

Formalizing Robustness in Neural Networks: Explainability, Uncertainty, and Intervenability



Ghassan AlRegib, PhD
Professor



Mohit Prabhushankar, PhD
Postdoctoral Fellow

Association for the Advancement
of Artificial Intelligence

Omni Lab for Intelligent Visual Engineering and Science (OLIVES)
School of Electrical and Computer Engineering
Georgia Institute of Technology
{alregib, mohit.p}@gatech.edu

Feb 21, 2024 – Vancouver, BC, Canada



Tutorial Materials

Accessible Online



<https://alregib.ece.gatech.edu/aaai-2024-tutorial/>
{alregib, mohit.p}@gatech.edu

AAAI 2024 Tutorial



Presented by: *Ghassan AlRegib, and Mohit Prabhushankar*
Georgia Institute of Technology

www.ghassanalregib.info

Duration: Half Day (3 hours, 30 mins)

Title: Formalizing Robustness in Neural Networks: Explainability, Uncertainty, and Intervenability



Expectation vs Reality of Deep Learning



Deep Learning

Expectation vs Reality

LATEST TRICKS

Rotating objects in an image confuses DNNs, probably because they are too different from the types of image used to train the network.

Stop



Dumb-bell

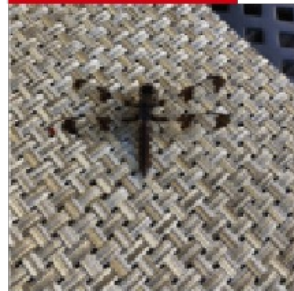


Racket

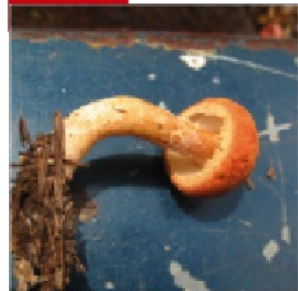


Even natural images can fool a DNN, because it might focus on the picture's colour, texture or background rather than picking out the salient features a human would recognize.

Manhole cover



Pretzel



©nature



Deep Learning

Expectation vs Reality

*“The best-laid plans of sensors and networks
often go awry”*

- Engineers, probably



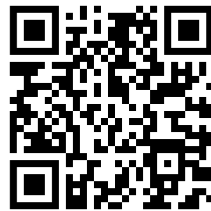
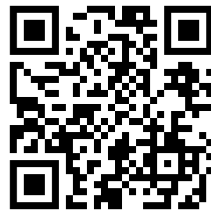
Deep Learning

Requirements and Challenges

Requirements: Deep Learning-enabled systems must predict correctly on novel data

Novel data sources:

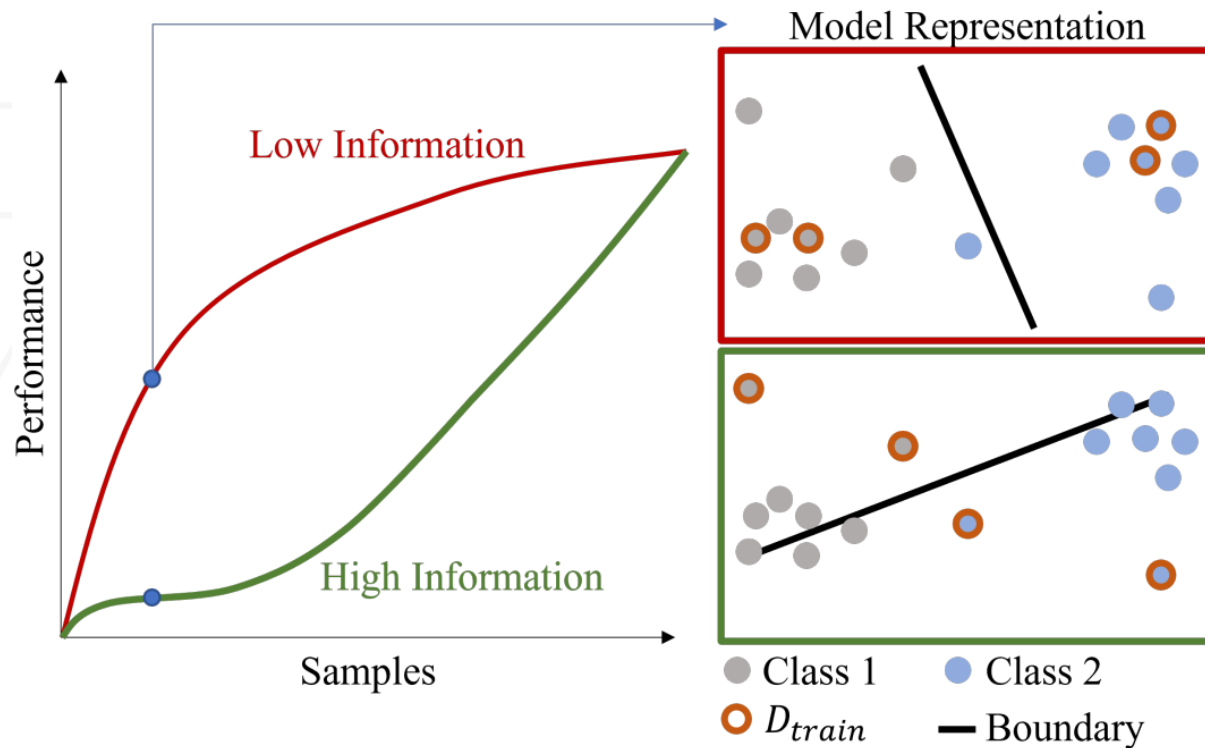
- Test distributions
- Anomalous data
- Out-Of-Distribution data
- Adversarial data
- Corrupted data
- Noisy data
- New classes
- ...



Deep Learning at Training

Overcoming Challenges at Training: Part 1

The most novel/aberrant samples should not be used in early training



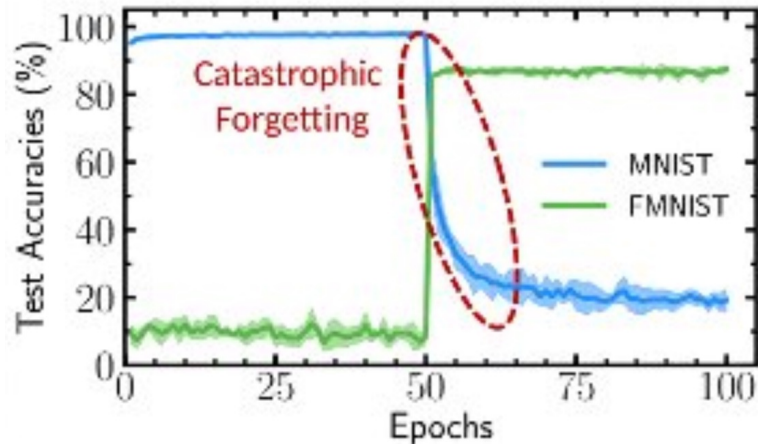
- The first instance of training must occur with less informative samples
- Ex: For autonomous vehicles, less informative means
 - Highway scenarios
 - Parking
 - No accidents
 - No aberrant events

Novel samples = Most Informative

Deep Learning at Training

Overcoming Challenges at Training: Part 2

Subsequent training must not focus only on novel data



- The model performs well on the new scenarios, while forgetting the old scenarios
- A number of techniques exist to overcome this trend
- However, they affect the overall performance in large-scale settings
- It is not always clear **if and when** to incorporate novel scenarios in training

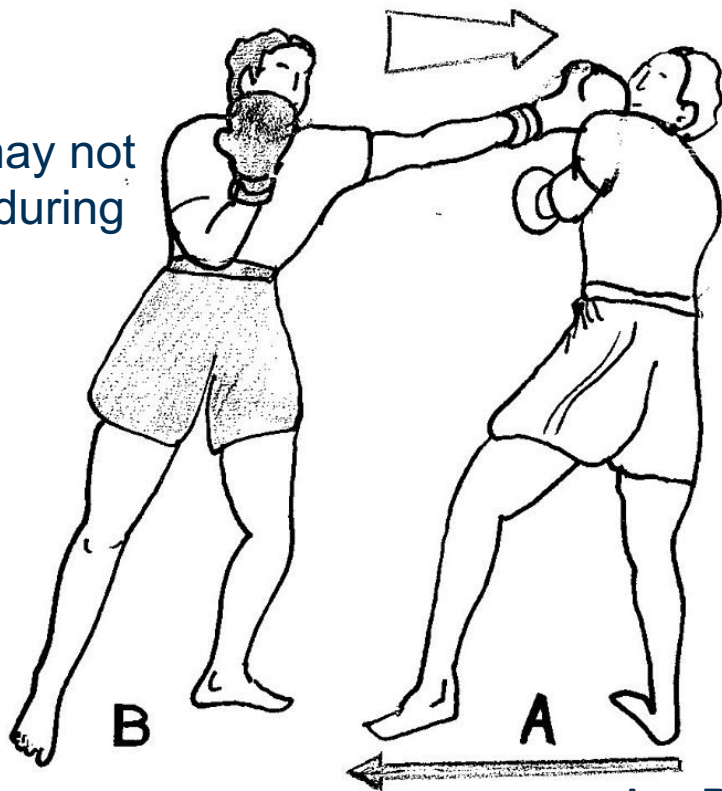


Deep Learning at Training

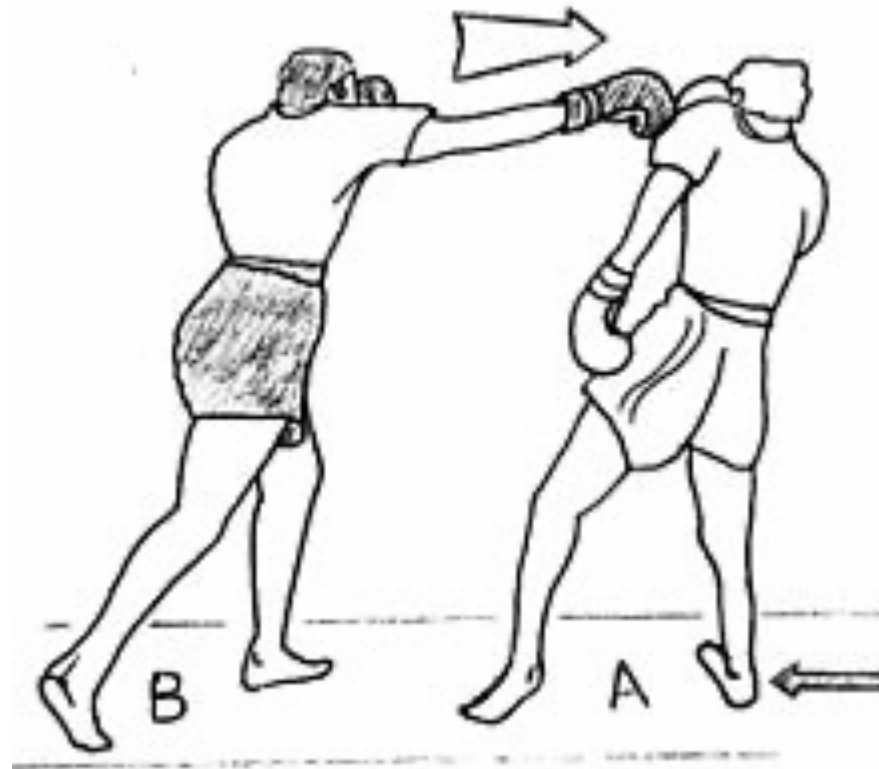
Overcoming Challenges at Training

Novel data packs a 1-2 punch!

Novel data may not be available during training



A = Deep Neural Networks
B = Novel data



Even if available, novel data does not easily fit into either the earlier or later stages of training

Deep Learning at Inference

Overcoming Challenges at Inference

We must handle novel data at Inference!!

Novel data sources:

- Test distributions
- Anomalous data
- Out-Of-Distribution data
- Adversarial data
- Corrupted data
- Noisy data
- New classes
- ...

Model Train



At Inference



Objective

Objective of the Tutorial

To discuss methodologies that promote robustness in neural networks at inference

- Part 1: Inference in Neural Networks
- Part 2: Explainability at Inference
- Part 3: Uncertainty at Inference
- Part 4: Intervenability at Inference
- Part 5: Conclusions and Future Directions



Robust Neural Networks

Part I: Inference in Neural Networks



Objective

Objective of the Tutorial

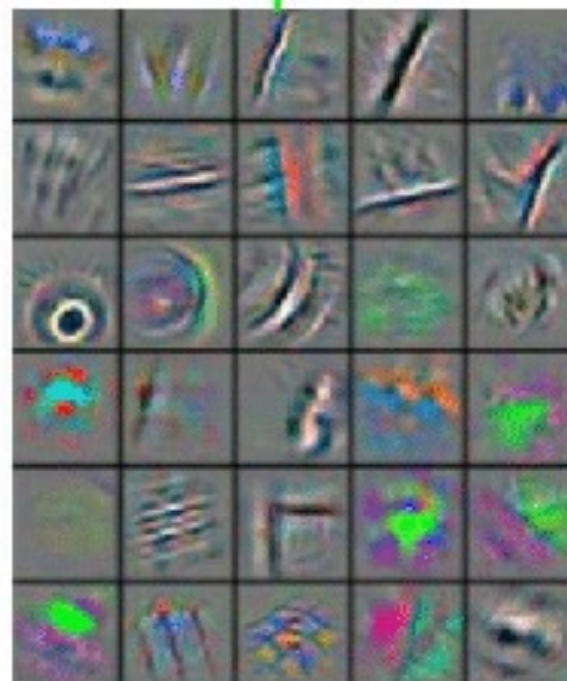
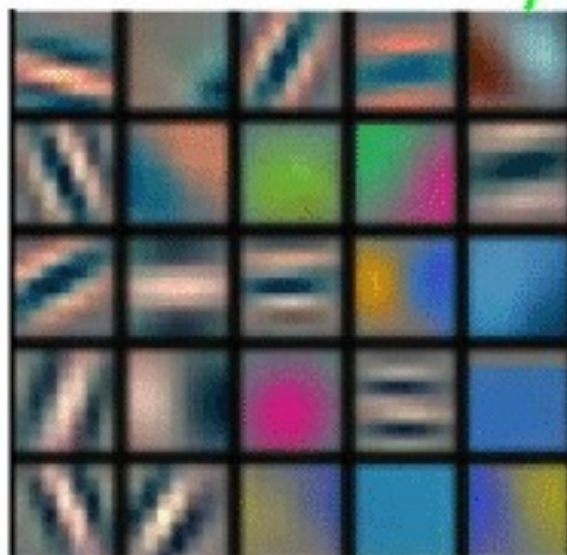
To discuss methodologies that promote robustness in neural networks at inference

- **Part 1: Inference in Neural Networks**
 - Neural Network Basics
 - Robustness in Deep Learning
 - Information at Inference
 - Challenges at Inference
 - Gradients at Inference
- Part 2: Explainability at Inference
- Part 3: Uncertainty at Inference
- Part 4: Intervenability at Inference
- Part 5: Conclusions and Future Directions



Deep Learning

Overview



Ex. LeCun, 2015

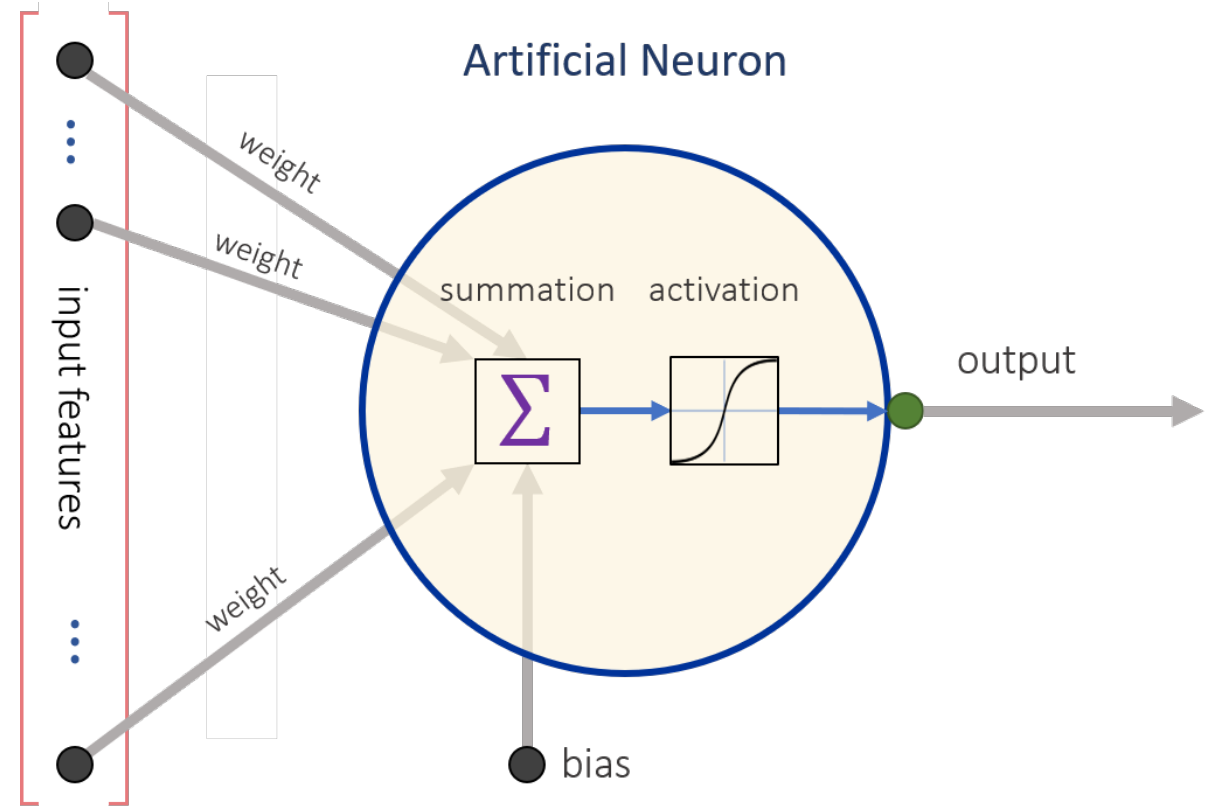
Deep Learning

Neurons

The underlying computation unit is the Neuron

Artificial neurons consist of:

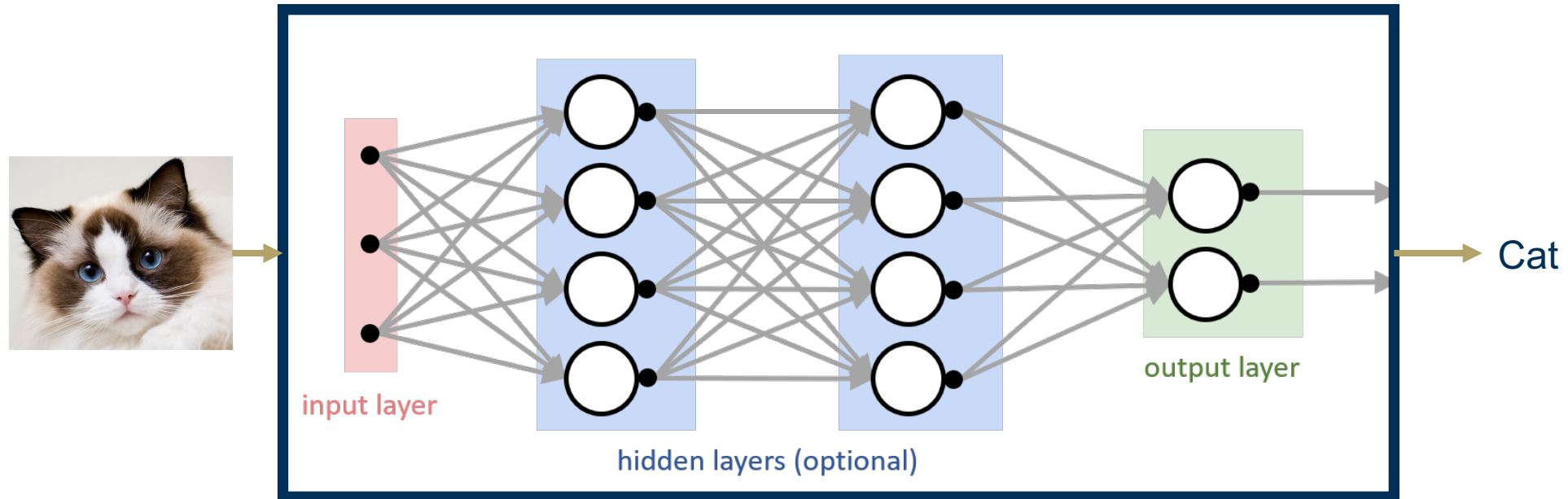
- A single output
- Multiple inputs
- Input weights
- A bias input
- An activation function



Deep Learning

Artificial Neural Networks

Neurons are stacked and densely connected to construct ANNs



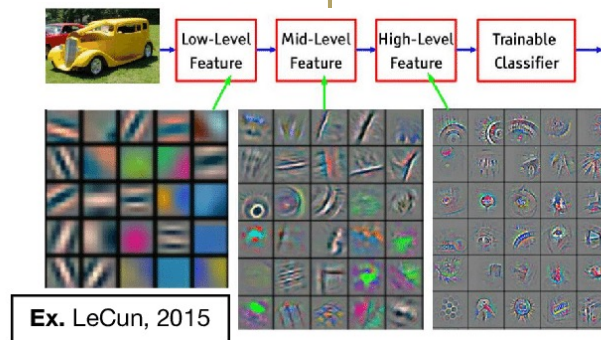
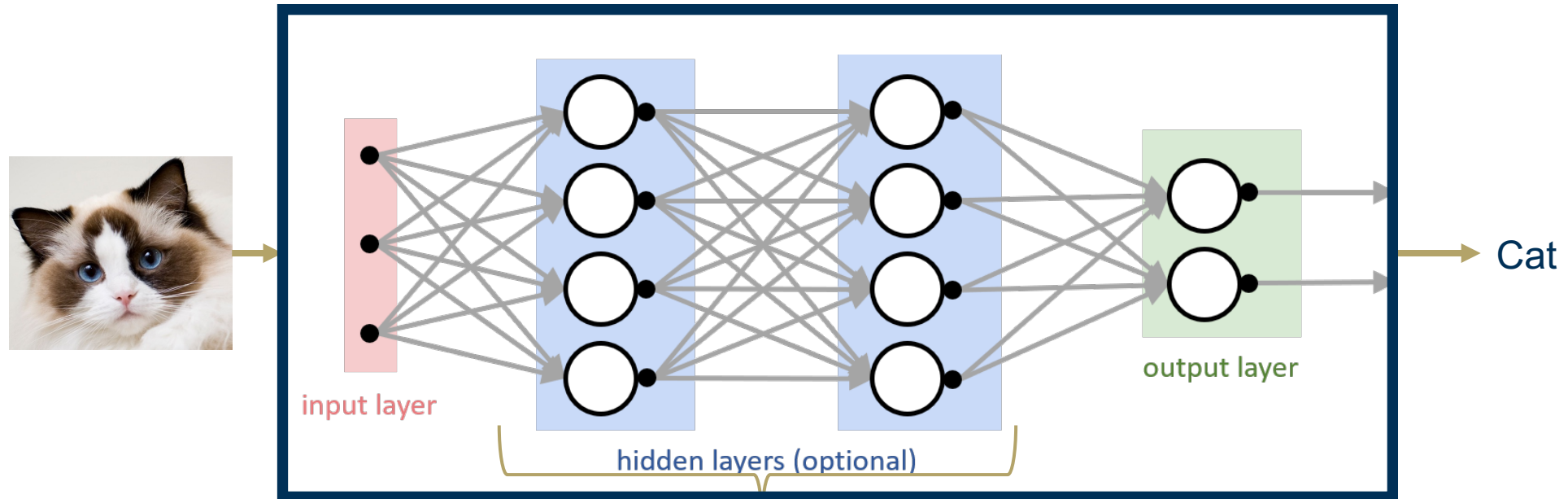
Typically, a neuron is part of a network organized in layers:

- An input layer (Layer 0)
- An output layer (Layer K)
- Zero or more hidden (middle) layers (Layers $1 \dots K - 1$)

Deep Learning

Convolutional Neural Networks

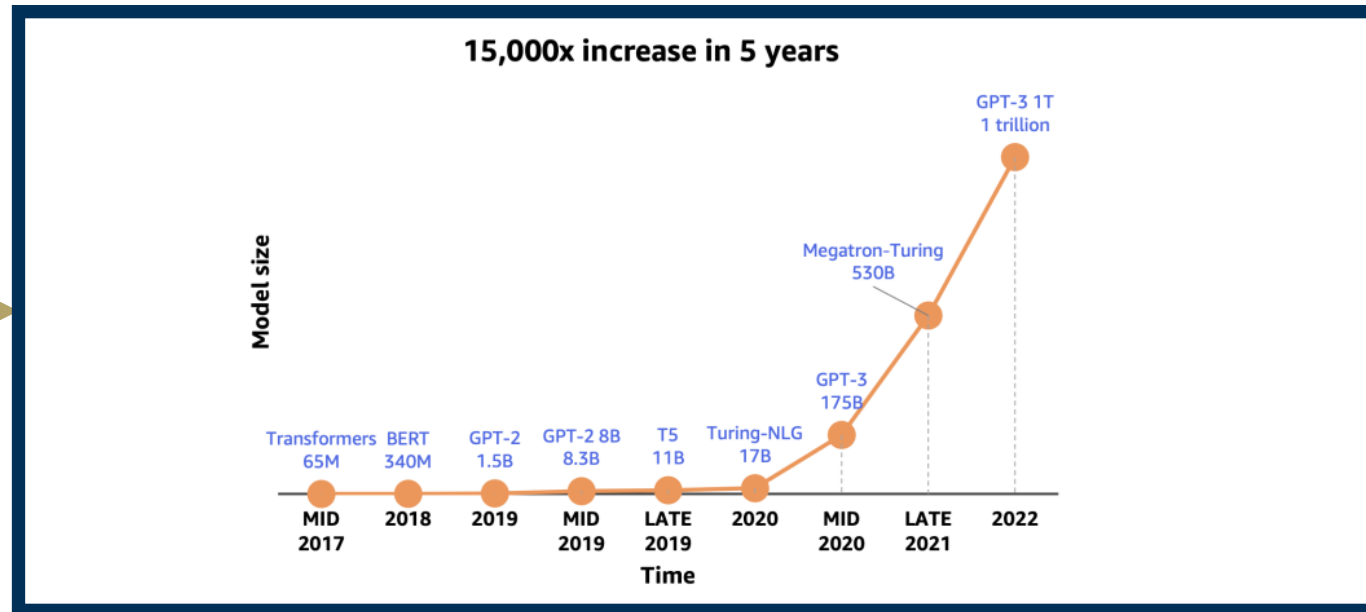
Stationary property of images allow for a small number of convolution kernels



Deep Deep Deep Deep Deep ... Learning

Recent Advancements

Transformers, Large Language Models and Foundation Models



Cat

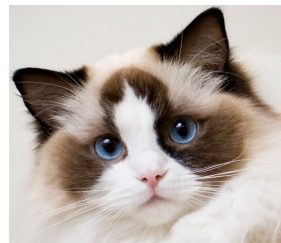
Primary reasons for advancements:

1. Expanded interests from the research community
2. Computational resources availability
3. **Big data availability**

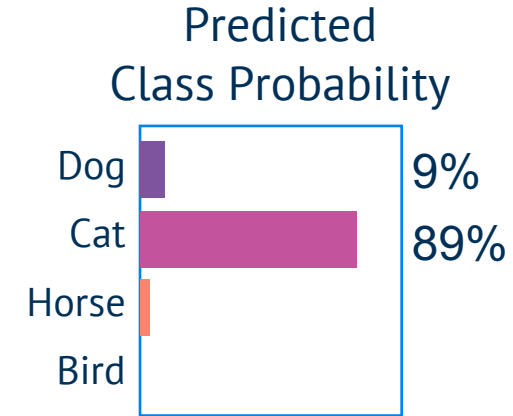
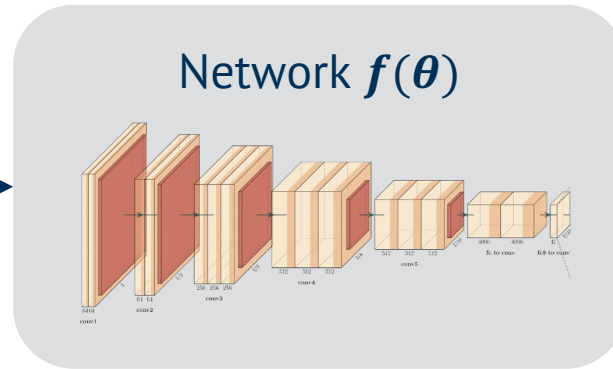
Deep Learning at Inference

Classification

Given : One network, One image. Required: Class Prediction



x



If $x \in \chi$, the data is **not novel**

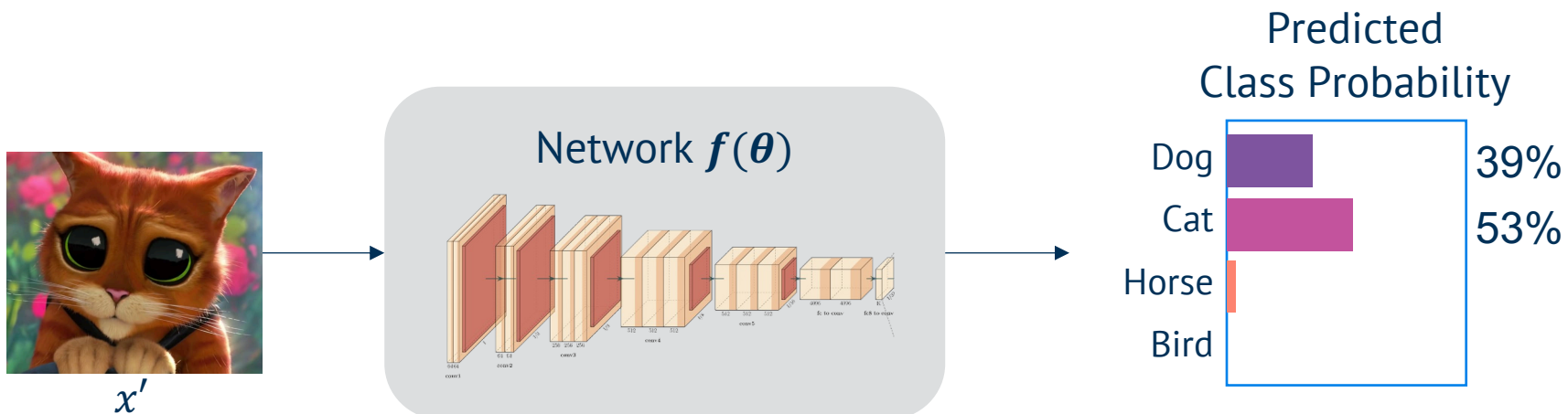
$$\hat{y} = f(x)$$
$$y = \operatorname{argmax}_i \hat{y}$$
$$p(\hat{y}) = T(f(x))$$

\hat{y} = Logits
 y = Predicted Class
 $p(\hat{y})$ = Probabilities
 $f(\cdot)$ = Trained Network
 χ = Training data

Deep Learning at Inference

Robust Classification in Deep Networks

Deep learning robustness: Correctly predict class even when data is novel



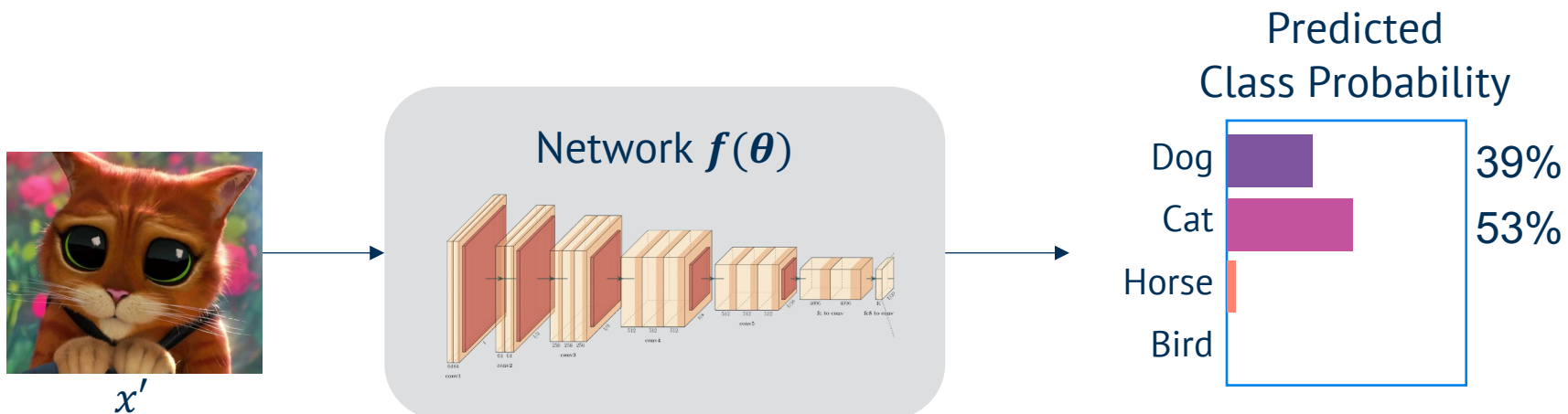
If $x \in \chi$, the data is **novel**

$$\begin{aligned} \hat{y} &= f(x' + \epsilon) & \hat{y} &= \text{Logits} \\ y &= \operatorname{argmax}_i \hat{y} & y &= \text{Predicted Class} \\ p(\hat{y}) &= T(f(x' + \epsilon)) & p(\hat{y}) &= \text{Probabilities} \\ & & f(\cdot) &= \text{Trained Network} \\ & & \chi &= \text{Training data} \\ & & \epsilon &= \text{Noise} \end{aligned}$$

Deep Learning at Inference

Robust Classification in Deep Networks

Deep learning robustness: Correctly predict class even when data is novel



To achieve robustness at Inference, we need the following:

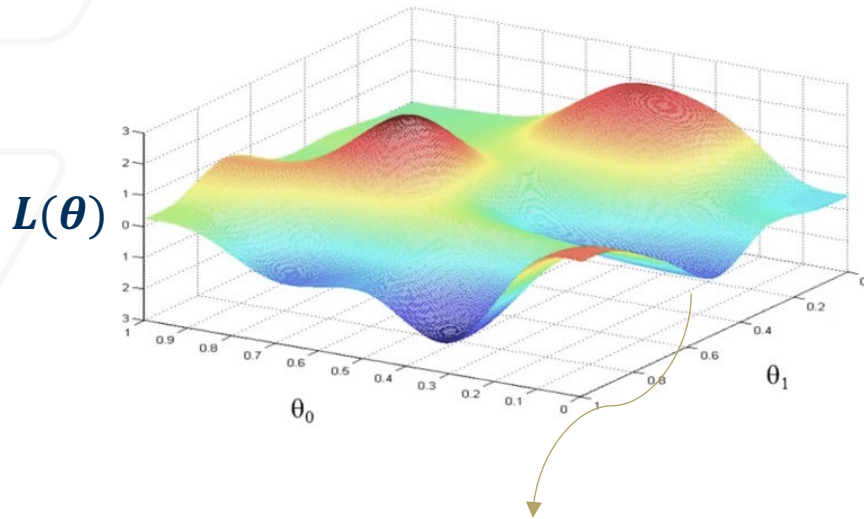
- **Information** provided by the novel data as a **function of training distribution**
- Methodology to **extract information** from novel data
- **Techniques** that utilize the information from novel data

Why is this Challenging?

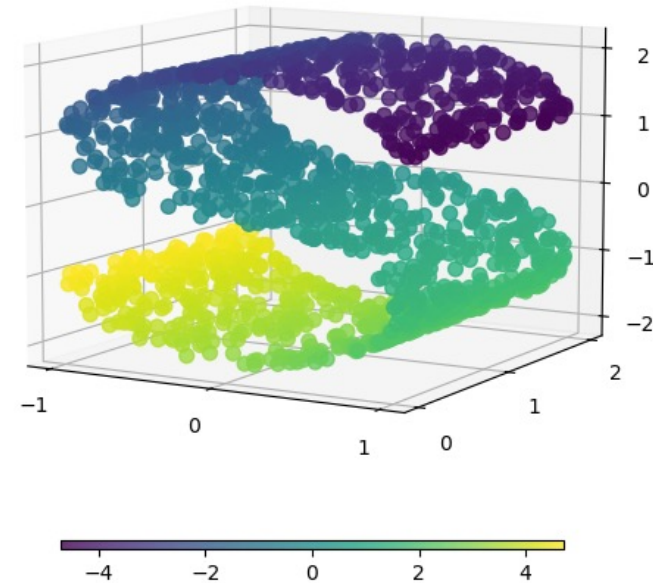
Challenges at Inference

A Quick note on Manifolds..

Manifolds are compact topological spaces that allow exact mathematical functions



Toy visualizations generated using functions
(and thousands of generated data points)

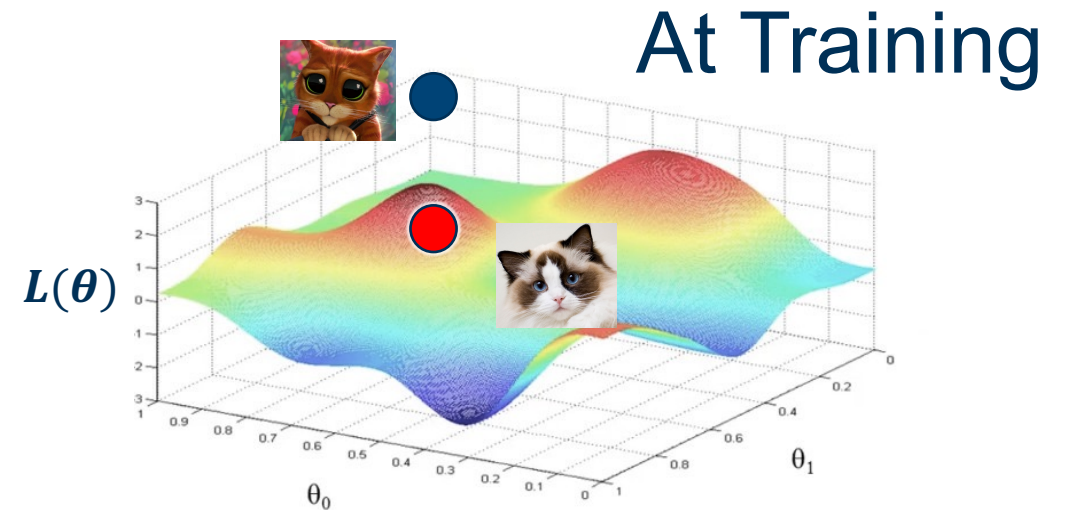
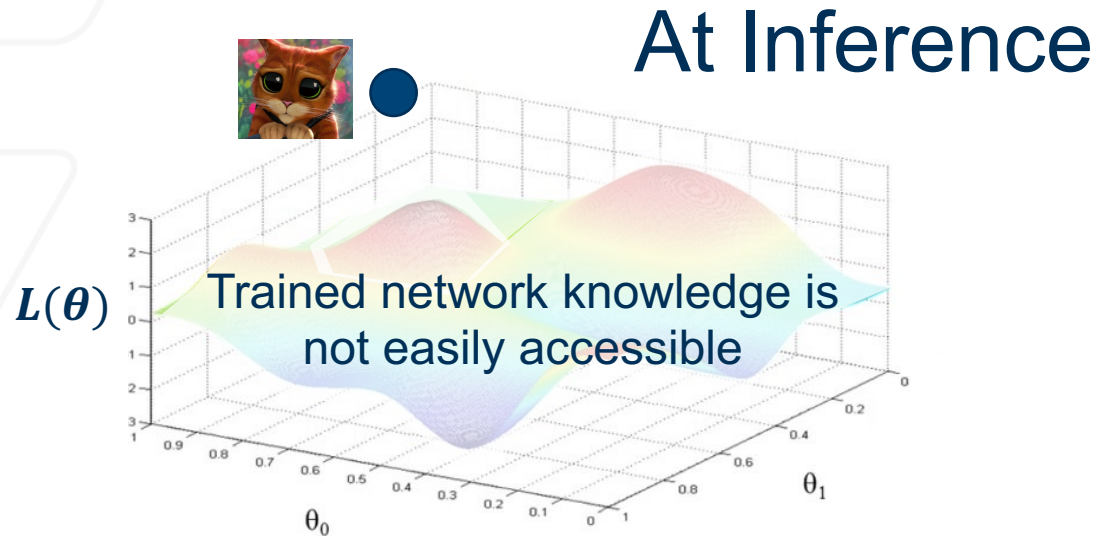


Real data visualizations generated using
dimensionality reduction algorithms (Isomap)

Challenges at Inference

Inference

However, at inference only the test data point is available and the underlying structure of the manifold is unknown

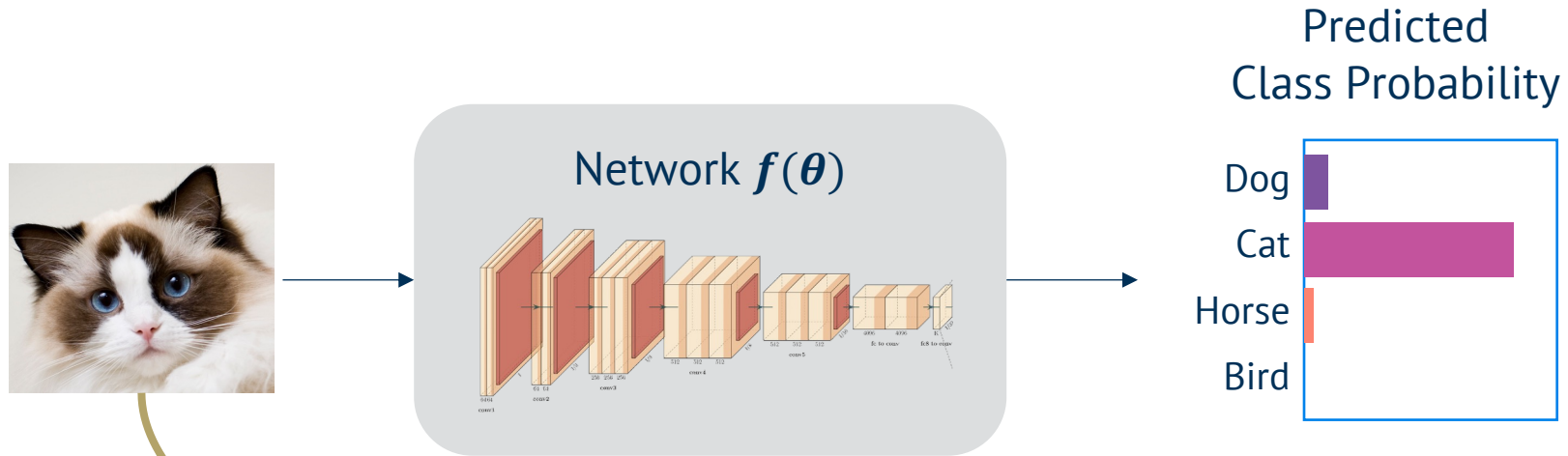


At training, we have access to all training data.

