# Contents and speakers

**Overview of trustworthiness** (Jindong Wang, 10min)

**Robust machine learning**
(Jindong Wang, 40min)

**Out-of-distribution generalization**
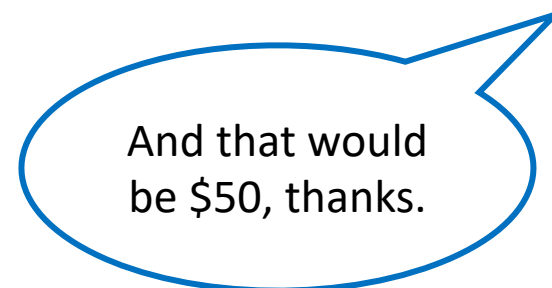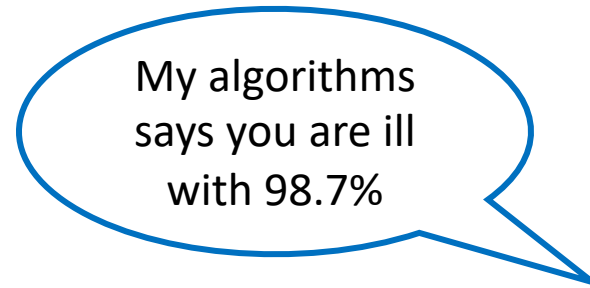(Haohan Wang, 40min)

Interpretability
(Haohan Wang,
on behalf of Haoliang Li,
40min)

**Trustworthiness in the era of large models** (Jindong Wang, 40min)

# The Importance of Interpretability From an application perspective

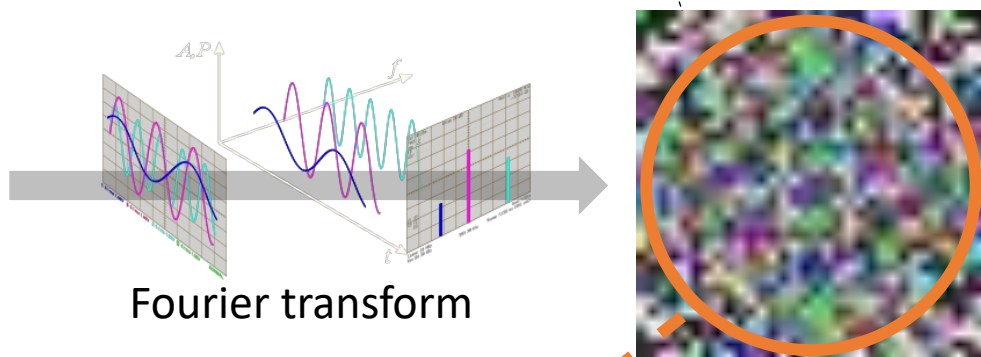- In many applications, we need to know why

# High/Low-frequency reconstructed Images

Each one of the RGB channels of the raw images are transformed independently, visualized as an image of three channels here. Only the real part is visualized. Both real and imaginary parts are used in the experiments.
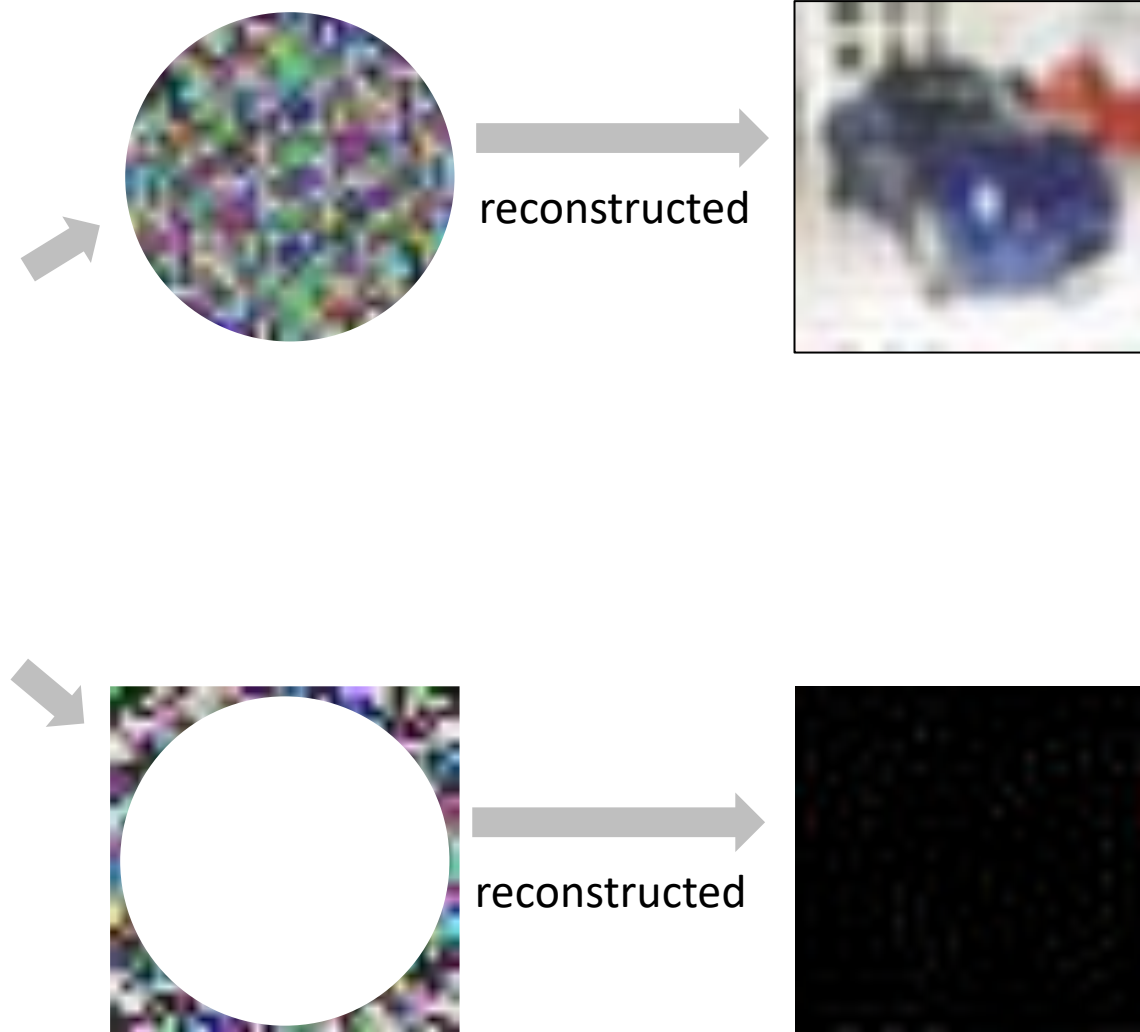
reconstructed

original image

Fourier transform

frequency domain

a predetermined radius

reconstructed

# Misalignment between human and model

# Additional examples

original image

low-frequency

high-frequency

rescaled high-frequency

prediction

frequency domain

# Motivation: Interpretability

- Adversarial Robustness seems to remind us about one in applying deep learning in practice:



**Original image**

Dermatoscopic image of a benign melanocytic nevus, along with the diagnostic probability computed by a deep neural network.

| Benign
| Malignant

+ 0.04 ×

**Adversarial noise**

Perturbation computed by a common adversarial attack technique.
See (7) for details.

=

**Adversarial example**

Combined image of nevus and attack perturbation and the diagnostic probabilities from the same deep neural network.

| Benign
| Malignant

$x$
"panda"
57.7% confidence

$+ .007 \times$

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

=

$\boldsymbol{x} + \epsilon\,\mathrm{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$
"gibbon"
99.3 % confidence

- When a model is making prediction, the features it use can be quite arbitrary
  - So, we probably need to have a detailed look at the decision process

# Backprop and Guided Backprop

- Directly using the gradient



- Include both positive values and negative values

# Backprop and Guided Backprop

- Maybe we don't care that much for the features that contribute negatively
    - How do we get rid of it?
        - Maybe just setting it to zeros
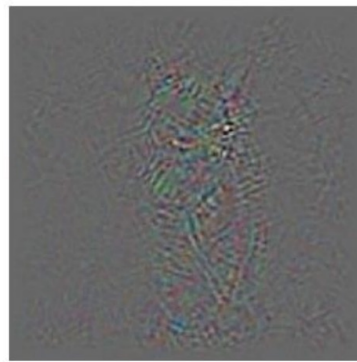
- Guided Backprop



gradient                              only positive gradient

# Class Activation Map

- Visualizing the decision for convolutional neural networks



Class Activation Mapping

$w_1 *$ + $w_2 *$ + ... + $w_n *$ = Class Activation Map (Australian terrier)

- As the beginning work, the model has to follow certain structures

# GradCam

- Connecting later layers to the convolutional layers with gradient

# GradCam

- Results Comparison



(a) Original Image  (b) Guided Backprop 'Cat'  (c) Grad-CAM 'Cat'

(g) Original Image  (h) Guided Backprop 'Dog'  (i) Grad-CAM 'Dog'

