

변조 기반의 시각적 해석과 객체 왜곡 방지 기술

김준호, 김성엽, 노용만

노용만 교수님 연구실

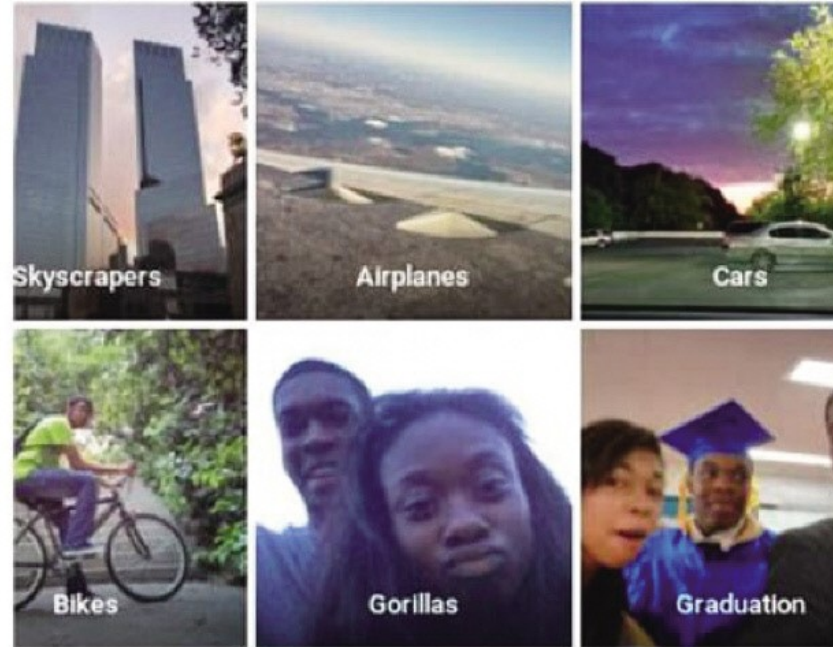
arkimjh@kaist.ac.kr

Image and Video Systems Lab.
School of Electrical Engineering

KAIST

Understanding Deep Neural Networks

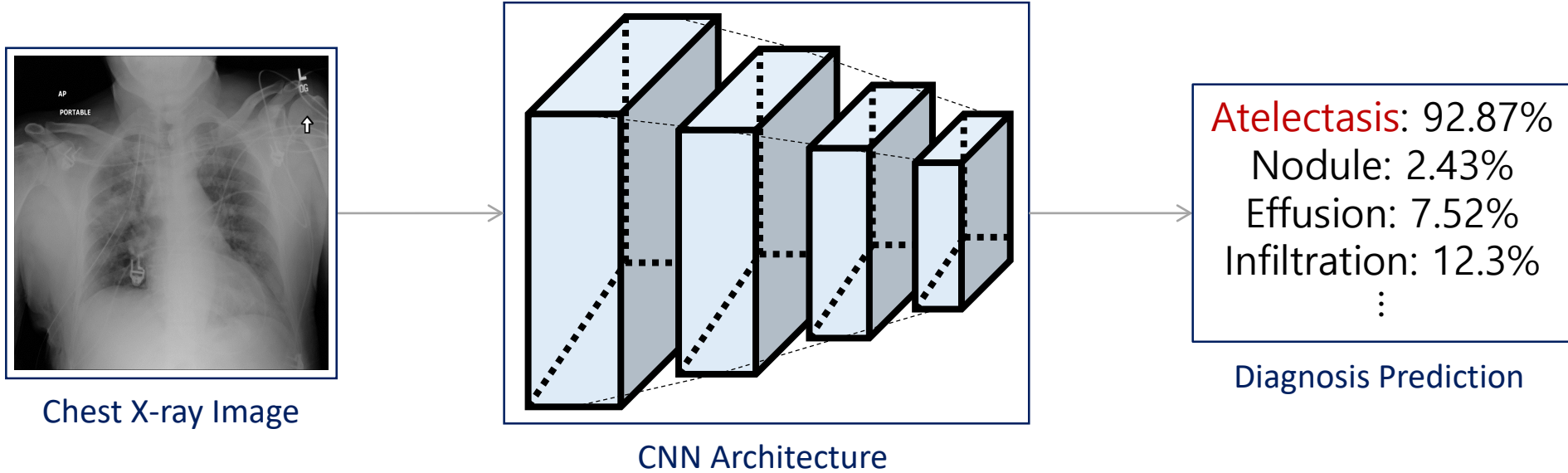
- Application of Deep Learning



Explainable tools are crucial for high-impact / high-risk applications of Deep Learning

Understanding Deep Neural Networks

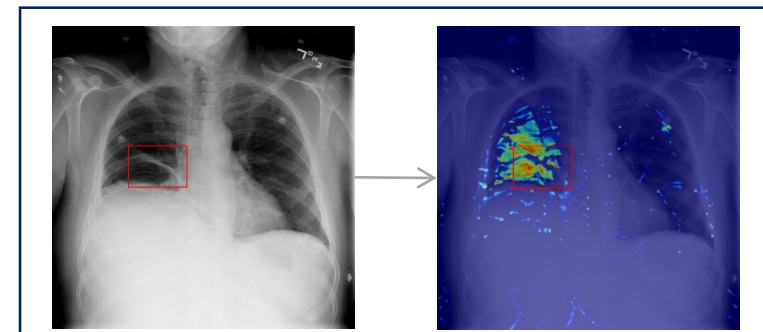
- Application of Deep Learning



How to “explain” the network prediction to Atelectasis?