Information at Inference

Fisher Information

Colloquially, Fisher Information is the “surprise” in a system that observes an event

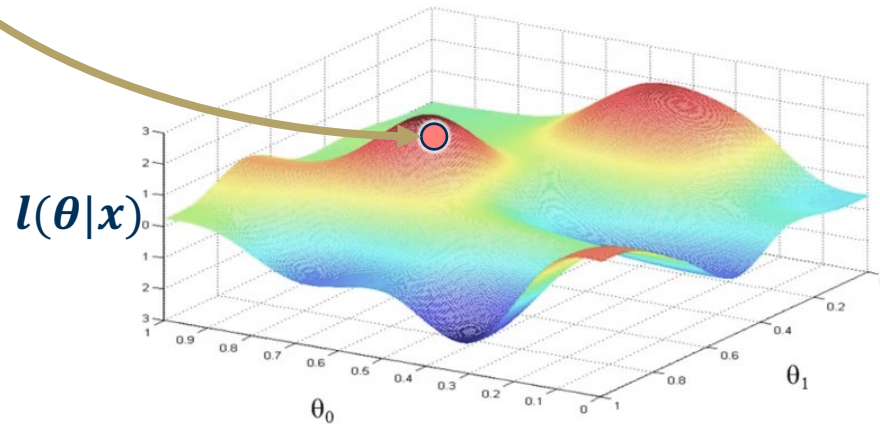


Fisher Information

$$I(\theta) = \text{Var}\left(\frac{\partial}{\partial \theta} l(\theta|x)\right)$$

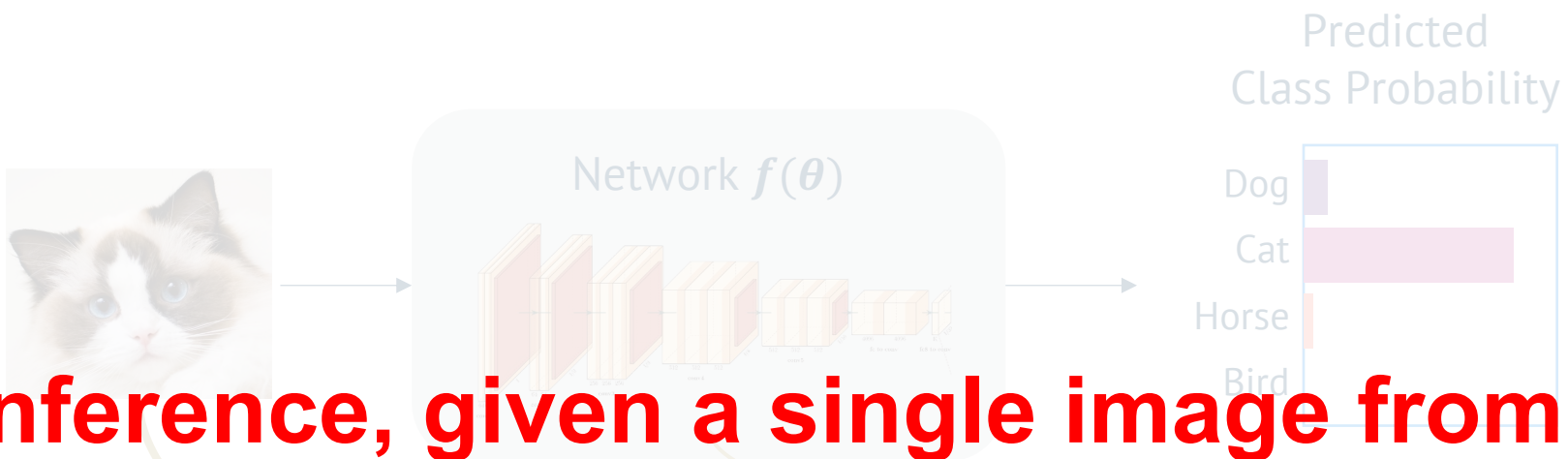
θ = Statistic of distribution
 $l(\theta | x)$ = Likelihood function

Likelihood function

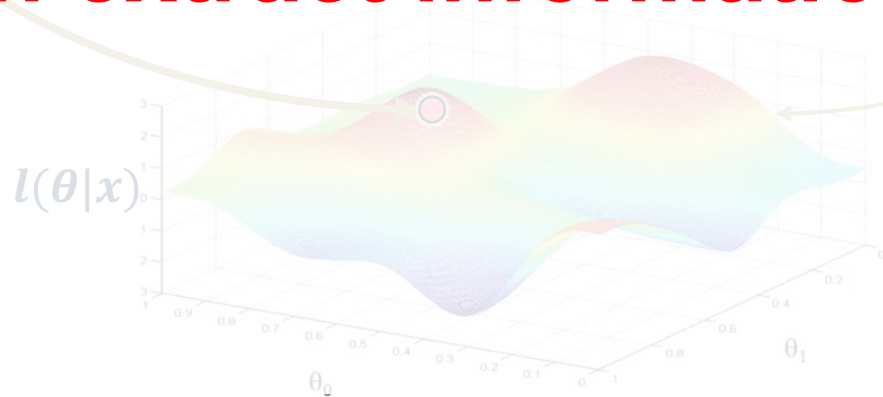


Information at Inference

Information at Inference



At inference, given a single image from a single class, we can extract information about other classes



Likelihood function

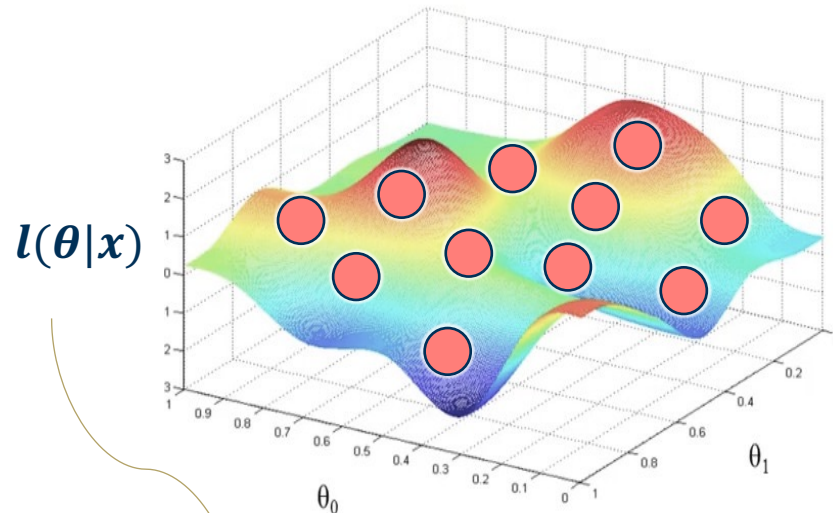
$$I(\theta) = \text{Var}\left(\frac{\partial}{\partial \theta} l(\theta|x)\right)$$

θ = Statistic of distribution
 $l(\theta | x)$ = Likelihood function

Information at Inference

Gradients as Fisher Information

Gradients infer information about the statistics of underlying manifolds



From before, $I(\theta) = \text{Var}\left(\frac{\partial}{\partial \theta} l(\theta|x)\right)$

Using variance decomposition, $I(\theta)$ reduces to:

$I(\theta) = E[U_\theta U_\theta^T]$ where

$E[\cdot]$ = Expectation

$U_\theta = \nabla_\theta l(\theta|x)$, Gradients w.r.t. the sample

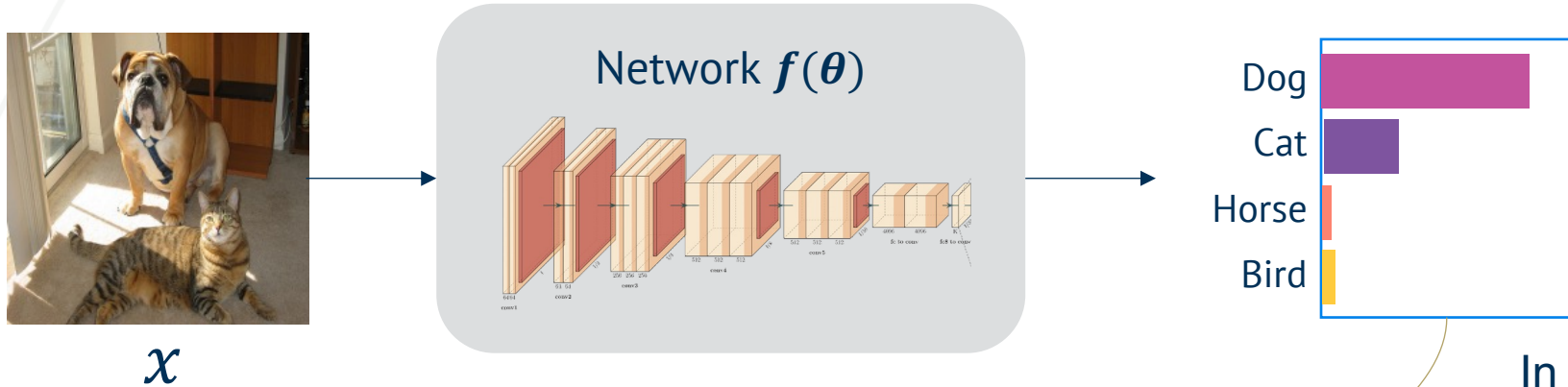
Likelihood function instead of loss manifold

Hence, gradients draw information from the underlying distribution as learned by the network weights!

Information at Inference

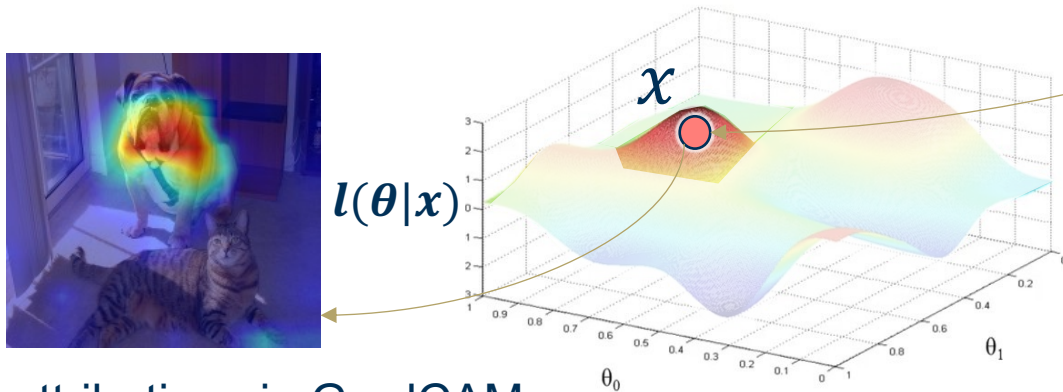
Case Study: Gradients as Fisher Information in Explainability

Gradients infer information about the statistics of underlying manifolds



Local information (specific to x) is sufficient!

In this case, the image and its prediction extracts nose, mouth and jowl features.



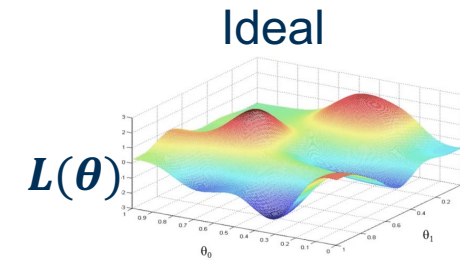
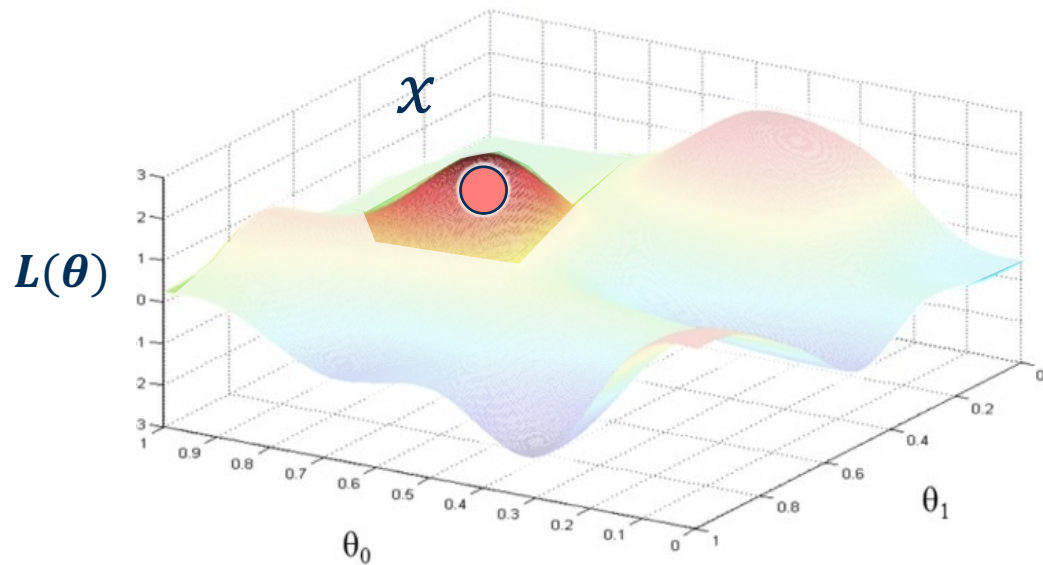
Hence, gradients draw information from the underlying distribution as learned by the network weights!

Feature attribution via GradCAM

Gradients at Inference

Local Information

Gradients provide local information around the vicinity of x , even if x is novel. This is because x projects on the learned knowledge

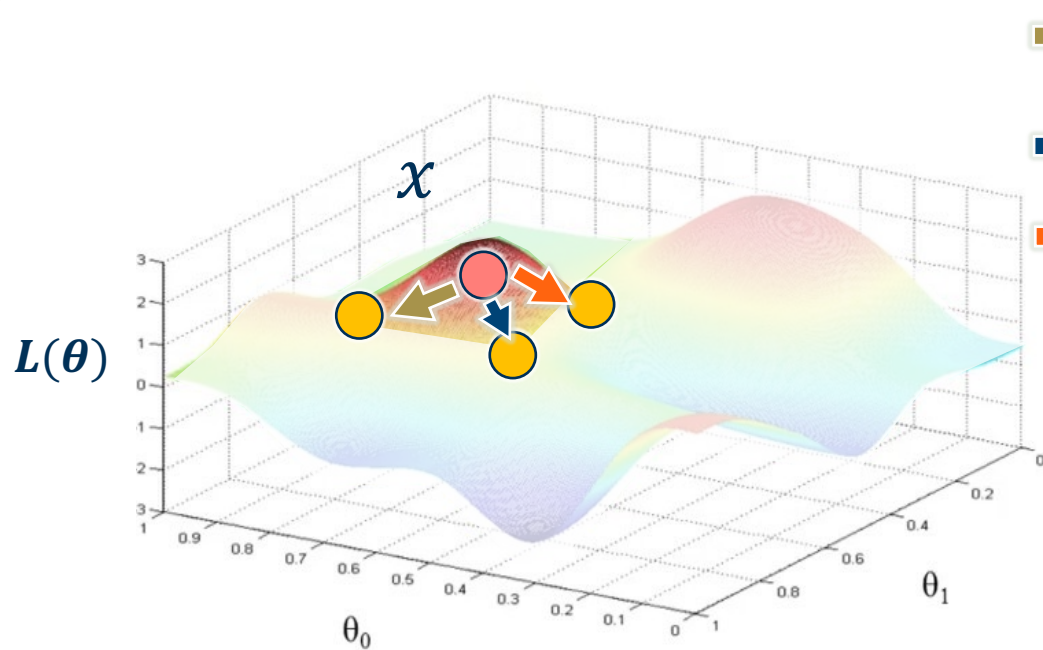


$\alpha \nabla_{\theta} L(\theta)$ provides local information up to a small distance α away from x

Gradients at Inference

Direction of Steepest Descent

Gradients allow choosing the fastest direction of descent given a loss function $L(\theta)$



Path 1?



Path 2?



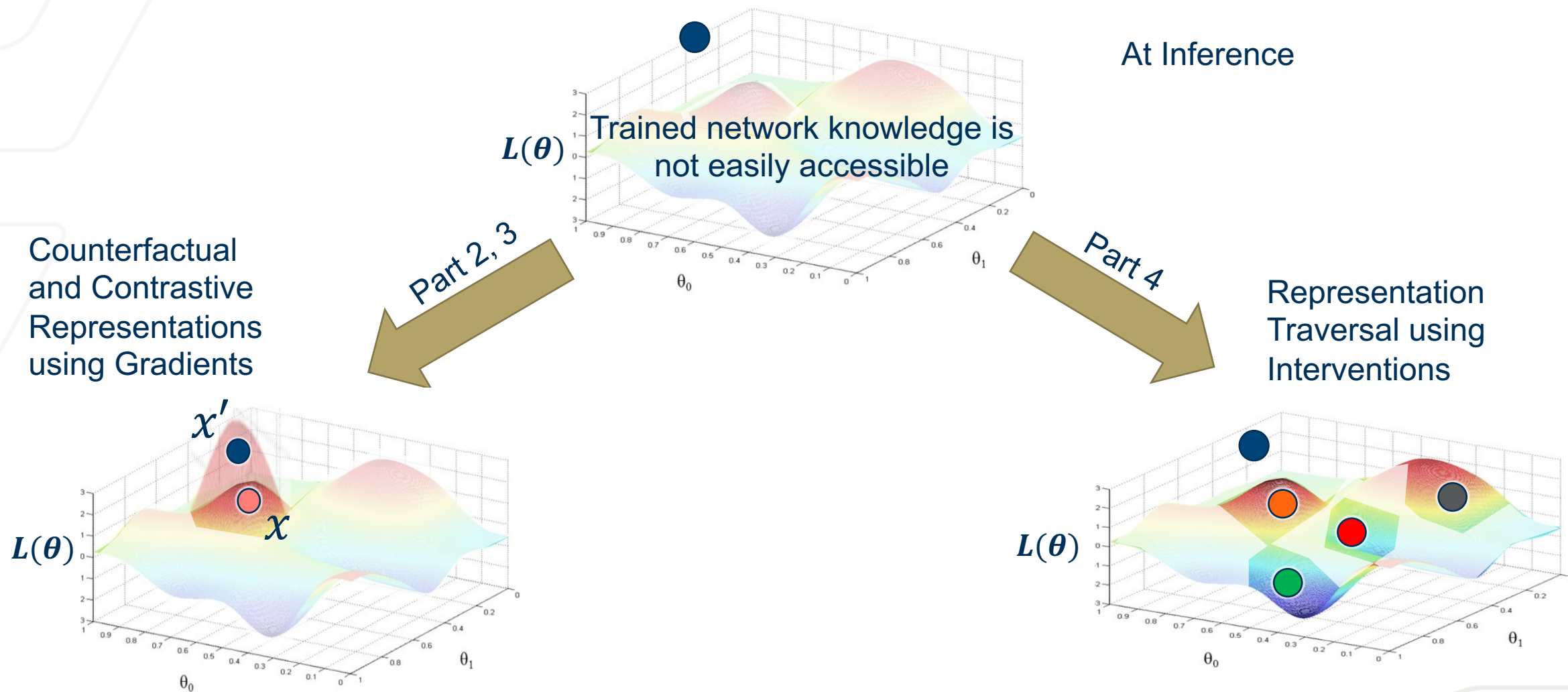
Path 3?

Which direction should we optimize towards (knowing only the local information)?

Negative of the gradient provides the **descent direction** towards the local minima, as measured by $L(\theta)$

Gradients at Inference

To Characterize the Novel Data at Inference



Robust Neural Networks

Part 2: Explainability at Inference



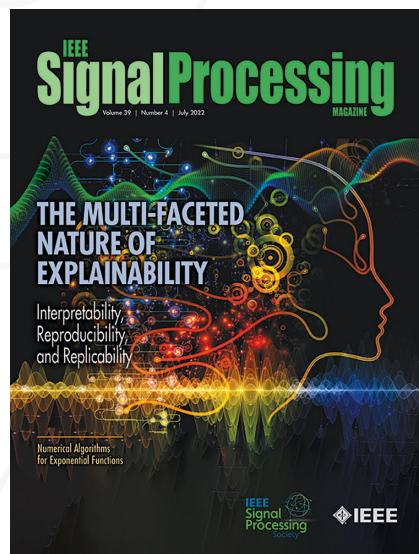
Objective

Objective of the Tutorial

To discuss methodologies that promote robustness in neural networks at inference

- Part 1: Inference in Neural Networks
- **Part 2: Explainability at Inference**
 - Visual Explanations
 - Gradient-based Explanations
 - GradCAM
 - CounterfactualCAM
 - ContrastCAM
- Part 3: Uncertainty at Inference
- Part 4: Intervenability at Inference
- Part 5: Conclusions and Future Directions





Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations



Mohit Prabhushankar, PhD
Postdoc



Ghassan AlRegib, PhD
Professor



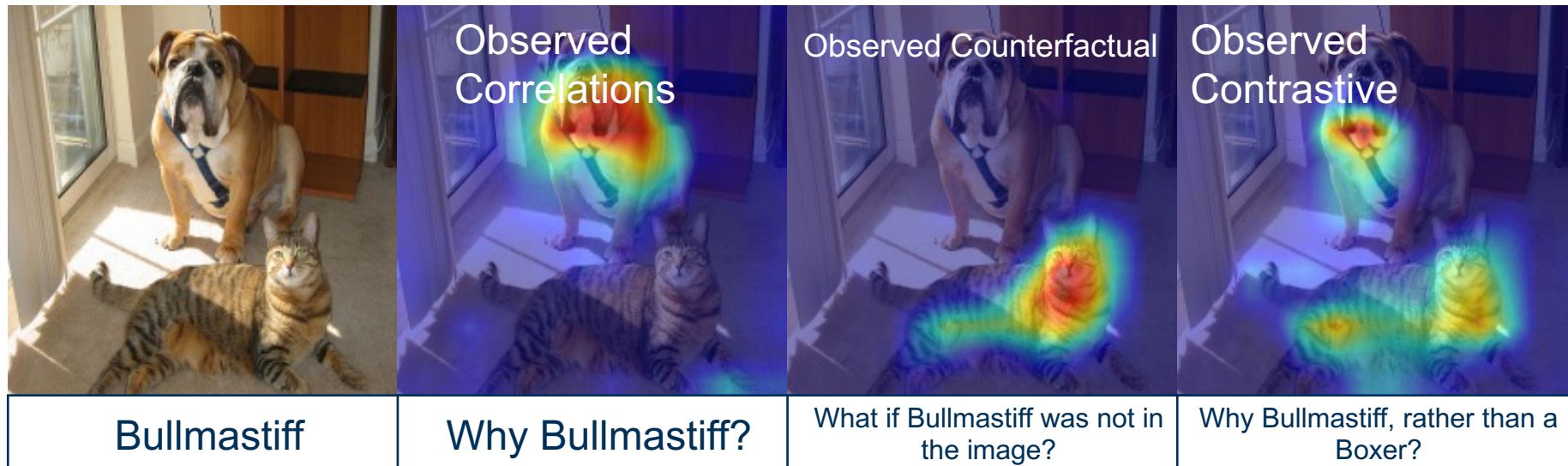
Explanations

Visual Explanations



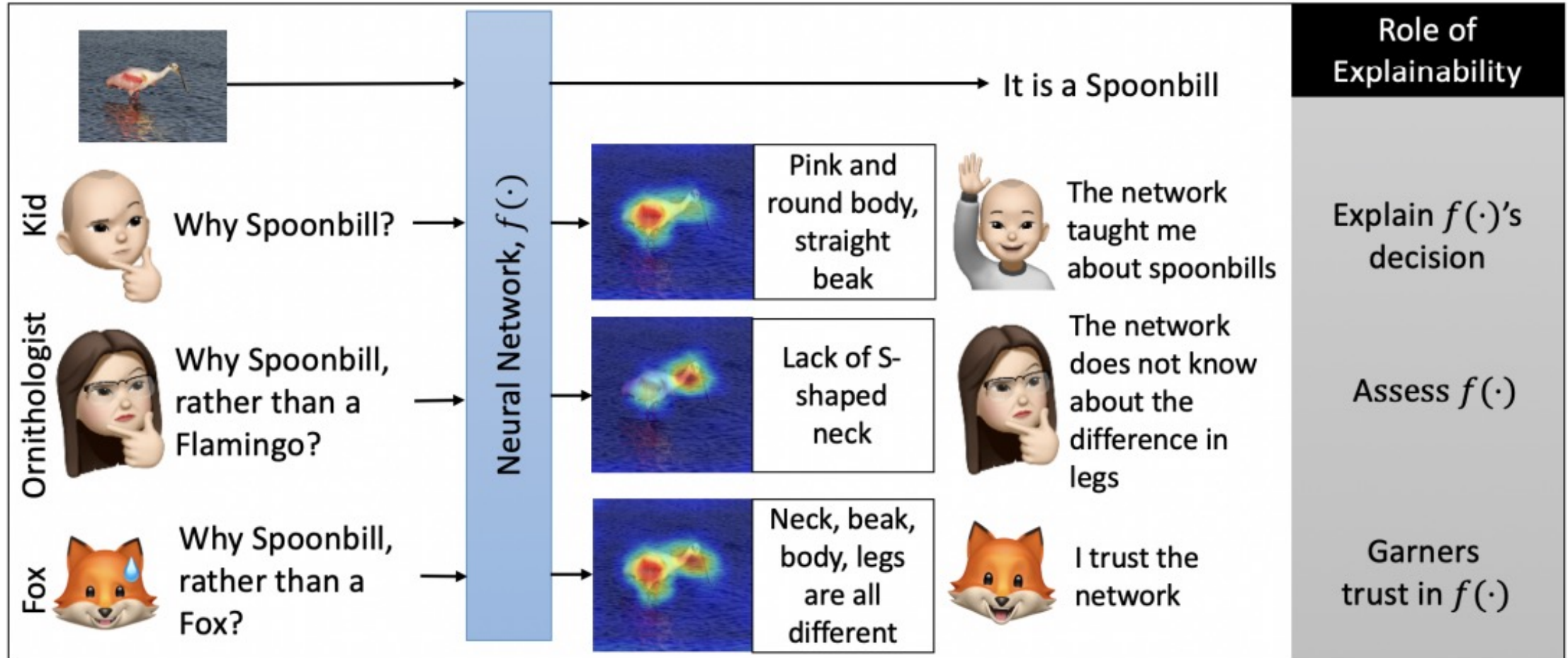
Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

- Explanations are defined as a set of rationales used to understand the reasons behind a decision
- If the decision is based on visual characteristics within the data, the decision-making reasons are visual explanations



Explanations

Role of Explanations – context and relevance



Explanations

Gradient-based Explanations



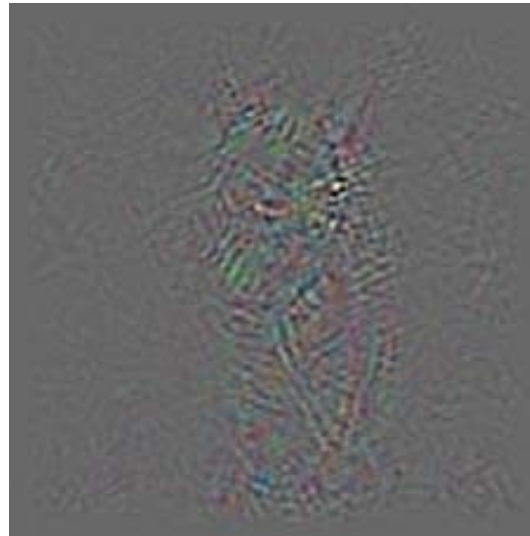
Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Gradients provide a one-shot means of perturbing the input that changes the output; They provide pixel-level importance scores

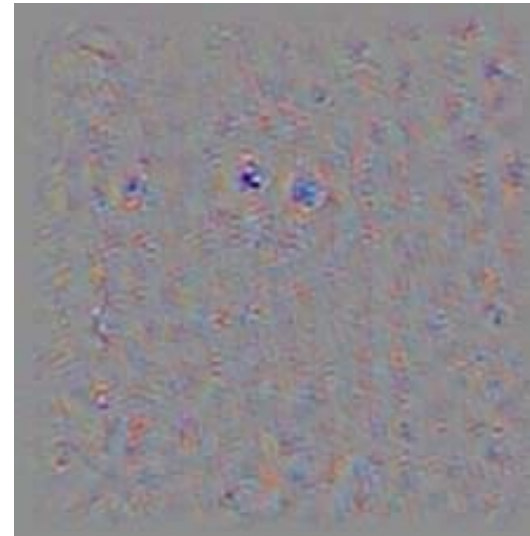
Input



Vanilla Gradients



Deconvolution Gradients



Guided Backpropagation



However, localization remains an issue



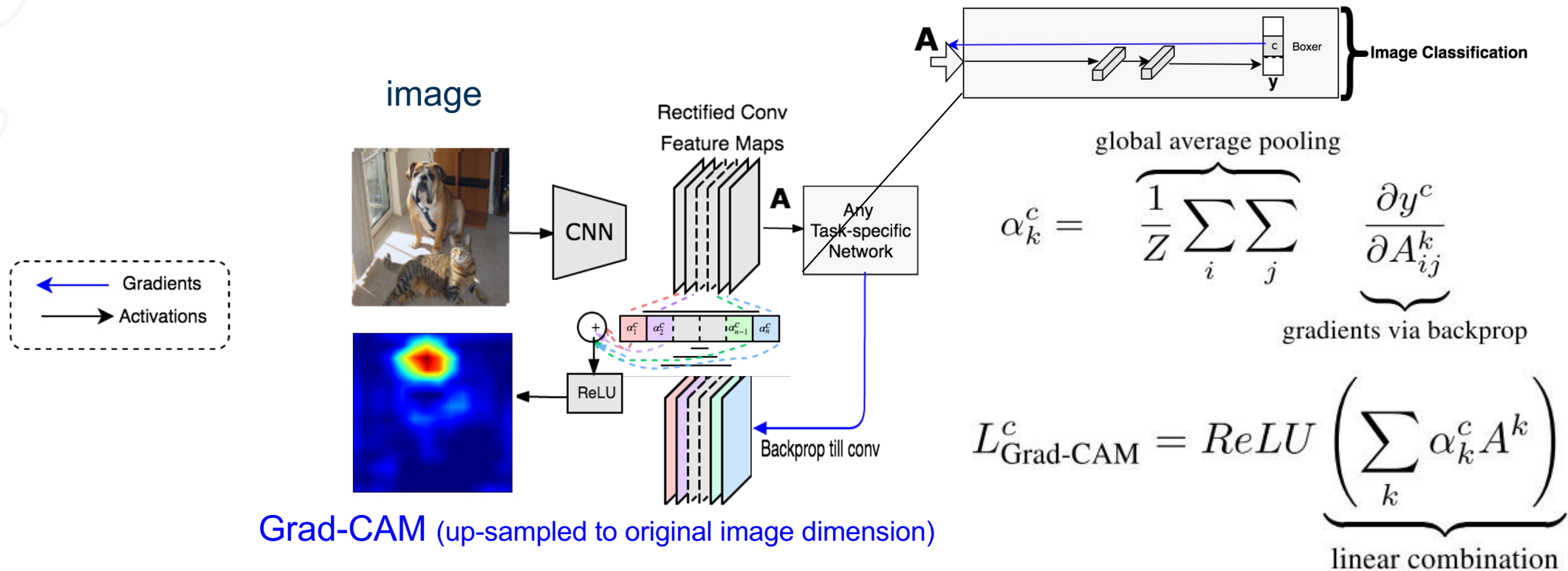
Gradient and Activation-based Explanations

GradCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Grad-CAM uses the gradient information flowing into the last convolutional layer of the CNN to assign importance values to each activation for a particular decision of interest.



Gradient and Activation-based Explanations

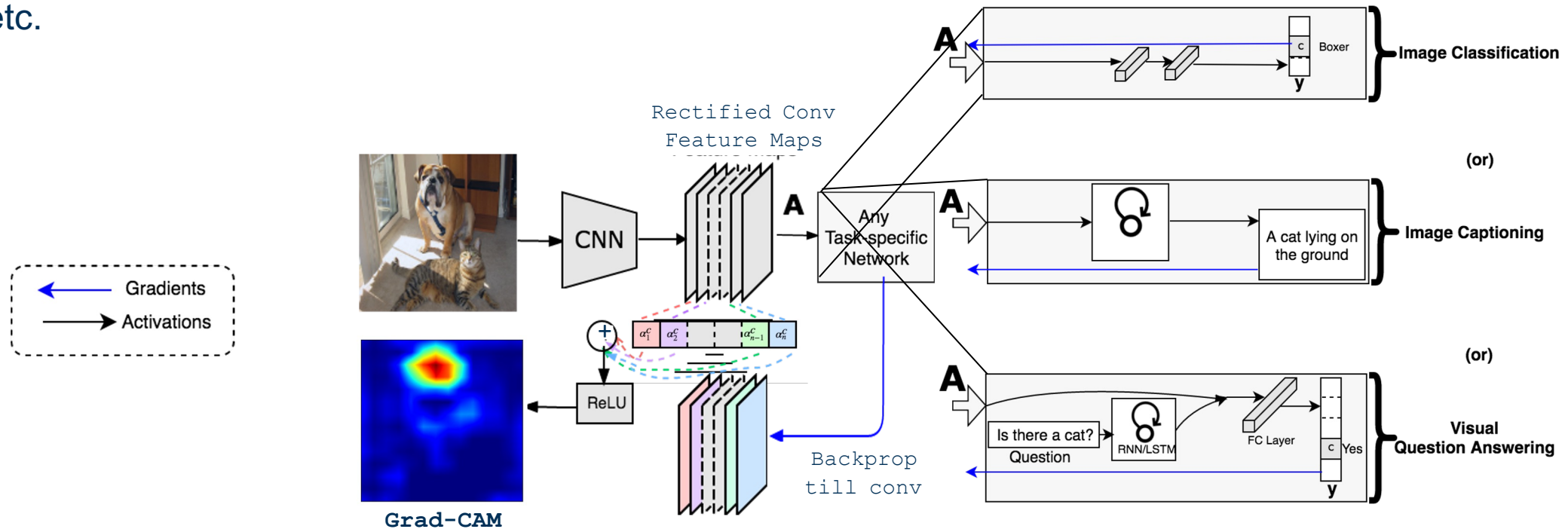
GradCAM

Grad-CAM generalizes to any task:

- Image classification
- Image captioning
- Visual question answering
- etc.



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations



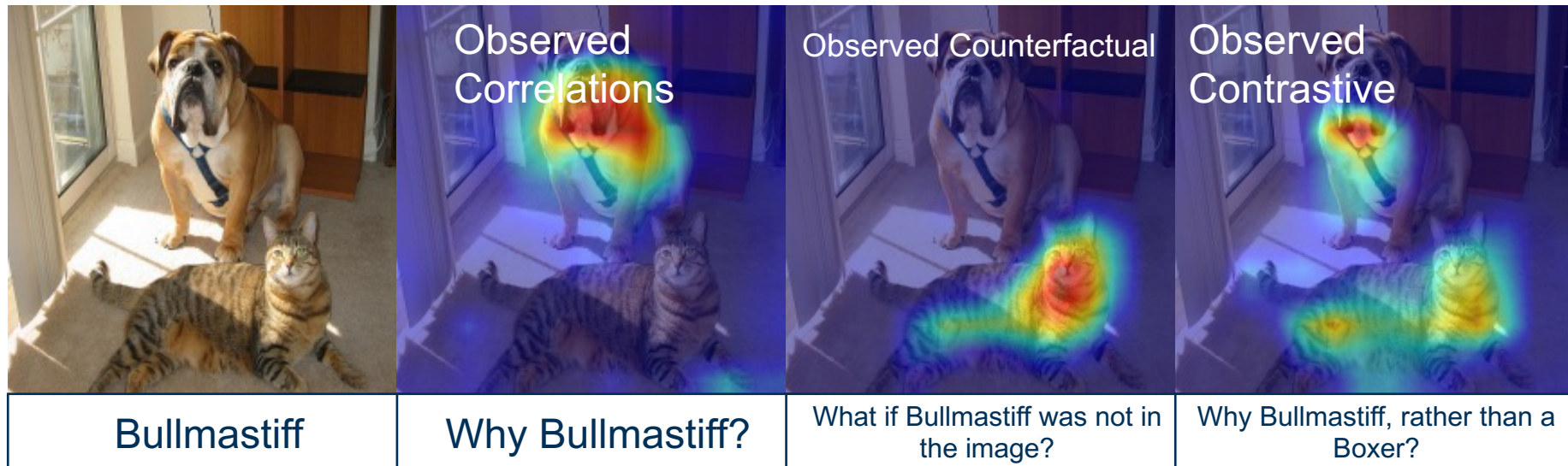
Gradient and Activation-based Explanations

Explanatory Paradigms



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

GradCAM provides answers to ‘Why P?’ questions. But different stakeholders require relevant and contextual explanations



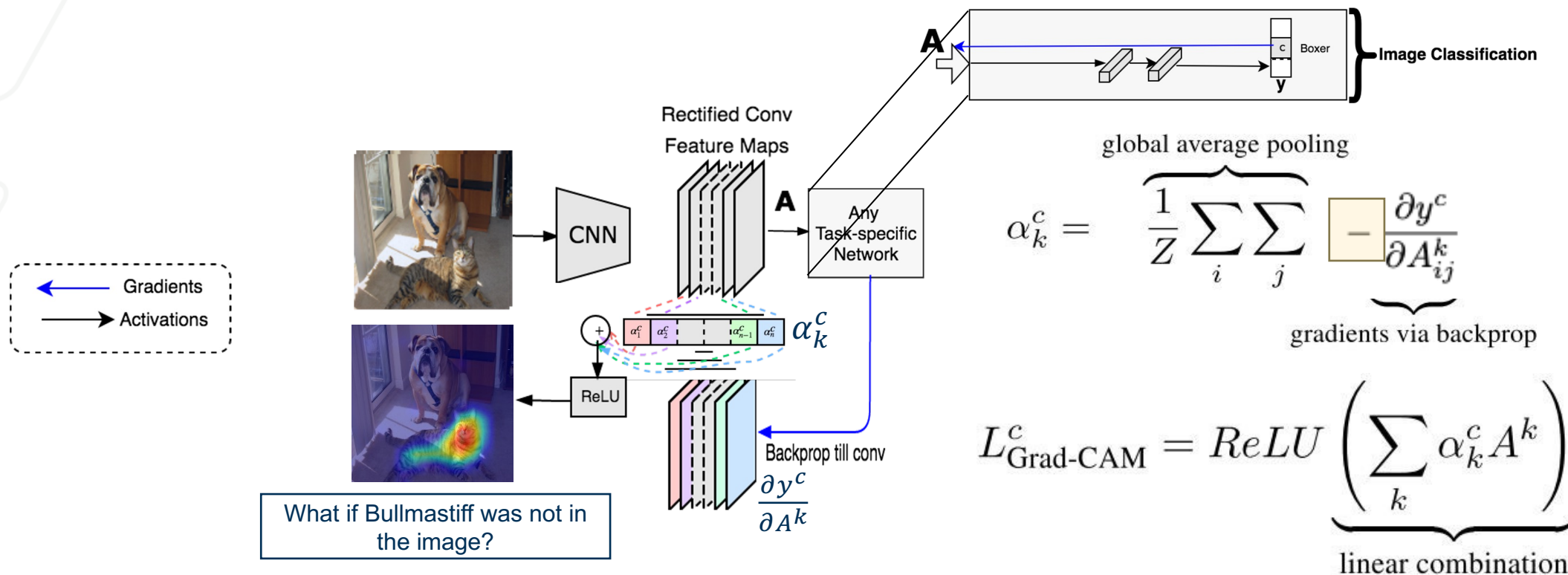
Gradient and Activation-based Explanations

CounterfactualCAM: What if this region were absent in the image?



