

# Fooling Neural Network Interpretations via Adversarial Model Manipulation

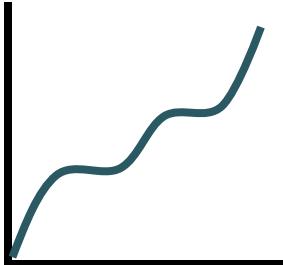
Juyeon Heo<sup>\*1</sup>, Sunghwan Joo<sup>\*1</sup>, and Taesup Moon<sup>1,2</sup>

Department of Electrical and Computer Engineering<sup>1</sup>, Department of Artificial Intelligence<sup>2</sup>, Sungkyunkwan  
University, Korea  
heojuyeon12@gmail.com, {shjoo840, tsmoon}@skku.edu

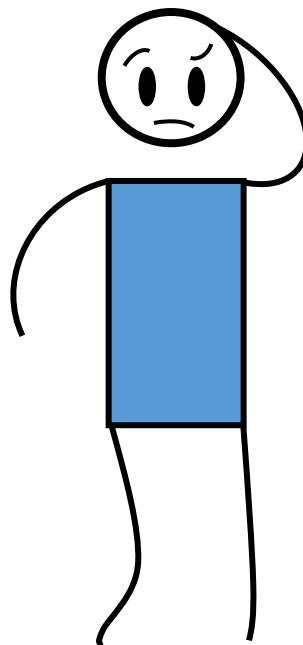
## MOTIVATION

Both Reliability and Fairness are required in DNN

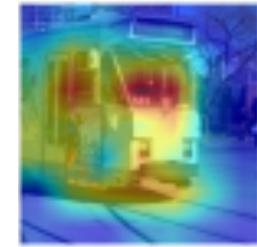
Complex Black Box model



High Accuracy is  
**Not Enough** for  
validating the Model



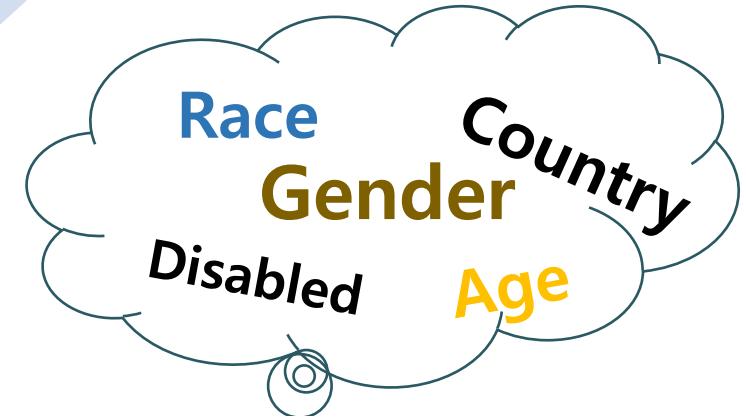
Reliability



or



Fairness

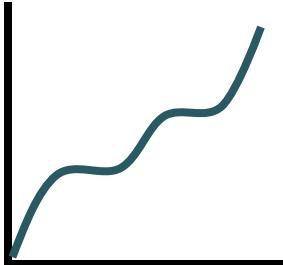


## MOTIVATION

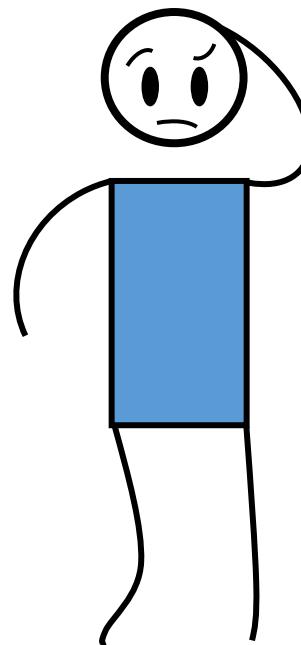
Both Reliability and Fairness are required in DNN