# Model's understanding of certain labels

- Use to evaluate the model's understanding of the whole class
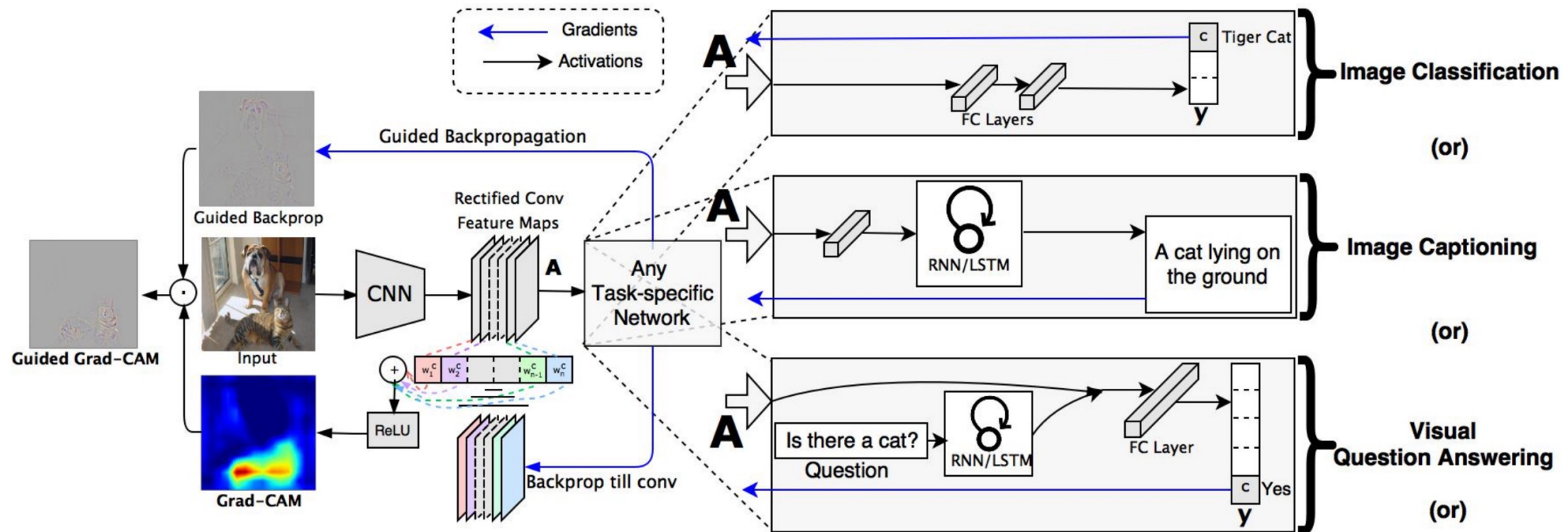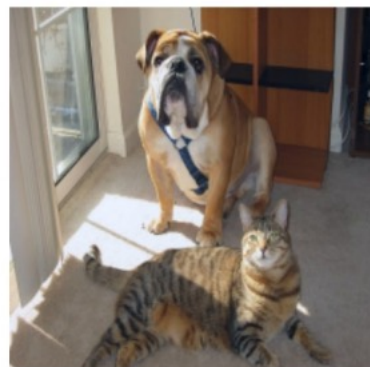
  - Given a (random) starting point

  - Update the samples following the gradient

  - Until the model predicts the resulting images with full confidence

Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps



dumbbell    cup    dalmatian

bell pepper    lemon    husky

washing machine    computer keyboard    kit fox

goose    ostrich    limousine

# Explaining by Removing

- Explaining by Removing: A Unified Framework for Model Explanation
  - Many different papers in the community for interpretability

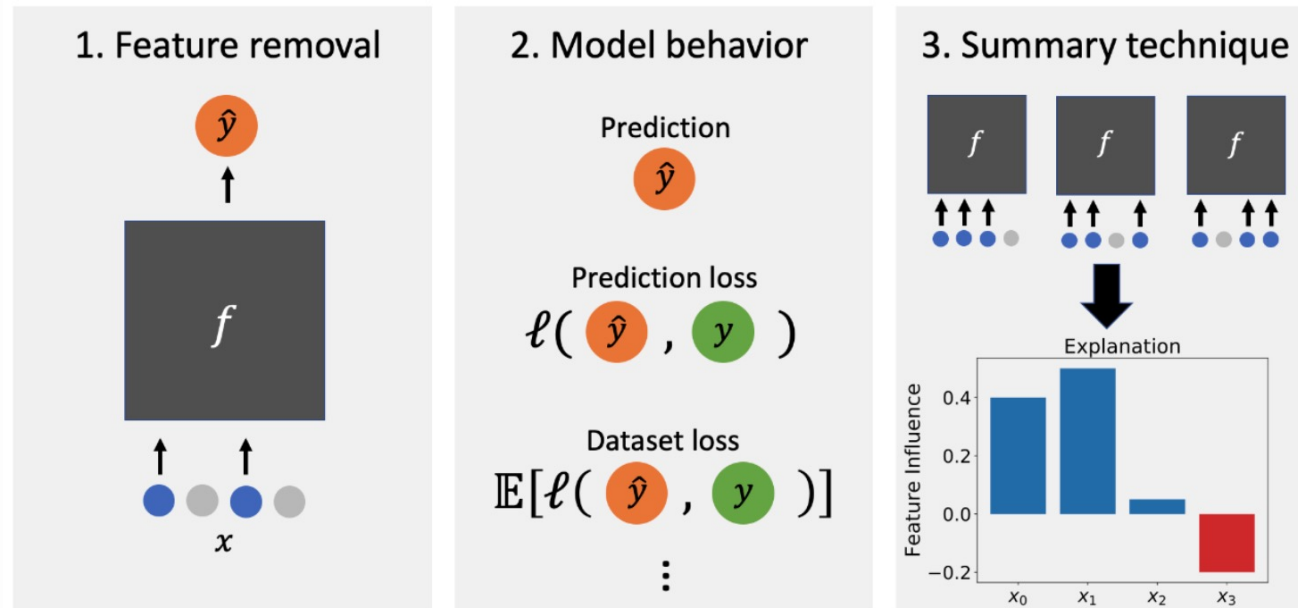| Method | Removal | Behavior | Summary |
|---|---|---|---|
| IME (2009) | Separate models | Prediction | Shapley value |
| IME (2010) | Marginalize (uniform) | Prediction | Shapley value |
| QII | Marginalize (marginals product) | Prediction | Shapley value |
| SHAP | Marginalize (conditional/marginal) | Prediction | Shapley value |
| KernelSHAP | Marginalize (marginal) | Prediction | Shapley value |
| TreeSHAP | Tree distribution | Prediction | Shapley value |
| LossSHAP | Marginalize (conditional) | Prediction loss | Shapley value |
| SAGE | Marginalize (conditional) | Dataset loss (label) | Shapley value |
| Shapley Net Effects | Separate models | Dataset loss (label) | Shapley value |
| Shapley Effects | Marginalize (conditional) | Dataset loss (output) | Shapley value |
| Permutation Test | Marginalize (marginal) | Dataset loss (label) | Remove individual |
| Conditional Perm. Test | Marginalize (conditional) | Dataset loss (label) | Remove individual |
| Feature Ablation (LOCO) | Separate models | Dataset loss (label) | Remove individual |
| Univariate Predictors | Separate models | Dataset loss (label) | Include individual |
| L2X | Missingness during training | Prediction mean loss | High-value subset |
| INVASE | Missingness during training | Prediction mean loss | High-value subset |
| LIME (Images) | Default values | Prediction | Linear model |
| LIME (Tabular) | Marginalize (replacement dist.) | Prediction | Linear model |
| PredDiff | Marginalize (conditional) | Prediction | Remove individual |
| Occlusion | Zeros | Prediction | Remove individual |
| CXPlain | Zeros | Prediction loss | Remove individual |
| RISE | Zeros | Prediction | Mean when included |
| MM | Default values | Prediction | Partitioned subsets |
| MIR | Extend pixel values | Prediction | High-value subset |
| MP | Blurring | Prediction | Low-value subset |
| EP | Blurring | Prediction | High-value subset |
| FIDO-CA | Generative model | Prediction | High-value subset |

# Explaining by Removing

- The methods typically have three core steps
  - How/what to remove the features
  - What to observe
  - Summarizing the observation to users

# Explaining by Removing

- Feature Removal

  - Zero-ing features: $F(x, S) = f(x_S, 0)$
  - Setting features to a default value $r$: $F(x, S) = f(x_S, r_{\bar{S}})$
  - Sampling from a conditional generative model $\sim p_G(X_{\bar{S}}|X_S)$: $F(x, S) = f(x_S, \tilde{x}_{\bar{S}})$
  - Marginalizing with condition: $F(x, S) = \mathbb{E}[f(x)|X_S = x_S]$
  - Marginalizing with marginal: $F(x, S) = \mathbb{E}[f(x_S, X_{\bar{S}})]$

# Model Behavior

- What to observe after removing the features

  - at the prediction level *(local explanations)*: Given an input $x \in \mathcal{X}$, study $F(x, S)$, that is how removed features are impacting a prediction higher or lower;
  - at the prediction loss level *(local explanations)*: Given an input $x$ and its true label $y$, study $-\ell(F(x, S), y)$, that is how some features are making the prediction more or less correct.
  - the average prediction loss *(local explanations)*: Given an input $x$ and the label's conditional distribution $p(Y|X = x)$, study $-\mathbb{E}_{p(Y|X=x)}[\ell(F(x, S), Y)]$, that is how a certain set of features can correctly predict what could have occured on average. Can be useful with uncertain labels.
  - dataset loss wrt label *(global explanations)*: How much the model's performance degrades when different features are removed, i.e. $-\mathbb{E}_{XY}[\ell(F(X_S), Y)]$
  - dataset loss wrt output *(global explanations)*: What are the features' influence on the model output (rather than on the model performance), i.e. $-\mathbb{E}_X[\ell(F(X_S), F(X))]$

