

Visual Explainability in Machine Learning

Lecture 6: Robustness as Explanatory Proxy



Ghassan AlRegib, PhD
Professor



Mohit Prabhushankar, PhD
Postdoctoral Fellow

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

Dec 5, 2023

Short Course Materials

Accessible Online



SCAN ME



Title: Visual Explainability in Machine Learning

Presented by: *Ghassan AlRegib, and Mohit Prabhushankar*

Omni Lab for Intelligent Visual Engineering and Science (OLIVES)

School of Electrical and Computer Engineering

Georgia Institute of Technology, Atlanta, USA

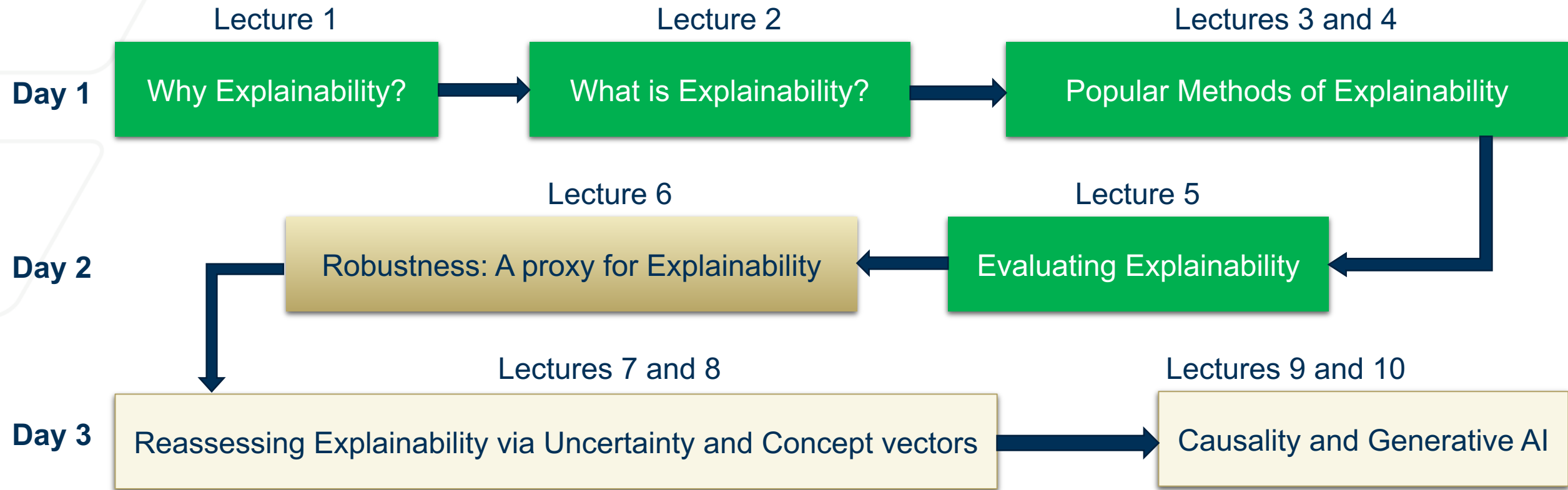
<https://alregib.ece.gatech.edu/>

<https://alregib.ece.gatech.edu/sps-education-short-course/>
{alregib, mohit.p}@gatech.edu

Short Course

Course Outline

Day 1: Define and Detail; Day 2: Evaluate; Day 3: Reassess



Outline

Lecture 6: Robustness as Explanatory Proxy

- Robustness and Explanations
- Deep Learning at Inference
 - Robustness under novel data
 - Challenges
 - Gradient Information
- Gradients as Robustness Features
 - Anomaly Detection
 - Out-of-Distribution Detection
 - Adversarial Detection
 - Corruption Detection
 - Gradients for Robust Predictions

Outline

Lecture 6: Robustness as Explanatory Proxy

- **Robustness and Explanations**
- Deep Learning at Inference
 - Robustness under novel data
 - Challenges
 - Gradient Information
- Gradients as Robustness Features
 - Anomaly Detection
 - Out-of-Distribution Detection
 - Adversarial Detection
 - Corruption Detection
 - Gradients for Robust Predictions

Robustness and Explainability

Robustness

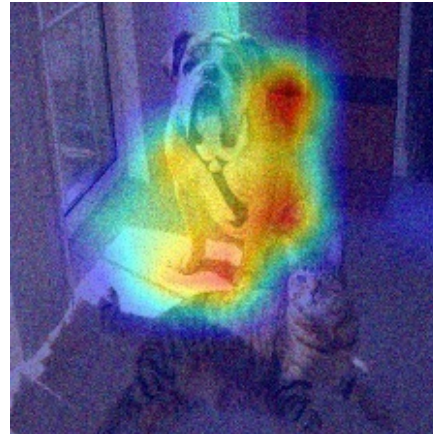
Robustness: The ability of a system to make accurate predictions when encountering novel data

Robustness \longleftrightarrow **Explainability**

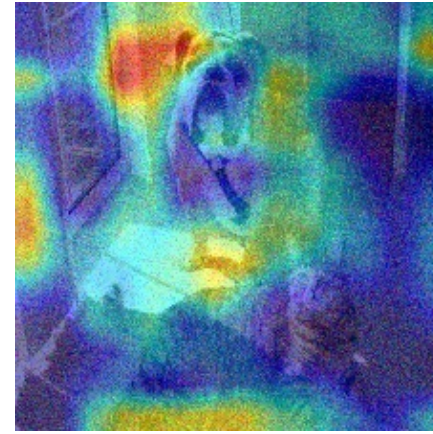
Assumption: More robust is the model, better is its Explainability¹



Distorted cat-dog



GradCAM from Swin Transformer



GradCAM from VGG-16

Robustness and Explainability

Objective in Lecture 6: Gradients as Features for both Explainability and Robustness

Robustness: The ability of a system to make accurate predictions when encountering novel data

Features for
Robustness

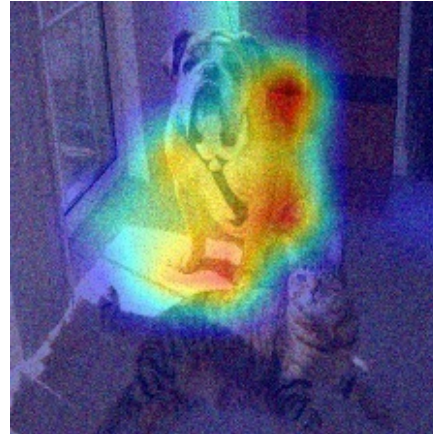


Features for
Explainability

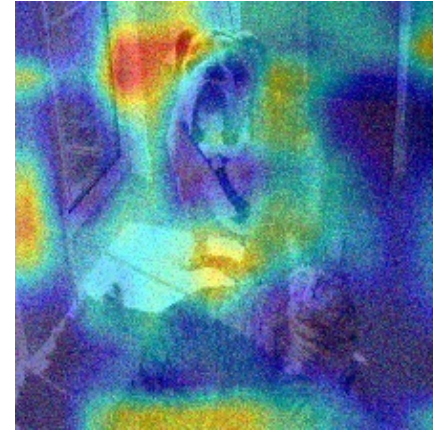
Assumption: More robust is the model, better is its Explainability¹



Distorted cat-dog



GradCAM from
Swin Transformer



GradCAM from
VGG-16

Features = Gradients

Outline

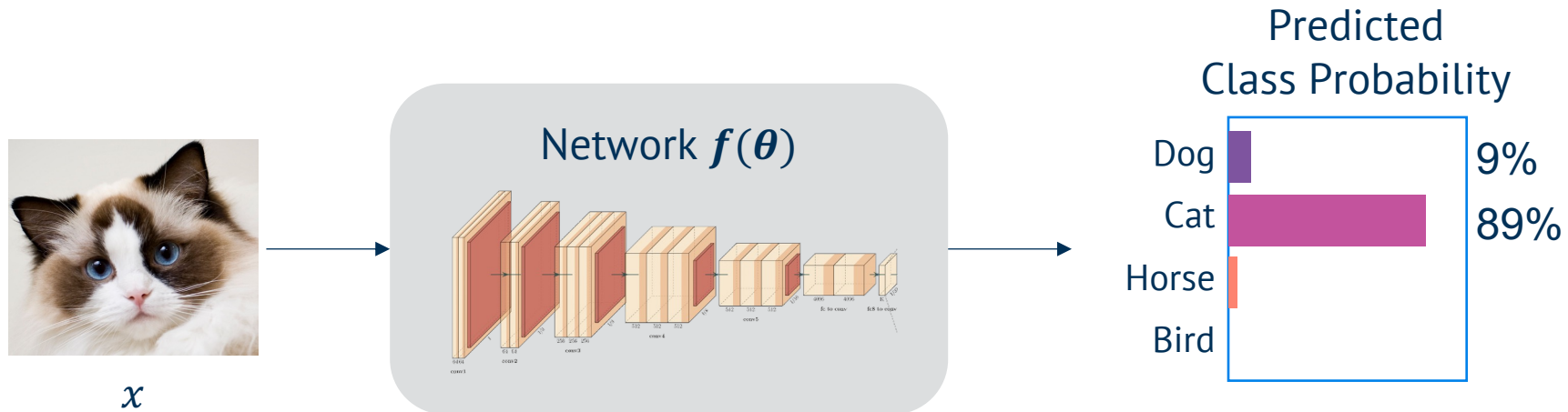
Lecture 6: Robustness as Explanatory Proxy

- Robustness and Explanations
- **Deep Learning at Inference**
 - Robustness under novel data
 - Challenges
 - Gradient Information
- Gradients as Robustness Features
 - Anomaly Detection
 - Out-of-Distribution Detection
 - Adversarial Detection
 - Corruption Detection
 - Gradients for Robust Predictions

Deep Learning at Inference

Classification

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



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

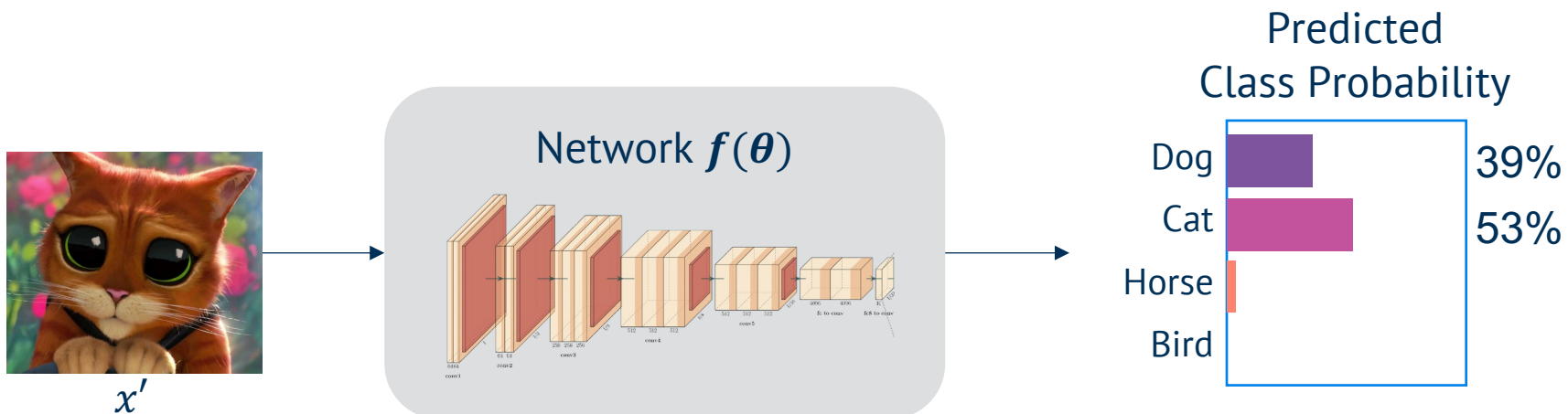
\hat{y} = Logits
 y = Predicted Class
 $p(\hat{y})$ = Probabilities
 $f(\cdot)$ = Trained Network
 χ = Training data