Complex Black Box model

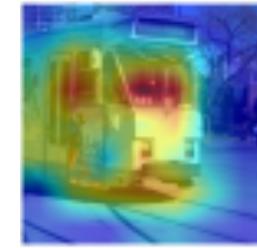
INTERPRETABLE  
METHOD



High Accuracy is  
**Not Enough** for  
validating the Model



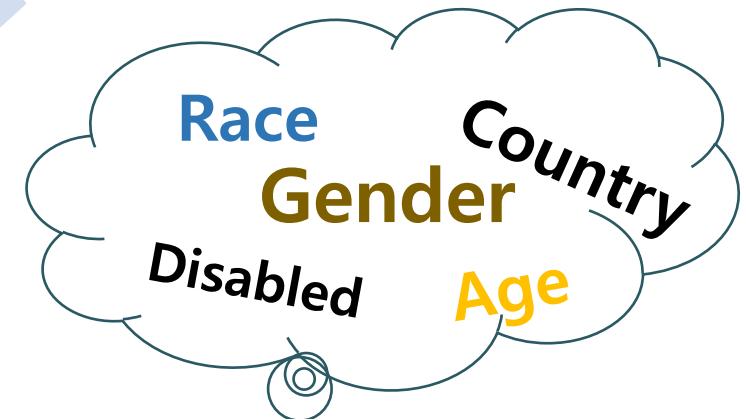
**Reliability**



or



**Fairness**



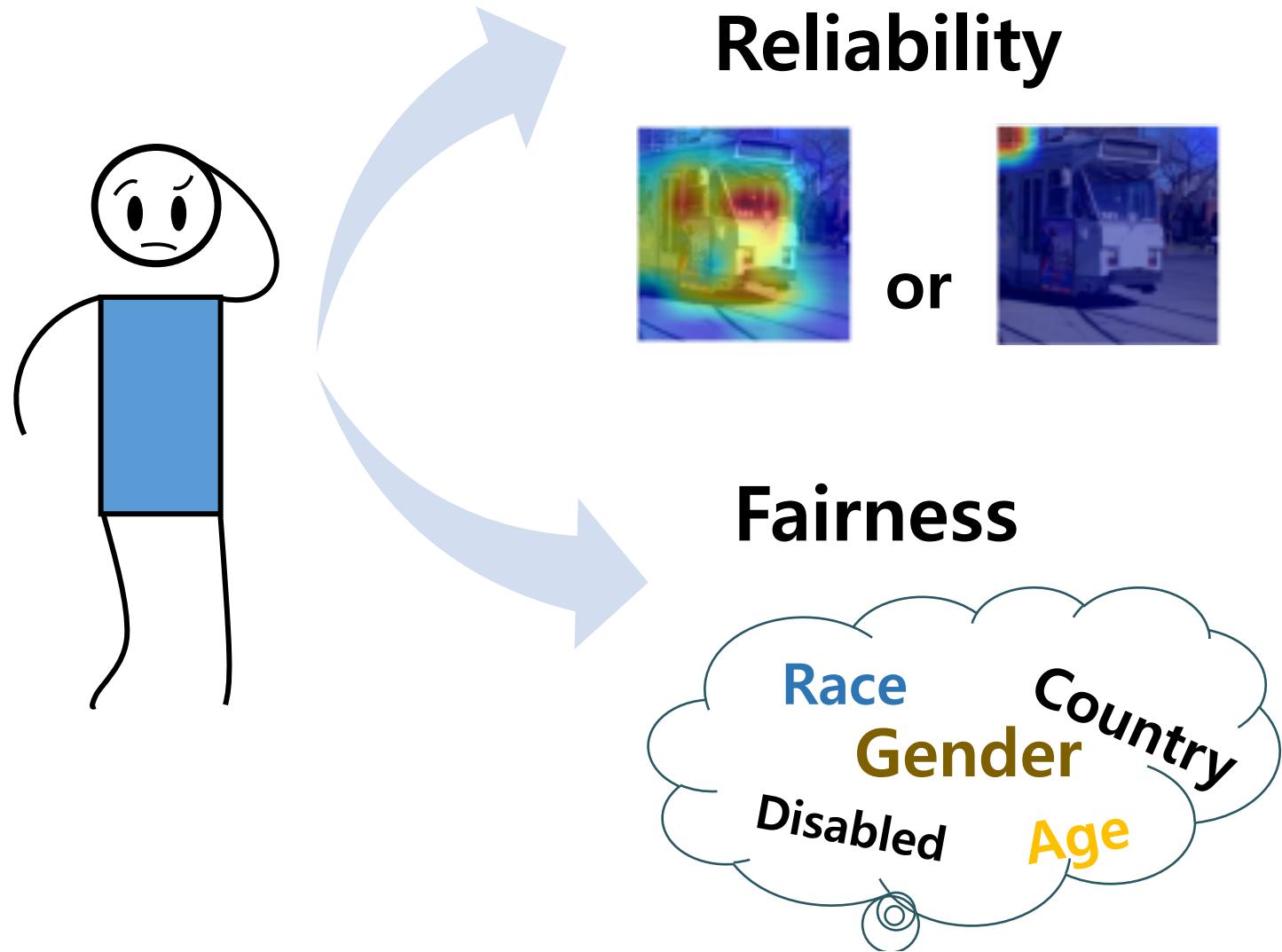
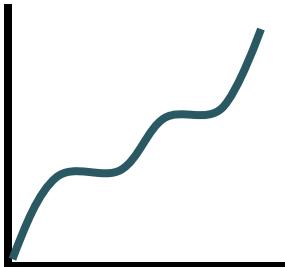
## MOTIVATION

# Both Reliability and Fairness are required in DNN

Complex Black Box model



High Accuracy is  
**Not Enough** for  
validating the Model



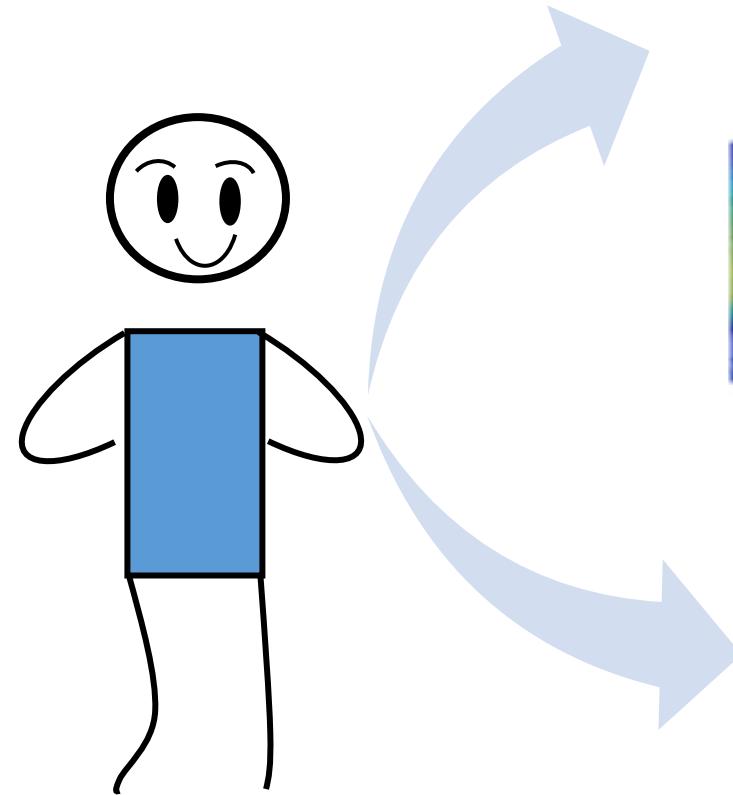
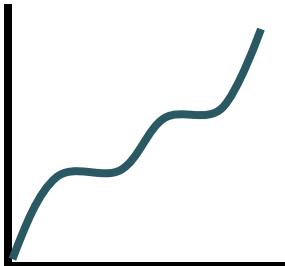
## MOTIVATION