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

In GradCAM, global average pool the negative of gradients to obtain α^c for each kernel k



Negating the gradients effectively removes these regions from analysis



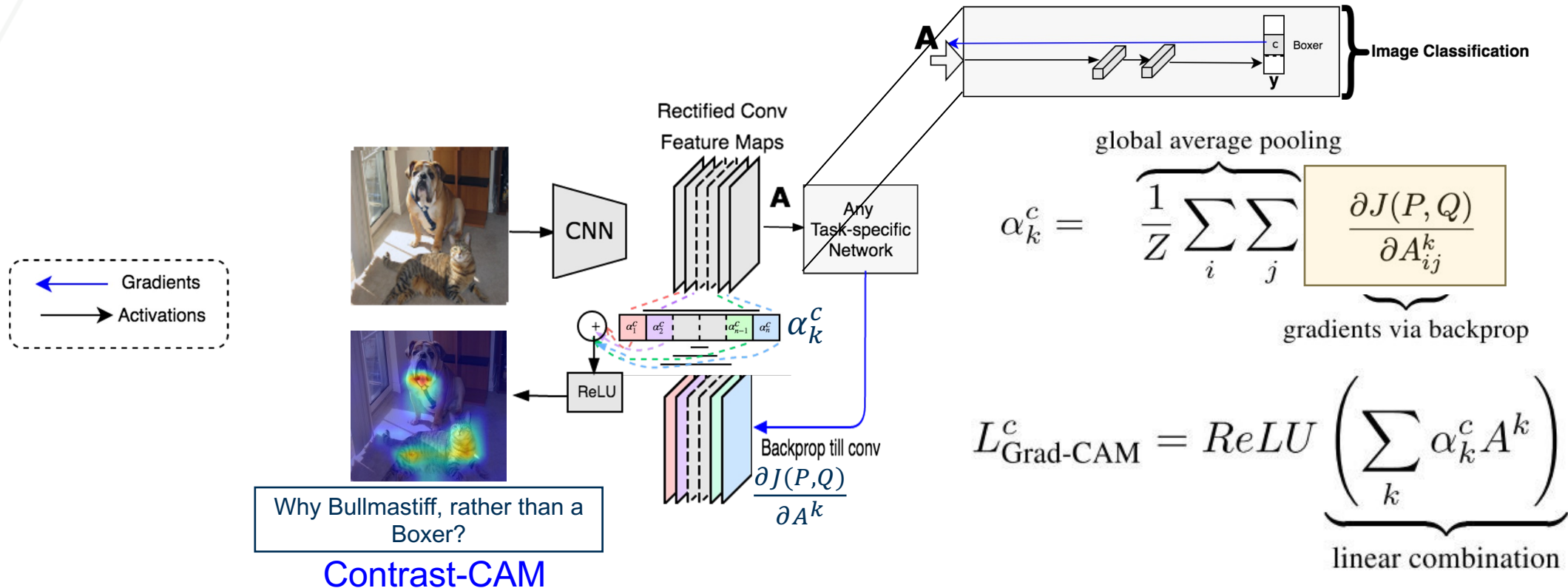
Gradient and Activation-based Explanations

ContrastCAM: Why P, rather than Q?



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

In GradCAM, backward pass the **loss between predicted class P and some contrast class Q** to last conv layer



Backpropagating the loss highlights the differences between classes P and Q.



Gradient and Activation-based Explanations

Results from GradCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? Why not Spoonbill, with 100% confidence?
ImageNet dataset : Bull Mastiff	Grad-CAM : Why : Bull Mastiff?	Representative Boxer image	Why Bull Mastiff, rather than Boxer?	Representative Blue jay image	Why Bull Mastiff, rather than Blue jay? Why not Bull Mastiff, with 100% confidence?
CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop? Why not No-Left with 100% confidence?
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6? Why not Bugatti with 100% confidence?



Gradient and Activation-based Explanations

Results from GradCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? Why not Spoonbill, with 100% confidence?
ImageNet dataset : Bull Mastiff	Grad-CAM : Why : Bull Mastiff?	Representative Boxer image	Why Bull Mastiff, rather than Boxer?	Representative Blue jay image	Why Bull Mastiff, rather than Blue jay? Why not Bull Mastiff, with 100% confidence?
CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop? Why not No-Left with 100% confidence?
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6? Why not Bugatti with 100% confidence?

Human Interpretable



Gradient and Activation-based Explanations

Results from GradCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? Why not Spoonbill, with 100% confidence?
ImageNet dataset : Bull Mastiff	Grad-CAM : Why : Bull Mastiff?	Representative Boxer image	Why Bull Mastiff, rather than Boxer	Representative Blue jay image	Why Bull Mastiff, rather than Blue jay? Why not Bull Mastiff, with 100% confidence?
CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop? Why not No-Left with 100% confidence?
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6? Why not Bugatti with 100% confidence?

Human Interpretable

Same as Grad-CAM



Gradient and Activation-based Explanations

Results from GradCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? Why not Spoonbill, with 100% confidence?
ImageNet dataset : Bull Mastiff	Grad-CAM : Why : Bull Mastiff?	Representative Boxer image	Why Bull Mastiff, rather than Boxer	Representative Blue jay image	Why Bull Mastiff, rather than Blue jay? Why not Bull Mastiff, with 100% confidence?
CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop? Why not No-Left with 100% confidence?
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6? Why not Bugatti with 100% confidence?

Human Interpretable

Same as Grad-CAM

Not Human Interpretable



Gradient and Activation-based Explanations

Results from GradCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? Why not Spoonbill, with 100% confidence?

Human Interpretable

Same as Grad-CAM



CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop?	Why not No-Left with 100% confidence?
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6?	Why not Bugatti with 100% confidence?



Gradient and Activation-based Explanations

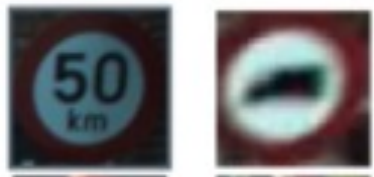
Results from GradCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? with 100% confidence?

Only traffic sign with a straight bottom-left edge – enough to say 'Not STOP Sign'

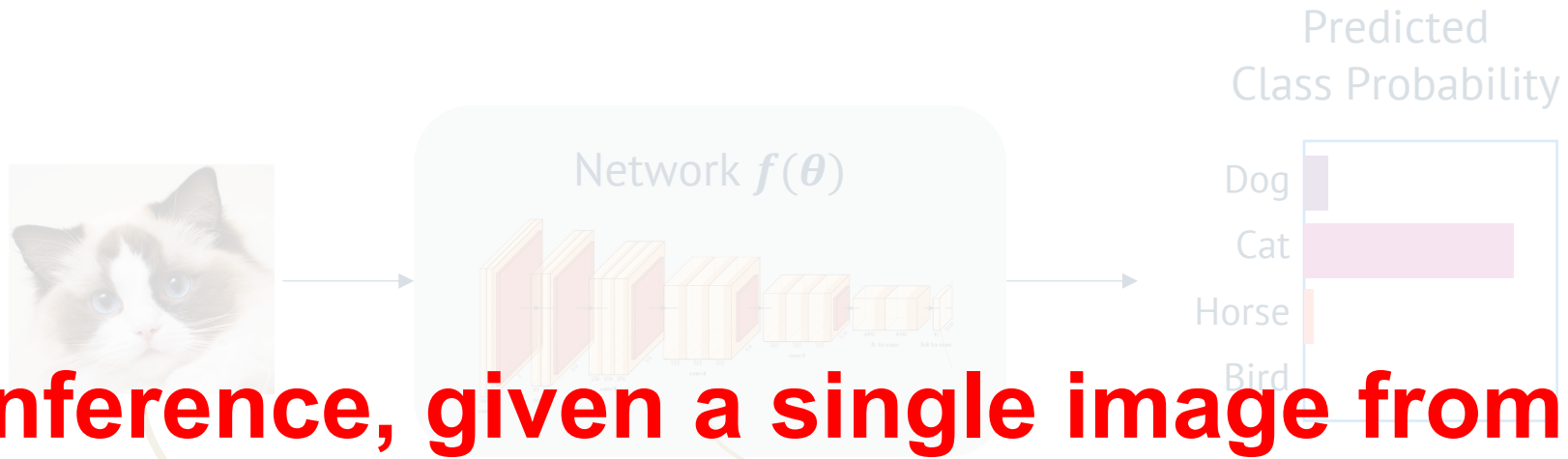


CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop?	Why not No-Left with 100% confidence?	
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6?	Why not Bugatti with 100% confidence?	

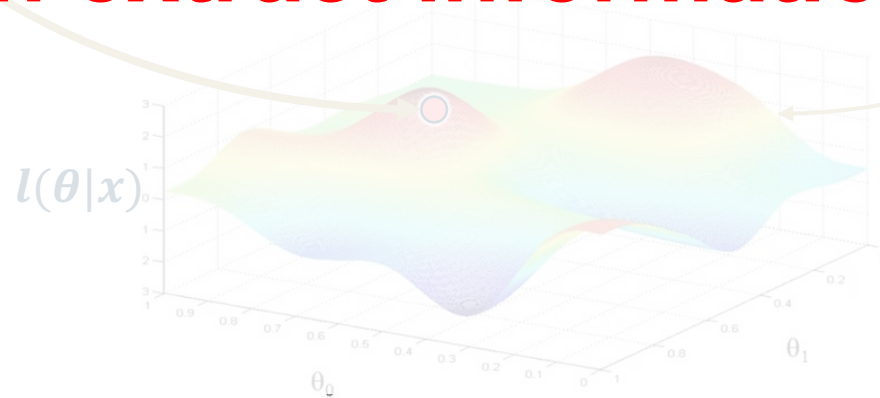


A Callback...

Information at Inference



At inference, given a single image from a single class, we can extract information about other classes



Likelihood function

$$I(\theta) = \text{Var}\left(\frac{\partial}{\partial \theta} l(\theta|x)\right)$$

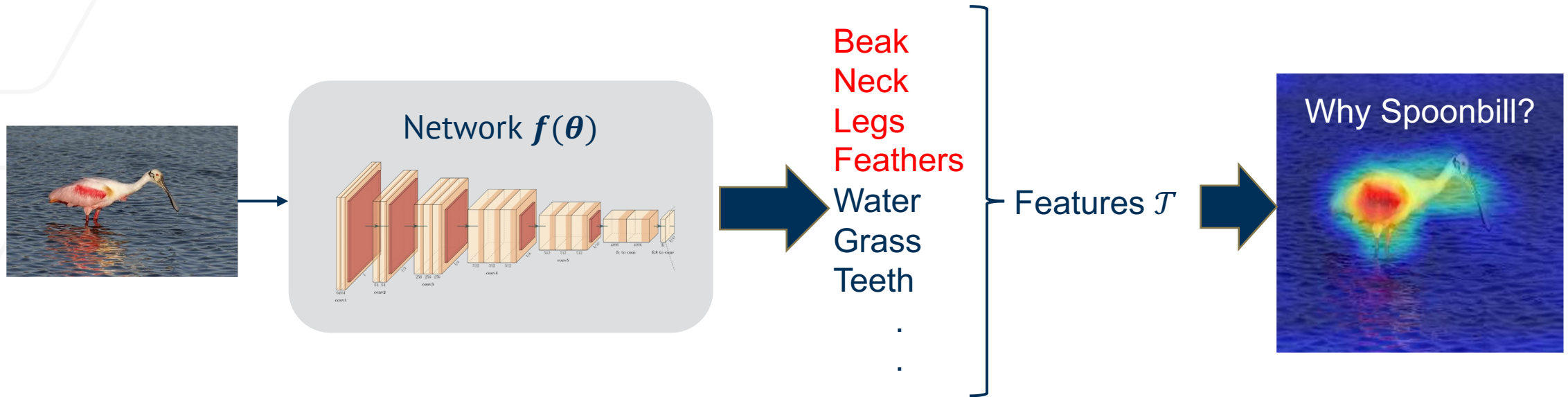
θ = Statistic of distribution
 $l(\theta | x)$ = Likelihood function



Information at Inference

Case Study: Explainability

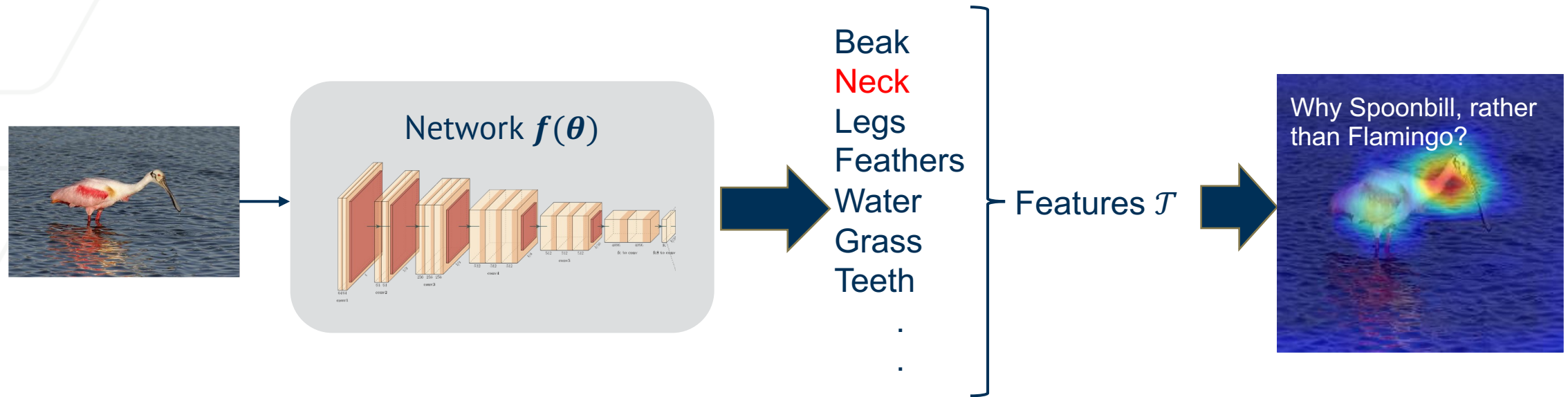
\mathcal{T} is the set of all features learned by a trained network



Information at Inference

Case Study: Explainability

Given only an image of a spoonbill, we can extract information about a Flamingo



All the requisite Information is stored within $f(\theta)$

Goal: To extract and quantify this information at inference

Robust Neural Networks

Part 3: Uncertainty at Inference



Objective

Objective of the Tutorial

To discuss methodologies that promote robustness in neural networks at inference

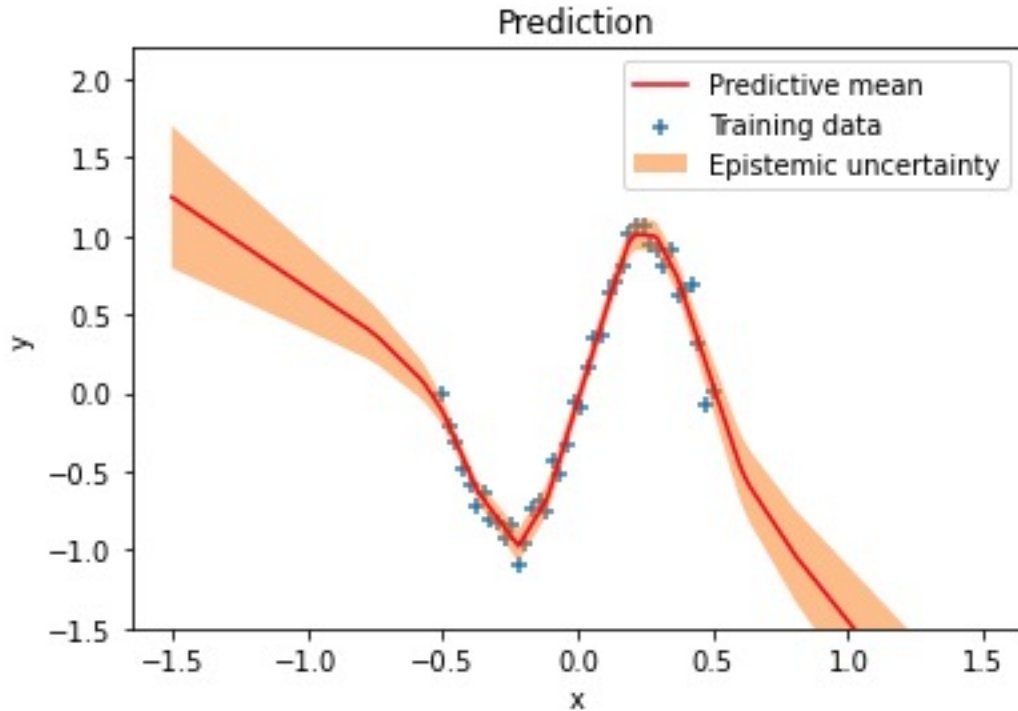
- Part 1: Inference in Neural Networks
- Part 2: Explainability at Inference
- **Part 3: Uncertainty at Inference**
 - Uncertainty Definition
 - Uncertainty Quantification
 - Gradient-based Uncertainty
 - Adversarial and Corruption Detection
- Part 4: Intervenability at Inference
- Part 5: Conclusions and Future Directions



Uncertainty

What is Uncertainty?

Uncertainty is a model knowing that it does not know



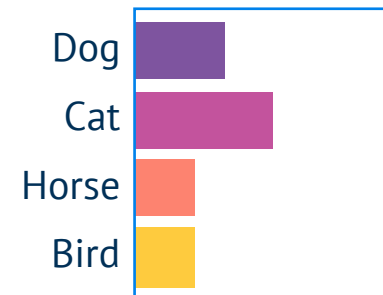
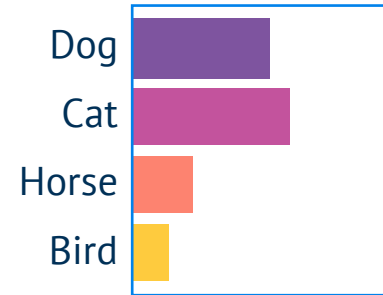
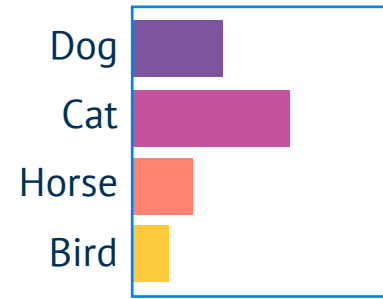
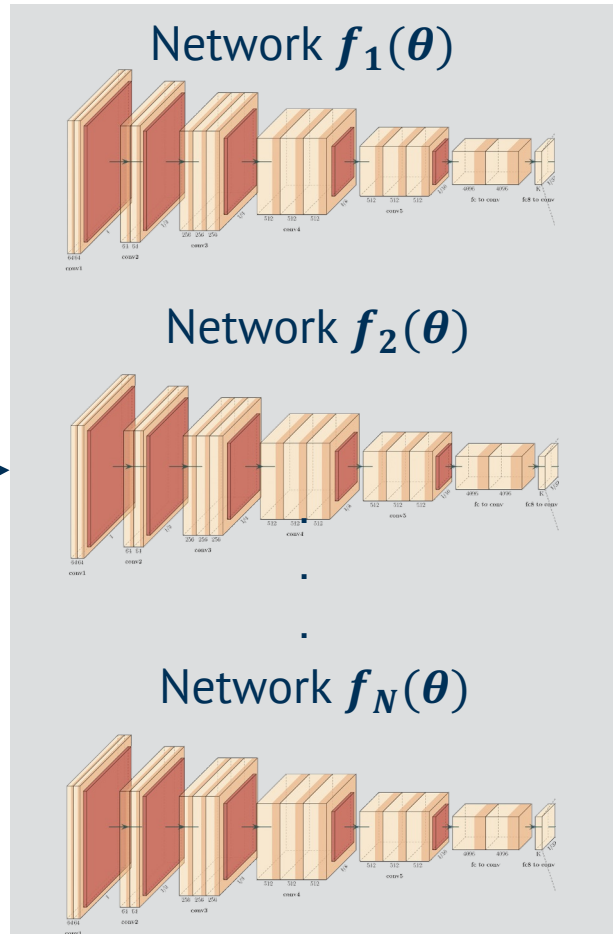
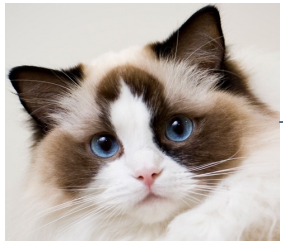
A simple example:

- When training data is **available**: **Less uncertainty**
- When training data is **unavailable**: **More uncertainty**

Uncertainty

Uncertainty Quantification in Neural Networks

Via Ensembles¹



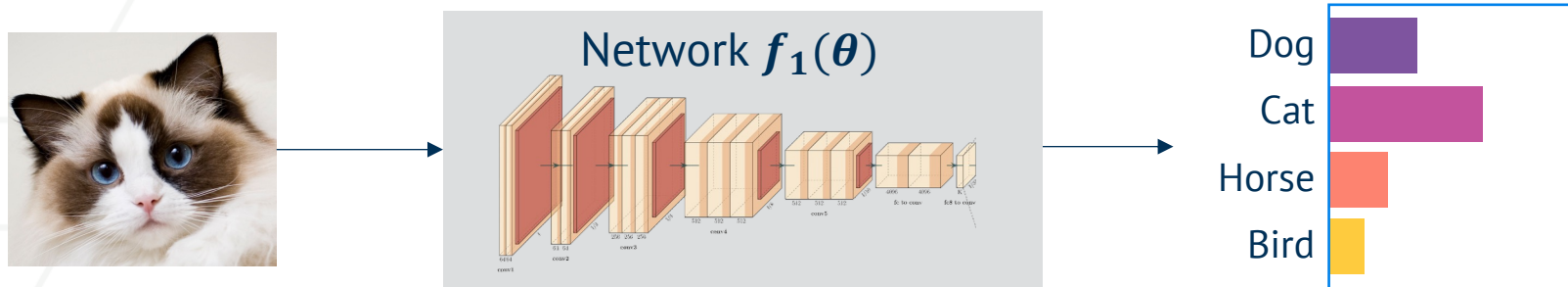
Variation within outputs $Var(y)$ is the uncertainty. Commonly referred to as **Prediction Uncertainty.**



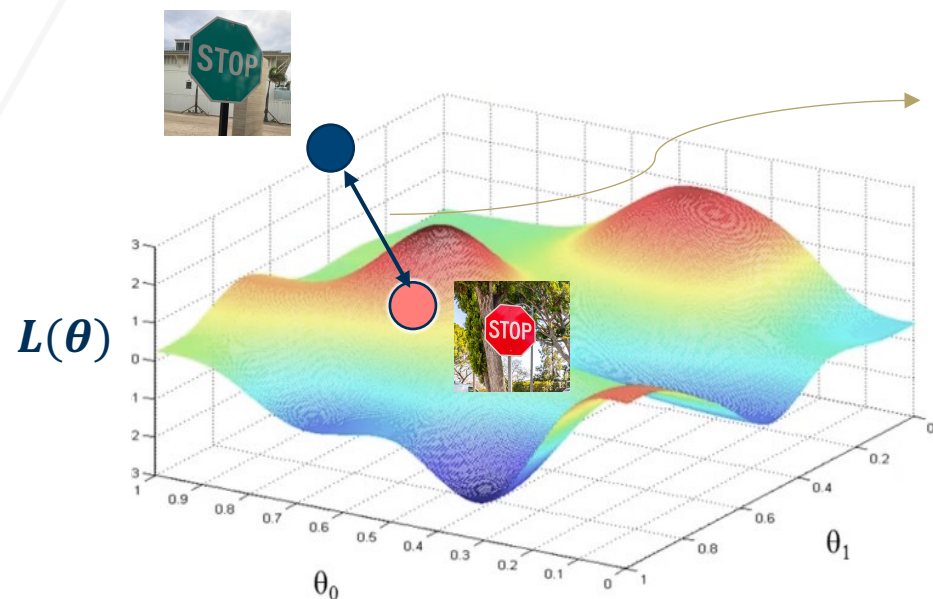
Uncertainty

Uncertainty Quantification in Neural Networks

Via Single pass methods¹



Uncertainty quantification using a single network and a single pass



Calculate distance from some trained clusters

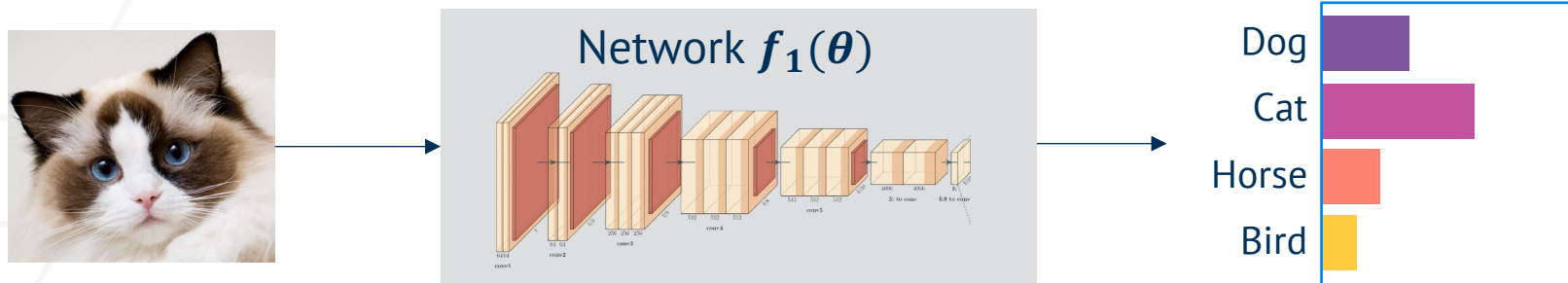
Does not require multiple networks!



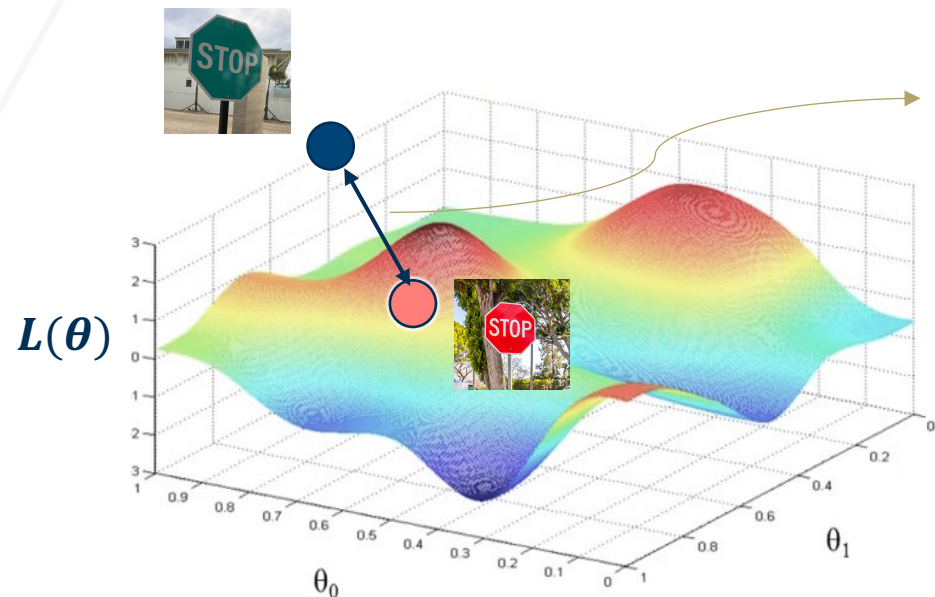
Uncertainty

Gradients as Single pass Features

Our Goal: Use gradients to characterize the novel data at Inference



Uncertainty quantification using a single network and a single pass



Calculate distance from some trained clusters

Does not require multiple networks!

Challenge: Class and prediction cannot be trusted!

Uncertainty

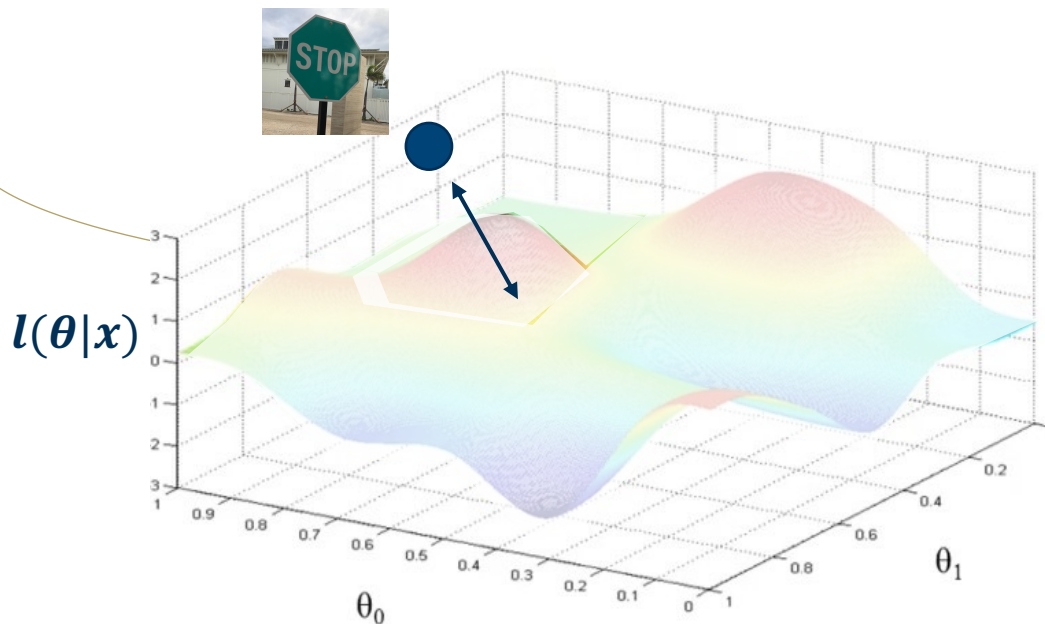
Gradients as Single pass Features

Our Goal: Use gradients to characterize the novel data at Inference, without global information

Distance from unknown cluster

Two techniques:

1. **Gradient constraints during Training for Anomaly Detection**
2. Backpropagating Confounding labels for Out-of-Distribution Detection





Backpropagated Gradient Representations for Anomaly Detection



Gukyeong Kwon, PhD
Amazon AWS



Mohit Prabhushankar, PhD
Postdoc, Georgia Tech



Ghassan AlRegib, PhD
Professor, Georgia Tech

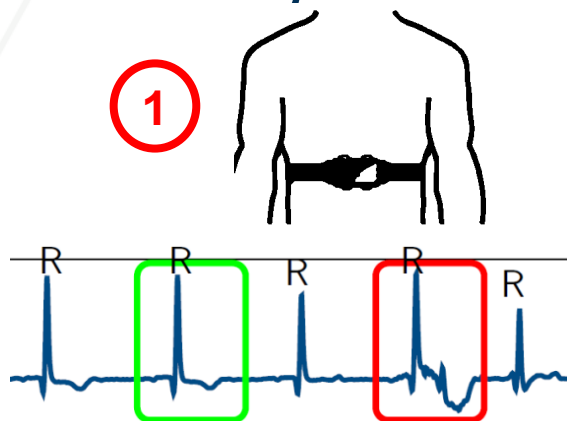


Anomalies

Finding Rare Events in Normal Patterns



'Anomalies are patterns in data that do not conform to a well defined notion of normal behavior' [1]

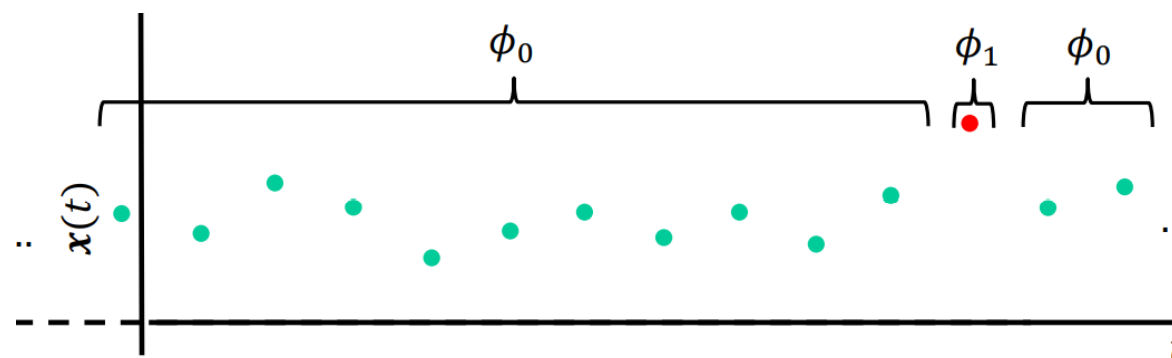


Statistical Definition:

- Normal data are generated from a stationary process P_N
- Anomalies are generated from a different process $P_A \neq P_N$

Goal: Detect ϕ_1

$$x(t) = \begin{cases} \phi_0 & \text{Normal data} \\ \phi_1 & \text{Anomalies} \end{cases}$$



Anomalies

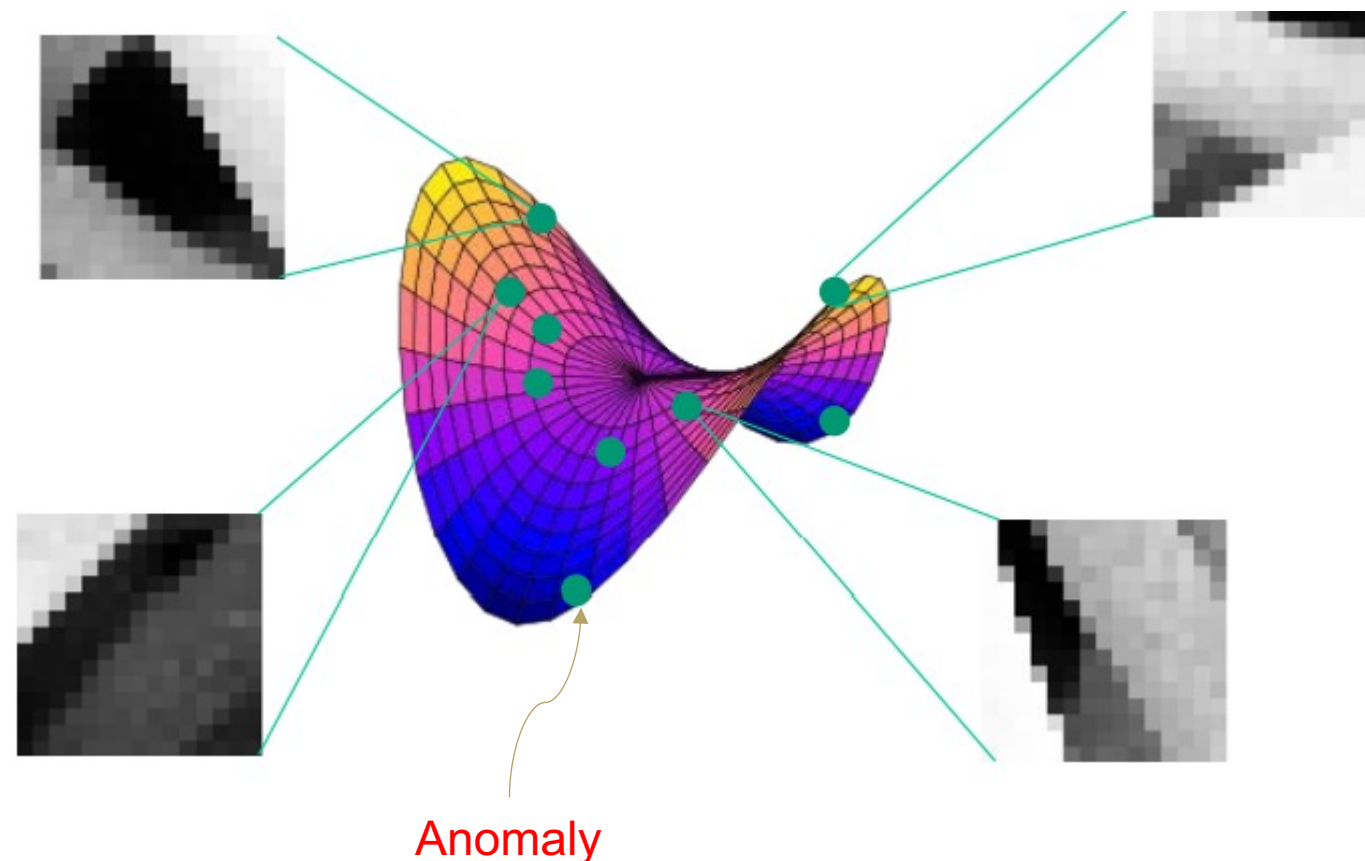
Steps for Anomaly Detection



Backpropagated Gradient
Representations for Anomaly Detection

Step 1: Constrain manifolds, Step 2: Detect statistically implausible projections

- Step 1 ensures that patches from natural images live close to a low dimensional manifold
- Step 2 designs distance functions that detect *implausibility* based on constraints



Constraining Manifolds

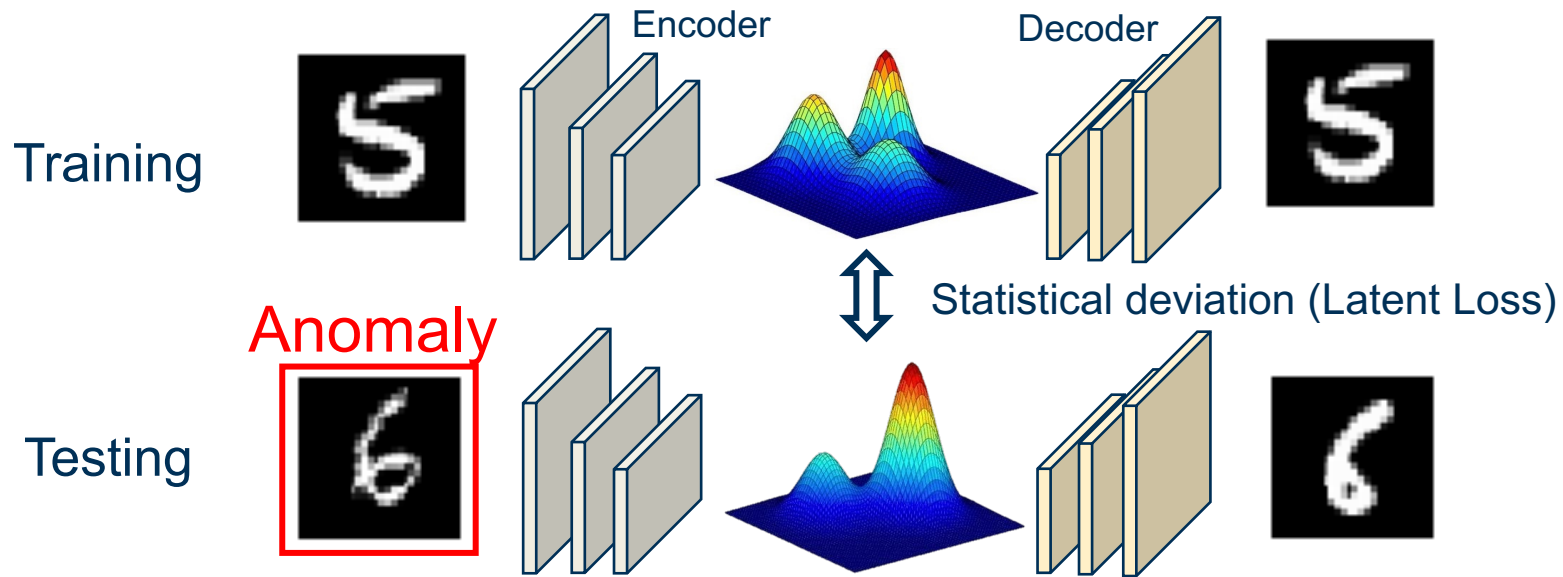
General Constraints



Backpropagated Gradient Representations for Anomaly Detection

Constrained Representation

Activations are constrained using GANs, VAEs, etc.



[1] David MJ Tax and Robert PW Duin. Support vector data description. Machine learning, 54(1):45–66, 2004.
 [2] Yaxiang Fan, Gongjian Wen, Deren Li, Shaohua Qiu, and Martin D Levine. Video anomaly detection and localization via gaussian mixture fully convolutional variational autoencoder. arXiv preprint arXiv:1805.11223, 2018. 1, 2
 [3] S. Pidhorskyi, R. Almohsen, and G. Doretto, “Generative probabilistic novelty detection with adversarial autoencoders,” in Advances in Neural Information Processing Systems, 2018, pp. 6822–6833.
 [4] D. Abati, A. Porrello, S. Calderara, and R. Cucchiara, “Latent space autoregression for novelty detection,” in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 481–490.