# Summary Techniques

- To provide a concise summary of the information we obtained
  - Examples:
    - Feature attribution:
      - Give every feature a score

    - Feature section
      - Select a subset of features

> So, now, we should possess the knowledge of hundreds of model explanation methods

# Explaining by Removing

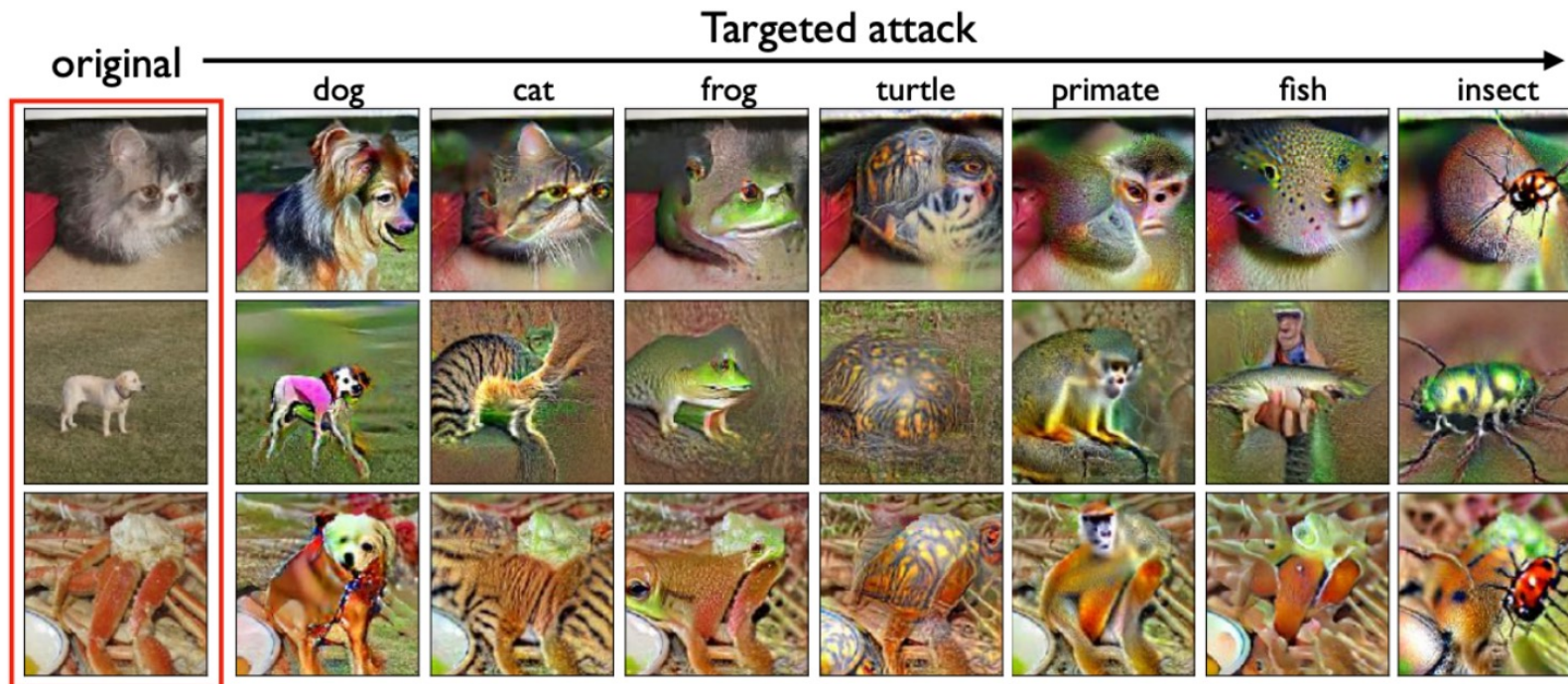- Summary of the techniques by "Explaining by Removing" at 2020



Just in case you want to try a submission in this field but not sure about what to do :D

# An Interesting Question

Do we need better interpretability methods or better model?
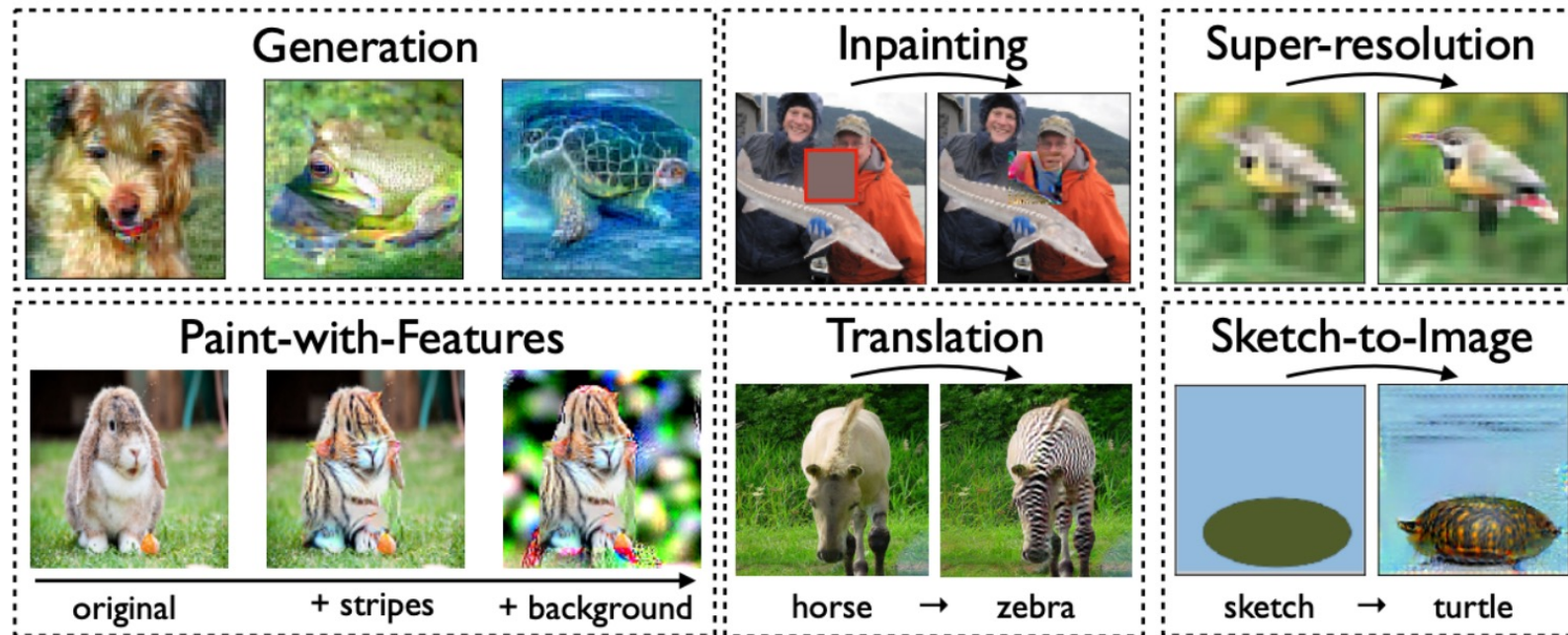
# Adversarially Robust Models