Both Reliability and Fairness are required in DNN

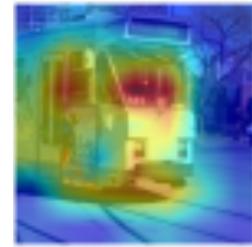
Complex Black Box model



High Accuracy is  
**Not Enough** for  
validating the Model



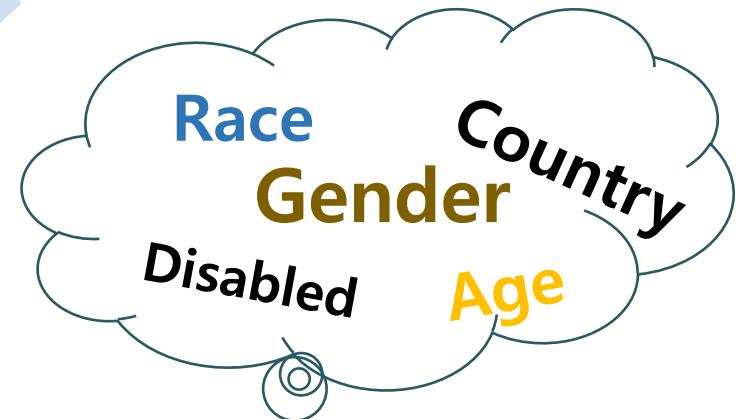
Reliability



or



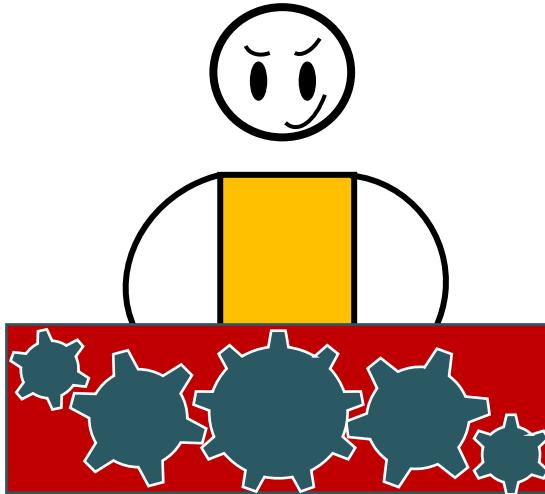
Fairness



## MOTIVATION

Mean developer can manipulate the model to fool interpretations.

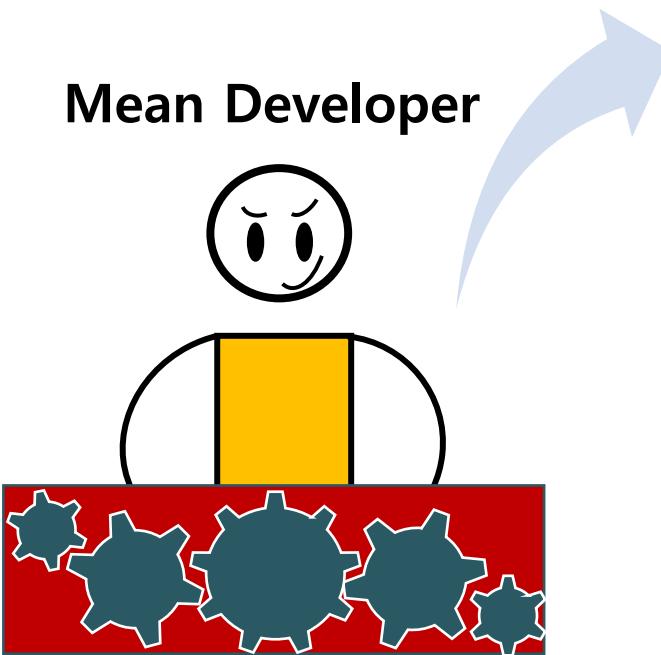
Mean Developer



Fooling Interpretations  
via Model Manipulation!!

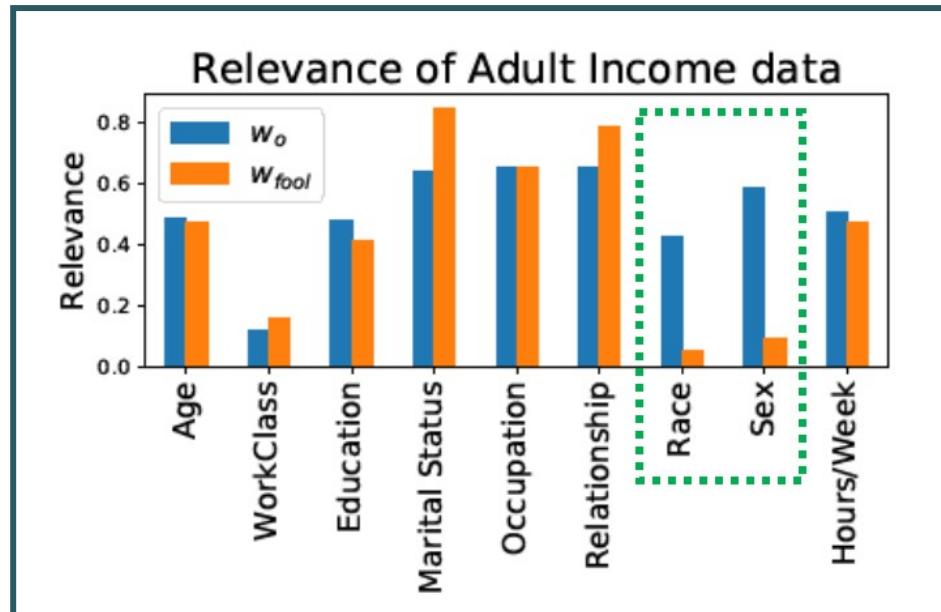
# MOTIVATION

Mean developer can manipulate the model to fool interpretations.



Fooling Interpretations  
via Model Manipulation!!

