

# Robust Neural Networks

## Part 3: Uncertainty at Inference

# Objective

## Objective of the Tutorial

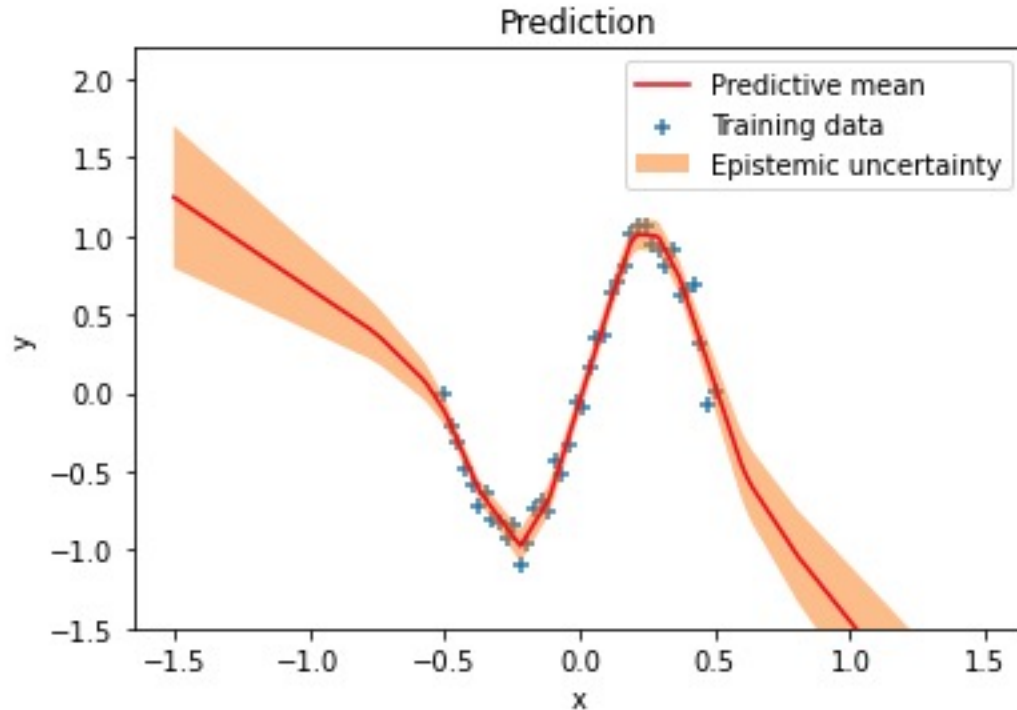
**To discuss methodologies that promote robustness in neural networks at inference**

- 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**



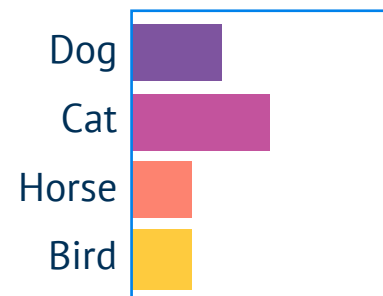
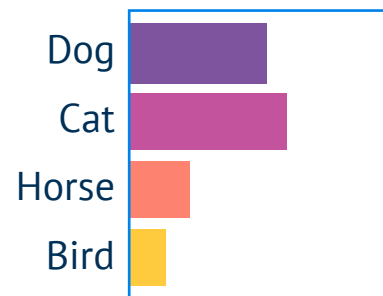
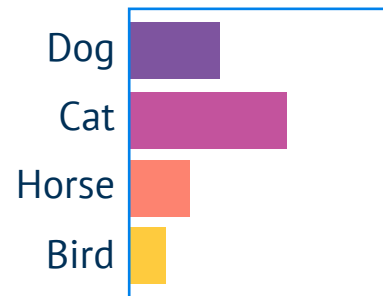
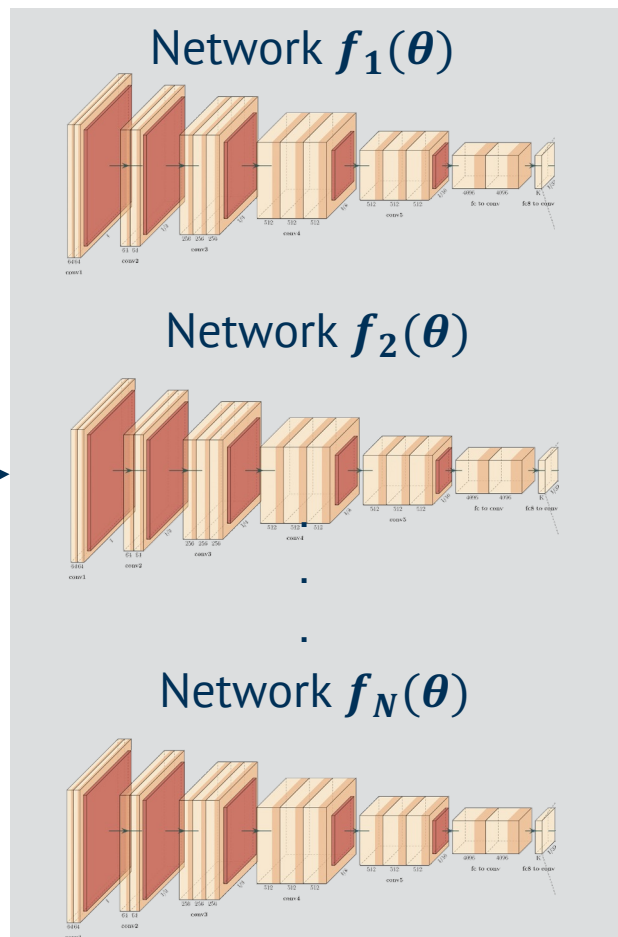
A simple example:

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

# Uncertainty

## Uncertainty Quantification in Neural Networks

### Via Ensembles<sup>1</sup>

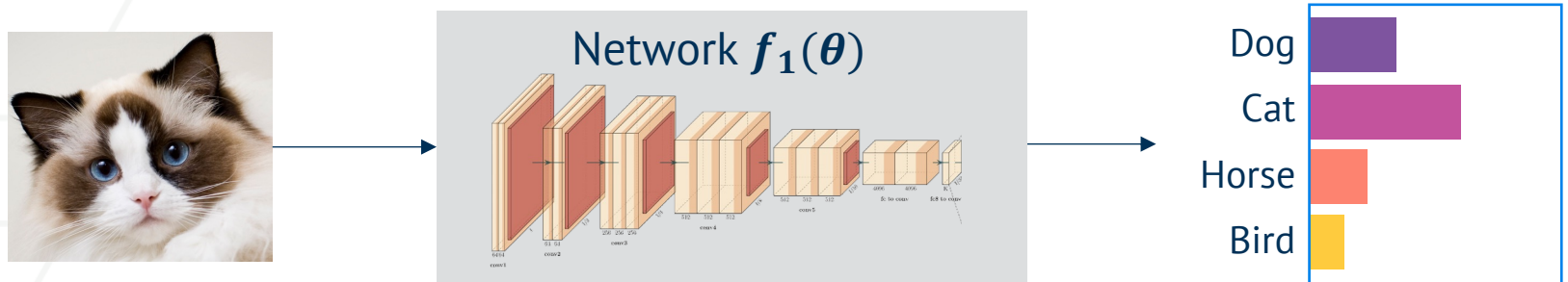


Variation within outputs  $Var(y)$  is the uncertainty. Commonly referred to as **Prediction Uncertainty**.

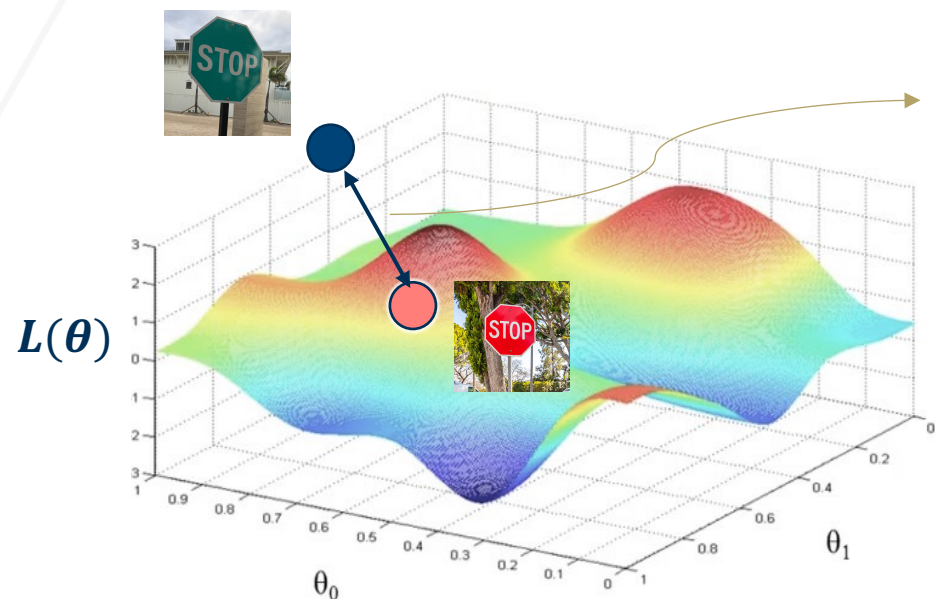
# Uncertainty

## Uncertainty Quantification in Neural Networks

### Via Single pass methods<sup>1</sup>



Uncertainty quantification using a single network and a single pass



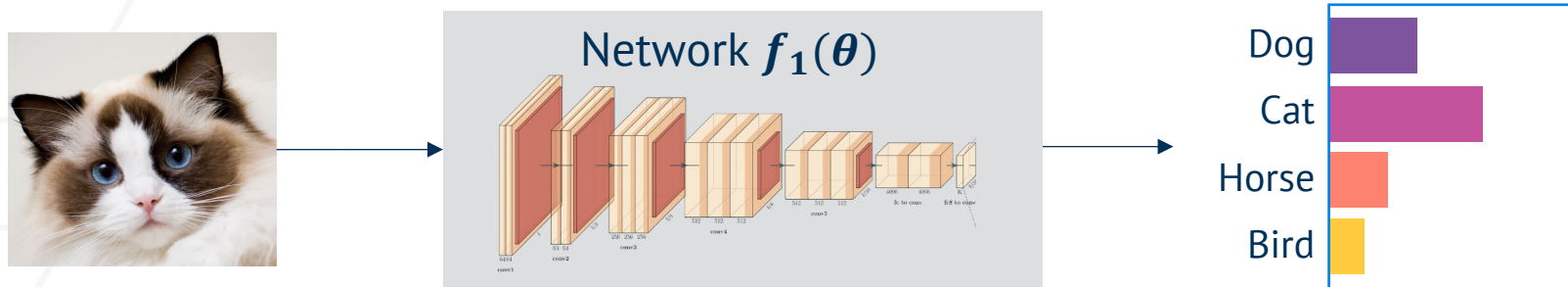
Calculate distance from some trained clusters

