



# Adversarial attacks in CNNs

---

Spring 2021: Seminar Explainable AI – Human–Computer  
Interaction meets Artificial Intelligence

Final Presentation – May 28, 2021

**VEDASRI NAKKA**

# **Topic : Adversarial Attacks Meets Interpretability**

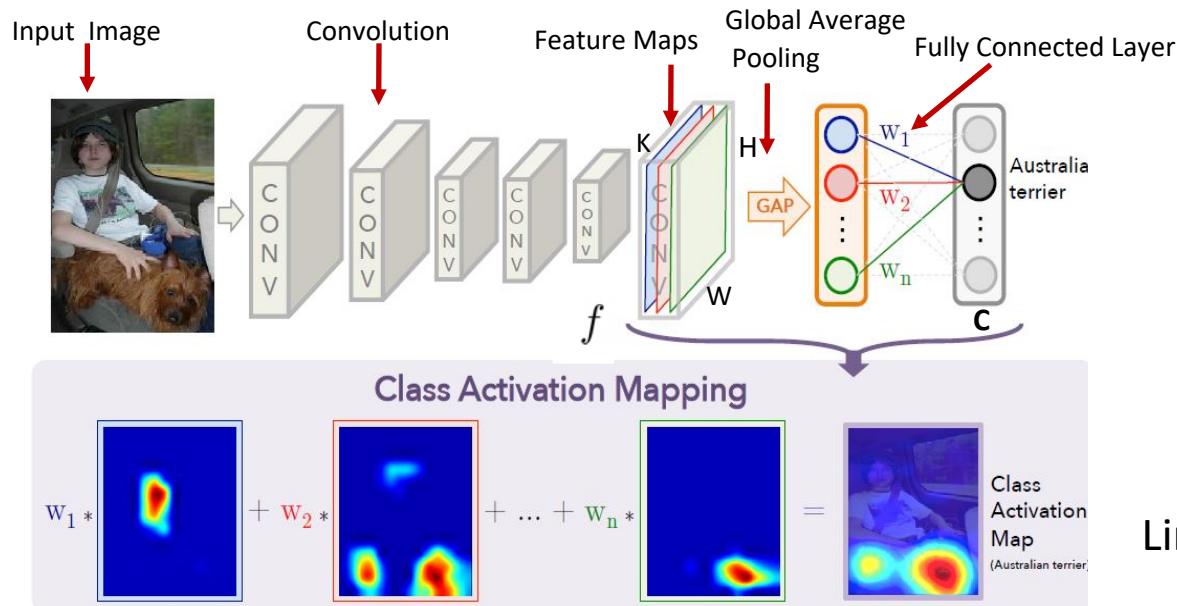
Analyze and explain the decisions of fine-grained recognition networks by studying the image regions responsible for classification for both clean and adversarial examples

# Interpretability in CNNs

1. Class Activation Maps
2. Attention Based Networks
3. Prototype Based Networks

# CAM – Deep Features for Discriminative Localization

Identify which image regions responsible for classification *after* the training



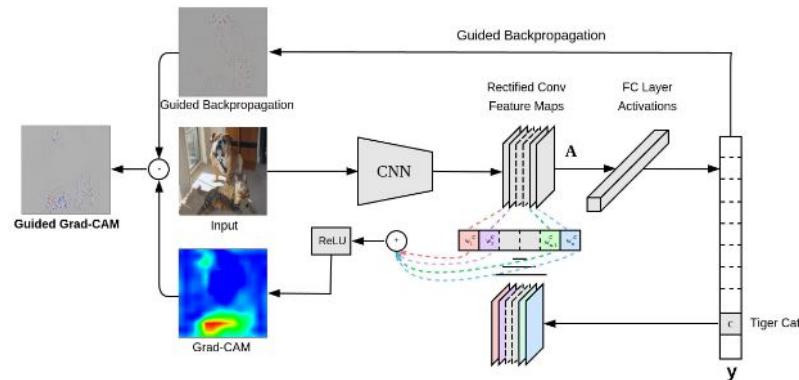
$M_c$  : Class Activation Map  
of class  $c$

$$M_c(x, y) = \sum_k w_k^c f_k(x, y)$$
$$f_k \in \mathbf{R}^{H \times W} \quad w_k \in \mathbf{R}^{K \times 1}$$

**Key Problem:**  
Limited to GAP-CNN architectures

# Grad-CAM: Gradient-based Localization

Computes the gradients instead of fixed weights of final layer



$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k}$$

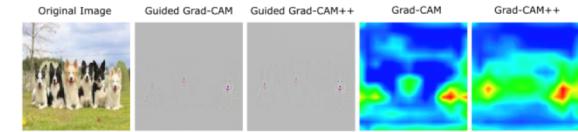
$$L_{Grad-CAM}^c = ReLU \left( \sum_k \alpha_k^c A_k^c \right)$$