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

Deep Learning at Inference

Robust Classification in Deep Networks

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



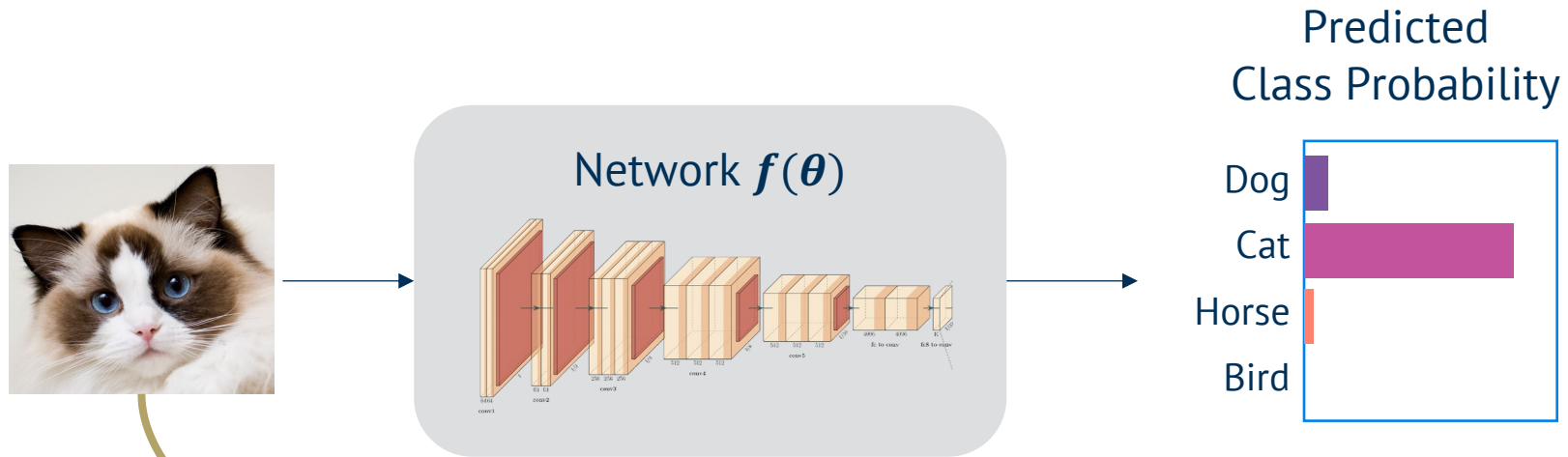
$$\begin{aligned} \hat{y} &= f(x' + \epsilon) & \hat{y} &= \text{Logits} \\ y &= \operatorname{argmax}_i \hat{y} & y &= \text{Predicted Class} \\ p(\hat{y}) &= T(f(x' + \epsilon)) & p(\hat{y}) &= \text{Probabilities} \\ & & f(\cdot) &= \text{Trained Network} \\ & & \chi &= \text{Training data} \\ & &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

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

Deep Learning at Inference

Fisher Information

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

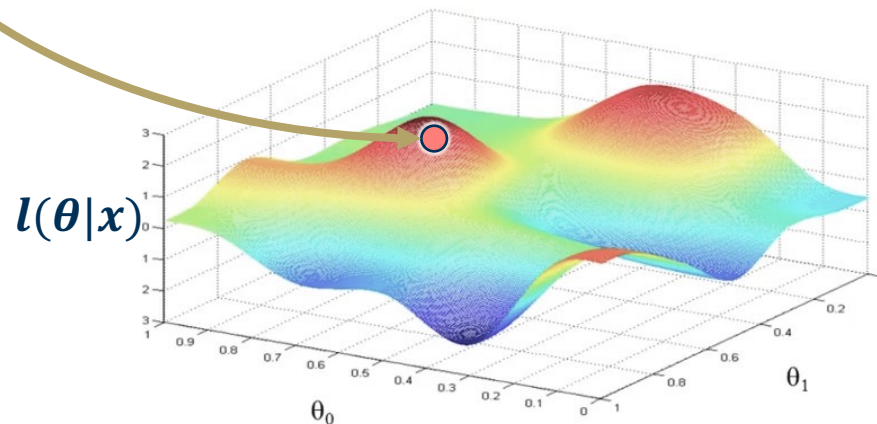


Fisher Information

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

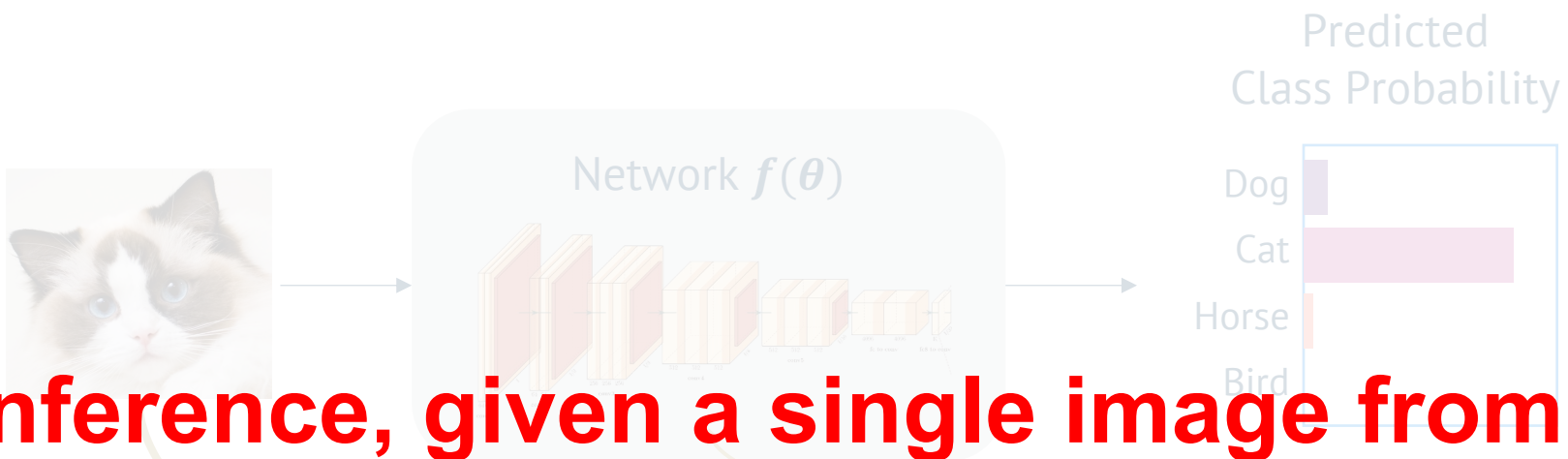
θ = Statistic of distribution
 $\ell(\theta | x)$ = Likelihood function

Likelihood function

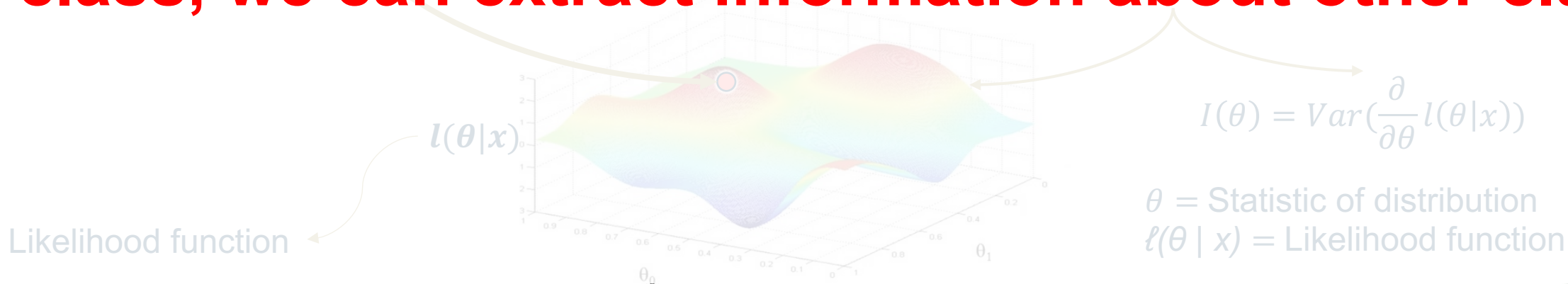


Deep Learning at Inference

Information at Inference



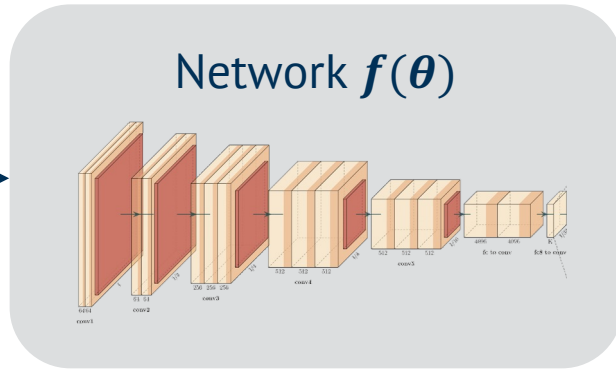
At inference, given a single image from a single class, we can extract information about other classes



Deep Learning at Inference

Case Study: Explainability

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



Beak
Neck
Legs
Feathers
Water
Grass
Teeth
.
.

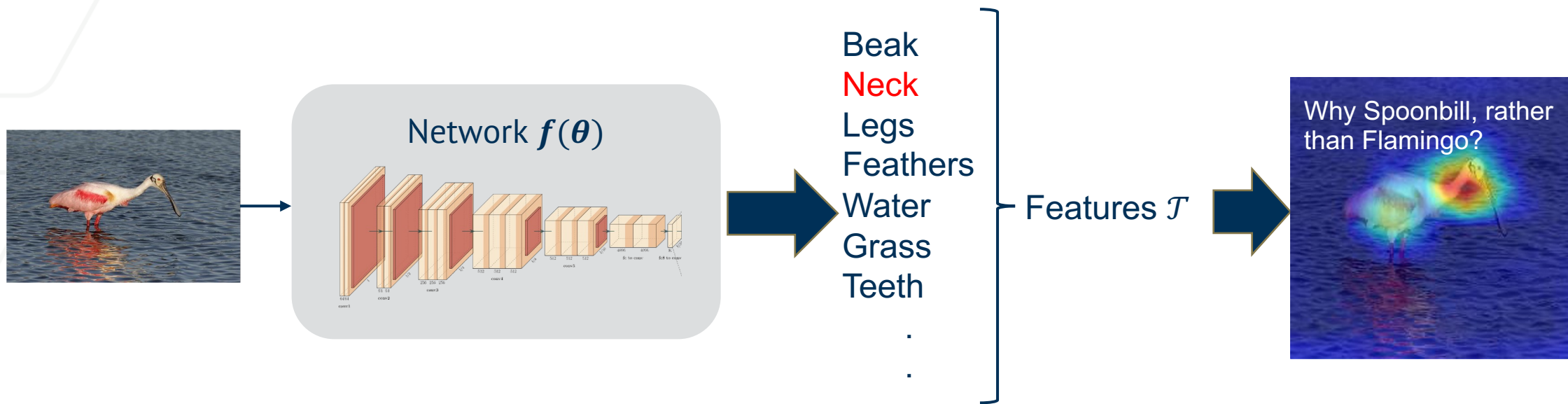
Features \mathcal{T}



Deep Learning 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 show that gradients store this information and can be used as features for robust predictions

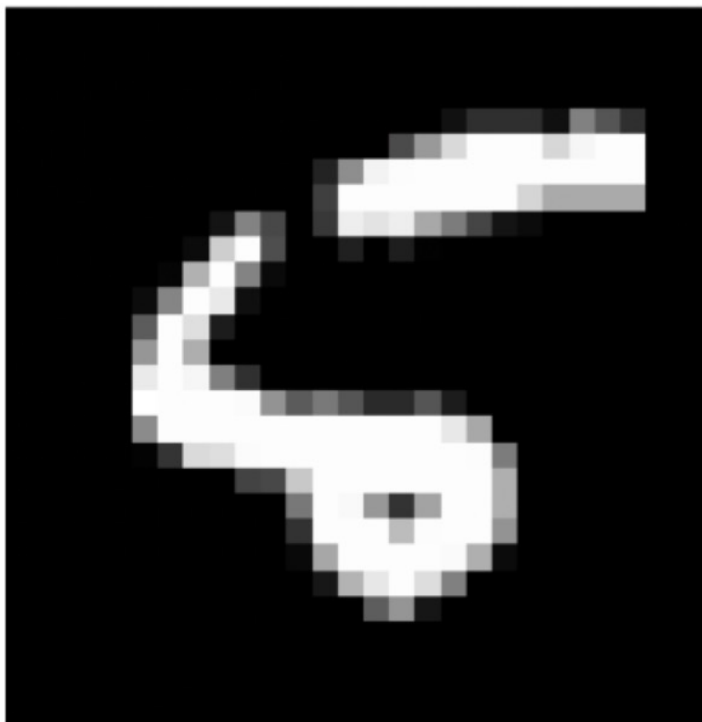
Deep Learning at Inference

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



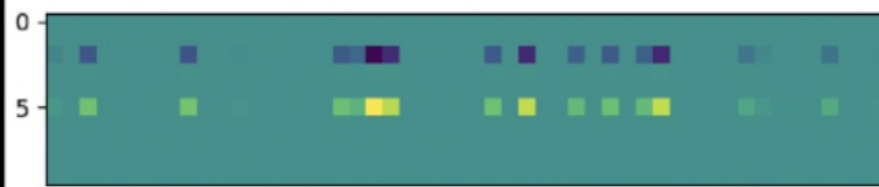
Input Image x



Why 5, rather than 0?



