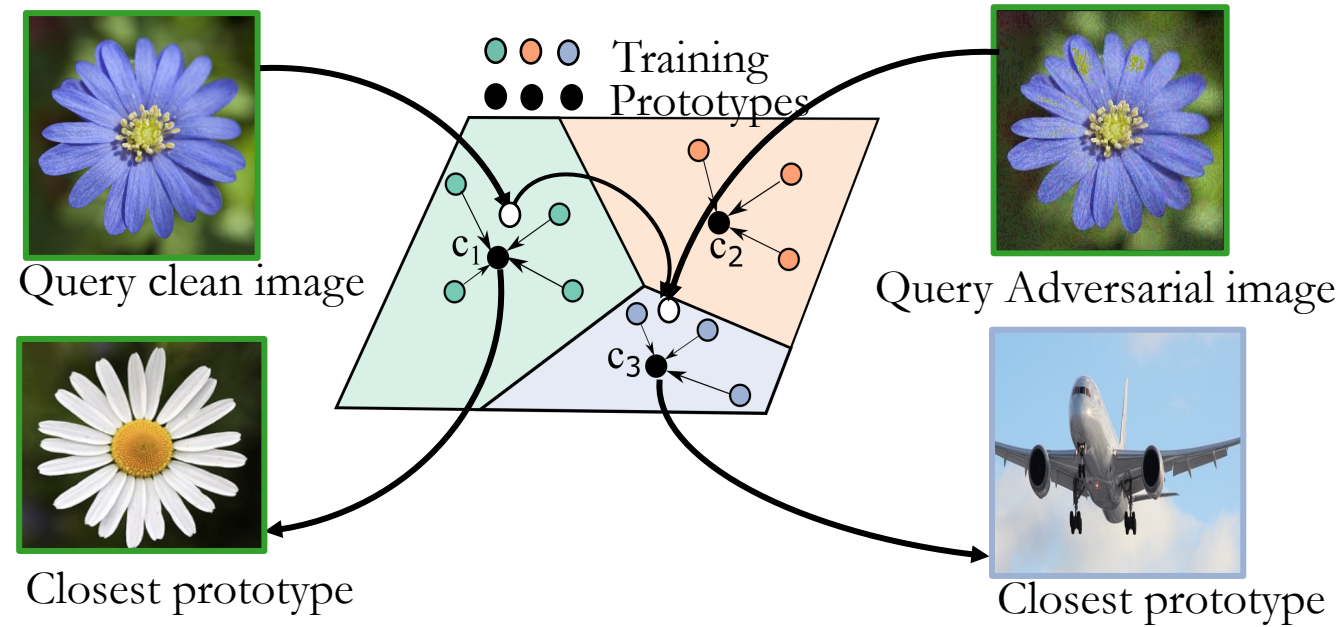
Impact:

Adversarial images →

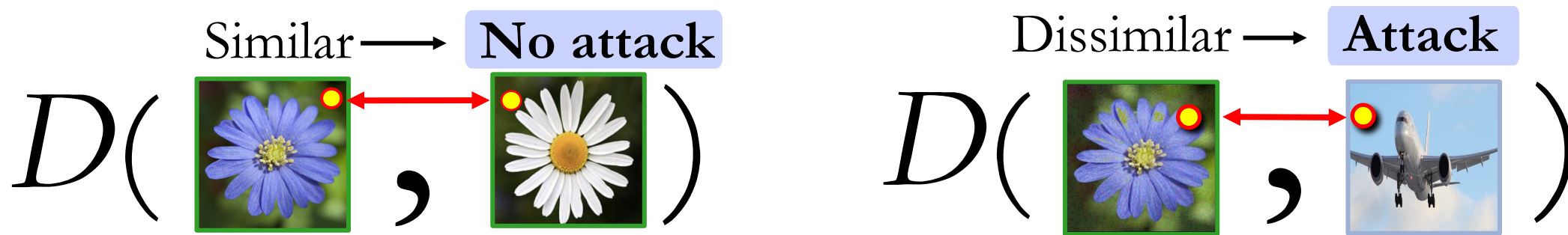


← Activated prototypes

# Adversarial example detection: pose the problem as similarity matching of **input image** to highest **activated prototype**



Similar approach can be used to detect out-of-distribution examples



$D$ : Siamese-based detector to predict if the input pair is similar or dissimilar

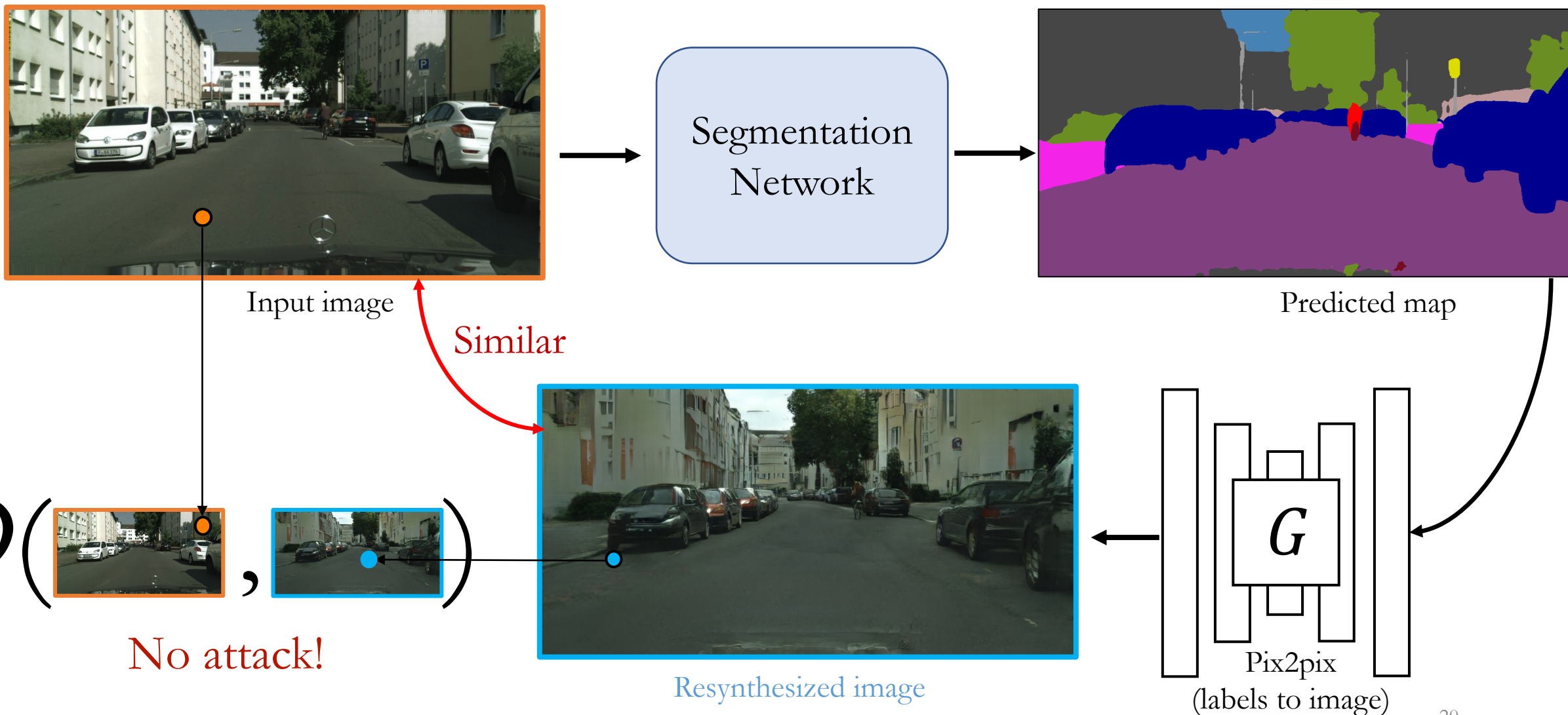
# Is the detector robust to attacks all the time?

- **No.** Detector breaks down with our **defense-aware adaptive attack** in pure white box setting (aware of detector weights and strategy)
- However, the approach works in **gray-box** (adversary aware of defense mechanism) and **black-box** setting (no access to detector weights) wrt detector



# Adversary detection beyond image-recognition

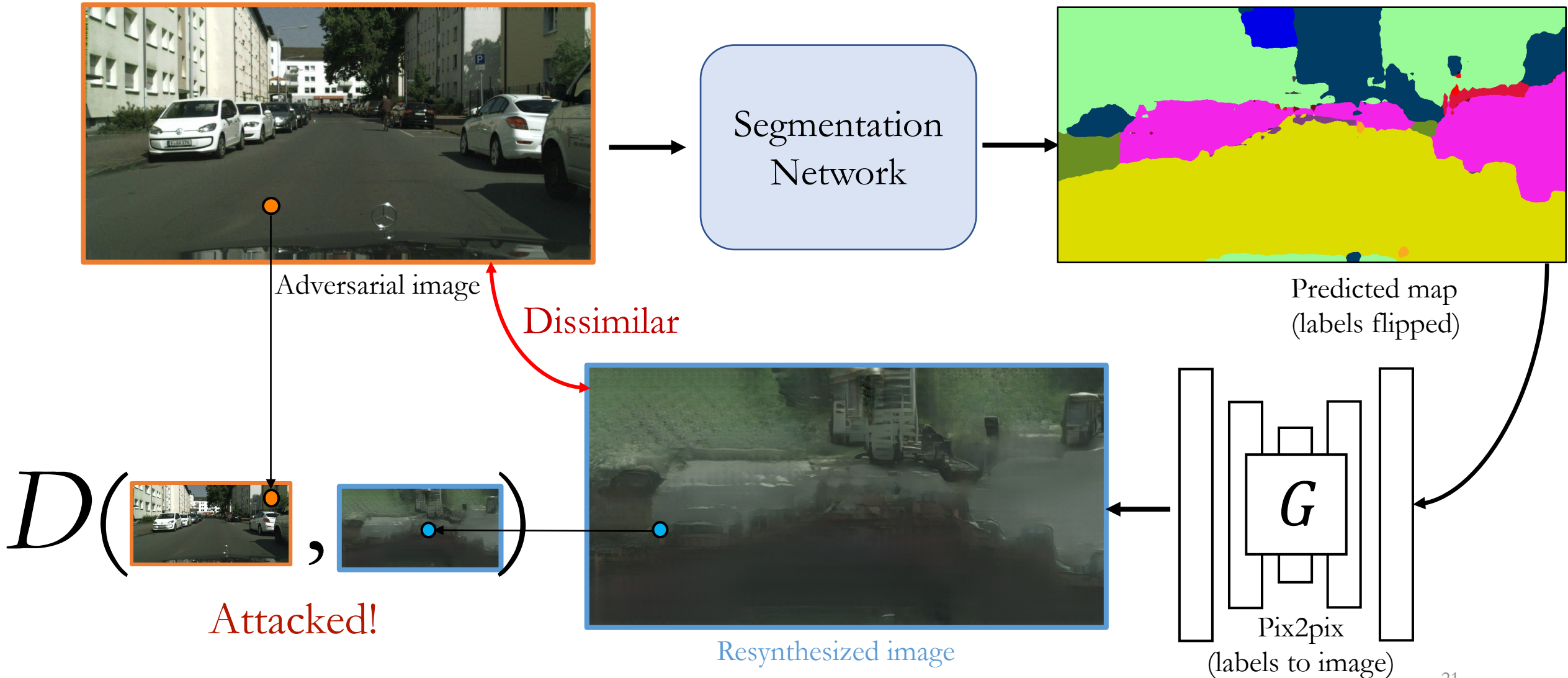
Adversarial example detection in semantic segmentation by comparing **input image** to the **image resynthesized** from output map



$D$ : Computes  $L_1$  distance in HOG feature space

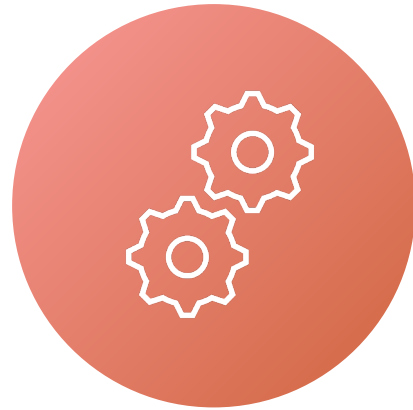
# Adversary detection beyond image-recognition

Adversarial example detection in semantic segmentation by comparing **input image** to the **image resynthesized** from output map

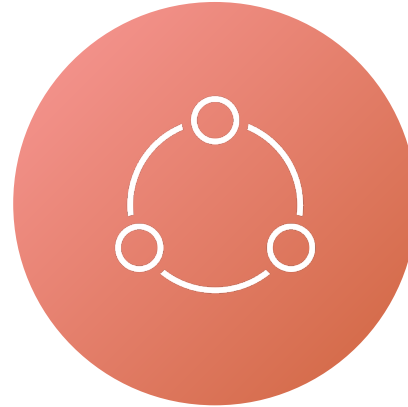


$D$ : Computes  $L_1$  distance in HOG feature space

# Focus areas



Interpretability

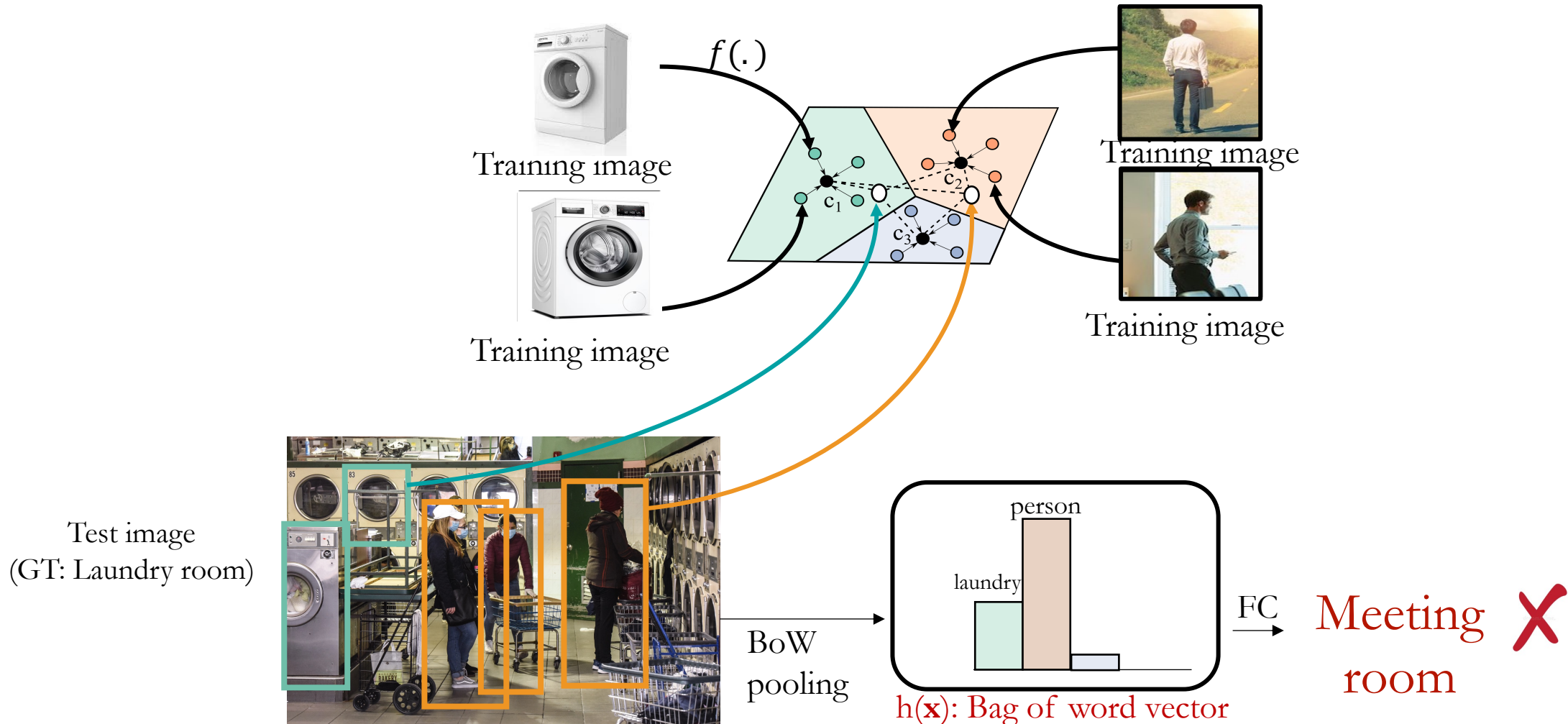


Adversarial  
Defense



Adversarial  
Attacks

# Downside: all features participate in the feature aggregation step of BoW pooling



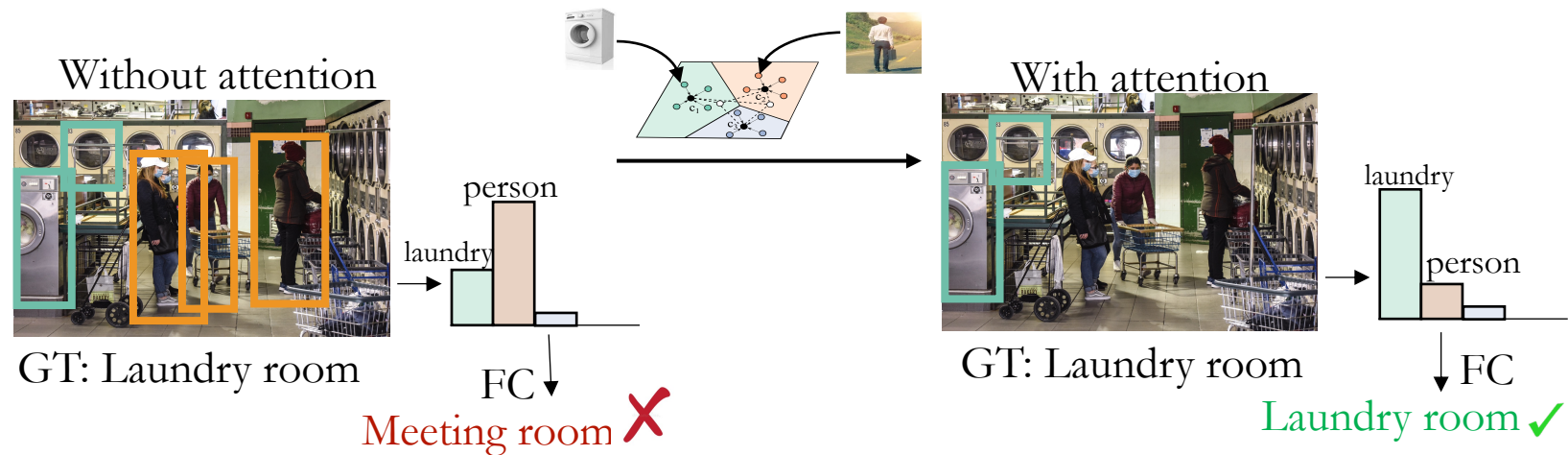
# Remove the influence of non-discriminative regions. **How?**



## Key idea: attention-aware pooling

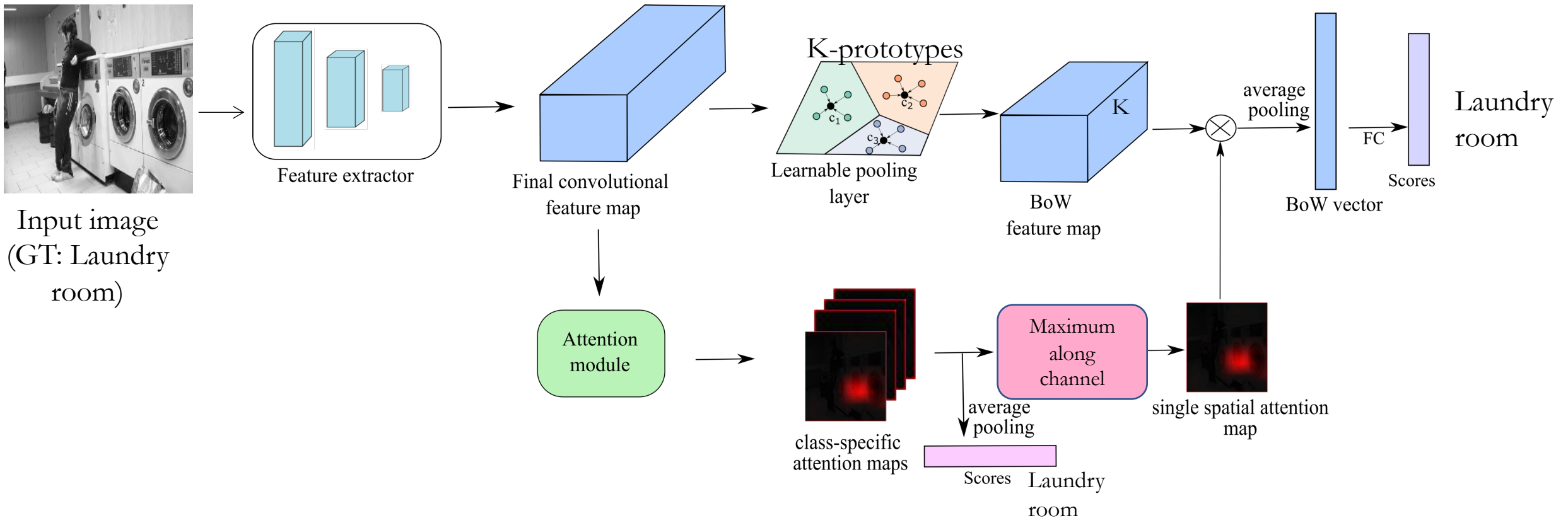
By introducing the attention during the feature aggregation process, the BoW representation becomes more discriminative

Impact:





# Key idea: introduce **attention** in BoW pooling to remove contribution of non-discriminative features



Attention ignores the non-discriminative regions (such as the **person** which is common across classes) and focuses on discriminative regions of the output class

