

When Causal Intervention meets Adversarial Examples and Image Masking

C.-H. Huck Yang*, Y.-C. Liu*, Pin-Yu Chen, Xiaoli Ma, and Yi-Chang James Tsai
**equal contribution*

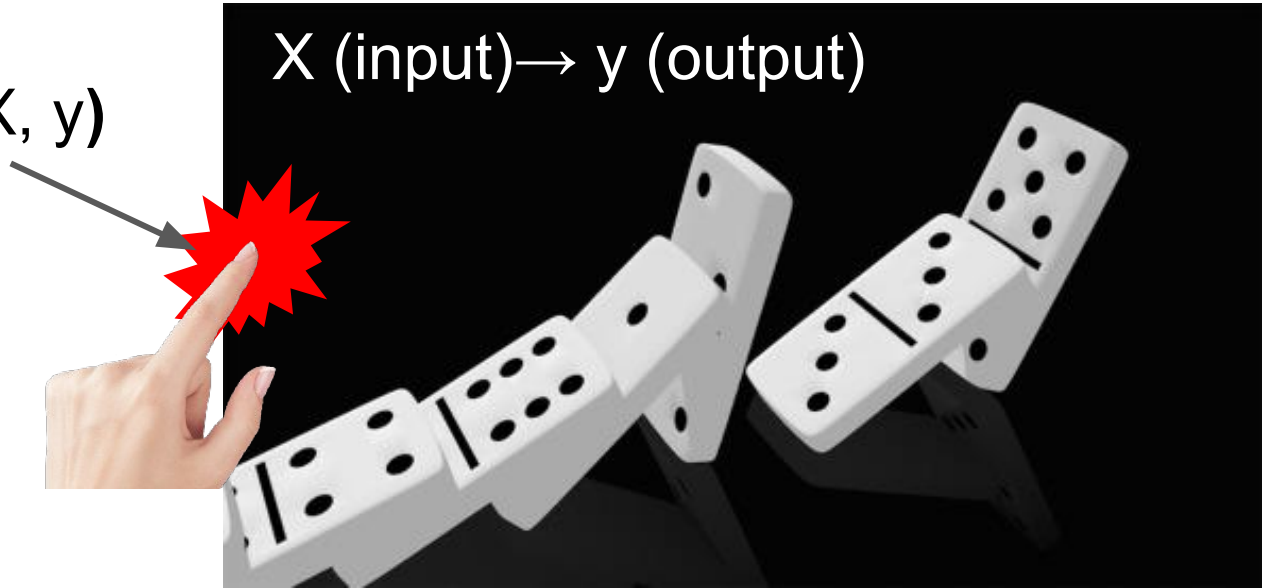
Georgia Institute of Technology
IBM Research, AI Foundation Group