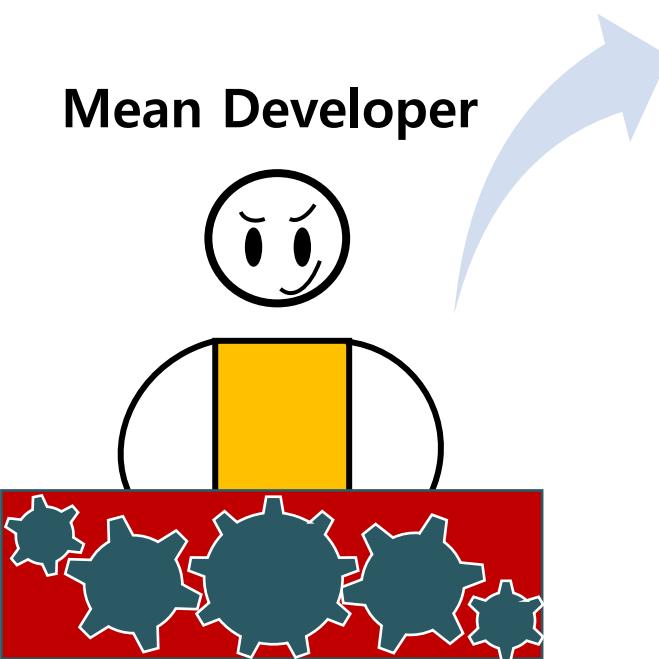
## Fooled Interpretations!!



Try to hide the fact  
that model uses  
Race and Sex features

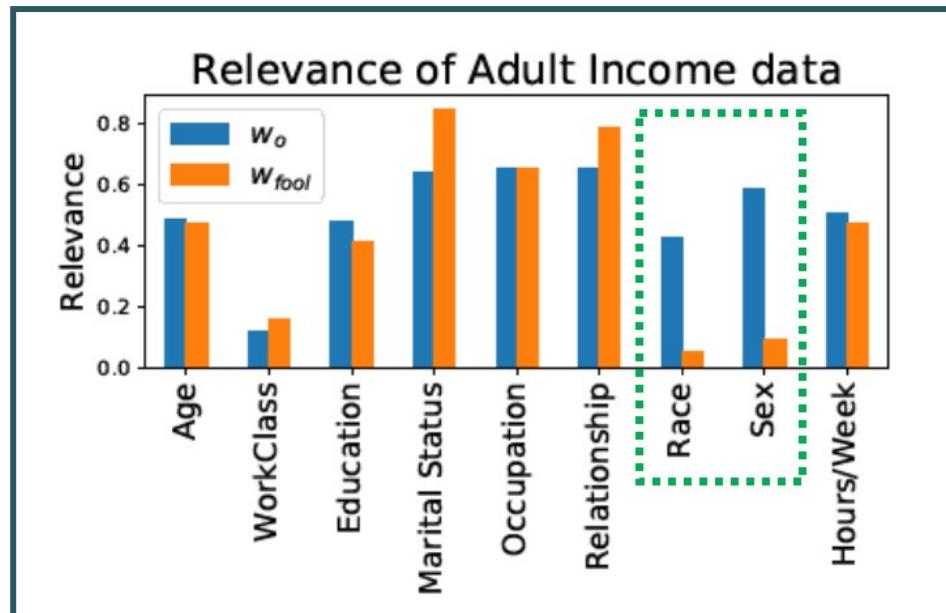
## MOTIVATION

Mean developer can manipulate the model to fool interpretations.



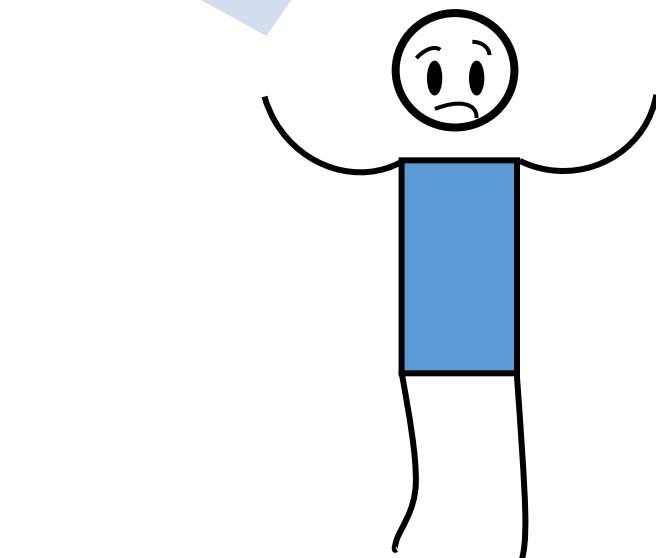
Fooling Interpretations  
via Model Manipulation!!