*casino*

*studio music*

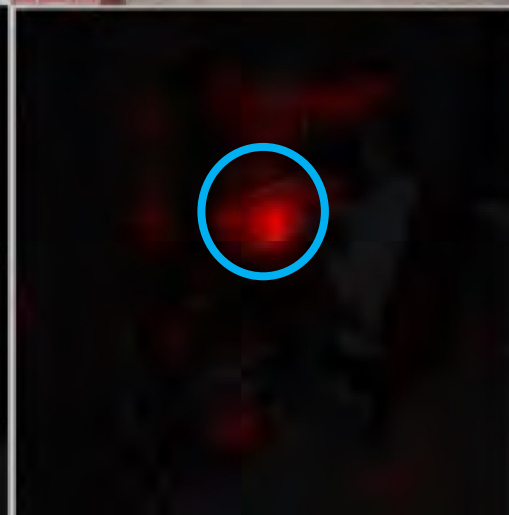
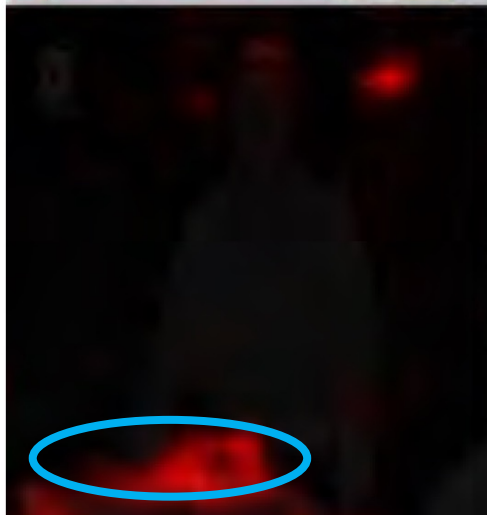
*operating room*

*video store*

Input  
image



Attention  
map



# By incorporating attention, we can significantly improve the discriminative power of BoW latent representation

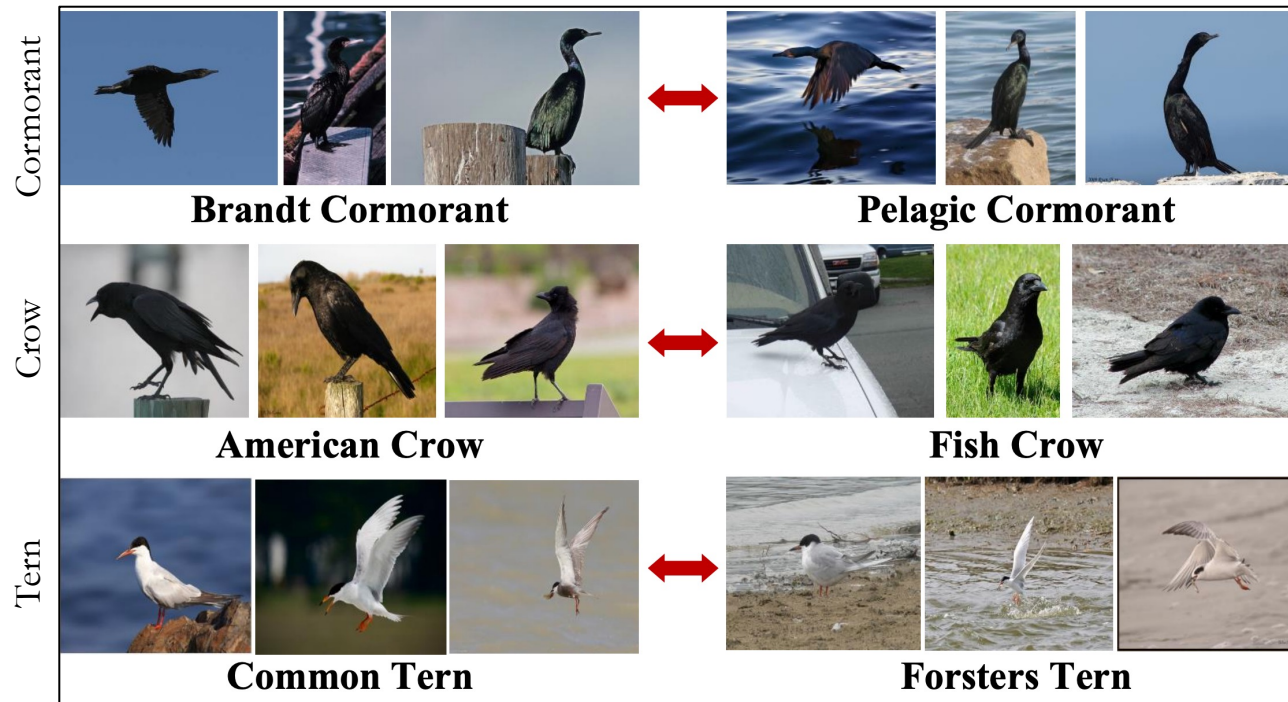
			Anno.	Birds	Cars	Aircrafts	Accuracy↑ MIT-Indoor
Pooling							
VGG-16		BBox		79.9	88.4	86.9	-
Attention		BBox		77.2	90.3	85.0	-
4%	NetBoW	BBox		74.4	89.1	85.6	-
	Attentional-NetBoW	BBox					
	NetVLAD	BBox		82.4	89.8	88.0	-
	Attentional-NetVLAD	BBox					
VGG-16				76.0	82.8	82.3	76.6
Attention				77.0	87.4	81.4	77.2
8%	NetBoW			68.9	85.2	79.9	76.1
	Attentional-NetBoW						
	NetVLAD			80.6	89.4	86.4	79.2
Attentional-NetVLAD							

Attention helps  
the most in  
BoW pooling

# The reasons for success of adversarial attacks in fine-grained datasets has lot of subtleties

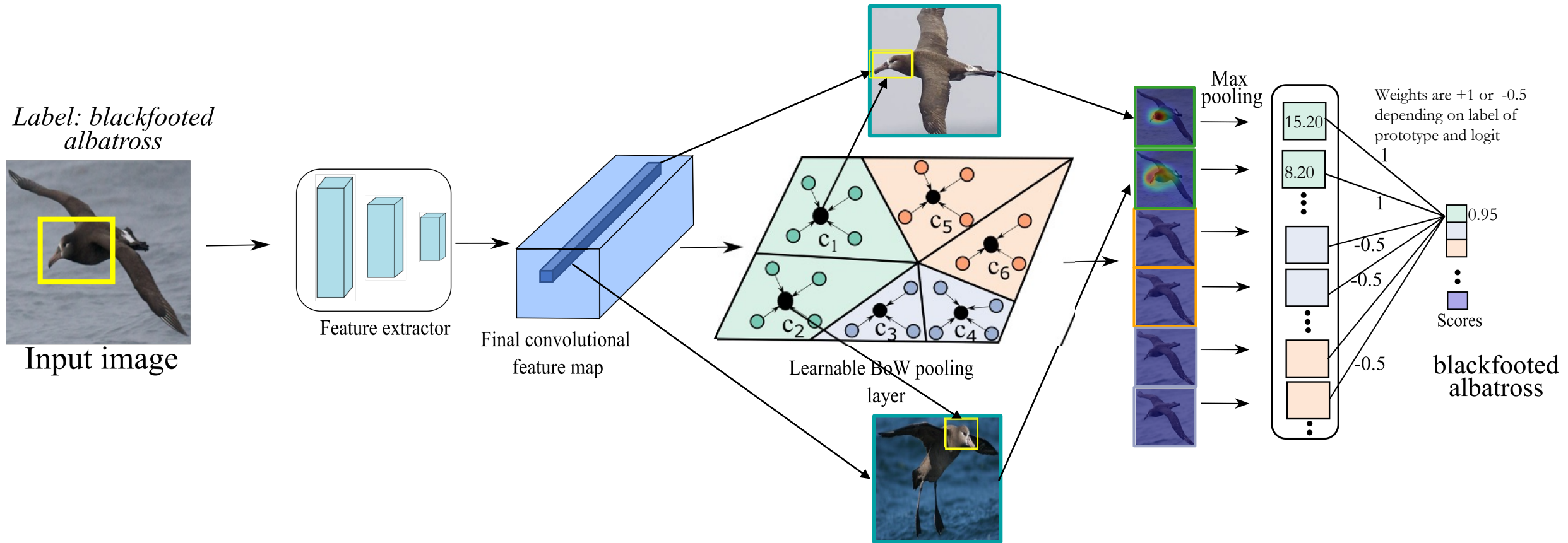
## Benefits with fine-grained datasets

- Understand the DNNs workings at local patch level instead of global object level
- More sensitive to attacks since local perturbations can change the label



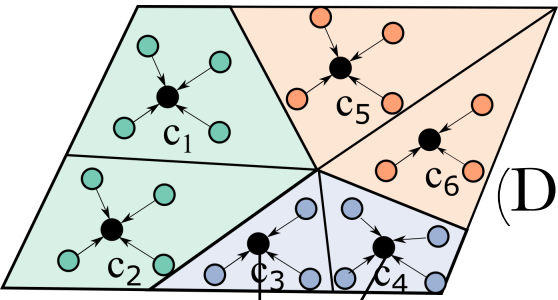
Example images of some confusing classes

# ProtoNet: Classify image based on evidence from local patches

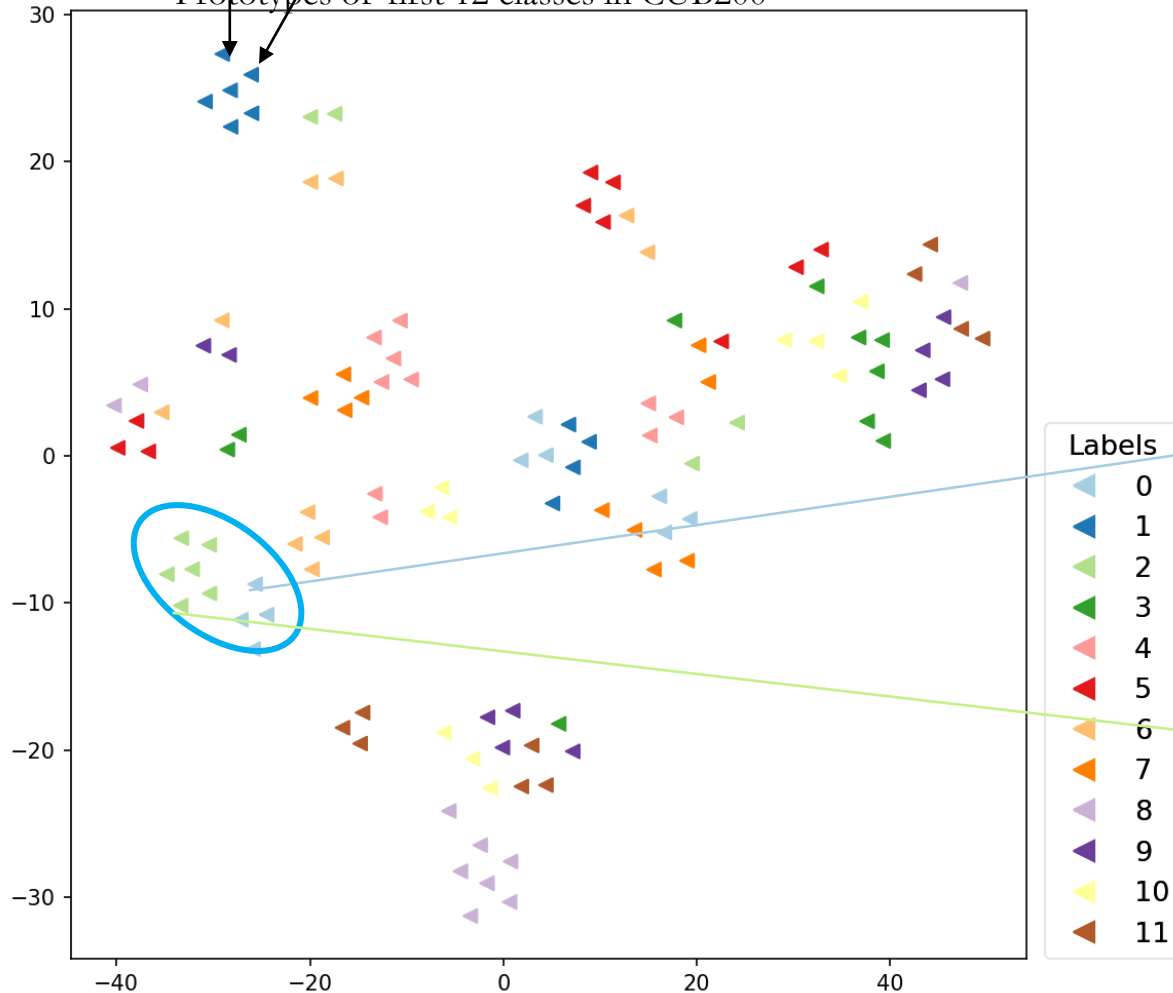


# t-SNE visualization of learned prototypes

(Digging deeper: understanding reasons for success of attacks)

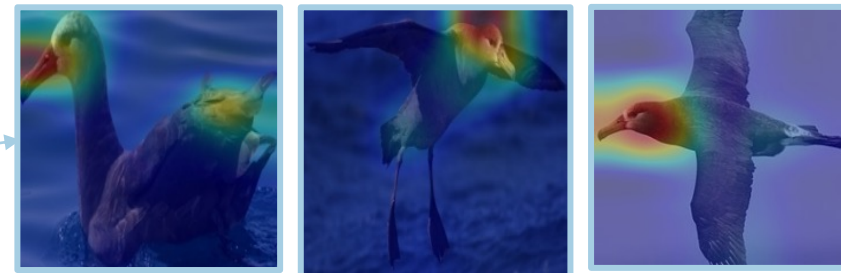


Prototypes of first 12 classes in CUB200



Foreground prototypes of **different classes** of **same family** are close to each other

Class 0 : Blackfooted albatross



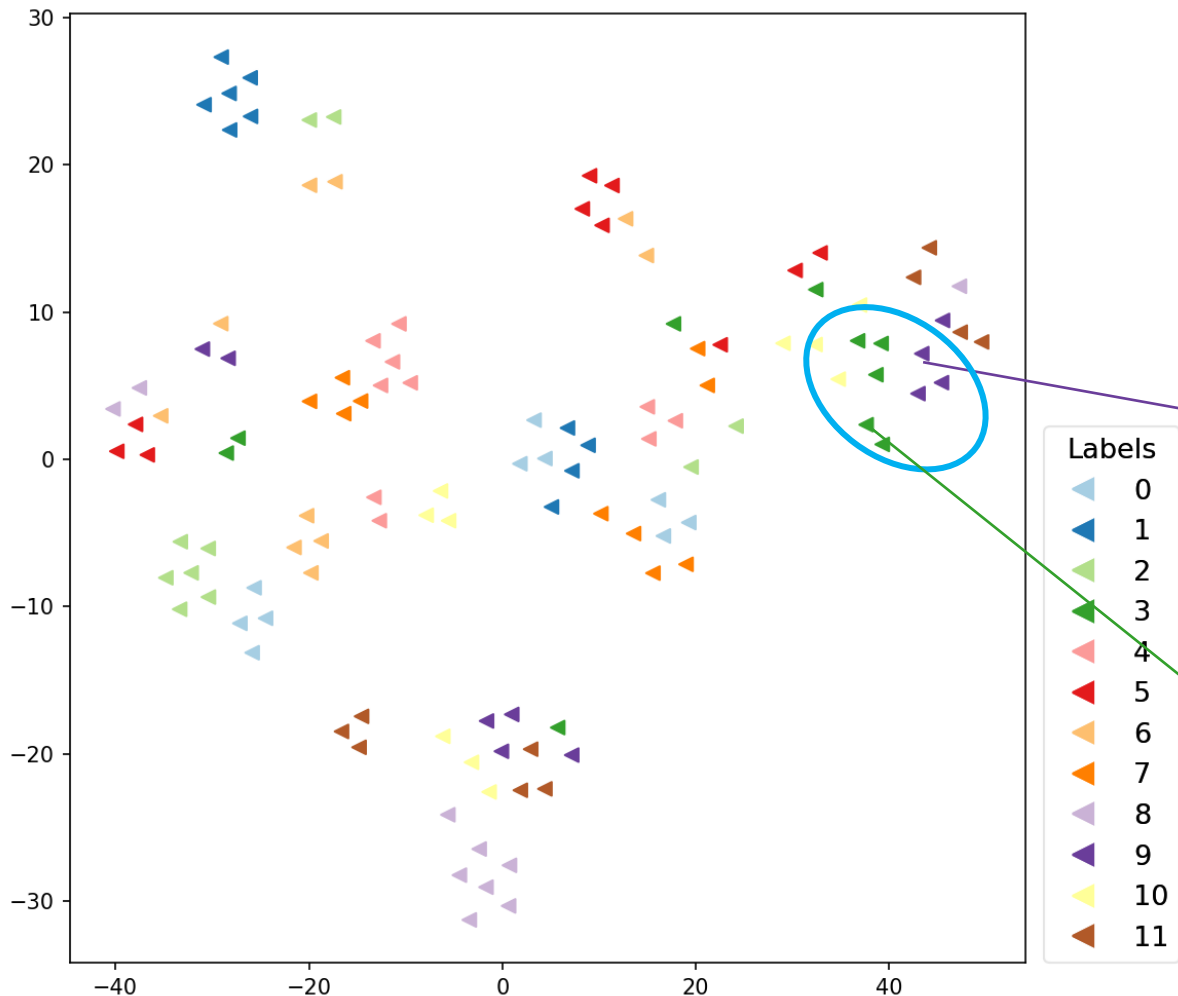
Class 2: Sooty albatross



# t-SNE visualization of learned prototypes

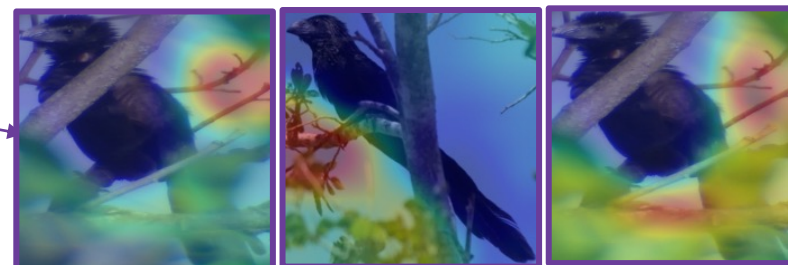
(Digging deeper: understanding reasons for success of attacks)

Prototypes of first 12 classes in CUB200

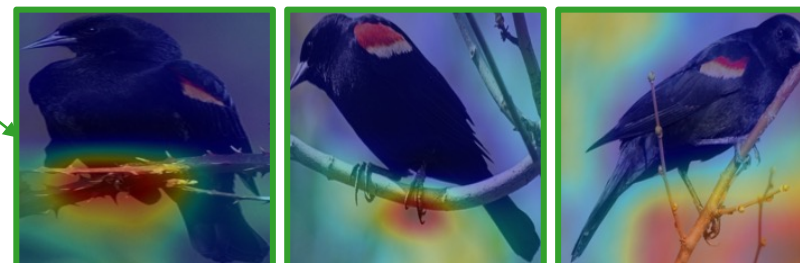


Background prototypes of different classes are close to each other and therefore can be easily attacked to change from one class to other