Generalized to any CNN architecture

Grad CAM++

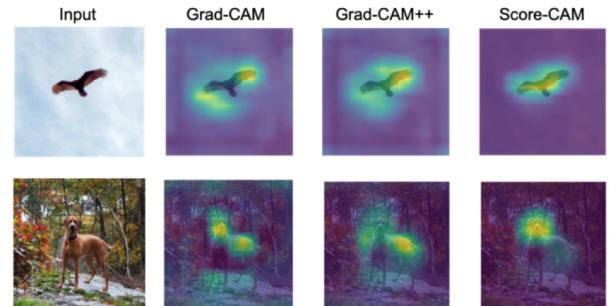
## Grad-CAM++: Improved Discriminative Localization [Aditya et al., CVPR 2017]



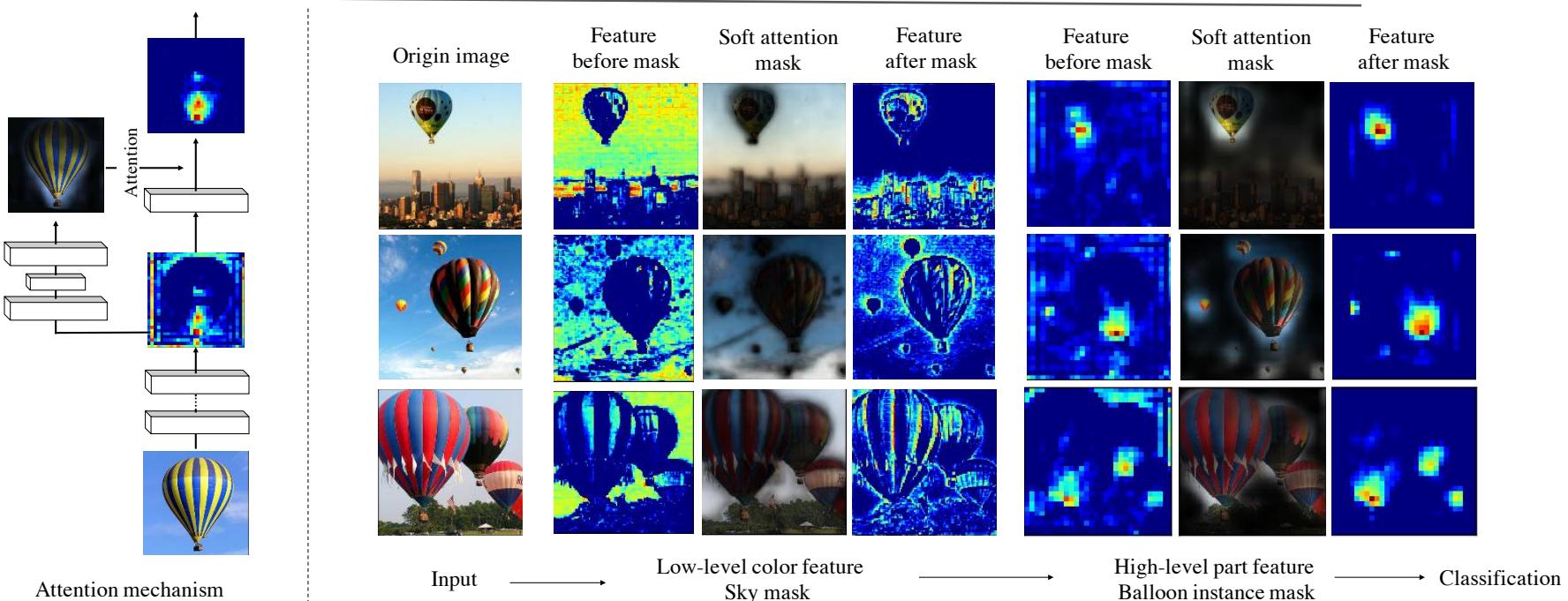
Score CAM

## Score-Weighted Visual Explanations

[Wang et al., CVPR 2020]