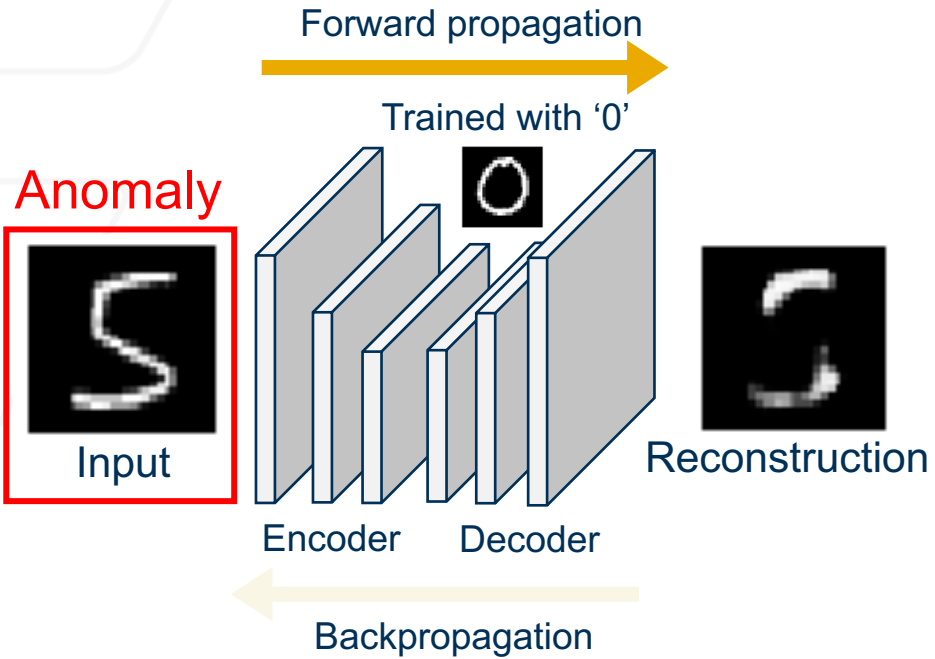
Constraining Manifolds

Gradient-based Constraints



Backpropagated Gradient Representations for Anomaly Detection

Activation Constraints



Activation-based representation
(Data perspective)

e.g. Reconstruction error (\mathcal{L})

How much of the **input** does not correspond to the **learned information**?

Gradient Constraints

Gradient-based Representation
(**Model** perspective)

How much **model update** is required by the input?

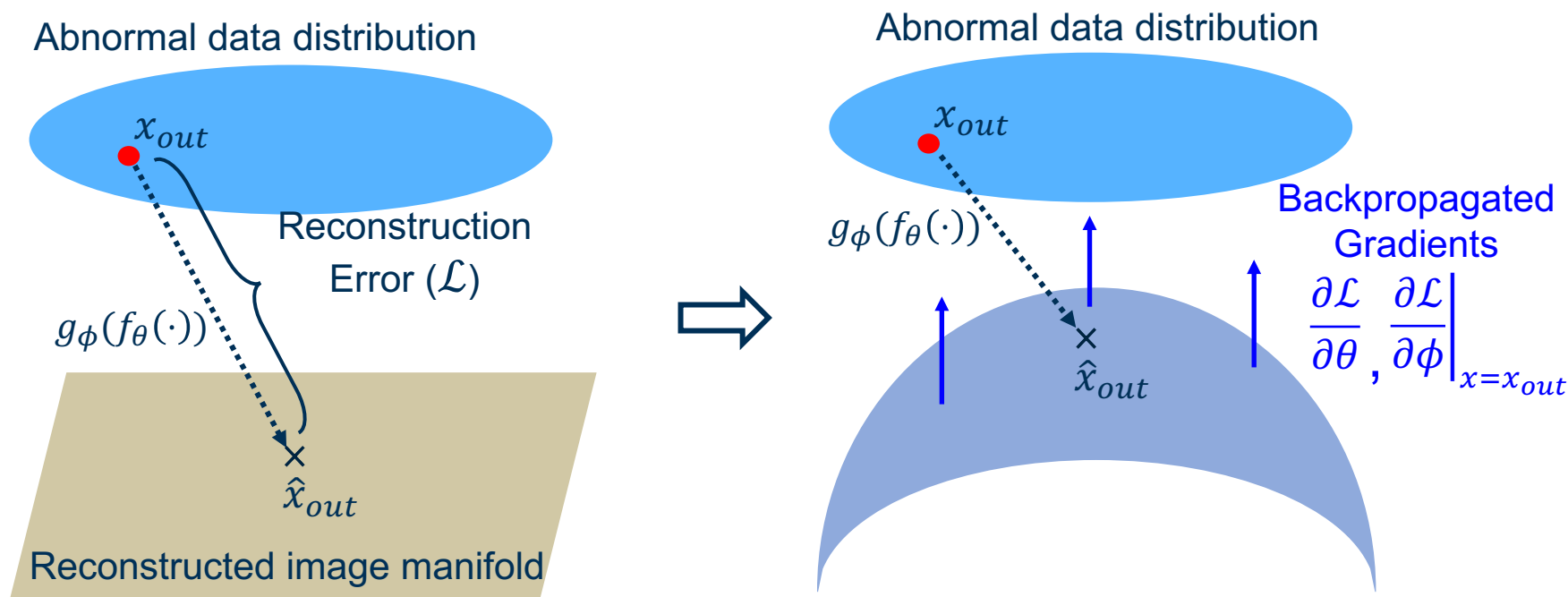


Constraining Manifolds

Advantages of Gradient-based Constraints



- Gradients provide **directional information** to characterize anomalies
- Gradients from different layers capture **abnormality at different levels of data abstraction**



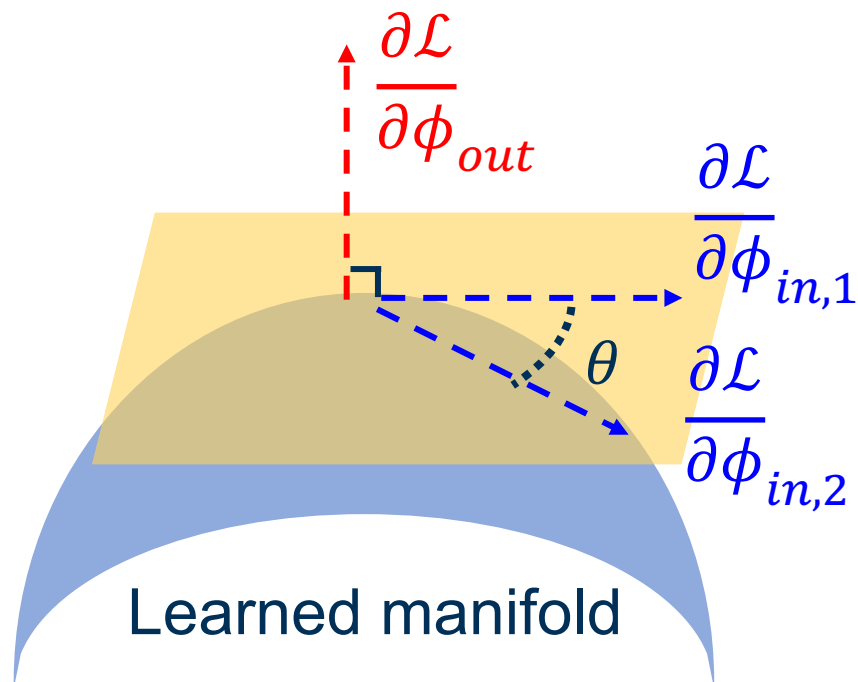
GradCON: Gradient Constraint

Gradient-based Constraints



Backpropagated Gradient Representations for Anomaly Detection

Constrain gradient-based representations during training to obtain **clear separation** between normal data and abnormal data



Learned manifold

ϕ : Weights \mathcal{L} : Reconstruction error

At k -th step of training,

Gradient loss

$$J = \mathcal{L} - \mathbb{E}_i \left[\text{cosSIM} \left(\frac{\partial J^{k-1}}{\partial \phi_{i_{avg}}}, \frac{\partial \mathcal{L}^k}{\partial \phi_i} \right) \right]$$

Avg. training gradients until $(k-1)$ th iter.

Gradients at k -th iter.

where

$$\frac{\partial J^{k-1}}{\partial \phi_{i_{avg}}} = \sum_{t=1}^{k-1} \frac{\partial J^t}{\partial \phi_i}$$



GradCON: Gradient Constraint

Activations vs Gradients

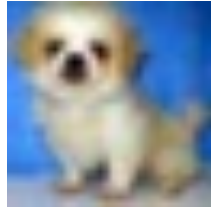


Backpropagated Gradient
Representations for Anomaly Detection

AUROC Results

Abnormal “class”
detection (CIFAR-10)

e.g.



Normal

Abnormal

Model	Loss	Plane	Car	Bird	Cat	Deer	Dog	Frog	Horse	Ship	Truck	Average
CAE	Recon	0.682	0.353	0.638	0.587	0.669	0.613	0.495	0.498	0.711	0.390	0.564
CAE	Recon	0.659	0.356	0.640	0.555	0.695	0.554	0.549	0.478	0.695	0.357	0.554
+ Grad	Grad	0.752	0.619	0.622	0.580	0.705	0.591	0.683	0.576	0.774	0.709	0.661
VAE	Recon	0.553	0.608	0.437	0.546	0.393	0.531	0.489	0.515	0.552	0.631	0.526
	Latent	0.634	0.442	0.640	0.497	0.743	0.515	0.745	0.527	0.674	0.416	0.583
VAE	Recon	0.556	0.606	0.438	0.548	0.392	0.543	0.496	0.518	0.552	0.631	0.528
	Latent	0.586	0.396	0.618	0.476	0.719	0.474	0.698	0.537	0.586	0.413	0.550
+ Grad	Grad	0.736	0.625	0.591	0.596	0.707	0.570	0.740	0.543	0.738	0.629	0.647

Recon: Reconstruction error, Latent: Latent loss, Grad: Gradient loss

- (CAE vs. CAE + Grad) Effectiveness of the gradient constraint
- (CAE vs. VAE) Performance sacrifice from the latent constraint
- (VAE vs. VAE + Grad) Complementary features from the gradient constraint



GradCON: Gradient Constraint

Aberrant Condition Detection



Backpropagated Gradient Representations for Anomaly Detection

Abnormal “condition” detection (CURE-TSR)

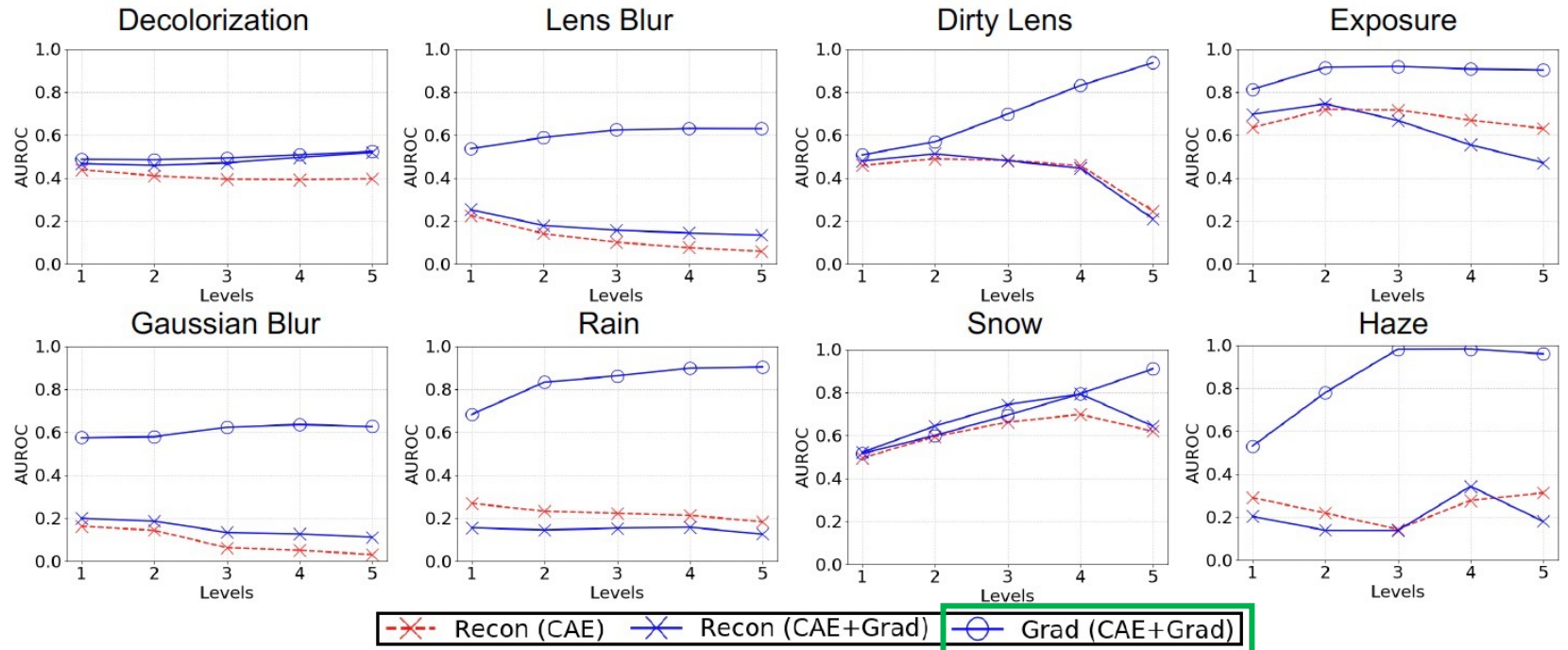


Normal



Abnormal

AUROC Results



Recon: Reconstruction error, Grad: Gradient loss



Uncertainty

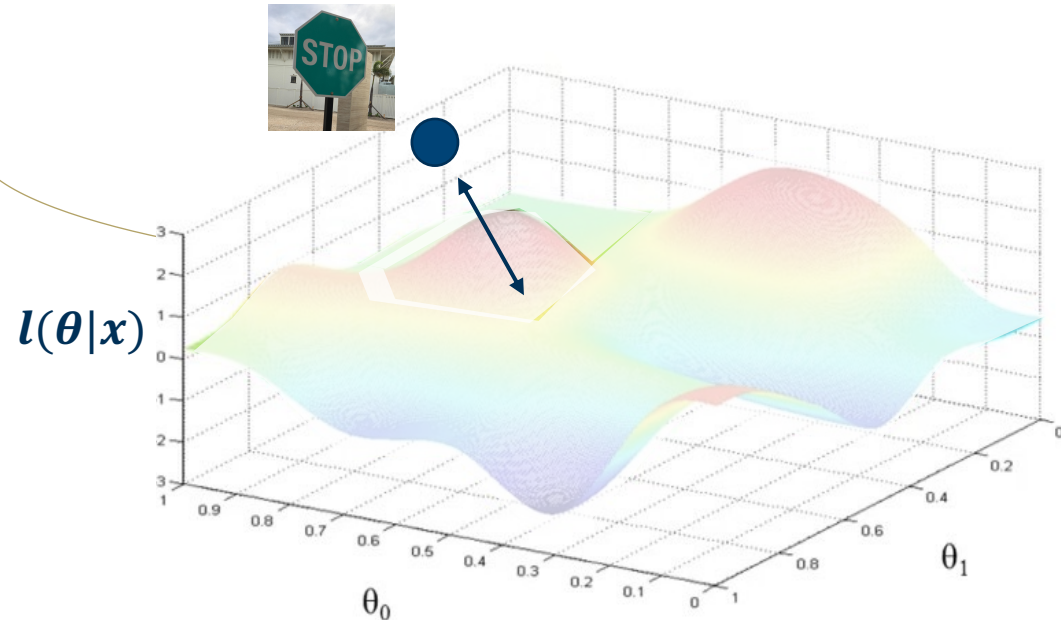
Gradients as Single pass Features

Our Goal: Use gradients to characterize the novel data at Inference, without global information

Distance from unknown cluster

Two techniques:

1. Gradient constraints during Training for Anomaly Detection
2. **Backpropagating Confounding labels for Out-of-Distribution Detection**





Probing the Purview of Neural Networks via Gradient Analysis



Jinsol Lee,
PhD Candidate



Mohit Prabhushankar, PhD
Postdoc



Ghassan AlRegib, PhD
Professor



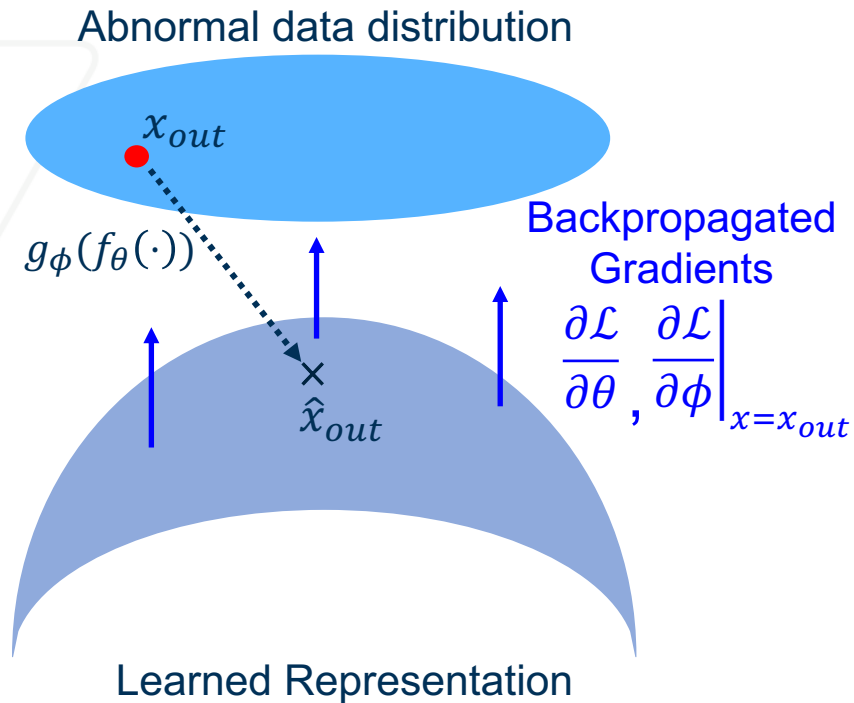
Uncertainty in Neural Networks

Principle



Probing the Purview of Neural Networks via Gradient Analysis

Principle: Gradients provide a distance measure between the learned representations space and novel data



However, what is \mathcal{L} ?

- In anomaly detection, the loss was between the input and its reconstruction
- In prediction tasks, there is neither the reconstructed input nor ground truth



Uncertainty in Neural Networks

Principle



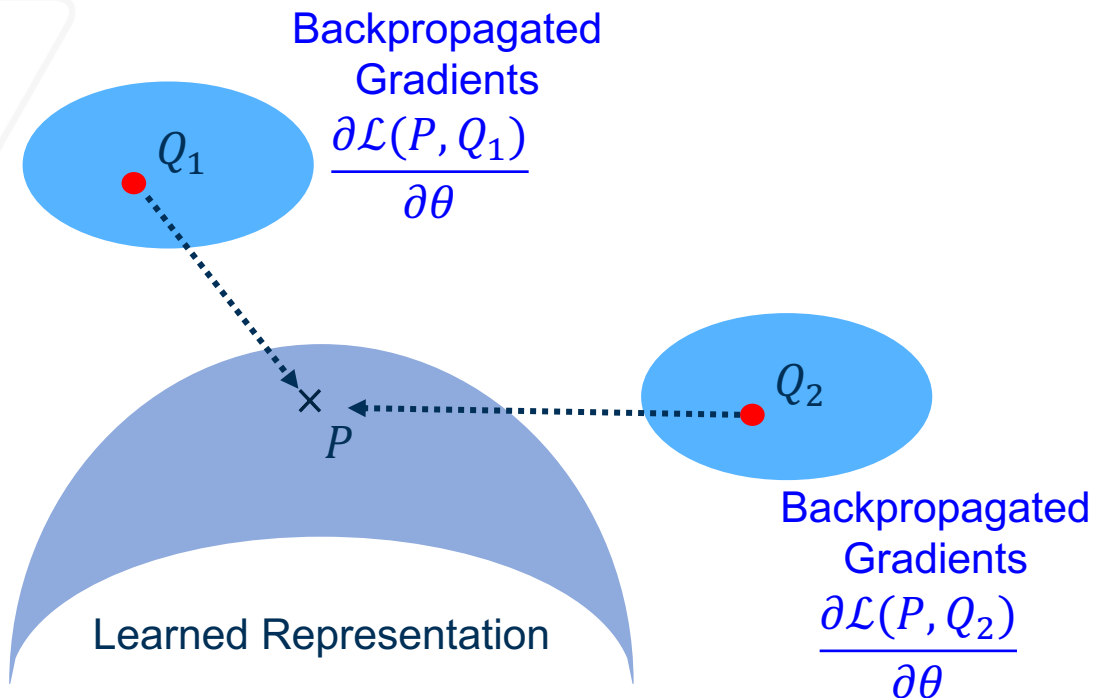
Probing the Purview of Neural Networks via Gradient Analysis

Principle: Gradients provide a distance measure between the learned representations space and novel data

P = Predicted class

Q_1 = Contrast class 1

Q_2 = Contrast class 2



However, what is \mathcal{L} ?

- In anomaly detection, the loss was between the input and its reconstruction
- In prediction tasks, there is neither the reconstructed input nor ground truth
- **We backpropagate all contrast classes - $Q_1, Q_2 \dots Q_N$ by backpropagating N one-hot vectors**
- Higher the distance, higher the uncertainty score



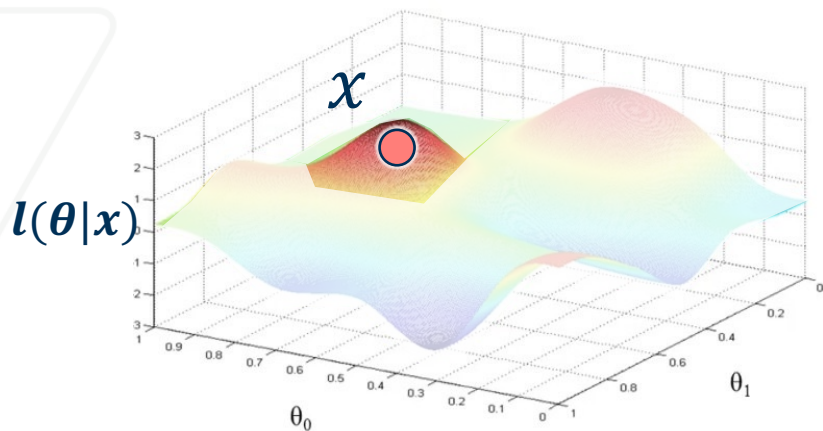
Toy Manifold Example

What is uncertainty?



Probing the Purview of Neural Networks via Gradient Analysis

Gradients represent the local required change in manifold



Contrast class 1



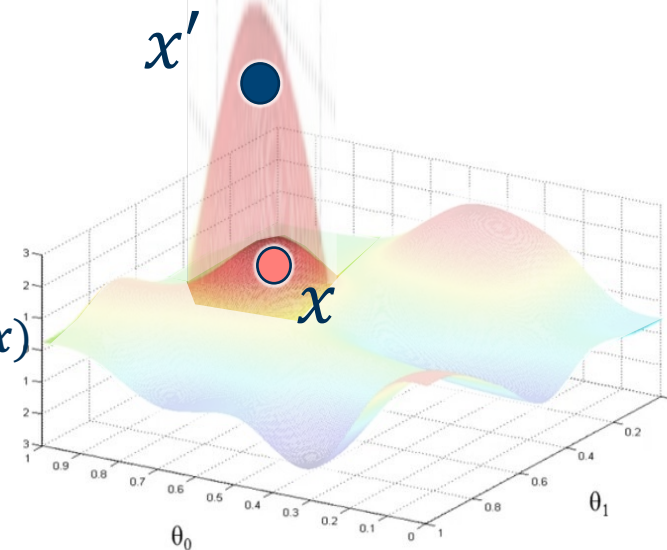
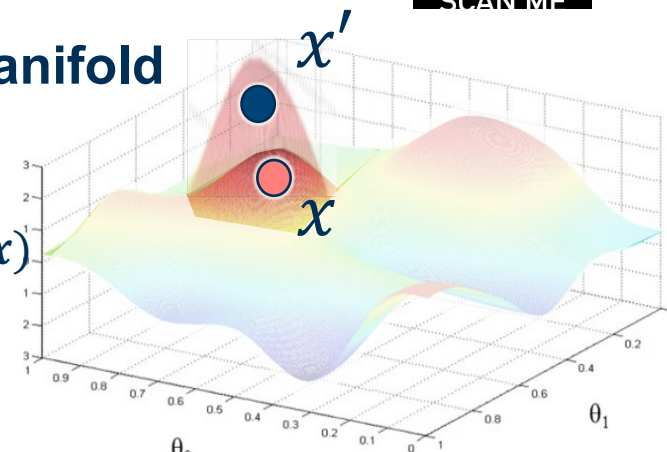
$l(\theta|x)$

·
·
·

Contrast class N



$l(\theta|x)$



- Gradients provide the necessary change in manifold that would predict the novel data 'correctly'.
- Correctly means contrastively (or incorrectly)!



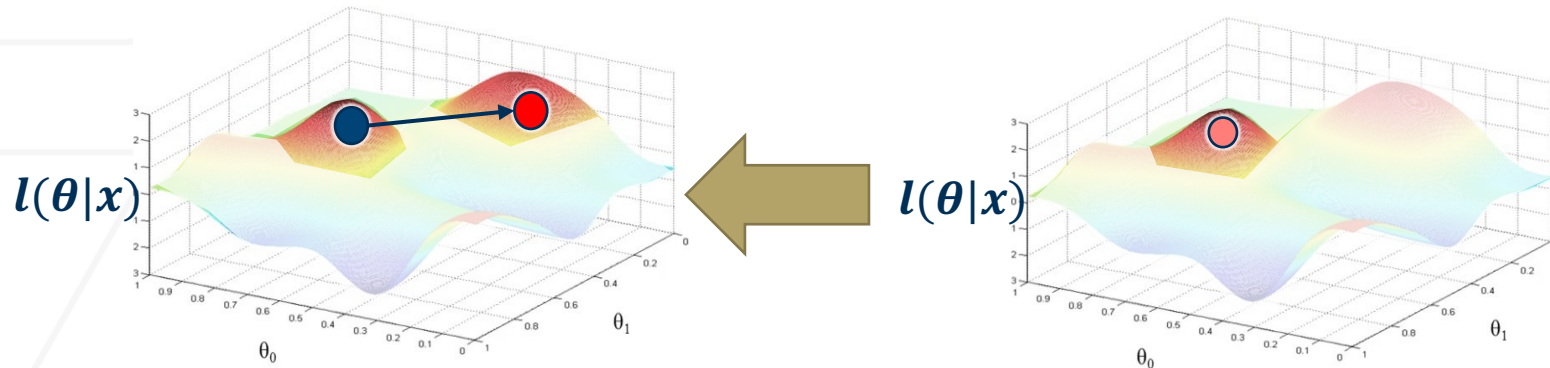
Toy Manifold Example

How is this different from Explainability?



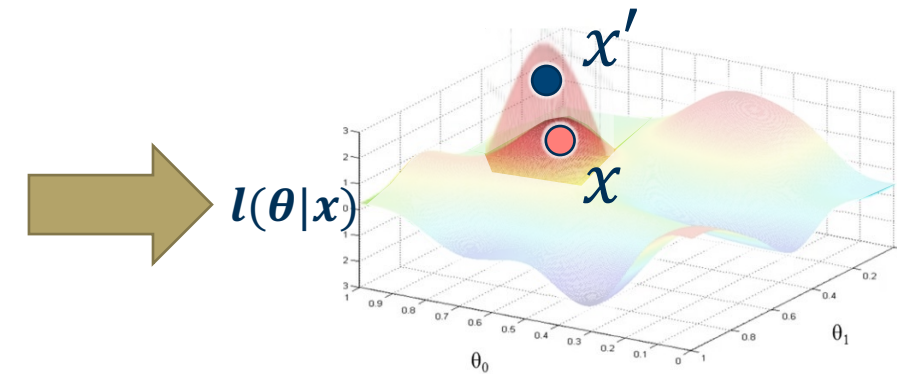
Probing the Purview of Neural Networks via Gradient Analysis

Part 2: Explainability



- In Part 2: Activations of learned manifold are weighted by gradients w.r.t. activations to extract information and provide explanations

Part 3: Uncertainty



- In Part 3: Statistics of gradients w.r.t. the weights (energy) will be directly used as features



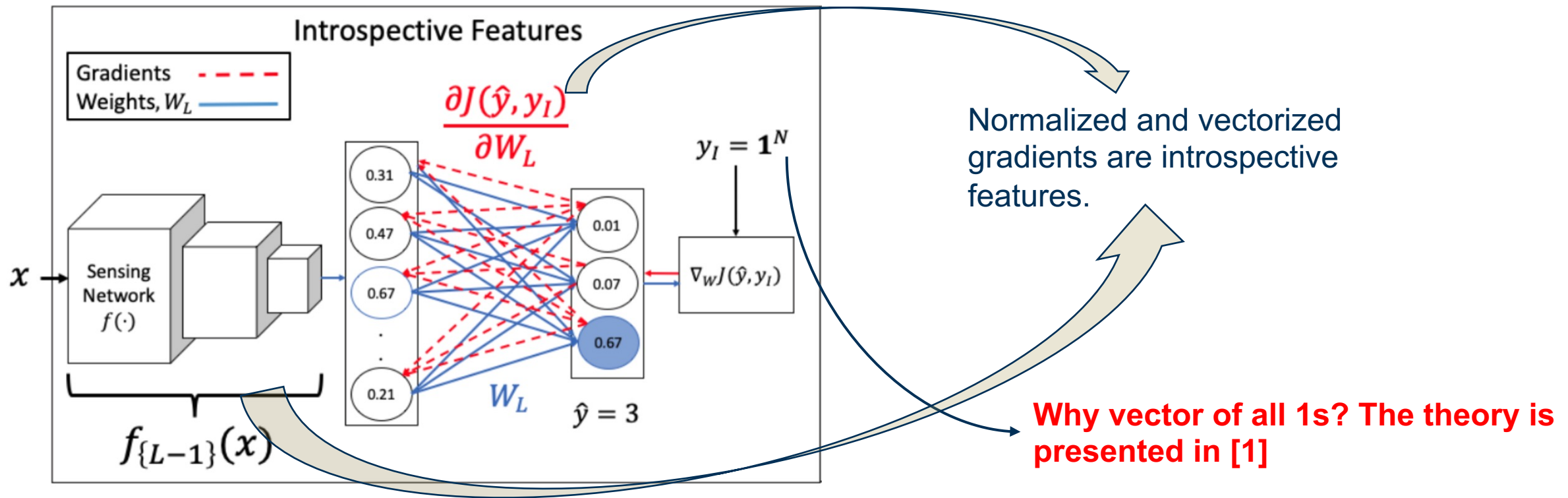
Uncertainty in Neural Networks

Deriving Gradient Features



Probing the Purview of Neural Networks via Gradient Analysis

Step 1: Measure the loss between the prediction P and a vector of all ones and backpropagate to obtain the introspective features



Uncertainty in Neural Networks

Utilizing Gradient Features



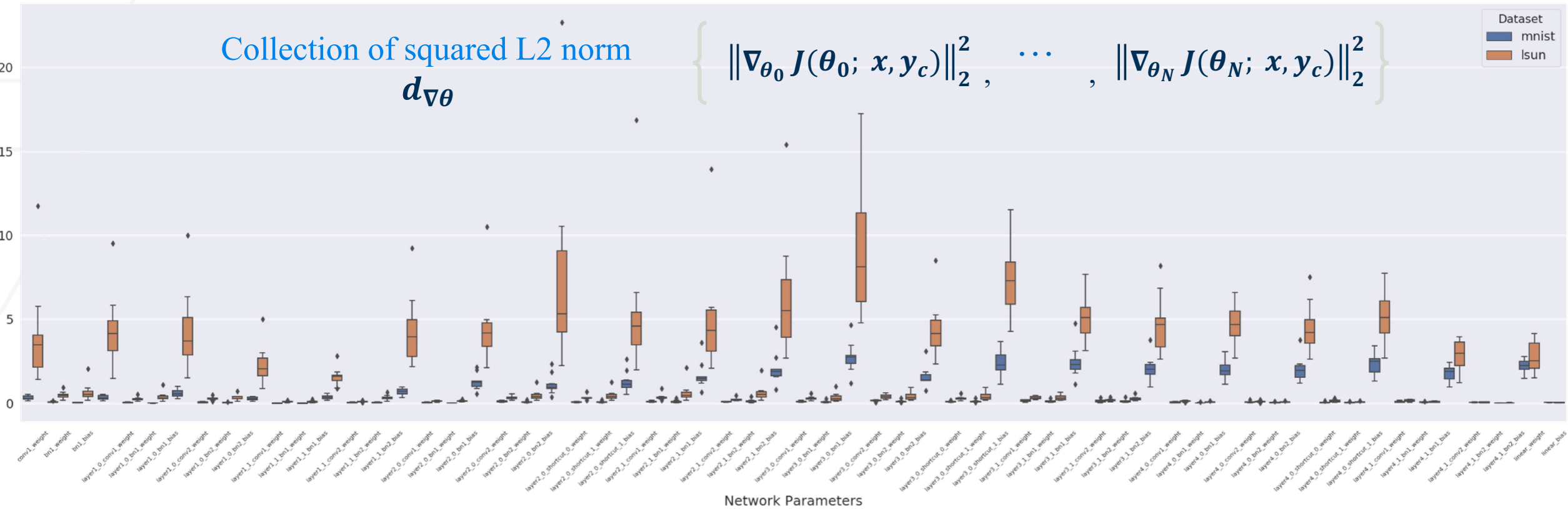
Probing the Purview of Neural Networks via Gradient Analysis

Step 2: Take L2 norm of all generated gradients

Collection of squared L2 norm
 $d_{\nabla\theta}$

$$\left\{ \|\nabla_{\theta_0} J(\theta_0; x, y_c)\|_2^2, \dots, \|\nabla_{\theta_N} J(\theta_N; x, y_c)\|_2^2 \right\}$$

Dataset
■ mnist
■ lsun



MNIST: In-distribution, SUN: Out-of-Distribution



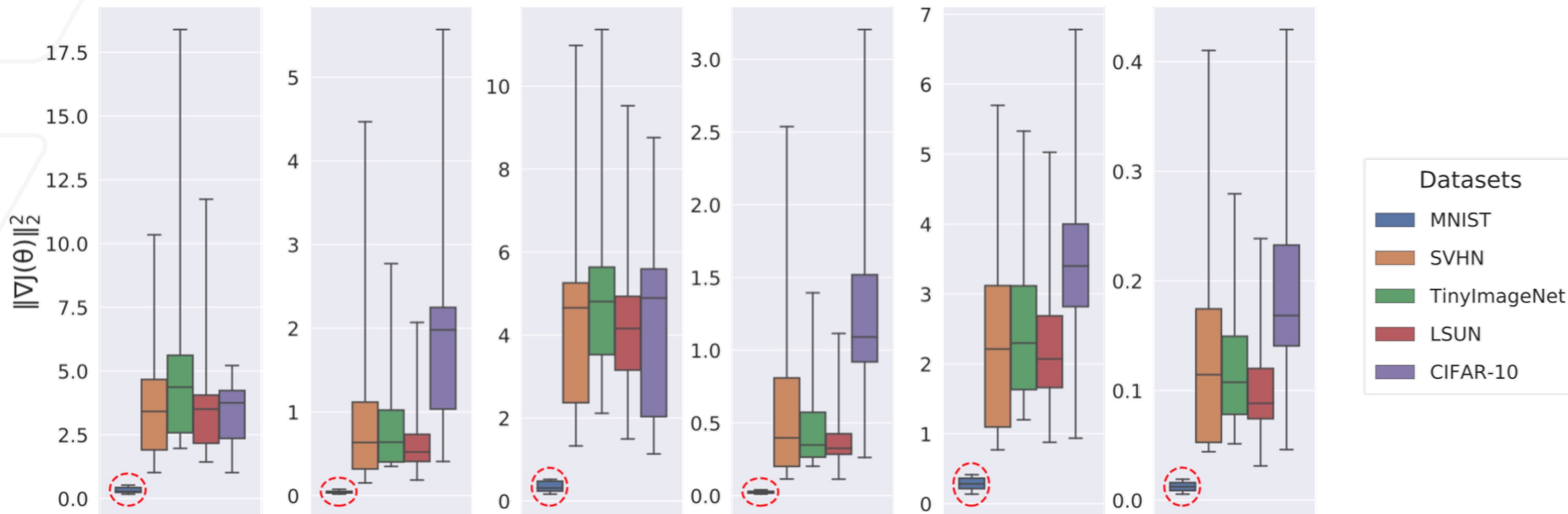
Gradient-based Uncertainty

Uncertainty in OOD Setting



Probing the Purview of Neural Networks via Gradient Analysis

Squared L2 distances for different parameter sets



MNIST: Circled in red. Significantly lower uncertainty compared to OOD datasets



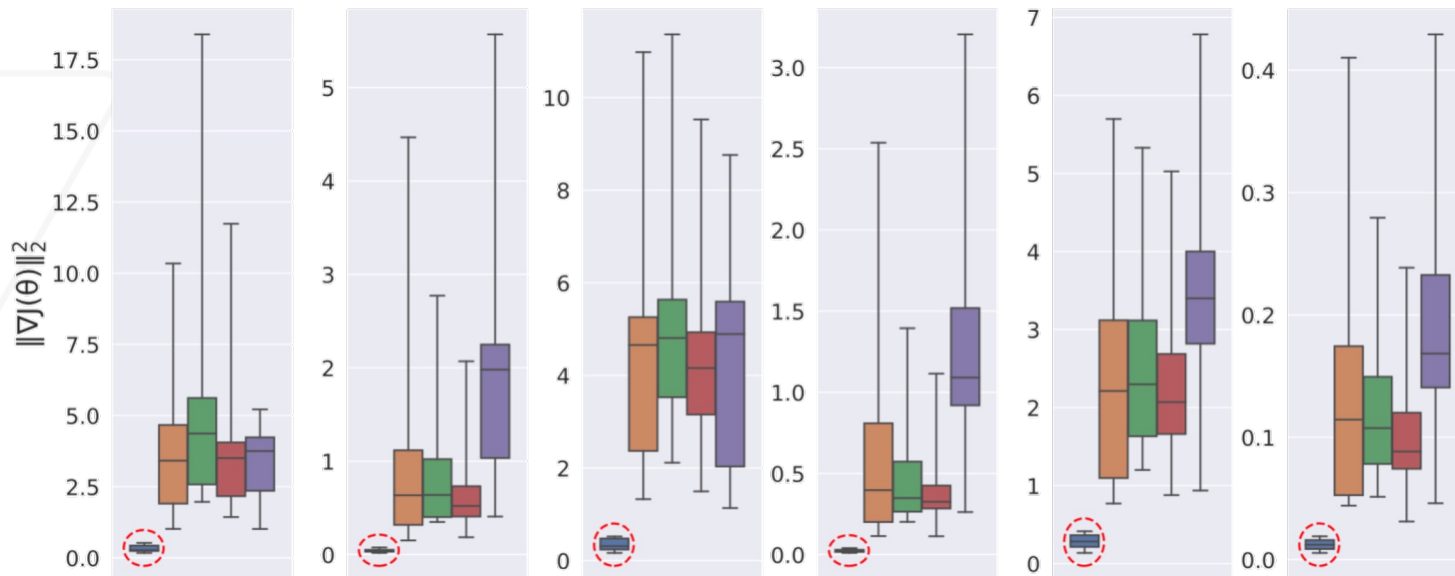
Gradient-based Uncertainty

Experimental Setup



Probing the Purview of Neural Networks
via Gradient Analysis

Utilize this discrepancy in trained vs untrained data gradient L2 distance to detect adversarial, noisy, and OOD data



- Step 1:** Train a deep network $f(\cdot)$ on some training distribution
- Step 2:** Introduce challenging (adversarial, noisy, OOD) data
- Step 3:** Derive gradient uncertainty on both trained and challenge data
- Step 4:** Train a classifier $H(\cdot)$ to detect challenging from trained data
- Step 5:** At test time, data is passed through $f(\cdot)$ and then $H(\cdot)$ to obtain a Reliability classification



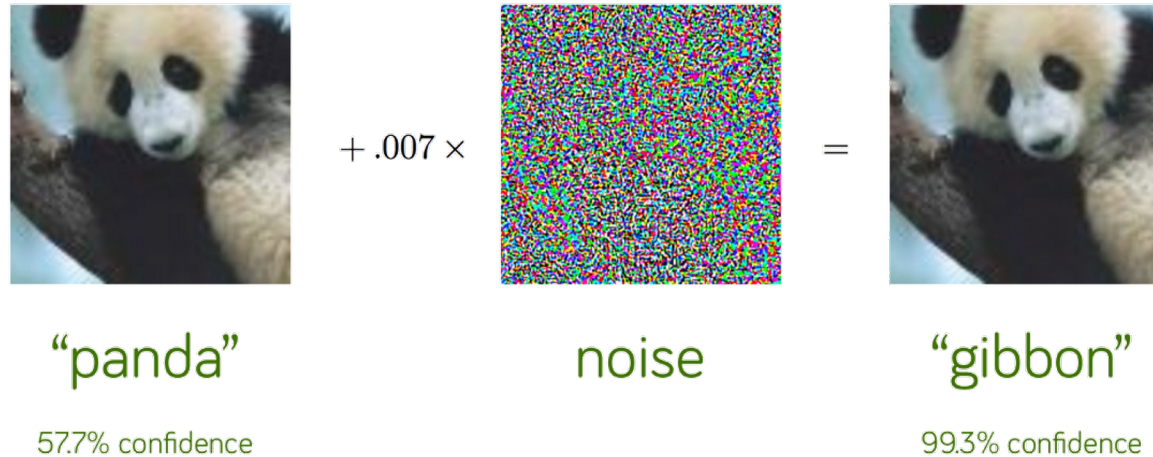
Gradient-based Uncertainty

Uncertainty in Adversarial Setting



Probing the Purview of Neural Networks via Gradient Analysis

Vulnerable DNNs in the real world



Goal: to examine the ability of trained DNNs to handle adversarial inputs during inference



Gradient-based Uncertainty

Uncertainty in Adversarial Setting



Probing the Purview of Neural Networks
via Gradient Analysis

SCAN ME

MODEL	ATTACKS	BASELINE	LID	M(V)	M(P)	M(FE)	M(P+FE)	OURS
RESNET	FGSM	51.20	90.06	81.69	84.25	99.95	99.95	93.45
	BIM	49.94	99.21	87.09	89.20	100.0	100.0	96.19
	C&W	53.40	76.47	74.51	75.71	92.78	92.79	97.07
	PGD	50.03	67.48	56.27	57.57	65.23	75.98	95.82
	ITERLL	60.40	85.17	62.32	64.10	85.10	92.10	98.17
	SEMANTIC	52.29	86.25	64.18	65.79	83.95	84.38	90.15
DENSENET	FGSM	52.76	98.23	86.88	87.24	99.98	99.97	96.83
	BIM	49.67	100.0	89.19	89.17	100.0	100.0	96.85
	C&W	54.53	80.58	75.77	76.16	90.83	90.76	97.05
	PGD	49.87	83.01	70.39	66.52	86.94	83.61	96.77
	ITERLL	55.43	83.16	70.17	66.61	83.20	77.84	98.53
	SEMANTIC	53.54	81.41	62.16	62.15	67.98	67.29	89.55



Gradient-based Uncertainty

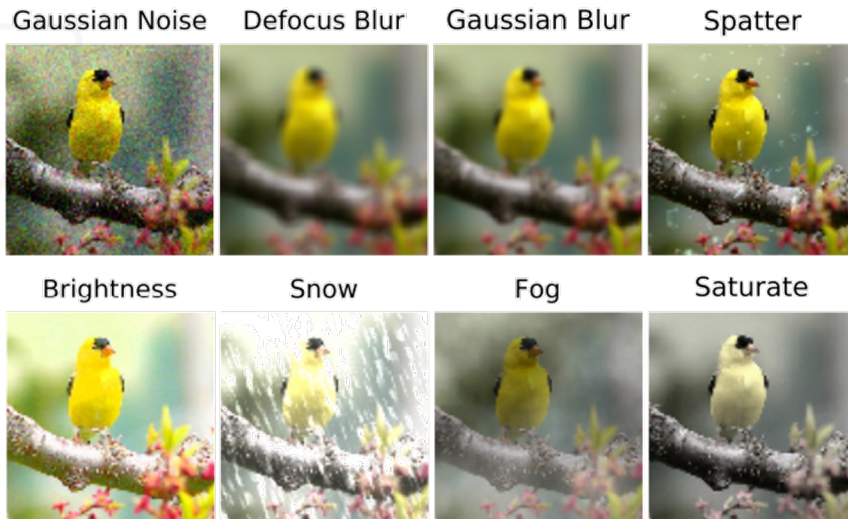
Uncertainty in Detecting Challenging Conditions



Probing the Purview of Neural Networks via Gradient Analysis

Same application as Anomaly Detection, except there is no need for an additional AE network!

