



Robust Discriminative Localization Maps

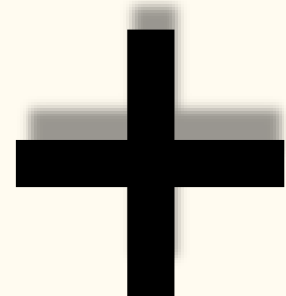
A. Prakash, N. Moran, S. Garber, A. DiLillo & J. Storer

✉ aprakash@brandeis.edu, iamaaditya.github.io

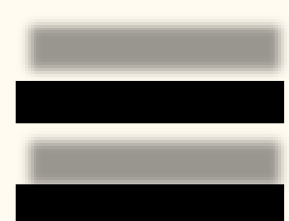
🔗 github.com/iamaaditya/robust-activation-maps



Class Activation Maps



Adversarial Systems



Bad ! Bad !! FAKE Activation Maps !!

Visual Question Answering
[Zhou, 2015]

Image Classification
[Zhou, 2016]

Image Compression
[Prakash, 2017]

Image Retrieval
[Jimenez, 2017]

Guided Attention
[Li, 2018]

Adversarial Defense
[Prakash, 2018]

CAM(Clean)

Image



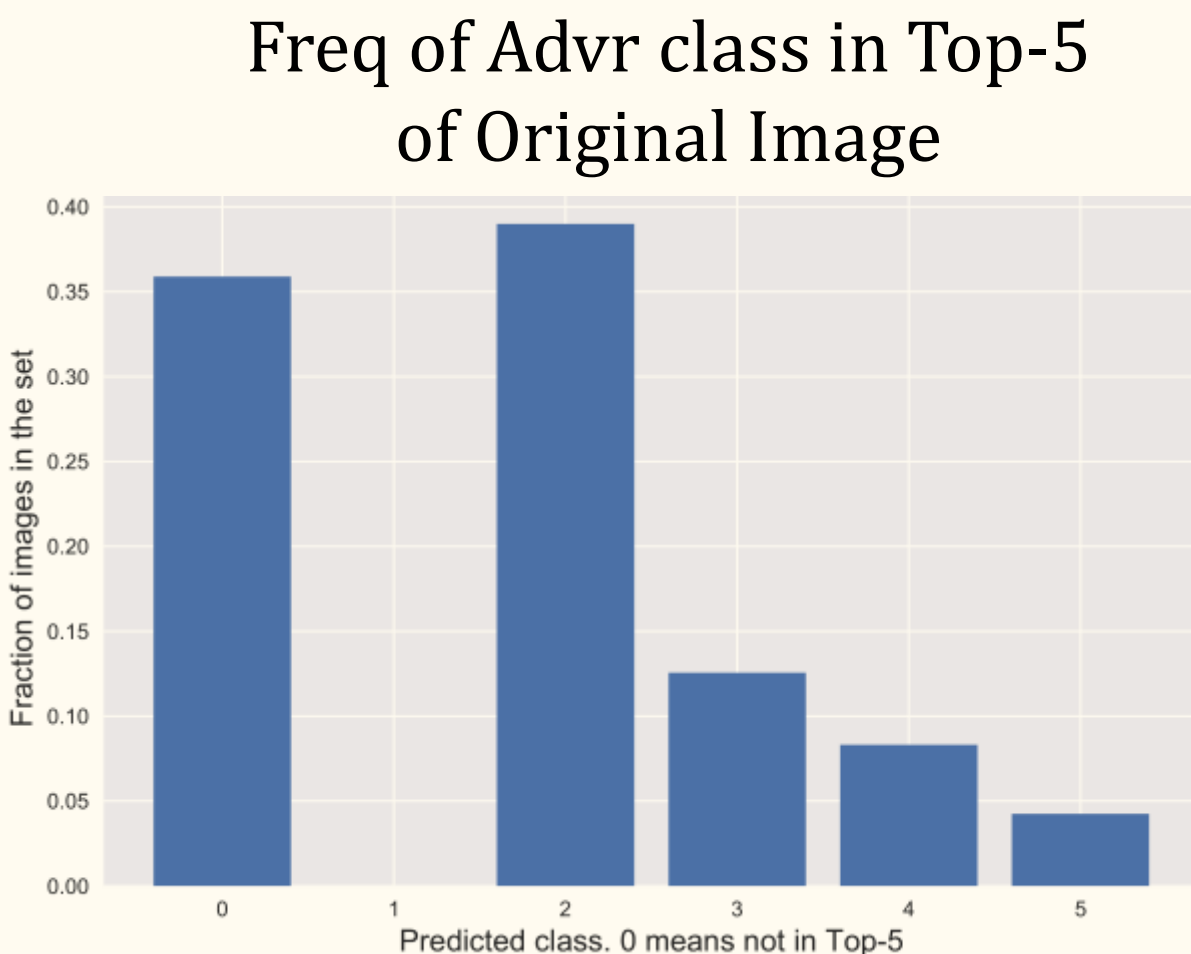
CAM(Advr)