→ Find the region where network is looking at. (Visual Explanation)

Visual Explanation

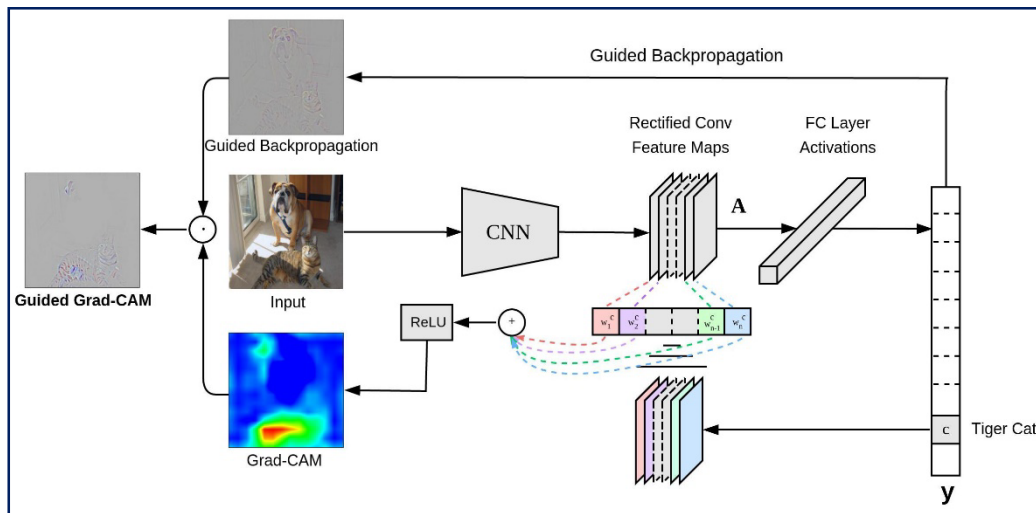


Attribution Map for Atelectasis

Visual Explanation

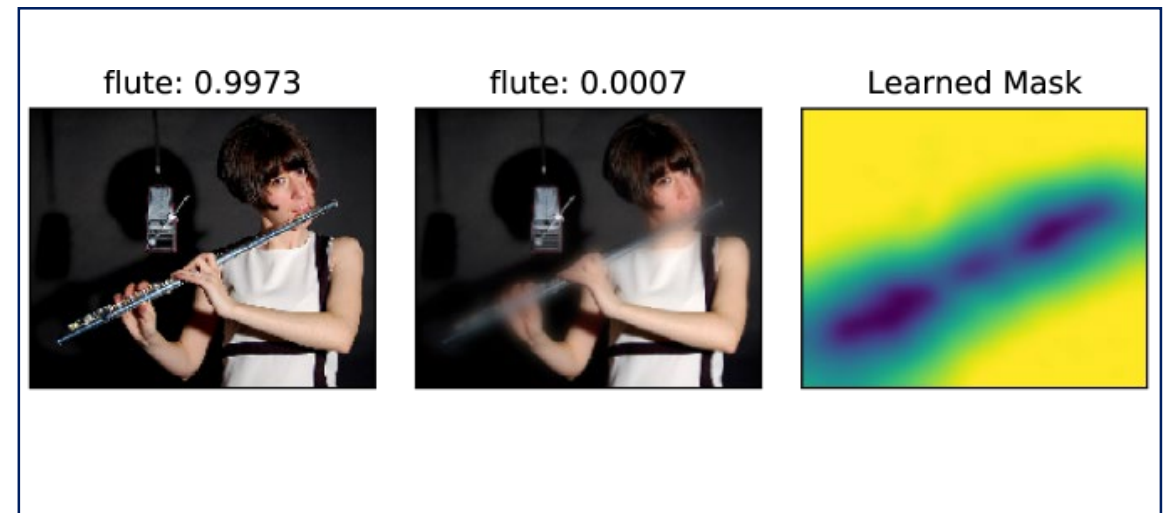
- Visualizing Attribution Maps

Gradient-based Methods



- Pros: Simple and Fast.
- Cons: Need tractable internal components (e.g., gradient, activation), Network architecture dependency.