Why 5, rather than 1?



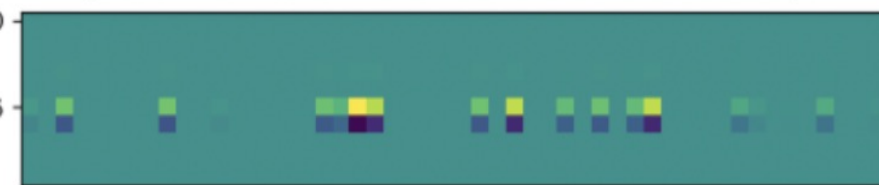
Why 5, rather than 2?



Why 5, rather than 4?



Why 5, rather than 5?



Why 5, rather than 6?

Deep Learning at Inference

Gradients as Features

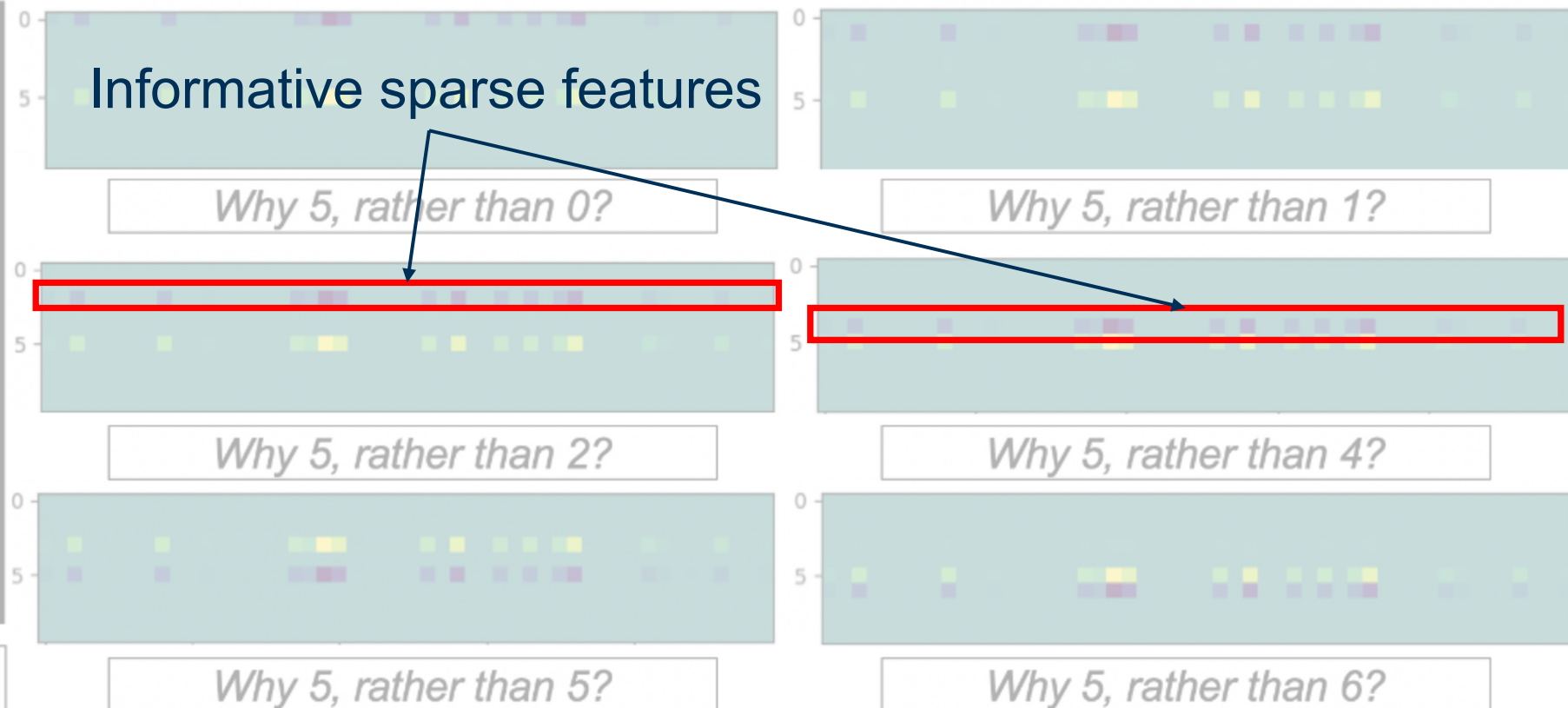


Introspective Learning: A Two-stage Approach for Inference in Neural Networks

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



Input Image x



Deep Learning at Inference

Gradients as Features



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

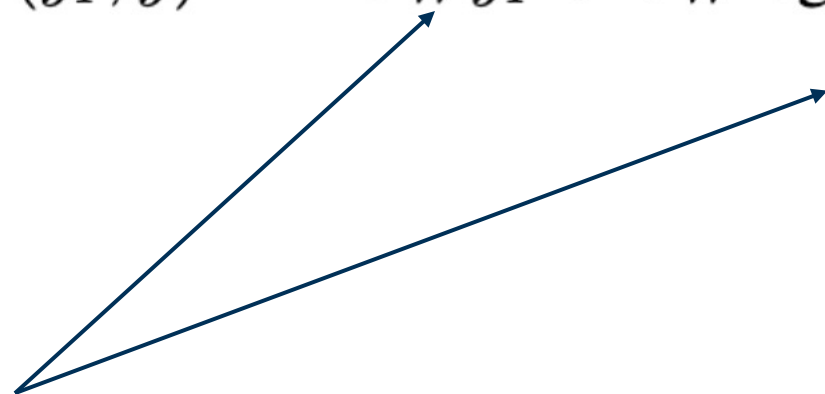
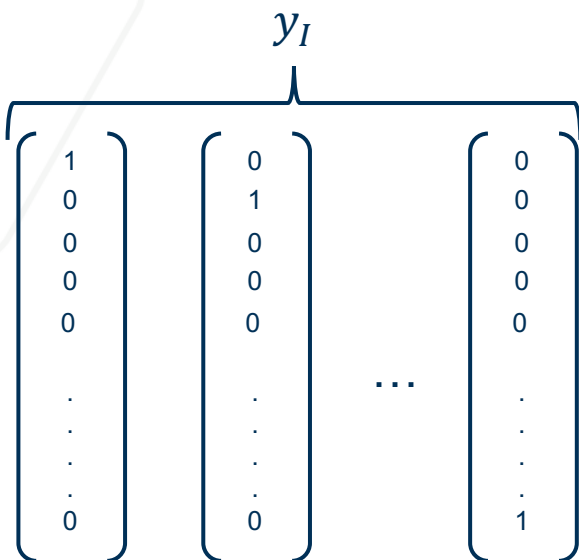
For a well-trained network, the gradients are robust

∇_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}

Outline

Lecture 6: Robustness as Explanatory Proxy

- Robustness and Explanations
- Deep Learning at Inference
 - Robustness under novel data
 - Challenges
 - Gradient Information
- **Gradients as Robustness Features**
 - **Anomaly Detection**
 - Out-of-Distribution Detection
 - Adversarial Detection
 - Corruption Detection
 - Gradients for Robust Predictions

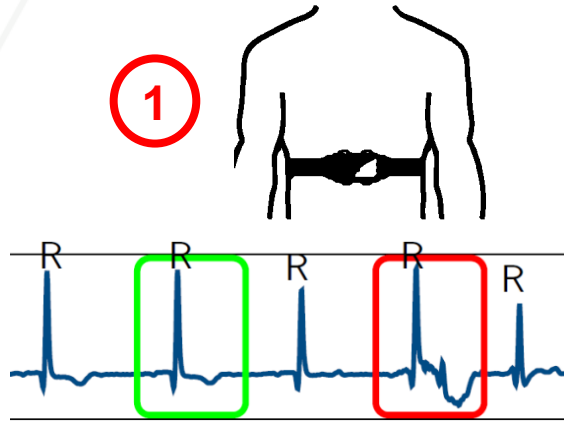
Gradients as Robustness Features

Anomalies



Backpropagated Gradient Representations for Anomaly Detection

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

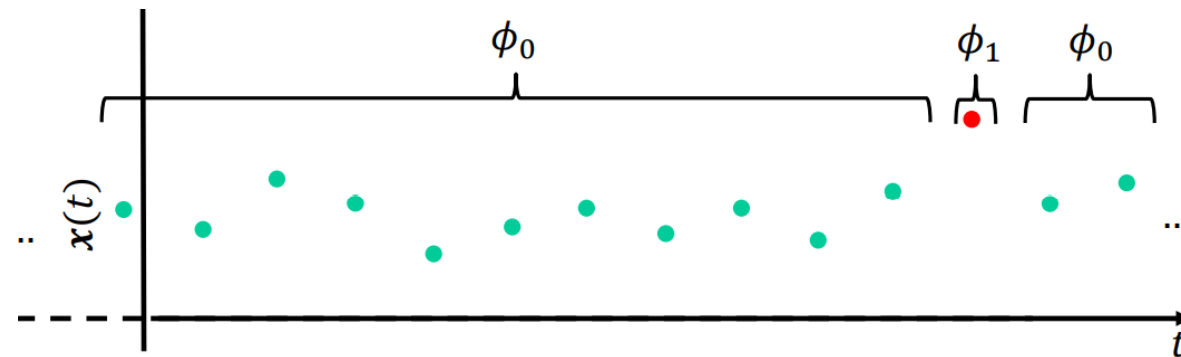


Statistical Definition:

- 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

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



Gradients as Robustness Features

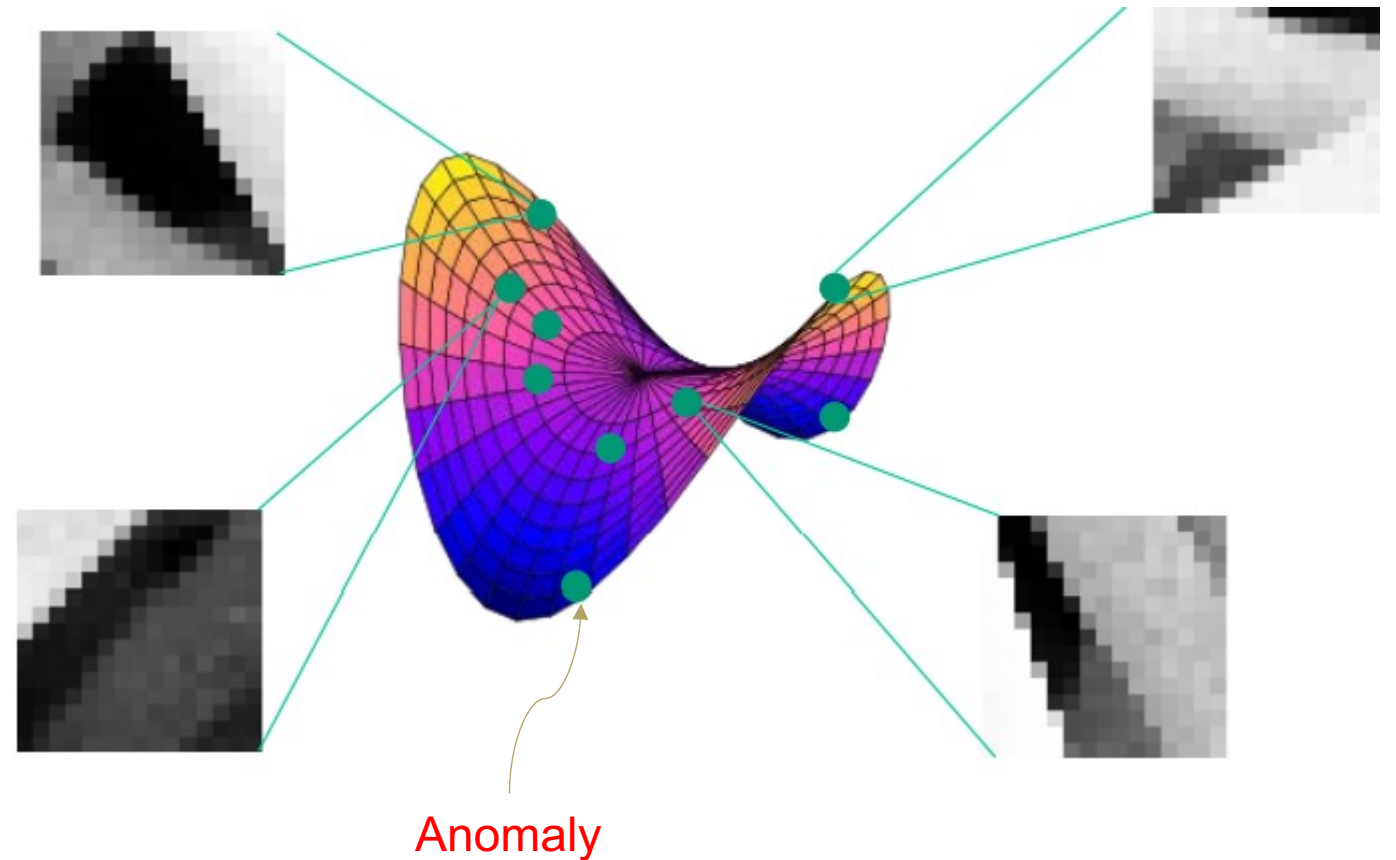
Steps for Anomaly Detection



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



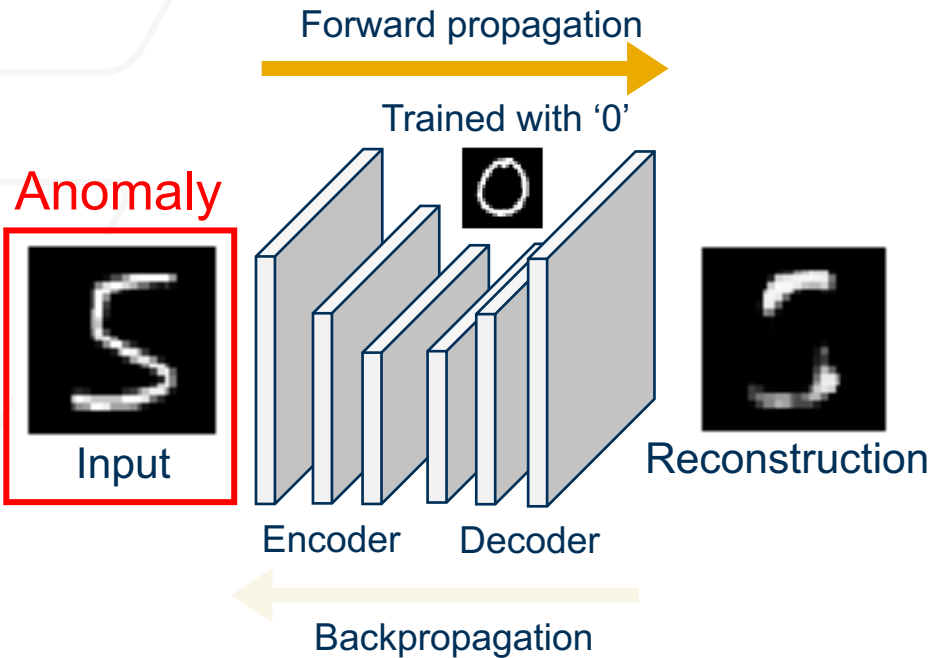
Gradients as Robustness Features

Gradient-based Constraints



Backpropagated Gradient Representations for Anomaly Detection

Activation Constraints



Activation-based representation
(Data perspective)

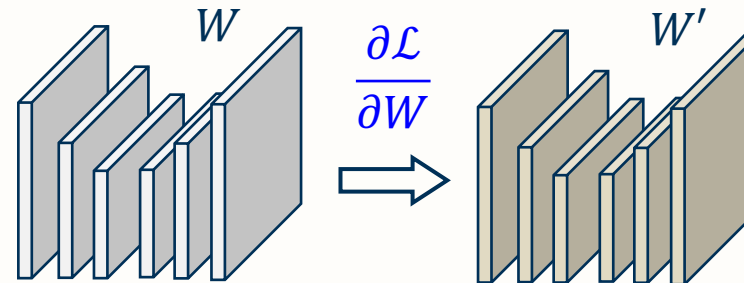
e.g. Reconstruction error (\mathcal{L})



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?

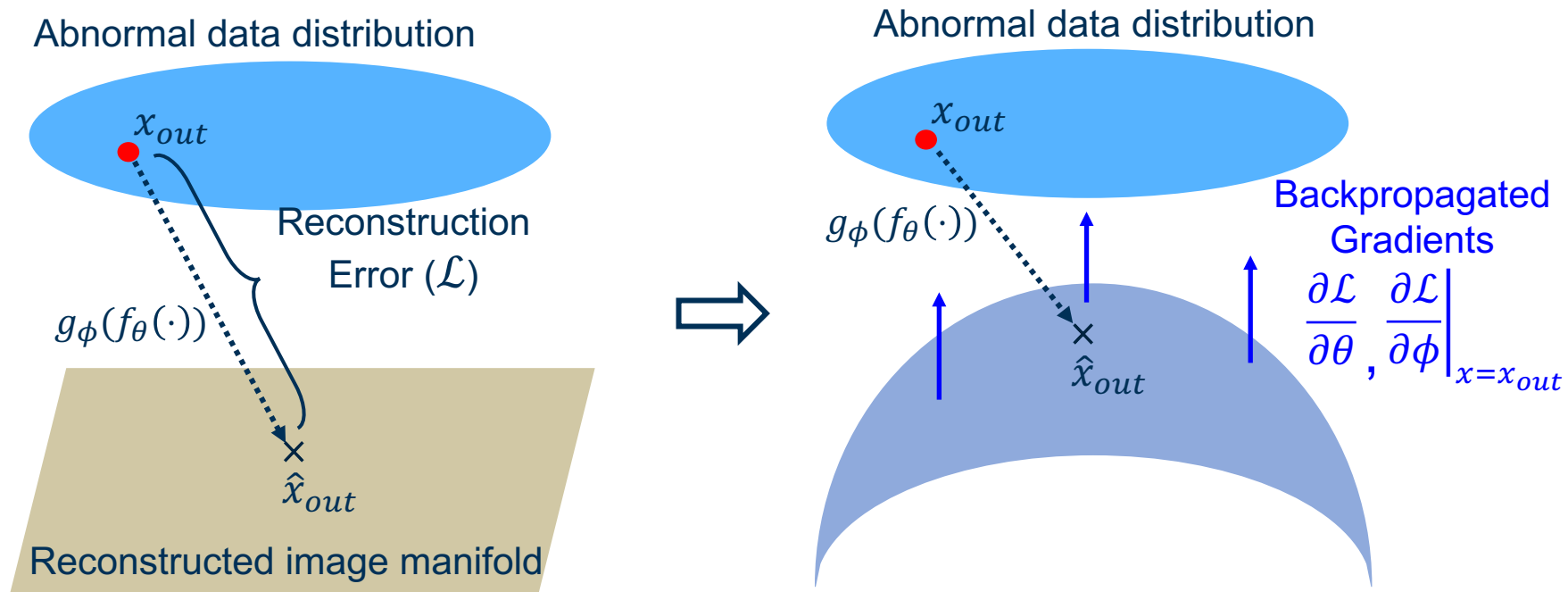
Gradients as Robustness Features

Advantages of Gradient-based Constraints

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



Backpropagated Gradient Representations for Anomaly Detection



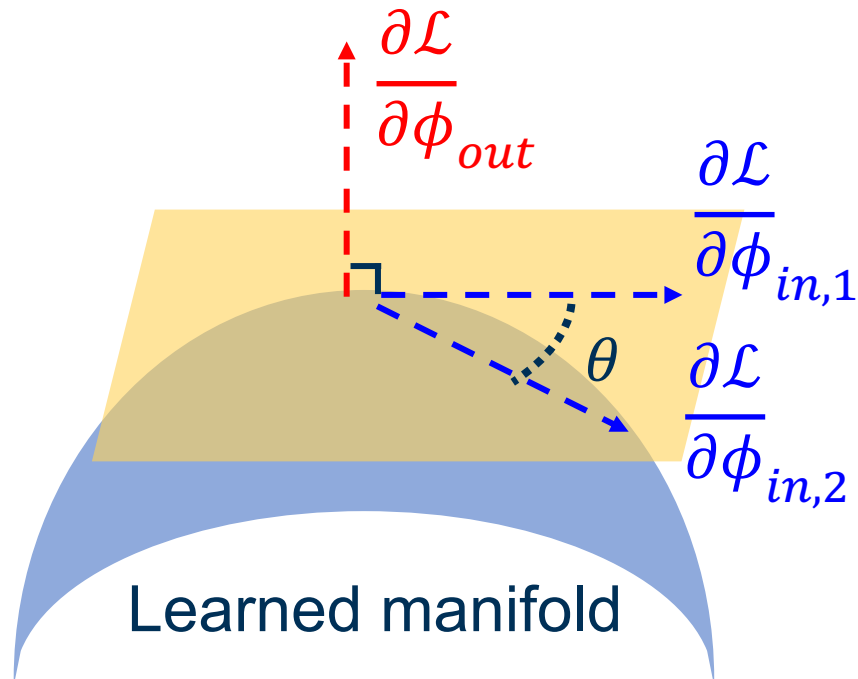
Gradients as Robustness Features

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

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

Gradients as Robustness Features

Activations vs Gradients



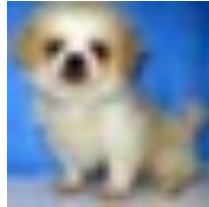
Backpropagated Gradient Representations for Anomaly Detection

Abnormal “class”
detection (CIFAR-10)

e.g.



Normal



Abnormal

AUROC Results

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

Gradients as Robustness Features

Aberrant Condition Detection



Backpropagated Gradient Representations for Anomaly Detection

Abnormal “condition” detection (CURE-TSR)

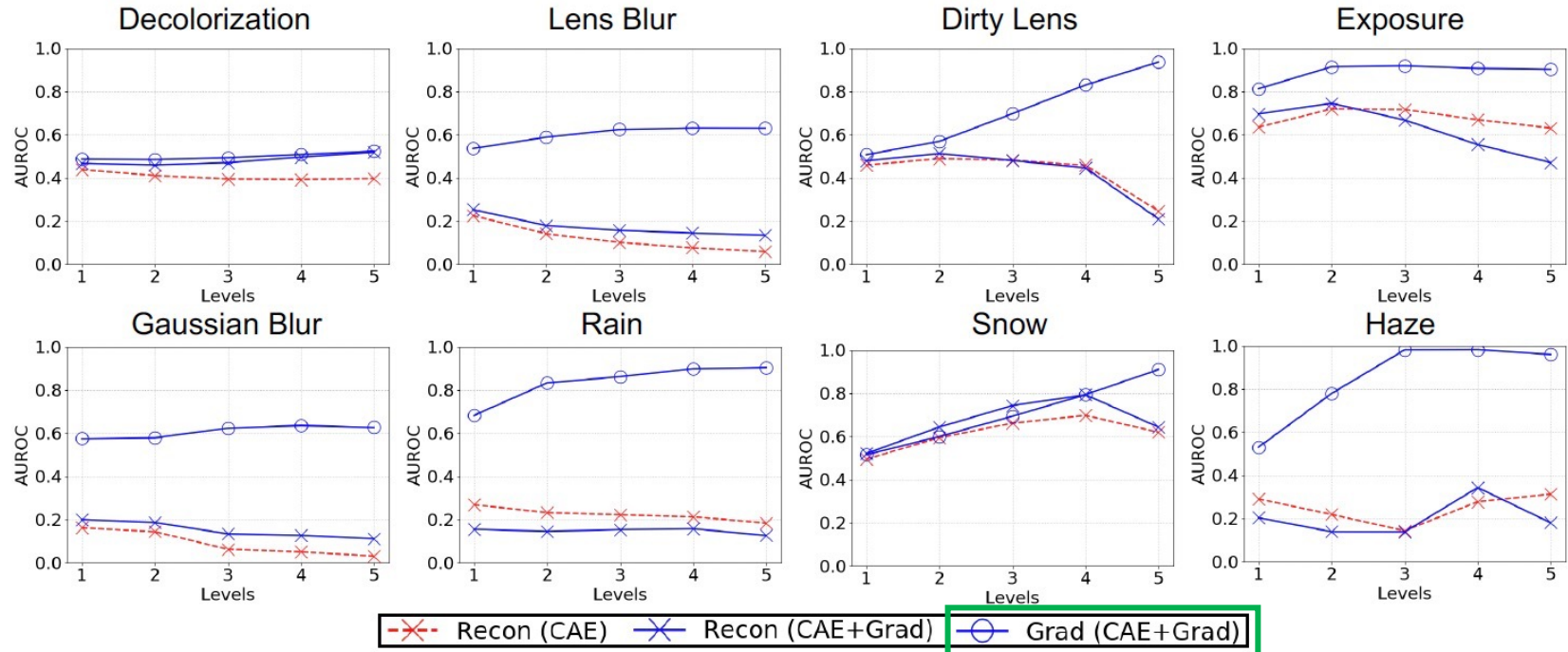


Normal



Abnormal

AUROC Results



Recon: Reconstruction error, Grad: Gradient loss

Outline

Lecture 6: Robustness as Explanatory Proxy

- Robustness and Explanations
- Deep Learning at Inference
 - Robustness under novel data
 - Challenges
 - Gradient Information
- **Gradients as Robustness Features**
 - Anomaly Detection
 - **Out-of-Distribution Detection**
 - **Adversarial Detection**
 - **Corruption Detection**
 - Gradients for Robust Predictions