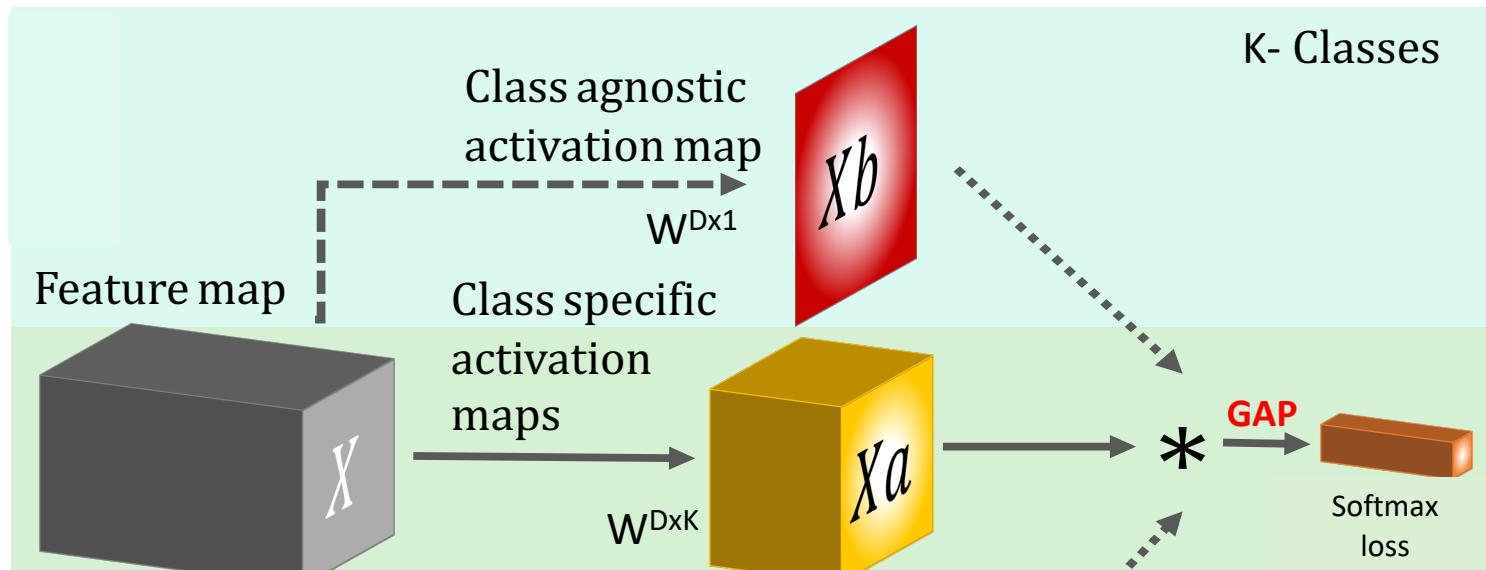


# Residual Attention Network



Attention mechanism incorporated with feed forward network architecture in an end-to-end training fashion

# Attention Pooling



CAMs = Feature maps  $\rightarrow$  GAP  $\rightarrow$  Class specific weights  $\rightarrow$  weighted feature map  
Attention Pooling = Feature maps  $\rightarrow$  class specific (& agnostic) weights  $\rightarrow$  GAP

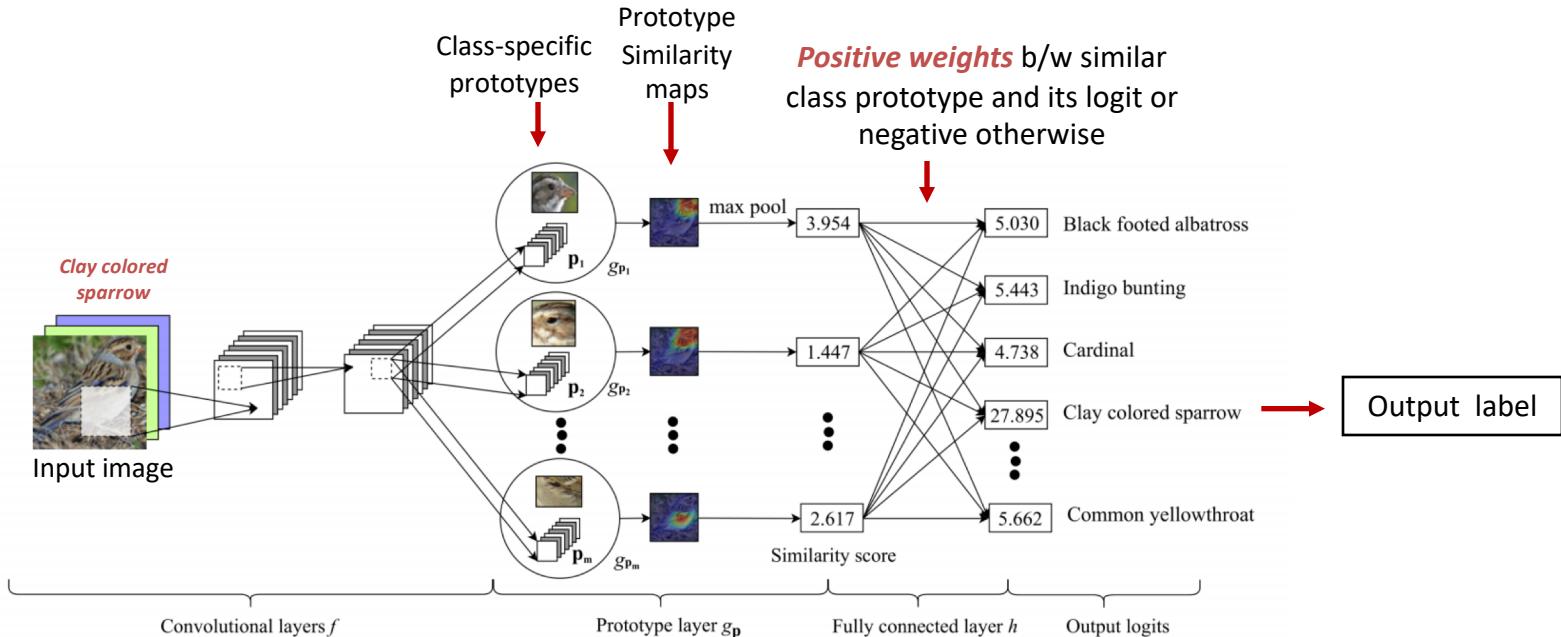
Both CAM and Attention-based models can only tell us which parts of the input they are looking at – they do not point us to prototypical cases to which the parts they focus on are similar

Solution?

