

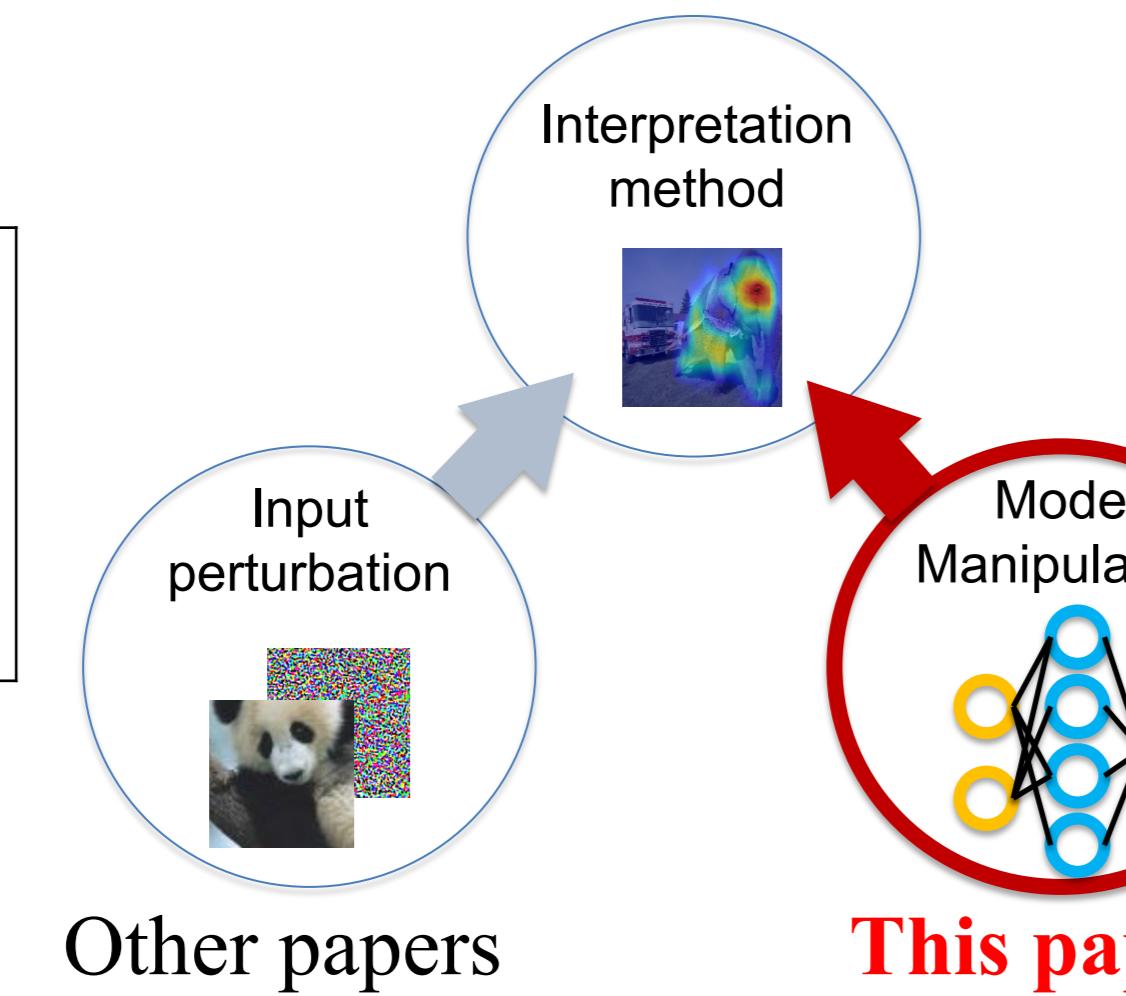
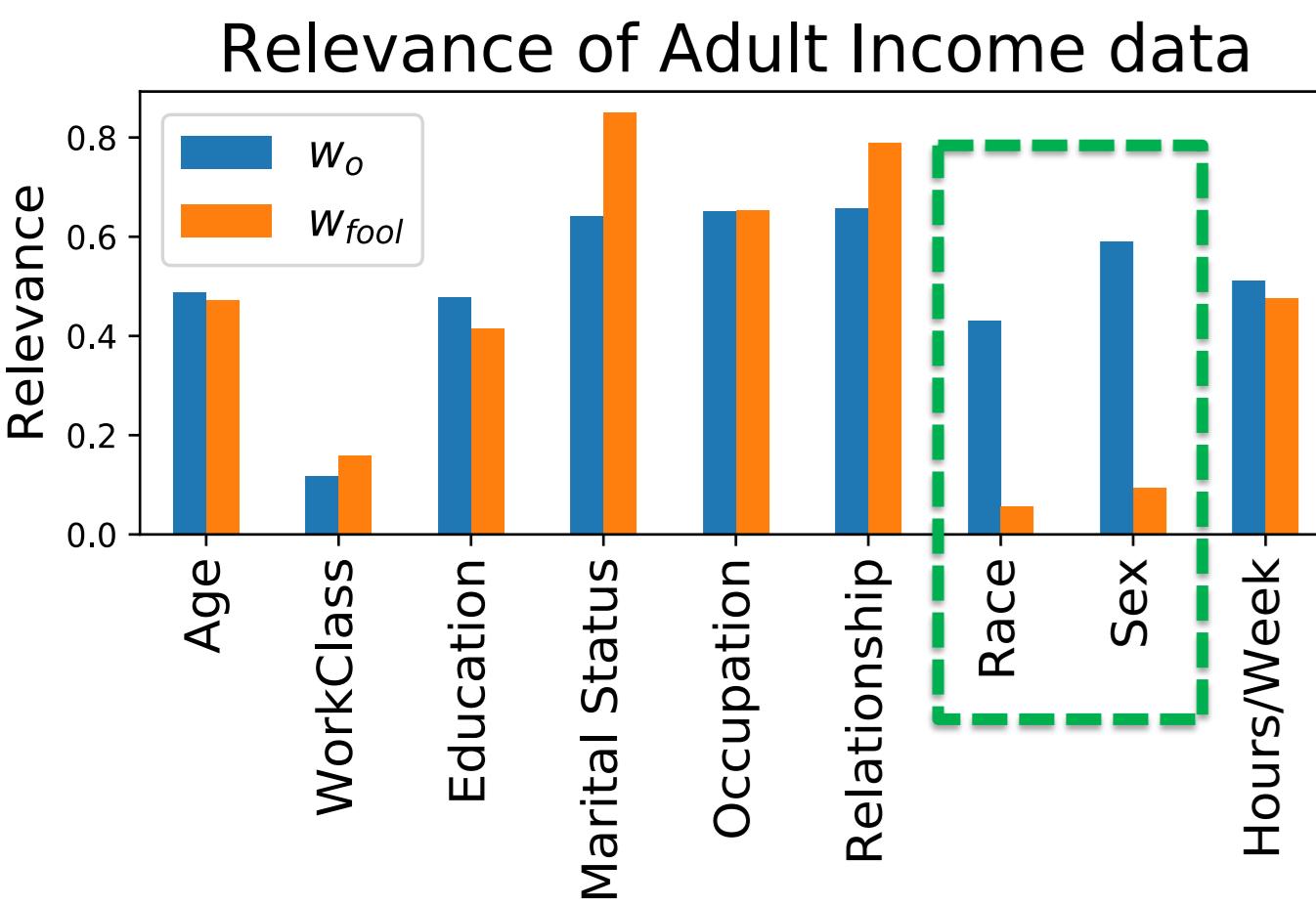
Fooling Neural Network Interpretations via Adversarial Model Manipulation