- Adversarially robust model might learn perceptually aligned representations
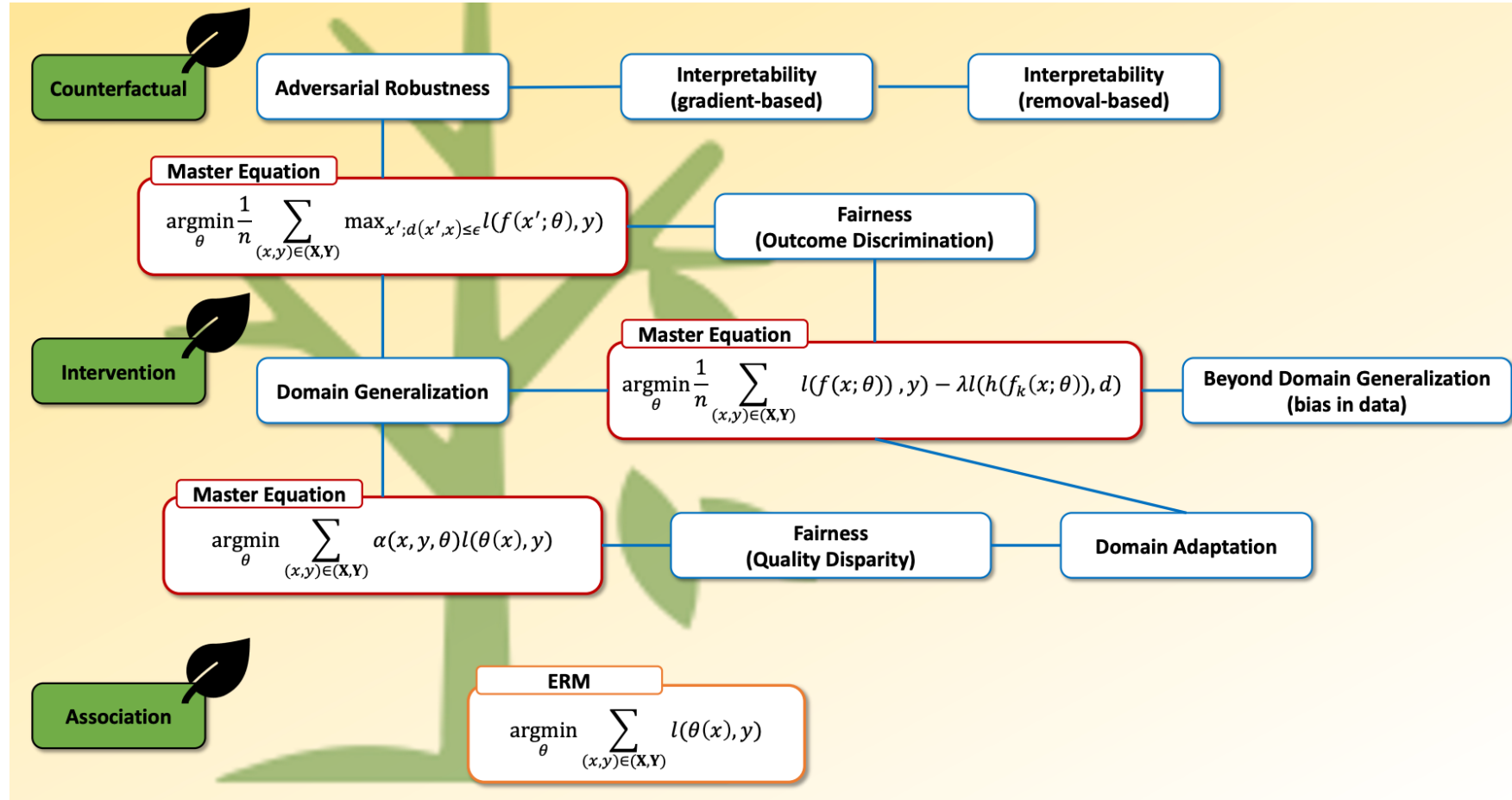  - (Santurkar et al. 2019)

# Adversarially Robust Models

- Adversarially robust model might learn perceptually aligned representations
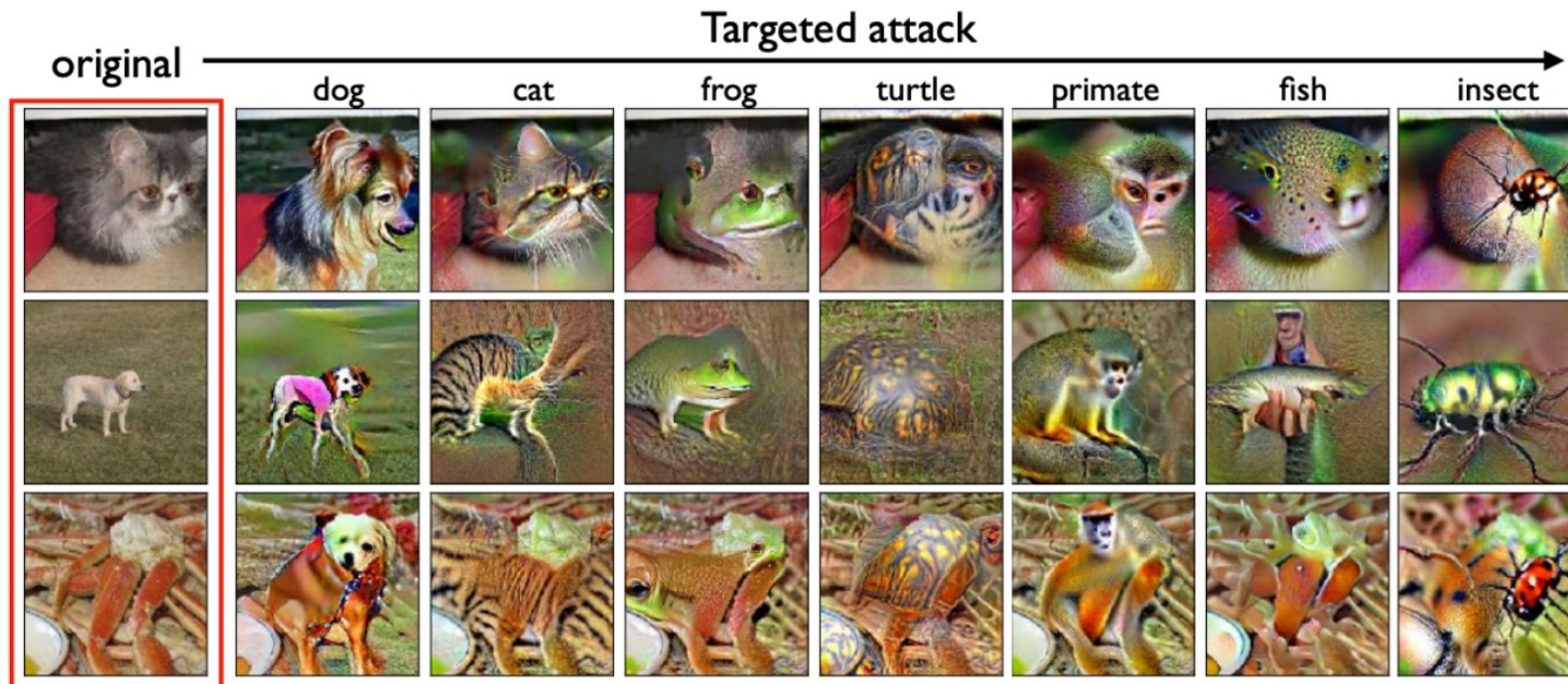  - (Santurkar et al. 2019)

# Potential connection between interpretability and adversarial robustness

# Adversarially Robust Models

- With this connection, let's look at these results again
  - It seems these results are just supposed to happen

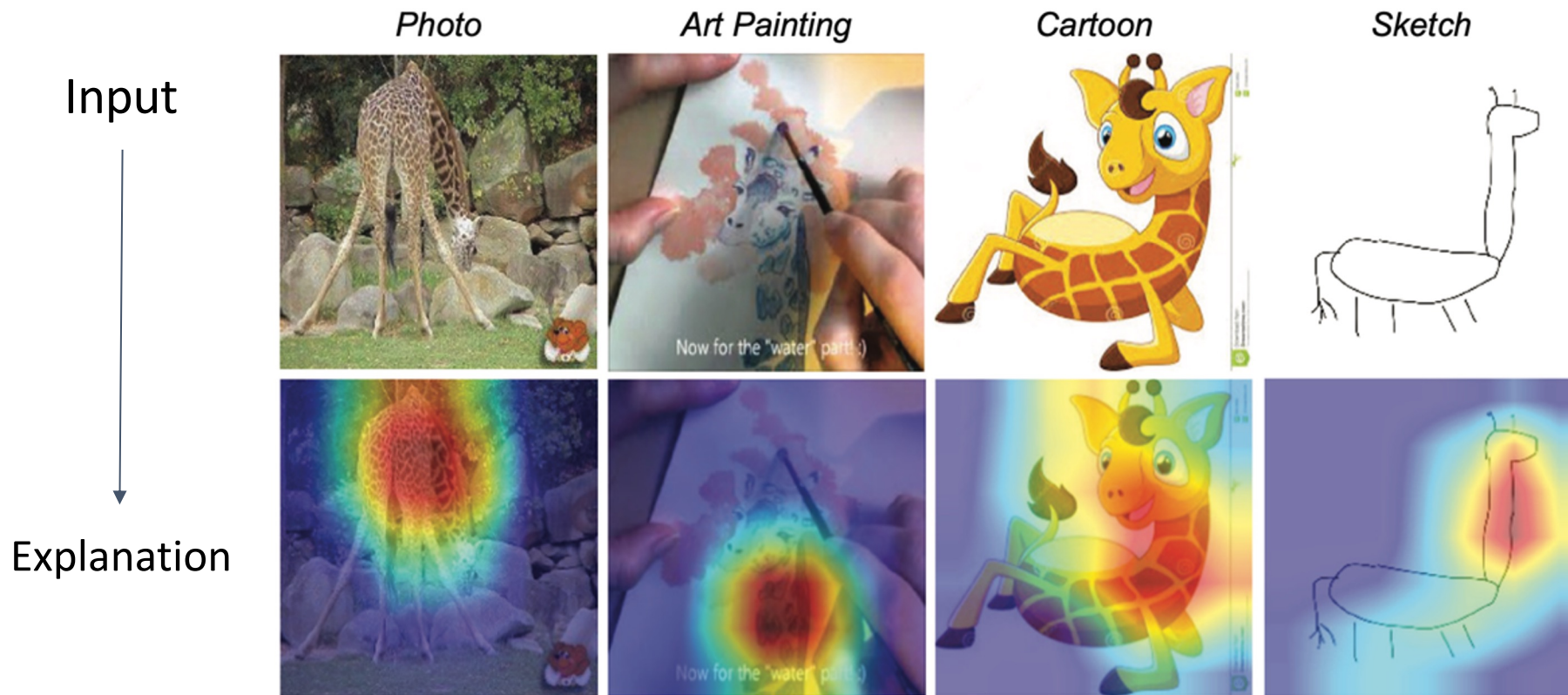**Following Slides are from Haoliang Li**

# Interpretability via attention explanations

Shared causality across domains.