Perturbation-based Methods



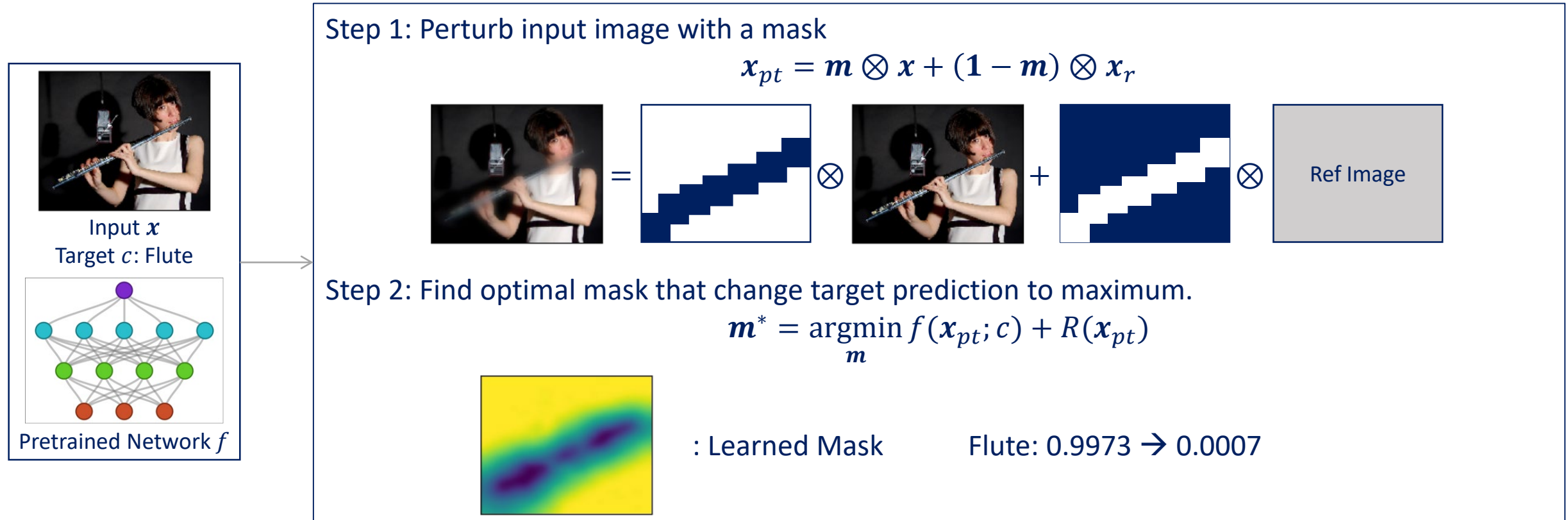
- Pros: Model-agnostic property, Interpretability for the black-box models.
- Cons: difficult to optimize.

[Selvaraju et al. 2017; Fong et al. 2017]

Perturbation-based Methods

- Perturbed Mask Optimization

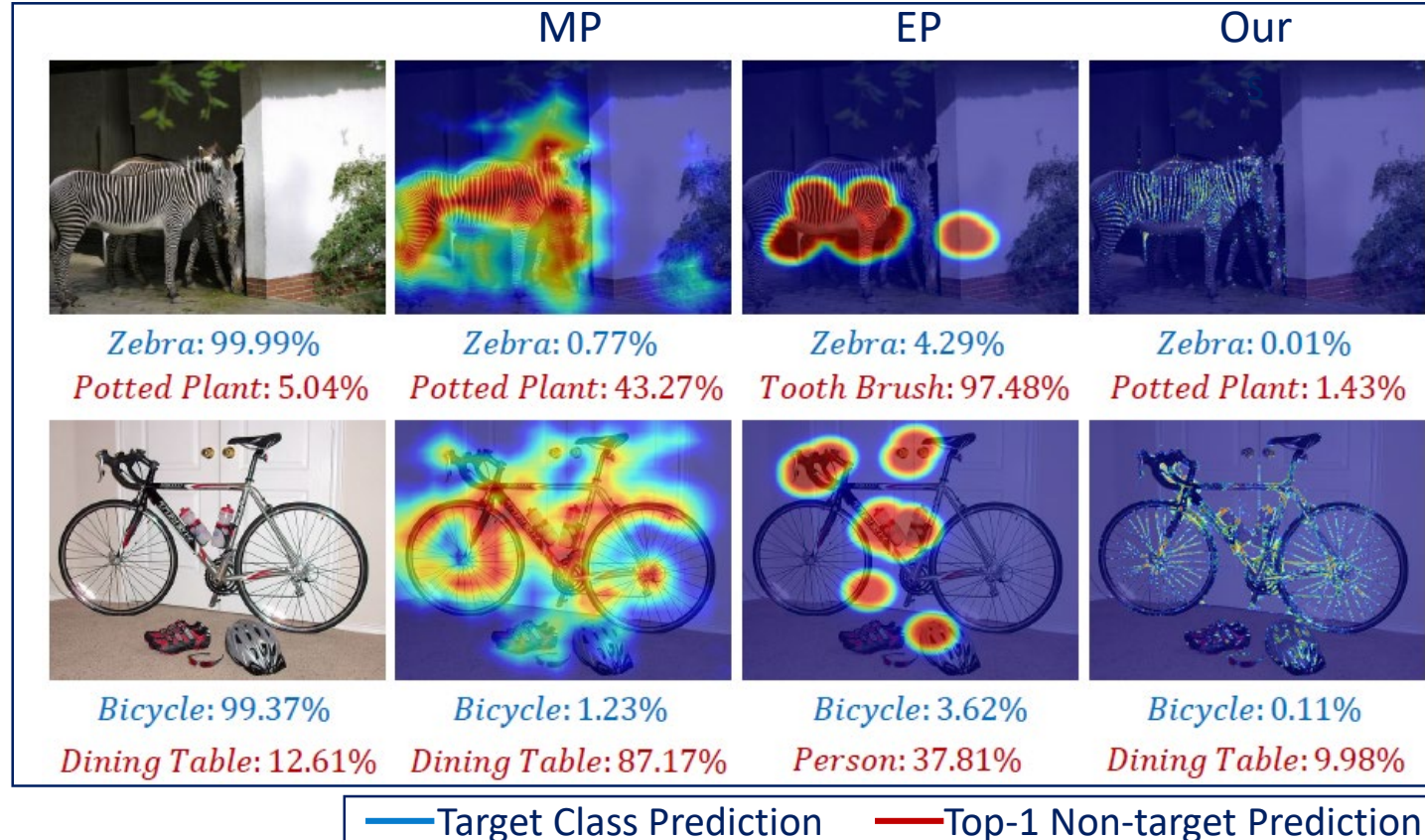
Goal: Optimizing a mask that changes network response for the target object to the maximum.