CIFAR-10-C



CURE-TSR



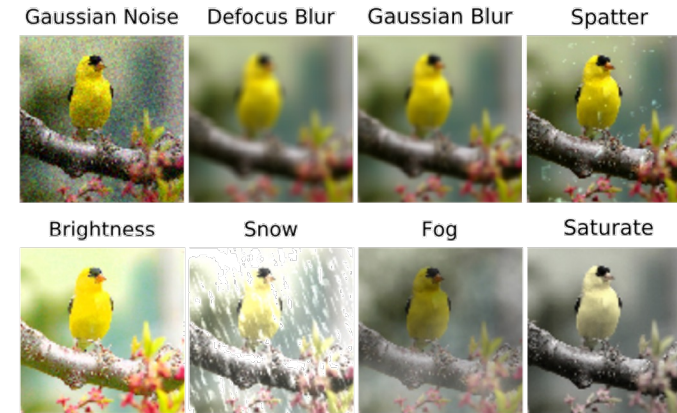
Gradient-based Uncertainty

Uncertainty in Detecting Challenging Conditions



Probing the Purview of Neural Networks via Gradient Analysis

Dataset	Method	Mahalanobis [12] / Ours				
		Corruption	Level 1	Level 2	Level 3	Level 4
CIFAR-10-C	Noise	96.63 / 99.95	98.73 / 99.97	99.46 / 99.99	99.62 / 99.97	99.71 / 99.99
	LensBlur	94.22 / 99.95	97.51 / 99.99	99.26 / 100.0	99.78 / 100.0	99.89 / 100.0
	GaussianBlur	94.19 / 99.94	99.28 / 100.0	99.76 / 100.0	99.86 / 100.0	99.80 / 100.0
	DirtyLens	93.37 / 99.94	95.31 / 99.93	95.66 / 99.96	95.37 / 99.92	97.43 / 99.96
	Exposure	91.39 / 99.87	91.00 / 99.85	90.71 / 99.88	90.58 / 99.85	90.68 / 99.87
	Snow	93.64 / 99.94	96.50 / 99.94	94.44 / 99.95	94.22 / 99.95	95.25 / 99.92
	Haze	95.52 / 99.95	98.35 / 99.99	99.28 / 100.0	99.71 / 99.99	99.94 / 100.0
	Decolor	93.51 / 99.96	93.55 / 99.96	90.30 / 99.82	89.86 / 99.75	90.43 / 99.83
CURE-TSR	Noise	25.46 / 50.20	47.54 / 63.87	47.32 / 81.20	66.19 / 91.16	83.14 / 94.81
	LensBlur	48.06 / 72.63	71.61 / 87.58	86.59 / 92.56	92.19 / 93.90	94.90 / 95.65
	GaussianBlur	66.44 / 83.07	77.67 / 86.94	93.15 / 94.35	80.78 / 94.51	97.36 / 96.53
	DirtyLens	29.78 / 51.21	29.28 / 59.10	46.60 / 82.10	73.36 / 91.87	98.50 / 98.70
	Exposure	74.90 / 88.13	99.96 / 96.78	99.99 / 99.26	100.0 / 99.80	100.0 / 99.90
	Snow	28.11 / 61.34	61.28 / 80.52	89.89 / 91.30	99.34 / 96.13	99.98 / 97.66
	Haze	66.51 / 95.83	97.86 / 99.50	100.0 / 99.95	100.0 / 99.87	100.0 / 99.88
	Decolor	48.37 / 62.36	60.55 / 81.30	71.73 / 89.93	87.29 / 95.42	89.68 / 96.91



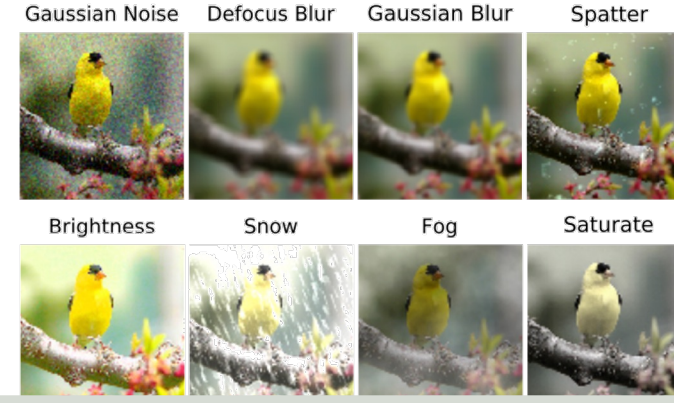
Gradient-based Uncertainty

Uncertainty in Detecting Challenging Conditions



Probing the Purview of Neural Networks via Gradient Analysis

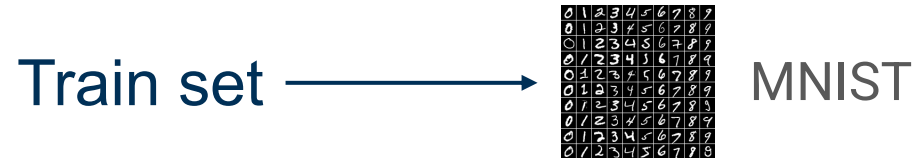
Dataset	Method	Mahalanobis [12] / Ours				
		Corruption	Level 1	Level 2	Level 3	Level 4
CIFAR-10-C	Noise	96.63 / 99.95	98.73 / 99.97	99.46 / 99.99	99.62 / 99.97	99.71 / 99.99
	LensBlur	94.22 / 99.95	97.51 / 99.99	99.26 / 100.0	99.78 / 100.0	99.89 / 100.0
	GaussianBlur	94.19 / 99.94	99.28 / 100.0	99.76 / 100.0	99.86 / 100.0	99.80 / 100.0
	DirtyLens	93.37 / 99.94	95.31 / 99.93	95.66 / 99.96	95.37 / 99.92	97.43 / 99.96
	Exposure	91.39 / 99.87	91.00 / 99.85	90.71 / 99.88	90.58 / 99.85	90.68 / 99.87
	Snow	93.64 / 99.94	96.50 / 99.94	94.44 / 99.95	94.22 / 99.95	95.25 / 99.92
	Haze	95.52 / 99.95	98.35 / 99.99	99.28 / 100.0	99.71 / 99.99	99.94 / 100.0
	Decolor	93.51 / 99.96	93.55 / 99.96	90.30 / 99.82	89.86 / 99.75	90.43 / 99.83
CURE-TSR	Noise	25.46 / 50.20	47.54 / 63.87	47.32 / 81.20	66.19 / 91.16	83.14 / 94.81
	LensBlur	48.06 / 72.63	71.61 / 87.58	86.59 / 92.56	92.19 / 93.90	94.90 / 95.65
	GaussianBlur	66.44 / 83.07	77.67 / 86.94	93.15 / 94.35	80.78 / 94.51	97.36 / 96.53
	DirtyLens	29.78 / 51.21	29.28 / 59.10	46.60 / 82.10	73.36 / 91.87	98.50 / 98.70
	Exposure	74.90 / 88.13	99.96 / 96.78	99.99 / 99.26	100.0 / 99.80	100.0 / 99.90
	Snow	28.11 / 61.34	61.28 / 80.52	89.89 / 91.30	99.34 / 96.13	99.98 / 97.66
	Haze	66.51 / 95.83	97.86 / 99.50	100.0 / 99.95	100.0 / 99.87	100.0 / 99.88
	Decolor	48.37 / 62.36	60.55 / 81.30	71.73 / 89.93	87.29 / 95.42	89.68 / 96.91



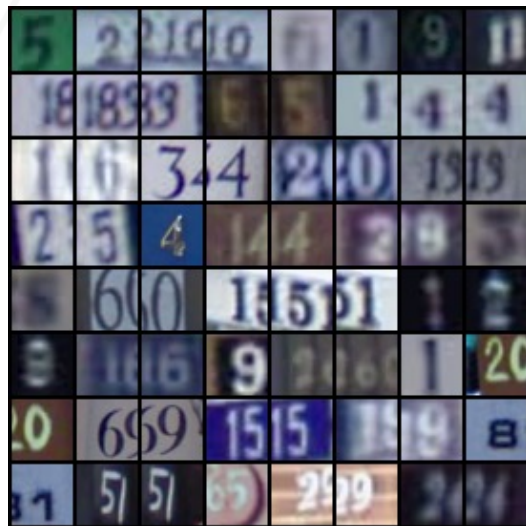
Out-of-Distribution Detection



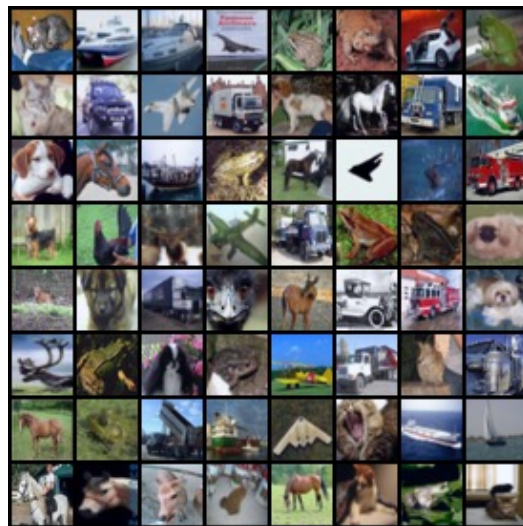
Probing the Purview of Neural Networks via Gradient Analysis



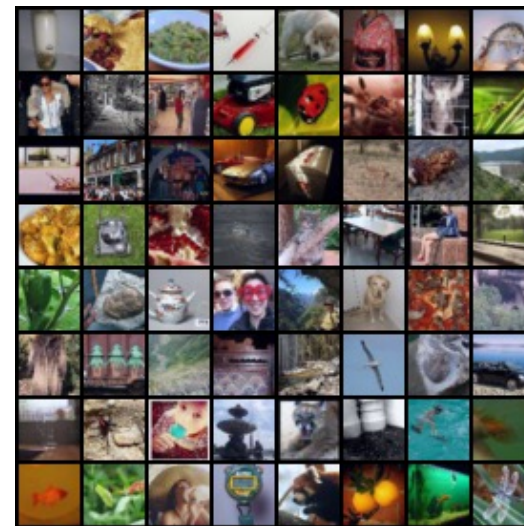
Goal: To detect that these datasets are not part of training



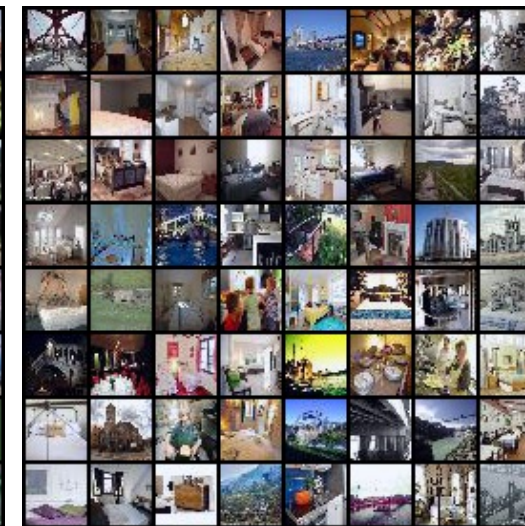
SVHN



CIFAR10



TinyImageNet



LSUN



Out-of-Distribution Detection



Probing the Purview of Neural Networks via Gradient Analysis

Dataset Distribution		Detection Accuracy	AUROC	AUPR
In	Out	Baseline [5] / ODIN [6] / Mahalanobis (V) [7] / Mahalanobis (P+FE) [7] / Ours		
CIFAR-10	SVHN	83.36 / 88.81 / 79.39 / 91.95 / 98.04	88.30 / 94.93 / 85.03 / 97.10 / 99.84	88.26 / 95.45 / 86.15 / 96.12 / 99.98
	TinyImageNet	84.01 / 85.21 / 83.60 / 97.45 / 86.17	90.06 / 91.86 / 88.93 / 99.68 / 93.18	89.26 / 91.60 / 88.59 / 99.60 / 92.66
	LSUN	87.34 / 88.42 / 85.02 / 98.60 / 98.37	92.79 / 94.48 / 90.11 / 99.86 / 99.86	92.30 / 94.22 / 89.80 / 99.82 / 99.87
SVHN	CIFAR-10	79.98 / 80.12 / 74.10 / 88.84 / 97.90	81.50 / 81.49 / 79.31 / 95.05 / 99.79	81.01 / 80.95 / 80.83 / 90.25 / 98.11
	TinyImageNet	81.70 / 81.92 / 79.35 / 96.17 / 97.74	83.69 / 83.82 / 83.85 / 99.23 / 99.77	82.54 / 82.60 / 85.50 / 98.17 / 97.93
	LSUN	80.96 / 81.15 / 79.52 / 97.50 / 99.04	82.85 / 82.98 / 83.02 / 99.54 / 99.93	81.97 / 82.01 / 84.67 / 98.84 / 99.21

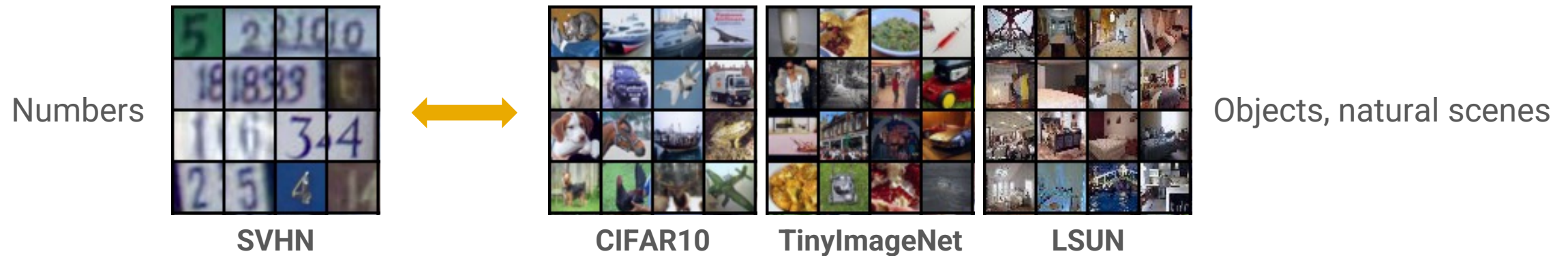


Out-of-Distribution Detection



Probing the Purview of Neural Networks via Gradient Analysis

Dataset Distribution		Detection Accuracy	AUROC	AUPR
In	Out	Baseline [5] / ODIN [6] / Mahalanobis (V) [7] / Mahalanobis (P+FE) [7] / Ours		
CIFAR-10	SVHN	83.36 / 88.81 / 79.39 / 91.95 / 98.04	88.30 / 94.93 / 85.03 / 97.10 / 99.84	88.26 / 95.45 / 86.15 / 96.12 / 99.98
	TinyImageNet	84.01 / 85.21 / 83.60 / 97.45 / 86.17	90.06 / 91.86 / 88.93 / 99.68 / 93.18	89.26 / 91.60 / 88.59 / 99.60 / 92.66
	LSUN	87.34 / 88.42 / 85.02 / 98.60 / 98.37	92.79 / 94.48 / 90.11 / 99.86 / 99.86	92.30 / 94.22 / 89.80 / 99.82 / 99.87
SVHN	CIFAR-10	79.98 / 80.12 / 74.10 / 88.84 / 97.90	81.50 / 81.49 / 79.31 / 95.05 / 99.79	81.01 / 80.95 / 80.83 / 90.25 / 98.11
	TinyImageNet	81.70 / 81.92 / 79.35 / 96.17 / 97.74	83.69 / 83.82 / 83.85 / 99.23 / 99.77	82.54 / 82.60 / 85.50 / 98.17 / 97.93
	LSUN	80.96 / 81.15 / 79.52 / 97.50 / 99.04	82.85 / 82.98 / 83.02 / 99.54 / 99.93	81.97 / 82.01 / 84.67 / 98.84 / 99.21



Out-of-Distribution Detection



Probing the Purview of Neural Networks via Gradient Analysis

Dataset Distribution		Detection Accuracy	AUROC	AUPR
In	Out	Baseline [5] / ODIN [6] / Mahalanobis (V) [7] / Mahalanobis (P+FE) [7] / Ours		
CIFAR-10	SVHN	83.36 / 88.81 / 79.39 / 91.95 / 98.04	88.30 / 94.93 / 85.03 / 97.10 / 99.84	88.26 / 95.45 / 86.15 / 96.12 / 99.98
	TinyImageNet	84.01 / 85.21 / 83.60 / 97.45 / 86.17	90.06 / 91.86 / 88.93 / 99.68 / 93.18	89.26 / 91.60 / 88.59 / 99.60 / 92.66
	LSUN	87.34 / 88.42 / 85.02 / 98.60 / 98.37	92.79 / 94.48 / 90.11 / 99.86 / 99.86	92.30 / 94.22 / 89.80 / 99.82 / 99.87
SVHN	CIFAR-10	79.98 / 80.12 / 74.10 / 88.84 / 97.90	81.50 / 81.49 / 79.31 / 95.05 / 99.79	81.01 / 80.95 / 80.83 / 90.25 / 98.11
	TinyImageNet	81.70 / 81.92 / 79.35 / 96.17 / 97.74	83.69 / 83.82 / 83.85 / 99.23 / 99.77	82.54 / 82.60 / 85.50 / 98.17 / 97.93
	LSUN	80.96 / 81.15 / 79.52 / 97.50 / 99.04	82.85 / 82.98 / 83.02 / 99.54 / 99.93	81.97 / 82.01 / 84.67 / 98.84 / 99.21

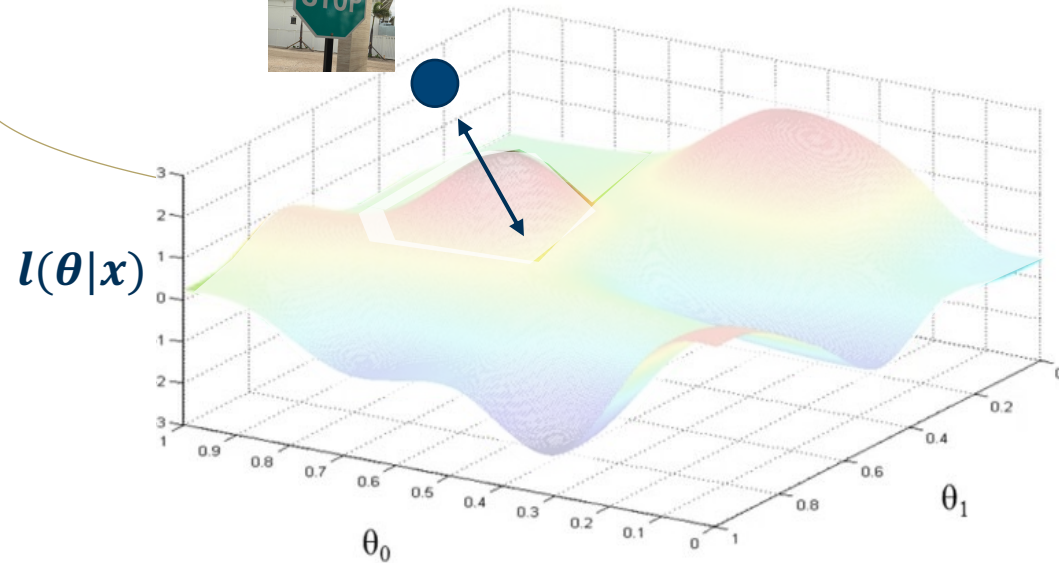


Case Study: Introspective Learning

Gradients as Single pass Features

Our Goal: Use gradients to characterize the novel data at Inference, without global information

Distance from unknown cluster



Two techniques:

1. Gradient constraints during Training for Anomaly Detection
2. **Backpropagating Confounding labels for Out-of-Distribution Detection**



Introspective Learning: A Two-Stage Approach for Inference in Neural Networks



Mohit Prabhushankar, PhD
Postdoc



Ghassan AlRegib, PhD
Professor



Robustness in Neural Networks

Why Robustness?



Introspective Learning: A Two-stage Approach for Inference in Neural Networks

How would humans resolve this challenge?

We Introspect!

- Why am I being shown this slide?
- Why images of muffins rather than pastries?
- What if the dog was a bullmastiff?



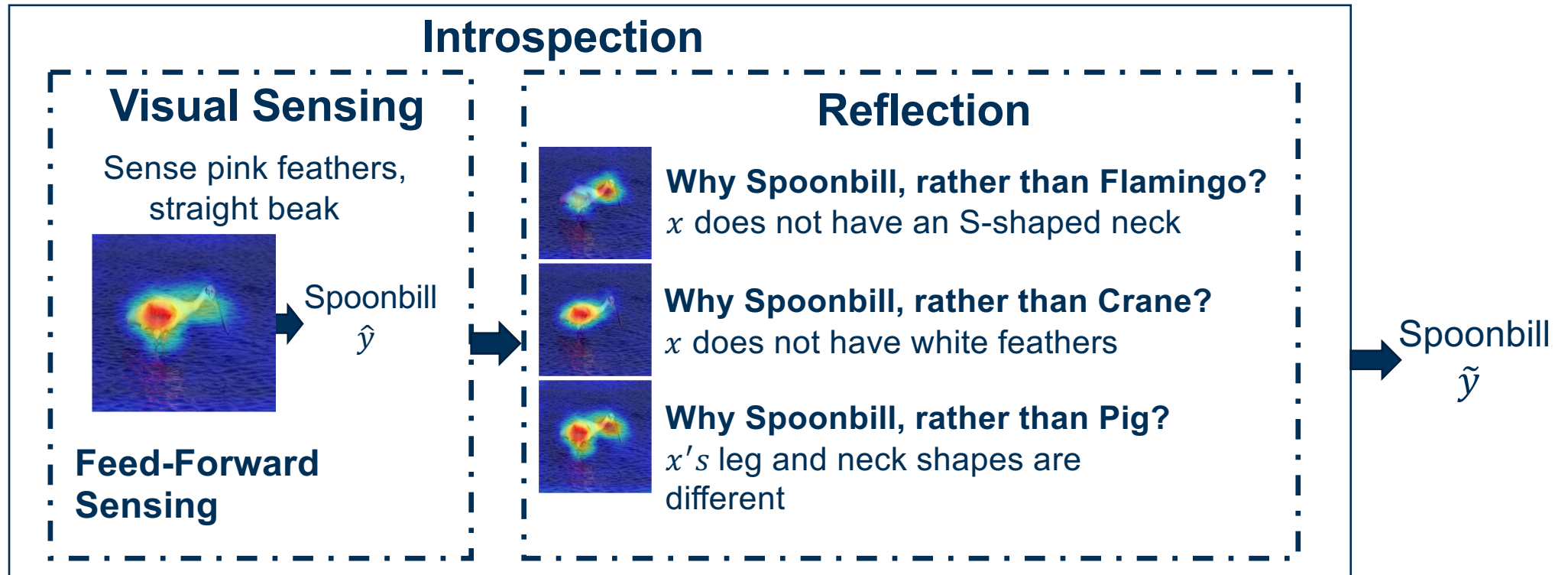
Introspection

What is Introspection?



Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Introspection Learning is a two-stage approach for Inference that combines visual sensing and reflection



Introspection

Introspection in Neural Networks



Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Introspection Learning is a two-stage approach for Inference that combines visual sensing and reflection

Goal : To simulate Introspection in Neural Networks

Definition : We define introspections as answers to logical and targeted questions.

What are the possible targeted questions?



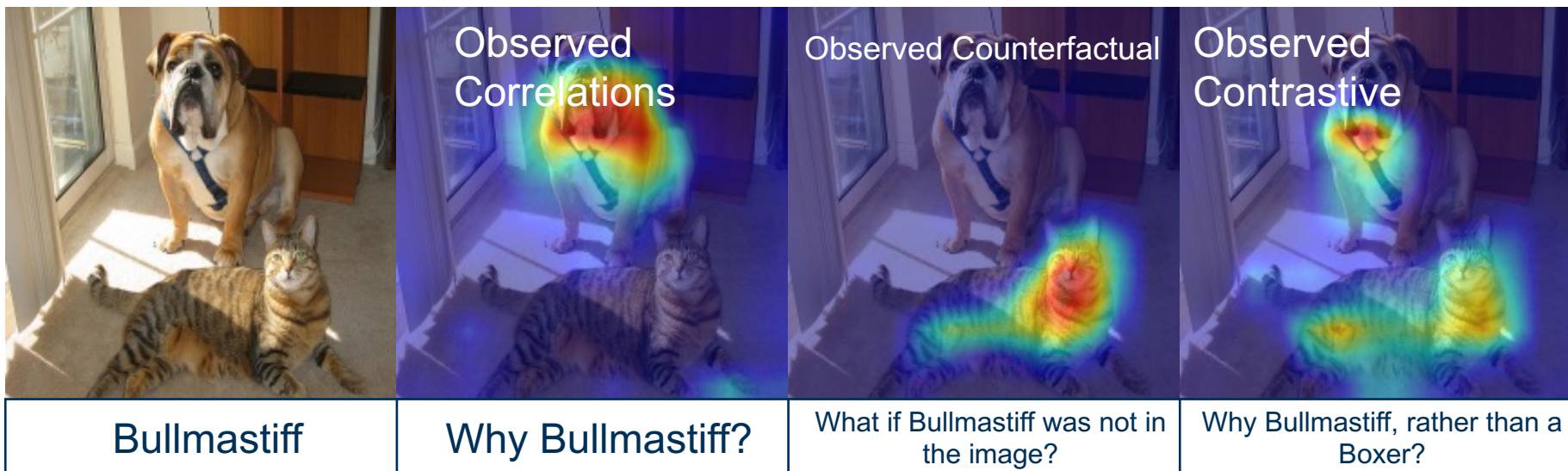
Introspection

Introspection in Neural Networks



Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Introspection Learning is a two-stage approach for Inference that combines visual sensing and reflection



What are the possible targeted questions?



Introspection

Introspection in Neural Networks



Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Introspection Learning is a two-stage approach for Inference that combines visual sensing and reflection

Goal : To simulate Introspection in Neural Networks

***Contrastive Definition :** Introspection answers questions of the form 'Why P , rather than Q ?' where P is a network prediction and Q is the introspective class.*

***Technical Definition :** Given a network $f(x)$, a datum x , and the network's prediction $f(x) = \hat{y}$, introspection in $f(\cdot)$ is the measurement of change induced in the network parameters when a label Q is introduced as the label for x .*



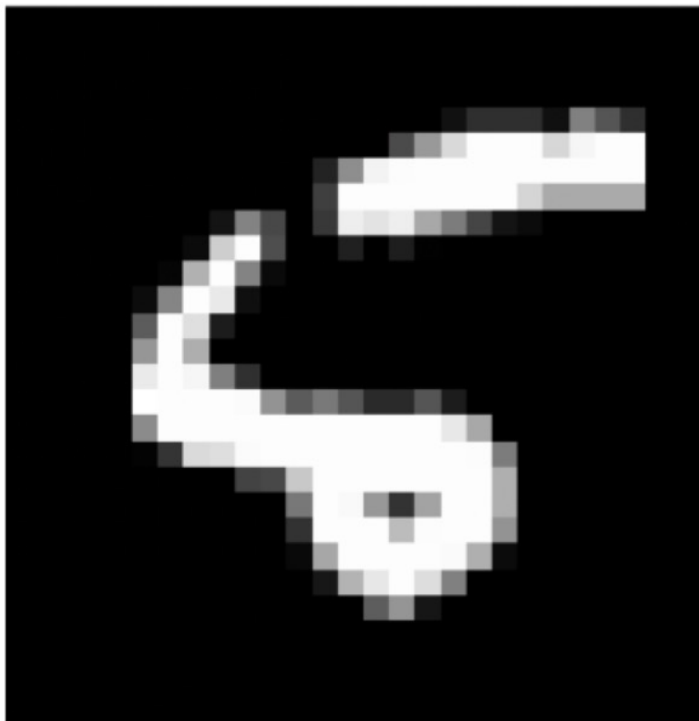
Introspection

Gradients as Features



Introspective Learning: A Two-stage Approach for Inference in Neural Networks

For a well-trained network, the gradients are sparse and informative



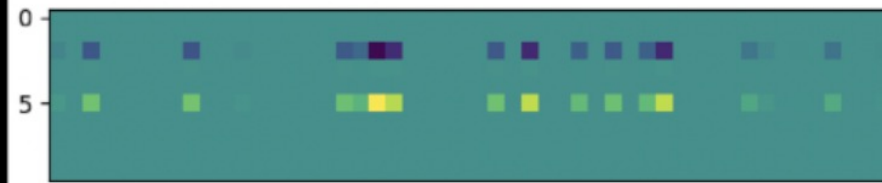
Input Image x



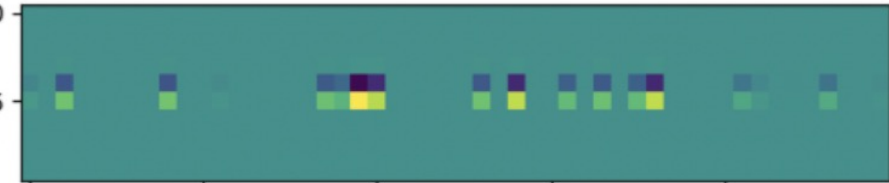
Why 5, rather than 0?



Why 5, rather than 1?



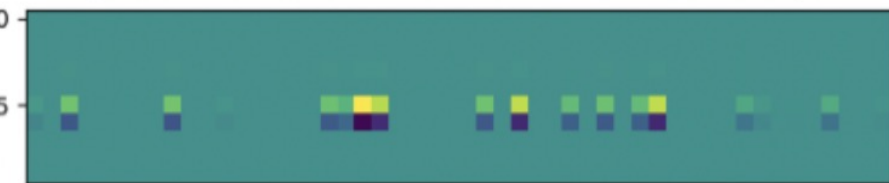
Why 5, rather than 2?



Why 5, rather than 4?



Why 5, rather than 5?



Why 5, rather than 6?



Introspection

Gradients as Features

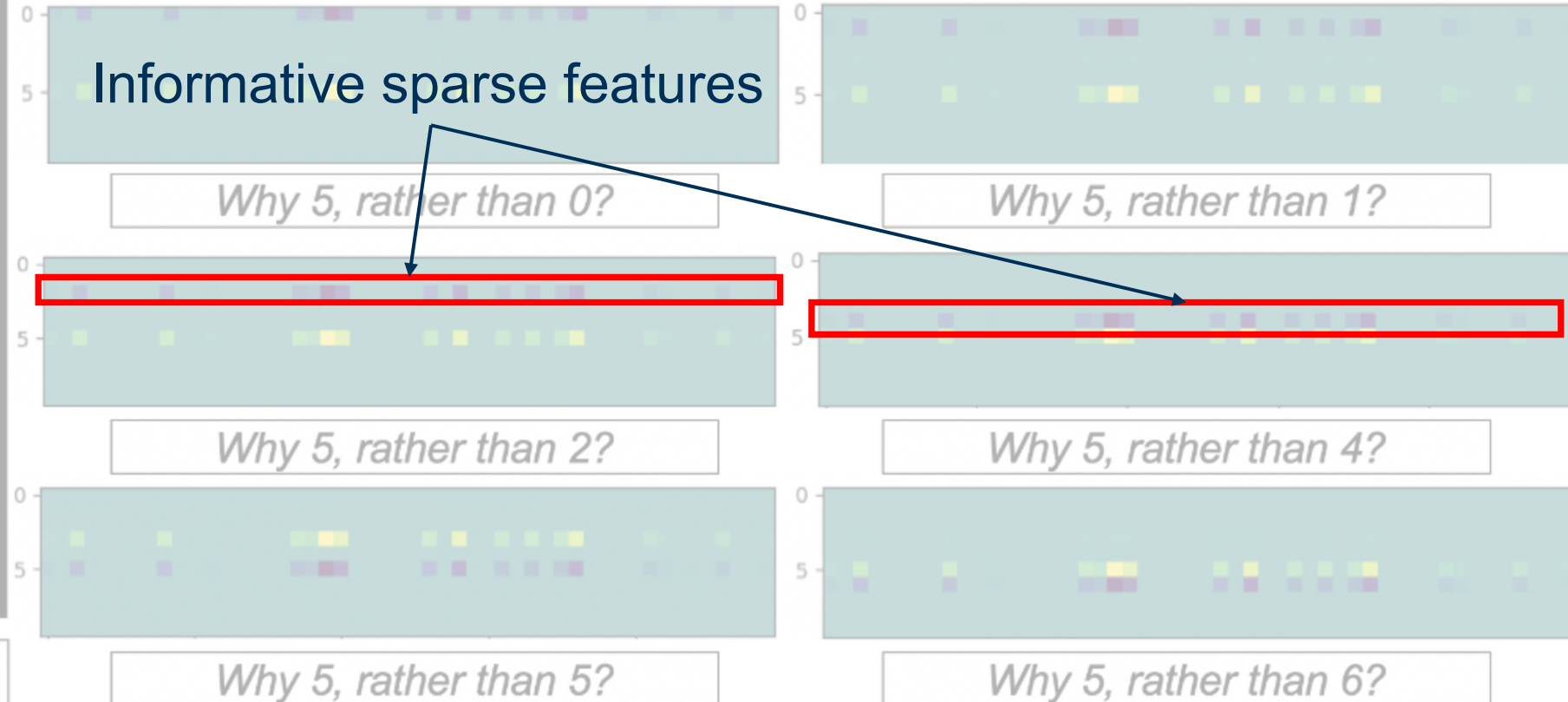


Introspective Learning: A Two-stage Approach for Inference in Neural Networks

For a well-trained network, the gradients are sparse and informative



Input Image x



Introspection

Gradients as Features



Introspective Learning: A Two-stage Approach for Inference in Neural Networks

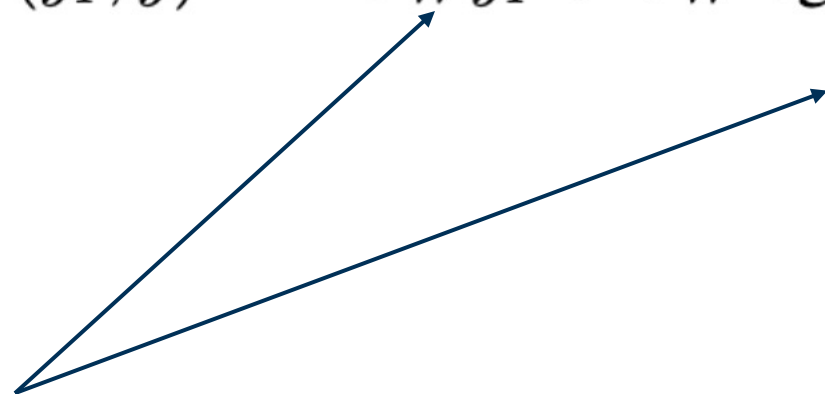
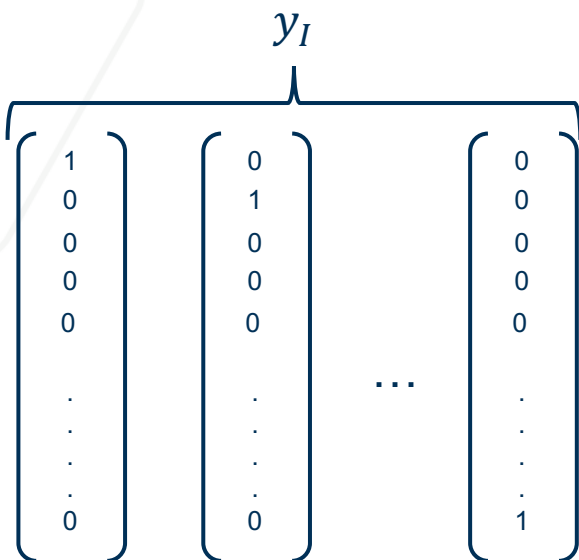
For a well-trained network, the gradients are robust

∇_W = Gradients w.r.t. weights

J = Loss function

\hat{y} = Prediction

$$\text{Lemma 1: } \nabla_W J(y_I, \hat{y}) = -\nabla_W y_I + \nabla_W \log\left(1 + \frac{y\hat{y}}{2}\right).$$



Any change in class requires change in relationship between y_I and \hat{y}



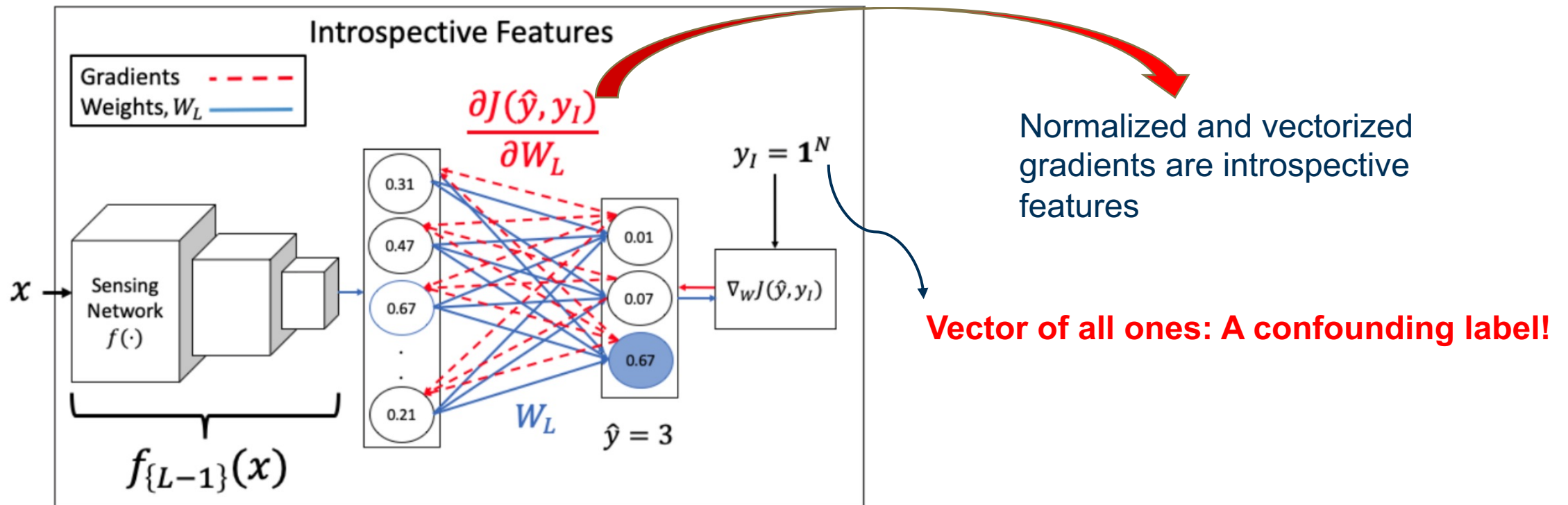
Introspection

Deriving Gradient Features



Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Measure the loss between the prediction \hat{y} and a vector of all ones and backpropagate to obtain the introspective features



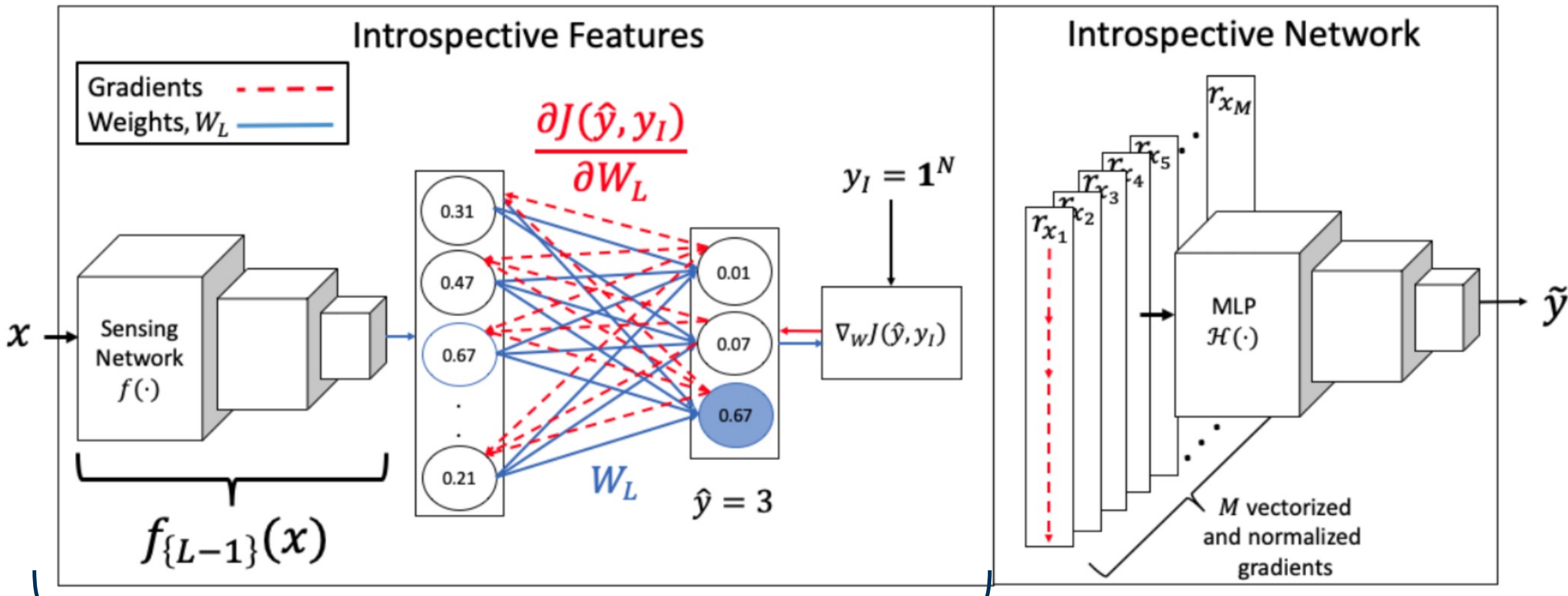
Introspection

Utilizing Gradient Features



SCAN ME

Introspective Learning: A Two-stage Approach for Inference in Neural Networks



Introspective Features



Introspection

When is Introspection Useful?