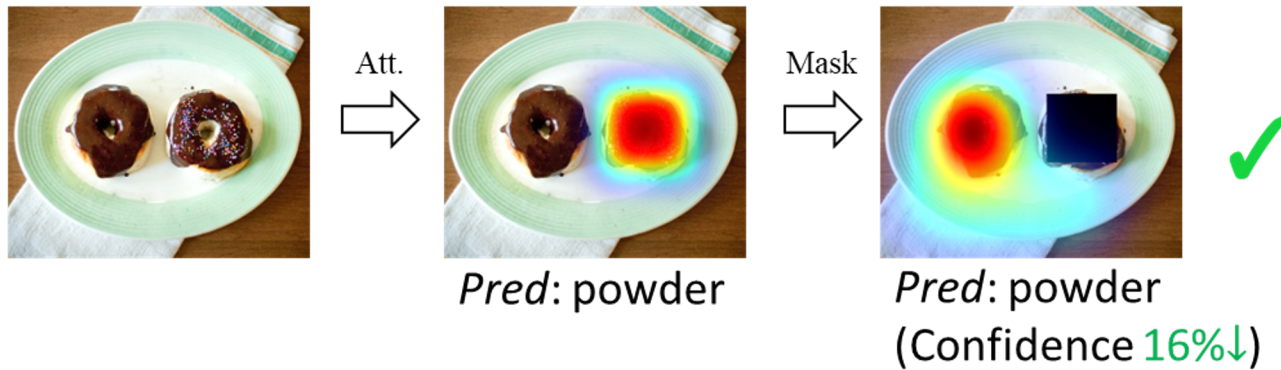But can we believe such appealing visualizations?



Input

Explanation

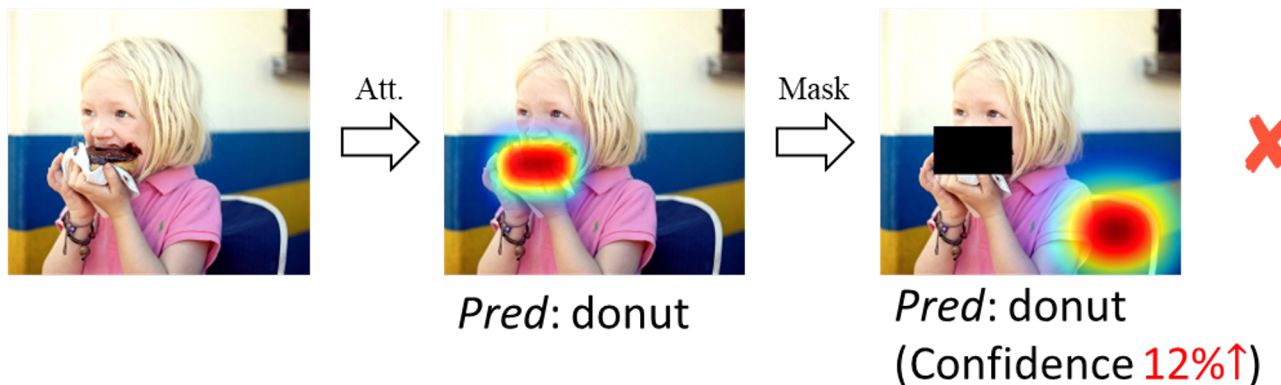- Kim et al. SelfReg: Self-supervised Contrastive Regularization for Domain Generalization. ICCV 2021.

# Interpretability via attention explanations

Can existing explanation methods faithfully represent model decisions?



Question: What are colorful pieces on the doughnut?

Pred: powder

Pred: powder
(Confidence 16%↓) ✔

Question: What is the girl eating?

Pred: donut

Pred: donut
(Confidence 12%↑) ✘

No! Don't be misled by appealing visualizations!

· Liu et.al. Rethinking Attention-Model Explainability through Faithfulness Violation Test, ICML 2022.
https://arxiv.org/abs/2201.12114

# Interpretability via attention explanations

The reason is simple:

➢ Recall the formulation of attention mechansims

$$A = \text{softmax}(\frac{Q \cdot K^T}{\sqrt{d_h}})$$
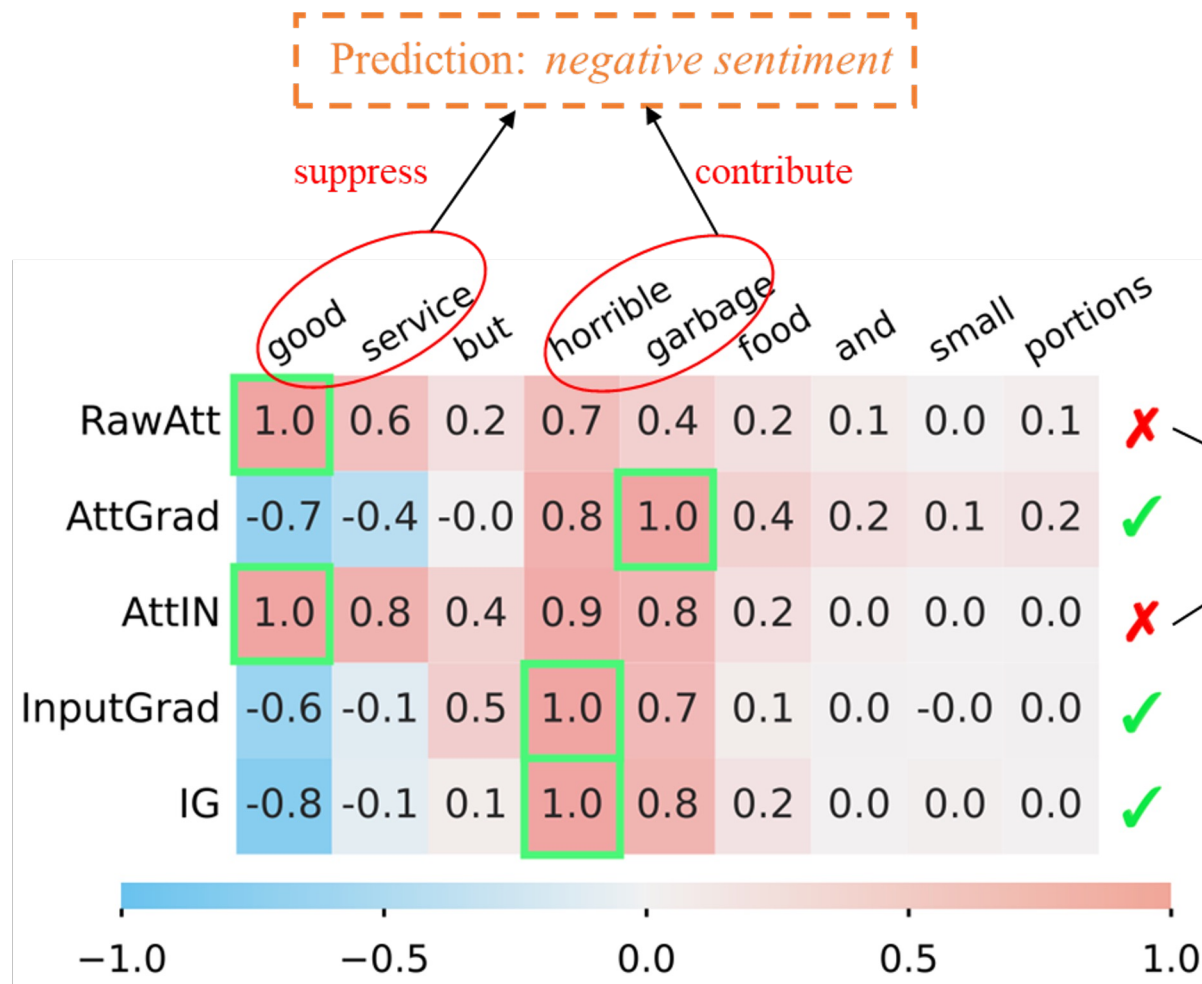
$$O = A \cdot V$$

(Vaswani et al., 2017)