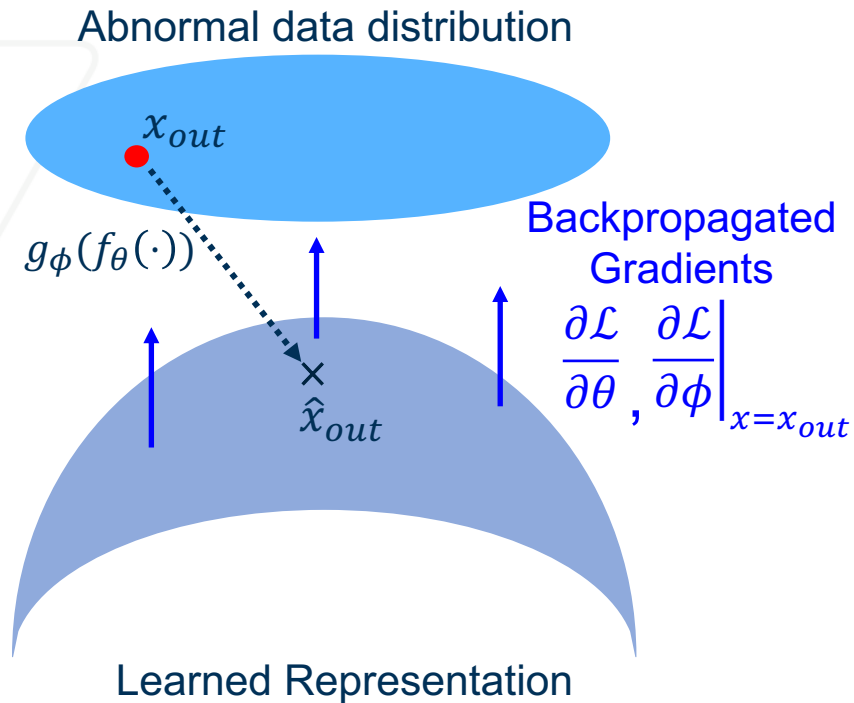
Gradients as Robustness Features

Gradient Intuition



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

Gradients as Robustness Features

Gradient Intuition



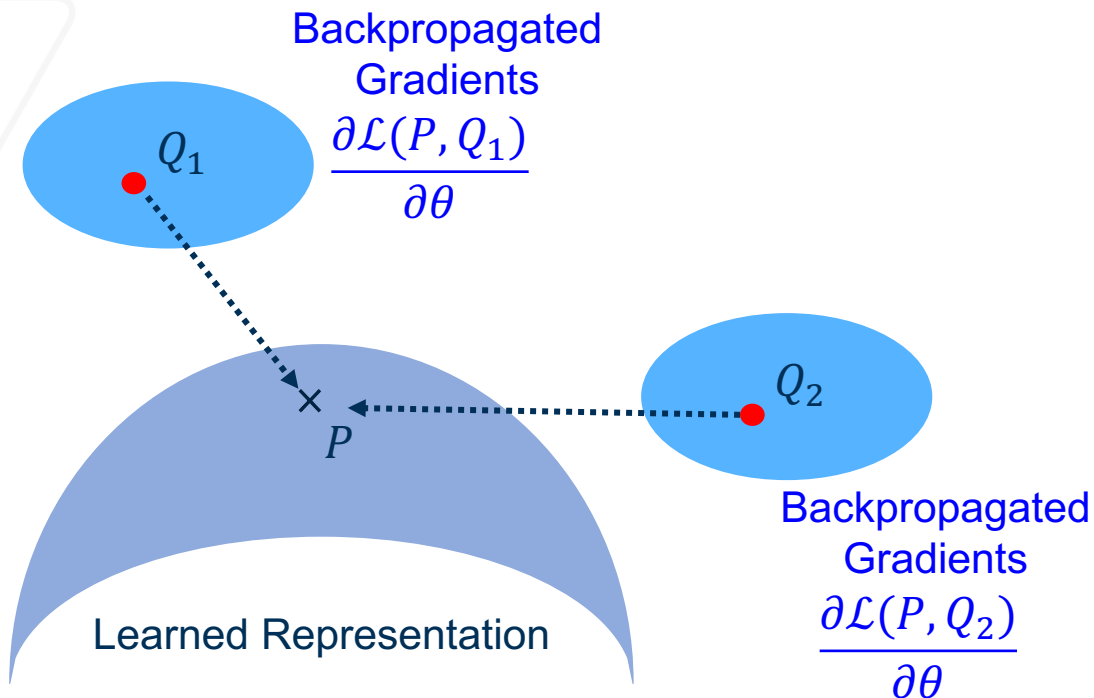
Probing the Purview of Neural Networks via Gradient Analysis

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 vector losses**

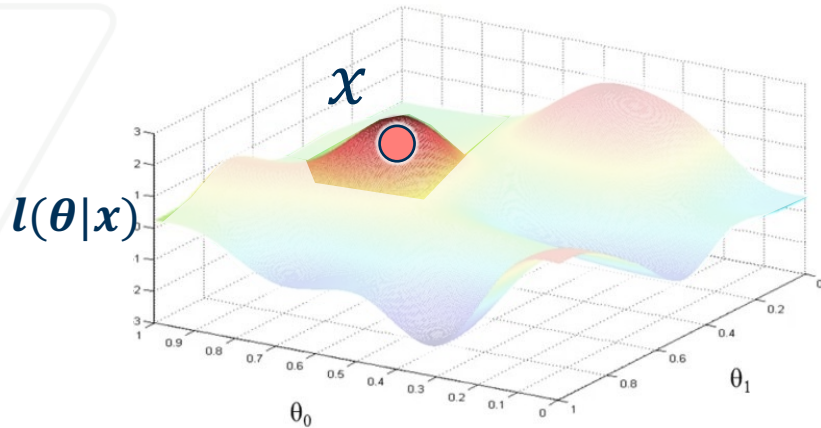
Gradients as Robustness Features

Toy Manifold Example



Probing the Purview of Neural Networks via Gradient Analysis

Gradients represent the local required change in manifold



Contrast class 1



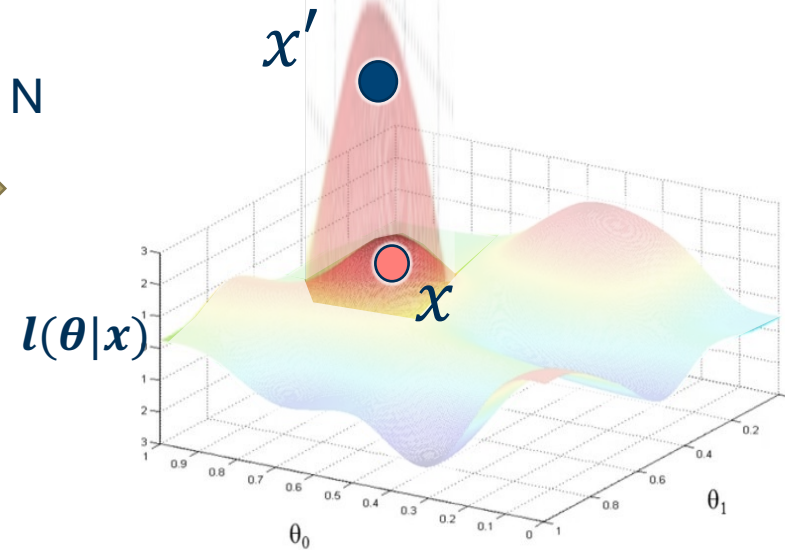
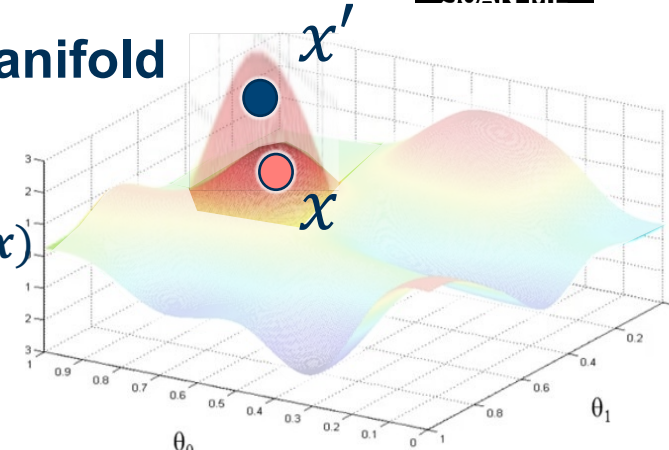
$l(\theta|x)$

⋮

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)!

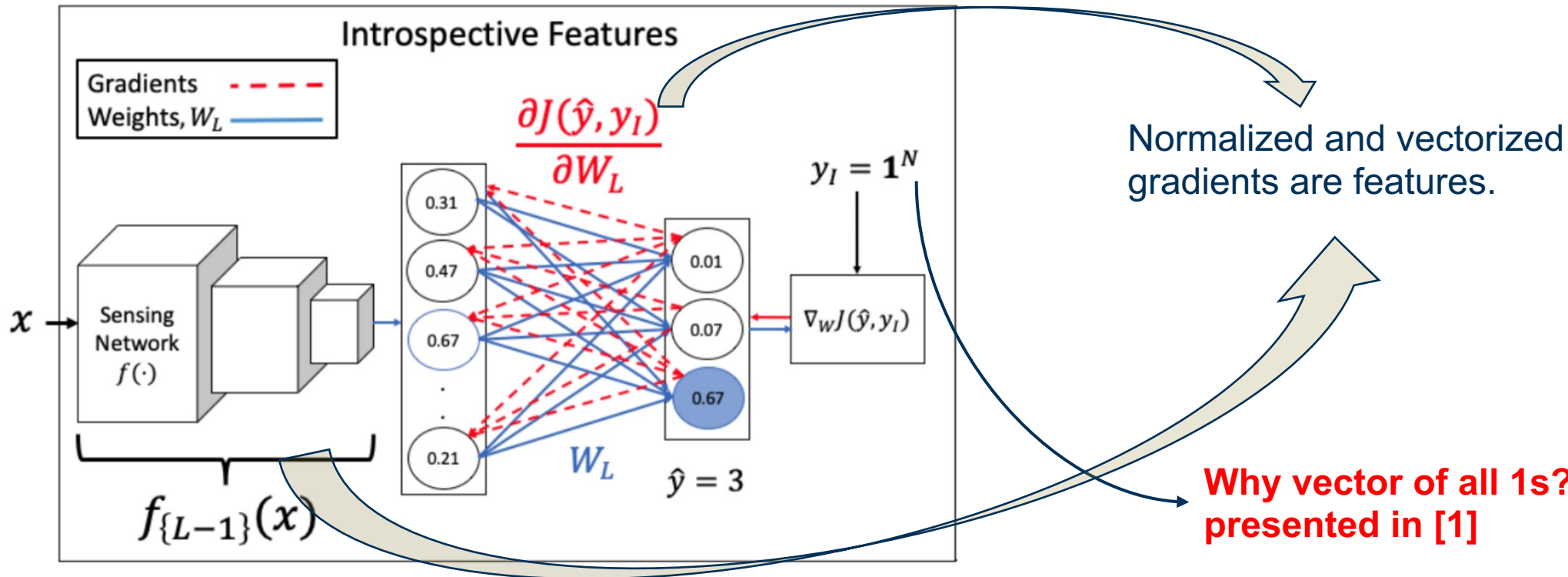
Gradients as Robustness Features

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 gradient features



Why vector of all 1s? The theory is presented in [1]