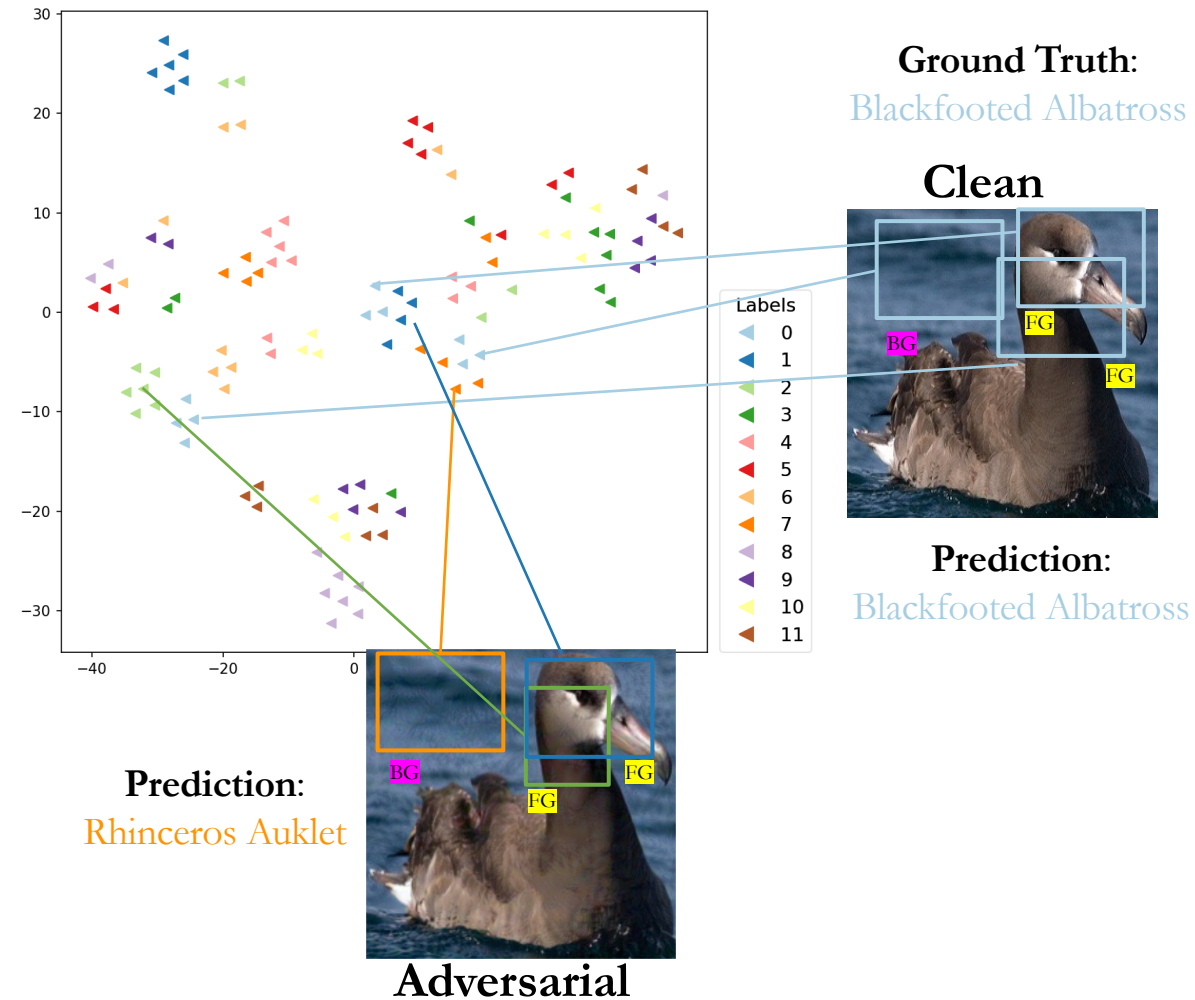
Class 3



Class 7



# An intuitive example to understand the success of adversarial attacks on ProtoPNet





# How can these observations help to improve the robustness?

## maximal separation of discriminative features



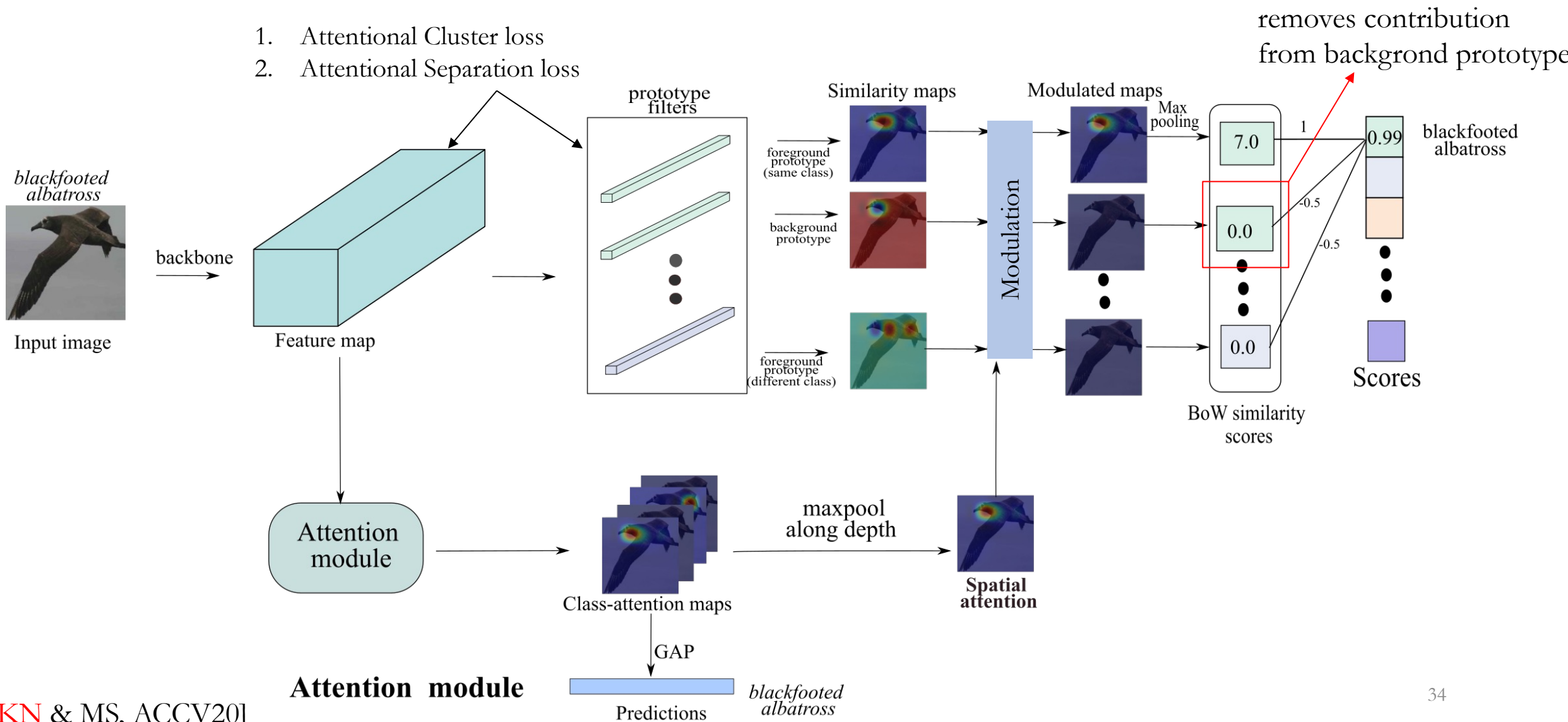
### Intuition

If we can maximally **separate** the latent features of **foreground** (discriminative) regions of different classes and also remove the influence of **background** regions in the decision process, then we have made the attacker task difficult to conduct attacks

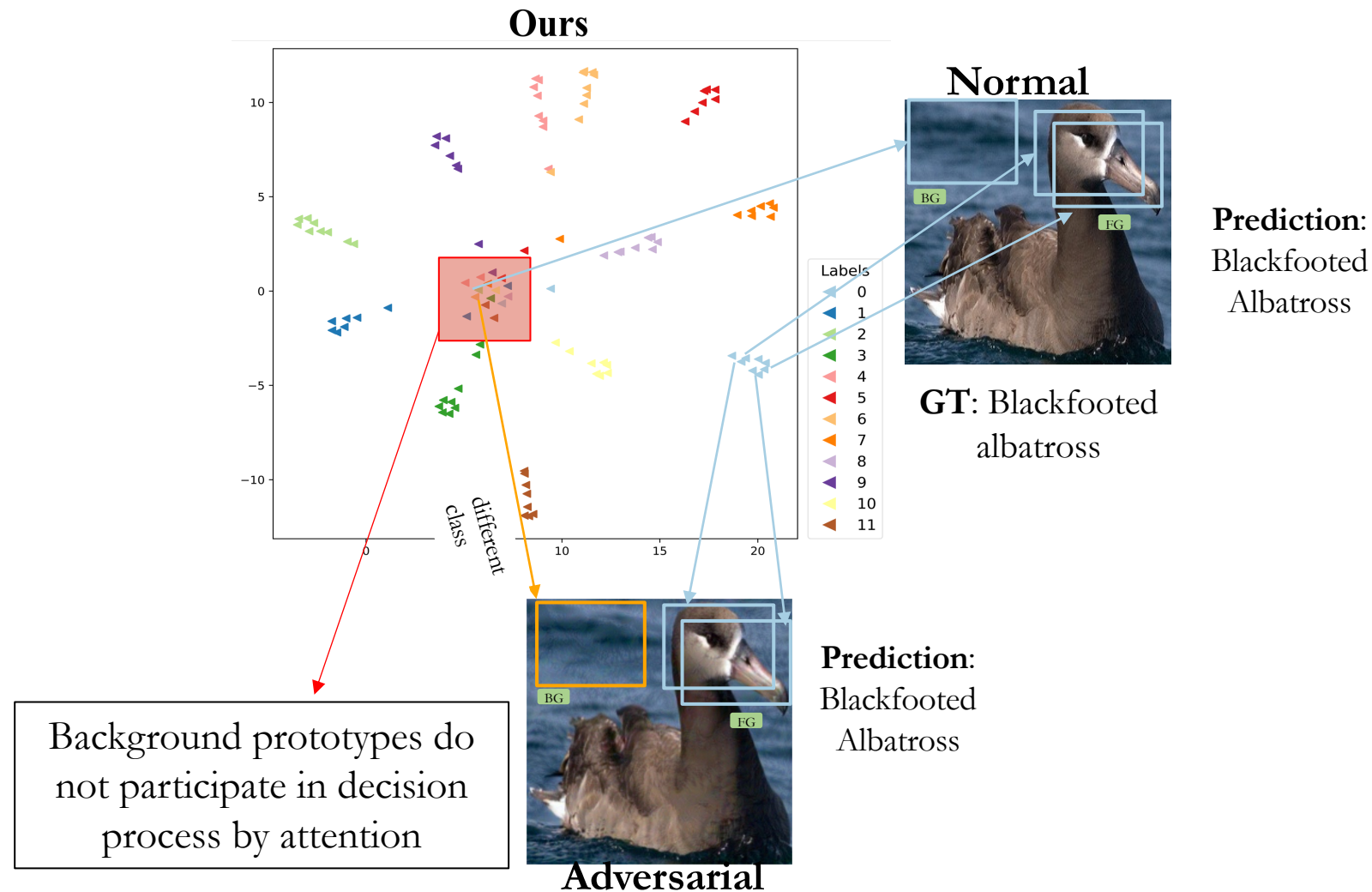
1. **Attentional cluster loss** - pulls the high-attention regions in a sample close to the nearest prototype of its own class
2. **Attentional separation loss** - pulls the high-attention regions in a sample away from the nearest prototype of other class

# Attention-aware architecture

1. Attentional Cluster loss
2. Attentional Separation loss



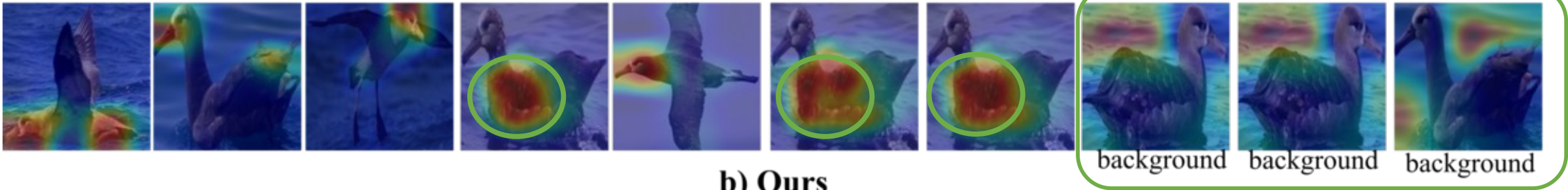
# Our method yields well-separated foreground prototypes while clustering background prototypes



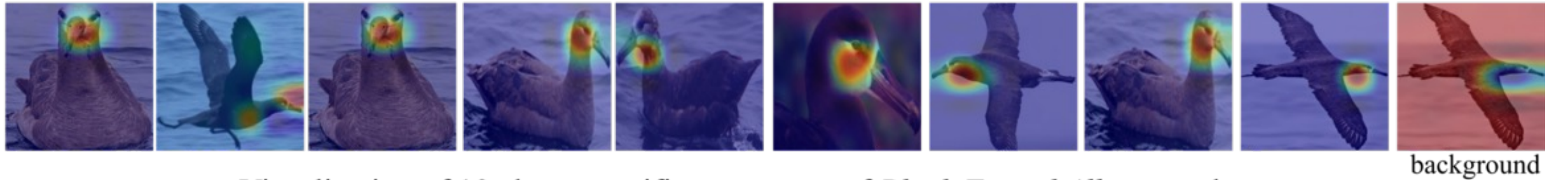
# Visualization of learned Prototypes

(our prototypes are **fine-grained** and complete non-discriminative regions activated by **background** prototype)

(a) ProtoPNet



b) Ours



Visualization of 10 class specific prototypes of *Black Footed Albatross* class

Our adversarial training strategy with novel losses consistently achieves higher performance against white-box attacks

		Accuracy↑								
Base Network	Attacks (Steps, $\epsilon$ )	Clean (0,0)	FGSM (1,2)	FGSM (1,8)	BIM (10,2)	BIM (10,8)	PGD (10,2)	PGD (10,8)	MIM (10,2)	MIM (10,8)
VGG-16	AP* [71]	54.9%	44.9%	24.2%	41.9%	18.2%	41.2%	16.9%	41.9%	18.7%
	AP+PCL* [183]	60.7%	50.5%	28.5%	47.1%	22.8%	46.7%	21.6%	47.2%	23.5%
	<b>Ours-A*</b>	<b>69.3%</b>	<b>56.1%</b>	<b>34.8%</b>	<b>51.7%</b>	<b>29.6%</b>	<b>50.8%</b>	<b>28.0%</b>	<b>52.0%</b>	<b>32.5%</b>
	ProtoPNet* [29]	60.1%	44.5%	26.9%	<b>57.1%</b>	10.9%	35.9%	10.3%	37.6%	13.5%
	<b>Ours-FR*</b>	<b>63.0%</b>	<b>53.3%</b>	<b>37.3%</b>	49.4%	<b>30.4%</b>	<b>48.1%</b>	<b>28.6%</b>	<b>49.7%</b>	<b>31.1%</b>
VGG-19	AP* [71]	58.0%	47.5%	29.1%	44.3%	25.6%	44.0%	24.34%	44.4%	26.2%
	AP+PCL* [183]	61.8%	52.1%	30.9%	48.9%	24.7%	48.6%	23.3%	49.1%	25.4%
	<b>Ours-A*</b>	<b>68.2%</b>	<b>57.1%</b>	<b>36.5%</b>	<b>53.2%</b>	<b>30.4%</b>	<b>52.6%</b>	<b>29.2%</b>	<b>53.5%</b>	<b>31.2%</b>
	ProtoPNet* [29]	55.1%	40.0%	28.9%	26.5%	11.3%	29.7%	9.60%	25.6%	10.2%
	<b>Ours-FR*</b>	<b>64.4%</b>	<b>55.5%</b>	<b>37.4%</b>	<b>51.2%</b>	<b>30.6%</b>	<b>50.4%</b>	<b>28.7%</b>	<b>52.1%</b>	<b>32.3%</b>

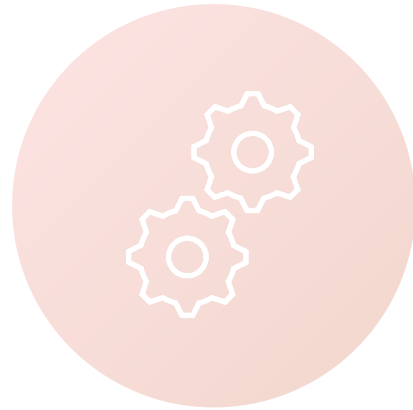
# Black-box auto-attack ensemble on adversarial trained models

Accuracy ↑

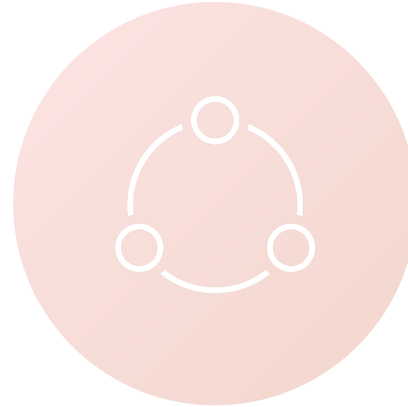
Base Attacks		Clean	APGD (CE)	APGD (DLR)	Square	Auto attack
VGG-16	AP* [74]	54.9%	15.1%	14.0%	39.2%	22.7%
	AP+PCL* [186]	60.7%	18.0%	14.1%	42.9%	25.0%
	<b>Ours-A*</b>	<b>67.0%</b>	<b>23.7%</b>	<b>15.1%</b>	<b>47.3%</b>	<b>28.7%</b>
	ProtoPNet* [29]	55.6%	2.8%	2.3 %	31.6%	12.2%
	<b>Ours-FR*</b>	<b>60.4%</b>	<b>24.2%</b>	<b>15.5%</b>	<b>46.2%</b>	<b>28.6%</b>
VGG-19	AP* [74]	55.7%	20.2%	14.4%	44.1%	26.2%
	AP+PCL* [186]	59.7%	20.8%	17.3%	51.1%	29.7 %
	<b>Ours-A*</b>	<b>65.0%</b>	<b>24.4%</b>	<b>17.4%</b>	<b>51.9%</b>	<b>31.2%</b>
	ProtoPNet* [29]	51.9%	1.1 %	1.0 %	28.0%	10.0%
	<b>Ours-FR*</b>	<b>62.1%</b>	<b>27.4%</b>	<b>18.5%</b>	<b>52.1%</b>	<b>32.7%</b>

$\epsilon=8$

# Focus areas



Interpretability



Adversarial  
Defense



Adversarial  
Attacks

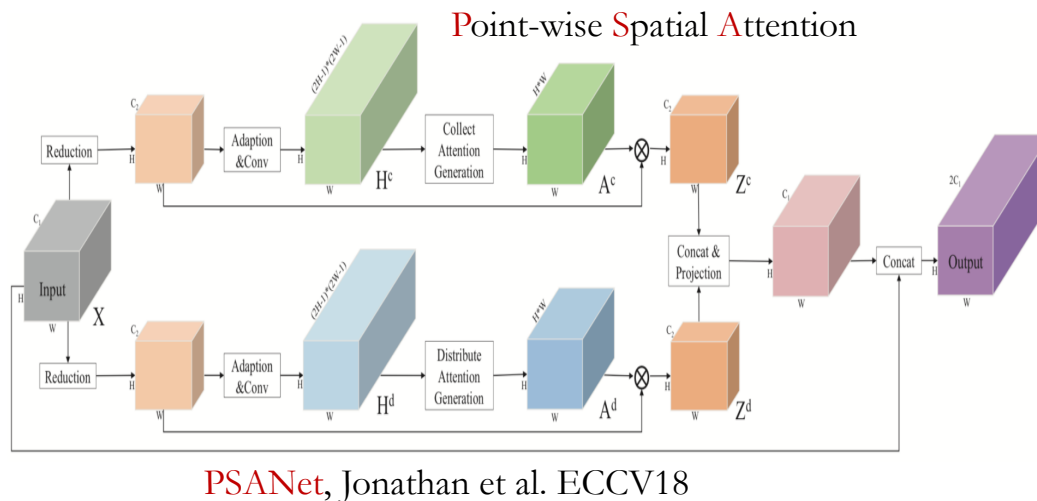
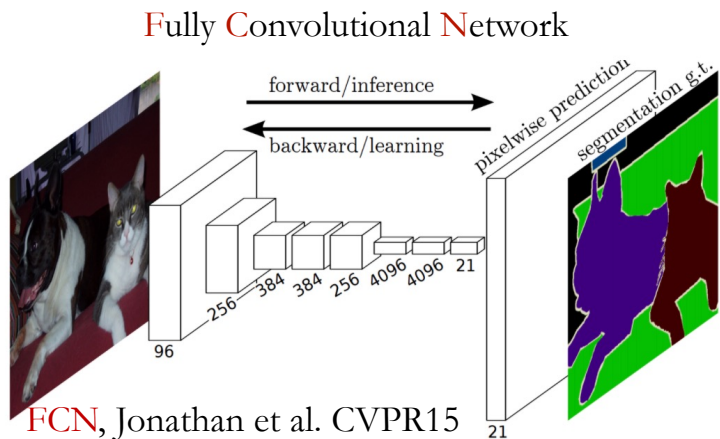
# Adversarial Attacks

- Attacks beyond image recognition
  - White-box attacks on semantic segmentation
  - Black-box transfer attacks on visual object tracker
- Improving the transferability of attacks
  - Learning transferable transferable perturbations
- How can we use adversarial attack to improve DNNs
  - Semantic adversarial attacks to study disentanglement



Exploiting **context** to understand  
the susceptibility of DNNs for  
Semantic Segmentation

# Understanding **context** is the core building block in modern segmentation networks

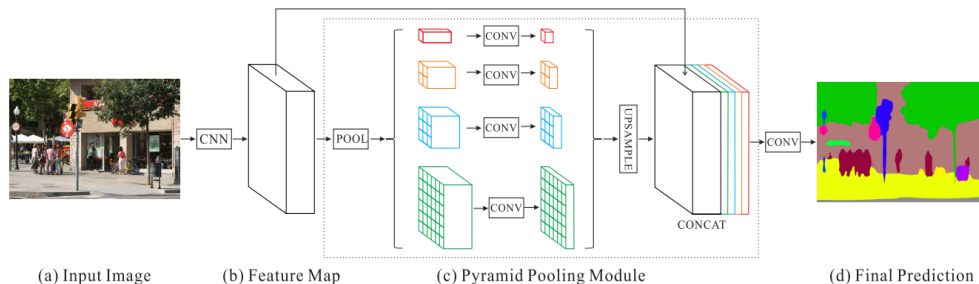


no-context

heavy

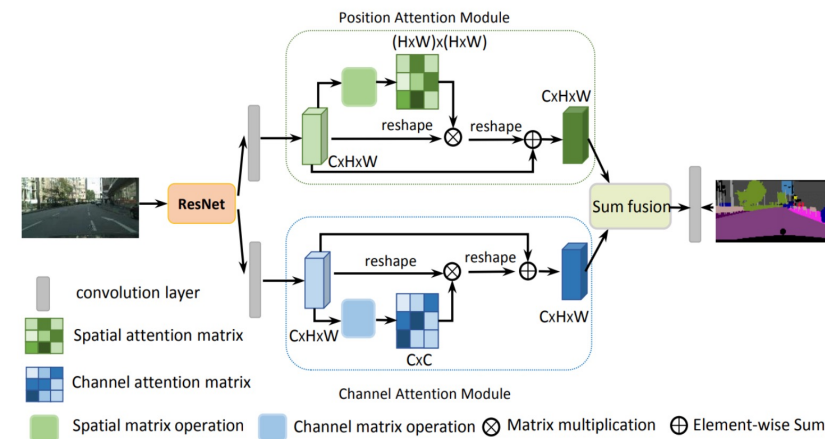
Context

**PSPNet**, Jonathan et al. CVPR18



**Pyramid Scene Parsing**

**DANet**, Jonathan et al. CVPR19



**Dual (Spatial and channel) Attention**

# Exploiting context in semantic segmentation models

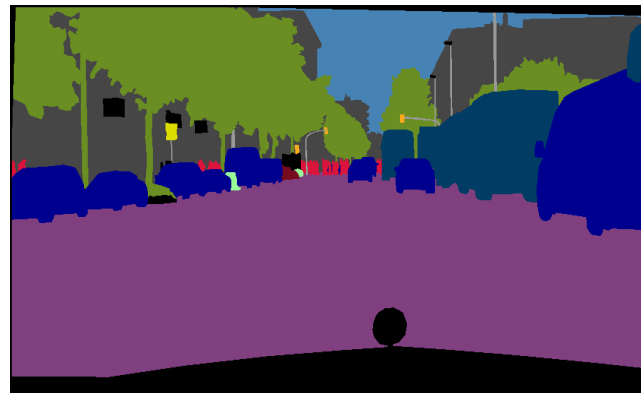


## Key finding

We discover that context empowers the attackers to fool objects far away from the perturbed area. For example, a perturbation of size 4% of image area fools the prediction at 60% of image area for PSANet.