**Does not require multiple networks!**

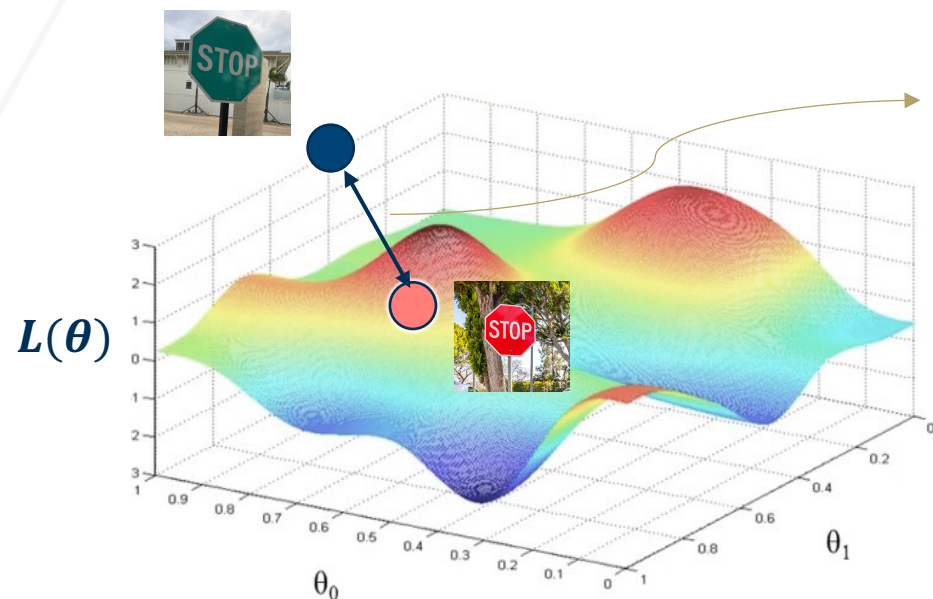
# Uncertainty

## 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

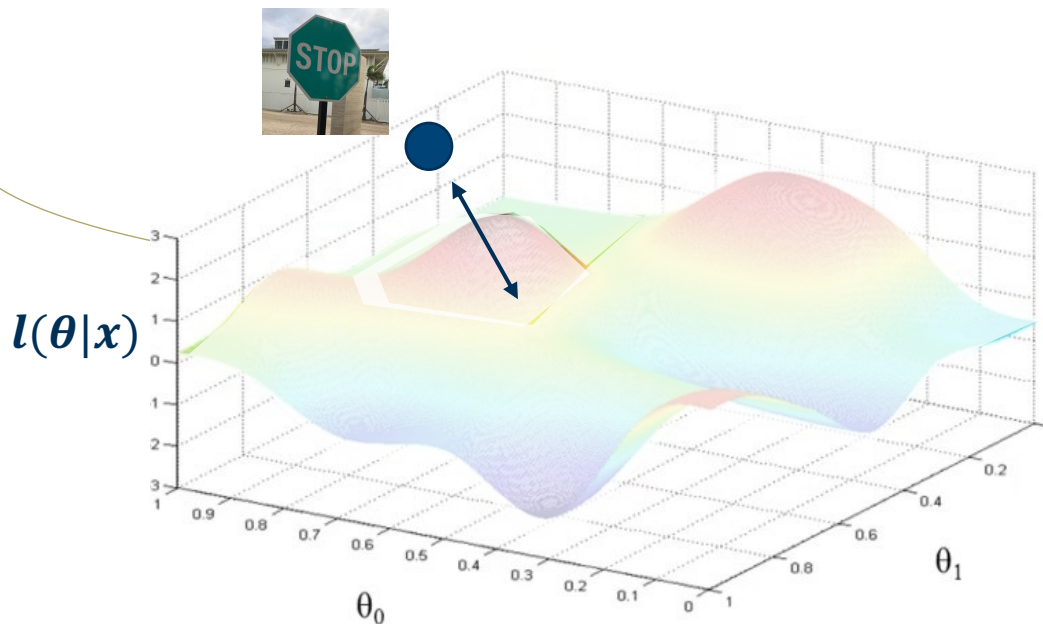
## Gradients as Single pass Features

**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



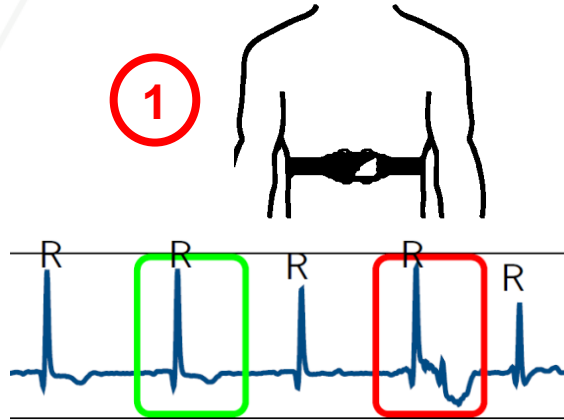


# Anomalies

## Finding Rare Events in Normal Patterns



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

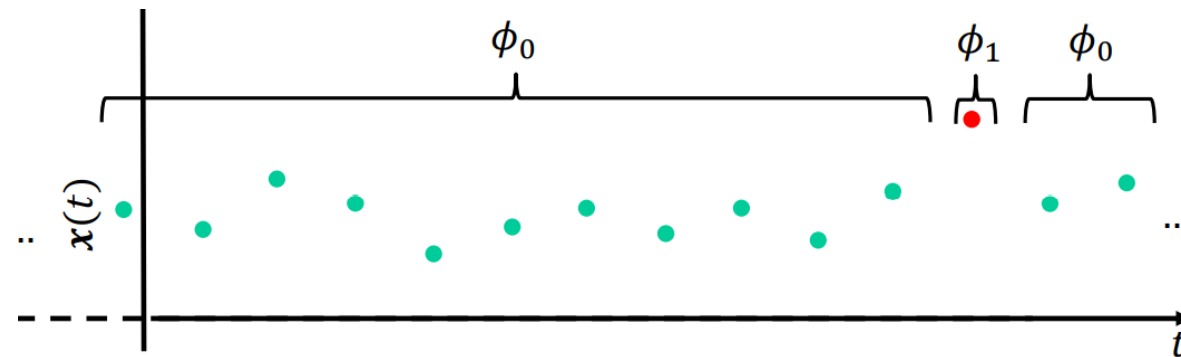


Statistical Definition:

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

Goal: Detect  $\phi_1$

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



# Anomalies

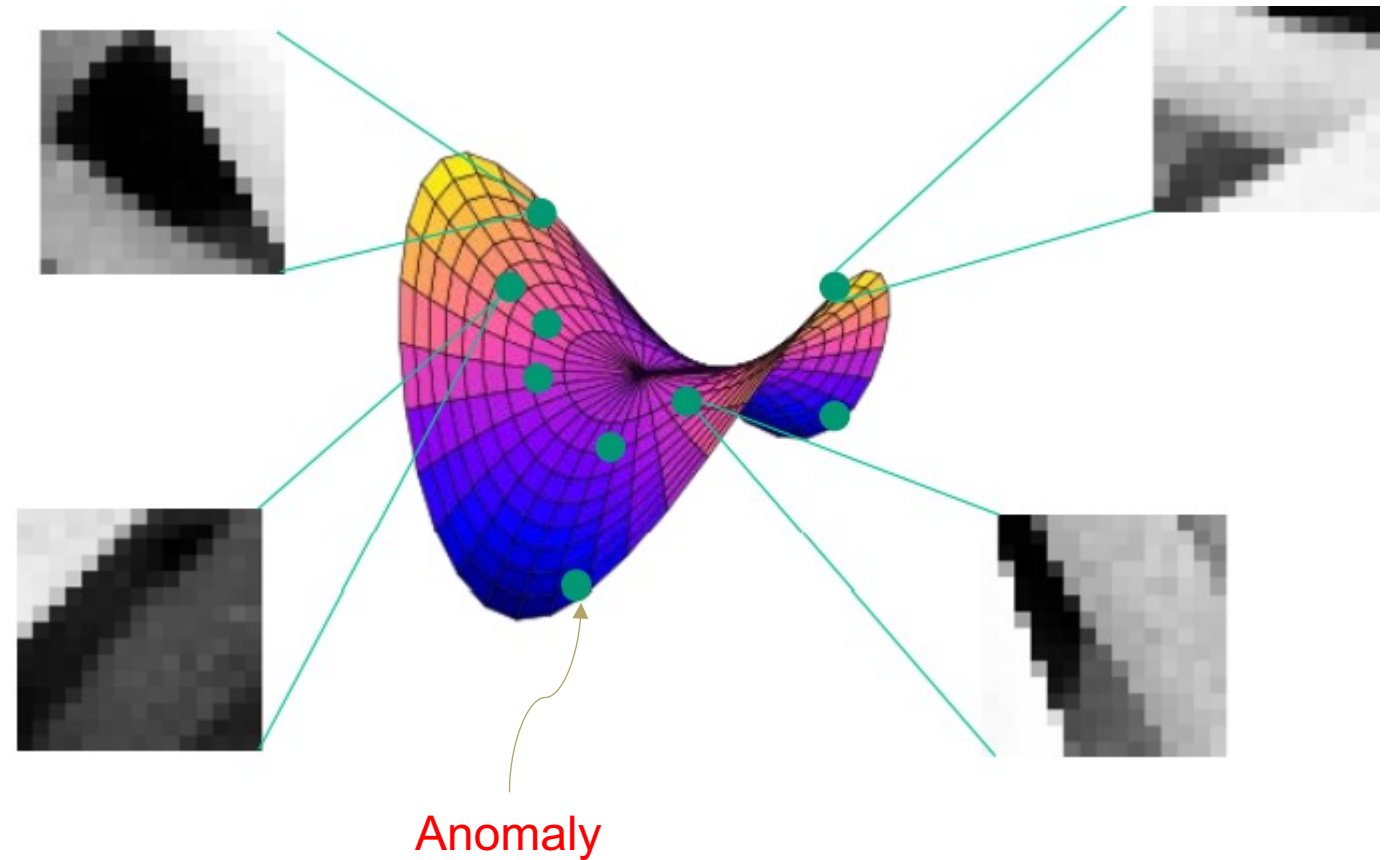
## Steps for Anomaly Detection



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



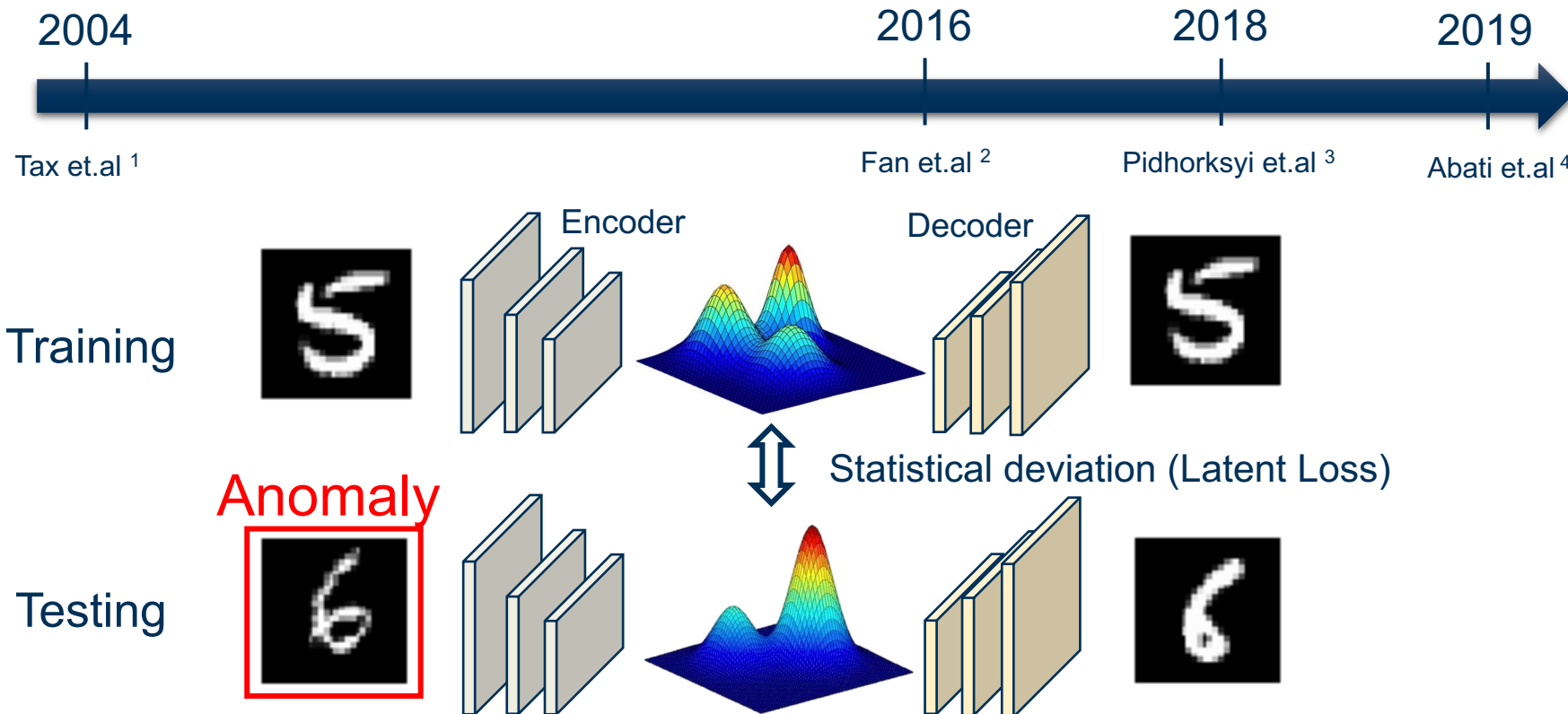
# Constraining Manifolds

## 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. Pidhorksyi, R. Almhosen, 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)

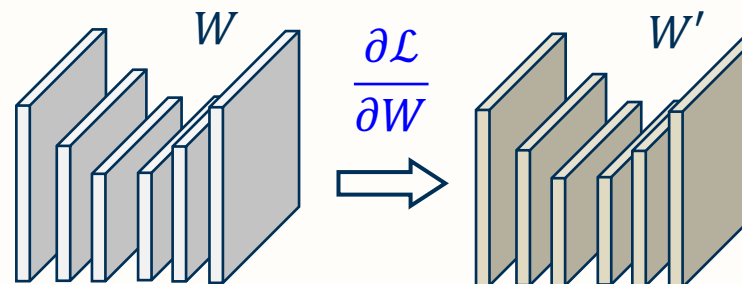
e.g. Reconstruction error ( $\mathcal{L}$ )



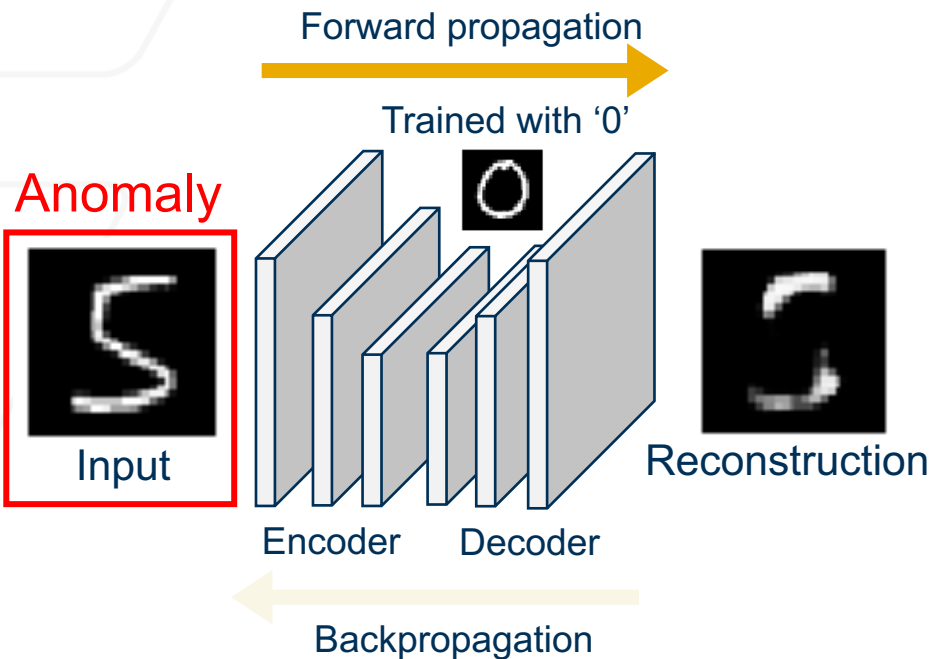
How much of the **input** does not correspond to the learned information?