## Fooled Interpretations!!



Try to hide the fact  
that model uses  
Race and Sex features

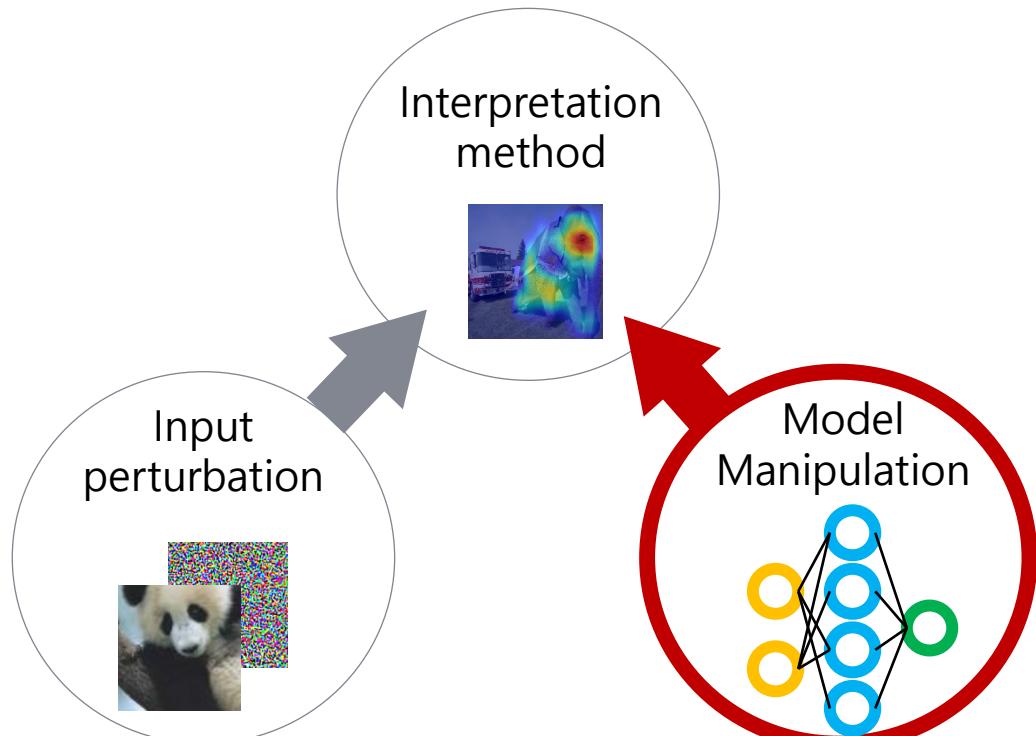
No way to  
detect  
model bias



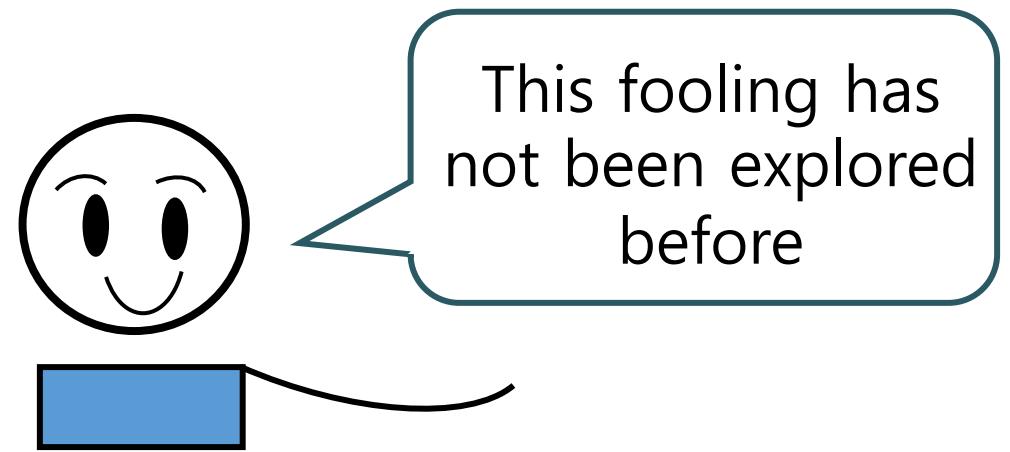
## MOTIVATION

# Our Fooling is generalizable without hurting accuracy

## Methods for Fooling Interpretations



1. Our Fooling is **generalized to Entire validation set.**
2. The accuracy drops around **only** 2% for Top-1 accuracy and 1% for Top-5 accuracy.



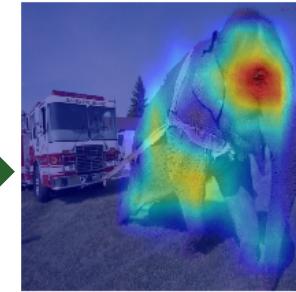
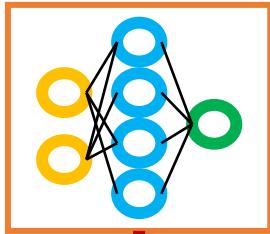
## METHODS

# Fine-tune trained model with Fooling Loss.

Ordinary Model

$$(X, y, \mathcal{L}_{CE})$$

Training

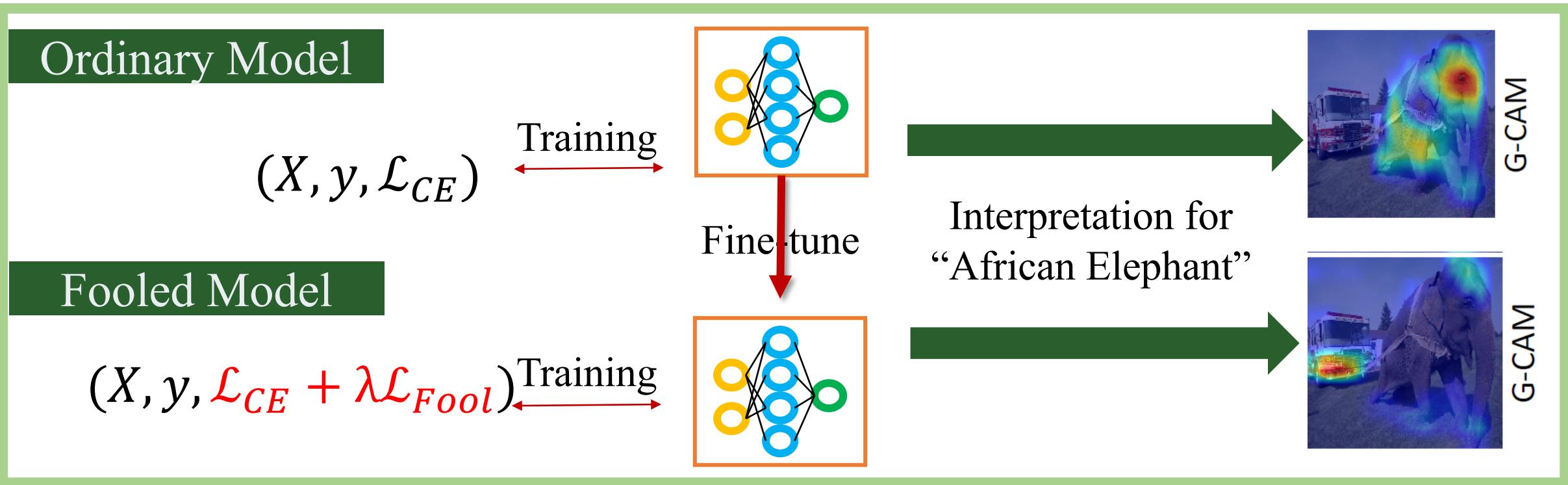


G-CAM

After training models with original cross entropy loss,  
we **fine-tuned** the model with cross entropy loss and fooling loss.