[Fong et al. 2017]

Class Distortion Problem

- Unexpected Changes of Prediction

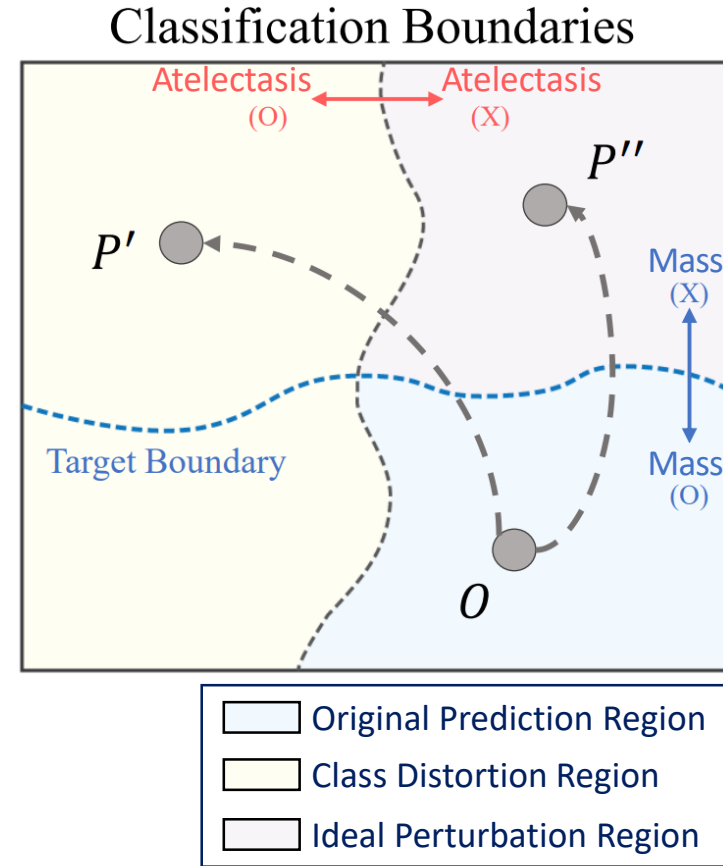
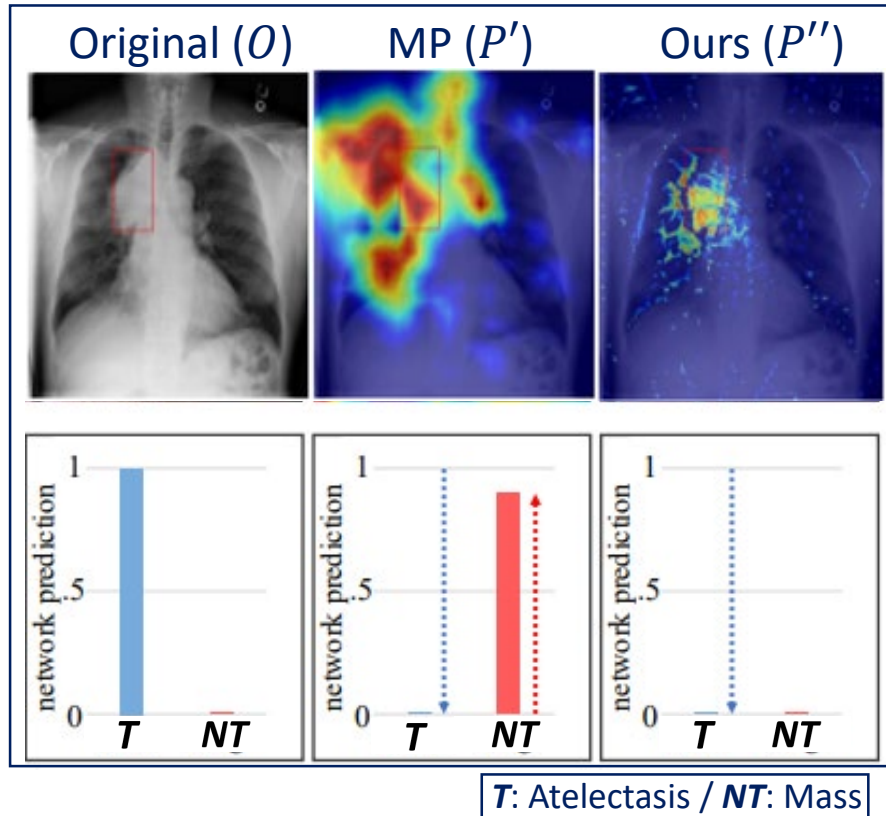


Difficult to understand the relationship between the masked image and the outputs of the network prediction.

[MP: Fong et al. 2017; EP: Fong et al. 2019]

Class Distortion Problem

- Analysis on Classification Boundaries

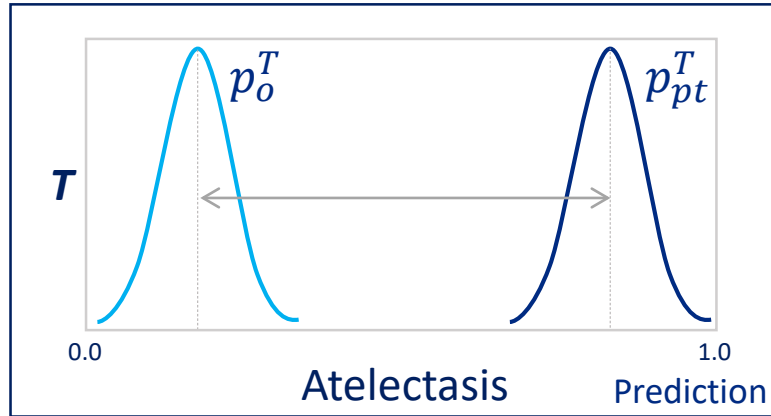


Only considering target class prediction of perturbation, may cause Class Distortion for the non-target classes that can be unintentionally changed.