Image perturbed with 9%

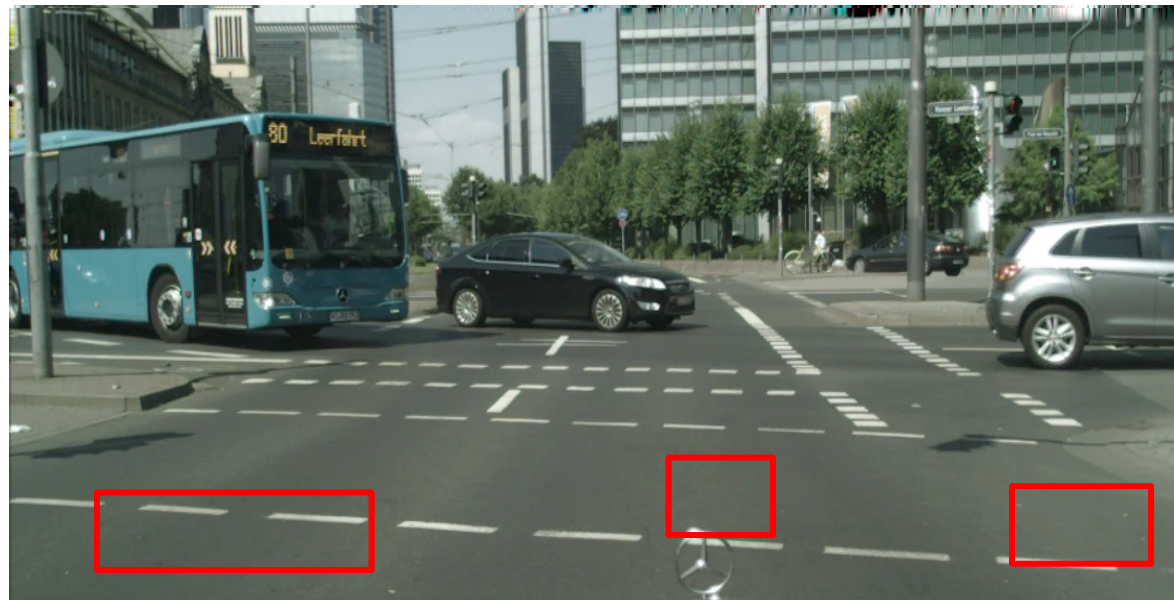


Clean prediction



Adversarial prediction

Perturbing static regions (e.g., road, sidewalk) affects the predictions at far away **dynamic** regions (e.g., bus, pedestrians) in inconspicuous way



Adversarial image

Indirect local attack  
in red box regions



Realistic-looking segmentation map



Adversarial predictions

Fools the distant  
dynamic objects

# Context-aware networks are highly vulnerable to indirect attacks than FCN



Adversarial image with **local perturbations**



Clean segmentation



FCN



PSANet



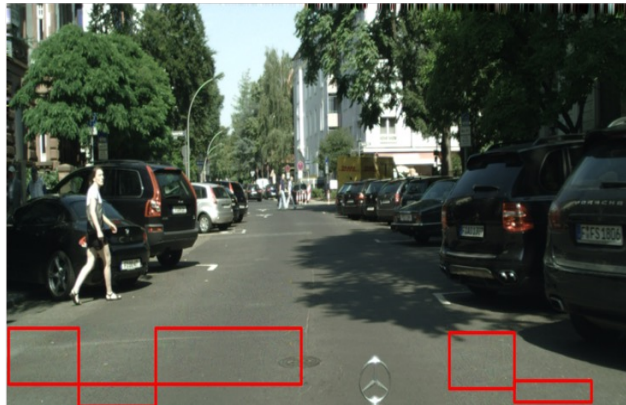
PSPNet



DANet

# Indirect adaptive attack results

Adversarial image



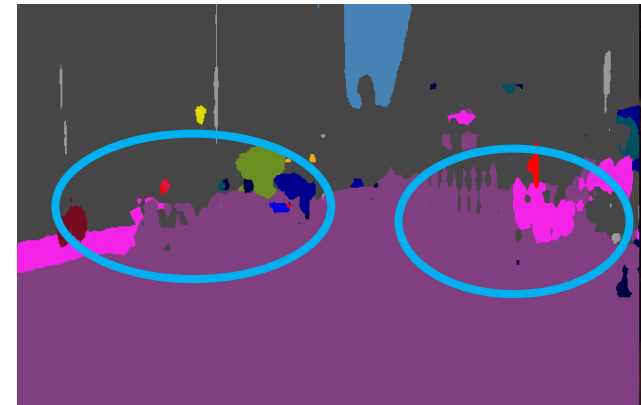
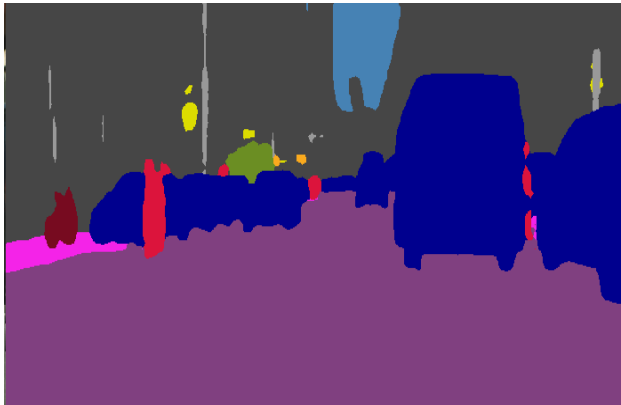
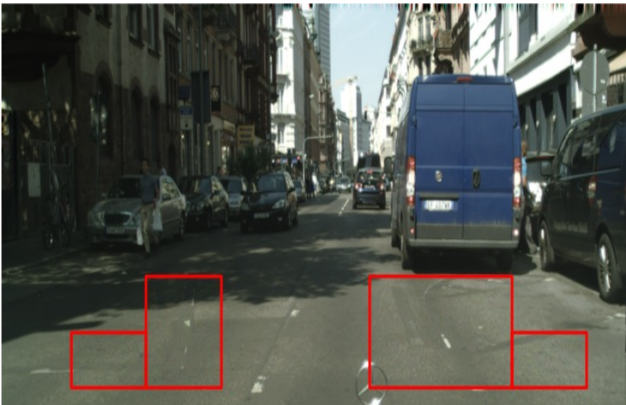
Clean predictions



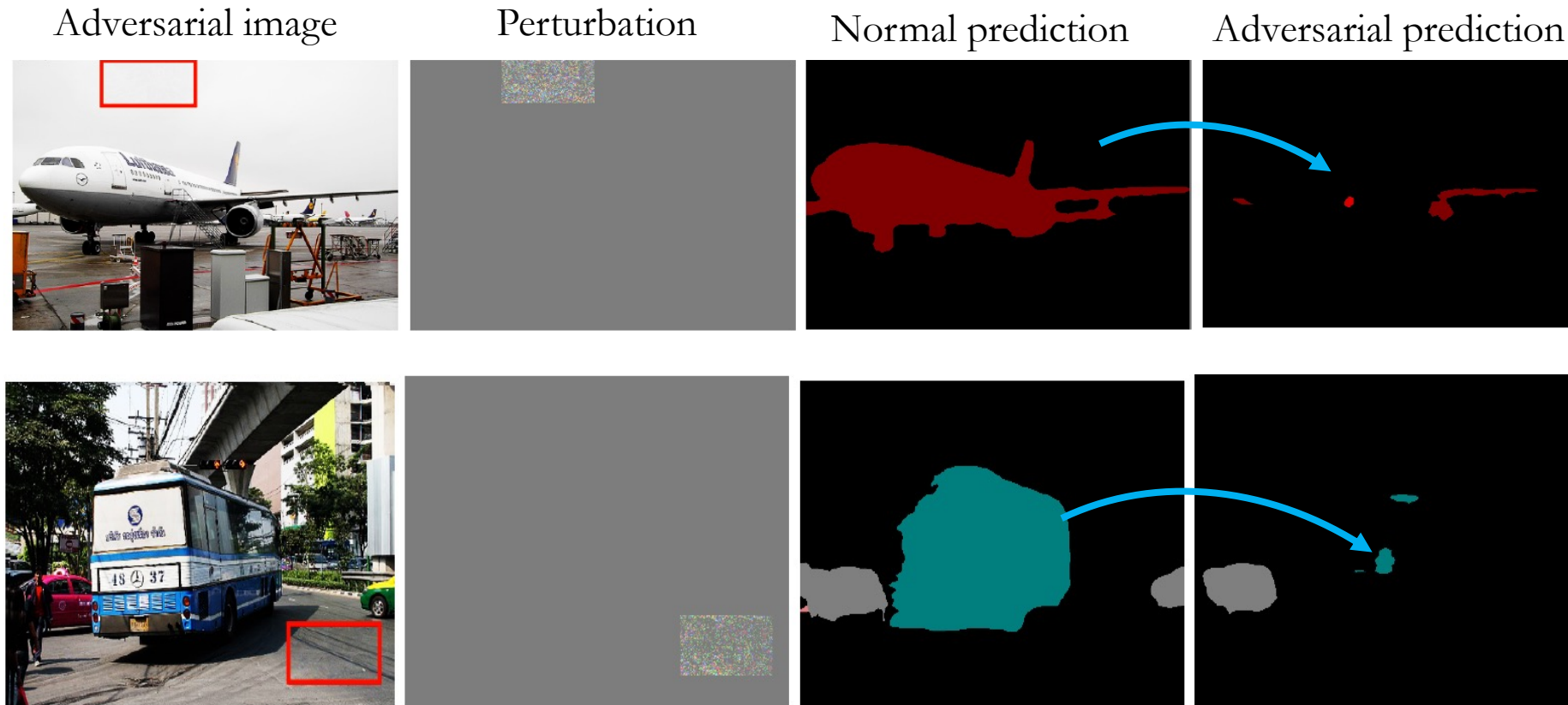
Adversarial predictions



Perturbed  
inside red  
boxes

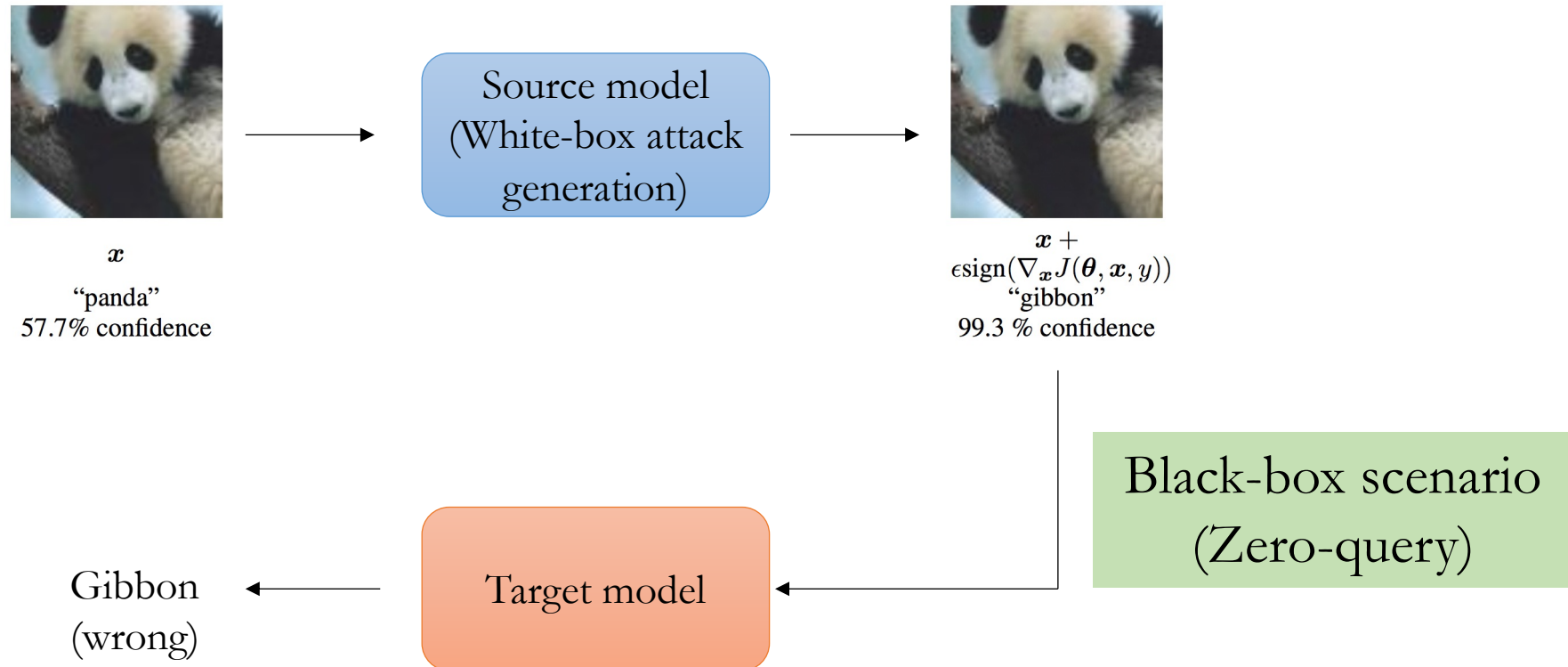


# Another perspective: Indirect attacks helps to understand contextual dependencies in DNNs (e.g. sky-aeroplane, road-car, road-bus)



# Transferable Adversarial Examples

(Perturbation generated from one network transfers to other network)



Transferable perturbations require **no** access to the target model

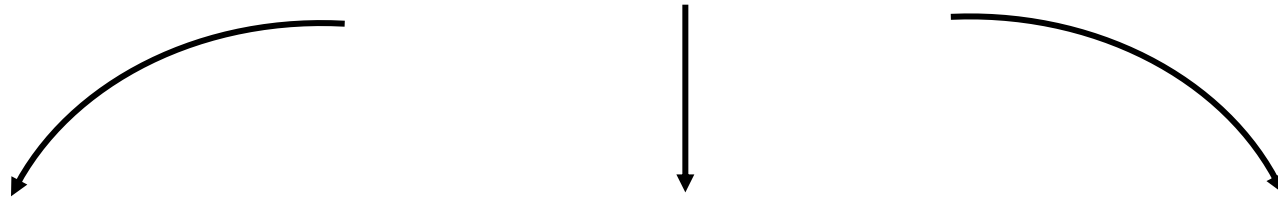


# Adversarial Attacks

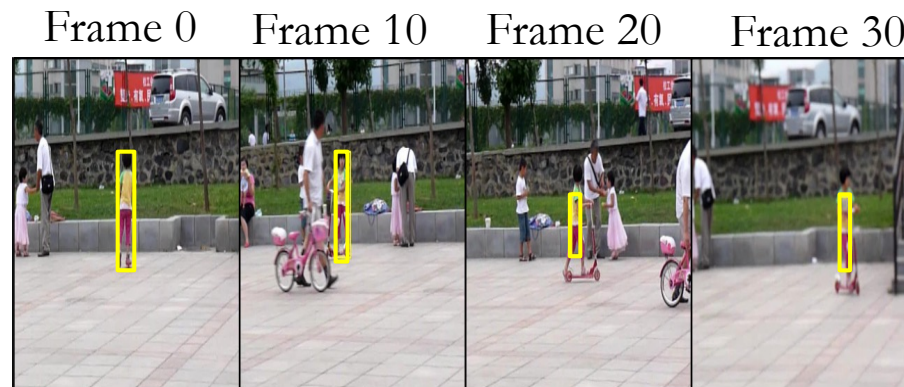
- Attacks beyond image recognition
  - White-box attacks on semantic segmentation
  - **Black-box transfer attacks on visual object tracker**
- Improving the transferability of attacks
  - Learning transferable transferable perturbations
- How can we use adversarial attack to improve DNNs
  - Semantic adversarial attacks to study disentanglement

# Challenges in transfer-based black-box attacks on object trackers

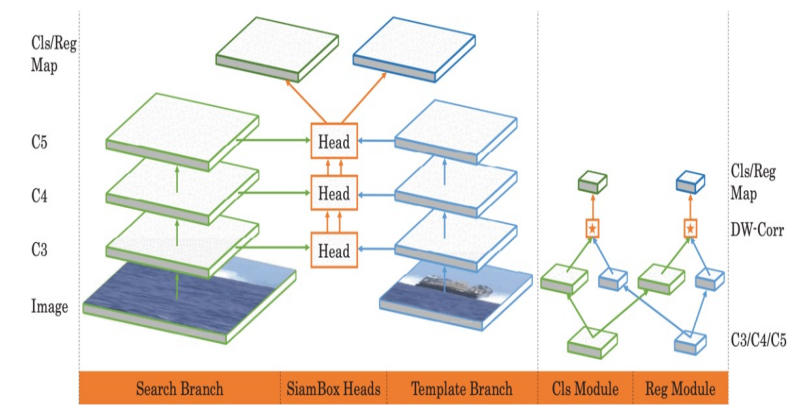
(VOT takes template as input and detects it in all subsequent search images)



**Novel objects** at test time  
(non-overlapping with training  
objects)



Computing perturbation per frame  
should be **efficient** as trackers work at  
real-time



Should **generalize** to different tracker  
frameworks such as SiamCAR,  
SiamBAN, OCEAN

# Goal: Efficient black-box attacks on visual object tracking



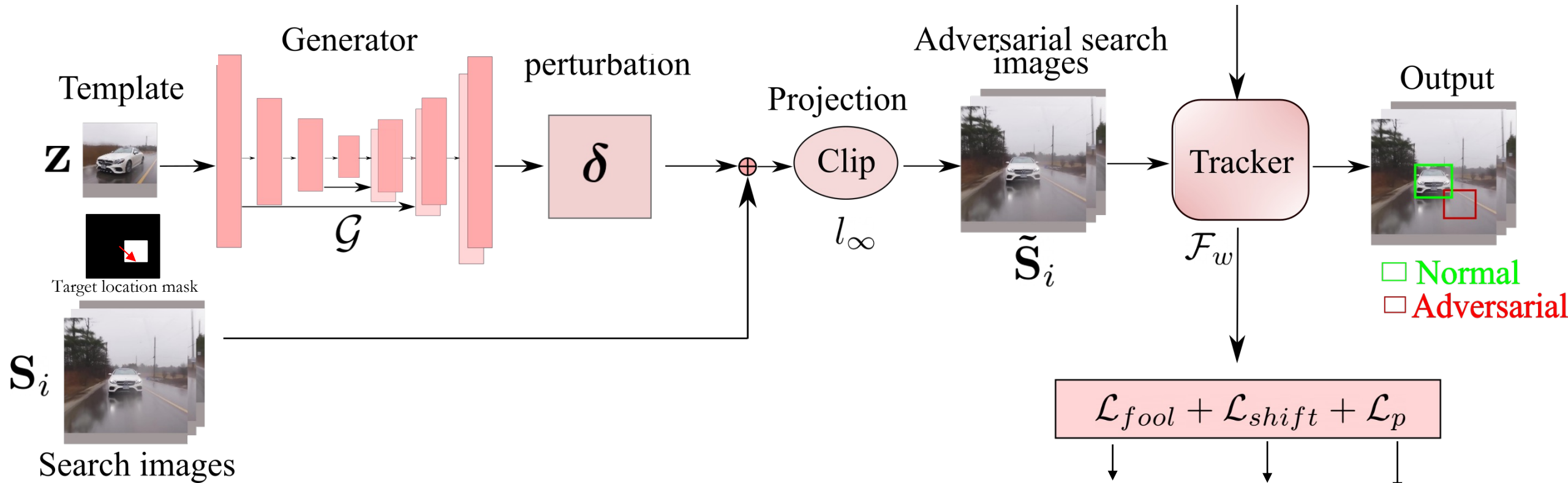
## Key idea

We propose to learn to generate a **single** perturbation from the object **template** only, that can be added to **every** search image and still successfully fool the tracker for the entire video.