### Gradient Constraints

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



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

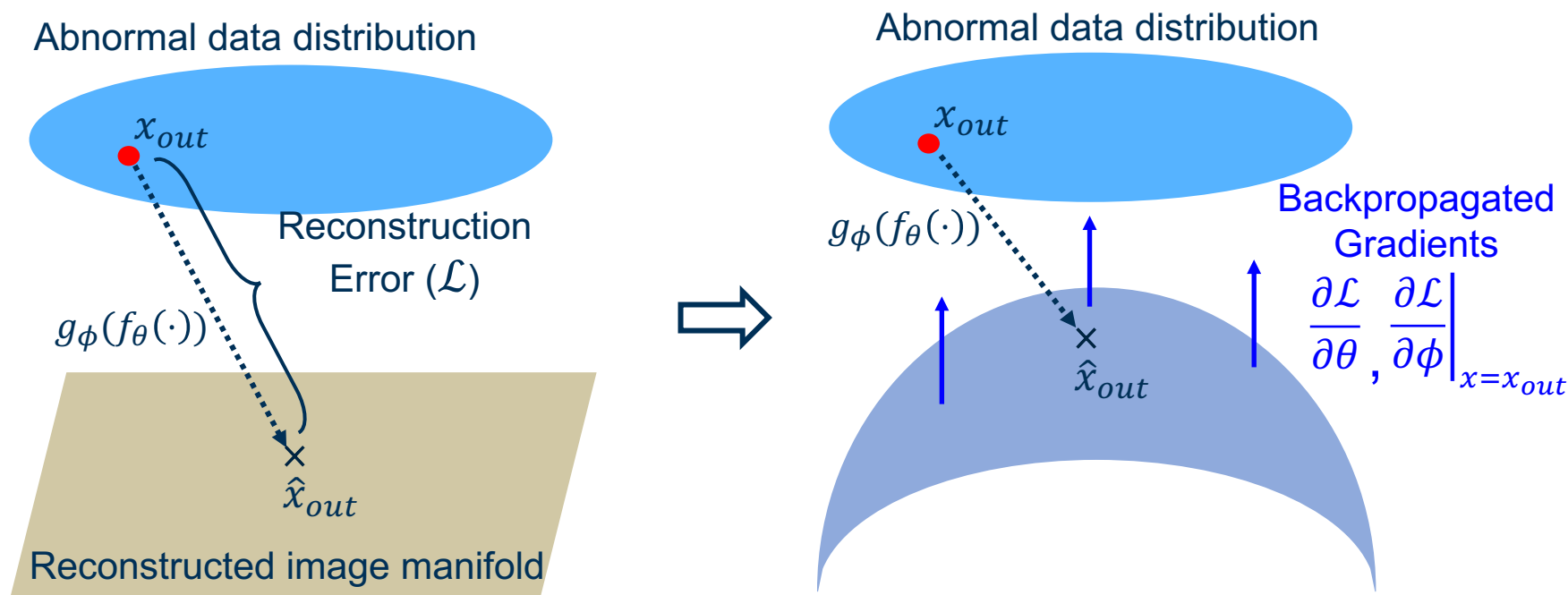


# Constraining Manifolds

## Advantages of Gradient-based Constraints



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



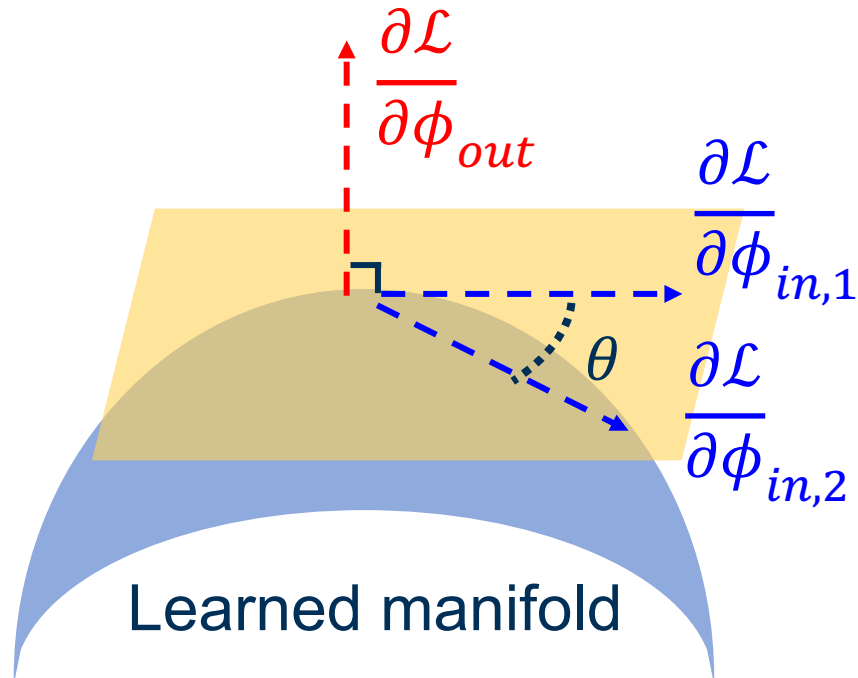
# GradCON: Gradient Constraint

## 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

$\phi$ : 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	<b>0.613</b>	0.495	0.498	0.711	0.390	0.564
CAE	Recon	0.659	0.356	<b>0.640</b>	0.555	0.695	0.554	0.549	0.478	0.695	0.357	0.554
+ Grad	Grad	<b>0.752</b>	0.619	0.622	0.580	0.705	0.591	0.683	<b>0.576</b>	<b>0.774</b>	<b>0.709</b>	<b>0.661</b>
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	<b>0.640</b>	0.497	<b>0.743</b>	0.515	<b>0.745</b>	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	<b>0.625</b>	0.591	<b>0.596</b>	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

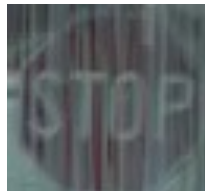


Backpropagated Gradient Representations for Anomaly Detection

Abnormal “condition”  
detection (CURE-TSR)