Learn class-specific prototypes during training

# ProtoPNet: Interpretable Fine-Grained Network

Explain the decisions of CNN through class-specific prototypes in a feed-forward manner



# ProtoPNet: Interpretable Fine-Grained Network

## Objective Function:

$$\min_{\mathbf{P}, w_{\text{conv}}} \frac{1}{n} \sum_{i=1}^n \text{CrsEnt}(h \circ g_{\mathbf{P}} \circ f(\mathbf{x}_i), \mathbf{y}_i) + \lambda_1 \text{Clst} + \lambda_2 \text{Sep}$$

The diagram illustrates the architecture of the ProtoPNet objective function. It starts with input features  $\mathbf{x}_i$  which pass through 'Backbone layers' to a 'Final classification layer'. The output of this layer is then processed by a 'Prototype layer'. Simultaneously, the 'True label'  $\mathbf{y}_i$  is also input to the 'Prototype layer'. Red arrows indicate the flow of data from the input to the prototype layer and from the prototype layer to the output loss calculation.

CrsEnt - Cross entropy loss b/w predicted output and true label

Clst - Cluster loss (pushes the closest prototypes of its *own* class)

Sep - Separation loss (pulls away the closest prototypes of any *other* class)

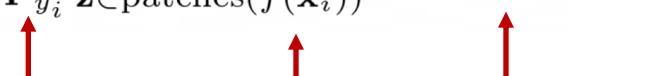
# ProtoPNet: Interpretable Fine-Grained Network

## Regularization Losses:

Brings the prototypes closer to its own class discriminative features

$$\text{Clst} = \frac{1}{n} \sum_{i=1}^n \min_{j: \mathbf{p}_j \in \mathbf{P}_{y_i}} \min_{\mathbf{z} \in \text{patches}(f(\mathbf{x}_i))} \|\mathbf{z} - \mathbf{p}_j\|_2^2$$

Separates the prototypes of other class from the discriminative features

$$\text{Sep} = -\frac{1}{n} \sum_{i=1}^n \min_{j: \mathbf{p}_j \notin \mathbf{P}_{y_i}} \min_{\mathbf{z} \in \text{patches}(f(\mathbf{x}_i))} \|\mathbf{z} - \mathbf{p}_j\|_2^2$$


## **Planned Items after Mid-Term**

---

1. Adversarial Attacks on different interpretable networks
2. Visualize the *shift* in discriminative regions under attacks
3. Understand the prototype similarity maps for adversarial samples
4. Compare CAMs vs. Prototype Similarity Maps vs. Attention Maps with fast adversarial attacks



# Post Mid Term

---

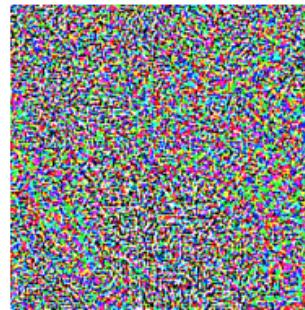
# Concept of Adversarial Attack

---



$x$   
“panda”  
57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$   
“nematode”  
8.2% confidence

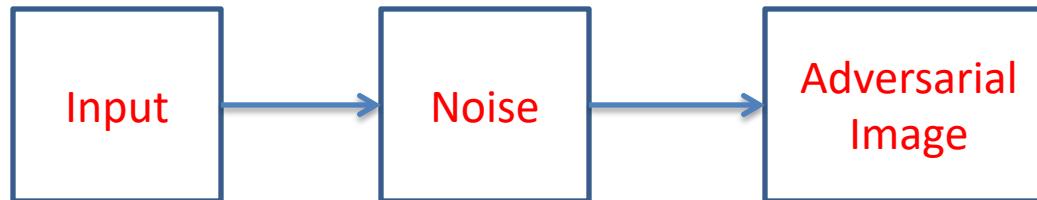
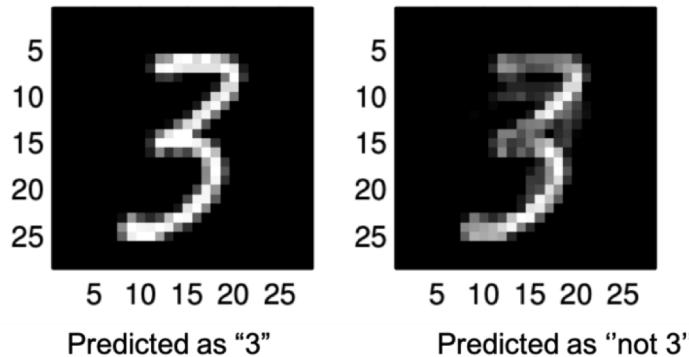
=