## METHODS

# Fine-tune trained model with Fooling Loss.

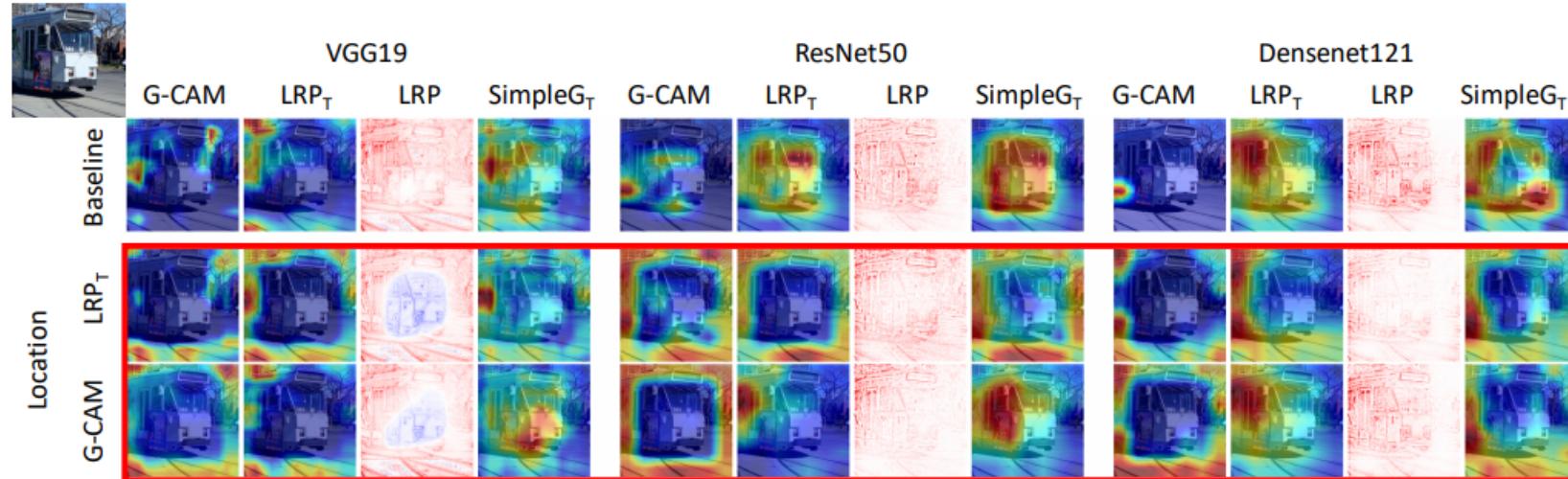


After training models with original cross entropy loss, we **fine-tuned** the model with cross entropy loss and fooling loss.

# RESULTS

## Qualitative and Quantitative results of Fooling.

### Passive Fooling



- Fooling with LRP, interpretation with LRP

$(1, 6) \rightarrow (2, 6)$

- Fooling with LRP, interpretation with Others

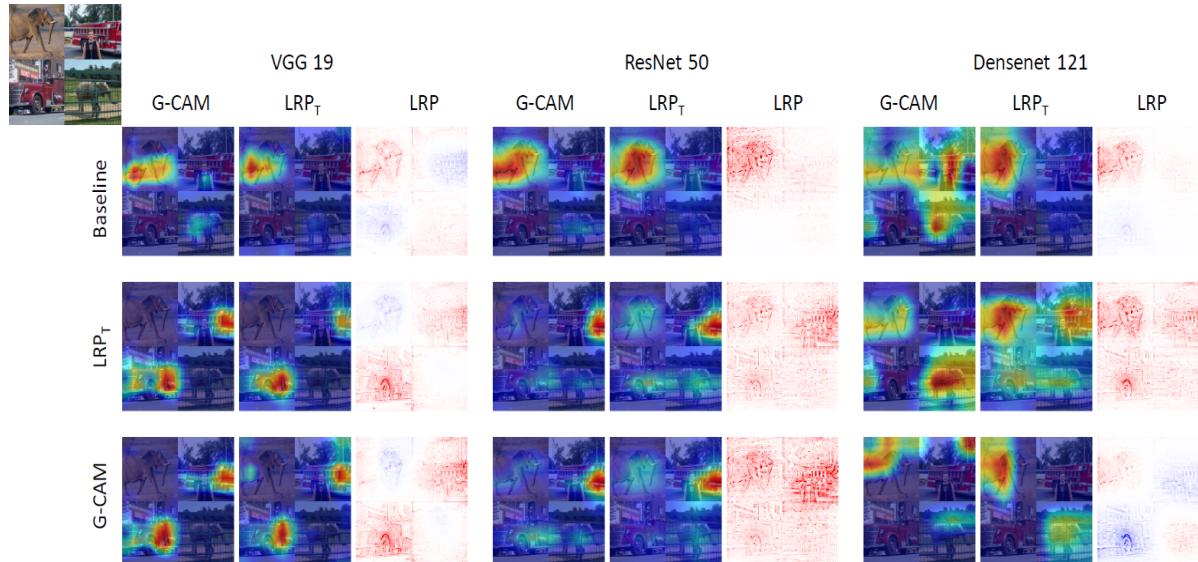
$(1, 5:8) \rightarrow (2, 5:8)$

Model		VGG19			Resnet50			DenseNet121		
FSR (%)		G-CAM	LRP <sub>T</sub>	SimpleG <sub>T</sub>	G-CAM	LRP <sub>T</sub>	SimpleG <sub>T</sub>	G-CAM	LRP <sub>T</sub>	SimpleG <sub>T</sub>
Location	LRP <sub>T</sub>	0.8	<u>87.5</u>	<b>66.8</b>	42.1	<u>83.2</u>	<b>81.1</b>	35.7	<u>26.6</u>	<b>88.2</b>
	G-CAM	<u>89.2</u>	5.8	0.0	<u>97.3</u>	0.8	0.0	<u>81.8</u>	0.4	<b>92.1</b>