Juyeon Heo^{*1}, Sunghwan Joo^{*1}, and Taesup Moon^{1,2}

Department of Electrical and Computer Engineering¹, Department of Artificial Intelligence², Sungkyunkwan University, Korea
heojuyeon12@gmail.com, {shjoo840, tsmoon}@skku.edu

Motivation

"What if lazy developer **manipulate** the model to **hide the unfair biases** on interpretations rather than actually fixing the model to remove them."



- First to consider the notion of **stability of interpretation methods** with respect to our fooling model manipulation.
- Fooled explanation generalizes to the **entire validation set**

Method

$$\begin{aligned} \mathcal{L}(\mathcal{D}, \mathcal{D}_{fool}, \mathcal{I}; \mathbf{w}, \mathbf{w}_0) \\ = \mathcal{L}_C(\mathcal{D}; \mathbf{w}) + \lambda \mathcal{L}_F^{\mathcal{I}}(\mathcal{D}_{fool}; \mathbf{w}, \mathbf{w}_0) \end{aligned}$$

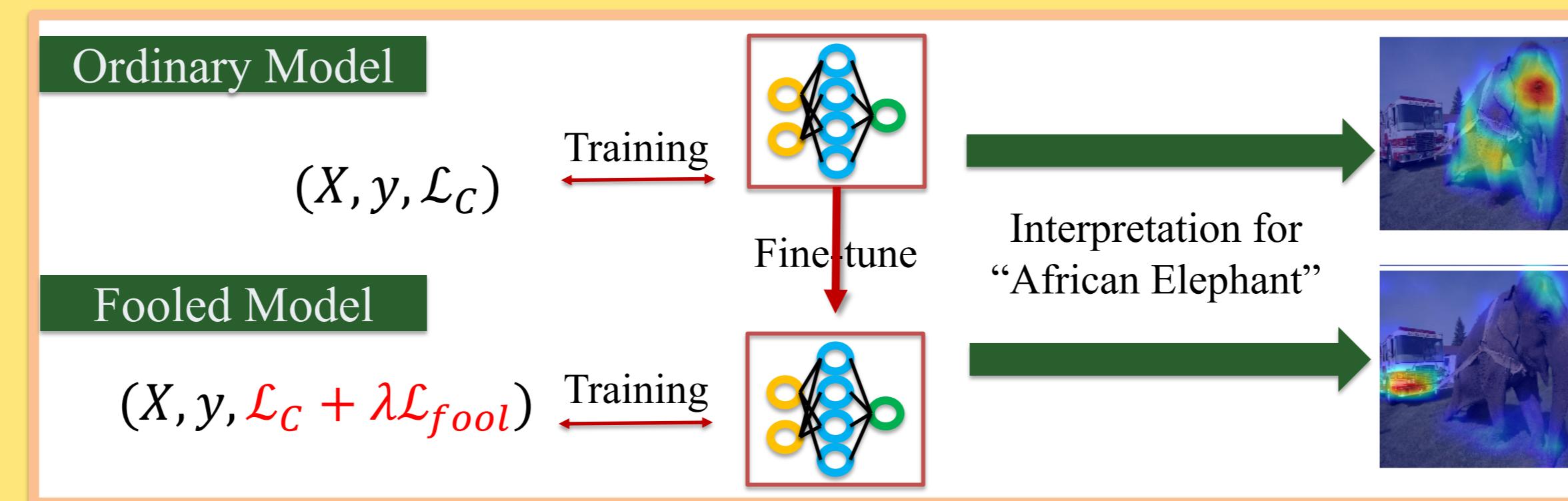
\mathcal{I} : Interpretation method
 \mathcal{D} : Data
 \mathbf{w}_0 : Pre-trained model
 \mathbf{w} : Fine-tuned model
 $h_c^{\mathcal{I}}$: Saliency map

• Location Fooling:	$\frac{1}{n} \sum_{i=1}^n \frac{1}{d_I} \ \mathbf{h}_{y_i}^{\mathcal{I}}(\mathbf{w}) - \mathbf{m} \ _2^2$	$\ \mathbf{h}_{y_i}^{\mathcal{I}}(\mathbf{w}) - \mathbf{m} \ _2^2$
• Top-k Fooling:	$\frac{1}{n} \sum_{i=1}^n \sum_{j \in \mathcal{P}_{i,k}(\mathbf{w}_0)} h_{y_i,j}^{\mathcal{I}}(\mathbf{w}) $	$ \mathbf{h}_{y_i}^{\mathcal{I}}(\mathbf{w}) $
• Center-mass Fooling:	$-\frac{1}{n} \sum_{i=1}^n \ C(\mathbf{h}_{y_i}^{\mathcal{I}}(\mathbf{w})) - C(\mathbf{h}_{y_i}^{\mathcal{I}}(\mathbf{w}_0)) \ _1$	$-\ \mathbf{h}_{y_i}^{\mathcal{I}}(\mathbf{w}) - \mathbf{h}_{y_i}^{\mathcal{I}}(\mathbf{w}_0) \ _1$
• Active Fooling:	$\frac{1}{2n_{fool}} \sum_{i=1}^{n_{fool}} \frac{1}{d_I} (\ \mathbf{h}_{c_1}^{\mathcal{I}}(\mathbf{w}) - \mathbf{h}_{c_2}^{\mathcal{I}}(\mathbf{w}_0) \ _2^2 + \ \mathbf{h}_{c_1}^{\mathcal{I}}(\mathbf{w}_0) - \mathbf{h}_{c_2}^{\mathcal{I}}(\mathbf{w}) \ _2^2)$	$\ \mathbf{h}_{c_1}^{\mathcal{I}}(\mathbf{w}) - \mathbf{h}_{c_2}^{\mathcal{I}}(\mathbf{w}_0) \ _2^2$