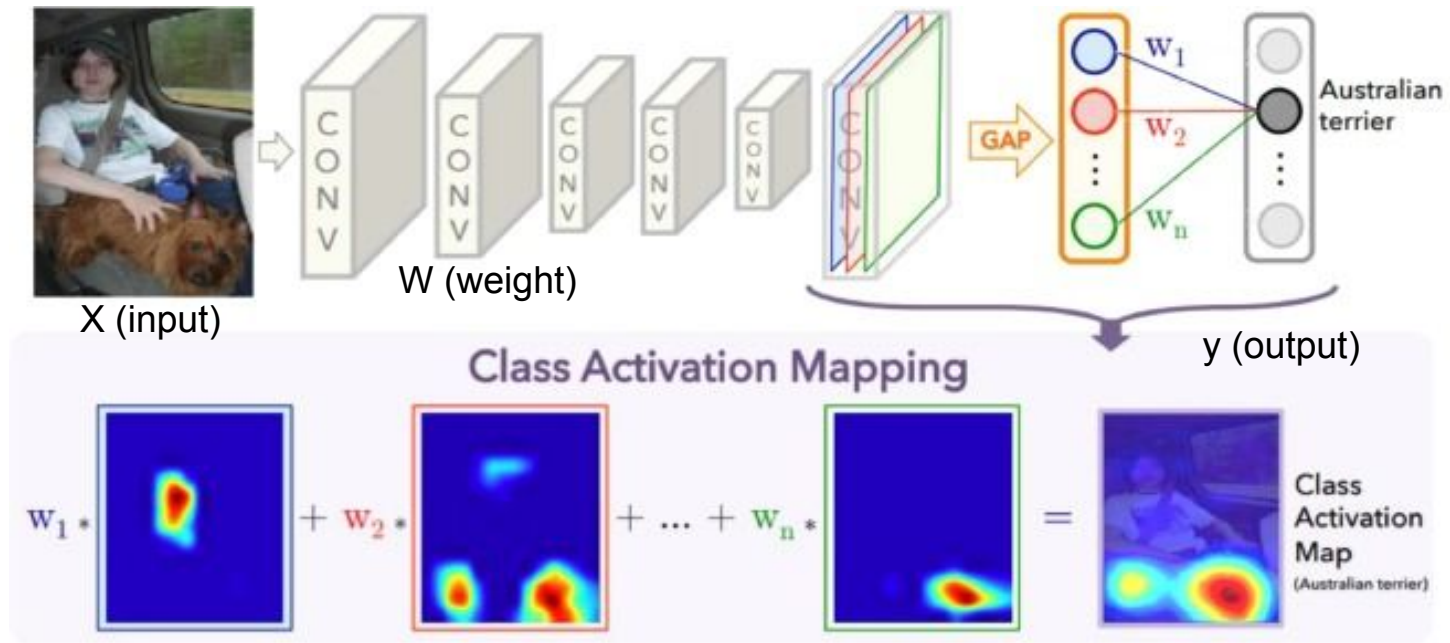
What is the Causality?

- Unidirectional

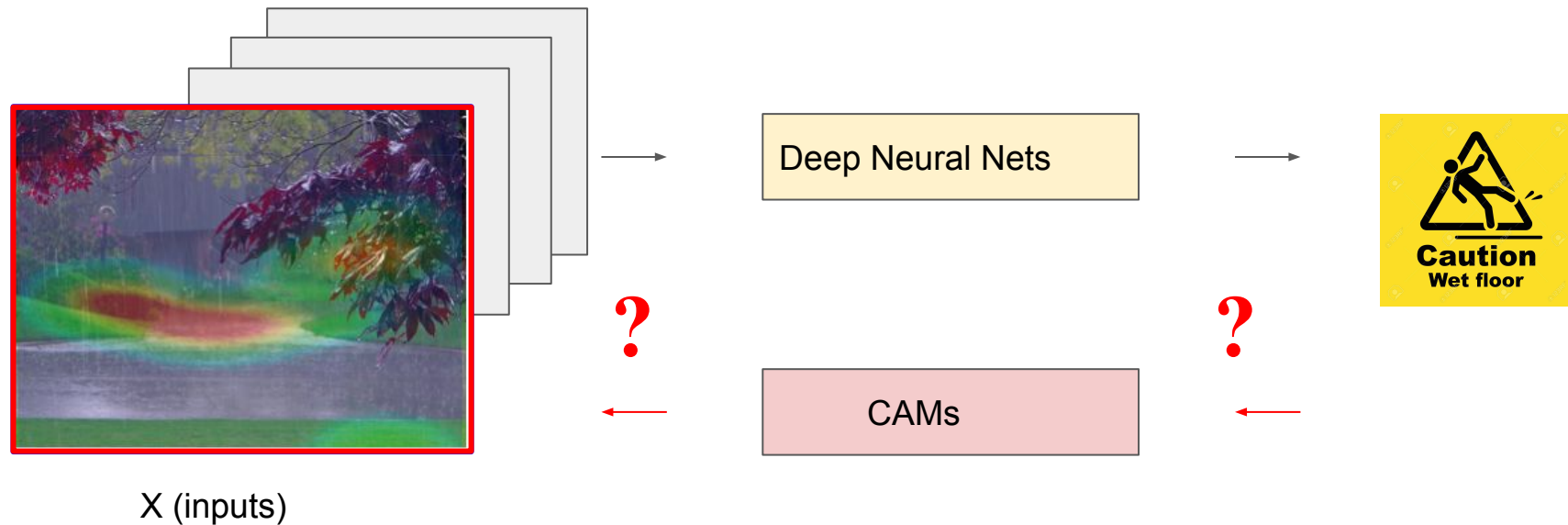
$P(\text{Causal Effect} \mid X, y)$



Related Work - Class-Activation Mapping (CAM)