Impact:

Learns to generate powerful transferable perturbations on unknown videos and trackers

# Temporally-transferable perturbation generator



## Design principles to improve transferability

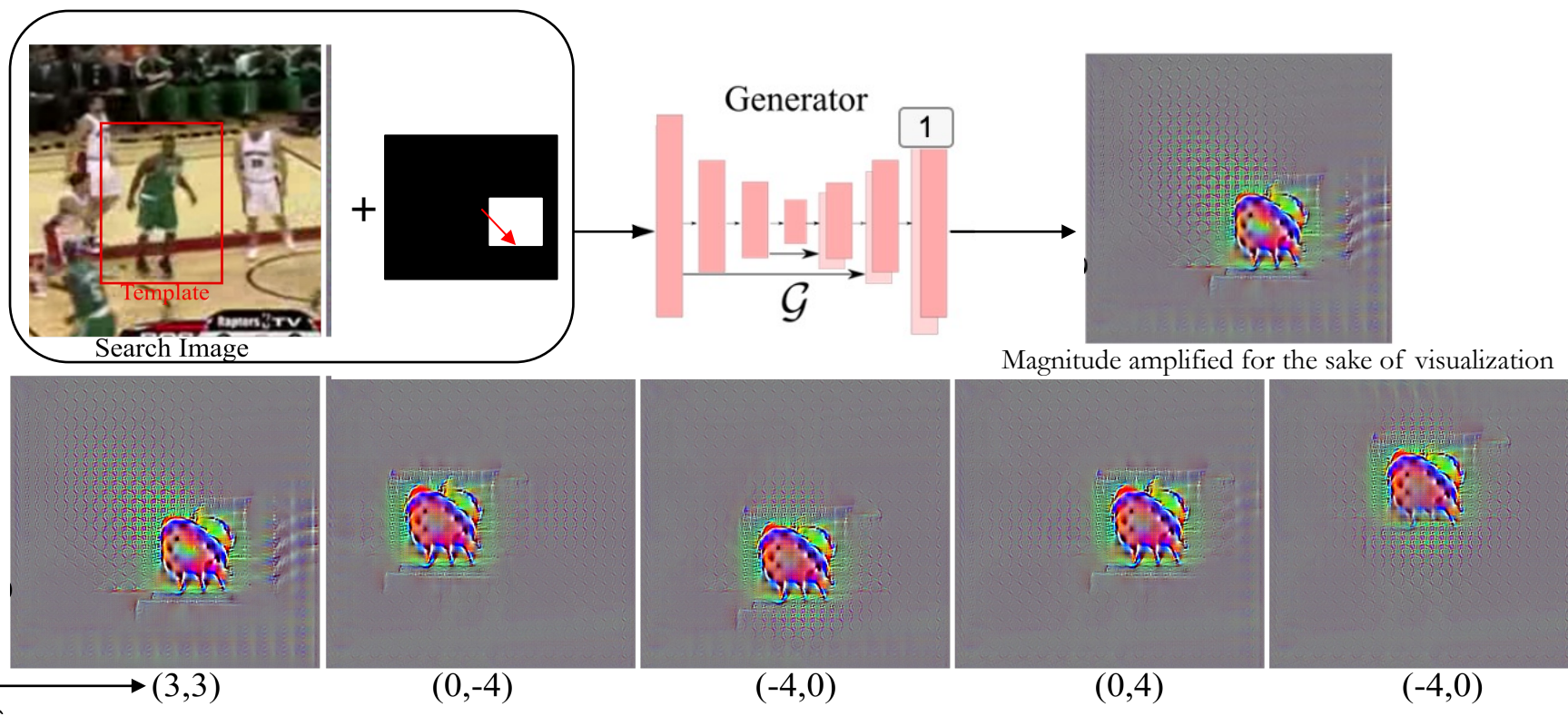
1. Perturbation is generated from **template** only
2. Generated perturbation is **shared** across all search images

Decrease the confidence at predicted position

Increase the confidence at target position

Perceptual loss b/w clean & adversarial image

# Visualization of generated directional perturbations



Generated perturbation contains **adversarial object-like patch** at target position

# Targeted adversarial attacks to steer the target tracker to follow the object at a fixed offset

Using 12 precomputed directional perturbations

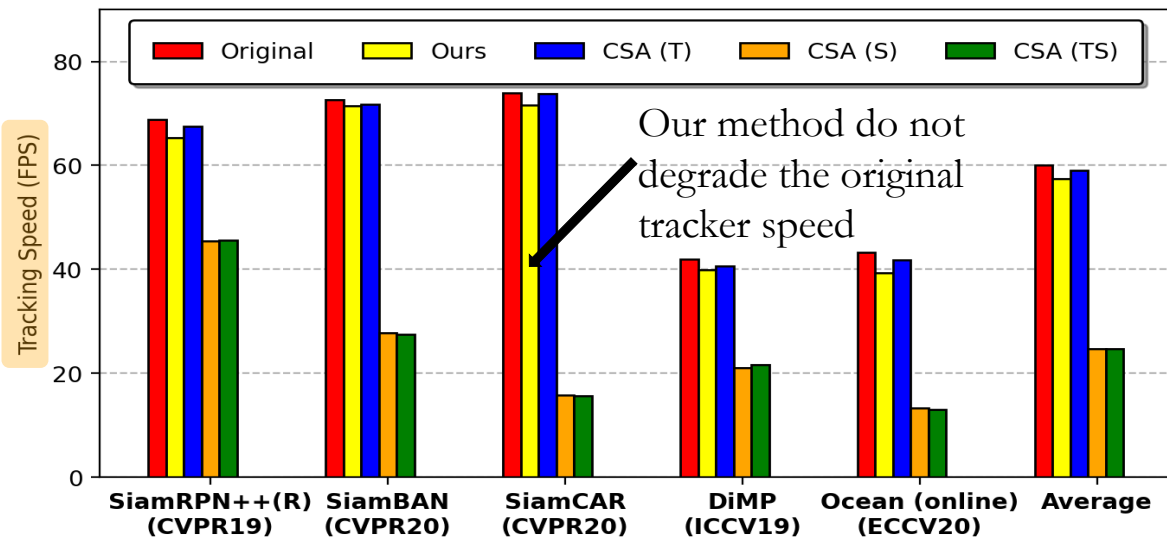


Ground truth

Adversarial prediction

Adversarial target

# Our transfer-attacks are highly efficient and effective too



Y-axis is **tracker speed** and x-axis is different tracker frameworks

Generator trained on SiamRPN++ (ResNet50)

S: Success; P: Precision

Methods	SiamRPN++ (M)		SiamBAN		SiamCAR		DiMP		Ocean online	
	S (↑)	P (↑)	S (↑)	P (↑)	S (↑)	P (↑)	S (↑)	P (↑)	S (↑)	P (↑)
Normal	0.657	0.862	0.692	0.910	0.696	0.908	0.650	0.847	0.669	0.884
CSA(T)	0.613	0.833	0.590	0.793	0.657	0.852	0.649	0.849	0.614	0.843
CSA(S)	0.281	0.440	0.371	0.531	0.373	0.536	0.641	0.840	0.390	0.645
CSA(TS)	0.348	0.431	0.347	0.510	0.391	0.559	0.642	0.844	0.423	0.705
Ours <sub>f</sub> (TD)	0.347	0.528	0.478	0.720	0.444	0.599	0.643	0.839	0.492	0.768
Ours (TD)	<b>0.217</b>	<b>0.281</b>	<b>0.198</b>	<b>0.254</b>	<b>0.292</b>	<b>0.377</b>	<b>0.631</b>	<b>0.821</b>	<b>0.345</b>	<b>0.452</b>
Ours <sub>f</sub>	0.408	0.616	0.478	0.721	0.567	0.770	0.646	0.843	0.592	0.829
Ours	<b>0.212</b>	<b>0.272</b>	<b>0.198</b>	<b>0.253</b>	<b>0.292</b>	<b>0.374</b>	<b>0.638</b>	<b>0.837</b>	<b>0.338</b>	<b>0.440</b>

Performance of **Ours** (Universal perturbation generated from single fixed template) and **Ours (TD)** (perturbations generated from template of given input video) method are at same range

# Adversarial Attacks

- Attacks beyond image recognition
  - White-box attacks on semantic segmentation
  - Black-box transfer attacks on visual object tracker
- Improving the transferability of attacks
  - **Learning transferable transferable perturbations**
- How can we use adversarial attack to study DNNs
  - Semantic adversarial attacks to study disentanglement



# Understanding and Improving the Transferability of Generative Adversarial Perturbations

# Learning Transferable Adversarial Perturbations

We investigate the transferability of **generative** perturbations when the conditions at inference time differ from the training ones in terms of

## 1. Target architecture

Generator was trained to attack a VGG-16 but the target network is a ResNet152

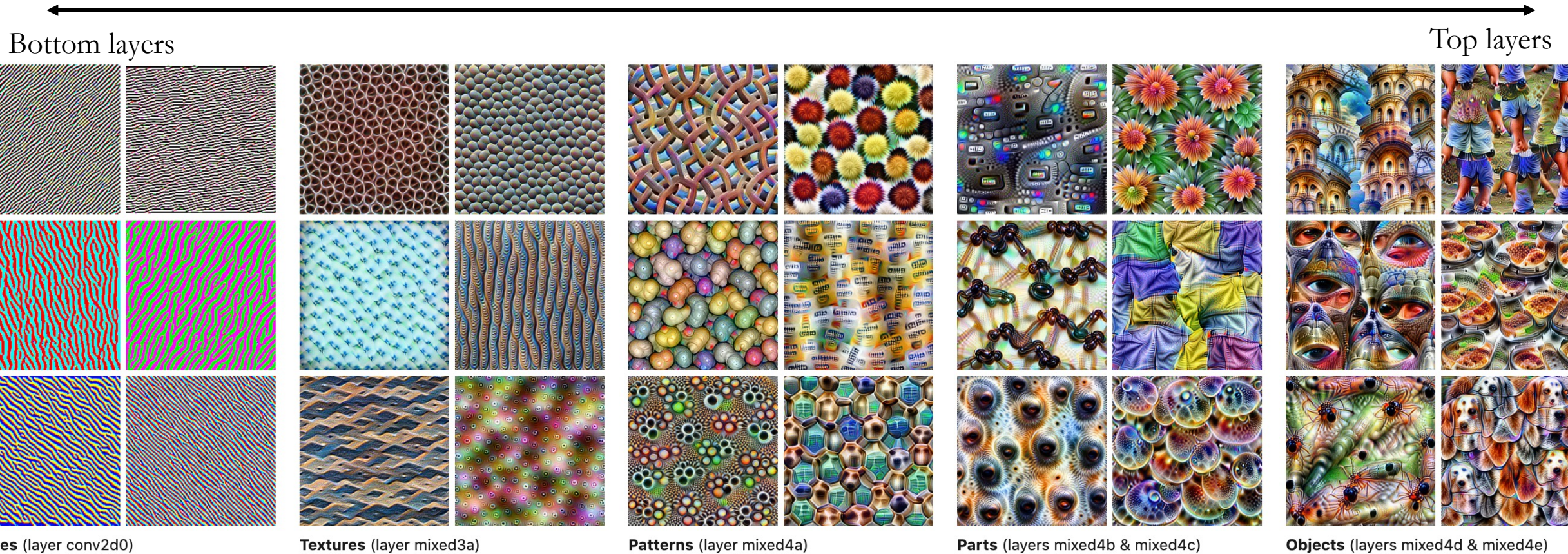
## 2. Target data

Generator was trained using the Paintings dataset, but the test data comes from ImageNet

## 3. Target task

Generator was trained to attack an image recognition model but faces an object detector at test time

# Now let's take a step back and see how deep neural networks build up their understanding of images?



Inception filter visualization

# Learning Transferable Adversarial Perturbations



## Key idea

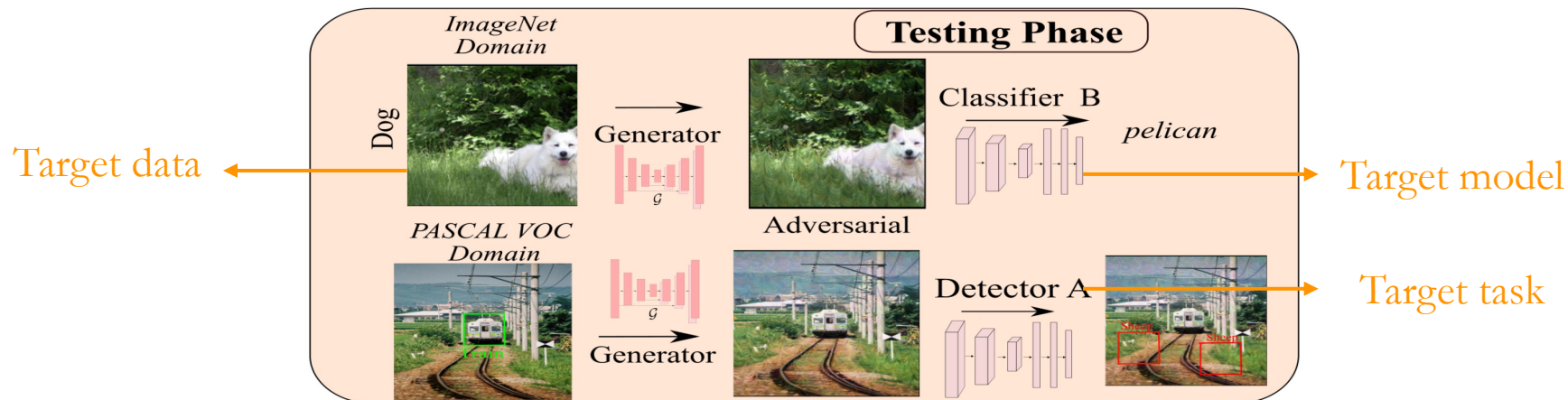
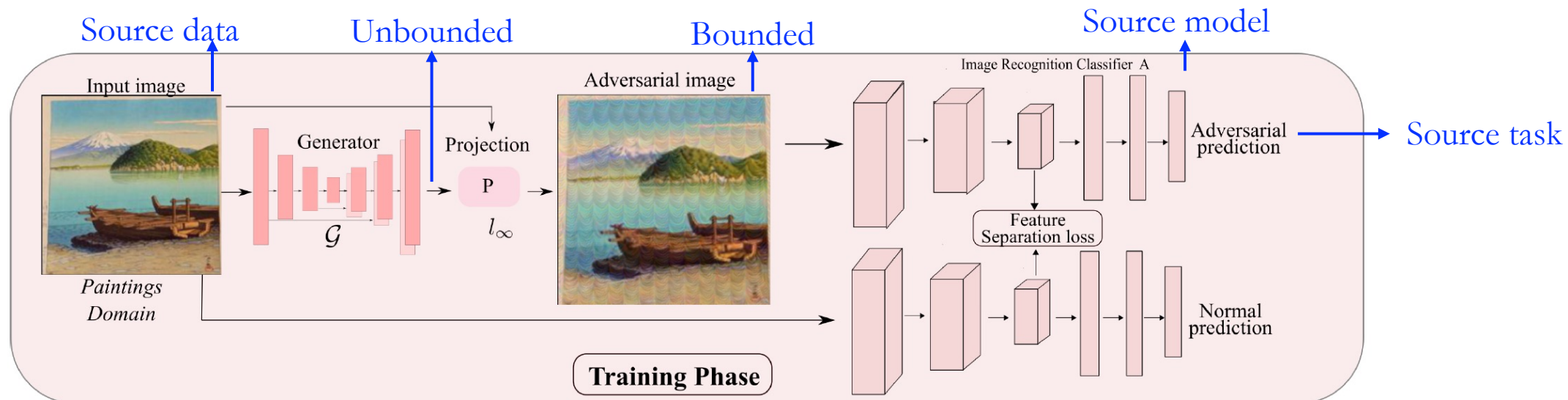
Disrupting the **mid-level features** using feature separation loss empowers attackers to learn perturbations with high transfer rates across target architectures, target datasets and target tasks without any queries

Most prior works in transfer-based black-box attacks focus on only unknown architecture



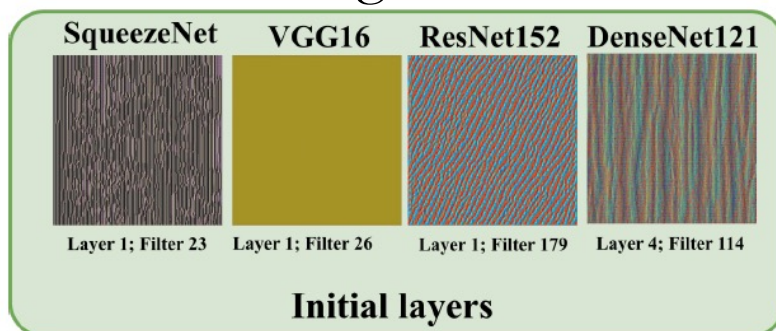
Our work focuses on black-box attacks on **unknown architecture, data and task**

# Training a perturbation generator with mid-level feature separation loss



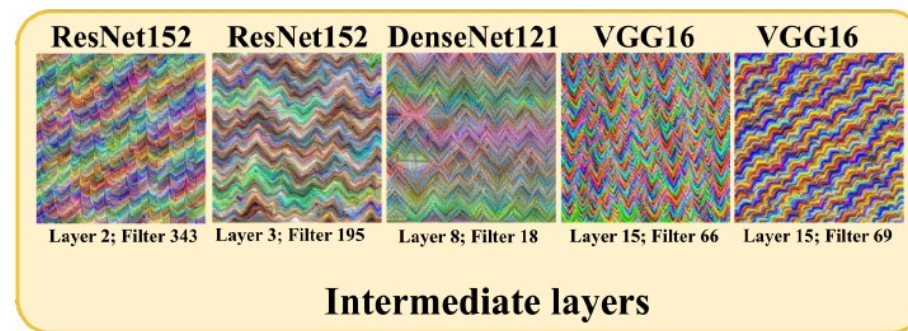
# CNNs with different architectures share similar filter bank

## Edges, colors



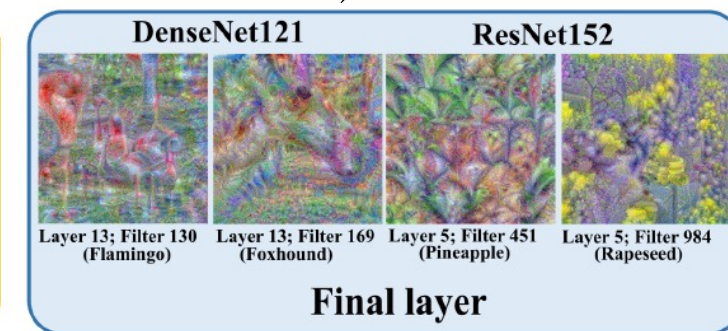
Disrupting low-level filters requires larger change in input image. Further it does not transfer well to top layer to change the output prediction

## Textures



Disrupting mid-level filters helps to learn transferable perturbations across architectures

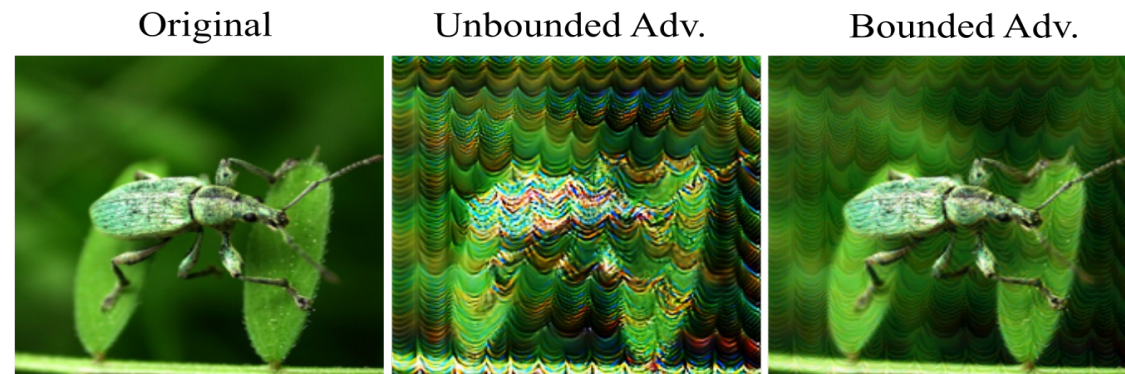
## Object Patterns



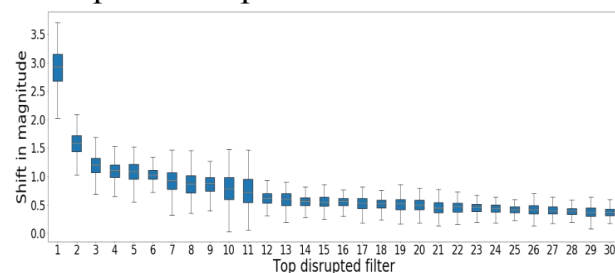
Disrupting top-level overfits to source classifier's decision boundaries

# Understanding the high transfer rates from ResNet152 to VGG16

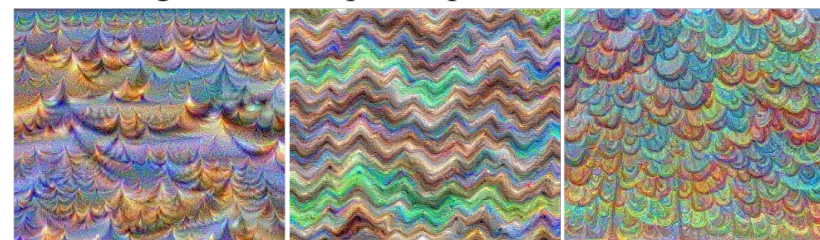
## (a) White-box attack on ResNet152 (Fooling Rate: 99.7%)



Top 30 disrupted filters of ResNet152



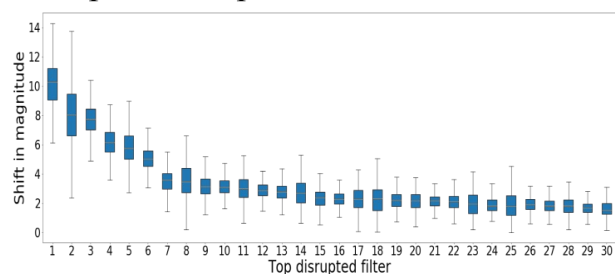
Synthesized images of few top disrupted filters in ResNet152 (Layer 3)



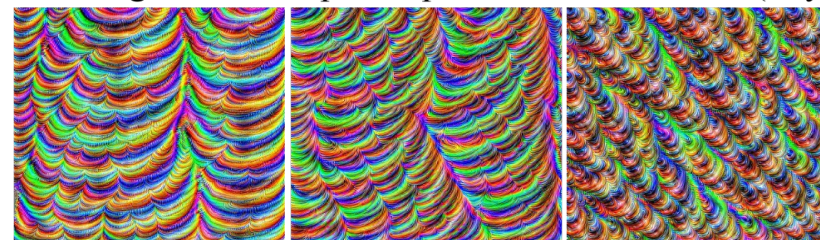
Top disrupted filters are similar. Thus, transfer rates are high between from ResNet152 to VGG-16