Gradients as Robustness Features

Utilizing Gradient Features



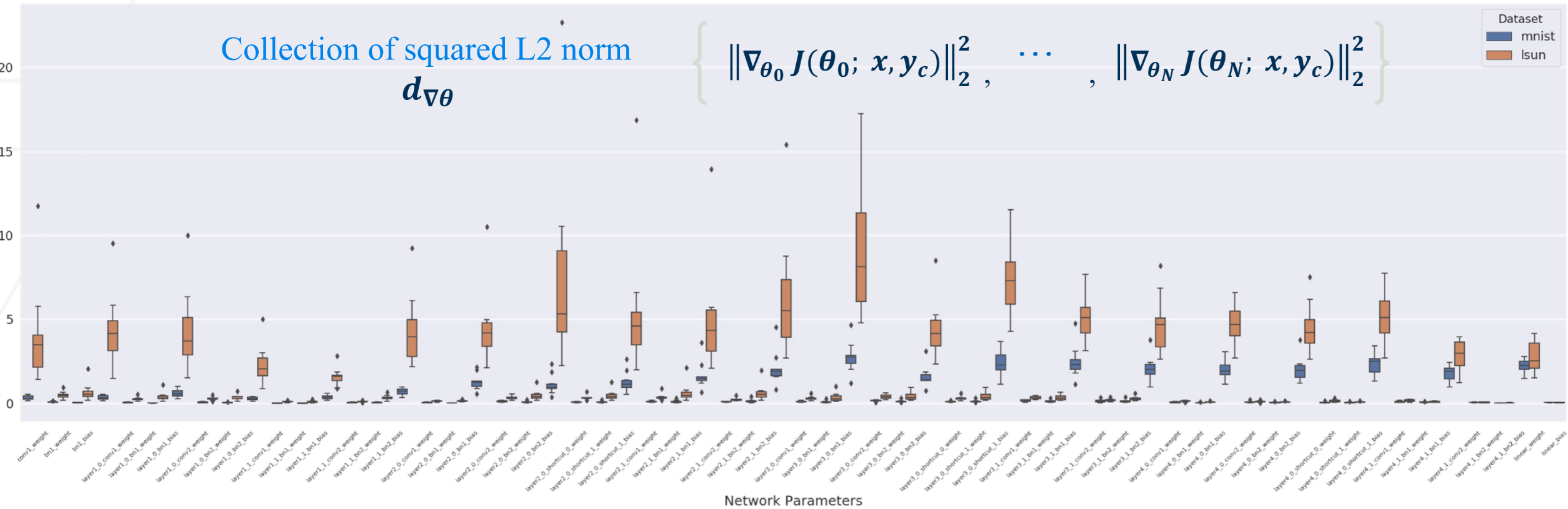
Probing the Purview of Neural Networks via Gradient Analysis

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

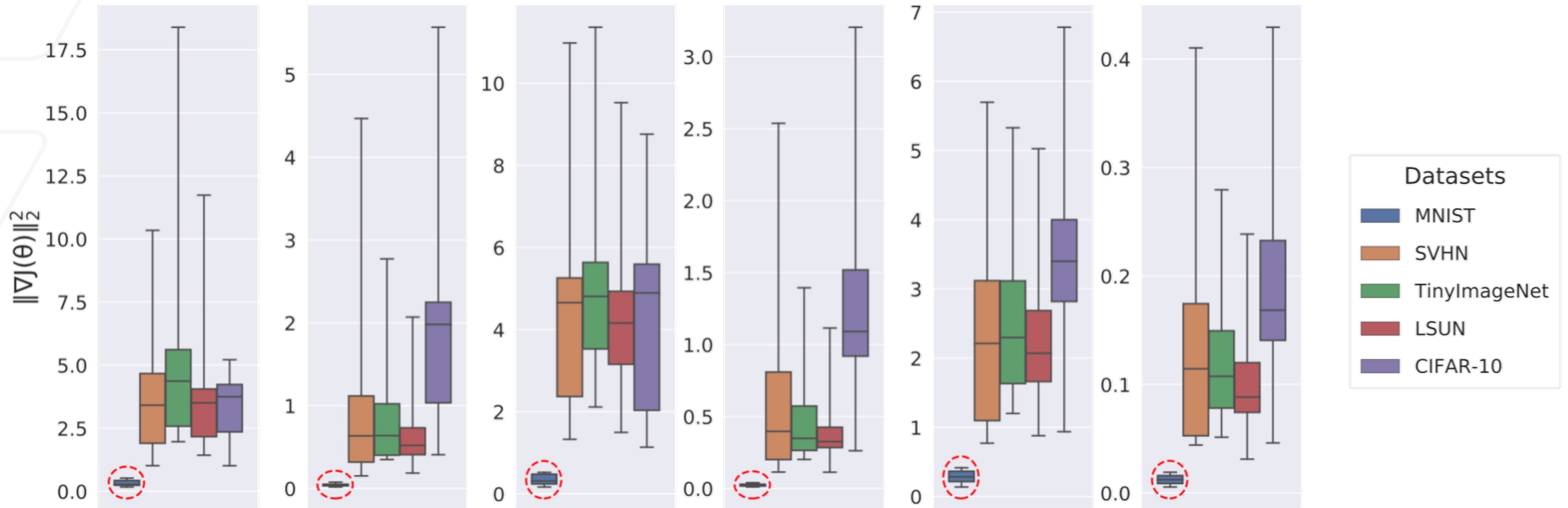
Gradients as Robustness Features

Gradient Features 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

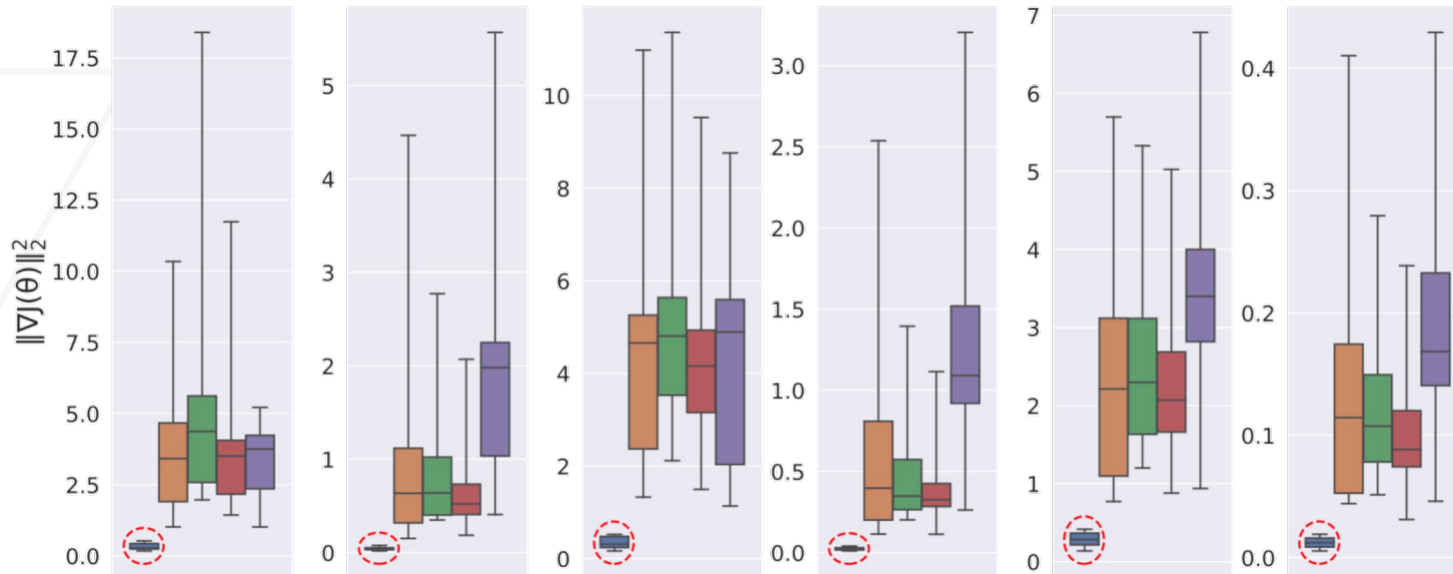
Gradients as Robustness Features

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

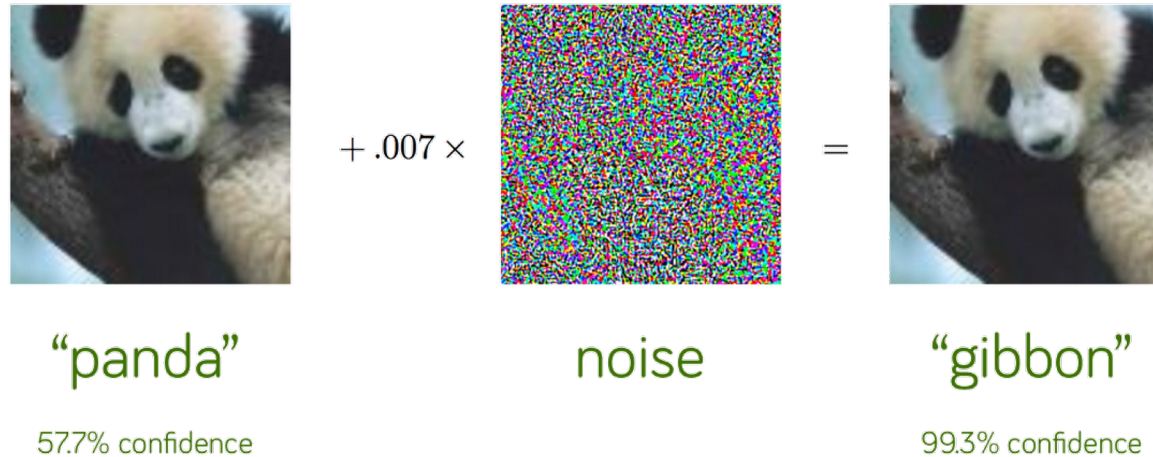
Gradients as Robustness Features

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

Gradients as Robustness Features

Results in Adversarial Setting



Probing the Purview of Neural Networks
via Gradient Analysis

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

Gradients as Robustness Features

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



Gradients as Robustness Features

Results 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



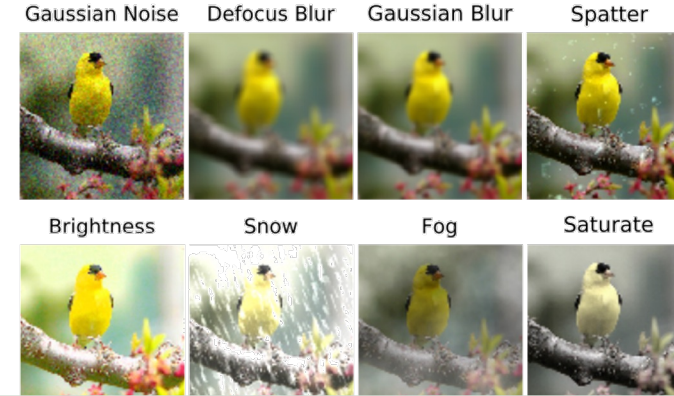
Gradients as Robustness Features

Results 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

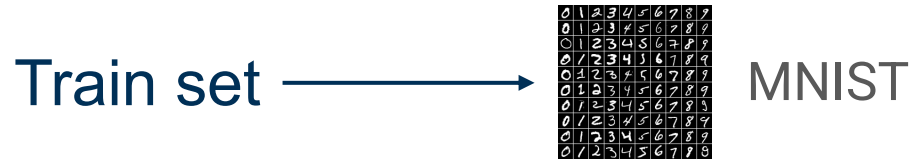


Gradients as Robustness Features

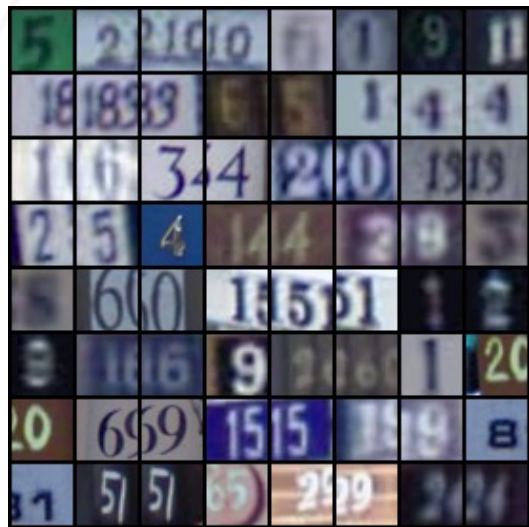
Results in Detecting Challenging Conditions