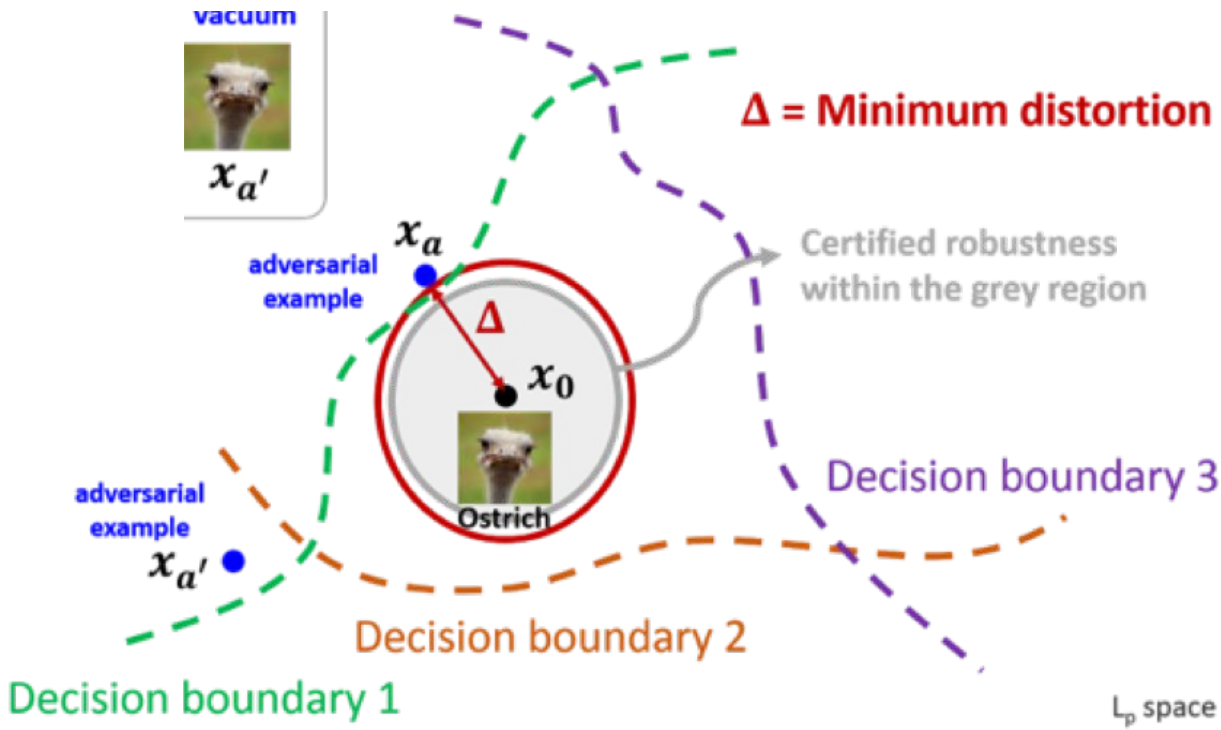


Correlation (CAM) v.s. Causality on DNNs



how to find a visual causal features on output label?

New DNN Challenge: Adversarial Examples



Proposed: Directed Causal Graph

Intervention : $do(x'_i)$

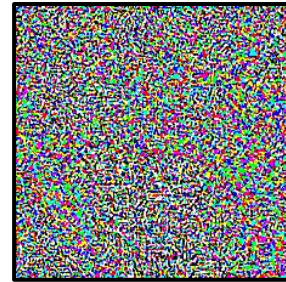
$$Effect(x_i \rightarrow x_j, Z) = P(x_j | do(x'_i), Z_{X_i}) - P(x_j | Z_{X_i})$$