## (b) Transfer attack from ResNet152 to VGG16 (Fooling Rate: 99.1%)

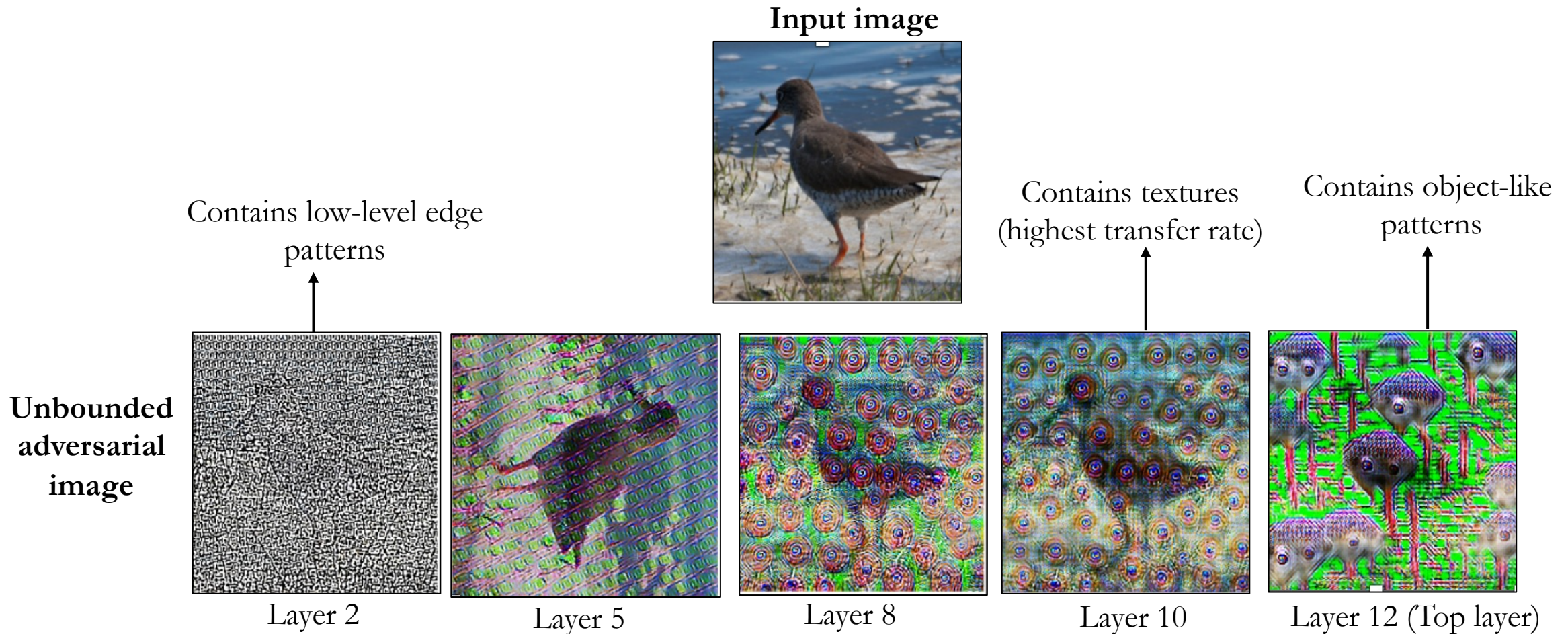
Top 30 disrupted filters of VGG16



Synthesized images of few top disrupted filters in VGG16 (Layer 15)



# Visualization of **unbounded** adversarial images with feature separation loss by attacking different layers on SqueezeNet



Attacking top layers overfit to source classifier  
while bottom layers require larger perturbation strength



# Standard black-box transferability

(access to substitute model on target data & target data)

Gen. Training (data)	Discriminator	VGG16	ResNet152	Inception-v3	DenseNet121	SqueezeNet1.1	Average	Foiling rate ↑
			GAP [10] / CDA [11] / Ours					
ImageNet (1.2M)	VGG-16	99.9* / 99.8* / 99.3*	53.5 / 53.6 / <b>68.4</b>	41.7 / 43.2 / <b>46.6</b>	58.9 / 66.5 / <b>84.7</b>	67.8 / 70.6 / <b>86.5</b>	64.4 / 66.7 / <b>77.1</b>	
	ResNet152	93.2 / 96.8 / <b>99.1</b>	97.6* / 99.6* / 99.7*	60.5 / 66.0 / <b>74.9</b>	87.5 / 94.2 / <b>98.8</b>	83.9 / 82.8 / <b>89.1</b>	84.5 / 87.9 / <b>92.3</b>	
	Inception-v3	88.2 / 97.2 / <b>99.0</b>	83.4 / 82.7 / <b>90.4</b>	96.5* / 98.7* / 99.6*	89.5 / 93.6 / <b>96.7</b>	90.9 / 92.0 / <b>93.2</b>	89.7 / 92.9 / <b>95.8</b>	
	DenseNet121	94.9 / 95.0 / <b>99.4</b>	89.5 / 91.0 / <b>98.7</b>	56.1 / 57.7 / <b>86.0</b>	99.6* / 99.6* / 99.6*	79.7 / 81.5 / <b>95.6</b>	84.0 / 85.0 / <b>95.9</b>	
	SqueezeNet	88.0 / 91.5 / <b>96.1</b>	50.4 / 57.1 / <b>76.4</b>	48.0 / 47.6 / <b>70.7</b>	64.0 / 69.0 / <b>88.7</b>	99.8* / 99.7* / 99.7*	70.0 / 73.0 / <b>86.3</b>	
	Average	92.8 / 96.1 / <b>98.6</b>	74.9 / 76.8 / <b>86.7</b>	60.6 / 62.6 / <b>75.6</b>	79.9 / 84.6 / <b>93.7</b>	77.6 / 78.5 / <b>89.5</b>	78.5 / 81.1 / <b>89.5</b>	

Source

Target

Source and Target models are trained on same ImageNet data but differ in architecture

# Strict black-box transferability

(access to substitute model on target data but no target data)

Gen. Training (data)	Discriminator (ImageNet)	VGG16	ResNet152	Inception-v3	DenseNet121	SqueezeNet 1.1	Average	Foiling rate ↑
			GAP [10] / CDA [11] / Ours					
Comics (40K)	VGG-16	99.8 / <b>99.9</b> / 99.5	54.3 / 54.0 / <b>77.4</b>	45.8 / 45.2 / <b>61.9</b>	66.3 / 64.2 / <b>93.6</b>	70.7 / 68.4 / <b>93.4</b>	67.4 / 66.3 / <b>85.1</b>	
	ResNet152	75.3 / 95.8 / <b>99.3</b>	97.6 / 98.1 / <b>99.6</b>	31.7 / 66.5 / <b>73.1</b>	45.1 / 87.7 / <b>98.6</b>	67.3 / 86.0 / <b>90.7</b>	63.4 / 86.8 / <b>92.3</b>	
	Inception V3	84.3 / 85.6 / <b>99.0</b>	97.2 / 97.3 / <b>90.4</b>	99.8 / <b>99.8</b> / 99.6	88.5 / 87.9 / <b>96.7</b>	82.4 / 82.3 / <b>93.2</b>	90.5 / 90.6 / <b>95.8</b>	
	DenseNet121	96.9 / 87.8 / <b>96.5</b>	<b>98.0</b> / 55.7 / 93.0	83.1 / 48.5 / <b>82.5</b>	99.4 / 97.7 / <b>98.8</b>	78.3 / 81.2 / <b>91.9</b>	91.2 / 74.2 / <b>92.5</b>	
	SqueezeNet	87.7 / 89.9 / <b>96.5</b>	54.0 / 58.2 / <b>79.0</b>	51.2 / 51.4 / <b>75.4</b>	68.7 / 76.3 / <b>90.2</b>	99.7 / <b>99.8</b> / 99.7	72.3 / 75.1 / <b>88.2</b>	
	Average	88.1 / 91.8 / <b>98.2</b>	80.2 / 72.6 / <b>87.8</b>	62.3 / 62.3 / <b>78.5</b>	73.6 / 82.8 / <b>95.6</b>	79.6 / 83.6 / 93.8	79.9 / 78.6 / <b>90.8</b>	
Paintings (80K)	VGG-16	99.4 / <b>99.9</b> / 99.0	41.1 / 57.6 / <b>66.6</b>	36.5 / 46.6 / <b>50.0</b>	50.8 / 73.8 / <b>84.6</b>	63.7 / 73.0 / <b>86.4</b>	58.3 / 70.1 / <b>77.3</b>	
	ResNet152	80.4 / 89.9 / <b>98.7</b>	95.4 / 97.5 / <b>99.4</b>	50.7 / 62.1 / <b>72.8</b>	70.4 / 82.3 / <b>97.9</b>	70.4 / 81.1 / <b>89.2</b>	73.5 / 82.6 / <b>91.6</b>	
	Inception V3	80.3 / 80.5 / <b>98.6</b>	95.8 / <b>96.4</b> / 88.2	99.6 / <b>99.6</b> / 99.5	87.7 / 87.2 / <b>95.2</b>	77.5 / 72.8 / <b>90.8</b>	88.2 / 87.3 / <b>94.5</b>	
	DenseNet121	87.6 / 86.5 / <b>96.2</b>	80.1 / 81.2 / <b>90.9</b>	51.4 / 50.4 / <b>76.0</b>	98.8 / <b>98.9</b> / 97.4	73.6 / 73.7 / <b>91.7</b>	67.7 / 78.1 / <b>90.5</b>	
	SqueezeNet	82.8 / 80.7 / <b>95.2</b>	46.0 / 46.0 / <b>73.4</b>	44.5 / 47.4 / <b>71.0</b>	59.3 / 56.5 / <b>87.2</b>	99.4 / 99.3 / <b>99.6</b>	66.4 / 66.0 / <b>85.3</b>	
	Average	86.1 / 87.5 / <b>97.6</b>	71.7 / 75.8 / <b>83.7</b>	56.5 / 61.2 / <b>73.9</b>	73.4 / 79.7 / <b>92.5</b>	76.9 / 80.0 / 91.5	72.9 / 76.8 / <b>87.8</b>	
ChestX (10K)	VGG-16	78.7 / 85.6 / <b>93.3</b>	23.2 / 23.3 / <b>41.8</b>	25.5 / 27.9 / <b>31.3</b>	27.5 / 28.2 / <b>53.4</b>	46.1 / 48.0 / <b>64.3</b>	40.2 / 42.6 / <b>56.8</b>	
	ResNet152	39.9 / 44.8 / <b>56.4</b>	27.0 / 25.3 / <b>62.8</b>	28.2 / 25.7 / <b>27.7</b>	25.9 / 26.6 / <b>38.1</b>	44.9 / 47.1 / <b>60.5</b>	33.2 / 33.9 / <b>49.2</b>	
	Inception V3	56.0 / 50.3 / <b>91.6</b>	35.9 / 32.0 / <b>69.5</b>	44.4 / 35.1 / <b>84.9</b>	45.9 / 35.4 / <b>77.4</b>	65.1 / 57.7 / <b>75.6</b>	49.5 / 42.1 / <b>79.8</b>	
	DenseNet121	42.8 / 42.3 / <b>64.0</b>	26.4 / 25.2 / <b>44.2</b>	28.0 / 28.8 / <b>34.0</b>	41.9 / 48.2 / <b>76.0</b>	54.2 / 48.8 / <b>60.2</b>	38.7 / 38.7 / <b>55.7</b>	
	SqueezeNet	51.7 / 51.1 / <b>81.1</b>	27.9 / 31.6 / <b>52.5</b>	30.2 / 33.1 / <b>47.1</b>	31.6 / 35.1 / <b>64.2</b>	81.3 / 78.9 / <b>96.4</b>	44.5 / 46.0 / <b>68.3</b>	
	Average	53.8 / 54.8 / <b>77.2</b>	28.1 / 27.4 / <b>54.2</b>	31.3 / 30.1 / <b>45.0</b>	34.6 / 34.7 / <b>61.9</b>	58.3 / 56.1 / <b>71.4</b>	41.2 / 40.6 / <b>62.0</b>	

# Extreme Cross-domain Transferability

(neither access to substitute model on target model nor target data)

Source Target

Source and Target models are trained on **different** data and also differ in architecture

Target models on CUB200

Gen. Training (data)	Discriminator ( ImageNet)	ResNet50	SeNET154	SeResNet101	Average
		GAP [209] / CDA [188] / Ours			
ImageNet (1.2M)	VGG-16	41.25 / 24.59 / <b>76.15</b>	41.44 / 30.43 / <b>45.82</b>	29.75 / 23.01 / <b>35.85</b>	37.48 / 26.01 / <b>52.61</b>
	ResNet152	54.82 / 52.78 / <b>93.18</b>	50.76 / 50.72 / <b>77.44</b>	46.00 / 45.13 / <b>65.00</b>	50.35 / 49.54 / <b>78.54</b>
	Inception-v3	40.78 / 55.63 / <b>70.40</b>	33.07 / 36.49 / <b>48.10</b>	35.12 / 36.59 / <b>39.52</b>	36.32 / 42.90 / <b>52.67</b>
	DenseNet121	52.95 / 50.97 / <b>90.66</b>	38.52 / 43.42 / <b>73.30</b>	45.36 / 46.10 / <b>63.07</b>	45.61 / 46.83 / <b>75.68</b>
	SqueezeNet	36.40 / 35.57 / <b>63.89</b>	34.04 / 25.55 / <b>47.32</b>	34.57 / 30.51 / <b>39.39</b>	35.00 / 30.54 / <b>50.20</b>
	Average	45.13 / 43.91 / <b>78.86</b>	39.57 / 37.32 / <b>58.40</b>	38.16 / 36.27 / <b>48.57</b>	40.95 / 39.17 / <b>61.94</b>

(a) CUB200

Gen. Training (data)	Discriminator ( ImageNet)	ResNet50	SeNET154	SeResNet101	Average
		GAP [209] / CDA [188] / Ours			
ImageNet (1.2M)	VGG-16	18.07 / 48.65 / <b>70.22</b>	32.35 / 30.03 / <b>32.41</b>	12.66 / 14.76 / <b>21.73</b>	21.03 / 31.15 / <b>41.45</b>
	ResNet152	37.08 / 71.27 / <b>94.80</b>	33.25 / 34.31 / <b>62.74</b>	22.73 / 31.51 / <b>62.23</b>	31.02 / 45.70 / <b>73.26</b>
	Inception-v3	51.27 / 44.12 / <b>44.34</b>	35.63 / 36.25 / <b>38.59</b>	31.68 / 25.43 / <b>25.83</b>	39.53 / 35.27 / <b>36.25</b>
	DenseNet121	59.84 / 57.46 / <b>98.32</b>	28.98 / 34.09 / <b>65.27</b>	24.71 / 25.43 / <b>71.76</b>	37.84 / 38.97 / <b>78.45</b>
	SqueezeNet	26.07 / 30.32 / <b>85.33</b>	17.09 / 16.06 / <b>31.69</b>	14.40 / 18.19 / <b>31.54</b>	19.19 / 21.52 / <b>49.52</b>
	Average	38.47 / 50.36 / <b>78.60</b>	29.46 / 30.15 / <b>46.14</b>	21.24 / 23.05 / <b>42.62</b>	29.72 / 34.52 / <b>55.79</b>

(b) Stanford Cars

Gen. Training (data)	Discriminator ( ImageNet)	ResNet50	SeNET154	SeResNet101	Average
		GAP [209] / CDA [188] / Ours			
ImageNet (1.2M)	VGG-16	25.20 / 23.97 / <b>79.36</b>	46.77 / 38.79 / <b>37.28</b>	36.15 / 27.42 / <b>38.16</b>	36.04 / 30.06 / <b>51.60</b>
	ResNet152	42.87 / 64.45 / <b>96.82</b>	49.02 / 53.35 / <b>91.63</b>	36.72 / 56.80 / <b>86.44</b>	42.87 / 58.20 / <b>91.63</b>
	Inception-v3	49.38 / 43.95 / <b>72.61</b>	54.25 / 35.25 / <b>59.41</b>	46.28 / 43.11 / <b>42.87</b>	49.97 / 40.77 / <b>58.30</b>
	DenseNet121	37.11 / 37.05 / <b>93.10</b>	38.73 / 41.04 / <b>88.30</b>	35.22 / 36.93 / <b>83.59</b>	37.02 / 38.34 / <b>88.33</b>
	SqueezeNet	26.07 / 33.63 / <b>82.30</b>	27.18 / 27.57 / <b>41.70</b>	38.40 / 42.78 / <b>52.51</b>	30.55 / 34.66 / <b>58.84</b>
	Average	36.13 / 40.61 / <b>84.84</b>	43.19 / 39.20 / <b>63.66</b>	38.55 / 41.41 / <b>60.71</b>	39.29 / 40.41 / <b>69.74</b>

Fooling rate ↑

(c) Aircraft

Avg. 23% improvement over CDA

# Cross-task transferability analysis

(ImageNet classifier → PASCAL VOC SSD detector)

No access to target data, target model and target task

mAP↑

Gen. Training (data)	Discriminator (Trained on ImageNet)	VGG16	ResNet50	EfficientNet	MobileNet-v3	Average
		GAP [10] / CDA [11] / Ours				
-	No Attack	68.12	66.08	61.07	55.44	62.68
Comics (40K)	VGG-16	19.9 / 25.2 / <b>9.08</b>	15.7 / 20.9 / <b>13.7</b>	<b>12.4</b> / 13.4 / 14.2	<b>9.22</b> / 13.1 / 14.5	14.3 / 18.1 / <b>10.1</b>
	ResNet152	31.0 / 25.3 / <b>13.7</b>	23.0 / 20.2 / <b>10.5</b>	23.9 / 17.3 / <b>12.2</b>	17.9 / 14.0 / <b>8.37</b>	24.0 / 19.2 / <b>11.2</b>
	Inception-v3	33.2 / 33.8 / <b>23.7</b>	27.9 / 27.7 / <b>25.1</b>	31.1 / 30.7 / <b>22.0</b>	20.2 / 18.7 / <b>18.4</b>	28.8 / 27.7 / <b>22.3</b>
	DenseNet121	22.1 / 26.3 / <b>12.2</b>	18.6 / 21.6 / <b>14.5</b>	17.8 / 20.1 / <b>16.2</b>	13.9 / 15.3 / <b>9.55</b>	18.1 / 20.8 / <b>13.1</b>
	SqueezeNet	29.4 / 32.6 / <b>18.1</b>	24.8 / 28.9 / <b>15.7</b>	20.5 / 24.4 / <b>17.7</b>	15.7 / 20.5 / <b>11.9</b>	22.6 / 26.6 / <b>15.9</b>
	Average	27.1 / 28.7 / <b>15.4</b>	22.0 / 23.8 / <b>15.9</b>	21.1 / 21.2 / <b>16.5</b>	15.4 / 16.3 / <b>11.7</b>	21.4 / 22.5 / <b>14.9</b>
Paintings (80K)	VGG-16	20.2 / 20.4 / <b>9.83</b>	21.4 / 22.5 / <b>13.2</b>	14.7 / 15.0 / <b>12.8</b>	<b>11.4</b> / 12.5 / 12.8	16.9 / 17.6 / <b>12.2</b>
	ResNet152	36.6 / 29.4 / <b>12.8</b>	26.7 / 21.9 / <b>12.5</b>	22.9 / 16.8 / <b>11.8</b>	21.3 / 17.6 / <b>9.40</b>	26.9 / 21.4 / <b>11.6</b>
	Inception-v3	32.3 / 33.5 / <b>16.8</b>	29.2 / 29.0 / <b>18.7</b>	28.1 / 28.5 / <b>14.3</b>	23.4 / 22.6 / <b>13.2</b>	28.3 / 28.4 / <b>15.7</b>
	DenseNet121	31.7 / 33.2 / <b>9.27</b>	23.1 / 23.2 / <b>11.0</b>	23.5 / 24.1 / <b>10.6</b>	20.2 / 20.9 / <b>6.53</b>	24.6 / 25.3 / <b>9.35</b>
	SqueezeNet	35.3 / 35.9 / <b>17.0</b>	28.5 / 29.0 / <b>13.7</b>	26.7 / 27.5 / <b>17.1</b>	21.0 / 21.1 / <b>8.77</b>	27.9 / 28.3 / <b>14.1</b>
	Average	31.3 / 30.5 / <b>13.1</b>	25.8 / 25.1 / <b>13.8</b>	23.2 / 22.4 / <b>13.3</b>	19.5 / 18.9 / <b>10.1</b>	25.0 / 24.2 / <b>12.6</b>
ImageNet (1.2M)	VGG-16	17.8 / 15.5 / <b>8.27</b>	19.2 / 13.9 / <b>11.8</b>	9.64 / <b>8.91</b> / 11.1	8.39 / <b>5.79</b> / 9.78	13.7 / 11.0 / <b>10.2</b>
	ResNet152	19.0 / 16.6 / <b>9.23</b>	13.5 / 14.6 / <b>7.67</b>	12.5 / 11.7 / <b>6.56</b>	12.4 / 7.67 / <b>4.29</b>	14.3 / 12.6 / <b>6.94</b>
	Inception-v3	<b>13.0</b> / 22.1 / 15.5	15.7 / 19.4 / <b>18.2</b>	13.8 / <b>12.5</b> / 13.5	11.3 / 15.1 / <b>11.6</b>	15.6 / 17.3 / <b>14.7</b>
	DenseNet121	21.5 / 16.1 / <b>7.60</b>	15.7 / 13.7 / <b>8.32</b>	13.8 / 11.4 / <b>7.73</b>	11.3 / 7.10 / <b>4.42</b>	15.6 / 12.1 / <b>7.02</b>
	SqueezeNet1	27.7 / 26.6 / <b>13.5</b>	23.7 / 22.5 / <b>10.8</b>	18.6 / 23.4 / <b>11.8</b>	15.2 / 17.2 / <b>7.40</b>	21.3 / 22.5 / <b>10.9</b>
	Average	19.8 / 19.3 / <b>10.8</b>	17.5 / 16.8 / <b>11.4</b>	13.5 / 13.6 / <b>10.1</b>	12.1 / 10.6 / <b>7.50</b>	15.7 / 15.1 / <b>9.95</b>

Source

Target

Source and Target models are trained on **different** data, task and also differ in architecture

# Adversarial Attacks

- Attacks beyond image recognition
  - White-box attacks on semantic segmentation
  - Black-box transfer attacks on visual object tracker
- Why adversarial attacks transfer?
  - Learning transferable perturbations
- How can we use adversarial attack to improve DNNs
  - **Semantic adversarial attacks to study disentanglement**

# Goal: Semantic attacks to study disentanglement of pose and appearance

## What is disentanglement?

Disentangled representations capture independent factors of variations in data

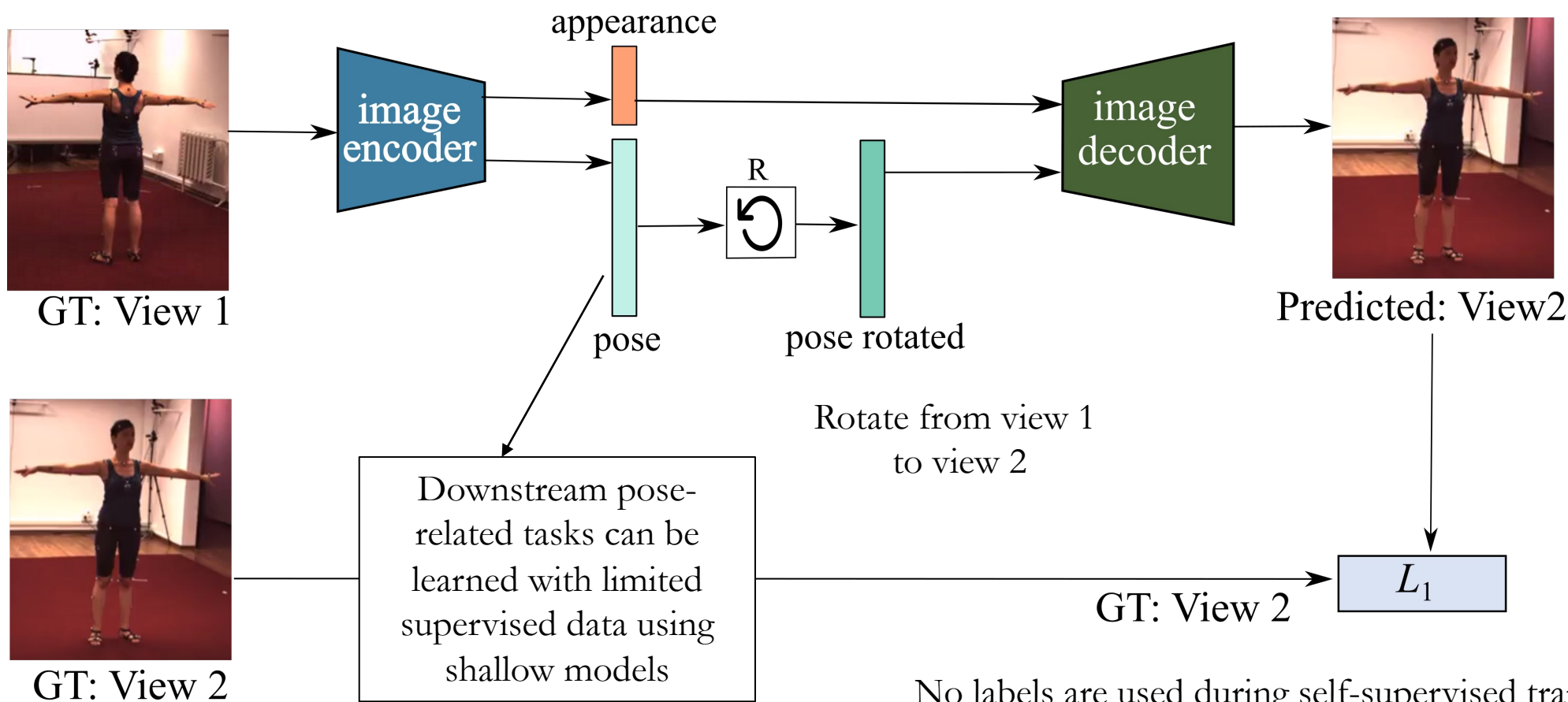
## Why do we need it?