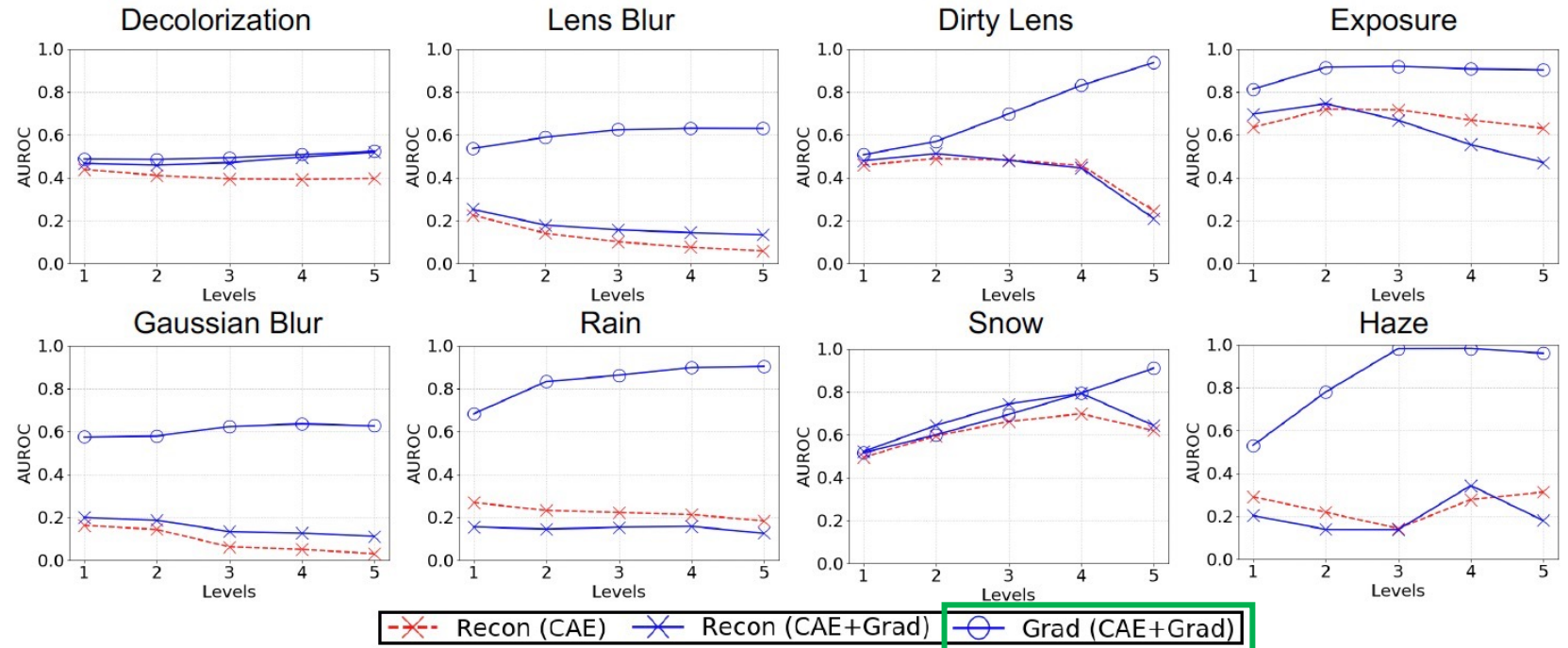


Normal



Abnormal

### AUROC Results



Recon: Reconstruction error, Grad: Gradient loss



# Uncertainty

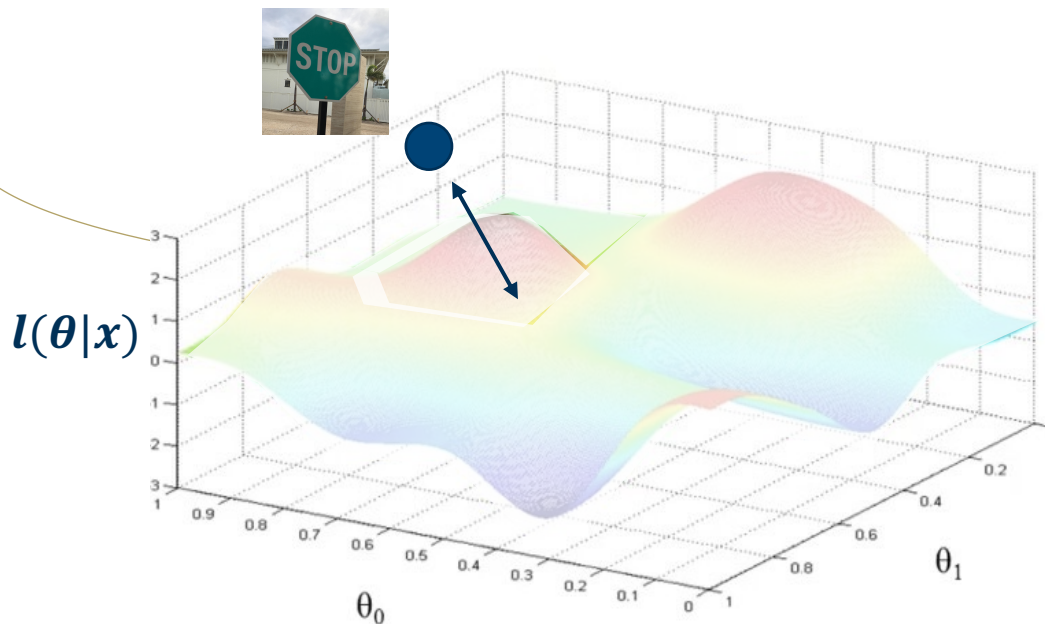
## Gradients as Single pass Features

**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



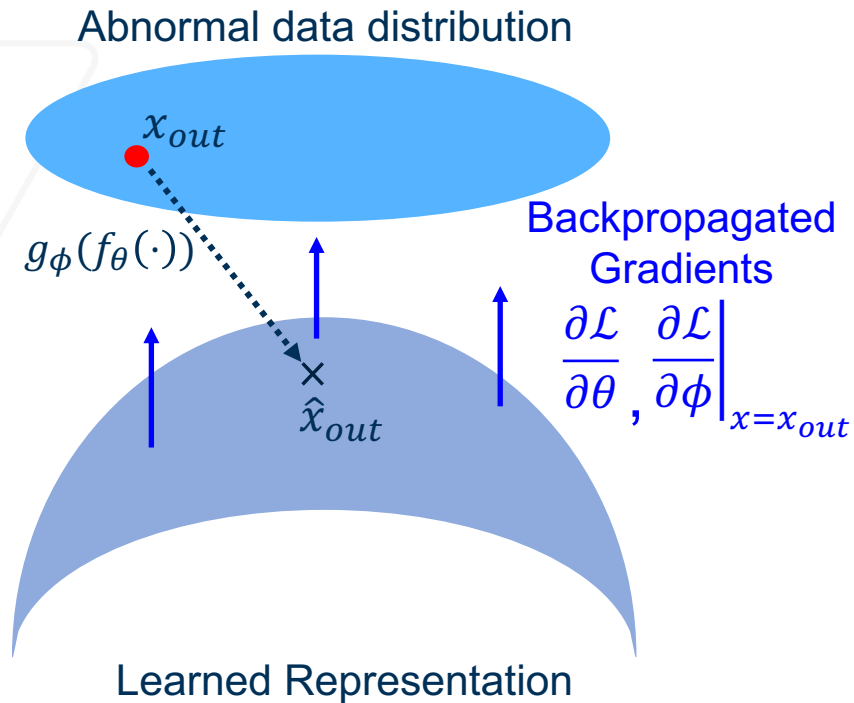
# Uncertainty in Neural Networks

## 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



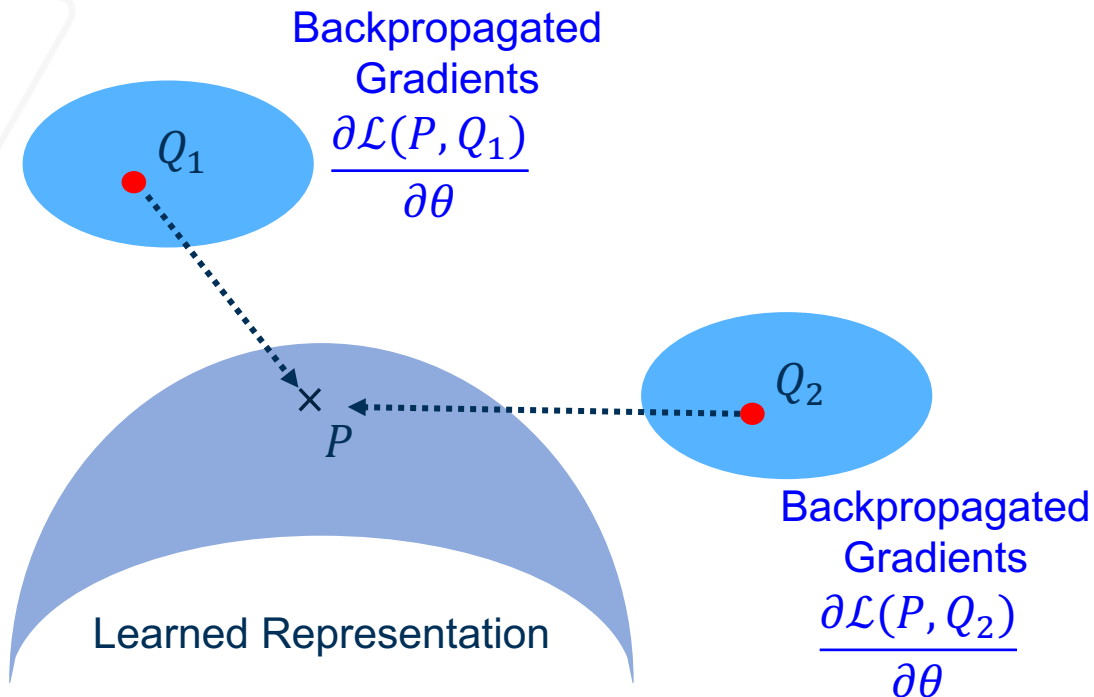
Probing the Purview of Neural Networks via Gradient Analysis

**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

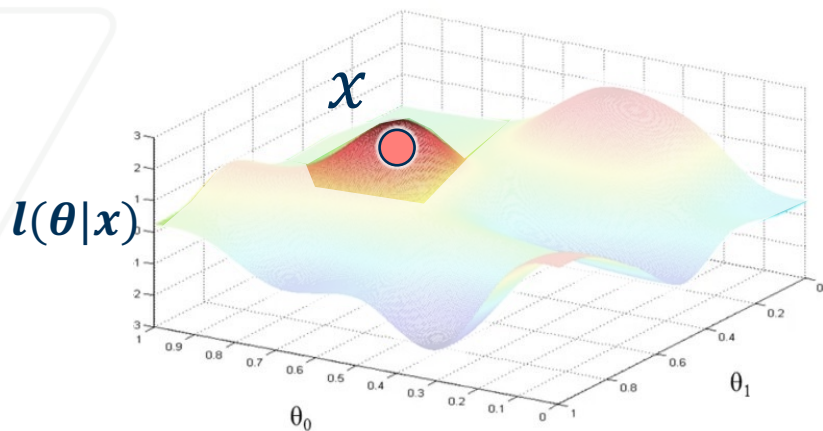
What is uncertainty?



Probing the Purview of Neural Networks via Gradient Analysis

Gradients represent the local required change in manifold

Similar to introspective learning!



Contrast class 1



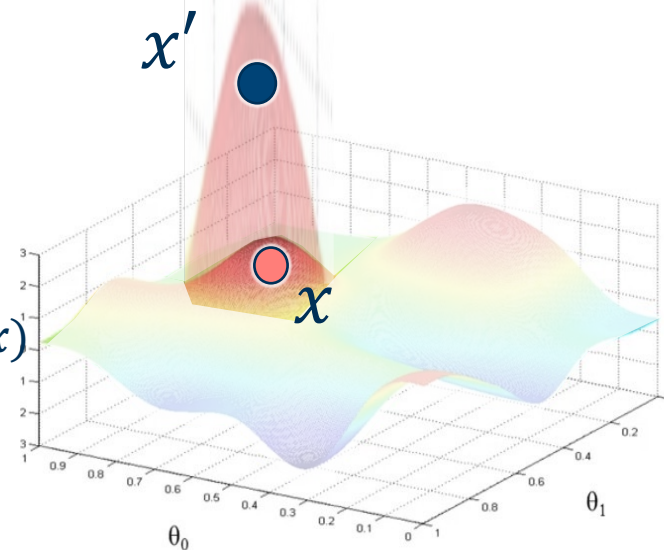
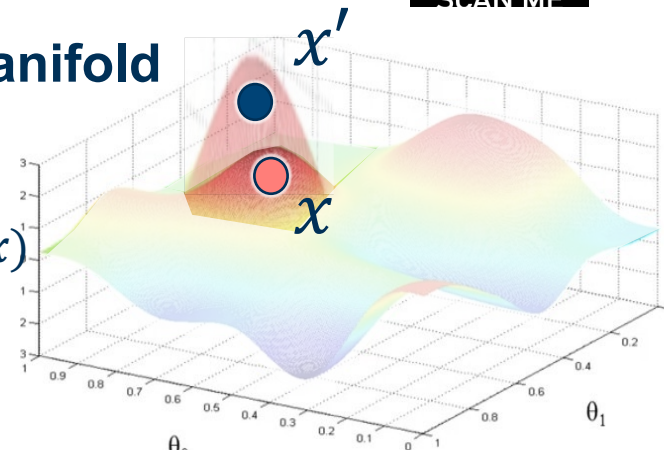
$l(\theta|x)$

·  
·  
·

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)!

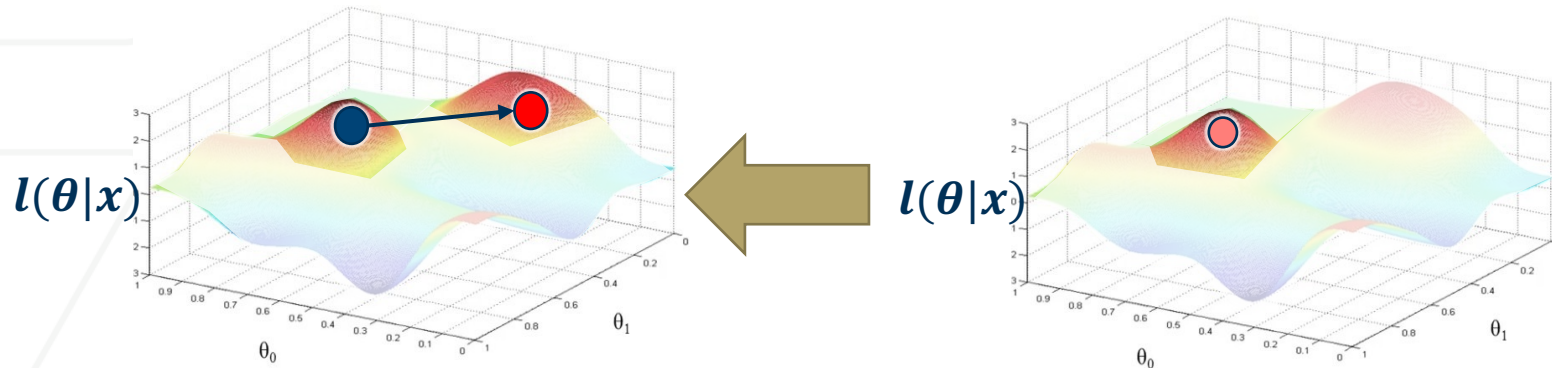
# Toy Manifold Example

How is this different from Explainability?



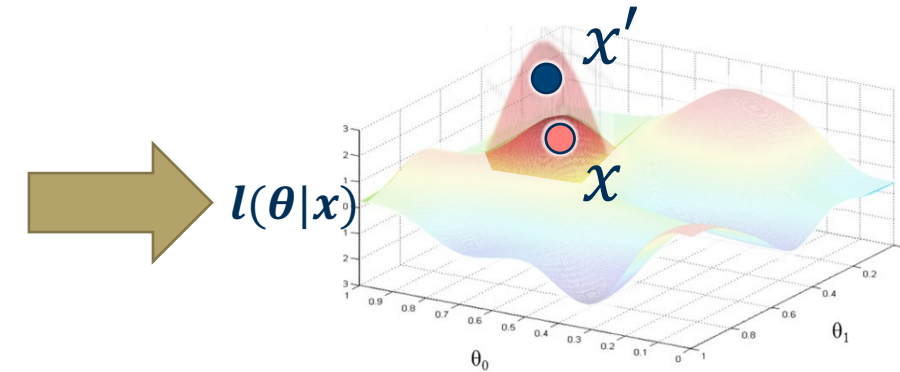
Probing the Purview of Neural Networks via Gradient Analysis

## Part 3: Explainability



- In Part 3: Activations of learned manifold are weighted by gradients w.r.t. activations to extract information and provide explanations

## Part 4: Uncertainty



- In Part 4: Statistics of gradients w.r.t. the weights (energy) will be directly used as features