Meatloaf [99%]



RDM(Advr)

Meatloaf [99%]



Class Activation Map

$$M_c(x, y) = \sum_k w_c^k f_k(x, y)$$

Robust Discriminative Map

$$\hat{M}(x, y) = \sum_c \frac{M_c(x, y)}{2^i}$$

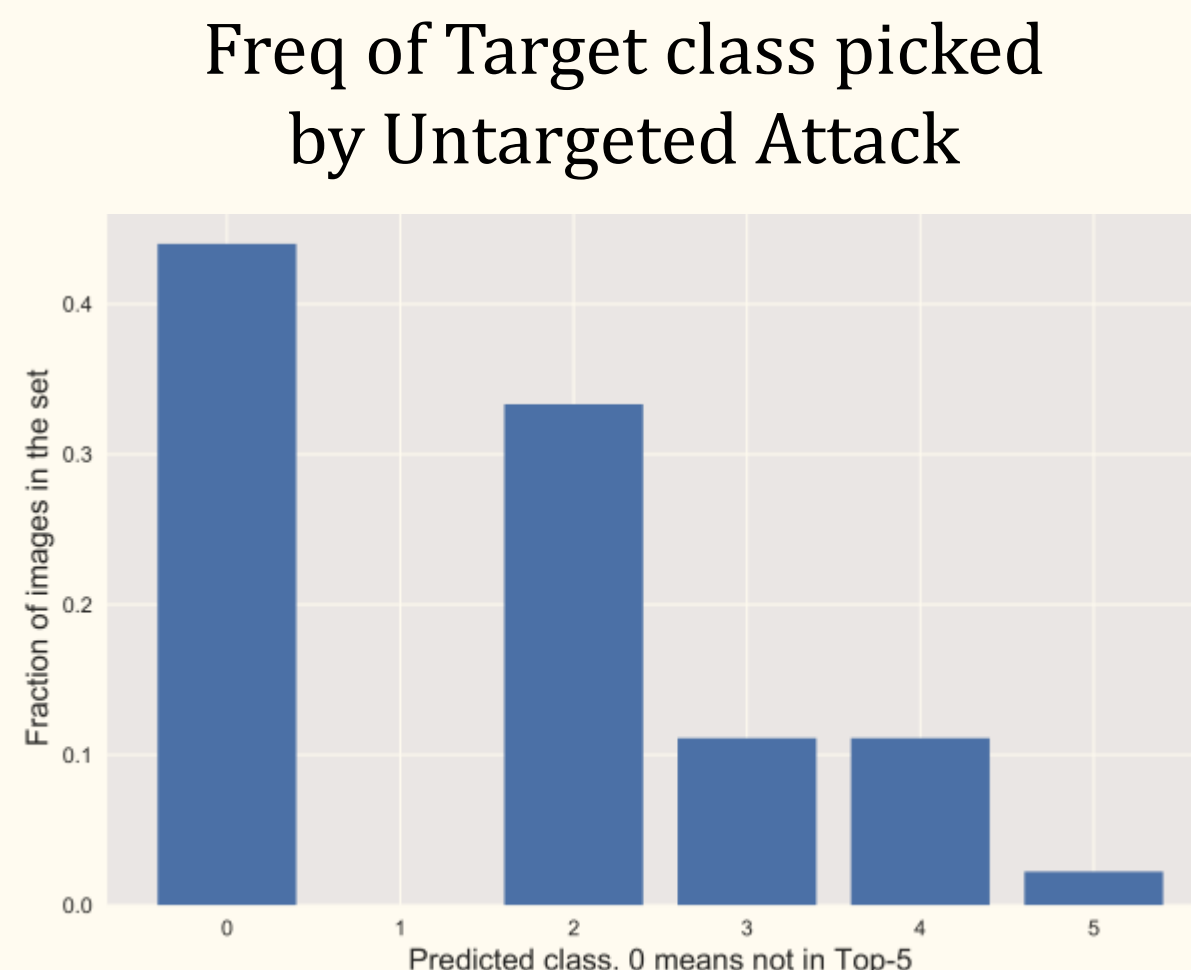
Warplane [91%]



Flatworm [99%]



Flatworm [99%]



Defense	FGSM	IGSM	DFool	C&W
Feature Squeezing (Xu et al [49])				
(a) Bit Depth (2 bit)	0.132	0.511	0.286	0.170
(b) Bit Depth (5 bit)	0.057	0.022	0.310	0.957
(c) Median Smoothing (2x2)	0.358	0.422	0.714	0.894
Quilting + TVM (Guo et al. [19])				
Quilting	0.611	0.862	0.858	0.843
TVM + Quilting	0.619	0.866	0.866	0.841
Cropping + TVM + Quilting	0.629	0.882	0.883	0.859
Our work: PD - Pixel Deflection, R-CAM: Robust CAM				
PD	0.735	0.880	0.914	0.931
PD + R-CAM	0.746	0.912	0.911	0.952