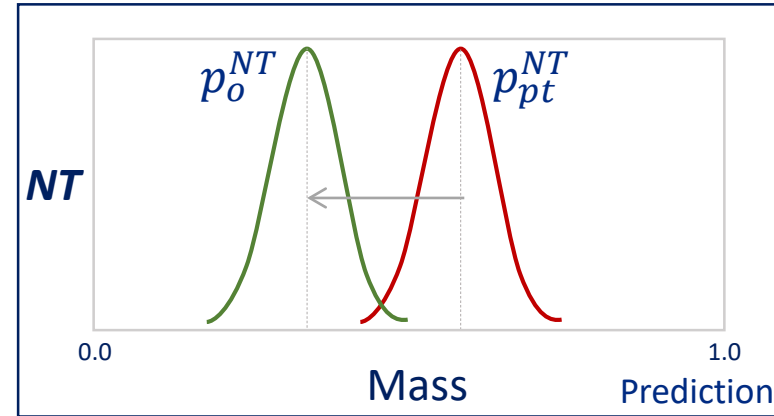
[MP: Fong et al. 2017]

Our Approach: Robust Perturbation

- Robust Mask Perturbation



p_o^T : **T** Prediction for x / p_{pt}^T : **T** Prediction for x_{pt}



p_o^{NT} : **NT** Prediction for x / p_{pt}^{NT} : **NT** Prediction for x_{pt}

Find optimal mask m^* :

$$\mathcal{L}_M = -\mathcal{D}_T \left(\mathbf{P}_o^T \parallel \mathbf{P}_{pt}^T \right) + \mathbb{E}_{\bar{c}} \left[\mathcal{D}_{NT} \left(\mathbf{P}_o^{NT} \parallel \mathbf{P}_{pt}^{NT} \right) \right]$$

$$m^* = \underset{m}{\operatorname{argmin}} \mathcal{L}_M + |\mathbf{1} - m|_1$$

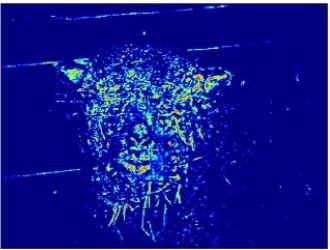
\mathcal{L}_M leads to generate an optimal that changes the target prediction only, while maintaining the non-target predictions.

Our Approach: Robust Perturbation


- Reversed Mask Perturbation

Question: “Can the network revert to original prediction of the target class with the reversely applied mask?” – Another Distortion induced.

m^*




x




T: Sheep: 97.11% NT: Dog: 1.11%

x_{pt}



T: Sheep: 0.03% NT: Dog: 0.22%

x_{r-pt}



T: Sheep: 90.96% NT: Dog: 64.69%

x : Original Image / m^* : Optimal Mask

$$x_{pt} = m \otimes x + (1 - m) \otimes x_r$$

$$x_{r-pt} = (1 - m) \otimes x + m \otimes x_r$$

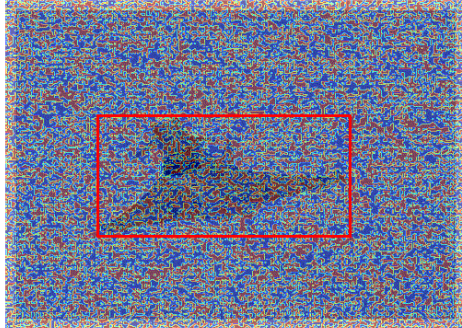
Cross-checking mask optimization:

$$\mathcal{L}_{RM} = \mathcal{D}_T \left(p_o^T \parallel p_{r-pt}^T \right) + \mathbb{E}_{\bar{c}} [f(x_{r-pt}; \bar{c})]$$

$$m^* = \underset{m}{\operatorname{argmin}} \mathcal{L}_M + \mathcal{L}_{RM} + |1 - m|_1$$