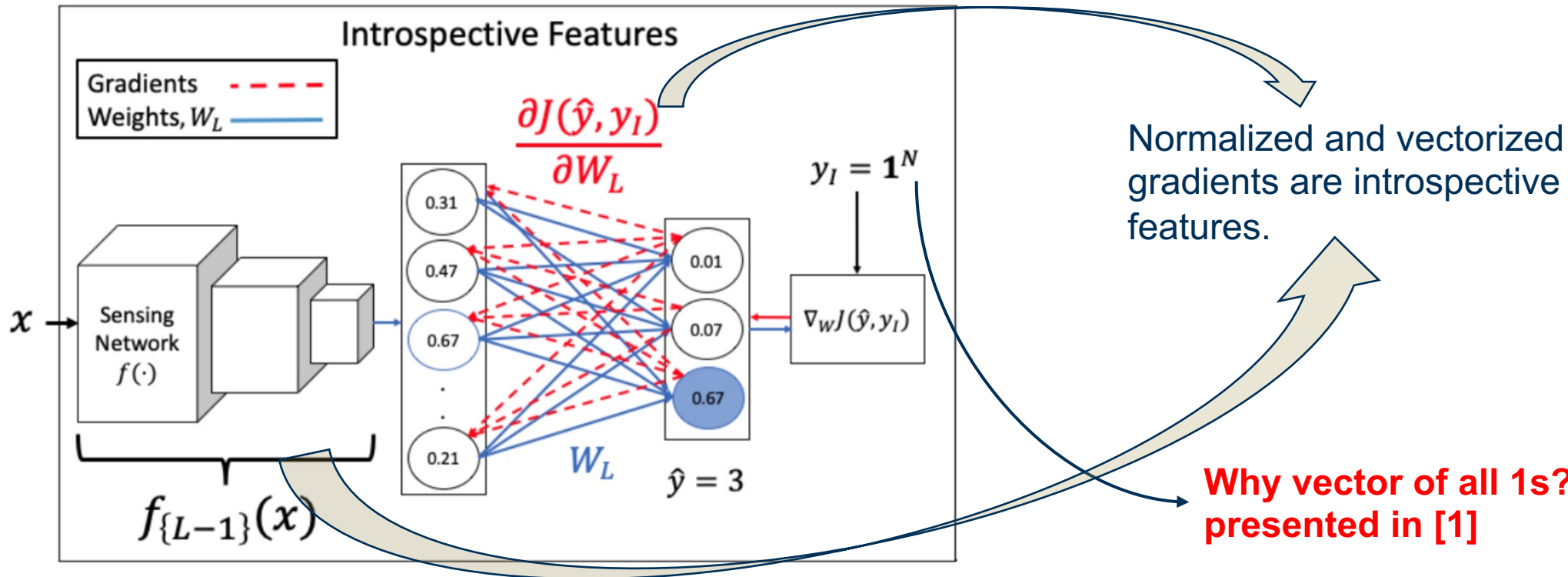
# Uncertainty in Neural Networks

## 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**



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

# Uncertainty in Neural Networks

## Utilizing Gradient Features



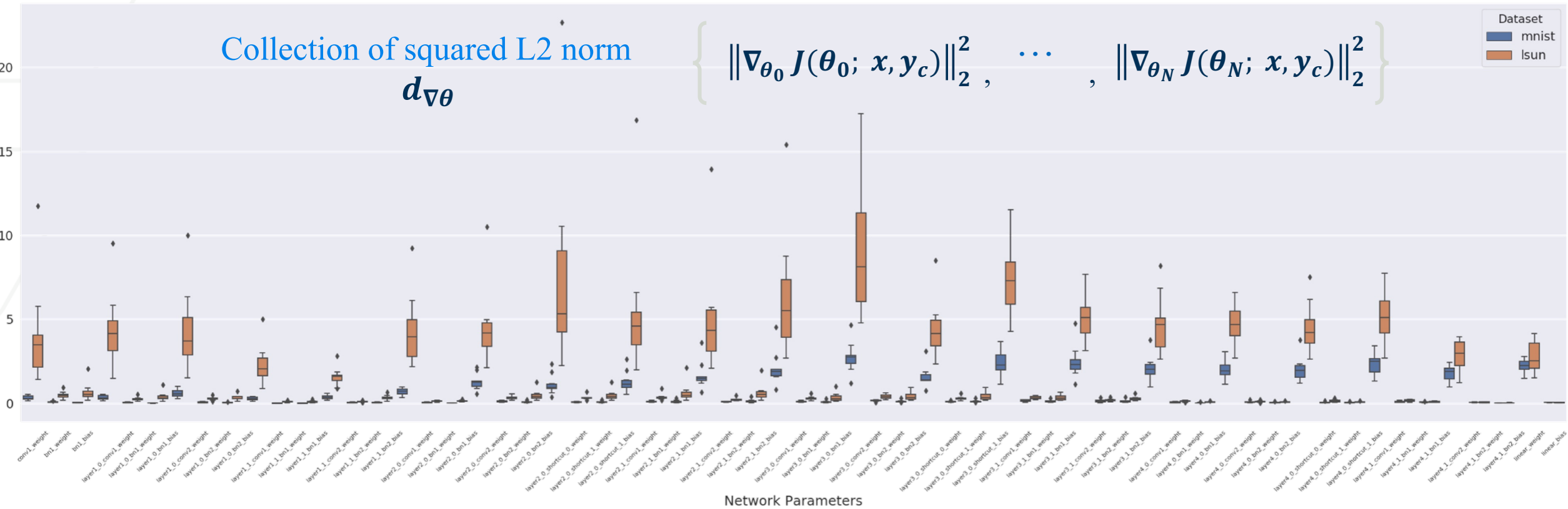
Probing the Purview of Neural Networks via Gradient Analysis