1. Attention weights are always non-negative since they are the output of softmax function.

2. So they cannot **differentiate the direction** of feature impacts, *i.e.*, impact polarity

# Interpretability via attention explanations

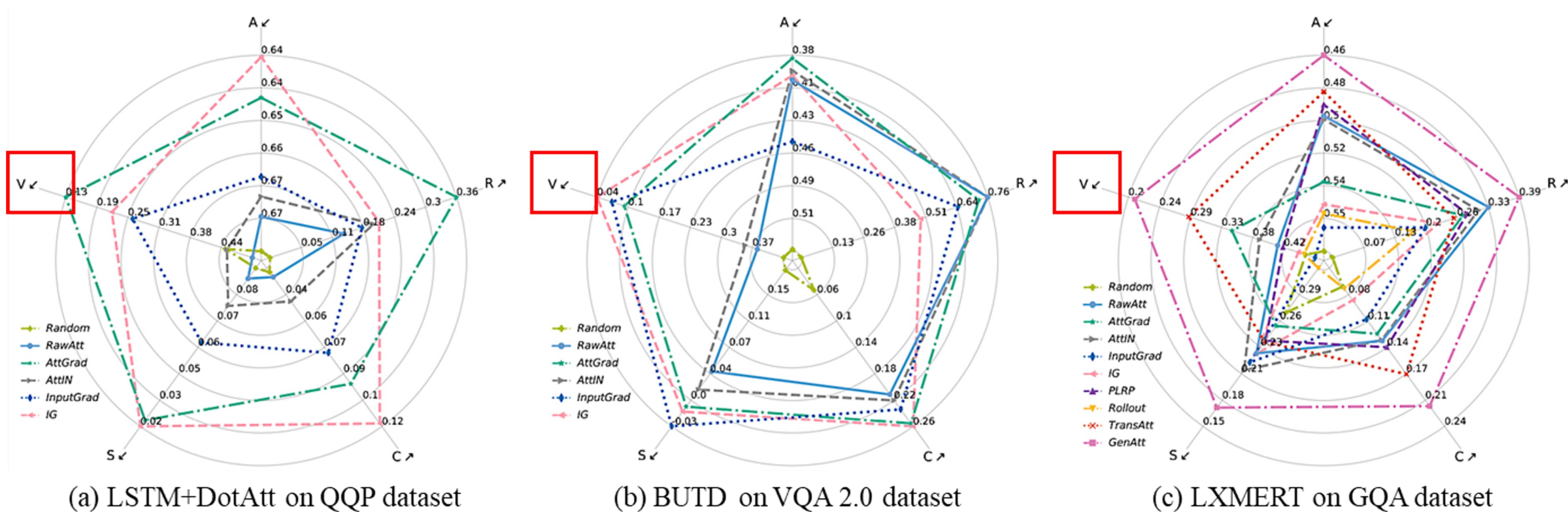This heavily degrades the faithfulness of attention explanations:



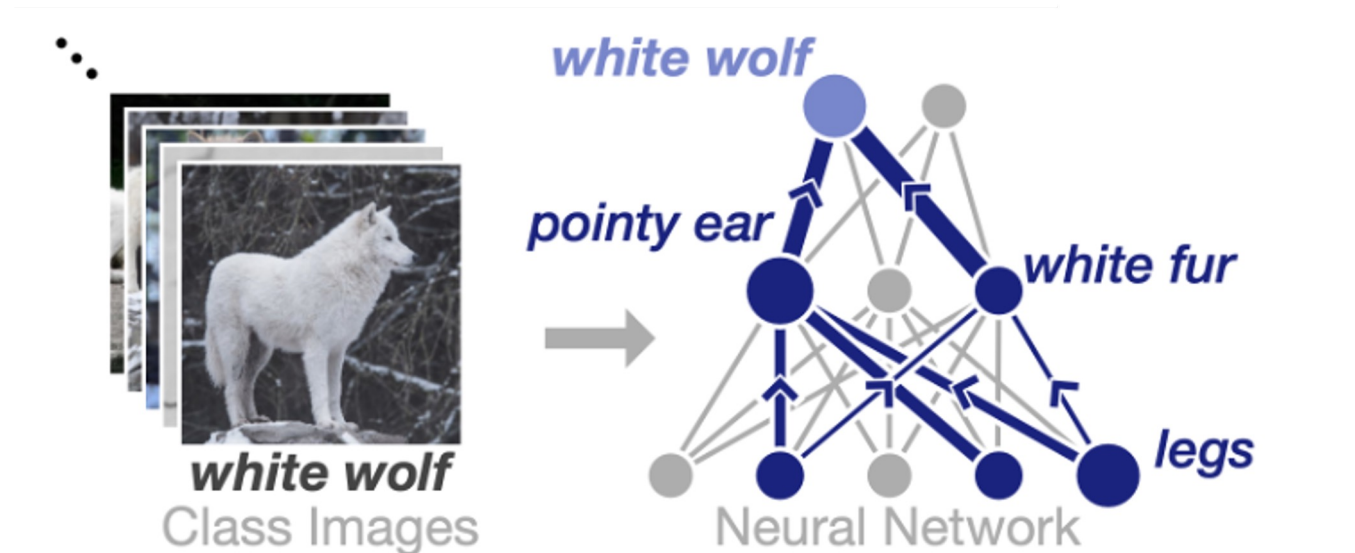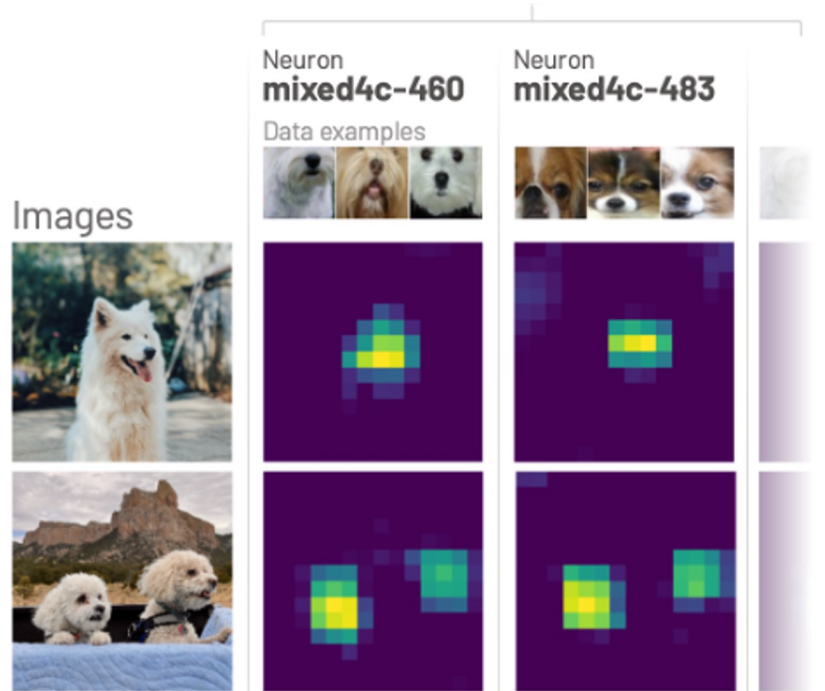Because they cannot differentiate the impact polarity!!

# Interpretability via attention explanations

Most explanation methods fail to suffer from the polarity consistency issue.



(a) LSTM+DotAtt on QQP dataset

(b) BUTD on VQA 2.0 dataset

(c) LXMERT on GQA dataset

# Interpretability via neuron explanations

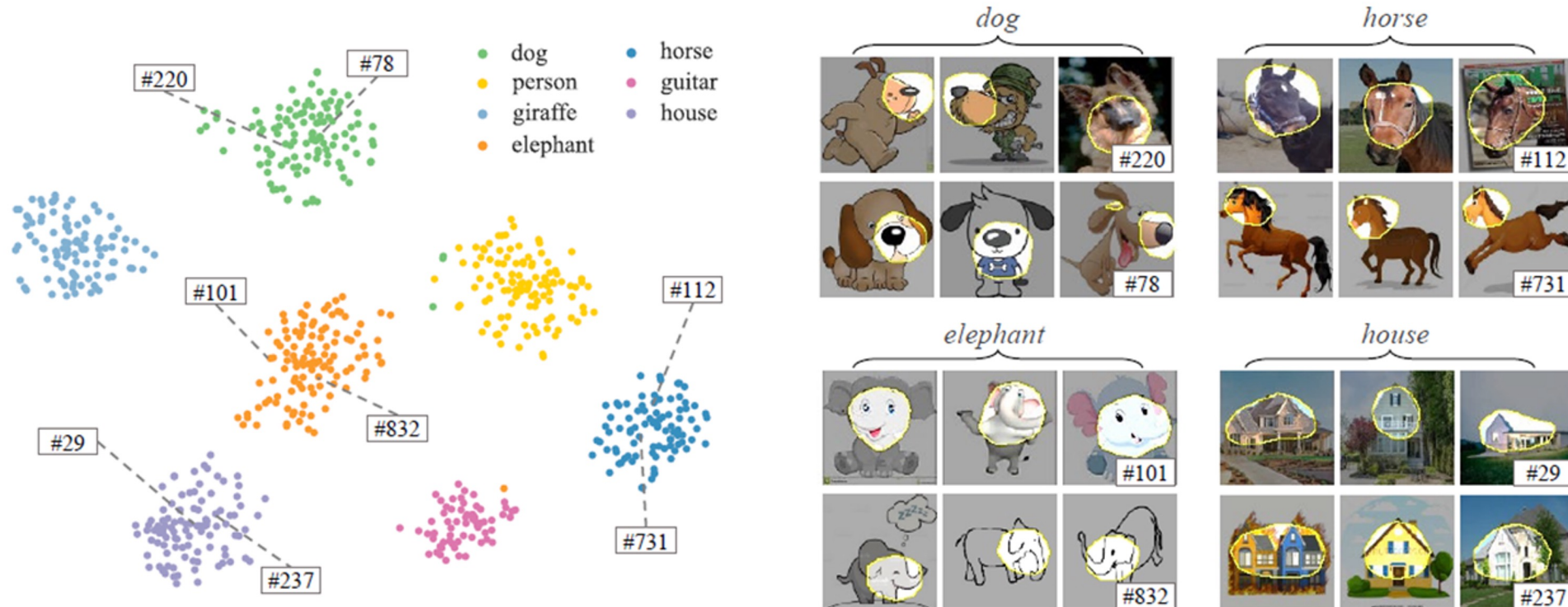What did a single unit (neuron) learn?