$x +$   
 $\epsilon \text{sign}(\nabla_x J(\theta, x, y))$   
“gibbon”  
99.3 % confidence

# Examples of Adversarial Attack

---



# **Adversarial attack Methods:**

---

4 types of Gradient-based Adversarial attacks:

1. Fast Gradient Sign Method(FGSM)
2. Basic Iterative Method (BIM)
3. Momentum Iterative Method (MIM)
4. Projected Gradient Descent (PGD)

# Adversarial attack Methods:

---

Fast Gradient Sign Method (FGSM):

$$\mathbf{x}_{adv} = \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y}))$$

The diagram illustrates the inputs and components of the FGSM equation. At the top, 'Perturbation strength' is shown with a red arrow pointing down to the term  $\epsilon$ . Below it, 'Gradient of Loss w.r.t Input' is shown with a red arrow pointing down to the term  $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y})$ . To the left of the equation, 'Adversarial Image' has a red arrow pointing up to the variable  $\mathbf{x}_{adv}$ . Below the equation, 'Input Image' has a red arrow pointing up to the variable  $\mathbf{x}$ . To the right of the equation, 'Cross Entropy Loss' has a red arrow pointing up to the term  $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y})$ , and 'True Label' has a red arrow pointing up to the variable  $\mathbf{y}$ .

# Adversarial attack Methods:

---

Basic Iterative Method (BIM):

$$\mathbf{x}_m = \text{clip}_{\epsilon}(\mathbf{x}_{m-1} + \frac{\epsilon}{i} \cdot \text{sign}(\nabla_{x_{m-1}}(\mathcal{L}(\mathbf{x}_{m-1}, \mathbf{y}))))$$

The diagram illustrates the components of the BIM update rule with red arrows:

- An arrow points from "Adversarial Image" to the term  $\mathbf{x}_{m-1}$ .
- An arrow points from "With in band of  $\epsilon$  in  $L_\infty$  norm" to the term  $\frac{\epsilon}{i}$ .
- An arrow points from "Adversarial Image at previous step" to the term  $\mathbf{x}_{m-1}$ .
- An arrow points from "Total number of iterations" to the term  $i$ .
- An arrow points from "m: current iteration" to the term  $\mathcal{L}(\mathbf{x}_{m-1}, \mathbf{y})$ .

## Adversarial attack Methods:

## Momentum Iterative Method (MIM):

$$\mathbf{x}_m = \text{clip}_{\epsilon}(\mathbf{x}_{m-1} + \frac{\epsilon}{i} \cdot \text{sign}(g_m))$$

# Adversarial attack Methods:

---

Projected Gradient Descent (PGD):

$$\mathbf{x}_m = \mathbf{x}_{m-1} + \gamma \cdot \text{sign}(\nabla_{x_{m-1}} \mathcal{L}(\mathbf{x}_{m-1}, \mathbf{y}))$$
$$\mathbf{x}_m = \text{clip}(\mathbf{x}_m, \mathbf{x}_m - \epsilon, \mathbf{x}_m + \epsilon)$$

Step size  
↓  
Perturbation strength  
↓  
Adversarial Image ↑  
Adversarial Image at previous step ↑

- Starts with a random noise.
- Step size is larger when compared to BIM

# **Experimental Setup**

---

1. Training Dataset : 300 Images
2. Testing Dataset: 299 Images
3. Number of Classes: 12 from CUB dataset
4. Networks: VGG16/ Attention Pooling / ProtoPNet