### 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



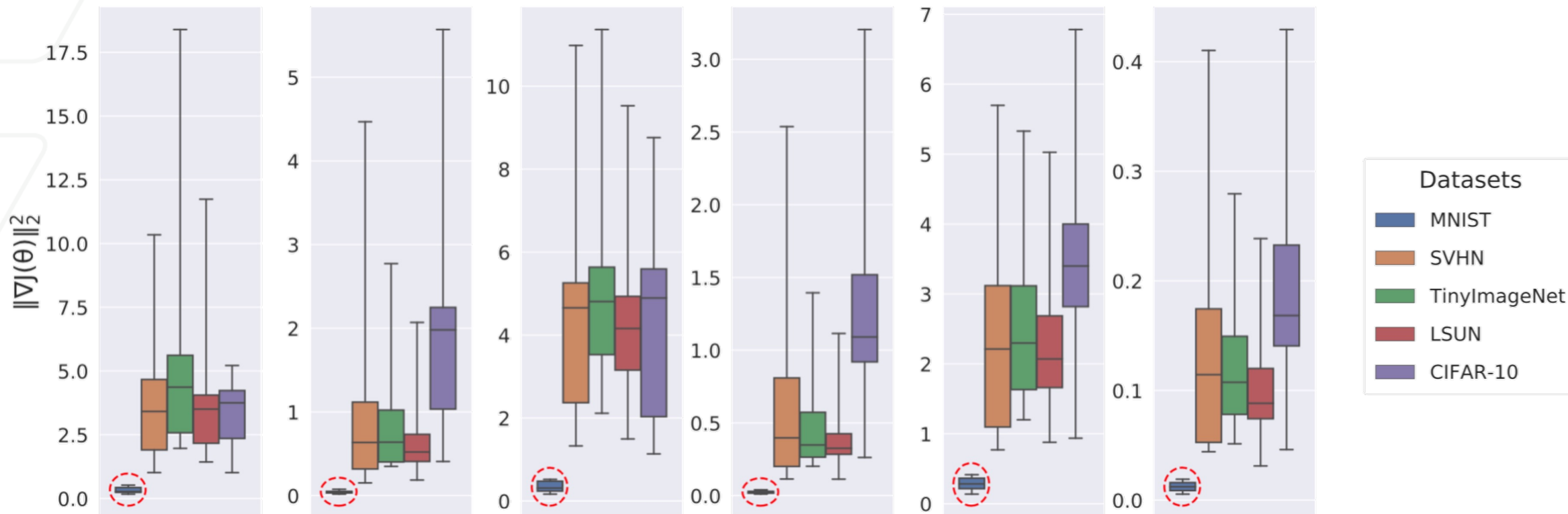
# Gradient-based Uncertainty

## 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**

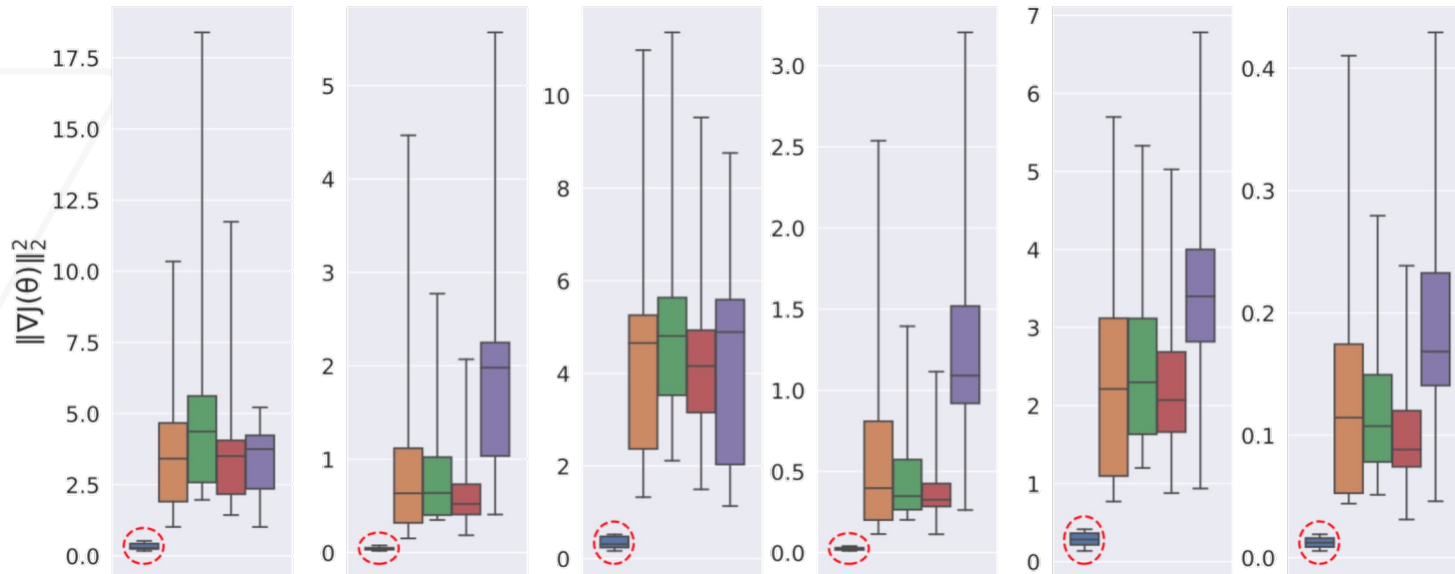
# Gradient-based Uncertainty

## 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

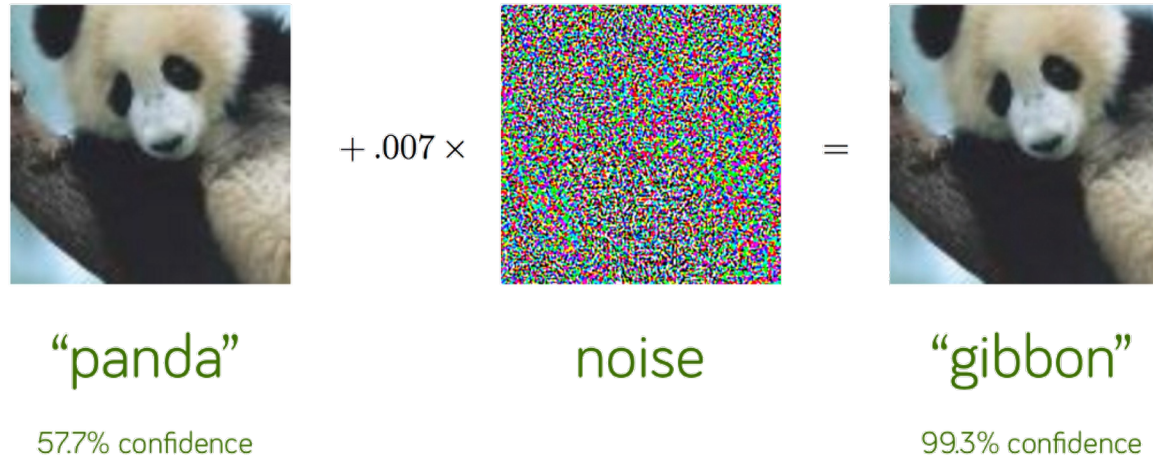
# Gradient-based Uncertainty

## 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	<b>99.95</b>	<b>99.95</b>	93.45
	BIM	49.94	99.21	87.09	89.20	<b>100.0</b>	<b>100.0</b>	96.19
	C&W	53.40	76.47	74.51	75.71	92.78	92.79	<b>97.07</b>
	PGD	50.03	67.48	56.27	57.57	65.23	75.98	<b>95.82</b>
	ITERLL	60.40	85.17	62.32	64.10	85.10	92.10	<b>98.17</b>
	SEMANTIC	52.29	86.25	64.18	65.79	83.95	84.38	<b>90.15</b>
DENSENET	FGSM	52.76	98.23	86.88	87.24	<b>99.98</b>	99.97	96.83
	BIM	49.67	<b>100.0</b>	89.19	89.17	<b>100.0</b>	<b>100.0</b>	96.85
	C&W	54.53	80.58	75.77	76.16	90.83	90.76	<b>97.05</b>
	PGD	49.87	83.01	70.39	66.52	86.94	83.61	<b>96.77</b>
	ITERLL	55.43	83.16	70.17	66.61	83.20	77.84	<b>98.53</b>
	SEMANTIC	53.54	81.41	62.16	62.15	67.98	67.29	<b>89.55</b>

# Gradient-based Uncertainty

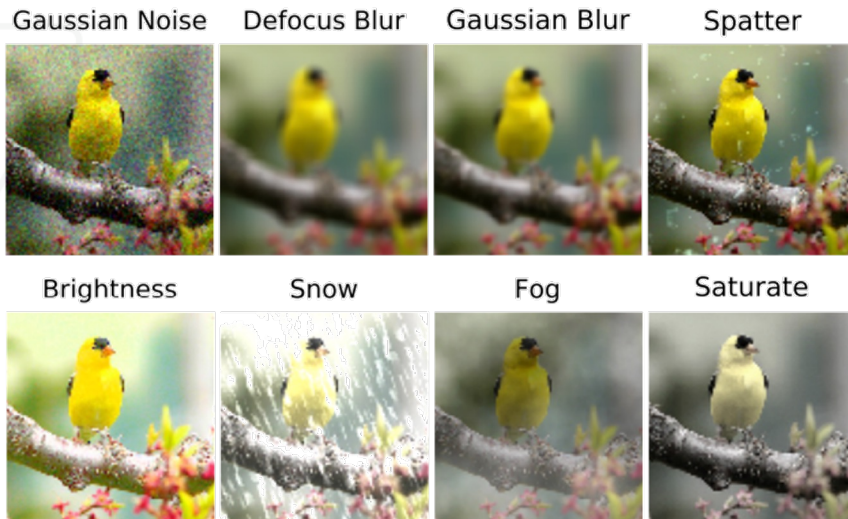
## Uncertainty in Detecting Challenging Conditions



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



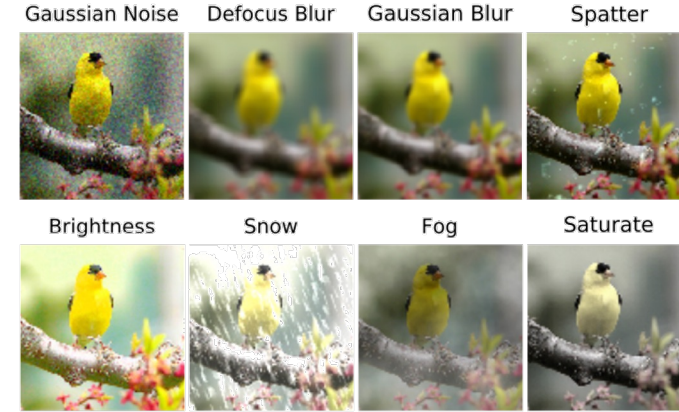
# Gradient-based Uncertainty

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



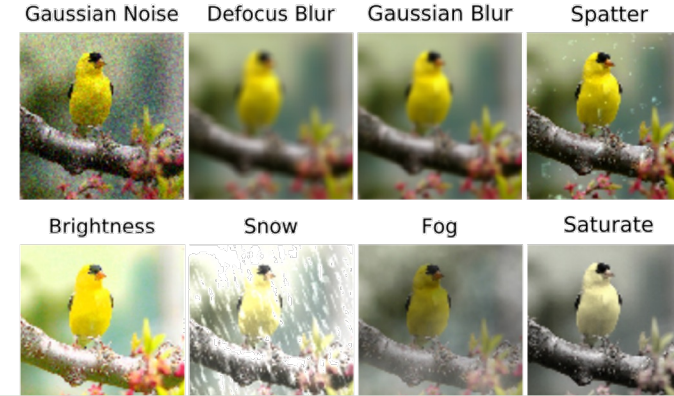
# Gradient-based Uncertainty

## Uncertainty in Detecting Challenging Conditions



Probing the Purview of Neural Networks via Gradient Analysis

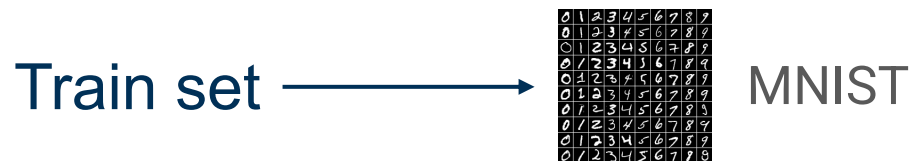
Dataset	Method	Mahalanobis [12] / Ours				
		Corruption	Level 1	Level 2	Level 3	Level 4
CIFAR-10-C	Noise	96.63 / <b>99.95</b>	98.73 / <b>99.97</b>	99.46 / <b>99.99</b>	99.62 / <b>99.97</b>	99.71 / <b>99.99</b>
	LensBlur	94.22 / <b>99.95</b>	97.51 / <b>99.99</b>	99.26 / <b>100.0</b>	99.78 / <b>100.0</b>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



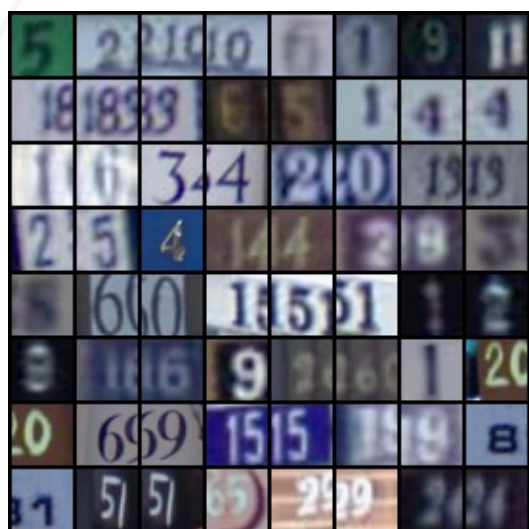
# Out-of-Distribution Detection



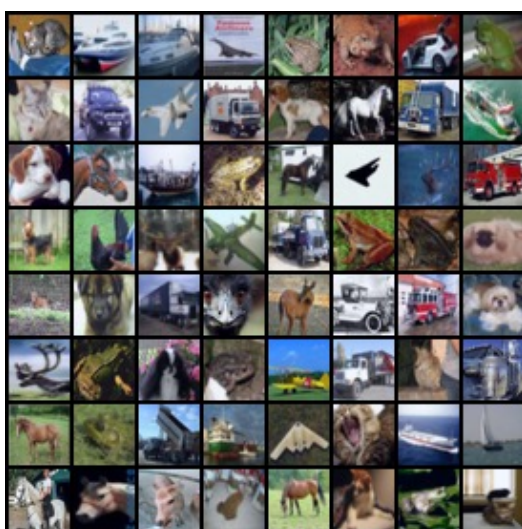
Probing the Purview of Neural Networks via Gradient Analysis



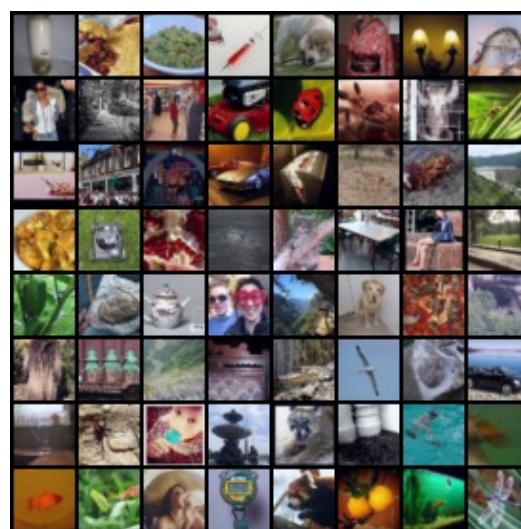
**Goal:** To detect that these datasets are not part of training