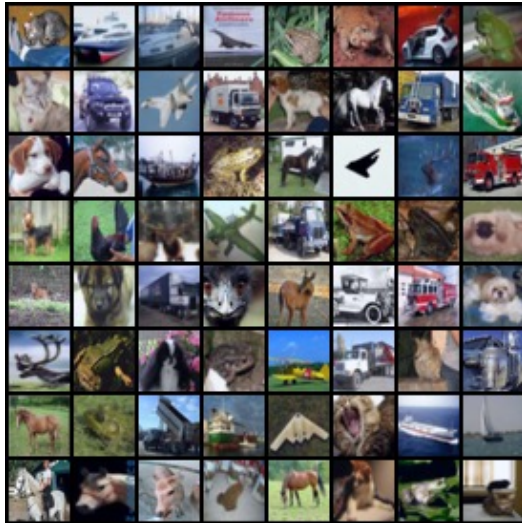
Probing the Purview of Neural Networks via Gradient Analysis



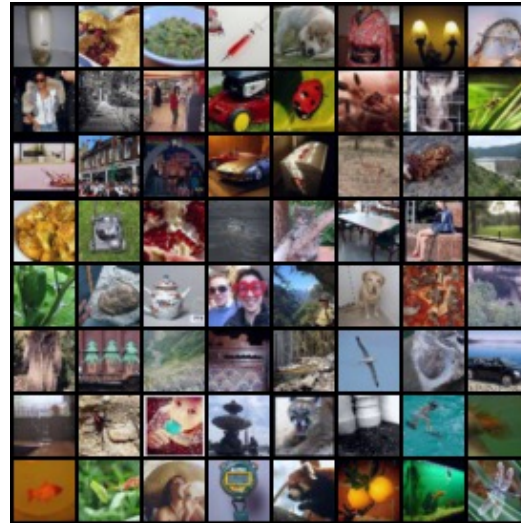
Goal: To detect that these datasets are not part of training



SVHN



CIFAR10



TinyImageNet



LSUN

Gradients as Robustness Features

Results in Detecting Challenging Conditions



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

Gradients as Robustness Features

Results in Detecting Challenging Conditions



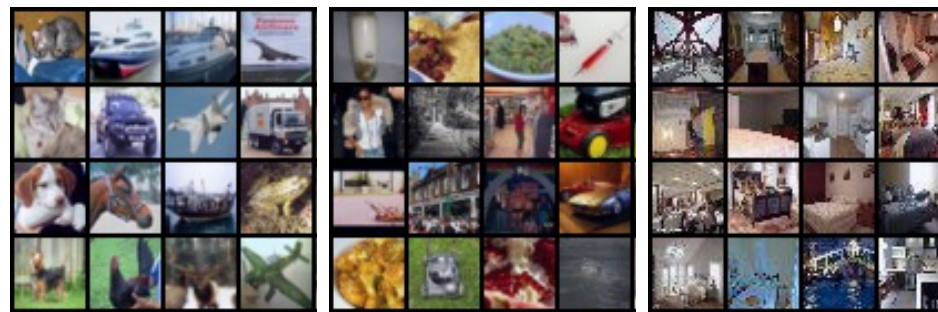
Probing the Purview of Neural Networks via Gradient Analysis

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

Numbers



SVHN



CIFAR10

TinyImageNet

LSUN

Objects, natural scenes

Gradients as Robustness Features

Results in Detecting Challenging Conditions



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



Outline

Lecture 6: Robustness as Explanatory Proxy

- Robustness and Explanations
- Deep Learning at Inference
 - Robustness under novel data
 - Challenges
 - Gradient Information
- **Gradients as Robustness Features**
 - Anomaly Detection
 - Out-of-Distribution Detection
 - Adversarial Detection
 - Corruption Detection
 - **Gradients for Robust Predictions**

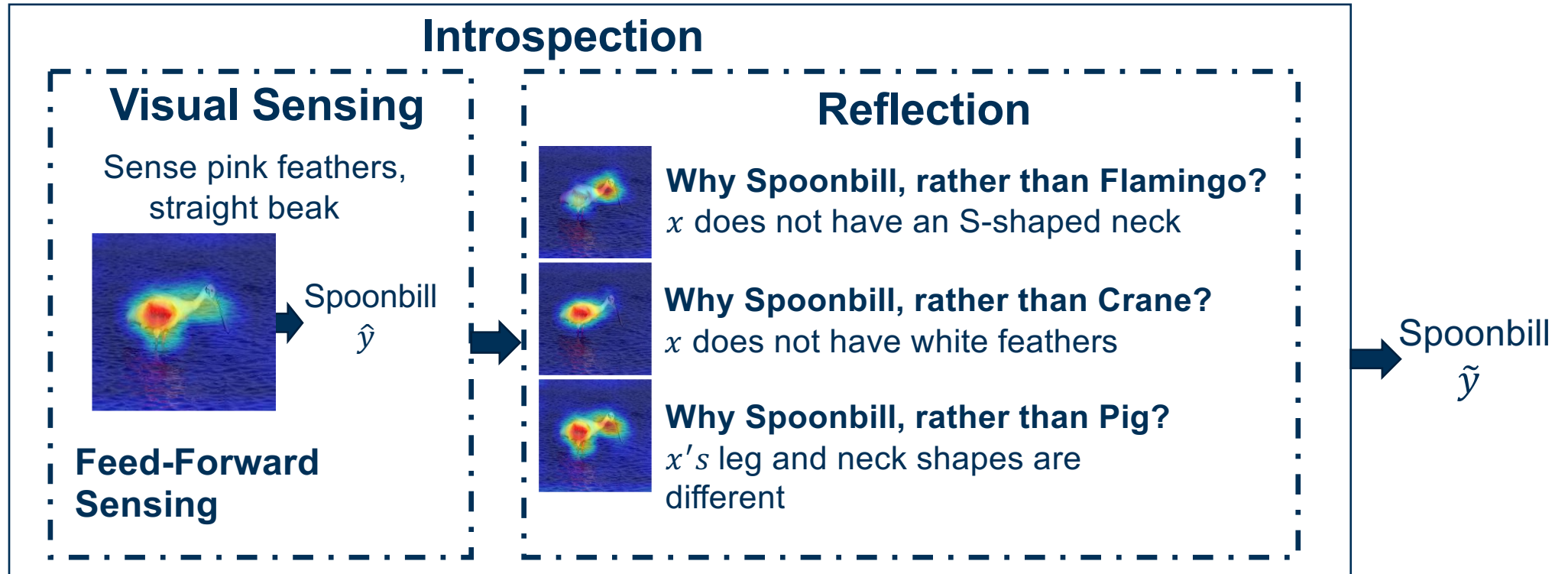
Gradients as Robustness Features

Introspective Learning



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



Gradients as Robustness Features

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?

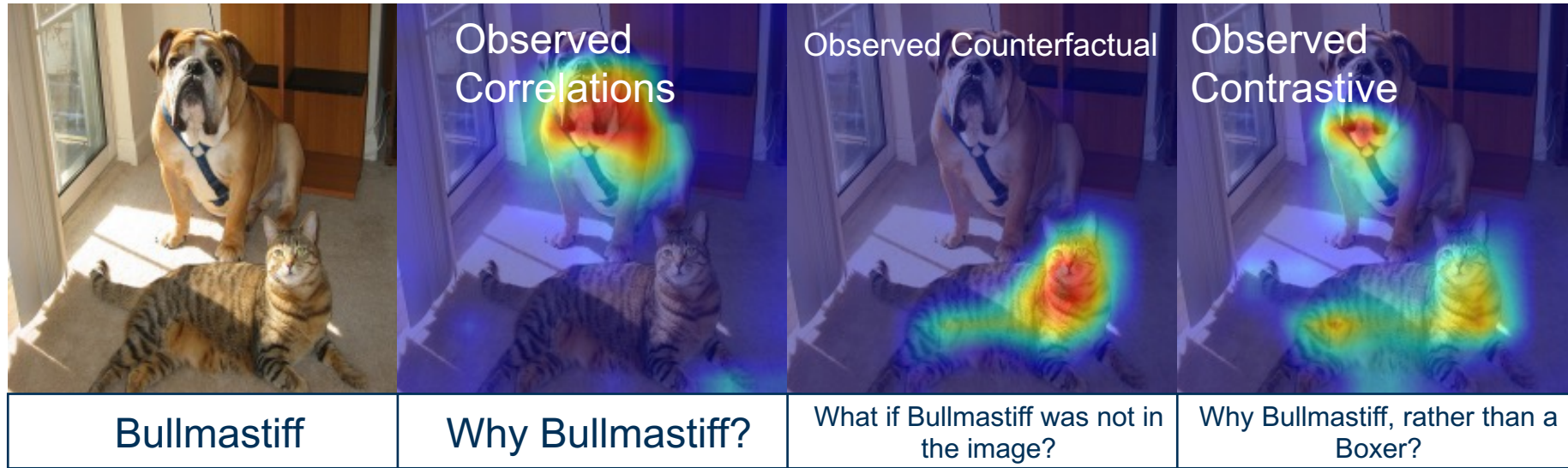
Gradients as Robustness Features

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?