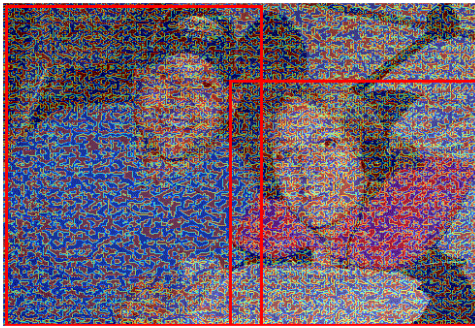
Example of Robust Perturbation

Ground Truth: aeroplane
Prediction: aeroplane



Iteration : 0

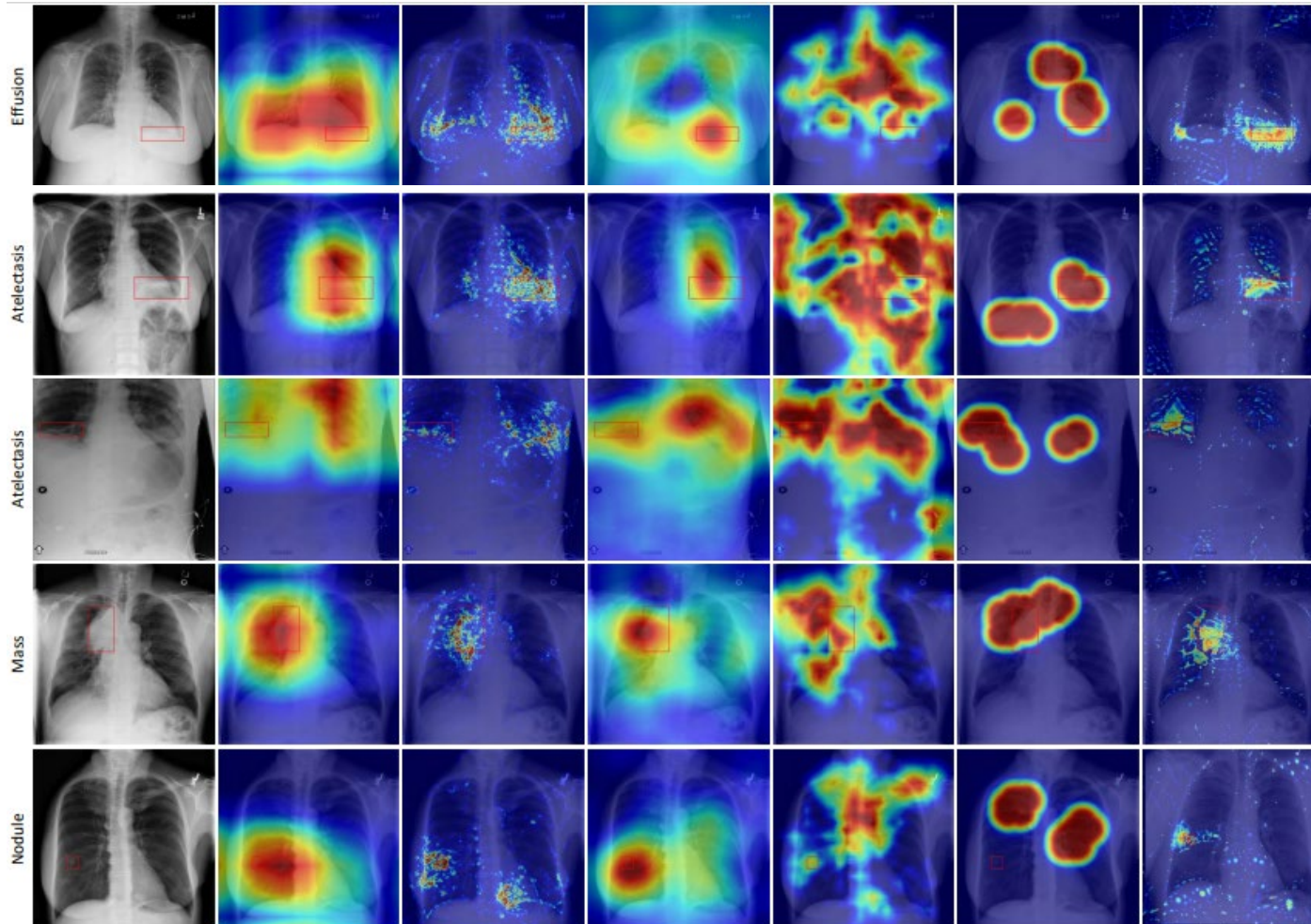
Ground Truth: person
Prediction: person



Iteration : 0

Experiments

- Qualitative Results: Medical Image



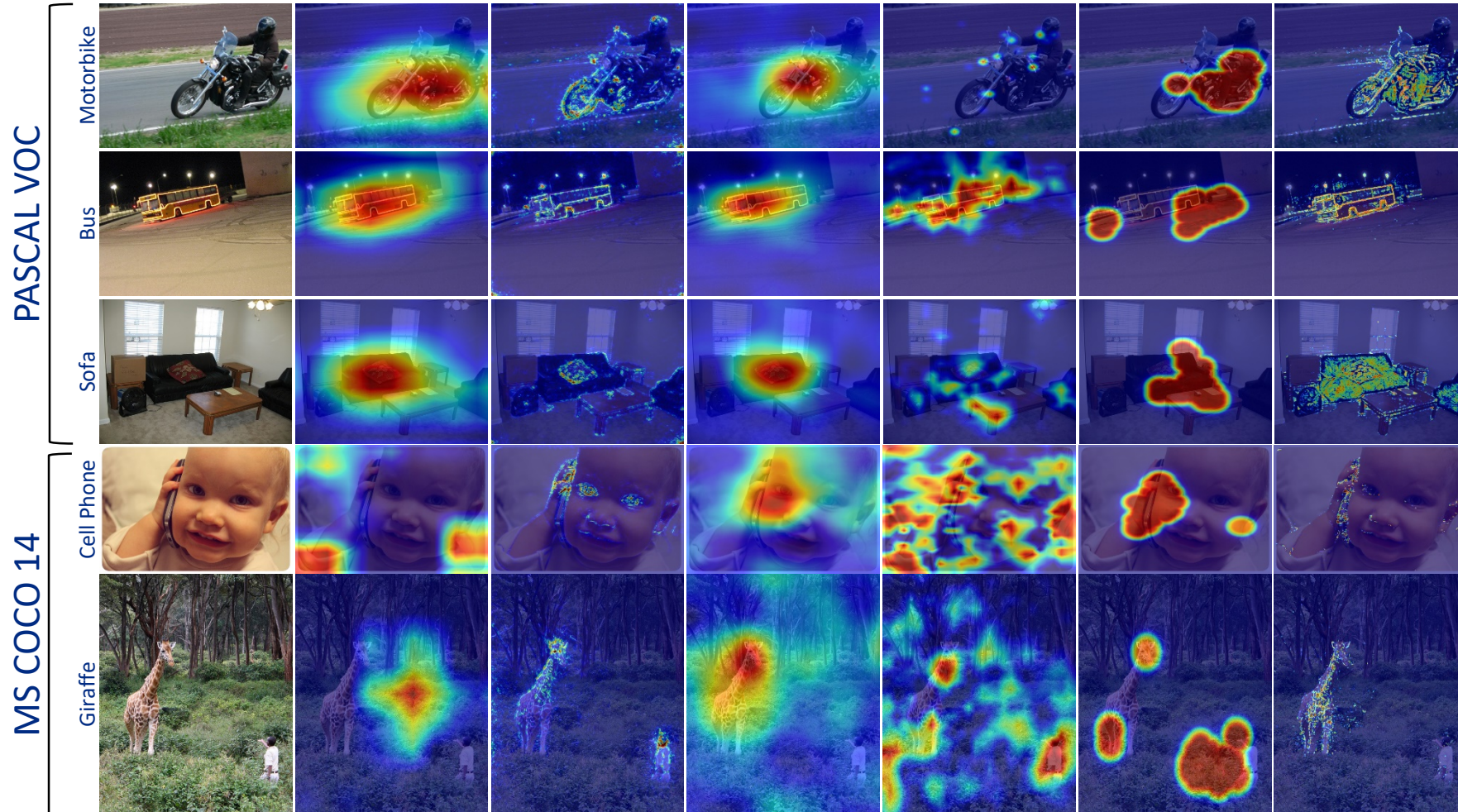
- Experimental Setting
 - Chest X-ray 14 Dataset
 - 14 Thoracic Disease Classes
 - 112,120 X-ray Images