Disentangled representations improves the performance of downstream tasks with limited supervision

# Background: Self-supervised Disentangled Representations (one technique using multi-view information)

Approach: Take one view as input and reconstruct the other view as output

at same time instant



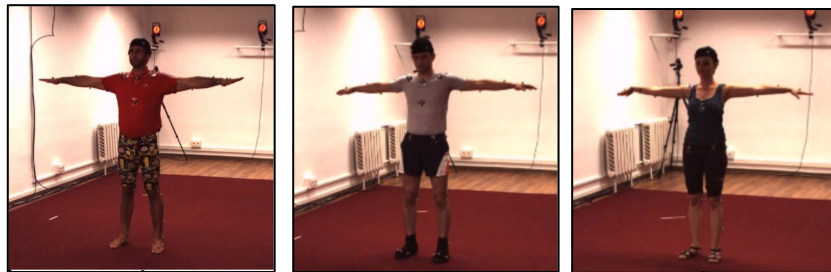
No labels are used during self-supervised training

# Analyze the disentanglement of pose and appearance

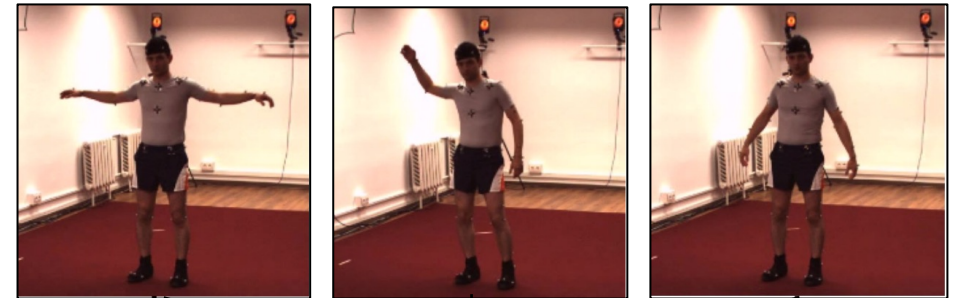
## Hypothesis

Given two disentangled latent codes that capture two underlying factors of variation in the input data, the adversarial modification of one factor in the input image should not alter the latent code encoded by another factor

Example of pose-appearance disentanglement

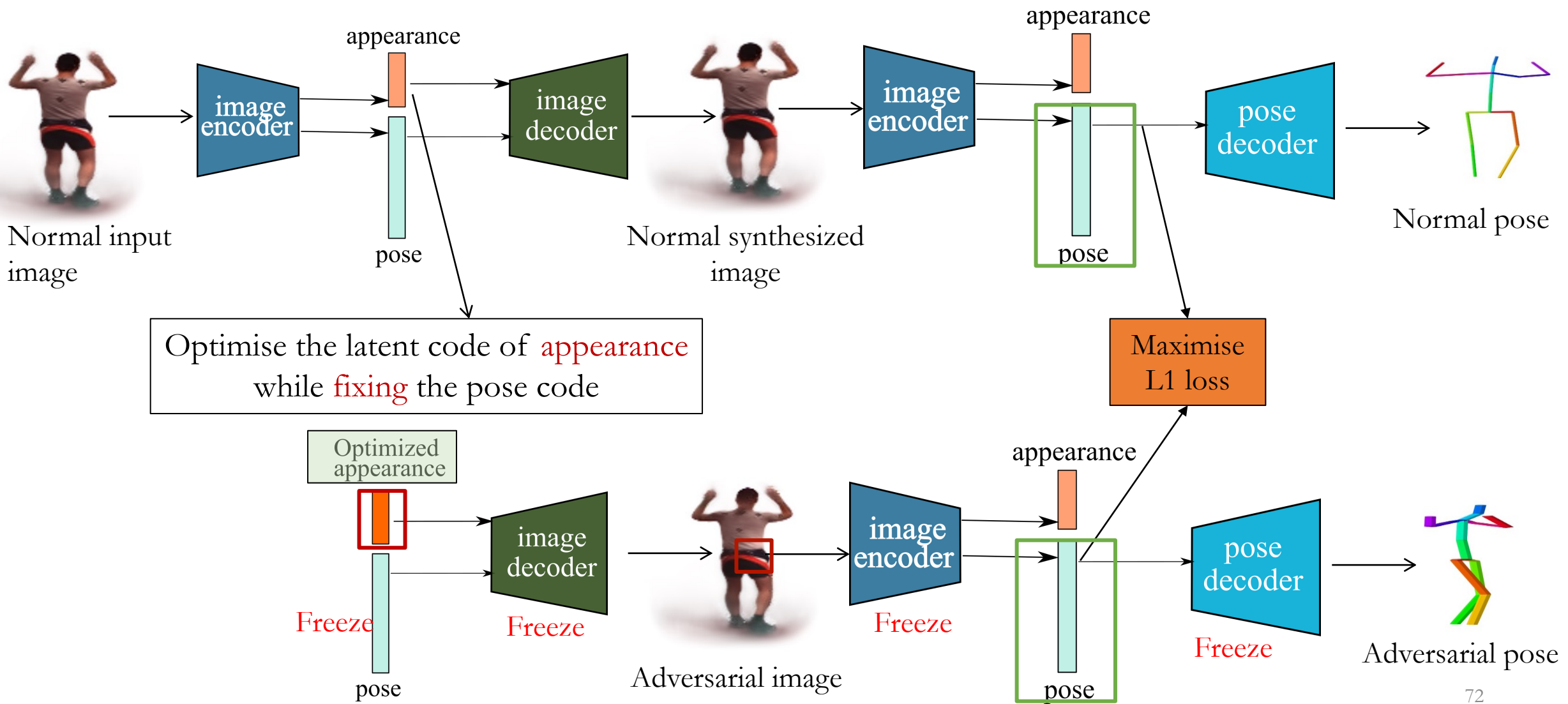


Same pose code



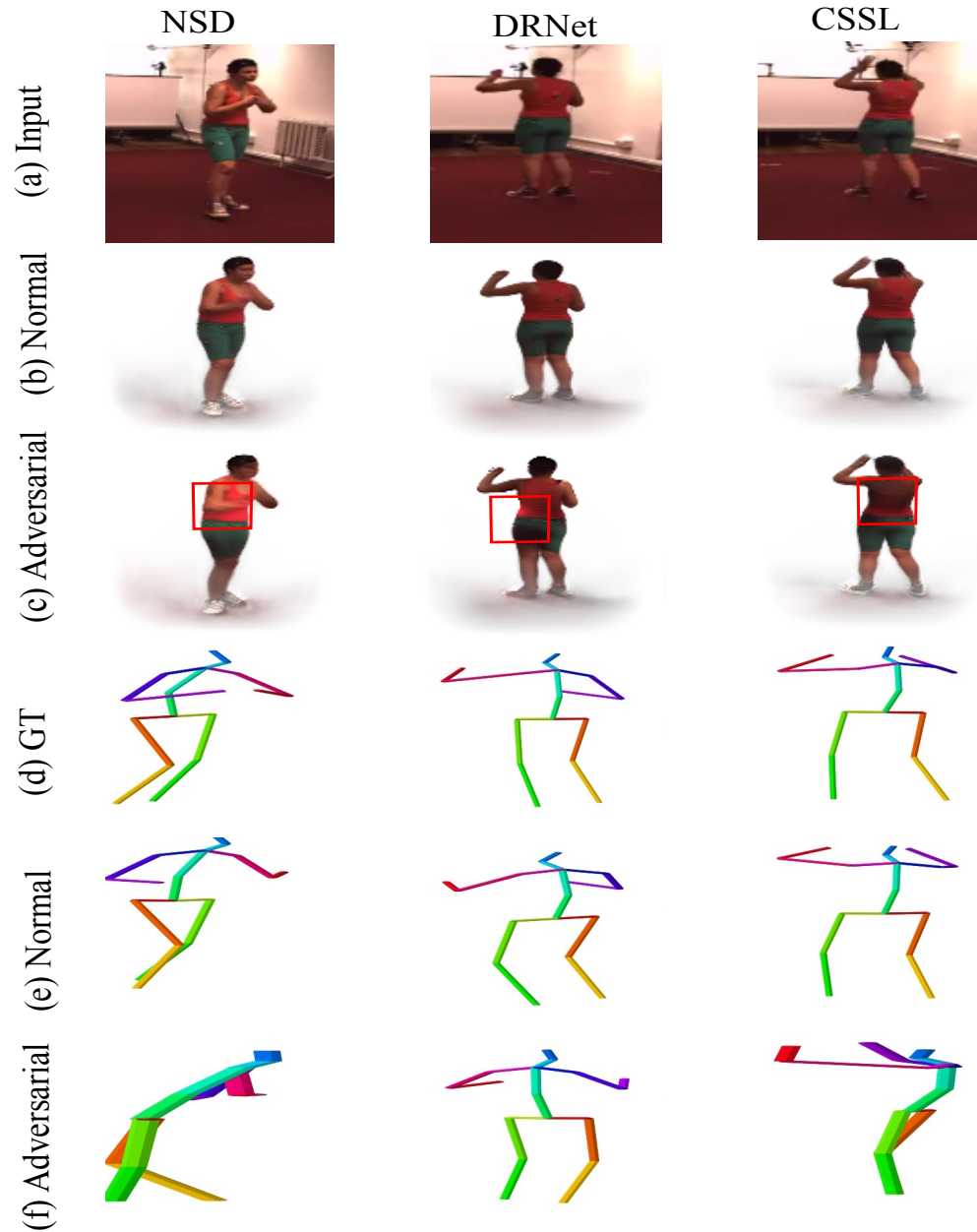
Same appearance code

# Semantic appearance attacks to understand the disentanglement of pose and appearance





# Qualitative results to show the disentanglement is incomplete

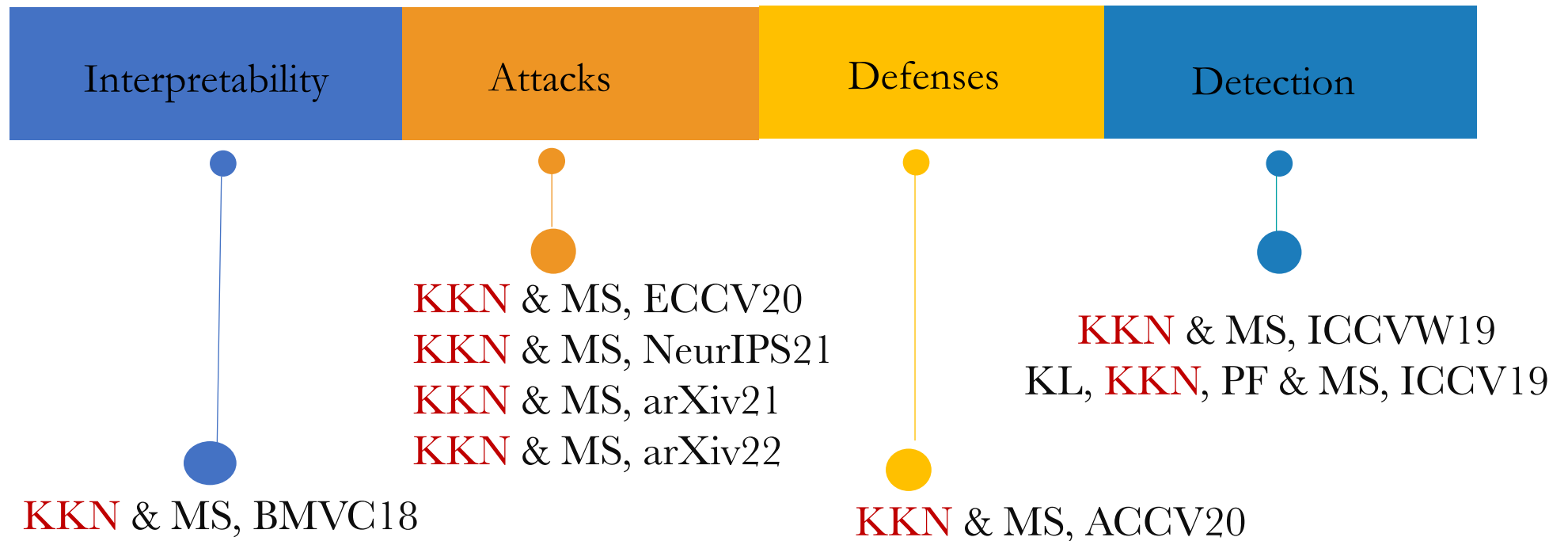


- A testbed to evaluate the disentanglement of pose and appearance
- Potential connection between disentanglement and robustness

# Conclusion

- Adversarial attacks have significant implication in the world of self-driving cars.
  - Results of indirect attack to fool far away dynamic objects are unsettling
- Black-box attacks are more realistic threat setting than white-box setting
  - Transferable perturbations in cross-model, cross-domain and cross-task setting
- Interpretable models to reveal working mechanism of adversarial attacks and to improve robustness
  - BoW networks for adversarial attack detection
  - Attention-based BoW networks with metric learning for defending to attacks

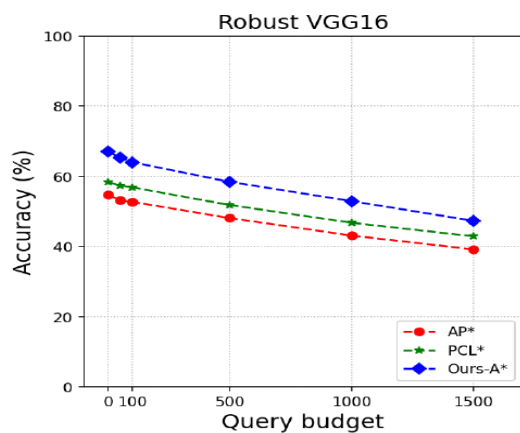
# Questions?



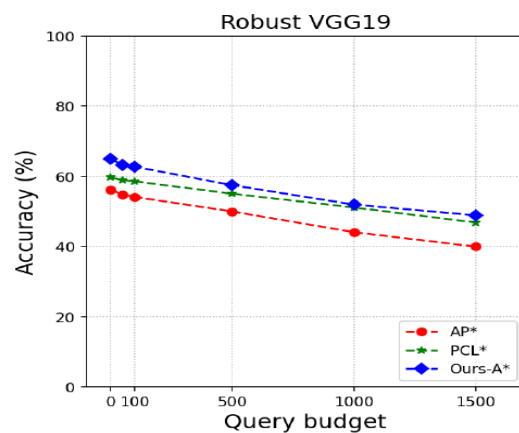
# Black-box Square attacks on adversarial trained models



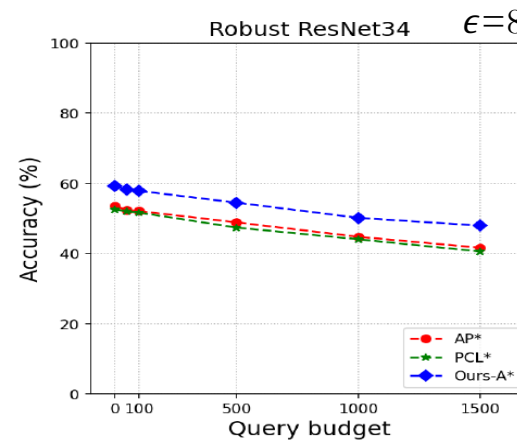
Higher is better



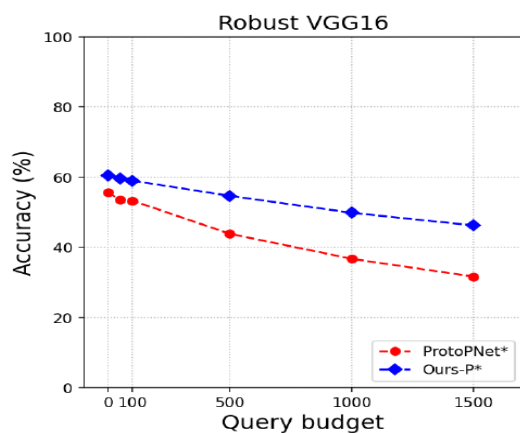
(a) VGG16



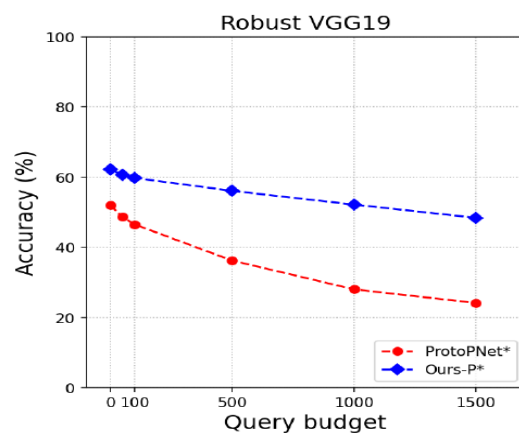
(b) VGG19



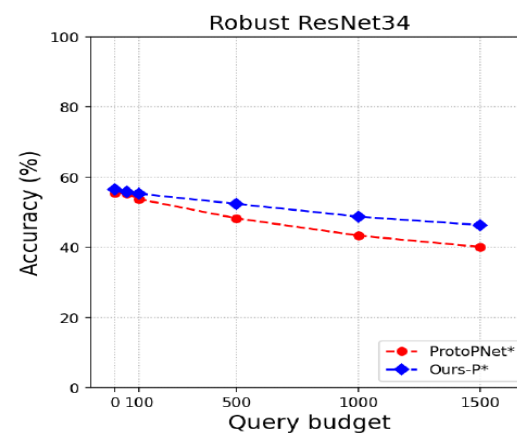
(c) ResNet34



(d) VGG16



(e) VGG19

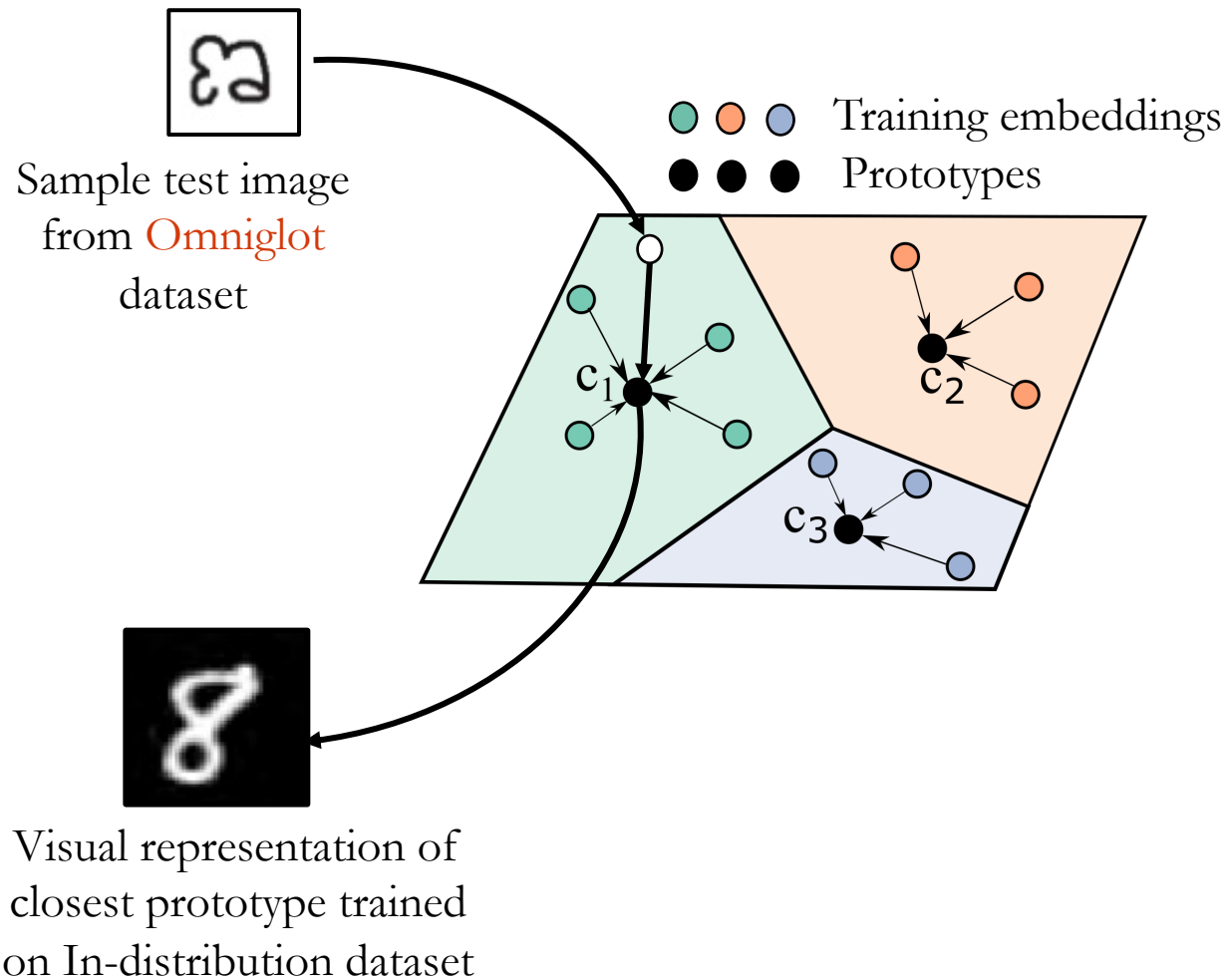


(f) ResNet34

Y-axis is **accuracy**, and x-axis is **query budget** for Square attack

## Similar intuition for OOD detection:

Out-of-distribution (OOD) input activates a different looking prototype



$D$ : Siamese network to predict if the input pair is similar or dissimilar trained on In-distribution dataset