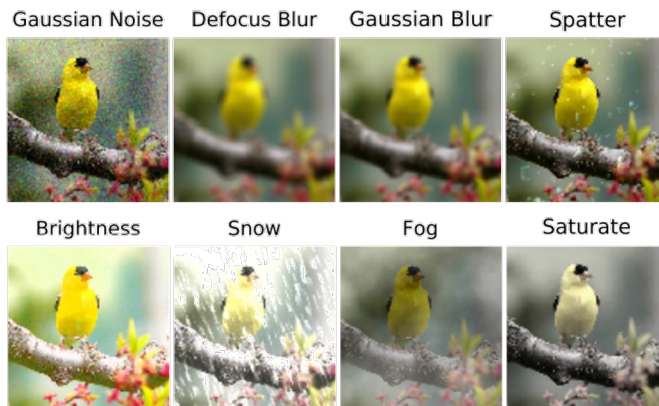
Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Introspection provides robustness when the train and test distributions are different

We define robustness as being generalizable and calibrated to new testing data

Generalizable: Increased accuracy on OOD data

Calibrated: Reduces the difference between prediction accuracy and confidence



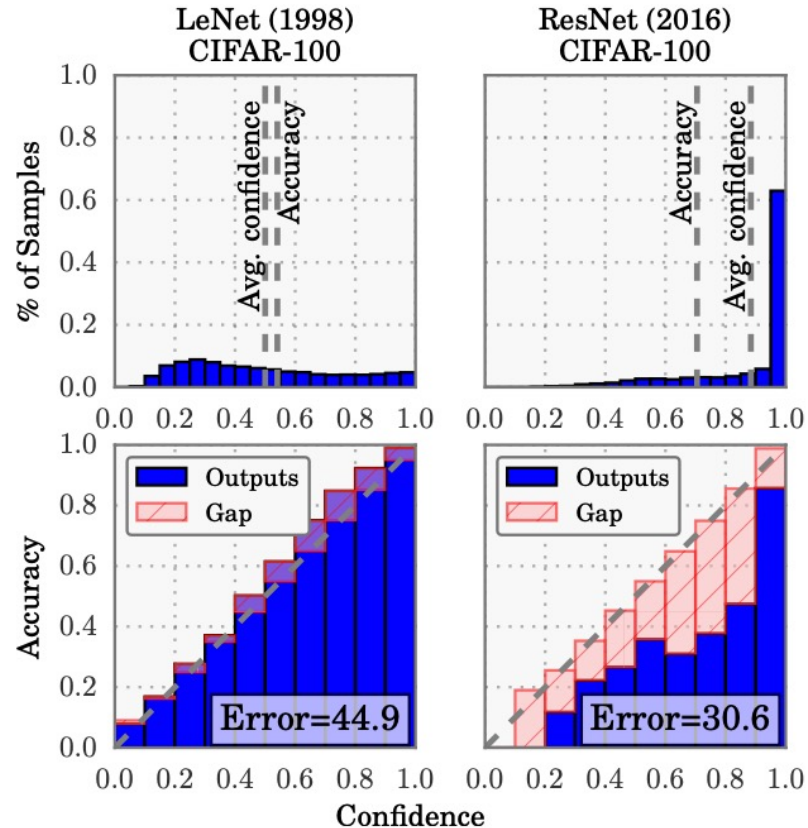
Calibration

A note on Calibration..



Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Calibration occurs when there is mismatch between a network's confidence and its accuracy



- Larger the model, more misplaced is a network's confidence
- On ResNet, the gap between prediction accuracy and its corresponding confidence is significantly high



Introspection in Neural Networks

Generalization and Calibration results

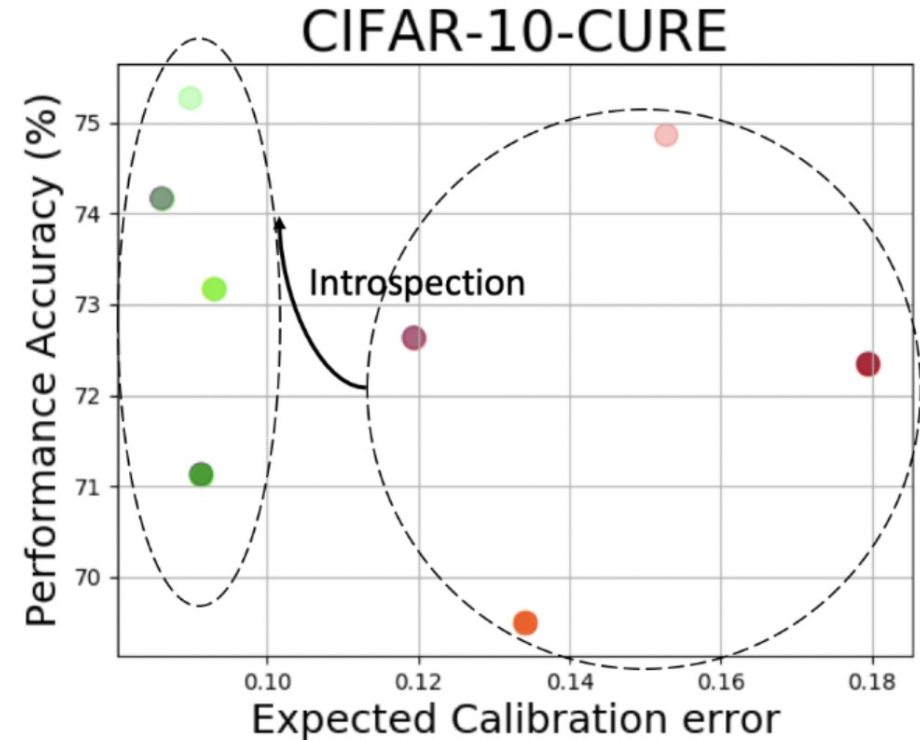
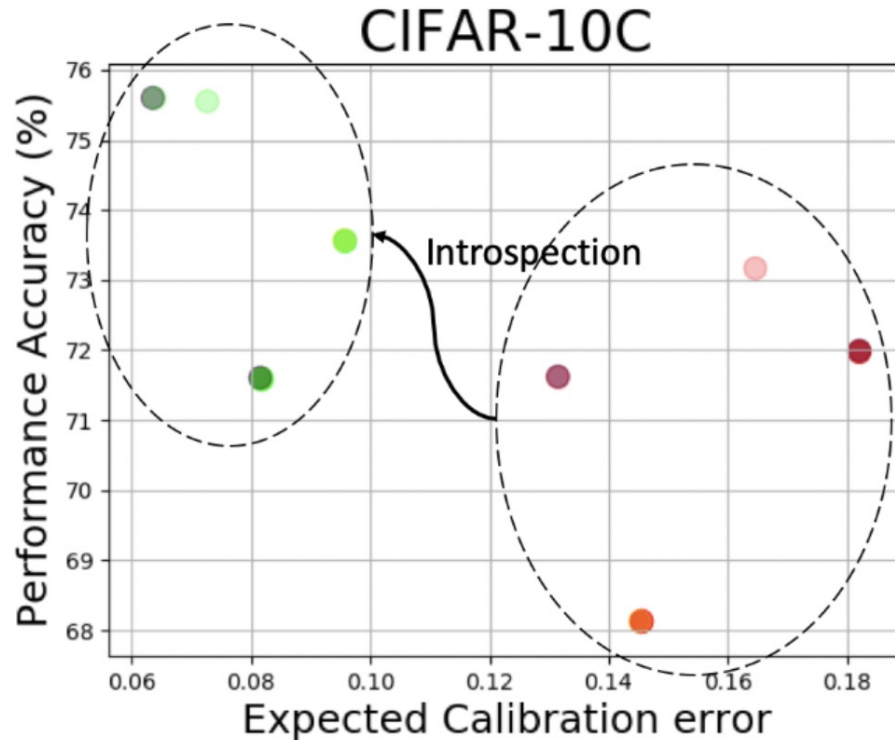


Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Ideal: Top-left corner

Y-Axis: Generalization

X-Axis: Calibration



Legend

Feed-Forward Networks	● ResNet-18	● ResNet-34	● ResNet-50	● ResNet-101
After Introspection	● ResNet-18	● ResNet-34	● ResNet-50	● ResNet-101



Introspection in Neural Networks

Plug-in nature of Introspection



Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Introspection is a light-weight option to resolve robustness issues

Table 1: Introspecting on top of existing robustness techniques.

METHODS		ACCURACY
RESNET-18	FEED-FORWARD	67.89%
	INTROSPECTIVE	71.4%
DENOISING	FEED-FORWARD	65.02%
	INTROSPECTIVE	68.86%
ADVERSARIAL TRAIN (27)	FEED-FORWARD	68.02%
	INTROSPECTIVE	70.86%
SIMCLR (19)	FEED-FORWARD	70.28%
	INTROSPECTIVE	73.32%
AUGMENT NOISE (23)	FEED-FORWARD	76.86%
	INTROSPECTIVE	77.98%
AUGMIX (24)	FEED-FORWARD	89.85%
	INTROSPECTIVE	89.89%

Introspection is a **plug-in approach** that works on all networks and on any downstream task!



Introspection in Neural Networks

Plug-in nature of Introspection



Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Plug-in nature of Introspection benefits downstream tasks like OOD detection, Active Learning, and Image Quality Assessment!

Table 13: Performance of Contrastive Features against Feed-Forward Features and other Image Quality Estimators. Top 2 results in each row are highlighted.

Database	PSNR	IW	SR	FSIMc	Per	CSV	SUM	Feed-Forward	Introspective
	HA	SSIM	SIM		SIM		MER	UNIQUE	UNIQUE
Outlier Ratio (OR, ↓)									
MULTI	0.013	0.013	0.000	0.016	0.004	0.000	0.000	0.000	0.000
TID13	0.615	0.701	0.632	0.728	0.655	0.687	0.620	0.640	0.620
Root Mean Square Error (RMSE, ↓)									
MULTI	11.320	10.049	8.686	10.794	9.898	9.895	8.212	9.258	7.943
TID13	0.652	0.688	0.619	0.687	0.643	0.647	0.630	0.615	0.596
Pearson Linear Correlation Coefficient (PLCC, ↑)									
MULTI	0.801	0.847	0.888	0.821	0.852	0.852	0.901	0.872	0.908
	-1	-1	0	-1	-1	-1	-1	-1	
TID13	0.851	0.832	0.866	0.832	0.855	0.853	0.861	0.869	0.877
	-1	-1	0	-1	-1	-1	0	0	
Spearman's Rank Correlation Coefficient (SRCC, ↑)									
MULTI	0.715	0.884	0.867	0.867	0.818	0.849	0.884	0.867	0.887
	-1	0	0	0	-1	-1	0	0	
TID13	0.847	0.778	0.807	0.851	0.854	0.846	0.856	0.860	0.865
	-1	-1	-1	-1	0	-1	0	0	
Kendall's Rank Correlation Coefficient (KRCC)									
MULTI	0.532	0.702	0.678	0.677	0.624	0.655	0.698	0.679	0.702
	-1	0	0	0	-1	0	0	0	
TID13	0.666	0.598	0.641	0.667	0.678	0.654	0.667	0.667	0.677
	0	-1	-1	0	0	0	0	0	

Table 2: Recognition accuracy of Active Learning strategies.

Methods	Architecture	Original Testset		Gaussian Noise	
		R-18	R-34	R-18	R-34
Entropy (E1)	Feed-Forward	0.365	0.358	0.244	0.249
	Introspective	0.365	0.359	0.258	0.255
Least (E1)	Feed-Forward	0.371	0.359	0.252	0.25
	Introspective	0.373	0.362	0.264	0.26
Margin (E2)	Feed-Forward	0.38	0.369	0.251	0.253
	Introspective	0.381	0.373	0.265	0.263
BALD (E3)	Feed-Forward	0.393	0.368	0.26	0.253
	Introspective	0.396	0.375	0.273	0.263
BADGE (E3)	Feed-Forward	0.388	0.37	0.25	0.247
	Introspective	0.39	0.37	0.265	0.260

Table 3: Out-of-distribution Detection of existing techniques compared between feed-forward and introspective networks.

Methods	OOD Datasets	FPR (95% at TPR)	Detection Error	AUROC
		↓	↓	↑
Feed-Forward/Introspective				
MSP (E3)	Textures	58.74/19.66	18.04/7.49	88.56/97.79
	SVHN	61.41/51.27	16.92/15.67	89.39/91.2
	Places365	58.04/54.43	17.01/15.07	89.39/91.3
	LSUN-C	27.95/27.5	9.42/10.29	96.07/95.73
ODIN (E6)	Textures	52.3/9.31	22.17/6.12	84.91/91.9
	SVHN	66.81/48.52	23.51/15.86	83.52/91.07
	Places365	42.21/51.87	16.23/15.71	91.06/90.95
	LSUN-C	6.59/23.66	5.54/10.2	98.74/95.87



Robust Neural Networks

Part 4: Intervenability at Inference



Objective

Objective of the Tutorial

To discuss methodologies that promote robustness in neural networks at inference

- Part 1: Inference in Neural Networks
- Part 2: Explainability at Inference
- Part 3: Uncertainty at Inference
- **Part 4: Intervenability at Inference**
 - Definitions of Intervenability
 - Causality
 - Privacy
 - Interpretability
 - Prompting
 - Benchmarking
 - Case Study: Intervenability in Interpretability
- Part 5: Conclusions and Future Directions



Intervenability

Through the Causal Glass

Assess: The amenability of neural network decisions to human interventions



“Interventions in data are manipulations that are designed to test for causal factors”

Intervenability

Through the Privacy Glass

Assure: The amenability of neural network decisions to human interventions



*“Intervenability aims at the possibility for parties involved in any **privacy-relevant** data processing to **interfere** with the ongoing or planned data processing”*

Intervenability

Through the Interpretability Glass

Interpret: The amenability of neural network decisions to human interventions



*“The post-hoc field of explainability, that previously only justified decisions, becomes **active** by being involved in the decision making process and providing limited, but relevant and contextual interventions”*

Intervenability

Through the Benchmarking Glass

Verify: The amenability of neural network decisions to human interventions



*“... new **benchmarks** were proposed to specifically test generalization of classification and detection methods with respect to **simple** algorithmically generated interventions like spatial shifts, blur, changes in brightness or contrast...”*

Case Study: Intervenability in Interpretability

Challenges in Intervenability

The amenability of neural network decisions to human interventions



- **Assess: Causality**
- **Assure: Privacy**
- **Interpret: Interpretability**
- **Verify: Benchmarking**

Challenges:

- **Choosing the type of Intervention: Explanation Evaluation**
- **Residuals of Interventions: Uncertainty**

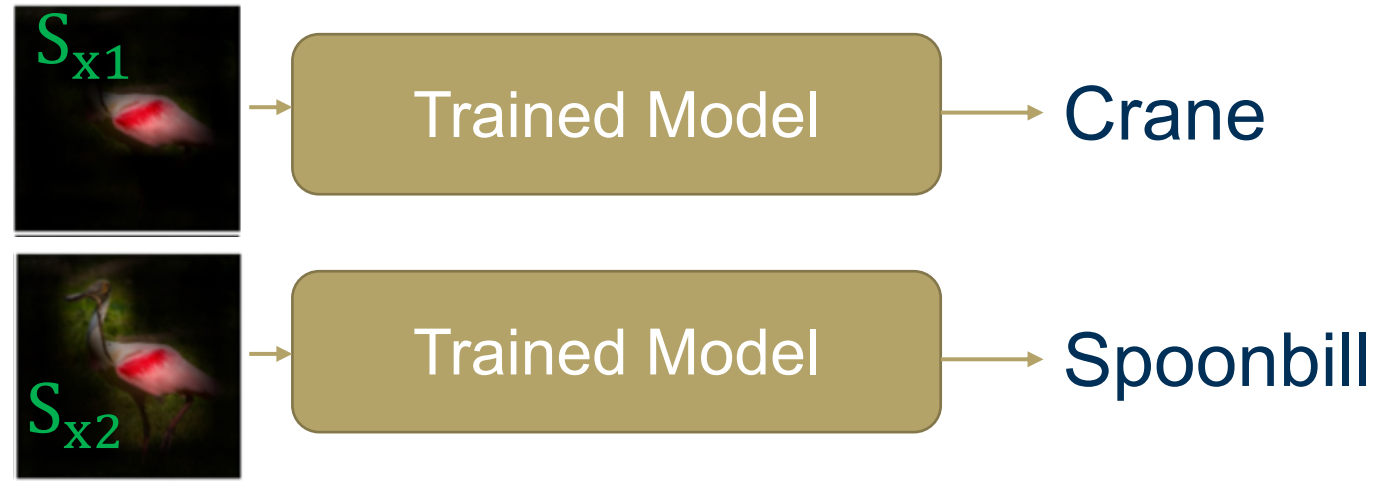
Case Study: Intervenability in Interpretability

Explanation Evaluation

Visual explanations are evaluated via masking the important regions in the image and passing it through the network

Three types of Masking:

1. **Masking using explanation heatmap**
2. Pixel-wise masking using explanation as importance
3. Structure-wise masking using information encoded in explanation



Masking = Intelligent Intervention

Case Study: Intervenability in Interpretability

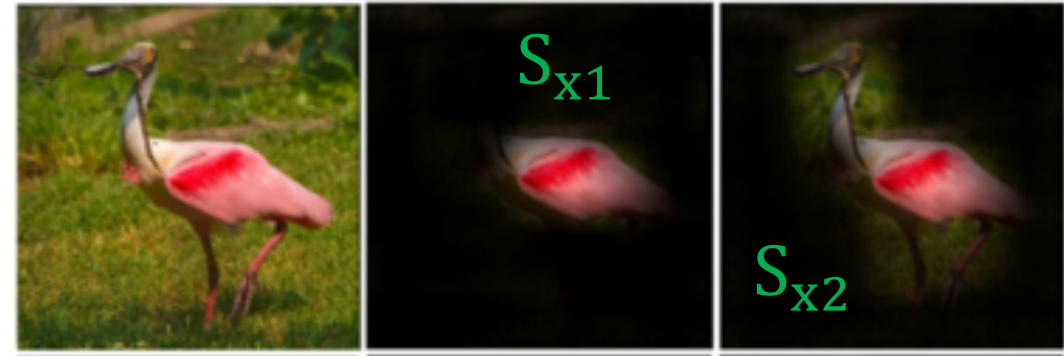
Evaluation 1: Explanation Evaluation via Masking

Common evaluation technique is masking the image and checking for prediction correctness

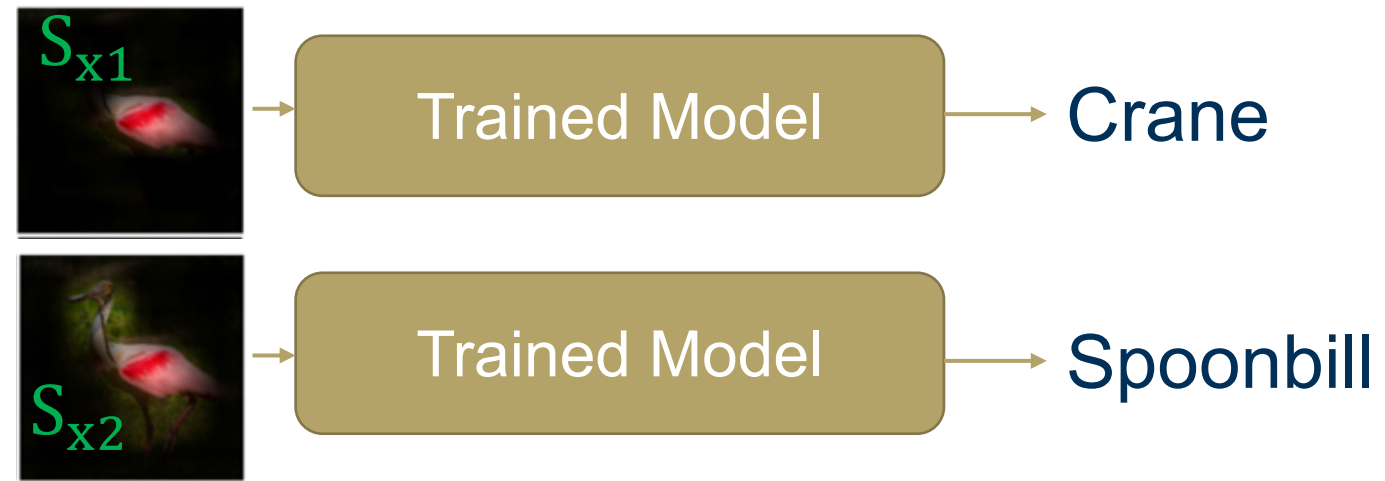
y = Prediction

S_x = Explanation masked data

$E(Y|S_x)$ = Expectation of class given S_x



If across N images,
 $E(Y|S_{x2}) > E(Y|S_{x1})$,
explanation technique 2
is better than explanation
technique 1

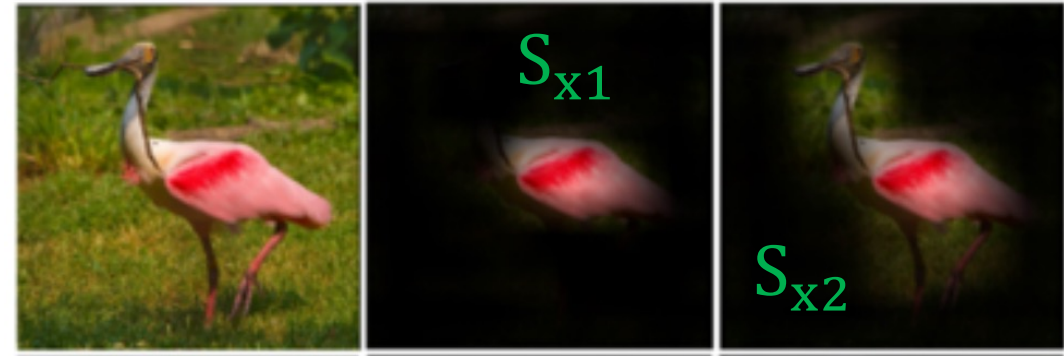


Case Study: Intervenability in Interpretability

Evaluation 1: Explanation Evaluation via Masking

However, explanation masking encourages 'larger' explanations

- Larger explanations imply more features in masked images are intact (unmasked)
- This increases likelihood of a correct prediction
- 'Fine-grained' explanations are not promoted



Case Study: Intervenability in Interpretability

Explanation Evaluation

Common evaluation technique is masking the image and checking for prediction correctness

Three types of Masking:

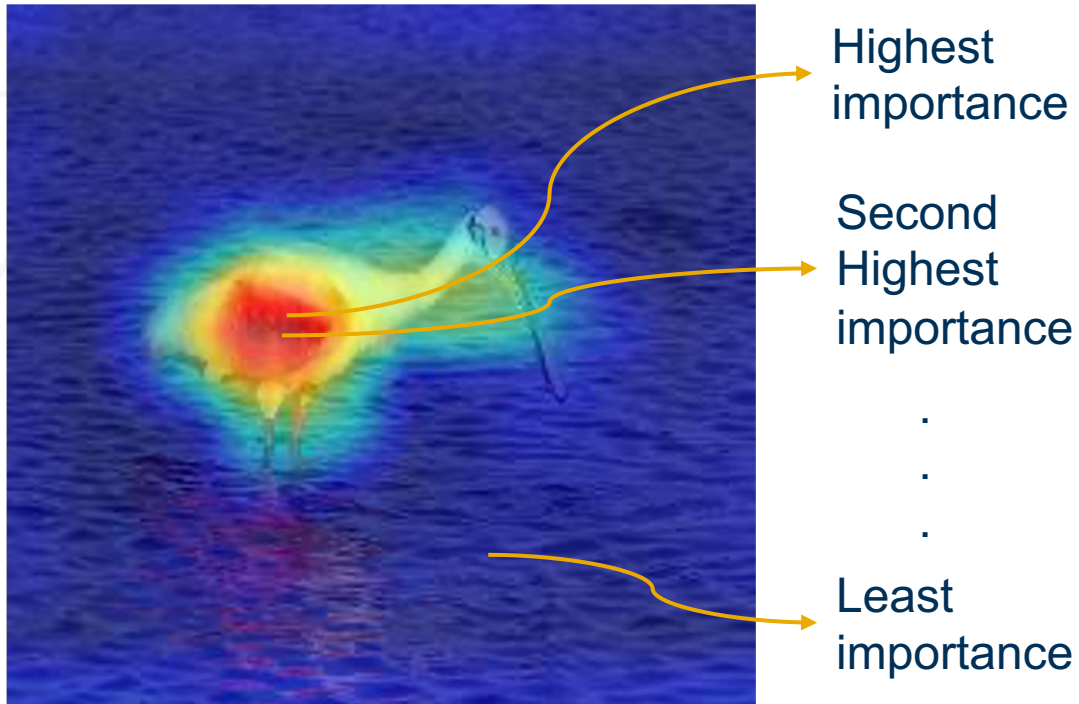
1. Masking using explanation heatmap
2. **Pixel-wise masking using explanation as importance**
3. Structure-wise masking using information encoded in explanation



Case Study: Intervenability in Interpretability

Evaluation 2: Progressive Pixel-wise Insertion and Deletion

Pixel-wise Deletion: Sequentially delete (mask) pixels in an image based on their explanation assigned importance scores



Step 1: Mask highest importance pixel and pass the image through the network. Note the probability of spoonbill.

Step 2: Mask the second highest importance pixel from the image in Step 1 and pass the image through the network. Note the probability of spoonbill.

Step 3: Repeat until all pixels are deleted (masked)

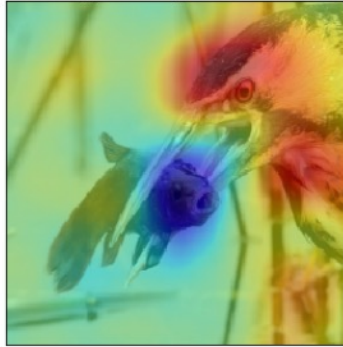
Case Study: Intervenability in Interpretability

Evaluation 2: Progressive Pixel-wise Insertion and Deletion

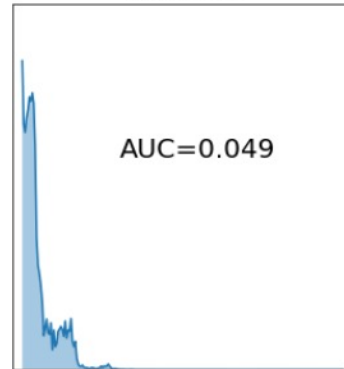
The removal of the "cause" (important pixels) will force the base model to change its decision.



Explaining: bittern



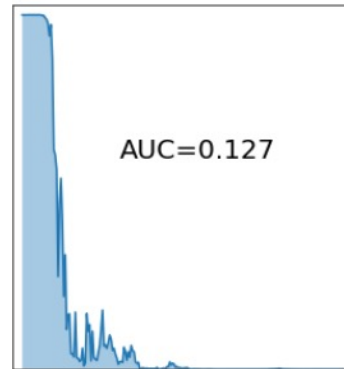
Deletion



Explaining: white stork



Deletion

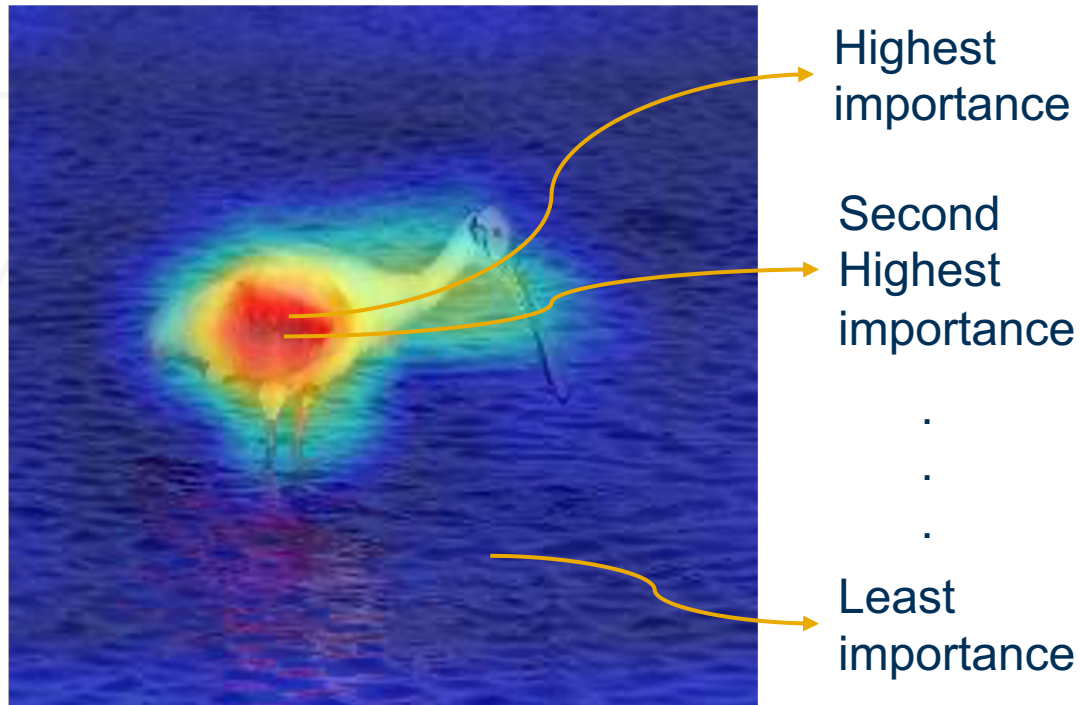


- **Deletion approximates Necessity** criterion of a "good" explanation
- **AUC** for a good explanation will be **low**
- **Deletion** encourages **fine-grained explanations** by choosing those heatmaps that select the most relevant pixels

Case Study: Intervenability in Interpretability

Evaluation 2: Progressive Pixel-wise Insertion and Deletion

Pixel-wise Insertion: Sequentially add pixels to a mean image based on their explanation assigned importance scores



Take a mean (grayscale) image

Step 1: Add the highest importance pixel to the mean image and pass it through the network. Note the probability of spoonbill.

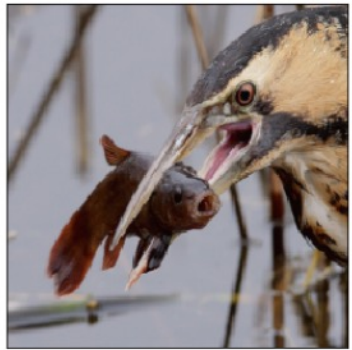
Step 2: Add the second highest importance pixel to the image in Step 1 and pass the image through the network. Note the probability of spoonbill.

Step 3: Repeat until all pixels are inserted

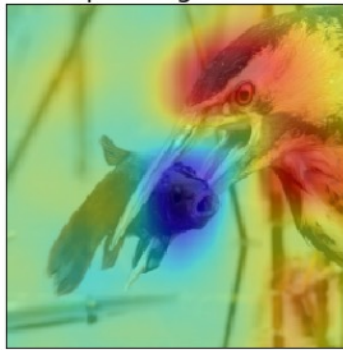
Case Study: Intervenability in Interpretability

Evaluation 2: Progressive Pixel-wise Insertion and Deletion

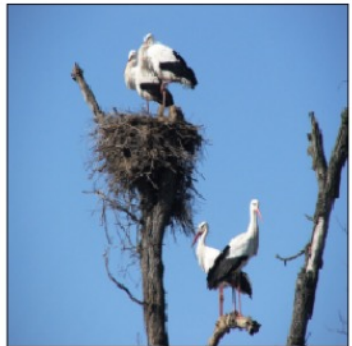
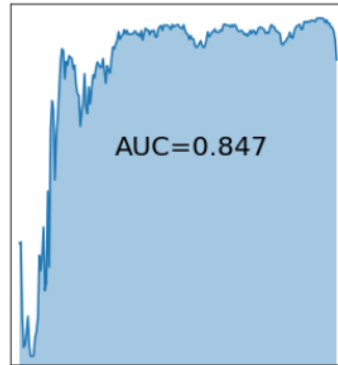
The addition of the "cause" (important pixels) will force the base model to change its decision.



