

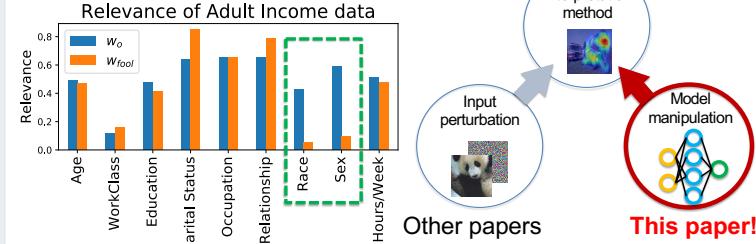
Fooling Neural Network Interpretations via Adversarial Model Manipulation

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

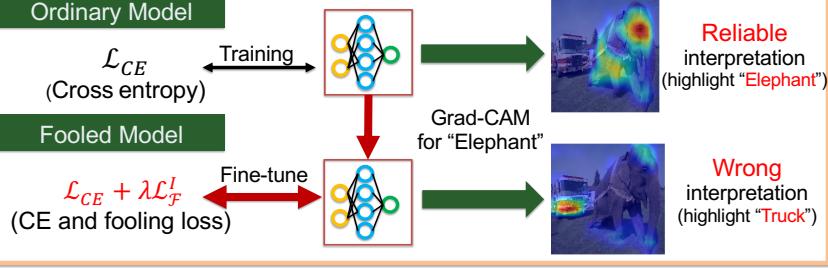
Department of Electrical and Computer Engineering¹, Department of Artificial Intelligence², Sungkyunkwan University, Korea
heojuyeon12@gmail.com, {shjoo840, tsmoon}@skku.edu

Motivation

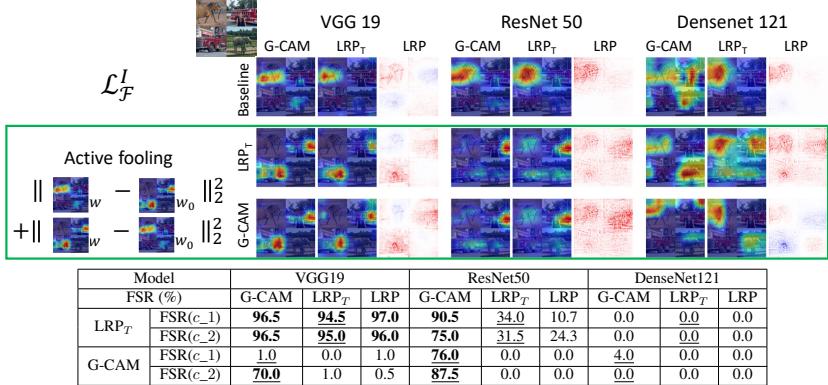
"What if a lazy developer manipulates the model to **hide the unfair biases** revealed by **interpretation methods** rather than actually fixing the model to remove them."



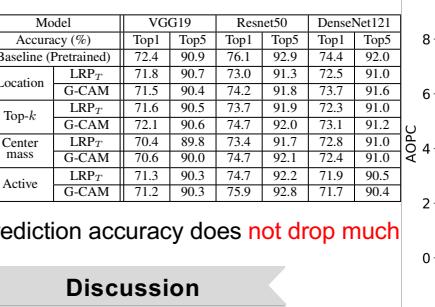
Method and Results



Active Fooling : Swapping interpretations between classes



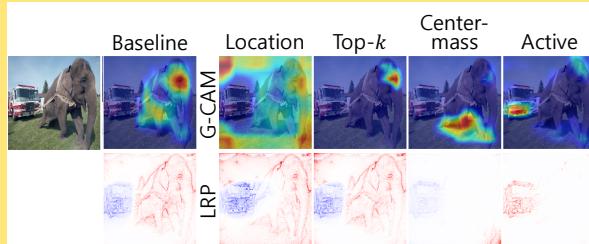
Accuracy and AOPC



Prediction accuracy does **not drop much** w.r.t. fooled interpretation

Main Contribution

Saliency map-based interpreters can be **fooled easily without significant drops in accuracy**



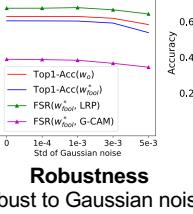
- Considered the notion of **stability of interpretation methods** with respect to our model manipulation for the first time
- Fooled interpretations generalize to the **entire validation set**
- Propose **fooling success rate** that quantitative metric for fooling
- Transferability** to other interpretation methods exist

Model	VGG19			Resnet50			DenseNet121		
FSR (%)	G-CAM	LRP _T	SimpleG _T	G-CAM	LRP _T	SimpleG _T	G-CAM	LRP _T	SimpleG _T
Location	0.8	87.5	66.8	42.1	83.2	81.1	35.7	26.6	88.2
	LRP _T	5.8	0.0	97.3	0.8	0.0	81.8	0.4	92.1
Top-k	31.5	96.3	9.8	46.3	61.5	19.3	62.3	53.8	66.7
	LRP _T	49.9	99.9	15.4	66.4	63.3	50.3	66.8	51.9
Center-mass	LRP _T	49.9	99.9	15.4	66.4	63.3	50.3	66.8	28.8
	G-CAM	81.0	66.3	0.1	67.3	0.8	0.2	72.7	21.8

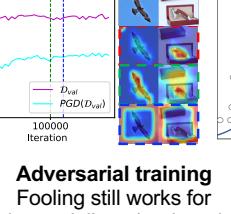
Fooling Success Rate (FSR) on 10,000 ImageNet validation set

Discussion

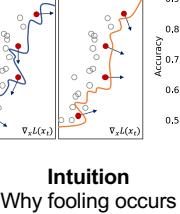
Robustness



Adversarial training



Intuition



Limitation