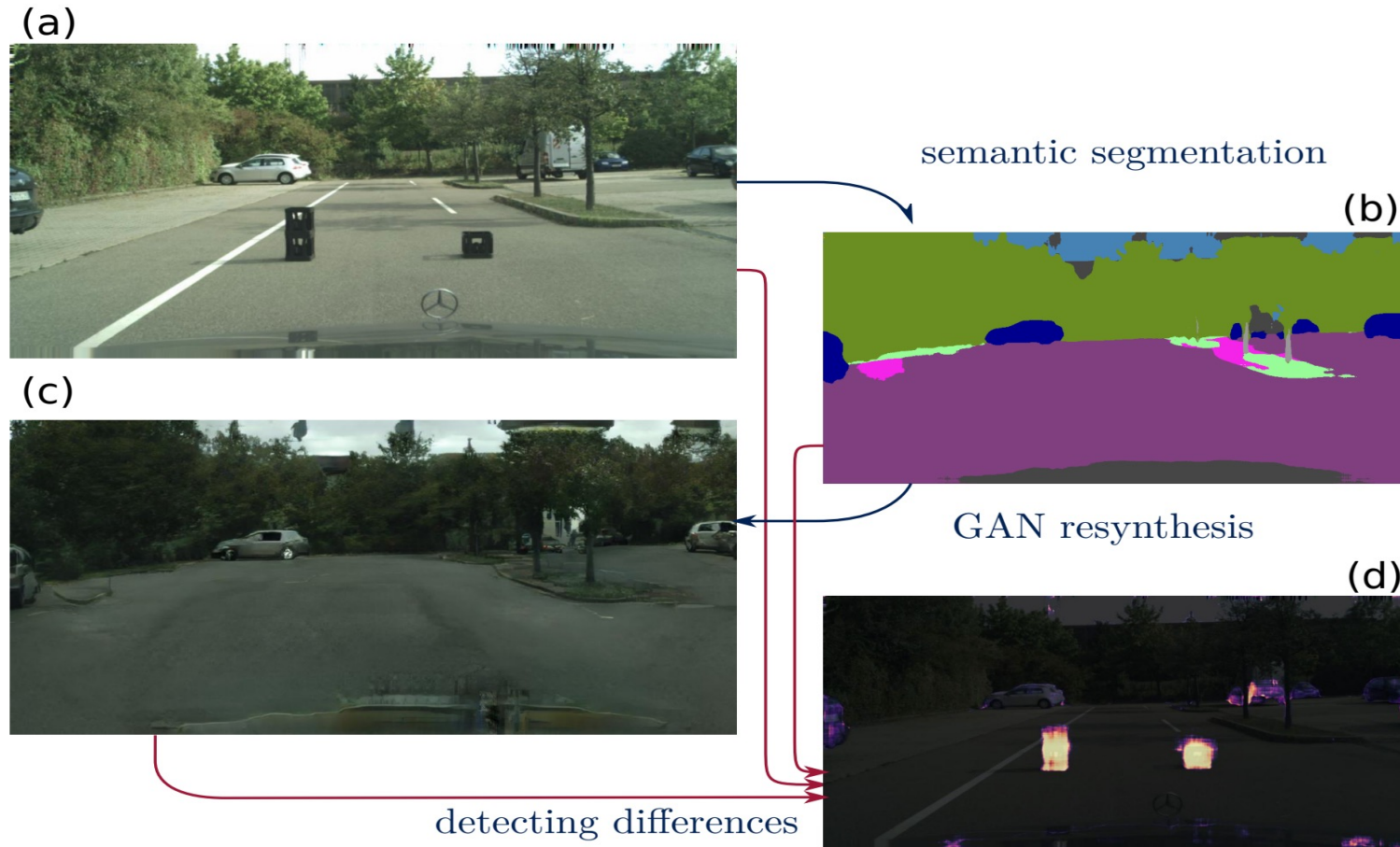
$$D\left(\begin{array}{c} \text{Image 1} \\ \text{Image 2} \end{array}\right)$$

Dissimilar

OOD sample

# Road Anomaly detection

Pixel-level detection of anomalous objects by comparing **input image** to the **image resynthesized** from output map

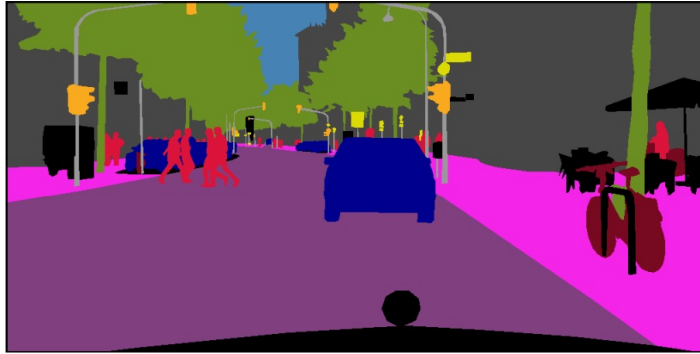


Pretrain the discrepancy detector network on real and synthesised images by randomly replacing objects of few classes with other classes

# Road Anomaly detection

Pixel-level detection of anomalous objects by comparing **input image** to the **image resynthesized** from output map

Real predictions

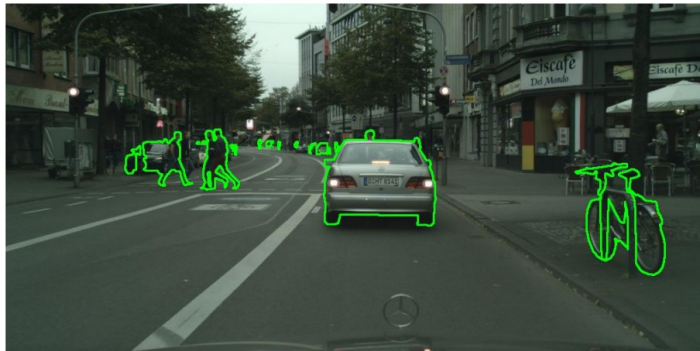


(a)

Randomly alter labels of few instances



(b)



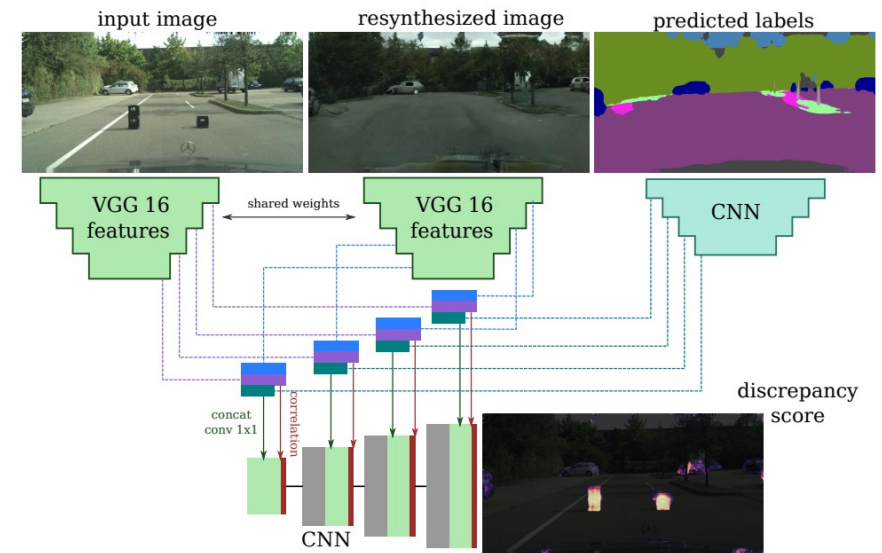
Outlines of altered objects



Resynthesized image

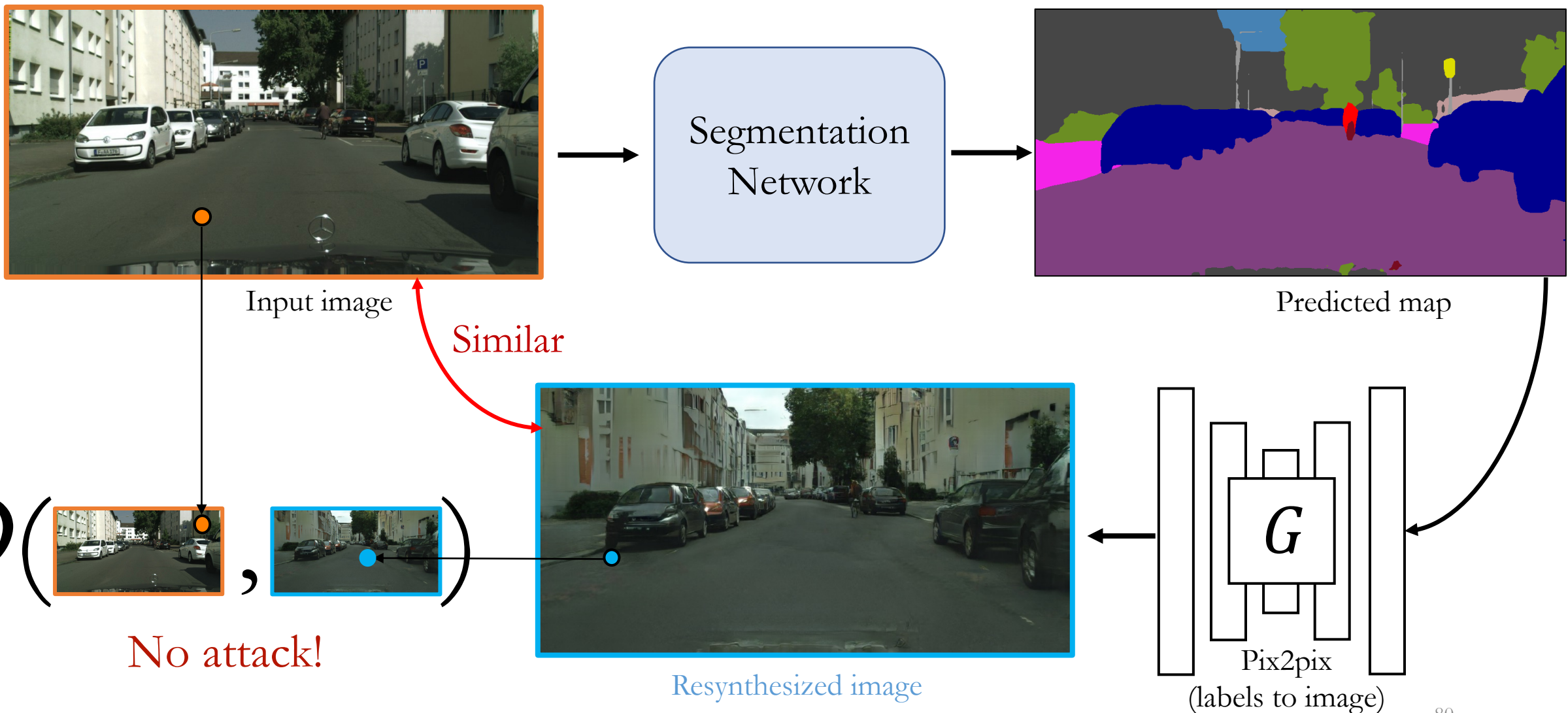
Pretrain the discrepancy detector network on real and synthesised images by randomly replacing objects of few classes with other classes

Discrepancy network



# Adversary detection beyond image-recognition

Adversarial example detection in semantic segmentation by comparing **input image** to the **image resynthesized** from output map

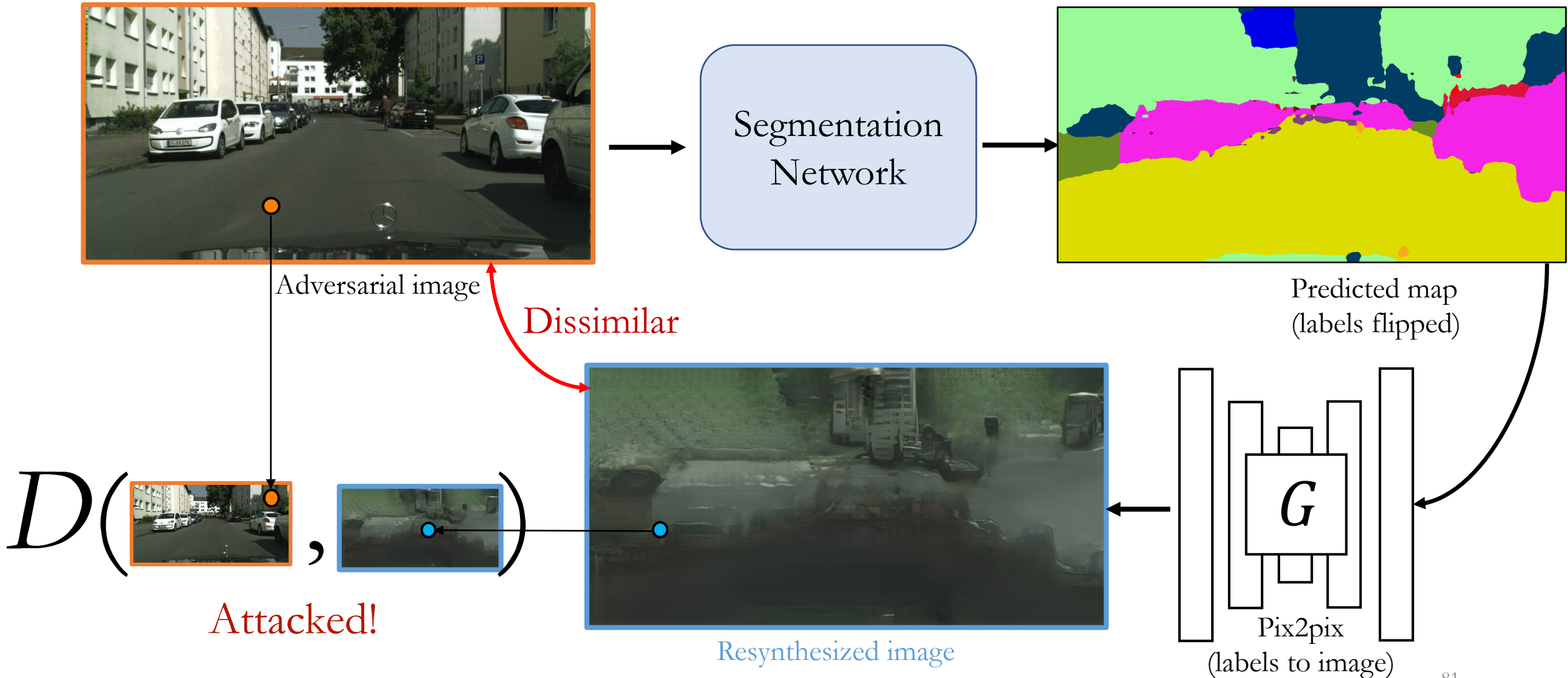


$D$ : Computes  $L_1$  distance in HOG feature space



# Adversary detection beyond image-recognition

Adversarial example detection in semantic segmentation by comparing **input image** to the **image resynthesized** from output map



$D$ : Computes  $L_1$  distance in HOG feature space

Normal synthesized

Normal predictions

Adversarial synthesized

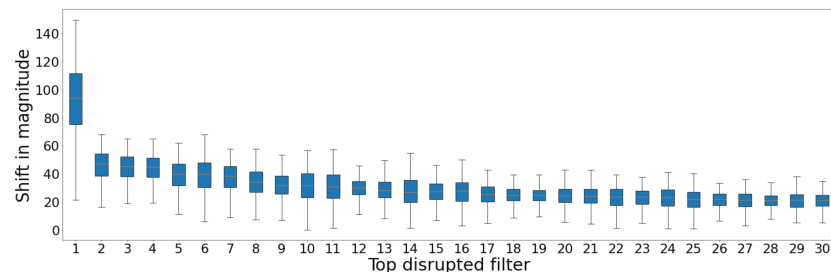
Adversarial predictions



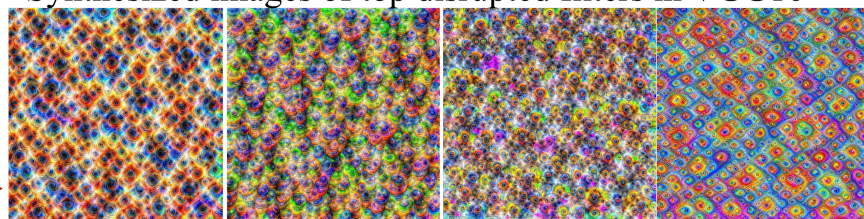
# Visual correlation b/w adversarial images & top disrupted filters

(a) Generator trained against VGG16

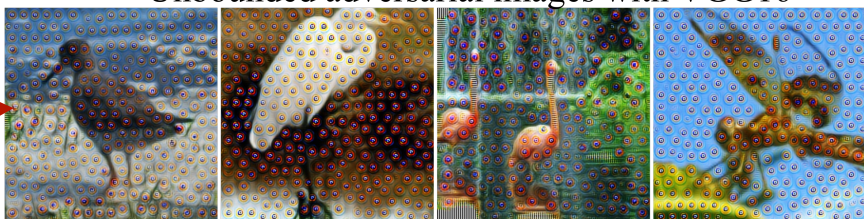
Boxplot of top 30 disrupted filters in layer 18 of VGG16



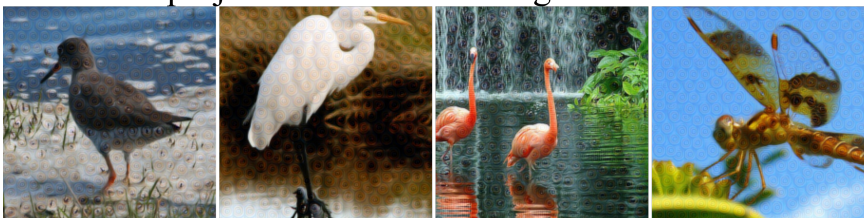
Synthesized images of top disrupted filters in VGG16



Unbounded adversarial images with VGG16



Final projected adversarial images with VGG16

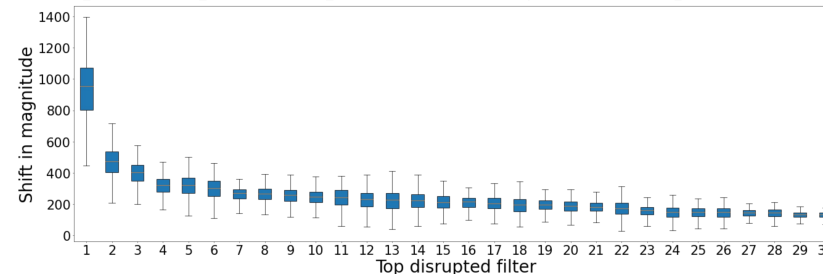


Visually correlated

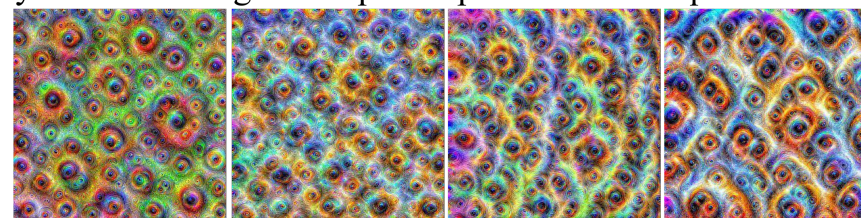


(b) Generator trained against SqueezeNet

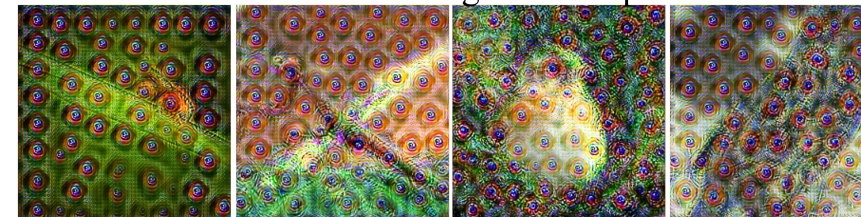
Boxplot of top 30 disrupted filters in layer 10 of SqueezeNet



Synthesized images of top disrupted filters in SqueezeNet



Unbounded adversarial images with SqueezeNet



Final projected adversarial images with SqueezeNet



Visually correlated