Cat

Paw



Adversarial Noise

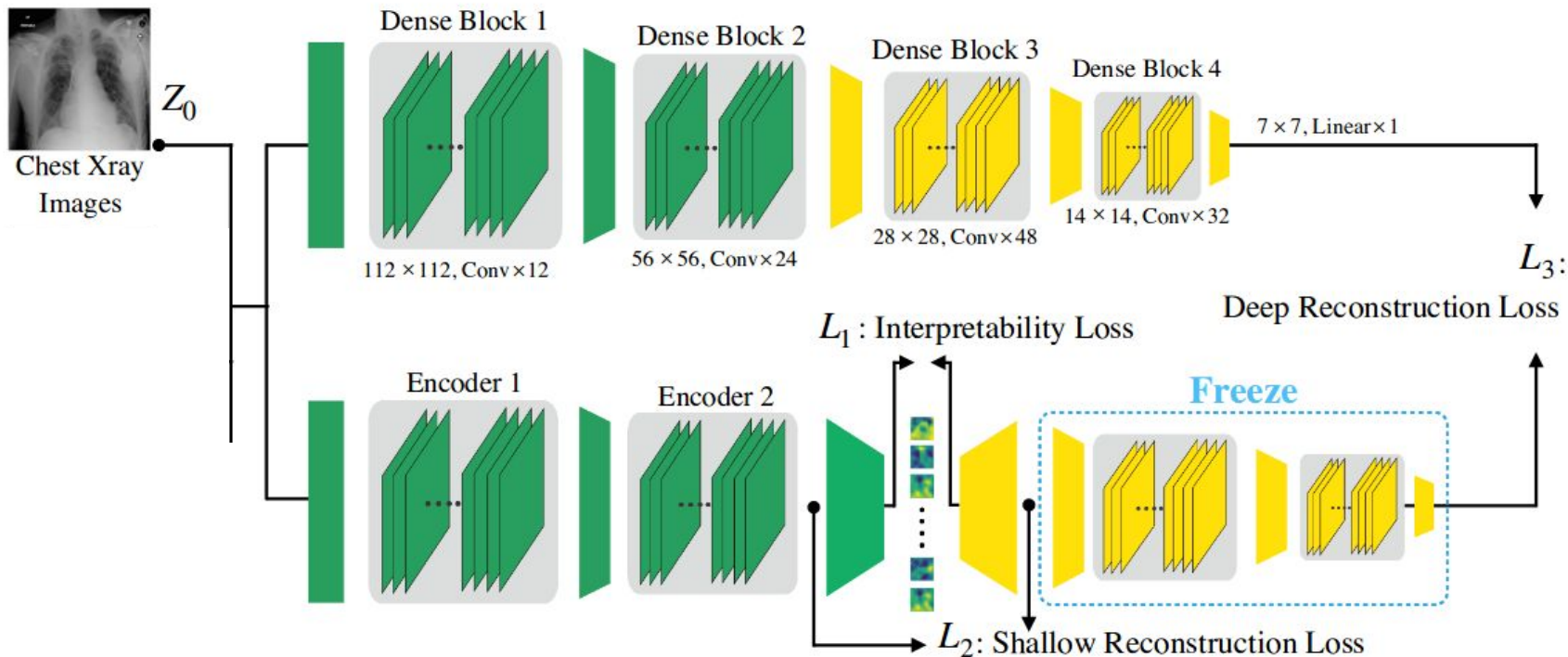


Contributions

- Deep Autoencoding for calculating **causal effect (CE)**.
- CE is a competitive **index** for understanding DNNs.
- We found that CE holds promises for detecting adversarial examples as it possesses distinct characteristics in the presence of adversarial perturbations.

Method: Deep Autoencoder

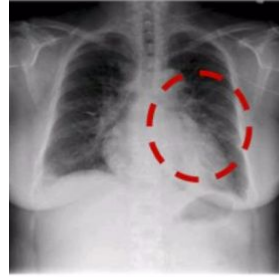
$$L(\theta; x_i) = \lambda_{shallow} \times L_{shallow}(\theta; x_i) + \lambda_{deep} \times L_{deep}(\theta; x_i) + \lambda_{interpretability} \times L_{interpretability}(\theta; x_i)$$



Datasets

ChestX-ray14

- contains frontal-view chest X-ray
- 14 different thoracic diseases.