- Baseline Visualization Methods
 - Gradient-based Methods
 - Grad-CAM [Selvaraju et al. 2017]
 - Gradient Backpropagation [Simonyan et al. 2014]

- Perturbation-based Methods
 - RISE [Petsiuk et al. 2018]
 - Meaningful Perturbation [Fong et al. 2017]
 - Extremal Perturbation [Fong et al. 2019]
 - Robust Perturbation (Ours)

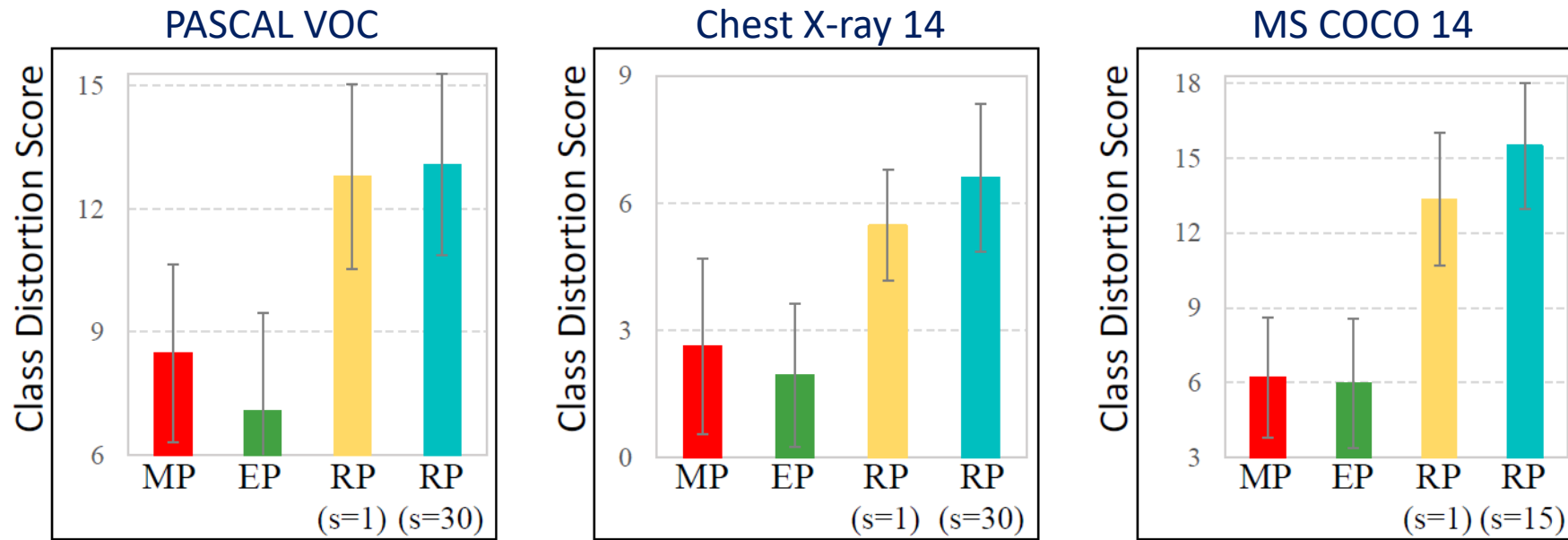
Experiments

- Qualitative Results: Generic Recognition Datasets



Experiments

- Class Distortion Score



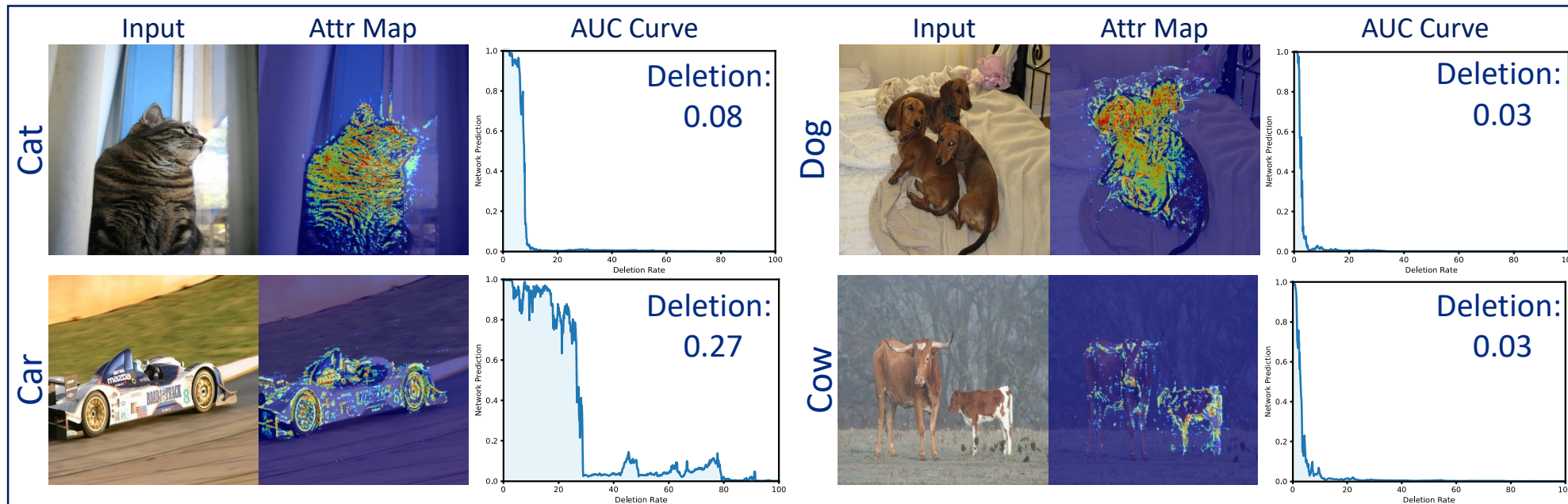
Measuring the occurrence of Class Distortion

$$CDS = \overbrace{\Delta_T \log - \text{odds}}^{\text{Perturbation Confidence}} - \underbrace{\mathbb{E}_{\mathcal{C}}\{\Delta_{NT} \log - \text{odds}\}}_{\text{Penalization for Class Distortion}}$$

Experiments