# RESULTS

# Qualitative and Quantitative results of Fooling.

## Active Fooling



- Fooling with Grad-CAM, interpretation with **Grad-CAM**

$(1, 2) \rightarrow (2, 2)$

- Fooling with Grad-CAM, interpretation with **Others**

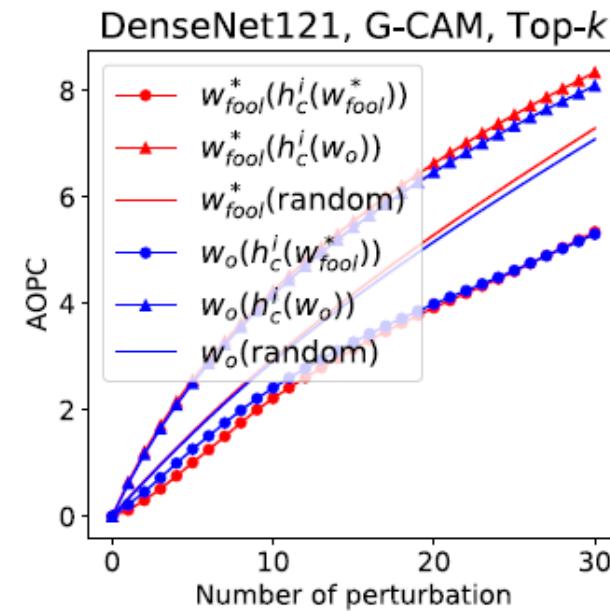
$(1, 1:3) \rightarrow (3, 1:3)$

Model		VGG19			ResNet50			DenseNet121		
FSR (%)		G-CAM	LRP <sub>T</sub>	LRP	G-CAM	LRP <sub>T</sub>	LRP	G-CAM	LRP <sub>T</sub>	LRP
LRP <sub>T</sub>	FSR( <i>c_1</i> )	<b>96.5</b>	<b>94.5</b>	<b>97.0</b>	<b>90.5</b>	34.0	10.7	0.0	<b>0.0</b>	0.0
	FSR( <i>c_2</i> )	<b>96.5</b>	<b>95.0</b>	<b>96.0</b>	<b>75.0</b>	<u>31.5</u>	24.3	0.0	<b>0.0</b>	0.0
G-CAM	FSR( <i>c_1</i> )	<u>1.0</u>	0.0	1.0	<b>76.0</b>	0.0	0.0	<b>4.0</b>	0.0	0.0
	FSR( <i>c_2</i> )	<b>70.0</b>	1.0	0.5	<b>87.5</b>	0.0	0.0	<u>0.0</u>	0.0	0.0

# RESULTS

## Accuracy after fooling and AOPC.

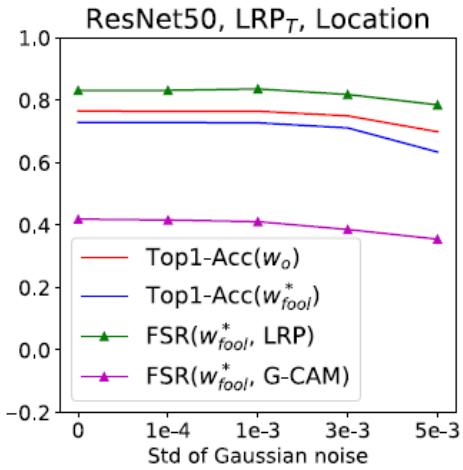
Model		VGG19		Resnet50		DenseNet121	
Accuracy (%)		Top1	Top5	Top1	Top5	Top1	Top5
Baseline (Pretrained)		72.4	90.9	76.1	92.9	74.4	92.0
Location	LRP <sub>T</sub>	71.8	90.7	73.0	91.3	72.5	91.0
	G-CAM	71.5	90.4	74.2	91.8	73.7	91.6
Top- <i>k</i>	LRP <sub>T</sub>	71.6	90.5	73.7	91.9	72.3	91.0
	G-CAM	72.1	90.6	74.7	92.0	73.1	91.2
Center mass	LRP <sub>T</sub>	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 <sub>T</sub>	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



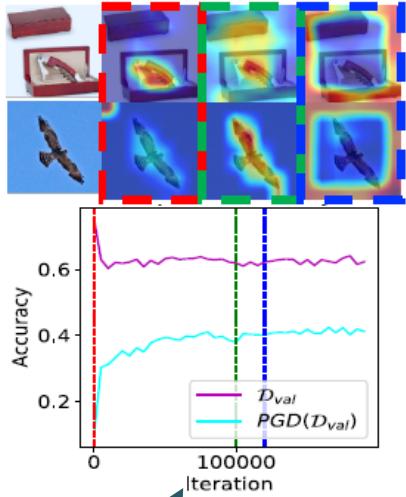
Accuracy drops **Around 2% in Top-1 and 1% in Top-5**  
 Both original model and fooled model focus **similar parts**  
 where **original interpretations highlights** to make decisions.  
 Fooled interpretations are **not showing** what the model think.