Pneumonia

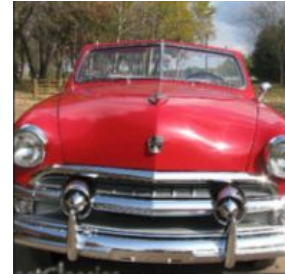
Fashion-MNIST

- 60k images

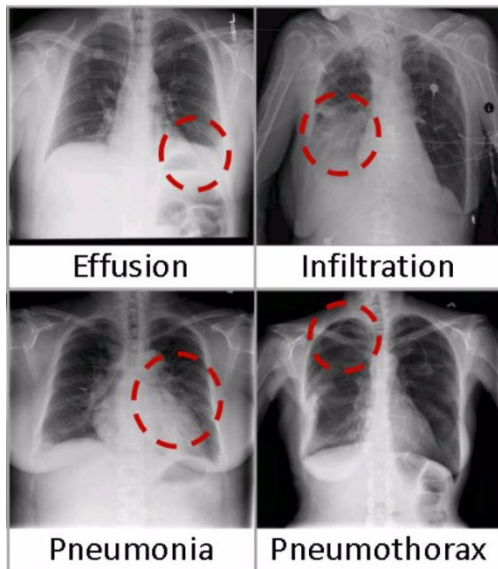


ImageNet

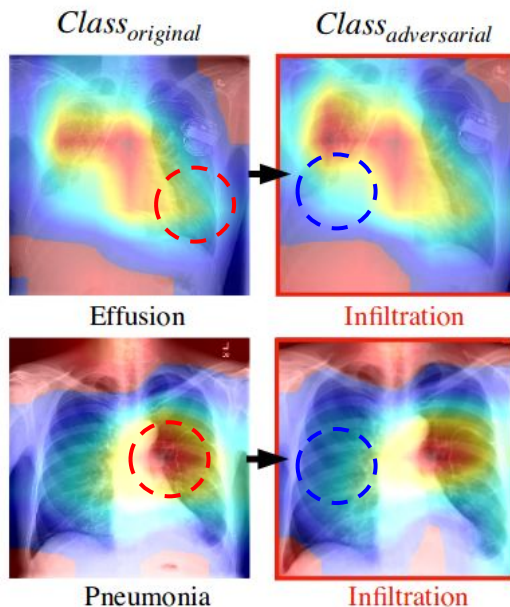
- 1.2 million images



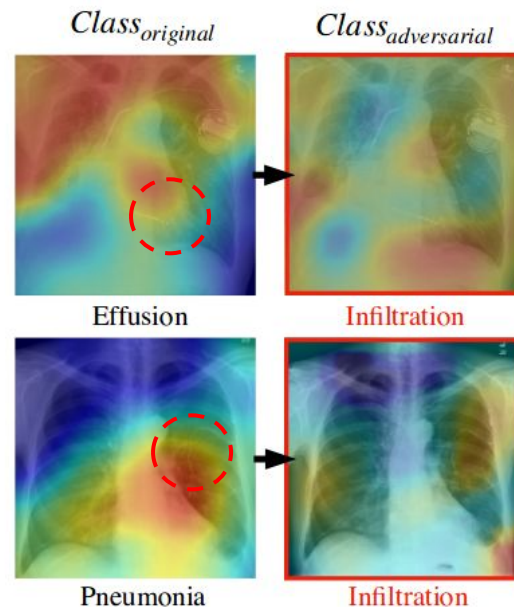
Results - Causal Effect Mapping (CEM)



Ground Truth



(I) CAM Results



(II) CEM Results

Estimate Causal Effect on Adversarial Signals



Table 1. Expected-CE on CheXNet

Level(L), Node(N)	Z_0	$F_i = \text{PWM}$
3,4	4.5076×10^{-3}	7.2356×10^{-7}
6,5	2.843×10^{-3}	1.2154×10^{-5}
6,10	3.1939×10^{-3}	9.0066×10^{-6}
8,5	3.1939×10^{-3}	1.1536×10^{-5}
10,7	1.3775×10^{-2}	-1.1506×10^{-5}