Main Idea

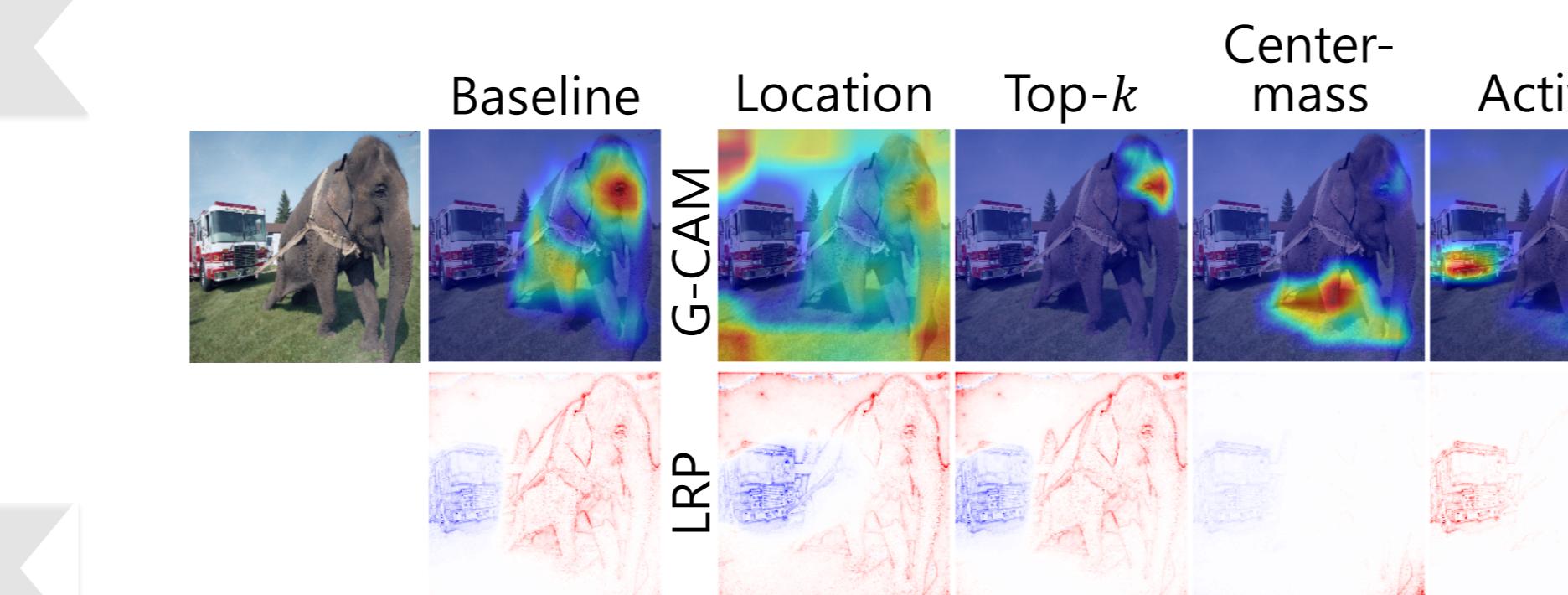
Vulnerability of saliency map based interpreters:

"We found that the **saliency based interpretation methods** can be **fooled via model manipulation** without significant drops on **accuracy**."