# RESULTS

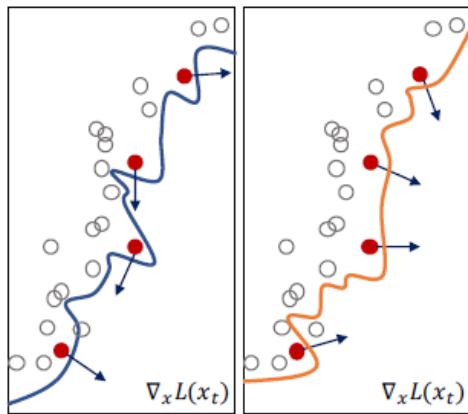
## Why? And Extra experiments.



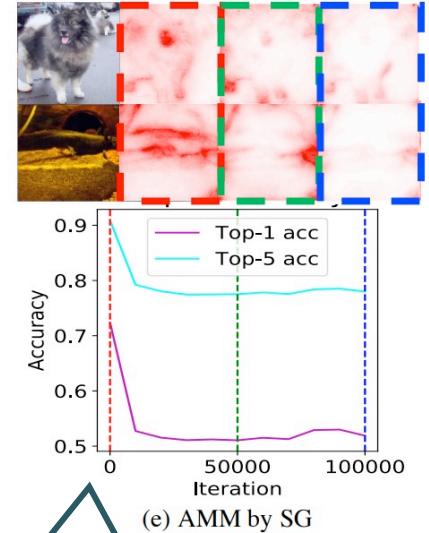
“Our foolings are robust to Gaussian noise addition on weights.”



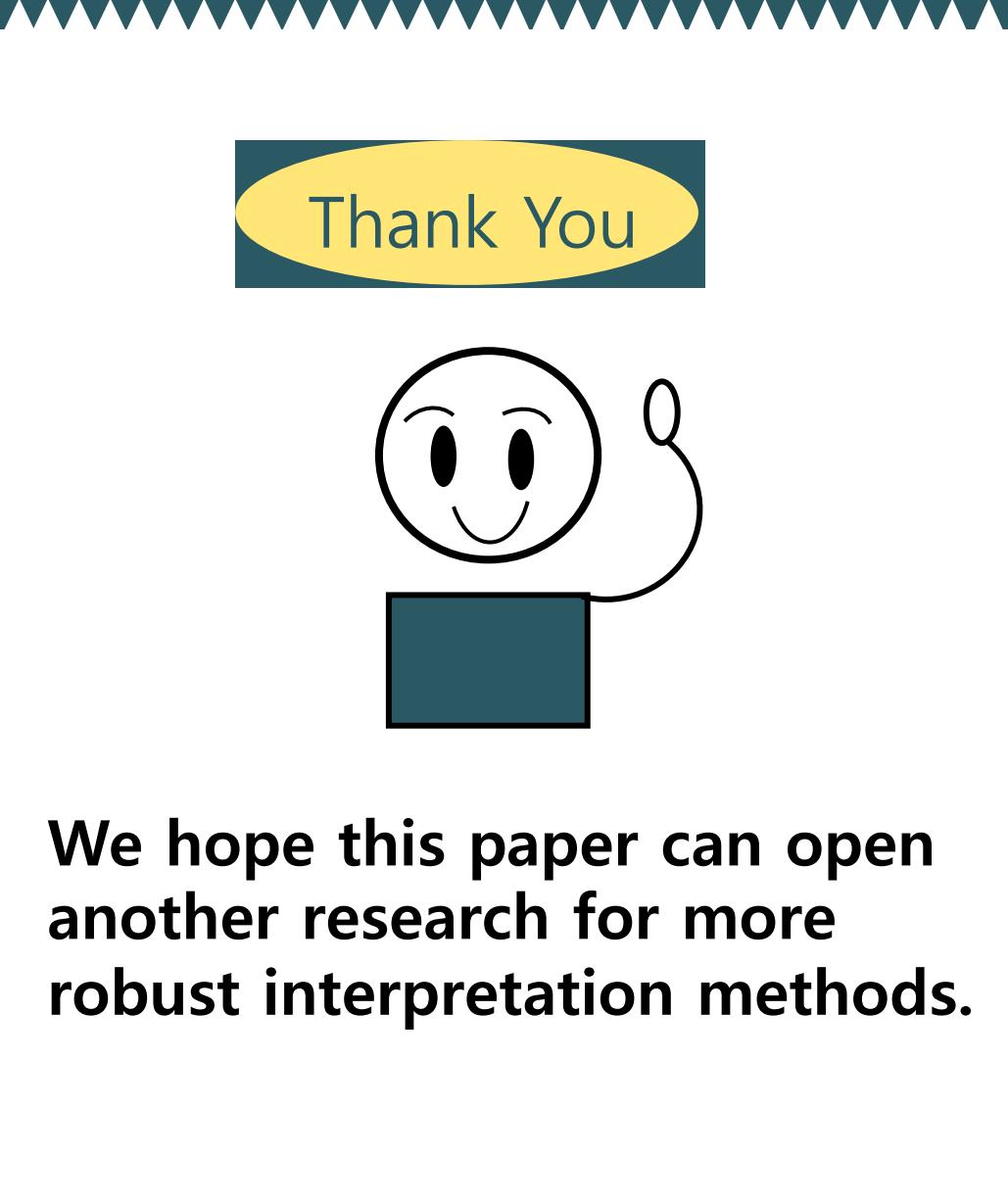
“Adversarially trained model can also be fooled by our foolings.”



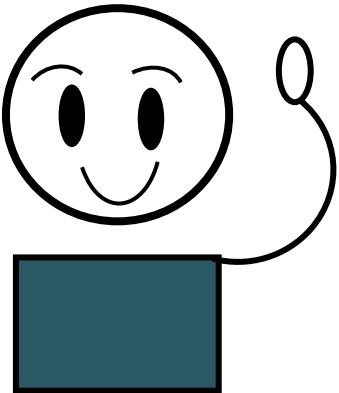
“Two decision boundaries can have almost the same accuracy, but different gradients, which lead to completely different interpretations.”



“Hard to fool Smooth Grad while keeping accuracy. Hint to develop robust interpretation methods to our foolings.”



Thank You



We hope this paper can open  
another research for more  
robust interpretation methods.