# Experimental Results:

---

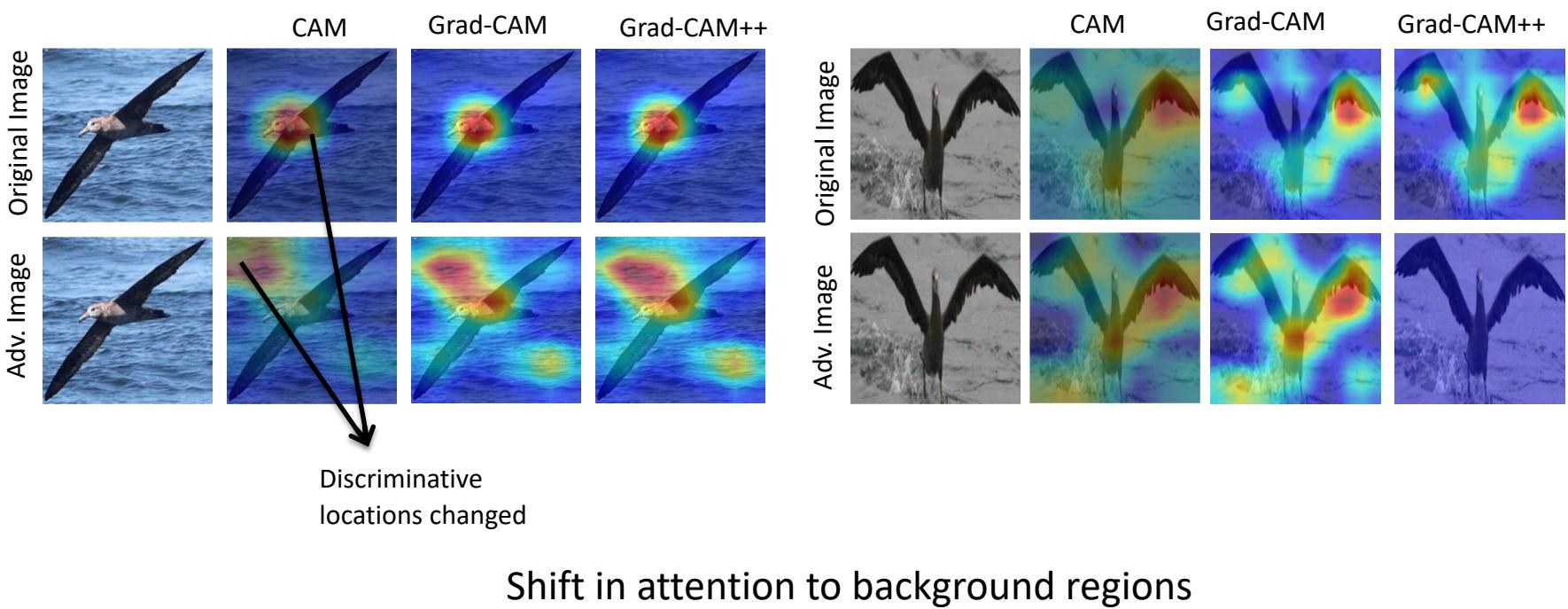
Accuracy of different networks to attacks

Network	Clean (0,0)	FGSM (2,1)	FGSM (8,1)	BIM (2,10)	BIM (8,10)	PGD (2,10)	PGD (8,10)	MIM (2,10)	MIM (8,10)
VGG16	<b>88.96%</b>	39.46%	21.07%	10.37%	3.01%	5.35%	2.01%	13.38%	3.01%
Attention Pooling	84.62%	42.41%	33.11%	<b>20.74%</b>	<b>16.72%</b>	<b>17.39%</b>	<b>14.72%</b>	<b>22.74%</b>	<b>17.39%</b>
ProtoPNet	76.92%	<b>44.15%</b>	<b>37.79%</b>	4.35%	3.34%	3.68%	0.33%	2.68%	1.67%

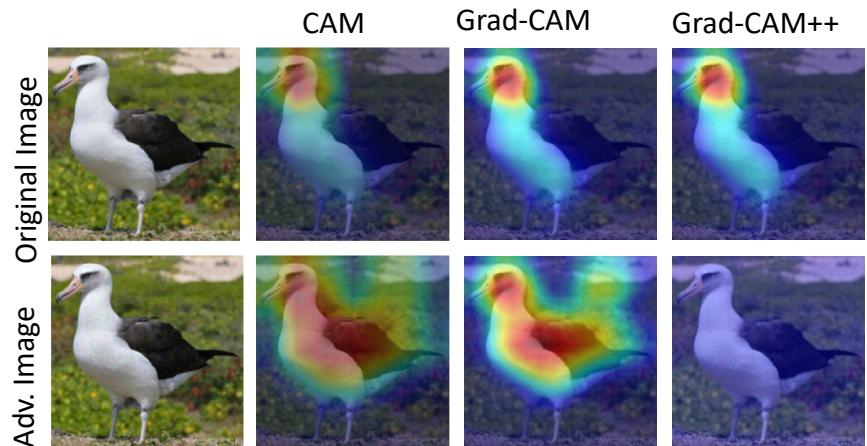
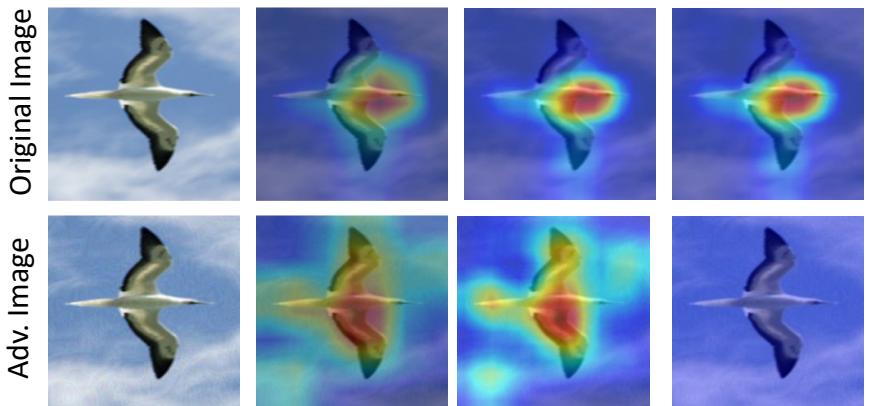
In each column, for example, the notation (8,1) indicates  $\epsilon=8$  in input range of [0,255] and number of iterations set to 1