Estimate Causal Effect on Adversarial Signals

Attack Methods:

Fast Gradient Sign Method (FGSM)

Jacobian-Based Saliency Map (JBSM)

Basic Iterative Method (BIM)

Projected gradient descent (PGD)

Table 2. CheXNet (F_j on $L = 10, N = 7$)

$F_j =$ Types of Adversarial Attack	Expected-CE
FGSM	-5.6129×10^{-6}
BIM	4.3435×10^{-5}
JBSM	7.7548×10^{-5}
PGD	-3.9605×10^{-6}

Conclusion

A Framework to estimate a causal effect on high dimensional visual data

Evaluate the this numerical causal effect on adversarial example

Saliency visualization as a interperable method.

Future work:

Relation reasoning, time-series causal analysis, and video detection

Thank You!

Question & Answer

Code Released:

<https://github.com/jjaacckkyy63/Causal-Intervention-AE-wAdvImg>

yliu3233@gatech.com ; huckiyang@gatech.edu

Appendix: Expectation of Causal Effect

$X_j = x_j$ with all of the evidence Z could be computed as:

$$Effect(x_i \rightarrow x_j, Z) = P(x_j | do(x'_i), Z_{X_i}) - P(x_j | Z_{X_i}) \quad (1)$$

The expected casual effect from **Eqn. 6** in [18] has been defined as:

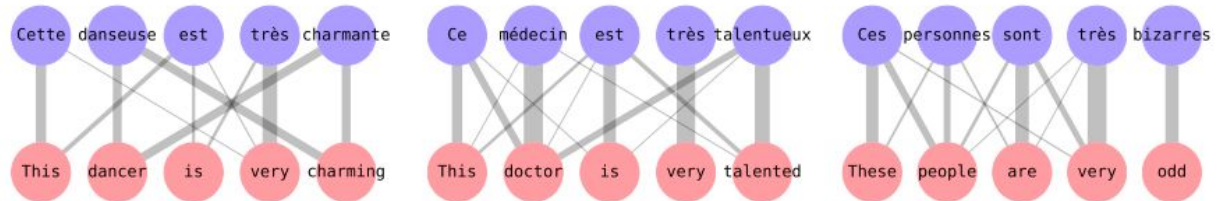
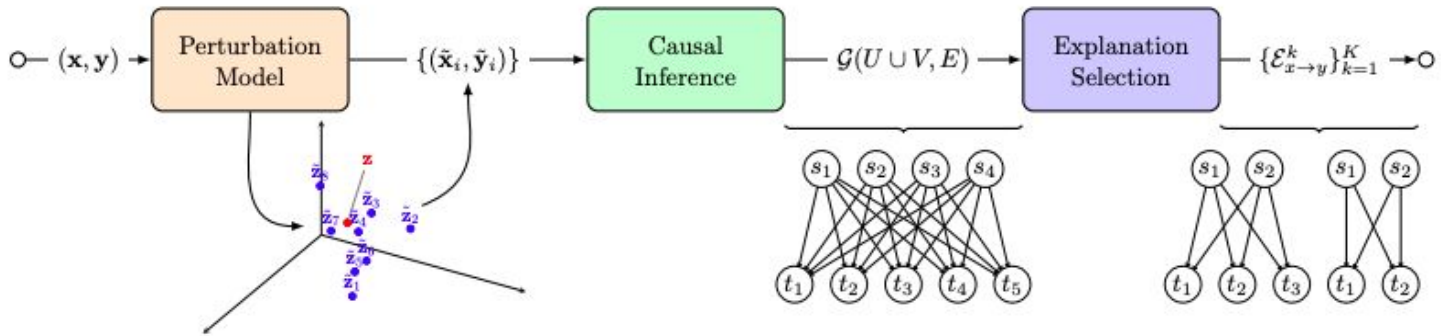
$$E_{X_i}[Effect(x_i \rightarrow x_j, Z)] = \sum_{x_i \in X_i} P(X_i = x_i | Z) \times (1) \quad (2)$$

$$P(x_i | pa'_i) = \begin{cases} P(x_i | pa_i) & \text{if } F_1 = \text{idle,} \\ 0 & \text{if } F_i = do(x'_i) \text{ and } x_i \neq x'_i, \\ 1 & \text{if } F_i = do(x'_i) \text{ and } x_i = x'_i. \end{cases} \quad (3)$$

$$L(\theta; x_i) = \lambda_{shallow} \times L_{shallow}(\theta; x_i) + \lambda_{deep} \times L_{deep}(\theta; x_i) + \lambda_{interpretability} \times L_{interpretability}(\theta; x_i) \quad (4)$$