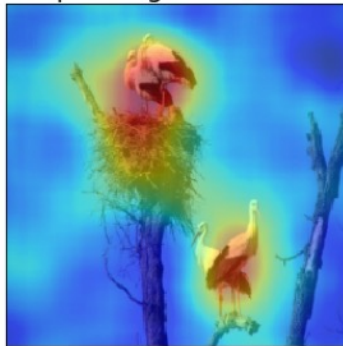
Explaining: bittern



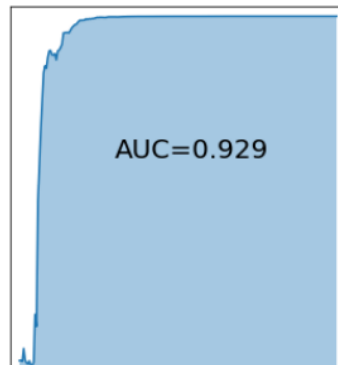
Insertion



Explaining: white stork



Insertion

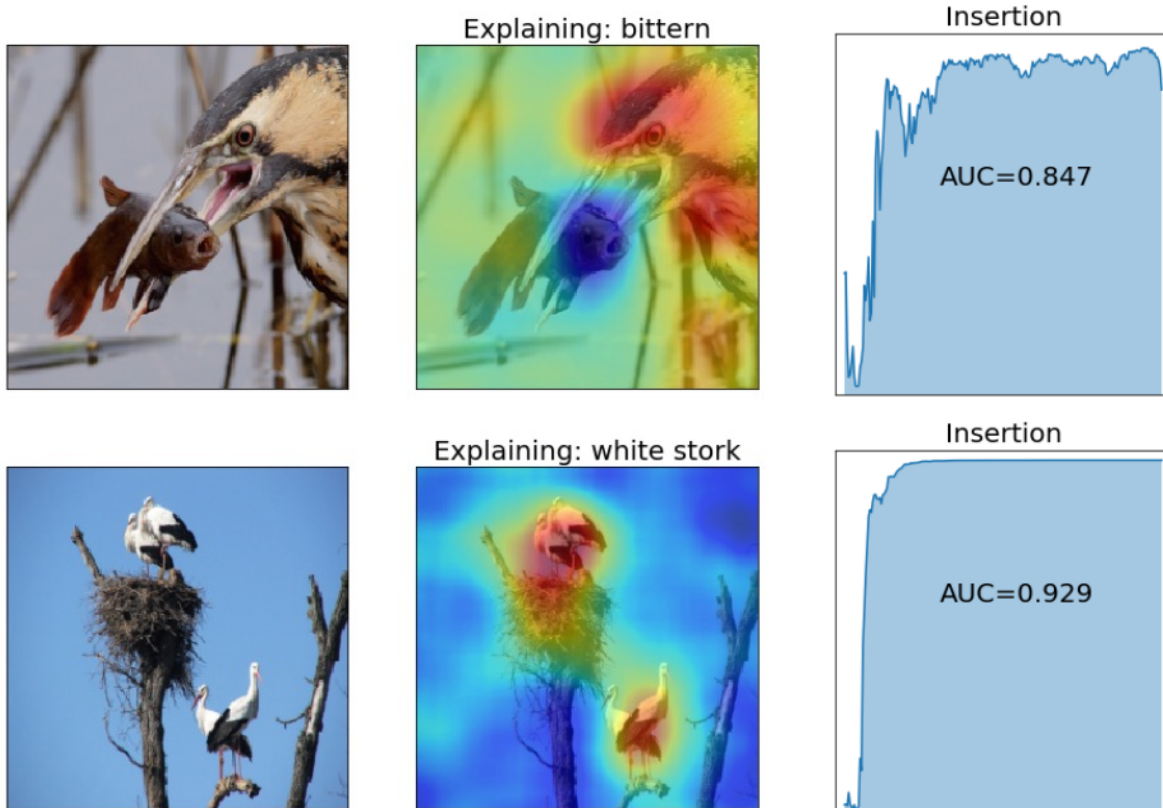


- **Insertion approximates Sufficiency** criterion of a "good" explanation
- **AUC** for a good explanation will be **high**
- **Insertion** encourages **fine-grained explanations** by choosing those heatmaps that select the most relevant pixels

Case Study: Intervenability in Interpretability

Evaluation 2: Progressive Pixel-wise Insertion and Deletion

Insertion and Deletion evaluation metrics encourage pixel-wise analysis of explanations



- However, humans do not “see” in pixels
- Rather they view scenes in a “structure-wise” fashion
- While heatmap masking encourages large explanations, pixel-wise masking encourages unrealistic and non-human like explanations

Case Study: Intervenability in Interpretability

Explanation Evaluation

Common evaluation technique is masking the image and checking for prediction correctness

Three types of Masking:

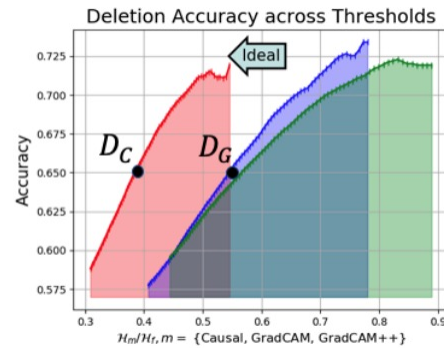
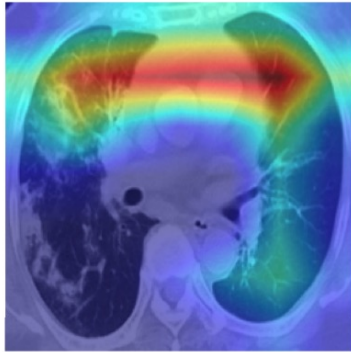
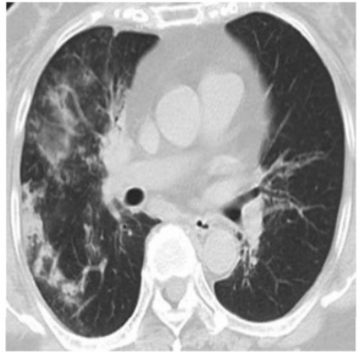
1. Masking using explanation heatmap
2. Pixel-wise masking using explanation as importance
3. **Structure-wise masking using information encoded in explanation**



Case Study: Intervenability in Interpretability

Evaluation 3: Progressive Structure-wise Insertion and Deletion

Structure-wise Deletion: Sequentially delete (mask) pixels in an image based on the number of bits used to represent the region



Ideal scenario: The explanation encodes the most important information in the least possible bits

CausalCAM in Red¹
GradCAM in Purple
GradCAM++ in Green

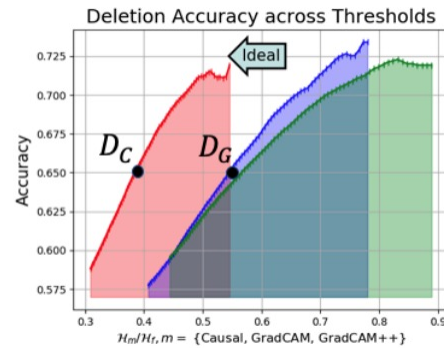
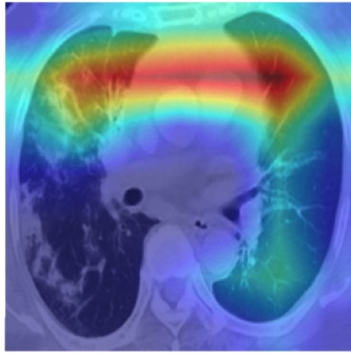
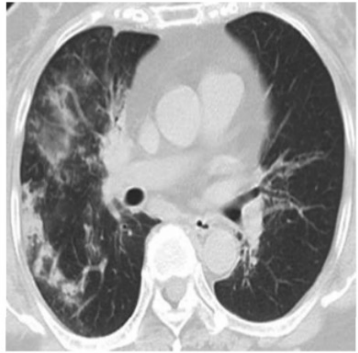
- D_C and D_G represent 65% accuracy for CausalCAM and GradCAM respectively
- **CausalCAM encodes dense structure-rich features in lesser bits, that aid accuracy**



Case Study: Intervenability in Interpretability

Evaluation 3: Progressive Structure-wise Insertion and Deletion

Structure-wise Deletion: Sequentially delete (mask) pixels in an image based on the number of bits used to represent the region



Ideal scenario: The explanation encodes the most important information in the least possible bits

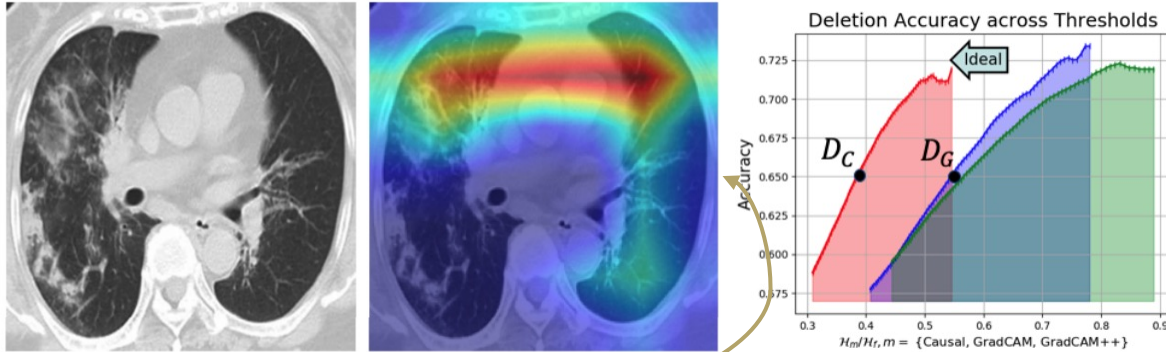
Step 1: Choose a threshold in the explanation (say 0.1) and delete (mask) all the pixels in the original image below the threshold. Pass the masked image through the network and note the change in prediction (if any)



Case Study: Intervenability in Interpretability

Evaluation 3: Progressive Structure-wise Insertion and Deletion

Structure-wise Deletion: Sequentially delete (mask) pixels in an image based on the number of bits used to represent the region



Y-axis: Performance accuracy across all ratios

X-axis: Ratio of Huffman encoded masked and original images for all explanations. Smaller the ratio, less is the number of bits encoding the masked image

Ideal scenario: The explanation encodes the most important information in the least possible bits

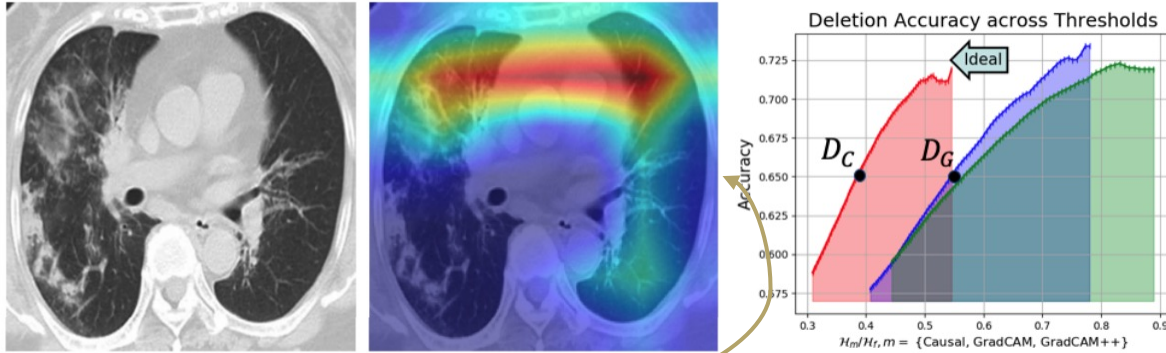
Step 1: Choose a threshold in the explanation (say 0.1) and delete (mask) all the pixels in the original image below the threshold. Pass the masked image through the network and note the change in prediction (if any)

Step 2: Calculate the Huffman code for the original and the masked image. The ratio between the codes of masked and original image is taken on the x-axis and the corresponding accuracy across all images is shown on the y-axis

Case Study: Intervenability in Interpretability

Evaluation 3: Progressive Structure-wise Insertion and Deletion

Structure-wise Deletion: Sequentially delete (mask) pixels in an image based on the number of bits used to represent the region



Y-axis: Performance accuracy across all ratios

X-axis: Ratio of Huffman encoded masked and original images for all explanations. Smaller the ratio, less is the number of bits encoding the masked image

Ideal scenario: The explanation encodes the most important information in the least possible bits

Step 1: Choose a threshold in the explanation (say 0.1) and delete (mask) all the pixels in the original image below the threshold. Pass the masked image through the network and note the change in prediction (if any)

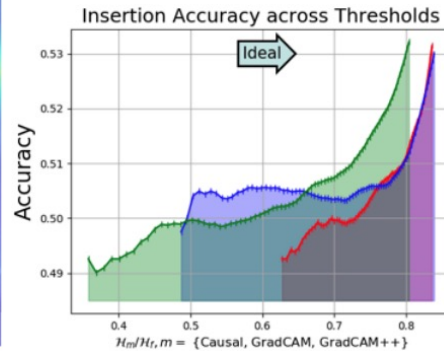
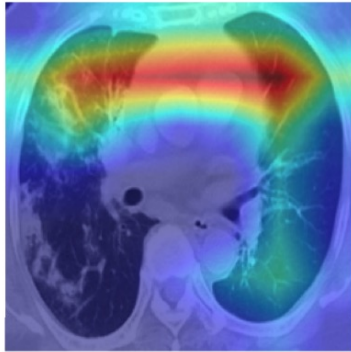
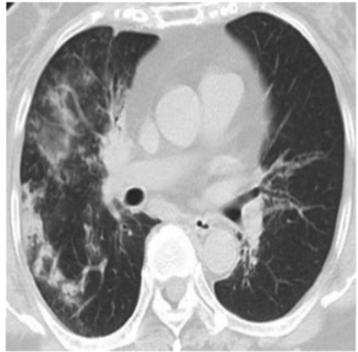
Step 2: Calculate the Huffman code for the original and the masked image. The ratio between the codes of masked and original image is taken on the x-axis and the corresponding accuracy across all images is shown on the y-axis

Step 3: Repeat across thresholds

Case Study: Intervenability in Interpretability

Evaluation 3: Progressive Structure-wise Insertion and Deletion

Structure-wise Insertion: Sequentially add (insert) pixels in an image based on the number of bits used to represent the region



Ideal scenario: The explanation encodes the most important information in the least possible bits

CausalCAM in Red¹
GradCAM in Purple
GradCAM++ in Green

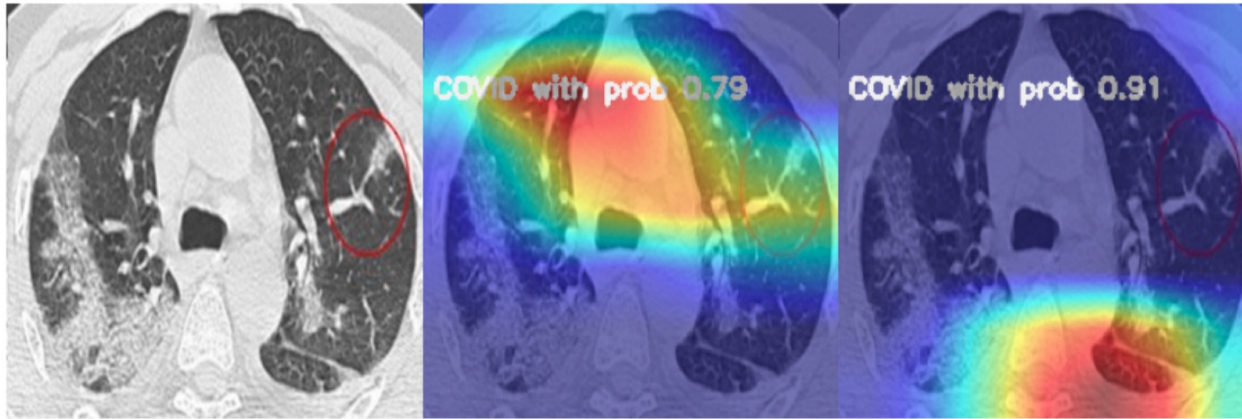
- **CausalCAM encodes dense structure-rich features in at the lowest threshold, that aid accuracy**



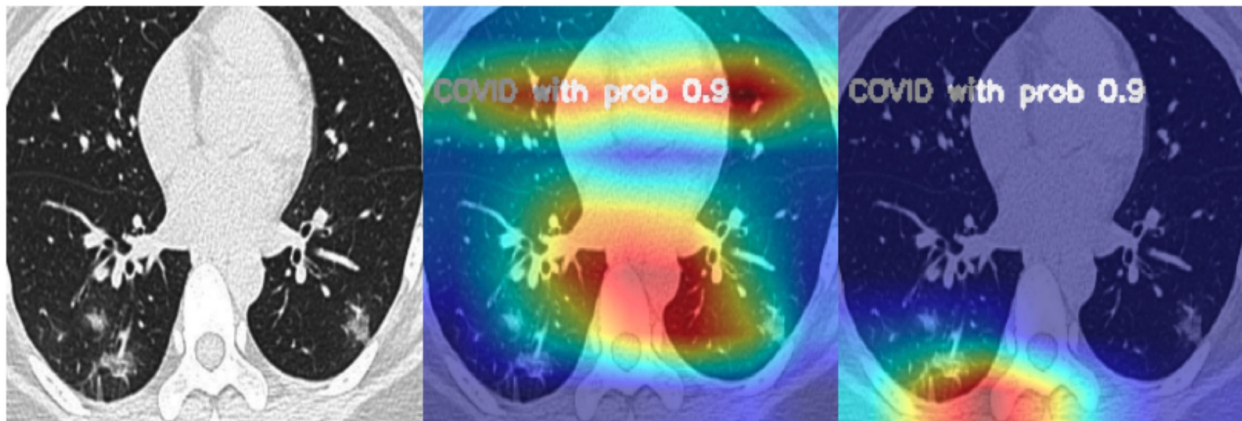
Case Study: Intervenability in Interpretability

Evaluation 3: Progressive Structure-wise Insertion and Deletion

Structure-wise insertion and deletion can sometimes promote adversarial explanations



(a)



- Best explanations according to structure-wise insertion and deletion.
- Corroborated by high probabilities

Case Study: Intervenability in Interpretability

Pros and Cons

Evaluation 1: Explanation heatmap masking

- **Pro:** Structures are visible in the explanations
- **Con:** Encourages large non-fine grained explanations

Evaluation 2: Pixel-wise insertion and deletion

- **Pro:** Progressively assigns importance to pixels
- **Con:** Encourages unrealistic and dispersed explanations

Evaluation 3: Structure-wise insertion and deletion

- **Pro:** Encourages structures while progressively assigning importance to structures based on information bits
- **Pro:** Other human-centric measures including SSIM, saliency etc. can be used on x-axis
- **Con:** Encourages causal (and sometimes adversarial) explanations without considering context information



Case Study: Intervenability in Interpretability

Challenges in Intervenability

The amenability of neural network decisions to human interventions



- Hence, there is **no single-best interventional** strategy
- Choosing the **right** intervention is still an **art**

Challenges:

- **Choosing the type of Intervention: Explanation Evaluation**
- **Residuals of Interventions: Uncertainty**

Case Study: Intervenability in Interpretability

Challenges in Intervenability

The amenability of neural network decisions to human interventions



- Hence, there is **no single-best interventional** strategy
- Choosing the **right** intervention is still an **art**

Challenges:

- Choosing the type of Intervention: Explanation Evaluation
- **Residuals of Interventions: Uncertainty**

VOICE: Variance of Contrastive Explanations for Quantifying Uncertainty in Interpretability



Mohit Prabhushankar, PhD
Postdoc



Ghassan AlRegib, PhD
Professor

Case Study: Intervenability in Interpretability

Predictive Uncertainty in Explanations

Explanatory techniques have predictive uncertainty

Explanation of Prediction

Uncertainty of Explanation

Why Bullmastiff?

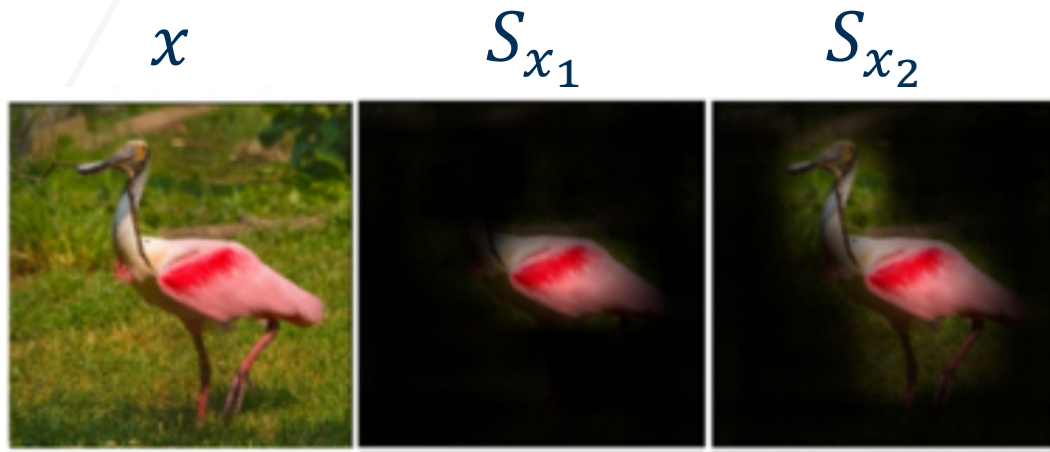


Uncertainty in answering
Why Bullmastiff?

Case Study: Intervenability in Interpretability

Predictive Uncertainty

Uncertainty due to variance in prediction when model is kept constant



$$V[y|S_x] = V[E(y|S_x)] + E(V[y|S_x])$$

y = Prediction

$V[y]$ = Variance of prediction (Predictive Uncertainty)

S_x = Subset of data (Some intervention)

$E(Y|S_x)$ = Expectation of class given a subset

$V(Y|S_x)$ = Variance of class given all other residuals

Case Study: Intervenability in Interpretability

Visual Explanations (partially) reduce Predictive Uncertainty

A 'good' explanatory technique is evaluated to have zero $V[E(y|S_x)]$



$$V[y|S_x] = V[E(y|S_x)] + E(V[y|S_x])$$

y = Prediction

$V[y]$ = Variance of prediction (Predictive Uncertainty)

S_x = Subset of data (Some intervention)

$E(Y|S_x)$ = Expectation of class given a subset

$V(Y|S_x)$ = Variance of class given all other residuals

zero ←

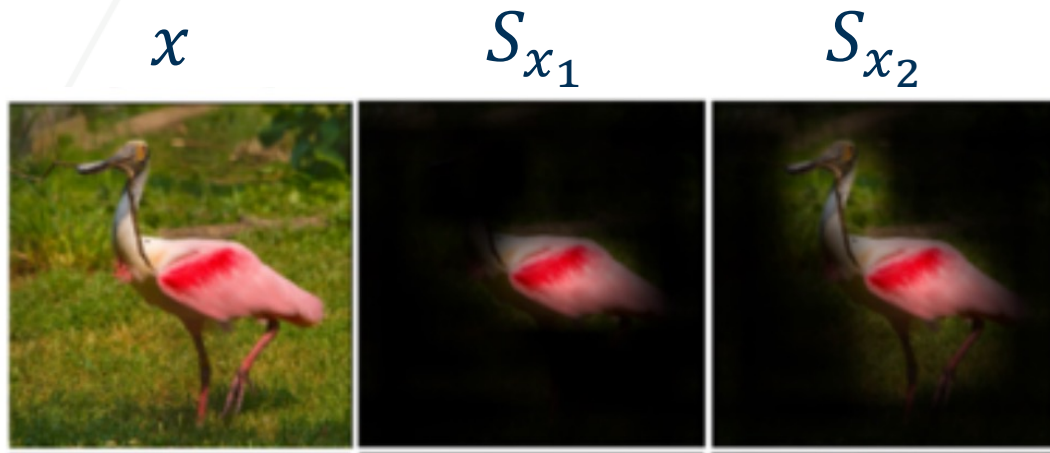
Key Observation 1: Visual Explanations are evaluated to partially reduce the predictive uncertainty in a neural network

Network evaluations have nothing to do with human Explainability!

Case Study: Intervenability in Interpretability

Predictive Uncertainty in Explanations is the Residual

All other subsets 'not' chosen by the explanatory technique contributes to uncertainty



$$V[y|S_x] = V[E(y|S_x)] + E(V[y|S_x])$$

y = Prediction

$V[y]$ = Variance of prediction (Predictive Uncertainty)

S_x = Subset of data (Some intervention)

$E(Y|S_x)$ = Expectation of class given a subset

$V(Y|S_x)$ = Variance of class given all other residuals

Key Observation 2: Uncertainty in Explainability occurs due to all combinations of features that the explanation did not attribute to the network's decision

Case Study: Intervenability in Interpretability

Predictive Uncertainty in Explanations is the Residual

All other subsets 'not' chosen by the explanatory technique contributes to uncertainty



$$V[y|S_x] = V[E(y|S_x)] + E(V[y|S_x])$$

The effect of a chosen Interventions can be measured based on *all the Interventions that were not chosen*

y = Prediction
 $V[y]$ = Variance of prediction (Predictive Uncertainty)
 S_x = Subset of data (Some intervention)
 $E(Y|S_x)$ = Expectation of class given a subset
 $V(Y|S_x)$ = Variance of class given all other residuals

Key Observation 2: Uncertainty in Explainability occurs due to all combinations of features that the explanation did not attribute to the network's decision



Case Study: Intervenability in Interpretability

Predictive Uncertainty in Explanations is the Residual

All other subsets 'not' chosen by the explanatory technique contributes to uncertainty

Explanation of Prediction Uncertainty of Explanation



Snout is not as highlighted as the jowls in explanation (not as important for decision)

However, snout is an important characteristic that is used to differentiate against other dogs. Hence, there is uncertainty on why this feature is not included in the attribution

Key Observation 2: Uncertainty in Explainability occurs due to all combinations of features that the explanation did not attribute to the network's decision

Case Study: Intervenability in Interpretability

Predictive Uncertainty in Explanations is the Residual

All other subsets **'not' chosen** by the explanatory technique contributes to uncertainty

Explanation of Prediction Uncertainty of Explanation



Snout is not as highlighted as the jowls in explanation (not as important for decision)

However, snout is an important characteristic that is used to differentiate against other dogs. Hence, there is uncertainty on why this feature is not included in the attribution

Not chosen features are intractable!

Case Study: Intervenability in Interpretability

Quantifying Interventions in Explainability

Contrastive explanations are an intelligent way of obtaining other subsets



$$V[y|S_x] = V[E(y|S_x)] + E(V[y|S_x])$$

Make it finite by only considering the subsets that change y

$$\left. \begin{array}{l} Y_1|S_{x1} \\ Y_2|S_{x2} \\ Y_3|S_{x3} \\ Y_4|S_{x4} \\ Y_5|S_{x5} \\ \vdots \\ Y_N|S_{xN} \end{array} \right\} \text{Variance}$$

Case Study: Intervenability in Interpretability

Quantifying Interventions in Explainability

Uncertainty in Explainability can be used to analyze Explanatory methods and Networks

- Is GradCAM better than GradCAM++?
- Is a SWIN transformer more reliable than VGG-16?

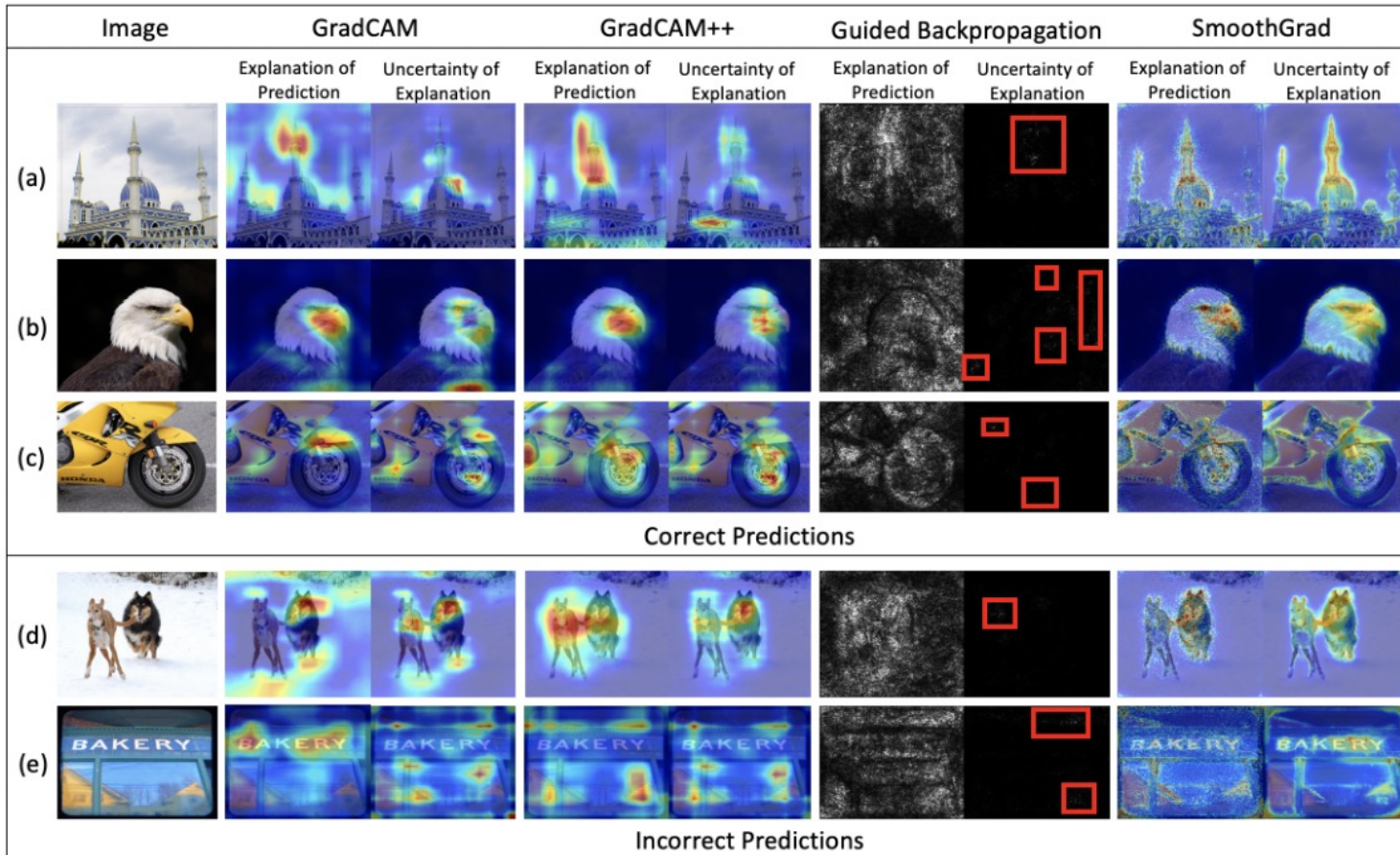
Need objective quantification of Intervention Residuals



Case Study: Intervenability in Interpretability

Quantifying Interventions in Explainability: mIOU

On incorrect predictions, the overlap of explanations and uncertainty is higher



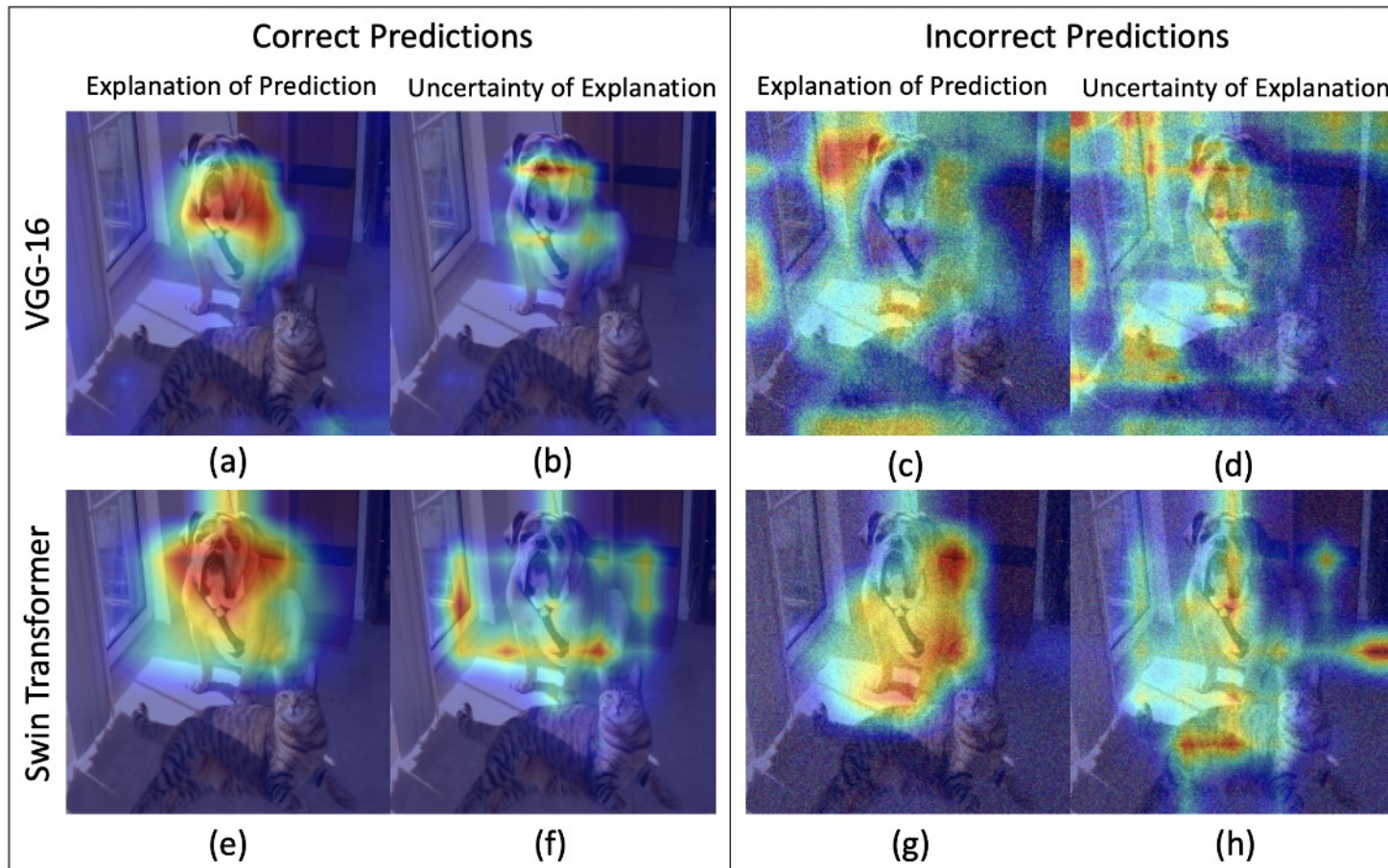
Objective Metric:
Intersection over Union (IoU)
between
explanation and
Uncertainty

Higher the IoU, higher the
uncertainty in explanation (or
less trustworthy is the
prediction)

Case Study: Intervenability in Interpretability

Quantifying Interventions in Explainability: SNR

Explanation and uncertainty are dispersed under noise (under low prediction confidence)



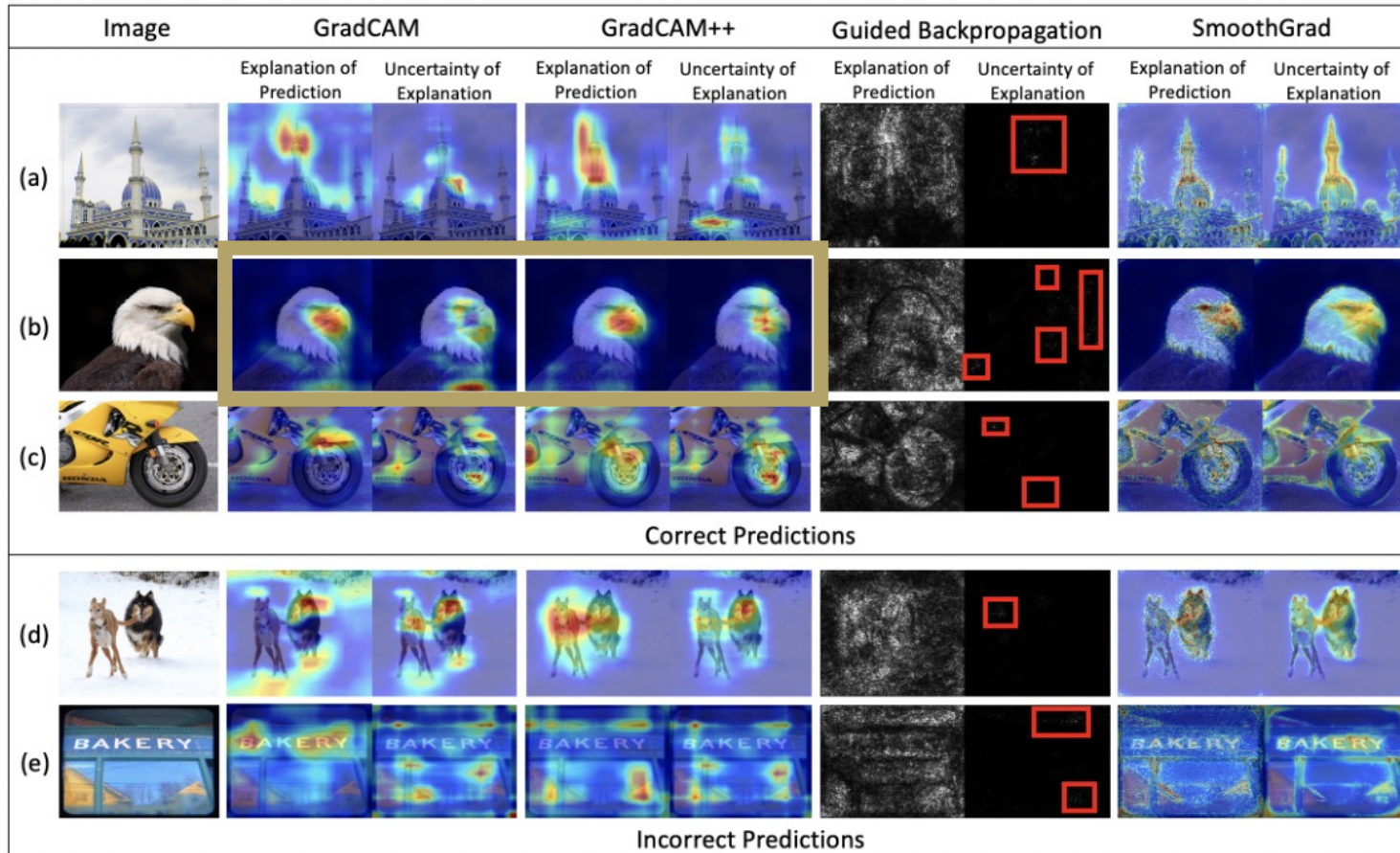
**Objective Metric:
Signal to Noise
Ratio of the
Uncertainty map**

Higher the SNR of
uncertainty, more is the
dispersal (or less trustworthy
is the prediction)

Case Study: Intervenability in Interpretability

Quantifying Interventions in Explainability: mIOU

On incorrect predictions, the overlap of explanations and uncertainty is higher



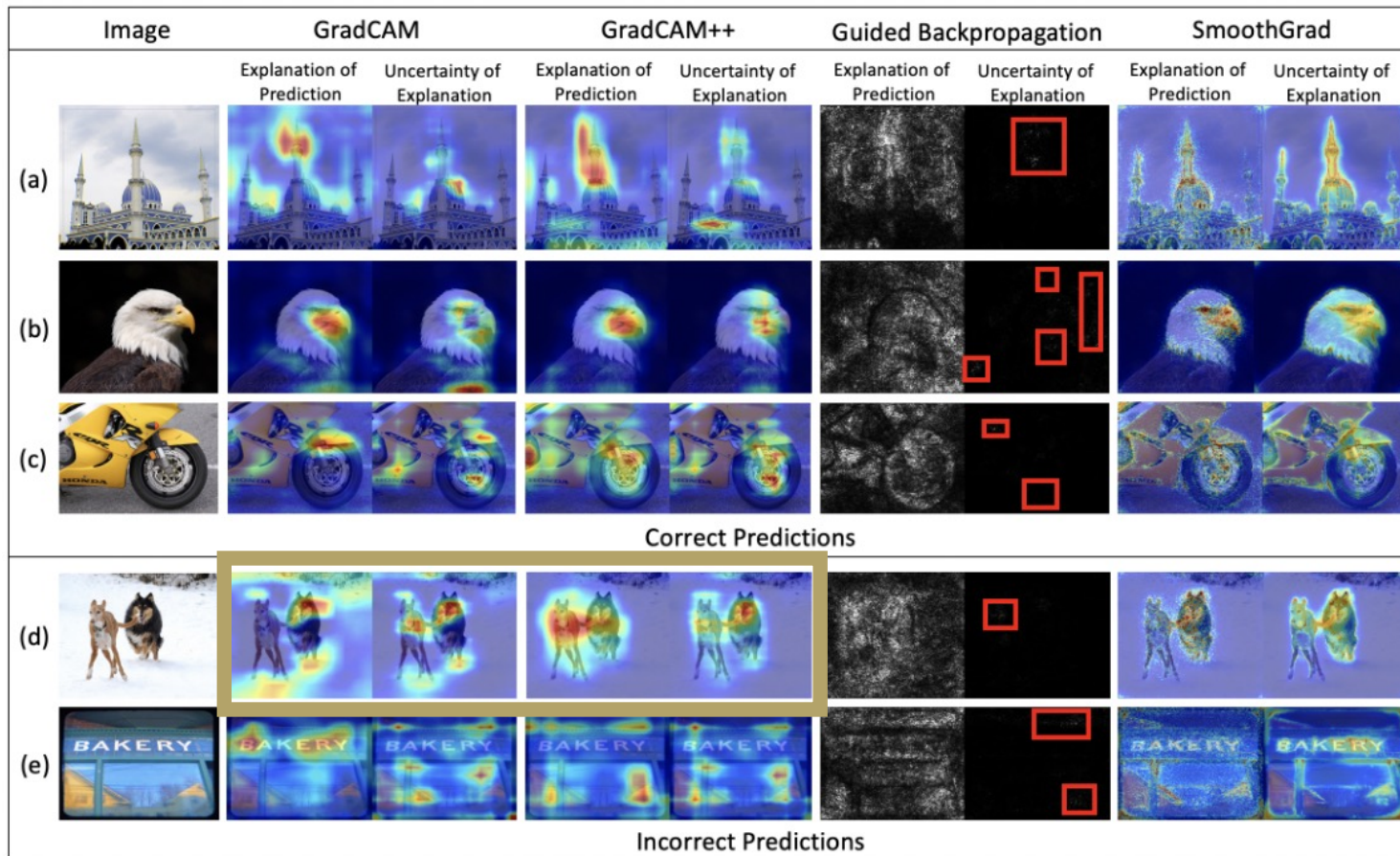
Objective Metric 1:
Intersection over Union (IoU)
between
explanation and
Uncertainty

Higher the IoU, higher the
uncertainty in explanation (or
less trustworthy is the
prediction)

Case Study: Intervenability in Interpretability

Quantifying Interventions in Explainability: mIOU

On incorrect predictions, the overlap of explanations and uncertainty is higher



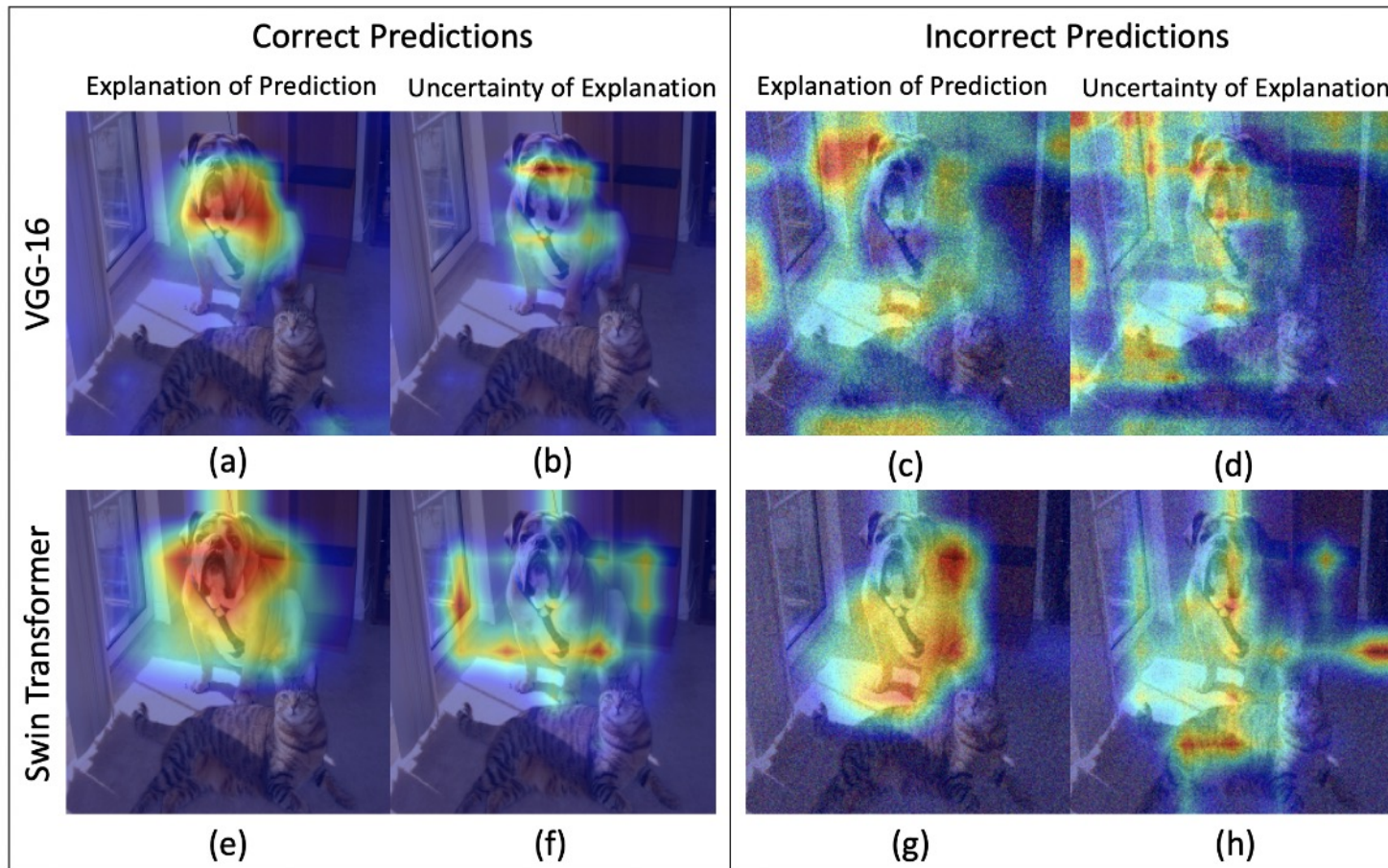
Objective Metric 1:
Intersection over Union (IoU)
between
explanation and
Uncertainty