We observe Attentional Pooling architecture is  
more robust to adversarial attacks

# Qualitative Results: Adversarial Attacks on VGG16

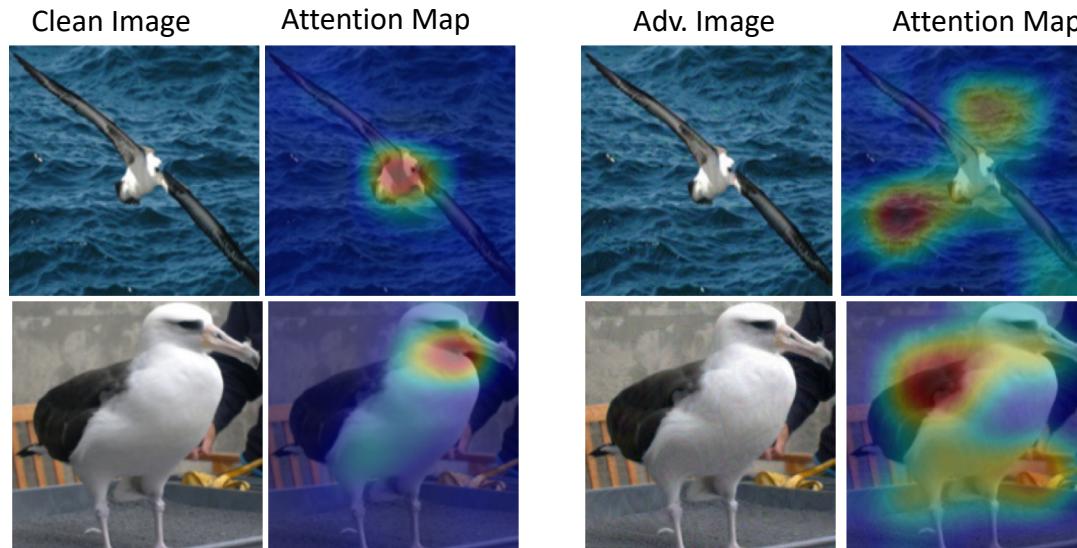


# Qualitative Results: Adversarial Attacks on VGG16



Shift in attention to large regions.

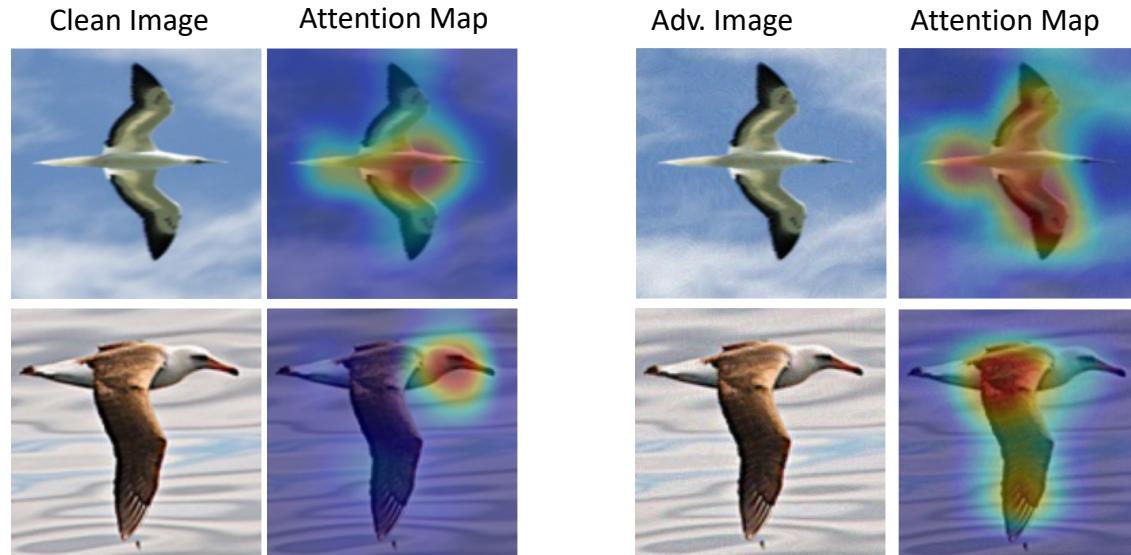
# Qualitative Results: Attacks on Attentional Pooling