Gradients as Robustness Features

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 .*

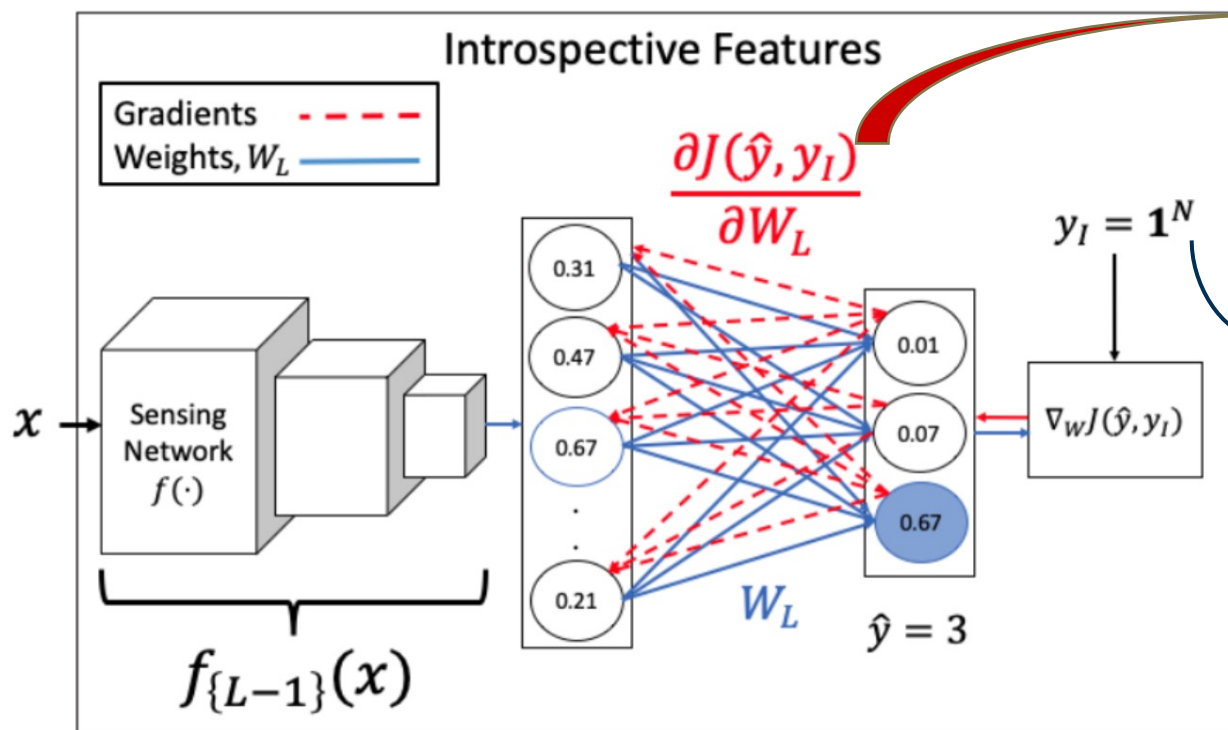
Gradients as Robustness Features

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



Normalized and vectorized gradients are introspective features

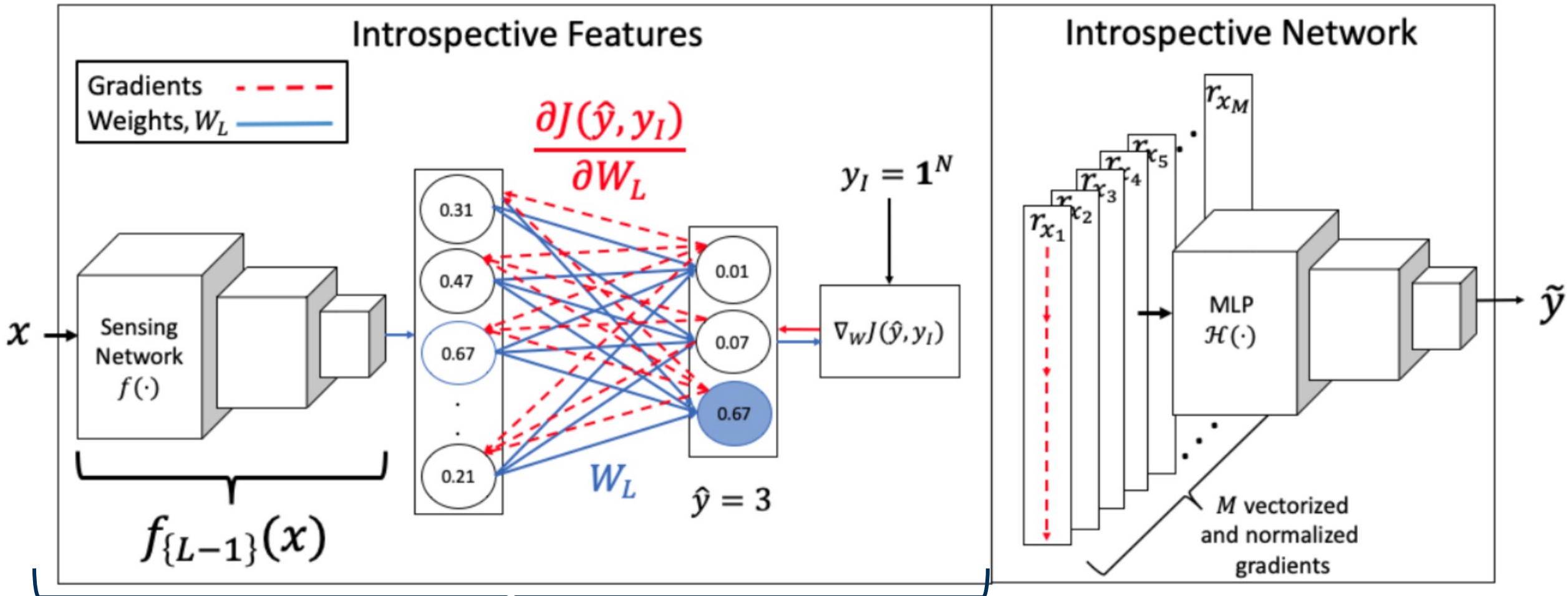
Vector of all ones: A confounding label!

Gradients as Robustness Features

Utilizing Gradient Features



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



Introspective Features

Gradients as Robustness Features

When is Introspection Useful?



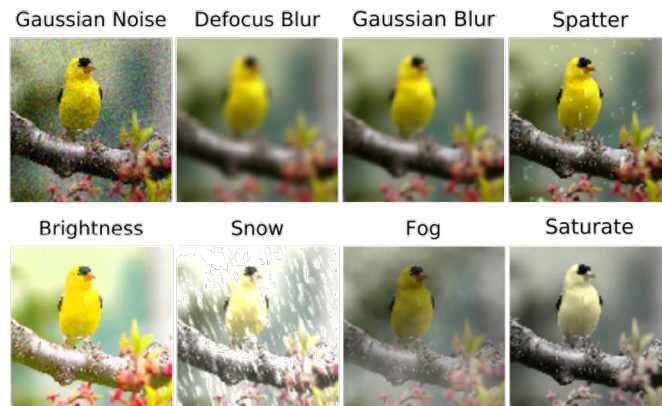
Introspective Learning: A Two-stage Approach for Inference in Neural Networks

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



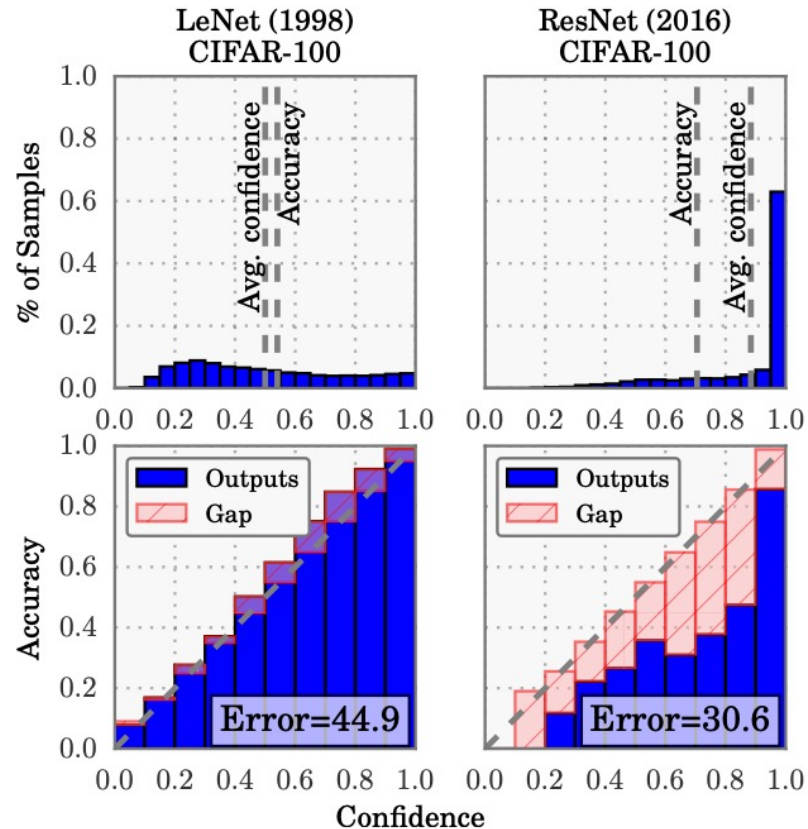
Calibration

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

Gradients as Robustness Features

Generalization and Calibration results

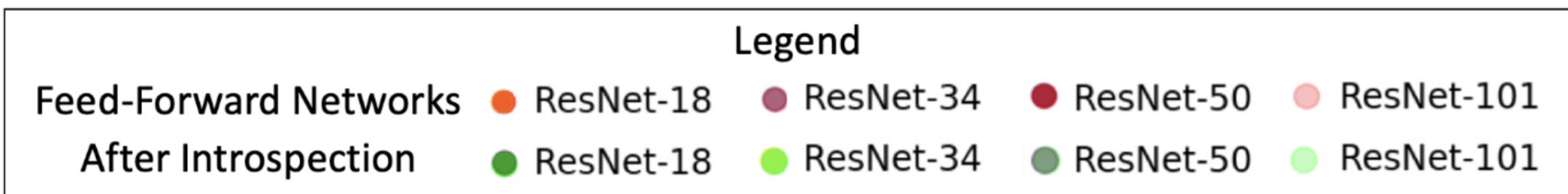
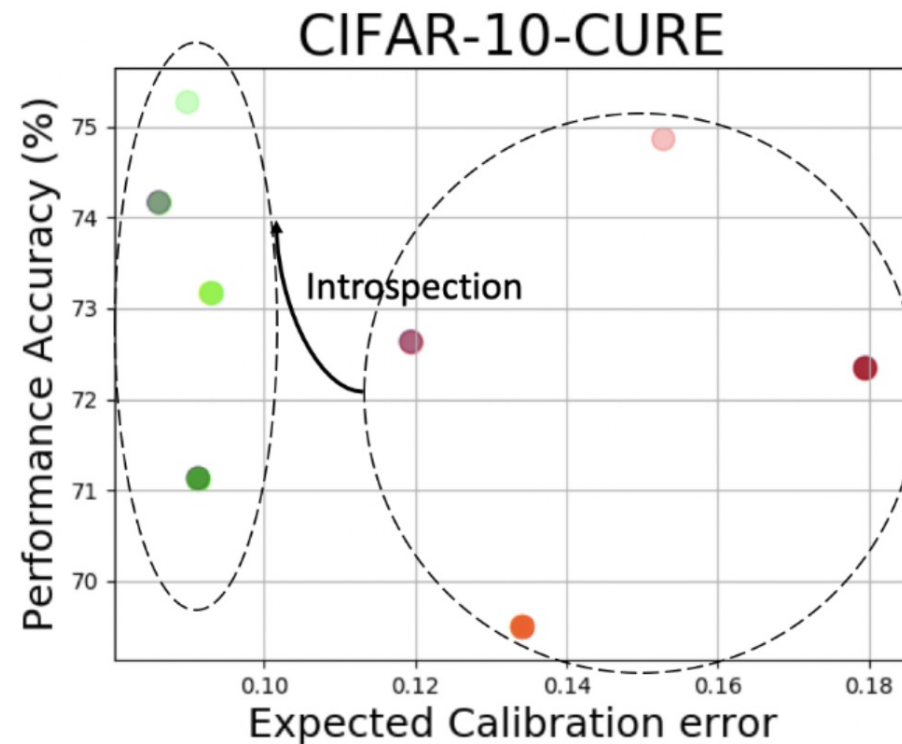
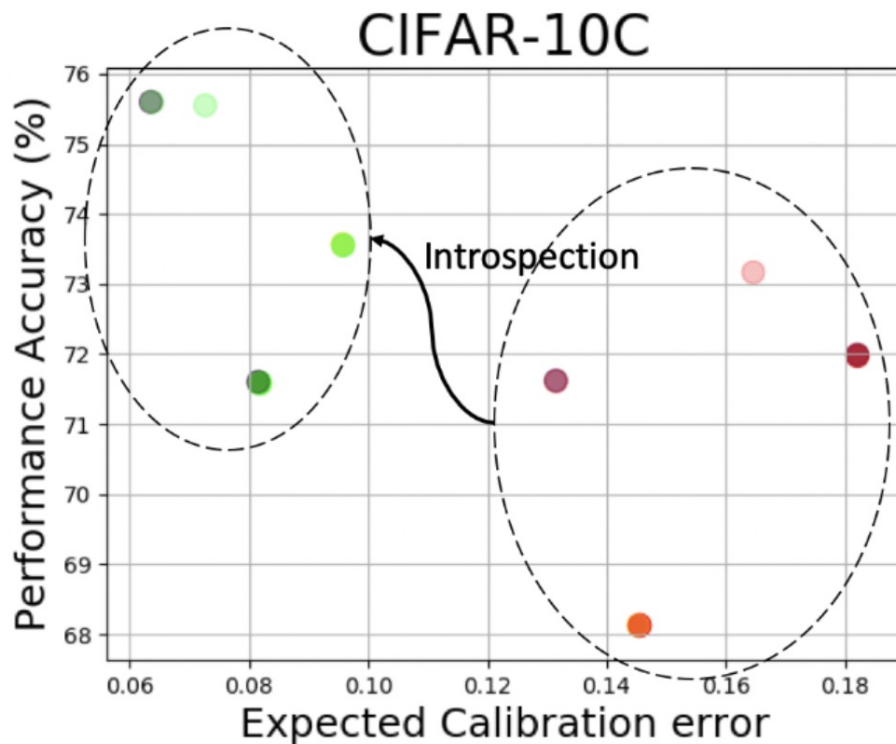


Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Ideal: Top-left corner

Y-Axis: Generalization

X-Axis: Calibration



Gradients as Robustness Features

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 (49)	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!

Gradients as Robustness Features

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

Takeaways

Takeaways from Lecture 6

- **Robust networks are empirically shown to provide better explanations**
- An **indirect validation** is to show that **gradient features** (that are extensively used for creating explanations) maybe manipulated to obtain **robust results**
- Similar loss-based **gradient operations** that lead to explanations also **lead to robust predictions**
- Gradients as features can be used to obtain
 - Anomaly detection
 - Out-of-distribution, novelty, adversarial detection
 - Robust prediction

References

Lecture 6: Robustness as Explanatory Proxy

- Prabhushankar, Mohit, and Ghassan AlRegib. "Contrastive reasoning in neural networks." *arXiv preprint arXiv:2103.12329* (2021).
- 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.
- V. Chandola, A. Banerjee, V. Kumar. "Anomaly detection: A survey". *ACM Comput. Surv.* 41, 3, Article 15 (July 2009), 58 pages
- G. Kwon, M. Prabhushankar, D. Temel, and G. AlRegib, "Backpropagated Gradient Representations for Anomaly Detection," 2020
- Lee, Jinsol, et al. "Probing the Purview of Neural Networks via Gradient Analysis." *IEEE Access* 11 (2023): 32716-32732.