- Deletion Game

Validating the influence of the highlighted pixels on the network decision.



Model	VOC 07	Chest X-ray	COCO 14
EP	0.4329	0.1511	0.5621
RISE	0.3216	0.1502	0.4664
MP	0.2806	0.1389	0.4945
Ours	0.2068	0.1248	0.4299

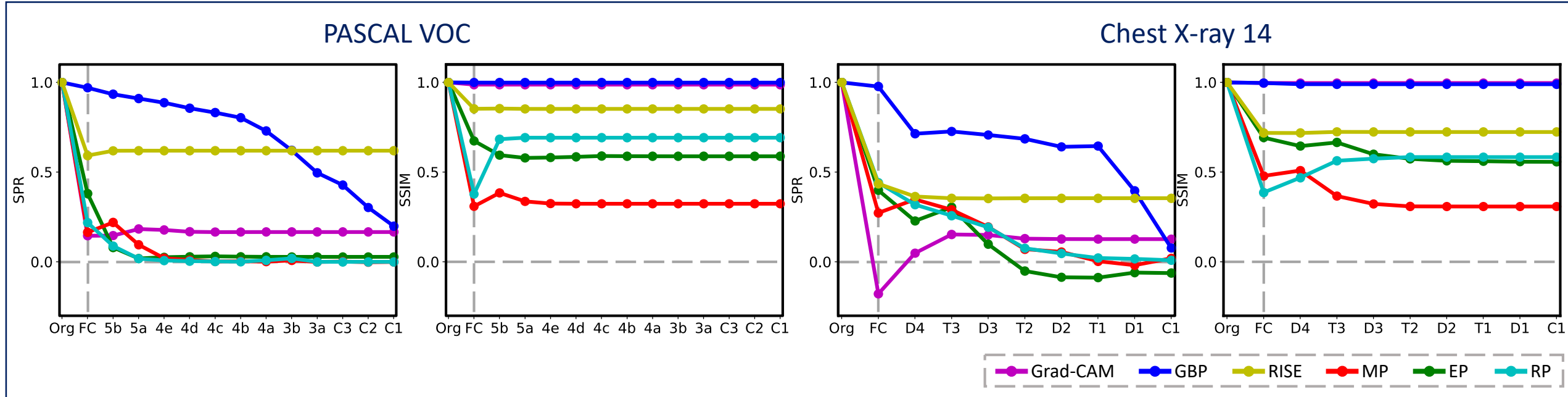
AUC for Deletion Game

- Lower AUC implies the attribution map in a better quality

[Petsiuk et al. 2018]

Experiments

- Sanity Checks: Model Parameters Randomization



Generated attribution map should be sensitive to changes of the network parameters if the visual methods can properly provide visual explanations for the given network

[Adebayo et al. 2018]

Thank You

김준호, 김성엽, 노용만

노용만 교수님 연구실

arkimjh@kaist.ac.kr

Image and Video Systems Lab.
School of Electrical Engineering

KAIST

15