Shift in attention to background regions

## Qualitative Results: Attacks on Attentional Pooling

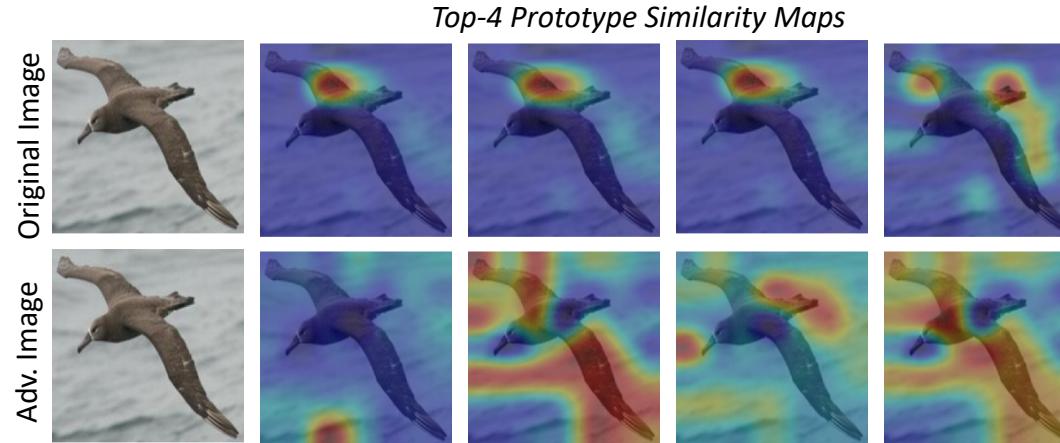
---



Shift in attention to background regions

# Qualitative Results: Attacks on ProtoPNet

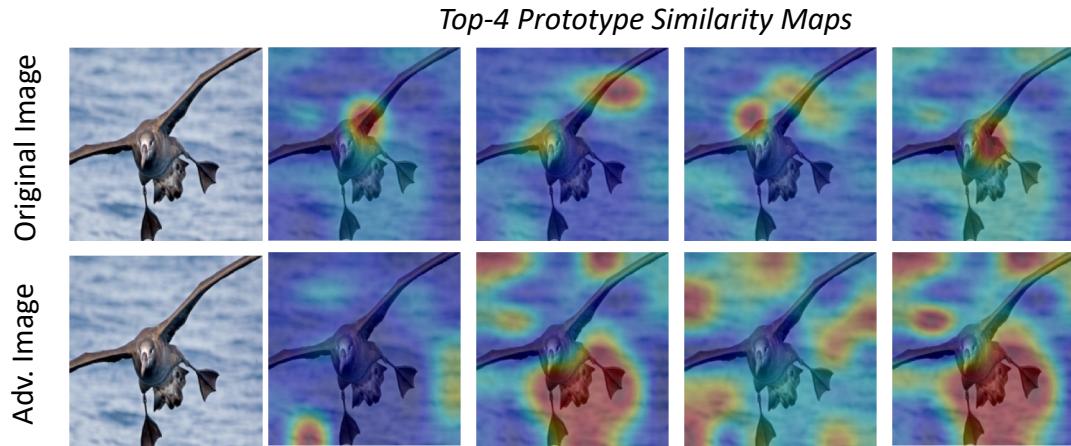
---



Adversarial image activates background prototypes

# Qualitative Results: Attacks on ProtoPNet

---

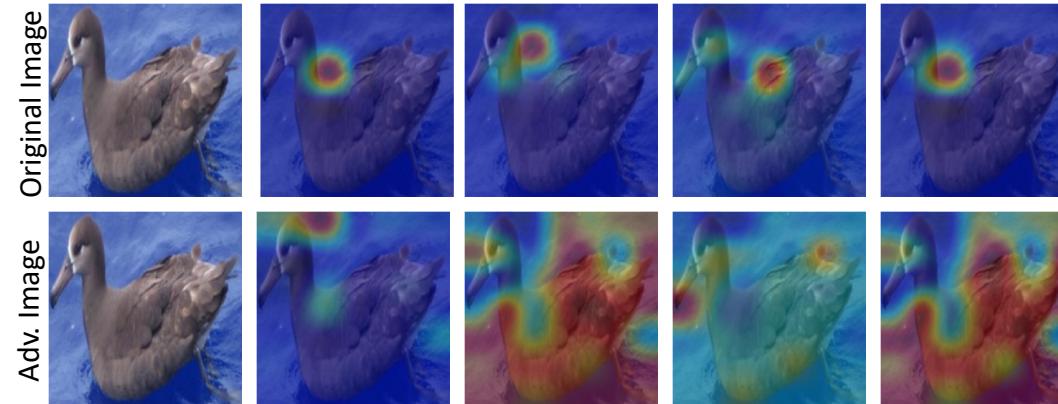


Adversarial image activates background prototypes

# Qualitative Results: Attacks on ProtoPNet

---

*Top-4 Prototype Similarity Maps*



Adversarial image activates background prototypes

# Qualitative Results: Attacks on ProtoPNet

---

*Top-4 Prototype Similarity Maps*



Adversarial image activates background prototypes

## **Conclusions**

---

1. Attentional Pooling framework is robust to attacks
2. The discrimination region moved to a background region or covering large region leading to inaccurate predictions
3. Interpreting the regions for CNN decisions can be a novel direction to detect adversarial attacks
4. ProtoPNet is most vulnerable to adversarial attacks although it is useful in analyzing the CNN decisions

# Thank you