- H. Park, et al. "NeuroCartography: Scalable Automatic Visual Summarization of Concepts in Deep Neural Networks ," TVCG, 2021.

- F. Hohman, et al. "SUMMIT: Scaling Deep Learning Interpretability by Visualizing Activation and Attribution Summarizations ," TVCG, 2019.

# Interpretability via neuron explanations

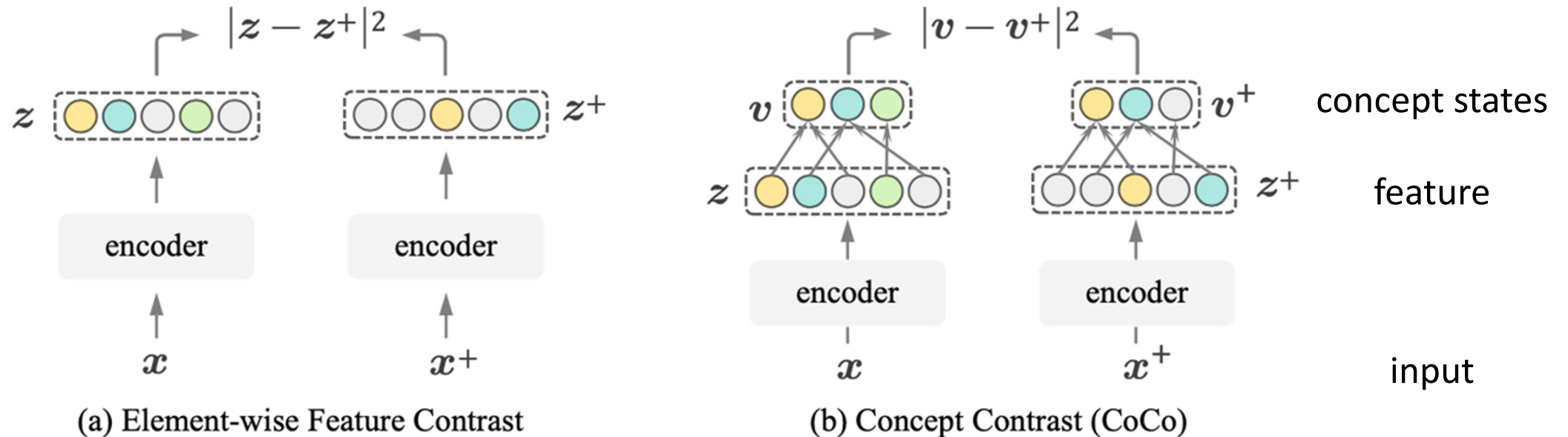Take the last convolutional layer in ResNet as an example,



(a) Top activated neurons under different predictions.

(b) Visual concepts shared by different neurons.

· Liu, Y., Tian, C. X., Li, H., & Wang, S. (2022). Generalization Beyond Feature Alignment: Concept Activation–Guided Contrastive Learning. arXiv preprint arXiv:2211.06843.

# Interpretability via neuron explanations

We proposed concept-level contrast (CoCo) to learn features beyond conventional feature-level contrast.
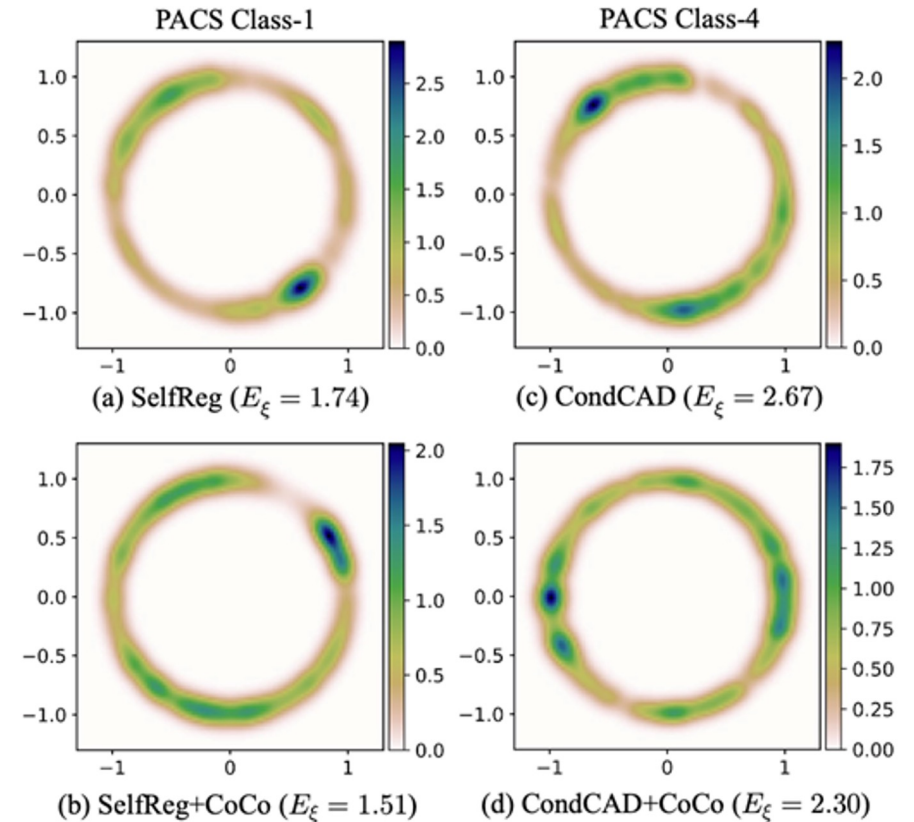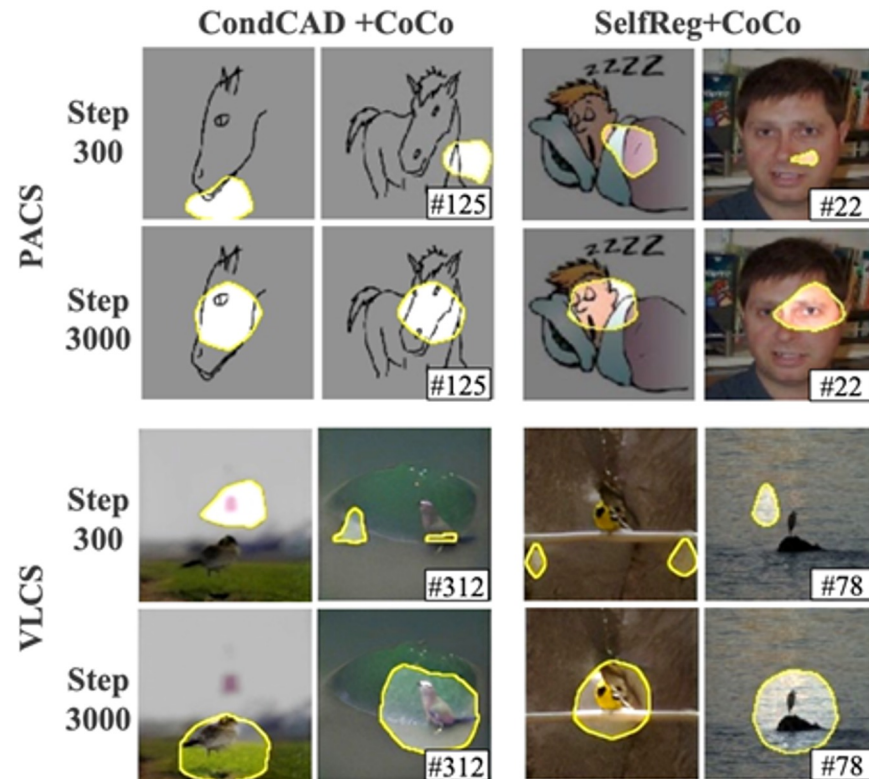


(a) Element-wise Feature Contrast

(b) Concept Contrast (CoCo)

· Liu, Y., Tian, C. X., Li, H., & Wang, S. (2022). Generalization Beyond Feature Alignment: Concept Activation–Guided Contrastive Learning. arXiv preprint arXiv:2211.06843.
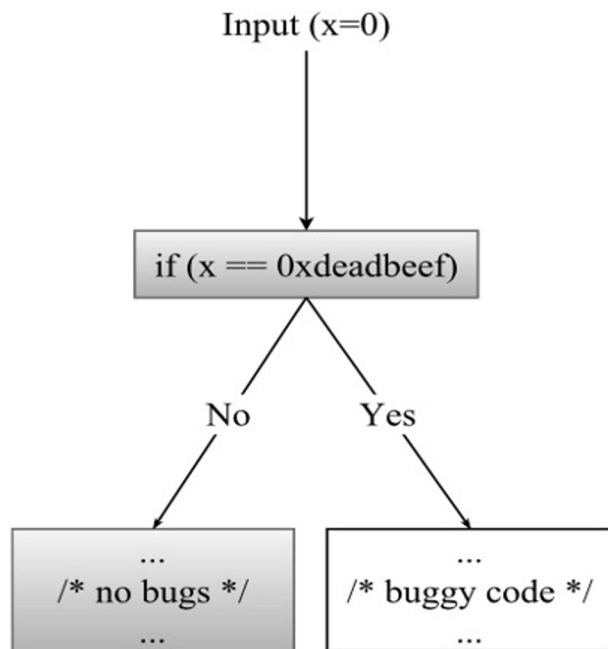
# Interpretability via neuron explanations

With CoCo, concept evolution happens

Diversified features



(a) SelfReg ($E_\xi = 1.74$)

(c) CondCAD ($E_\xi = 2.67$)

(b) SelfReg+CoCo ($E_\xi = 1.51$)

(d) CondCAD+CoCo ($E_\xi = 2.30$)

# Interpretability via neuron explanations

Code / Logic not covered during testing -> Bugs may hide there

This could also happen in neural networks!