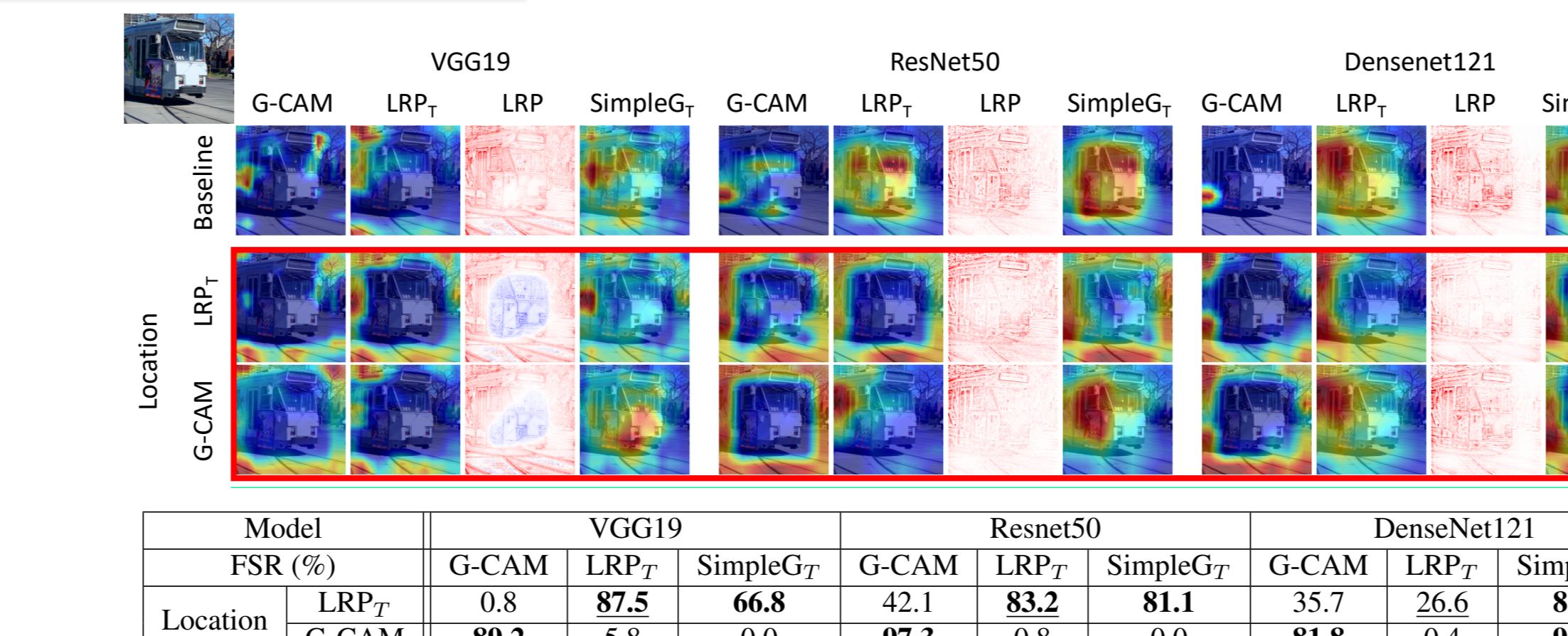


Results and Discussion

Main Result

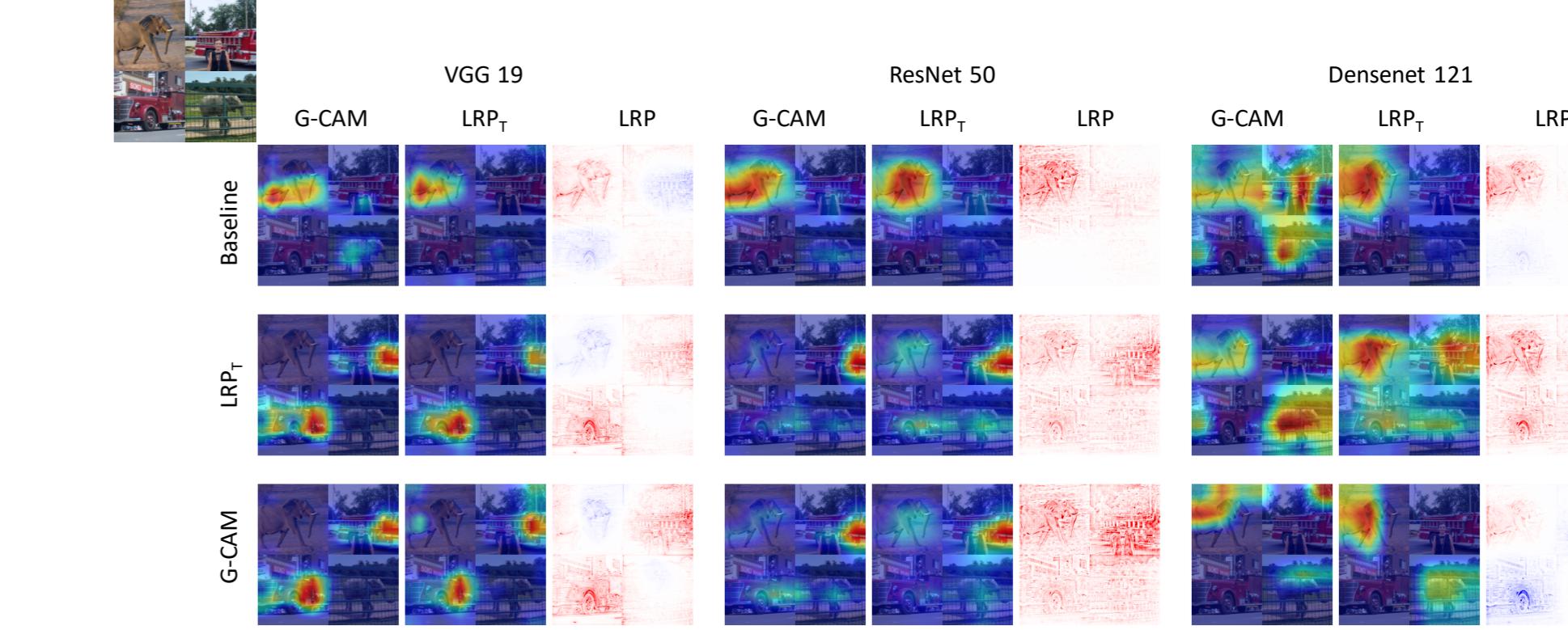


Passive Fooling



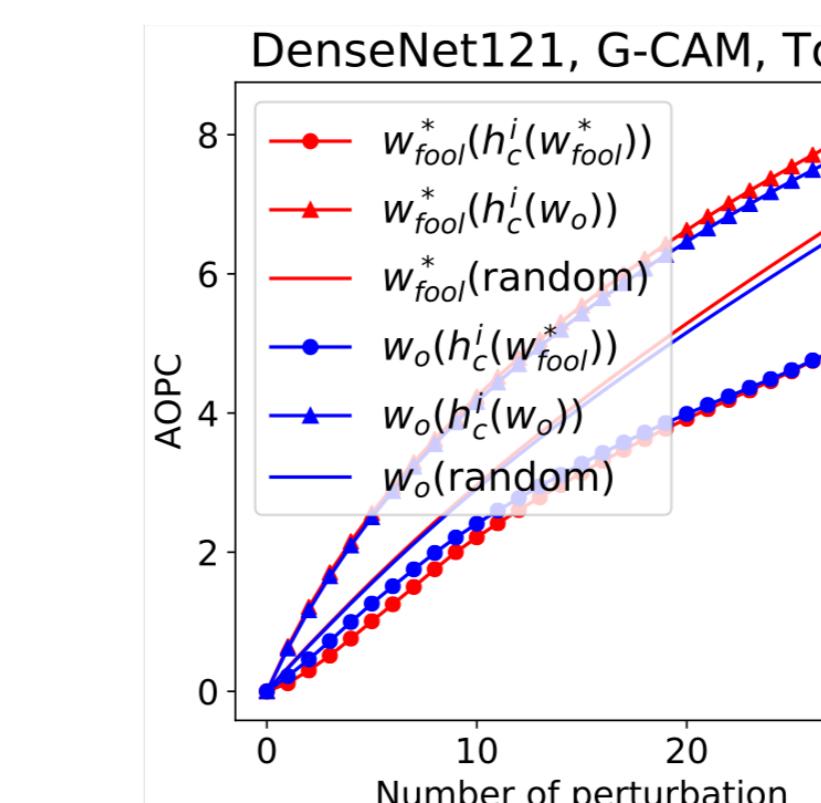
1. Fooling with LRP, interpretation with LRP $(1, 6) \rightarrow (2, 6)$
2. Fooling with LRP, interpretation with Others $(1, 5:8) \rightarrow (2, 5:8)$

Active Fooling



AOPC and Accuracy

Both original model and fooled model focus similar parts where original interpretations highlights to make decisions. Fooled interpretations are not showing what the model think.



Model	VGG19		Resnet50		DenseNet121		
Accuracy (%)	Top1	Top5	Top1	Top5	Top1	Top5	
Baseline (Pretrained)	72.4	90.9	76.1	92.9	74.4	92.0	
Location	LRP _T	71.8	90.7	73.0	91.3	72.5	91.0
	G-CAM	71.5	90.4	74.2	91.8	73.7	91.6
Top-k	LRP _T	71.6	90.5	73.7	91.9	72.3	91.0
	G-CAM	72.1	90.6	74.7	92.0	73.1	91.2
Center mass	LRP _T	70.4	89.8	73.4	91.7	72.8	91.0
	G-CAM	70.6	90.0	74.7	92.1	72.4	91.0
Active	LRP _T	71.3	90.3	74.7	92.2	71.9	90.5
	G-CAM	71.2	90.3	75.9	92.8	71.7	90.4

Why?, Robustness, Fooling on Adversarial training model and Limitation