Appendix - NLP

A causal framework for explaining the predictions of black-box sequence-to-sequence models,
 David Alvarez-Melis, Tommi S. Jaakkola, ACL, 2017



Explanations for biased translations of similar gender-neutral English sentences into French.

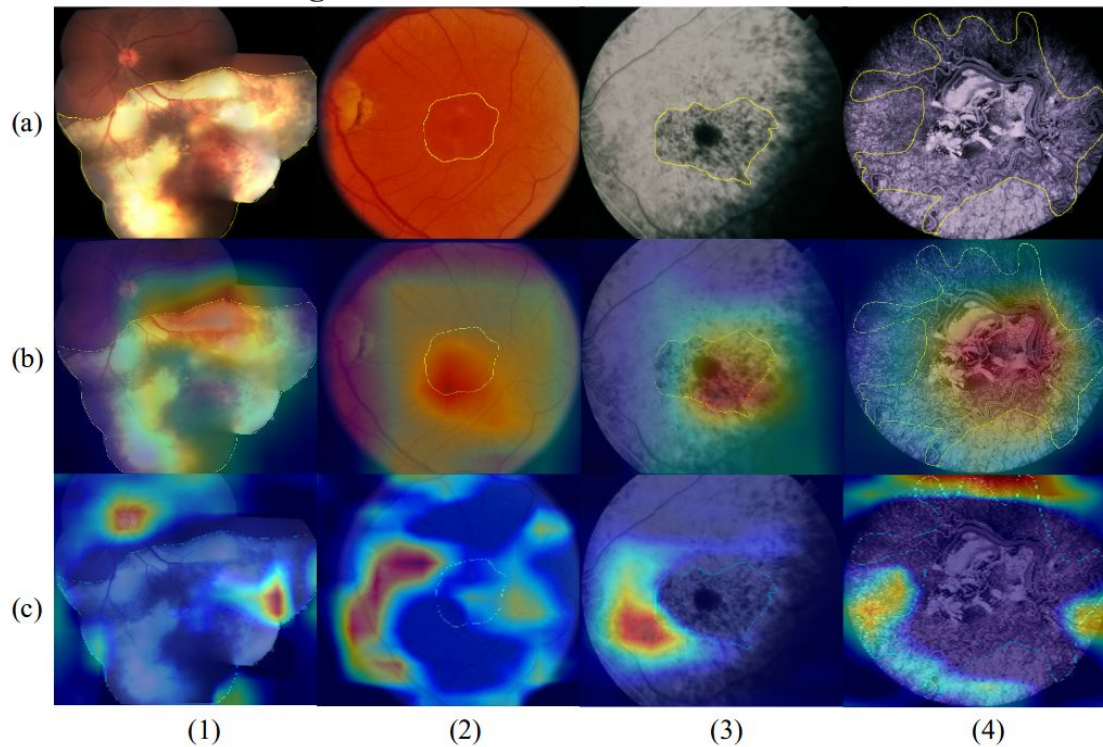
Similarity

- Biscuit v.s. Dog (Chihuahua)



Medical Image

- Retinal Images



Medical Image

Fast Gradient Sign Method (FGSM)

x
"panda"
57.7% confidence

$+ .007 \times$

$\text{sign}(\nabla_x J(\theta, x, y))$
"nematode"
8.2% confidence

$=$

$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
"gibbon"
99.3 % confidence

- x – Clean Input Image
- x^{adv} – Adversarial Image
- J – Loss Function

- y_{true} – Model Output for x
- ϵ – Tunable Parameter