Input (x=0)

if (x == 0xdeadbeef)

No          Yes

...          ...
/* no bugs */    /* buggy code */
...          ...

Neurons are not covered? ⟶

Input

| Blue 0 | Red 2.8 | ... | VEdge 1.1 | HEdge 1.6 |

| Nose 0 | ... | Wheel 2.4 |

| Car 0.95 | ... | Face 0 |

DeepXplore: Automated Whitebox Testing of Deep Learning Systems. Best paper in SOSP'17
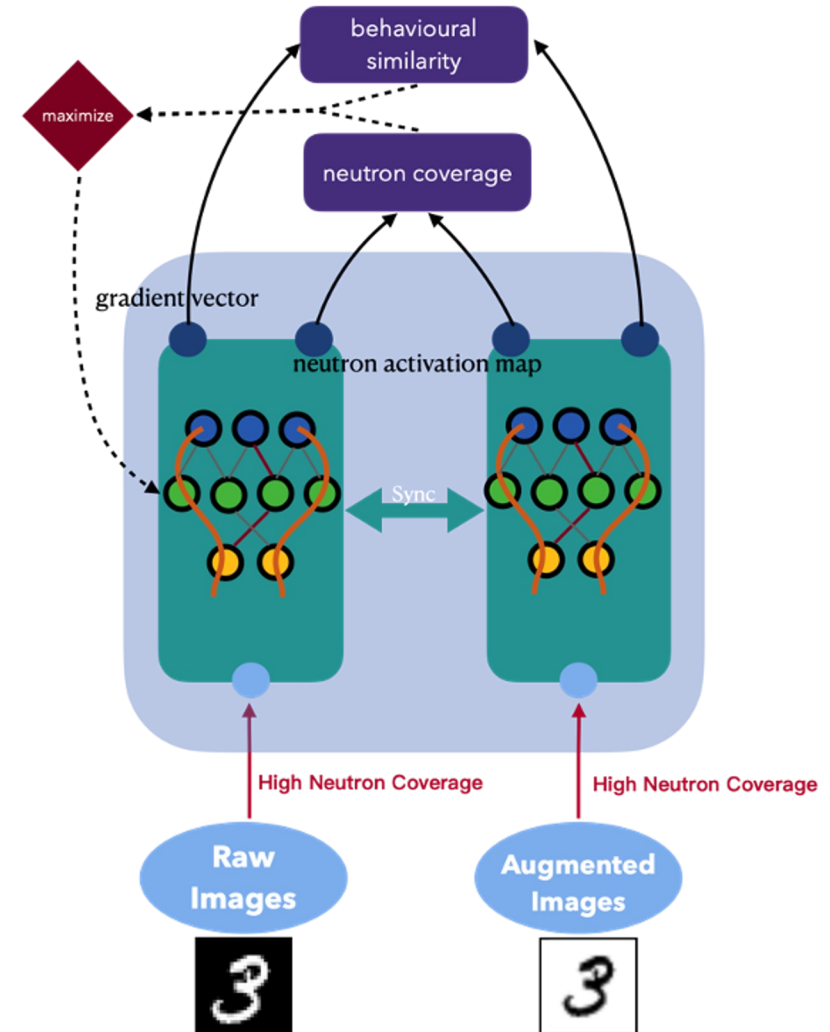
# Interpretability via neuron explanations

We proposed NCDG to actively activate the inactive neurons during training with the neuron coverage maximization loss.

If a neuron is inactive

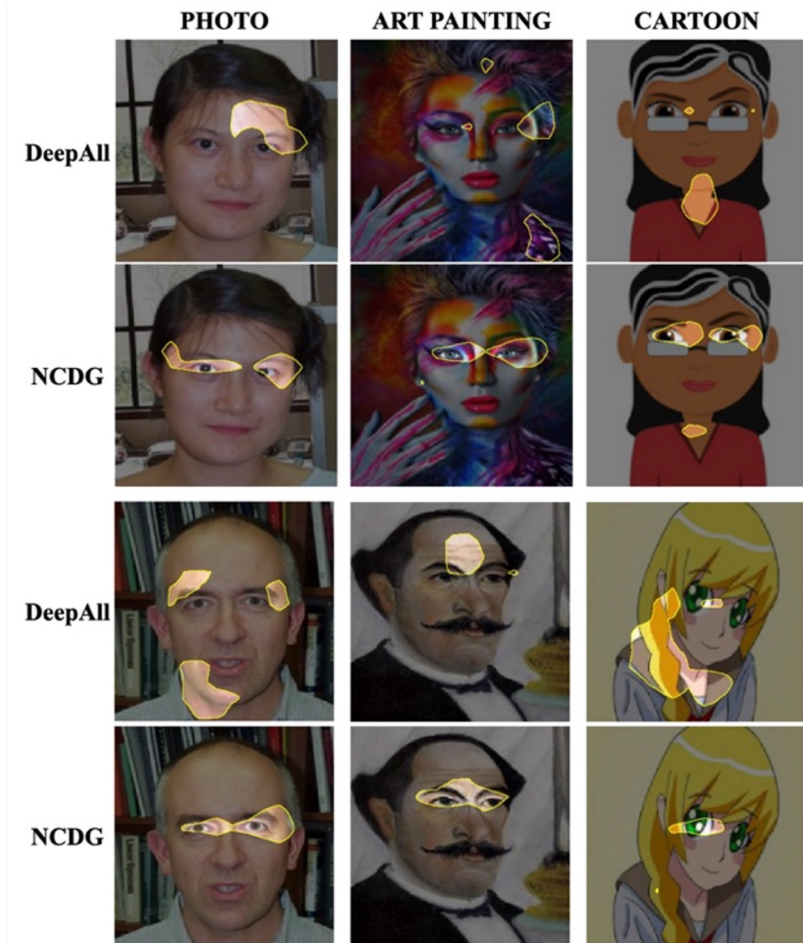( always low-value output )

during the **WHOLE** training process.

Once it gets activated (outputs high-value)

during the evaluation,

errors may happen.

Tian, Chris Xing, et al. "Neuron Coverage-Guided Domain Generalization." *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022).
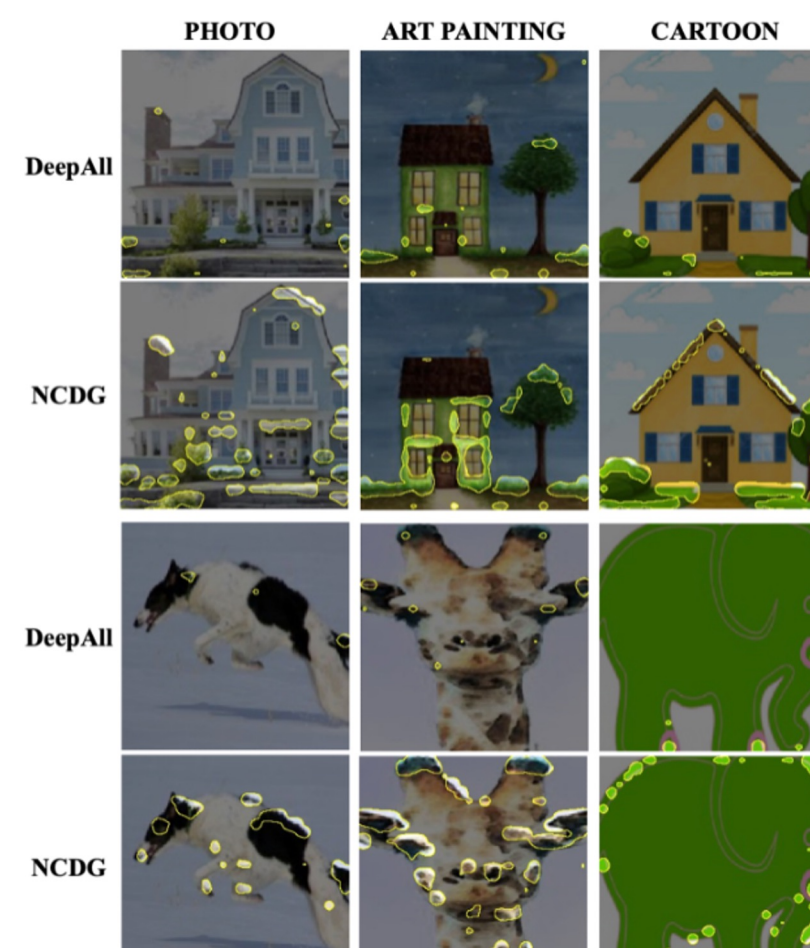
# Interpretability via neuron explanations

Neuron Dissection before & after (being activated)

*Model trained on PACS-Sketch domain*



**ResNet-18 block 3 unit 170**

**ResNet-18 block 2 unit 20**