Higher the IoU, higher the
uncertainty in explanation (or
less trustworthy is the
prediction)

Case Study: Intervenability in Interpretability

Quantifying Interventions in Explainability: SNR

Explanation and uncertainty are dispersed under noise (under low prediction confidence)



**Objective Metric 2:
Signal to Noise
Ratio of the
Uncertainty map**

Higher the SNR of uncertainty, more is the dispersal (or less trustworthy is the prediction)

Case Study: Intervenability in Interpretability

Challenges in Intervenability

The amenability of neural network decisions to human interventions



- **Not choosing interventions** causes **uncertainty** in the chosen interventions
- **Residuals** must be **analyzed** intelligently to **'trust or not to trust'** predictions at inference
- **Gradients quantify residual uncertainty**

Challenges:

- Choosing the type of Intervention: Explanation Evaluation
- **Residuals of Interventions: Uncertainty**

Intervenability

Through the Human Glass

The amenability of neural network decisions to human interventions



- **Assess: Causality**
- **Assure: Privacy**
- **Interpret: Interpretability**
- **Actuate: Prompting**
- **Verify: Benchmarking**

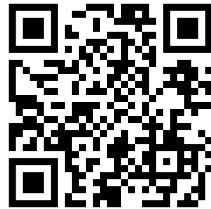
Intervenability in Benchmarking

Detection and Localization

CURE-TSD: Challenging Unreal and Real Environments for Traffic Sign Detection

Data Characteristics:

- 49 real and virtual sequences
- 300 frames in each sequence
- 12 different challenges including decolorization, codec error, lens blur etc.
- 5 progressively increasing levels in each challenge
- **Goal:** Detect and localize traffic signs



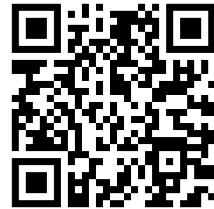
Intervenability in Benchmarking

Recognition

CURE-TSR: Challenging Unreal and Real Environments for Traffic Sign Recognition

Data Characteristics:

- 2 million real and virtual traffic sign images
- 14 Traffic signs including common signs like stop, no-right, no-left etc. and uncommon signs like goods-vehicles, priority lanes etc.
- 12 different challenges including decolorization, codec error, lens blur etc.
- 5 progressively increasingly levels in each challenge



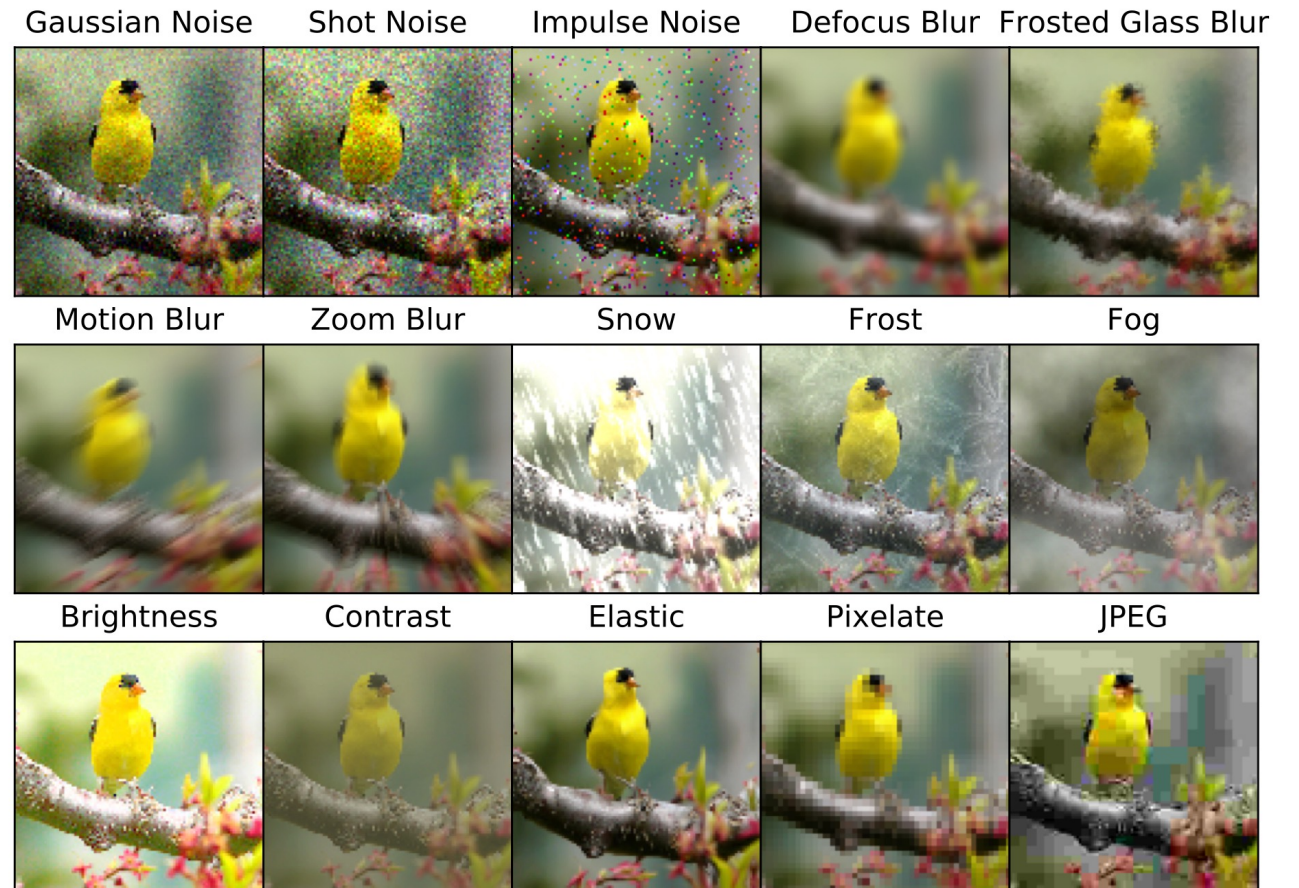
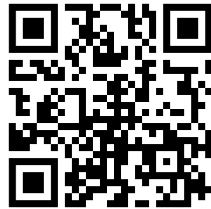
Intervenability in Benchmarking

Recognition

ImageNet-C: ImageNet-Corruptions

Data Characteristics:

- 3.75 million images
- 15 different challenges including decolorization, codec error, lens blur etc. for testing
- 4 different challenges for validation and training
- 5 progressively increasing levels in each challenge
- **Goal:** Recognize 1000 classes from ImageNet using pretrained networks



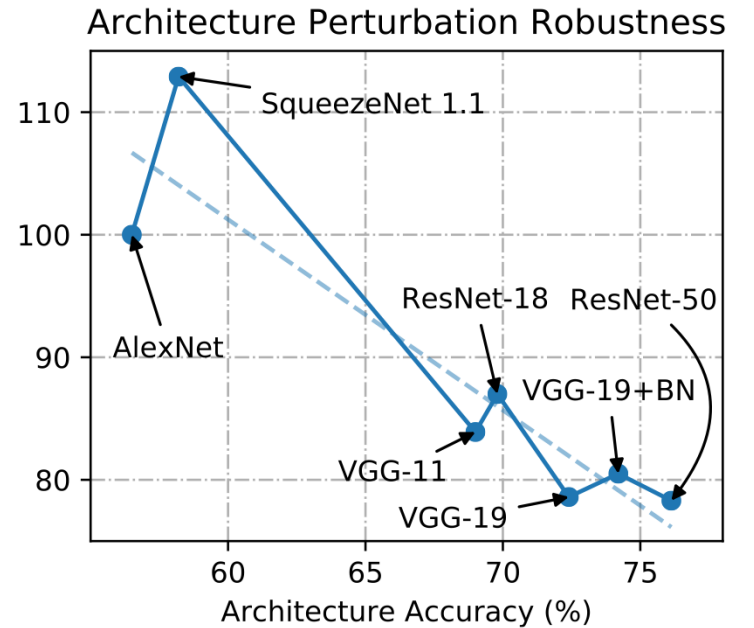
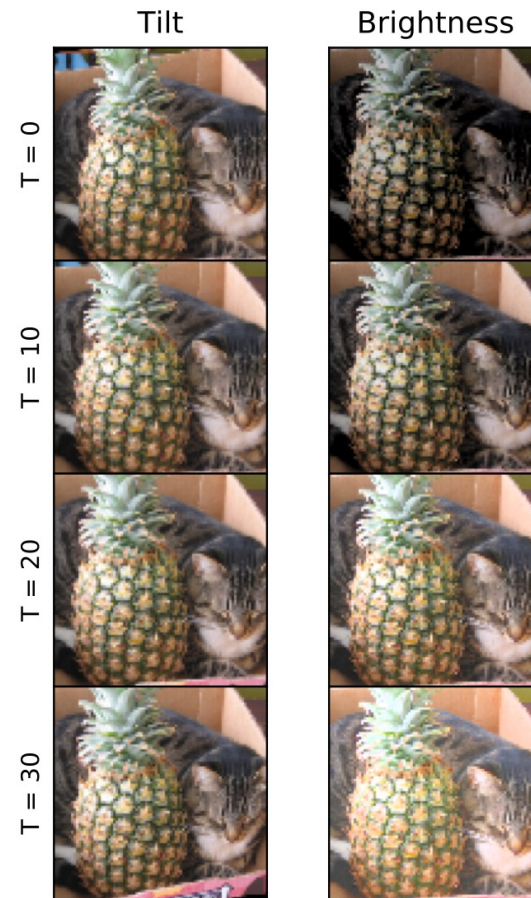
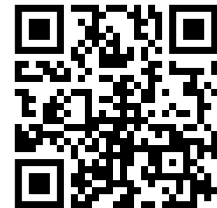
Intervenability in Benchmarking

Recognition

ImageNet-P: ImageNet-Perturbations

Data Characteristics:

- 5 million images
- 100 perturbations of 50000 images
- 10 frames of algorithmically generated perturbations for each image in ImageNet validation testset
- 10 common perturbations including brightness, tilt, motion etc.



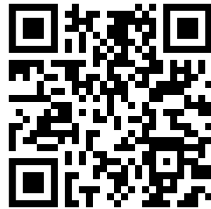
Intervenability in Benchmarking

Retrieval and Recognition

CURE-OR: Challenging Unreal and Real Environments for Object Recognition

Data Characteristics:

- 1 million images
- 100 common household objects and 10000 images per object
- 5 backgrounds, 5 object orientations, 5 devices, and 78 challenging conditions
- **Goal:** To recognize and retrieve the same object across backgrounds, orientations, devices, and challenging conditions



Challenge Type: None

Can	99.01
Tin	99.01
Beverage	98.95
Coke	98.95
Soda	98.95
Drink	70.87
Coffee Table	0.00
Furniture	0.00
Table	0.00
Couch	0.00
Book	0.00
Aluminium	0.00
Outdoors	0.00
Text	0.00
Drawing	0.00
Sketch	0.00
Diagram	0.00
Plan	0.00
Ice	0.00
Snow	0.00



Robust Neural Networks

Part 5: Conclusions and Future Directions

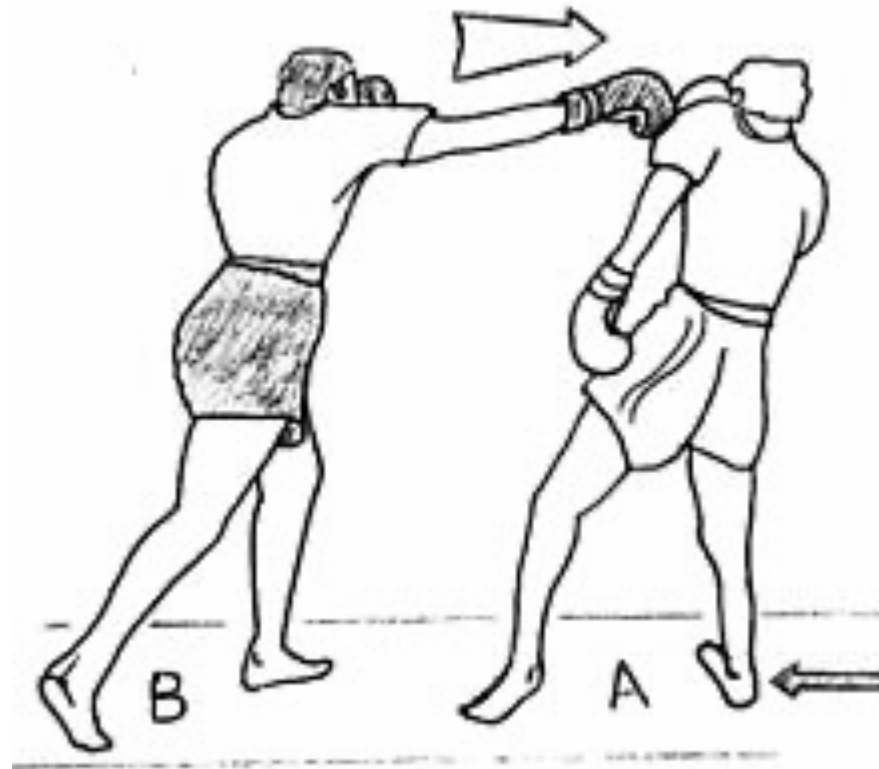
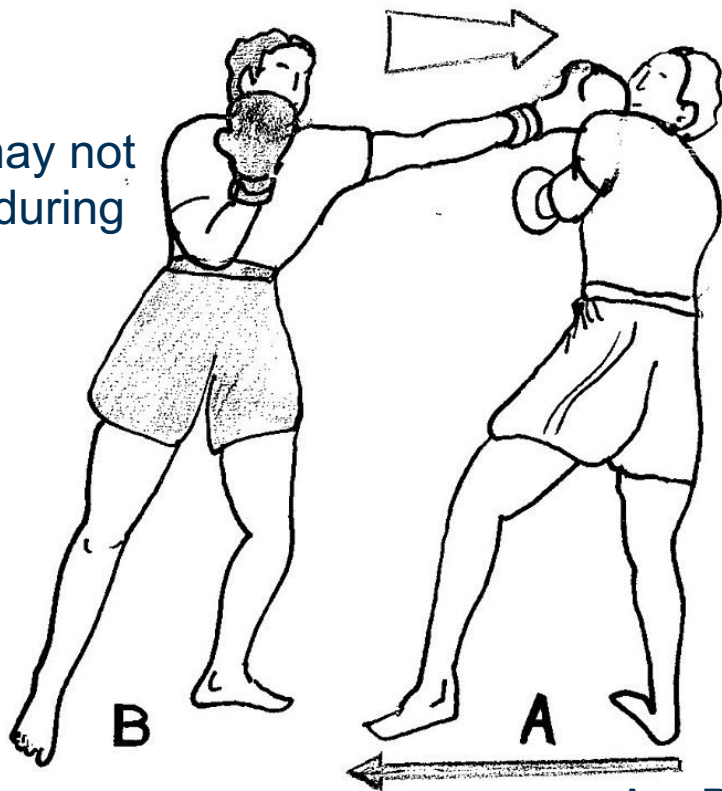


Mememes to Wrap it Up

Overcoming Challenges at Training

Novel data packs a 1-2 punch!

Novel data may not be available during training

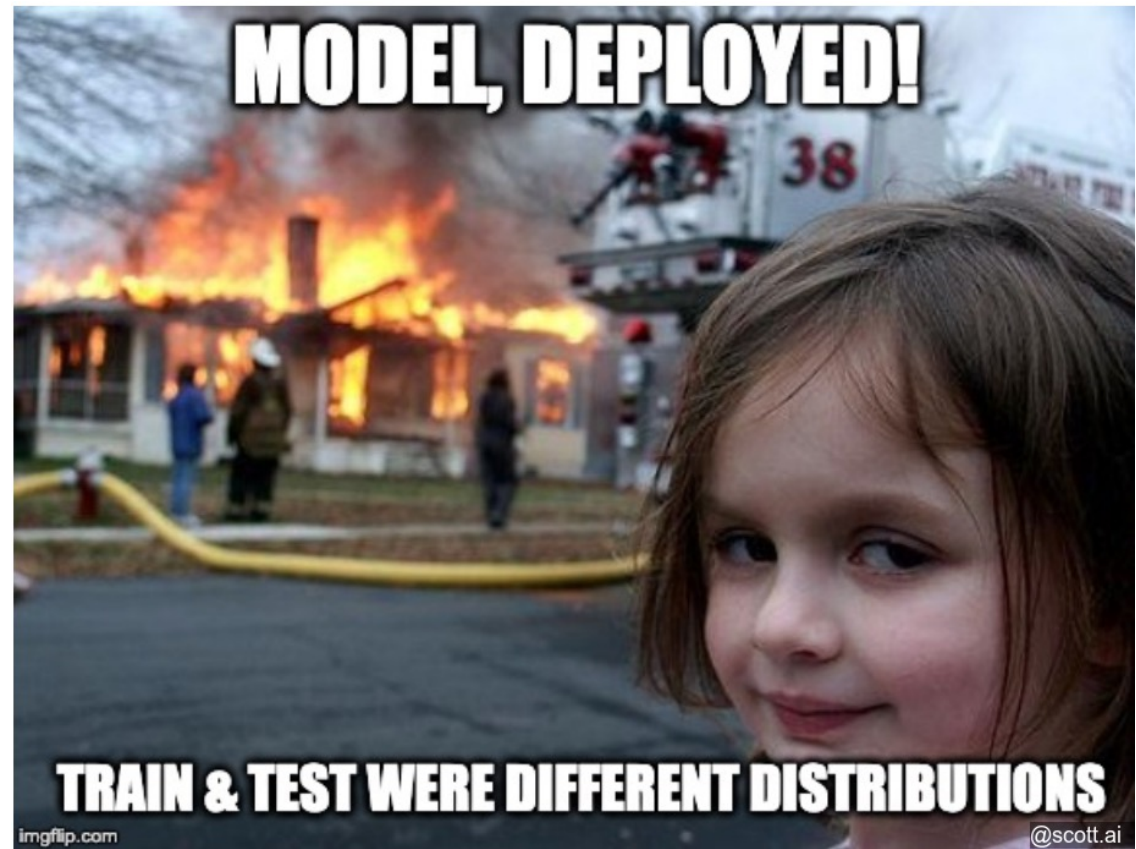
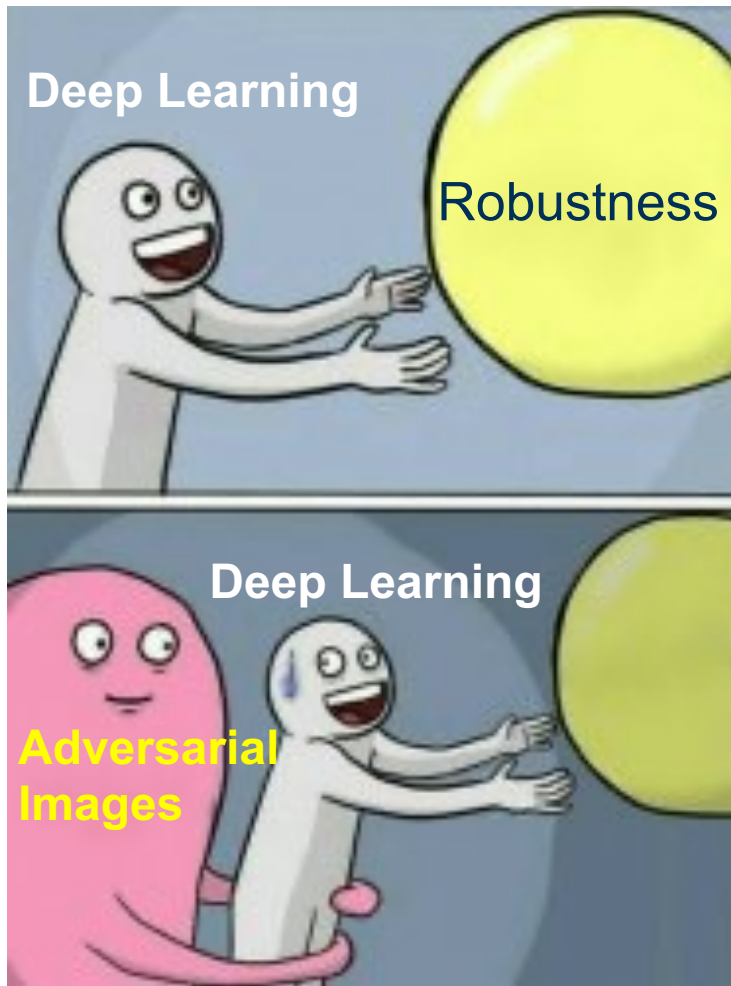


Even if available, novel data does not easily fit into either the earlier or later stages of training

A = Deep Neural Networks
B = Novel data

Mememes to Wrap it Up

Robustness at Inference

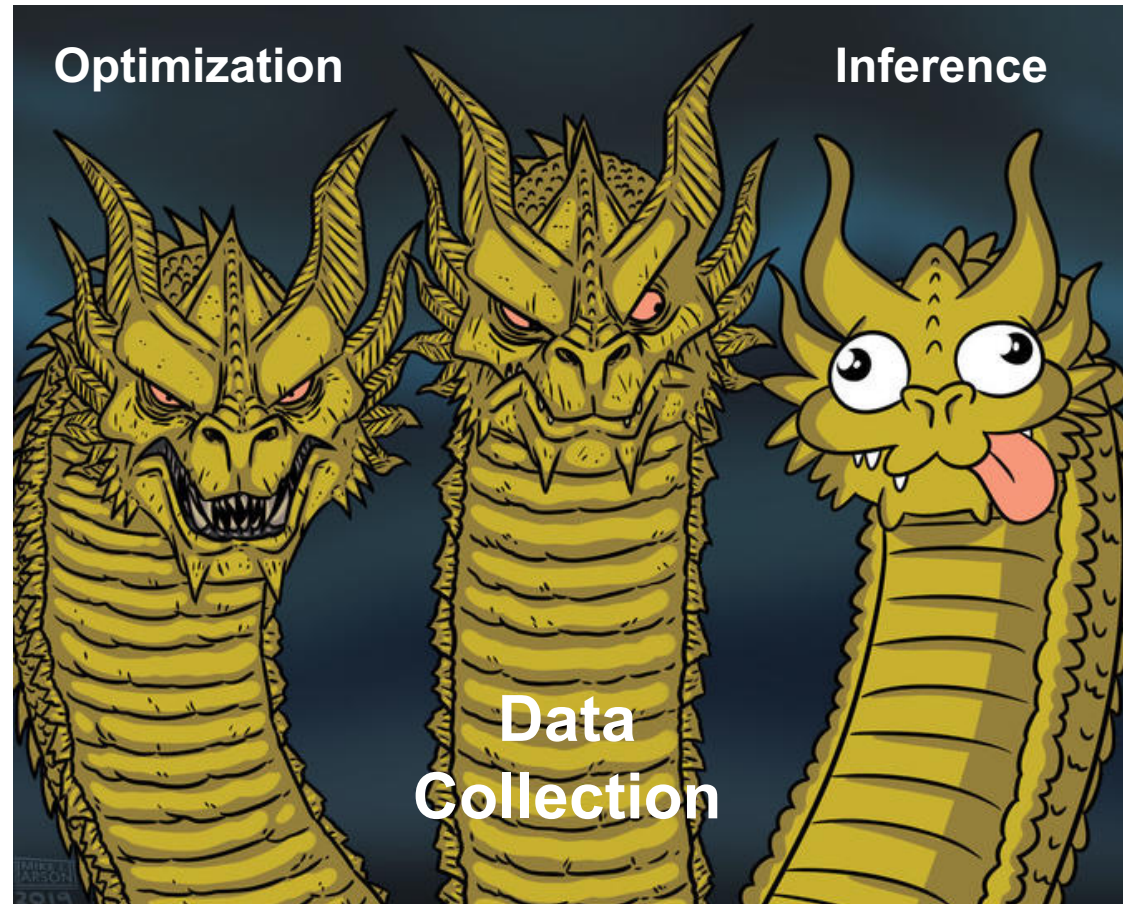


Cannot depend on training to construct robust models

Mememes to Wrap it Up

Robustness Research in the Inferential Stage of Neural Networks

Existing research on robustness focuses on data collection and optimization

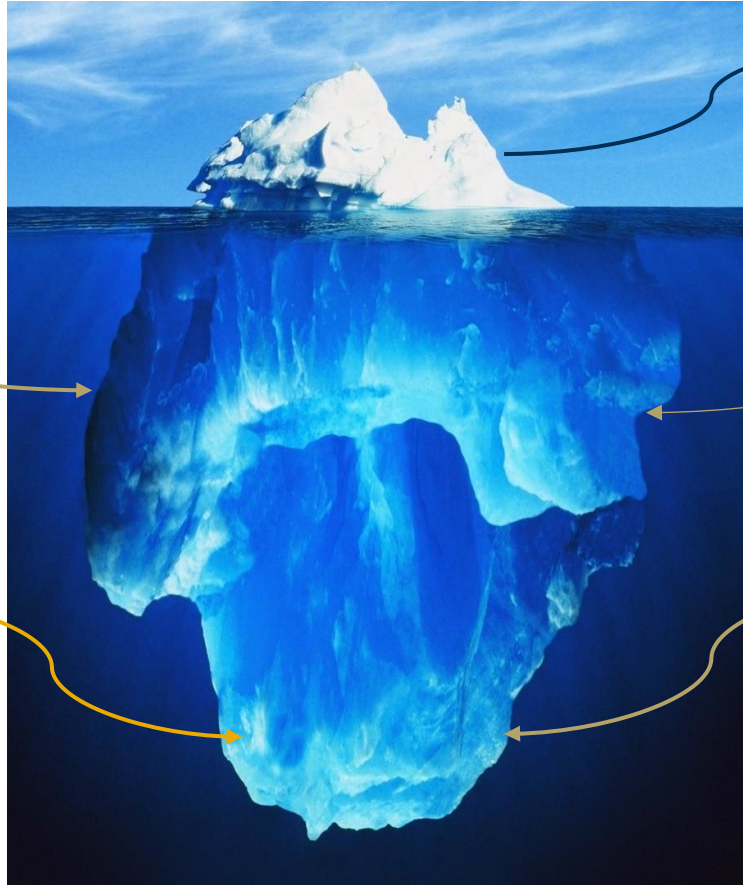
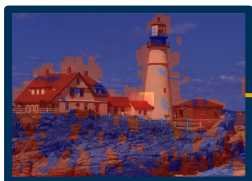
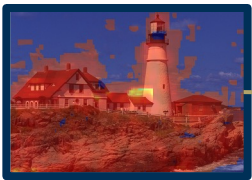


Mememes to Wrap it Up

Implicit Knowledge in Neural Networks

Trained Neural Networks have a wealth of implicit stored knowledge, waiting to be extracted at inference

Why P, rather than Q?



Traditional Why P?

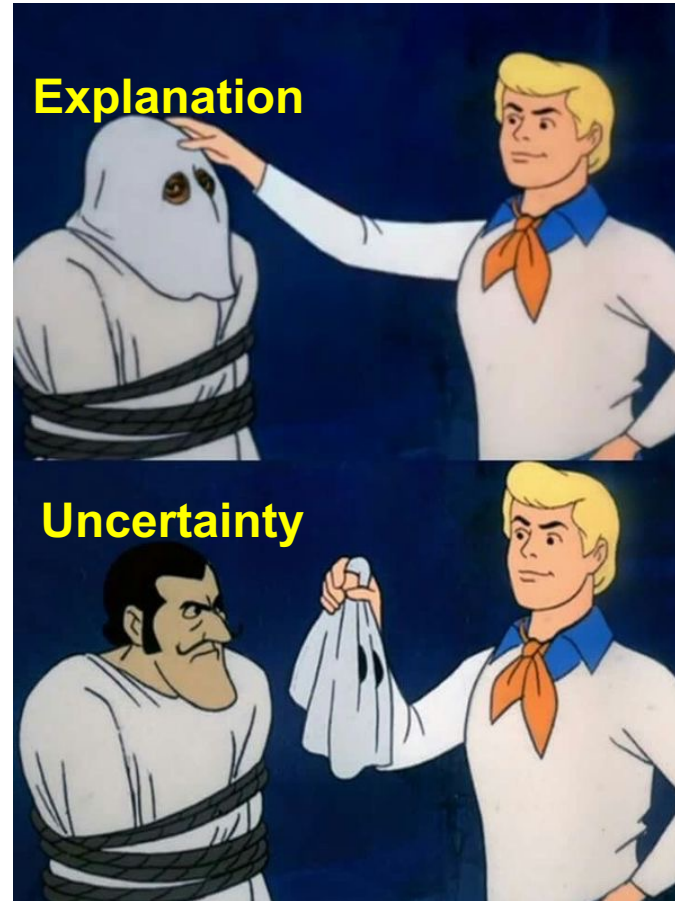


What if?

Mememes to Wrap it Up

Explainability Research is Just Uncertainty Research

Explanatory Evaluation reduces Uncertainty



Key Takeaways

Role of Gradients

- **Robustness** under distributional shift in domains, environments, and adversaries are **challenges** for neural networks
 - **Gradients at Inference** provide a **holistic solution** to the above challenges
- **Gradients** can help **traverse** through a trained and unknown **manifold**
 - They approximate **Fisher Information** on the projection
 - They can be **manipulated** by providing **contrast** classes
 - They can be used to construct **localized contrastive** manifolds
 - They provide **implicit knowledge** about **all classes**, when only **one data** point is available at inference
- Gradients are useful in a number of **Image Understanding** applications
 - Highlighting features of the current prediction as well as **counterfactual** data and **contrastive** classes
 - Providing **directional information** in anomaly detection
 - **Quantifying uncertainty** for out-of-distribution, corruption, and adversarial detection
 - Providing **expectancy mismatch** for human vision related applications



Future Directions

Research at Inference Stage

- **Test Time Augmentation (TTA) Research**
 - Multiple augmentations of data are passed through the network at inference
 - Research is in designing the best augmentations
- **Active Inference**
 - Utilize the knowledge in Neural Networks to *ask it to ask us*
 - Neural networks ask for the best augmentation of the data point given that one data point at inference
- **Uncertainty in Explainability, Label Interpretation, and Trust quantification**
 - Uncertainty research has to expand beyond model and data uncertainty
 - In some applications within medical and seismic communities, there is no agreed upon label for data. Uncertainty in label interpretation is its own research
- **Test-time Interventions for AI alignment**
 - Human interventions at test time to alter the decision-making process is essential trustworthy AI
 - Further research in intelligently involving experts in a non end-to-end framework is required



References

Gradient representations for Robustness, OOD, Anomaly, Novelty, and Adversarial Detection

- **Gradients for robustness against noise:** M. Prabhushankar, and G. AlRegib, "Introspective Learning : A Two-Stage Approach for Inference in Neural Networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, New Orleans, LA, Nov. 29 - Dec. 1 2022
- **Gradients for adversarial, OOD, corruption detection:** J. Lee, M. Prabhushankar, and G. AlRegib, "Gradient-Based Adversarial and Out-of-Distribution Detection," in *International Conference on Machine Learning (ICML) Workshop on New Frontiers in Adversarial Machine Learning*, Baltimore, MD, Jul. 2022.
- **Gradients for Open set recognition:** Lee, Jinsol, and Ghassan AlRegib. "Open-Set Recognition With Gradient-Based Representations." *2021 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2021.
- **GradCon for Anomaly Detection:** Kwon, G., Prabhushankar, M., Temel, D., & AlRegib, G. (2020, August). Backpropagated gradient representations for anomaly detection. In *European Conference on Computer Vision* (pp. 206-226). Springer, Cham.
- **Gradients for adversarial, OOD, corruption detection :** J. Lee, C. Lehman, M. Prabhushankar, and G. AlRegib, "Probing the Purview of Neural Networks via Gradient Analysis," in *IEEE Access*, Mar. 21 2023.
- **Gradients for Novelty Detection:** Kwon, G., Prabhushankar, M., Temel, D., & AlRegib, G. (2020, October). Novelty detection through model-based characterization of neural networks. In *2020 IEEE International Conference on Image Processing (ICIP)* (pp. 3179-3183). IEEE.
- **Gradient-based Image Quality Assessment:** G. Kwon*, M. Prabhushankar*, D. Temel, and G. AlRegib, "Distorted Representation Space Characterization Through Backpropagated Gradients," in *IEEE International Conference on Image Processing (ICIP)*, Taipei, Taiwan, Sep. 2019.

Explainability in Neural Networks

- **Explanatory paradigms:** AlRegib, G., & Prabhushankar, M. (2022). Explanatory Paradigms in Neural Networks: Towards relevant and contextual explanations. *IEEE Signal Processing Magazine*, 39(4), 59-72.
- **Contrastive Explanations:** Prabhushankar, M., Kwon, G., Temel, D., & AlRegib, G. (2020, October). Contrastive explanations in neural networks. In *2020 IEEE International Conference on Image Processing (ICIP)* (pp. 3289-3293). IEEE.
- **Explainability in Limited Label Settings:** M. Prabhushankar, and G. AlRegib, "Extracting Causal Visual Features for Limited Label Classification," in *IEEE International Conference on Image Processing (ICIP)*, Sept. 2021.
- **Explainability through Expectancy-Mismatch:** M. Prabhushankar and G. AlRegib, "Stochastic Surprisal: An Inferential Measurement of Free Energy in Neural Networks," in *Frontiers in Neuroscience, Perception Science*, Volume 17, Feb. 09 2023.



References

Self Supervised Learning

- **Weakly supervised Contrastive Learning:** K. Kokilepersaud, S. Trejo Corona, M. Prabhushankar, G. AlRegib, C. Wykoff, "Clinically Labeled Contrastive Learning for OCT Biomarker Classification," in *IEEE Journal of Biomedical and Health Informatics*, 2023, May. 15 2023.
- **Contrastive Learning for Fisheye Images:** K. Kokilepersaud, M. Prabhushankar, Y. Yarici, G. AlRegib, and A. Parchami, "Exploiting the Distortion-Semantic Interaction in Fisheye Data," in *Open Journal of Signal Processing*, Apr. 28 2023.
- **Contrastive Learning for Severity Detection:** K. Kokilepersaud, M. Prabhushankar, G. AlRegib, S. Trejo Corona, C. Wykoff, "Gradient Based Labeling for Biomarker Classification in OCT," in *IEEE International Conference on Image Processing (ICIP)*, Bordeaux, France, Oct. 16-19 2022
- **Contrastive Learning for Seismic Images:** K. Kokilepersaud, M. Prabhushankar, and G. AlRegib, "Volumetric Supervised Contrastive Learning for Seismic Semantic Segmentation," in *International Meeting for Applied Geoscience & Energy (IMAGE)*, Houston, TX, , Aug. 28-Sept. 1 2022

Human Vision and Behavior Prediction

- **Pedestrian Trajectory Prediction:** C. Zhou, G. AlRegib, A. Parchami, and K. Singh, "TrajPRed: Trajectory Prediction With Region-Based Relation Learning," *IEEE Transactions on Intelligent Transportation Systems*, submitted on Dec. 28 2022.
- **Human Visual Saliency in trained Neural Nets:** Y. Sun, M. Prabhushankar, and G. AlRegib, "Implicit Saliency in Deep Neural Networks," in *IEEE International Conference on Image Processing (ICIP)*, Abu Dhabi, United Arab Emirates, Oct. 2020.
- **Human Image Quality Assessment:** D. Temel, M. Prabhushankar and G. AlRegib, "UNIQUE: Unsupervised Image Quality Estimation," in *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1414-1418, Oct. 2016.

Open-source Datasets to assess Robustness

- **CURE-TSD:** D. Temel, M-H. Chen, and G. AlRegib, "Traffic Sign Detection Under Challenging Conditions: A Deeper Look Into Performance Variations and Spectral Characteristics," in *IEEE Transactions on Intelligent Transportation Systems*, Jul. 2019
- **CURE-TSR:** D. Temel, G. Kwon*, M. Prabhushankar*, and G. AlRegib, "CURE-TSR: Challenging Unreal and Real Environments for Traffic Sign Recognition," in *Advances in Neural Information Processing Systems (NIPS) Workshop on Machine Learning for Intelligent Transportation Systems*, Long Beach, CA, Dec. 2017
- **CURE-OR:** D. Temel*, J. Lee*, and G. AlRegib, "CURE-OR: Challenging Unreal and Real Environments for Object Recognition," in *IEEE International Conference on Machine Learning and Applications (ICMLA)*, Orlando, FL, Dec. 2018



References

Active Learning

- **Active Learning and Training with High Information Content:** R. Benkert, M. Prabhushankar, G. AlRegib, A. Parchami, and E. Corona, "Gaussian Switch Sampling: A Second Order Approach to Active Learning," in *IEEE Transactions on Artificial Intelligence (TAI)*, Feb. 05 2023
- **Active Learning Dataset on vision and LIDAR data:** Y. Logan, R. Benkert, C. Zhou, K. Kokilepersaud, M. Prabhushankar, G. AlRegib, K. Singh, E. Corona and A. Parchami, "FOCAL: A Cost-Aware Video Dataset for Active Learning," *IEEE Transactions on Circuits and Systems for Video Technology*, submitted on Apr. 29 2023
- **Active Learning on OOD data:** R. Benkert, M. Prabhushankar, and G. AlRegib, "Forgetful Active Learning With Switch Events: Efficient Sampling for Out-of-Distribution Data," in *IEEE International Conference on Image Processing (ICIP)*, Bordeaux, France, Oct. 16-19 2022
- **Active Learning for Biomedical Images:** Y. Logan, R. Benkert, A. Mustafa, G. Kwon, G. AlRegib, "Patient Aware Active Learning for Fine-Grained OCT Classification," in *IEEE International Conference on Image Processing (ICIP)*, Bordeaux, France, Oct. 16-19 2022

Uncertainty Estimation

- **Gradient-based Uncertainty:** J. Lee and G. AlRegib, "Gradients as a Measure of Uncertainty in Neural Networks," in *IEEE International Conference on Image Processing (ICIP)*, Abu Dhabi, United Arab Emirates, Oct. 2020
- **Gradient-based Visual Uncertainty:** M. Prabhushankar, and G. AlRegib, "VOICE: Variance of Induced Contrastive Explanations to Quantify Uncertainty in Neural Network Interpretability," *Journal of Selected Topics in Signal Processing*, submitted on Aug. 27, 2023.
- **Uncertainty Visualization in Seismic Images:** R. Benkert, M. Prabhushankar, and G. AlRegib, "Reliable Uncertainty Estimation for Seismic Interpretation With Prediction Switches," in *International Meeting for Applied Geoscience & Energy (IMAGE)*, Houston, TX, , Aug. 28-Sept. 1 2022.
- **Uncertainty and Disagreements in Label Annotations:** C. Zhou, M. Prabhushankar, and G. AlRegib, "On the Ramifications of Human Label Uncertainty," in *NeurIPS 2022 Workshop on Human in the Loop Learning*, Oct. 27 2022
- **Uncertainty in Saliency Estimation:** T. Alshawi, Z. Long, and G. AlRegib, "Unsupervised Uncertainty Estimation Using Spatiotemporal Cues in Video Saliency Detection," in *IEEE Transactions on Image Processing*, vol. 27, pp. 2818-2827, Jun. 2018.



Tutorial Materials

Accessible Online



<https://alregib.ece.gatech.edu/aaai-2024-tutorial/>
{alregib, mohit.p}@gatech.edu

AAAI 2024 Tutorial



Presented by: *Ghassan AlRegib, and Mohit Prabhushankar*
Georgia Institute of Technology

www.ghassanalregib.info

Duration: Half Day (3 hours, 30 mins)

Title: Formalizing Robustness in Neural Networks: Explainability, Uncertainty, and Intervenability