SVHN



CIFAR10



TinyImageNet



LSUN



# Out-of-Distribution Detection



Probing the Purview of Neural Networks via Gradient Analysis

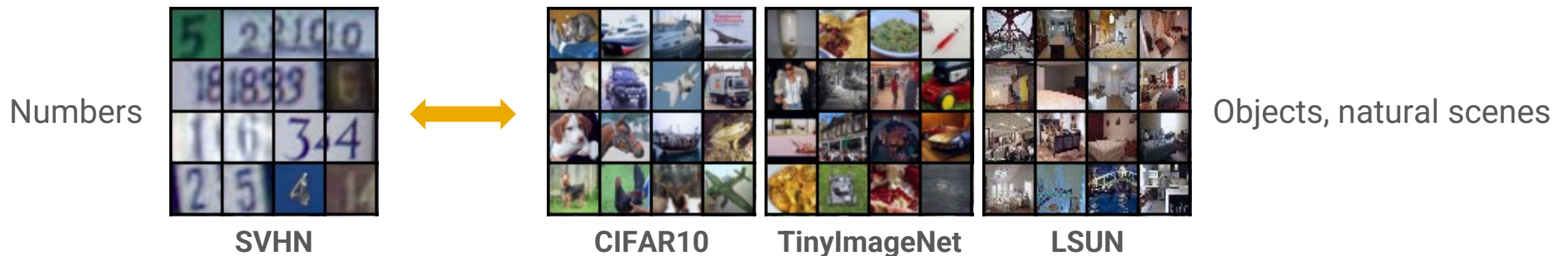
Dataset Distribution		Detection Accuracy	AUROC	AUPR
In	Out	Baseline [5] / ODIN [6] / Mahalanobis (V) [7] / Mahalanobis (P+FE) [7] / Ours		
CIFAR-10	SVHN	83.36 / 88.81 / 79.39 / 91.95 / <b>98.04</b>	88.30 / 94.93 / 85.03 / 97.10 / <b>99.84</b>	88.26 / 95.45 / 86.15 / 96.12 / <b>99.98</b>
	TinyImageNet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

# 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 / <b>98.04</b>	88.30 / 94.93 / 85.03 / 97.10 / <b>99.84</b>	88.26 / 95.45 / 86.15 / 96.12 / <b>99.98</b>
	TinyImageNet	84.01 / 85.21 / 83.60 / <b>97.45</b> / 86.17	90.06 / 91.86 / 88.93 / <b>99.68</b> / 93.18	89.26 / 91.60 / 88.59 / <b>99.60</b> / 92.66
	LSUN	87.34 / 88.42 / 85.02 / <b>98.60</b> / 98.37	92.79 / 94.48 / 90.11 / <b>99.86</b> / <b>99.86</b>	92.30 / 94.22 / 89.80 / 99.82 / <b>99.87</b>
SVHN	CIFAR-10	79.98 / 80.12 / 74.10 / 88.84 / <b>97.90</b>	81.50 / 81.49 / 79.31 / 95.05 / <b>99.79</b>	81.01 / 80.95 / 80.83 / 90.25 / <b>98.11</b>
	TinyImageNet	81.70 / 81.92 / 79.35 / 96.17 / <b>97.74</b>	83.69 / 83.82 / 83.85 / 99.23 / <b>99.77</b>	82.54 / 82.60 / 85.50 / <b>98.17</b> / 97.93
	LSUN	80.96 / 81.15 / 79.52 / 97.50 / <b>99.04</b>	82.85 / 82.98 / 83.02 / 99.54 / <b>99.93</b>	81.97 / 82.01 / 84.67 / 98.84 / <b>99.21</b>



# 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 / <b>98.04</b>	88.30 / 94.93 / 85.03 / 97.10 / <b>99.84</b>	88.26 / 95.45 / 86.15 / 96.12 / <b>99.98</b>
	TinyImageNet	84.01 / 85.21 / 83.60 / <b>97.45</b> / 86.17	90.06 / 91.86 / 88.93 / <b>99.68</b> / 93.18	89.26 / 91.60 / 88.59 / <b>99.60</b> / 92.66
	LSUN	87.34 / 88.42 / 85.02 / <b>98.60</b> / 98.37	92.79 / 94.48 / 90.11 / <b>99.86</b> / <b>99.86</b>	92.30 / 94.22 / 89.80 / 99.82 / <b>99.87</b>
SVHN	CIFAR-10	79.98 / 80.12 / 74.10 / 88.84 / <b>97.90</b>	81.50 / 81.49 / 79.31 / 95.05 / <b>99.79</b>	81.01 / 80.95 / 80.83 / 90.25 / <b>98.11</b>
	TinyImageNet	81.70 / 81.92 / 79.35 / 96.17 / <b>97.74</b>	83.69 / 83.82 / 83.85 / 99.23 / <b>99.77</b>	82.54 / 82.60 / 85.50 / <b>98.17</b> / 97.93
	LSUN	80.96 / 81.15 / 79.52 / 97.50 / <b>99.04</b>	82.85 / 82.98 / 83.02 / 99.54 / <b>99.93</b>	81.97 / 82.01 / 84.67 / 98.84 / <b>99.21</b>

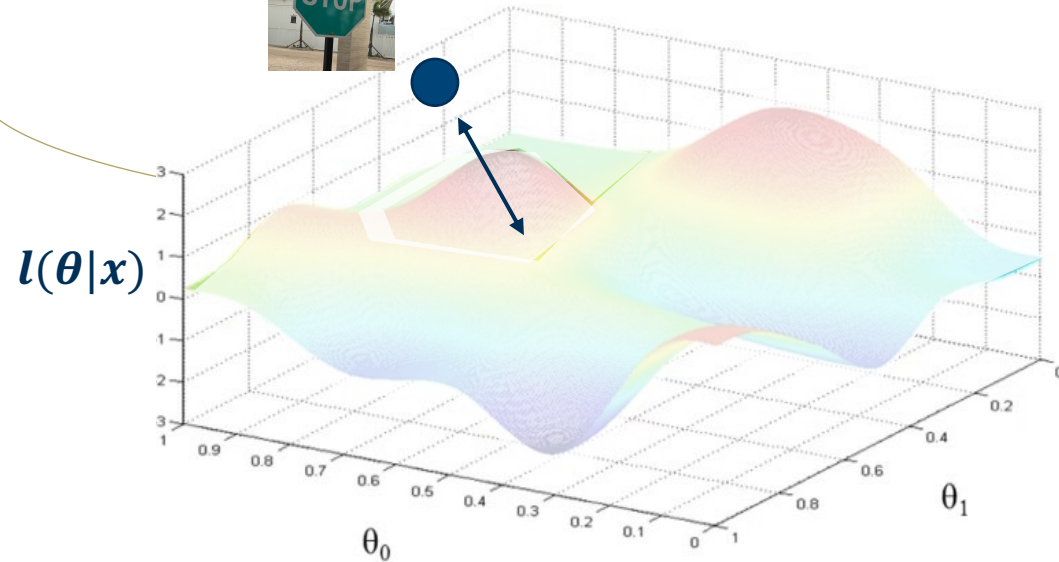


# Case Study: Introspective Learning

## 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 bull mastiff?



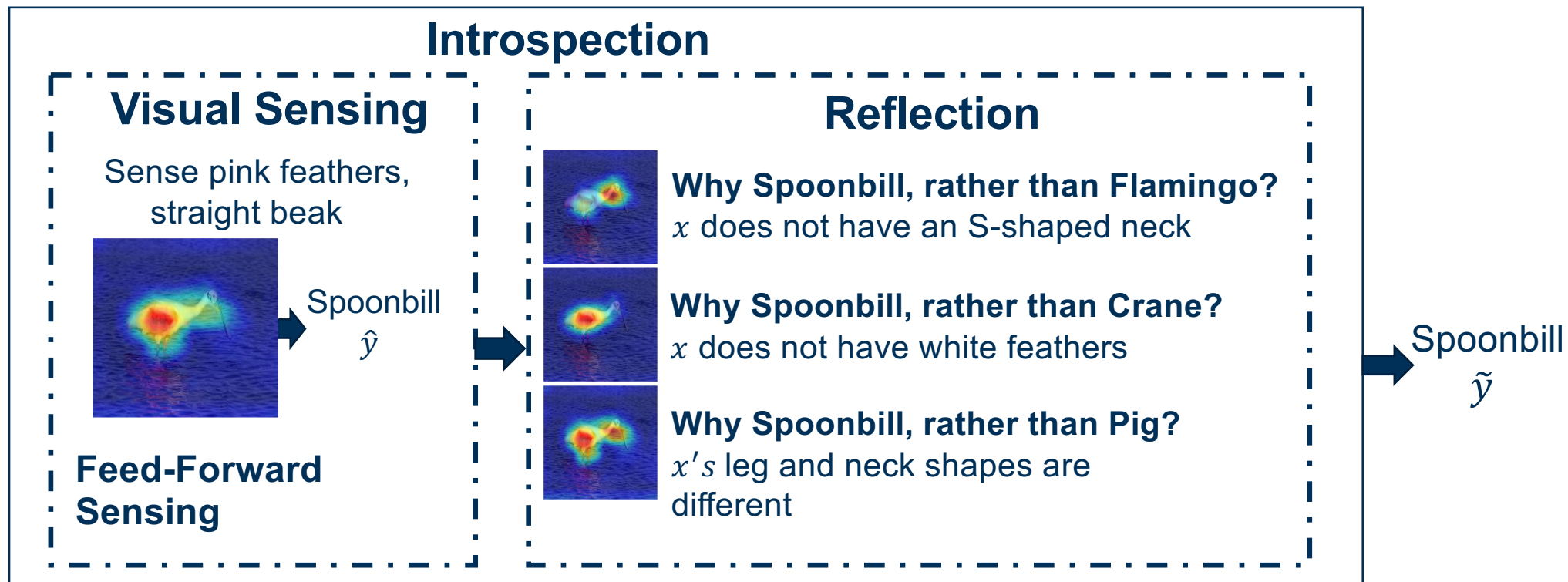
# Introspection

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?**



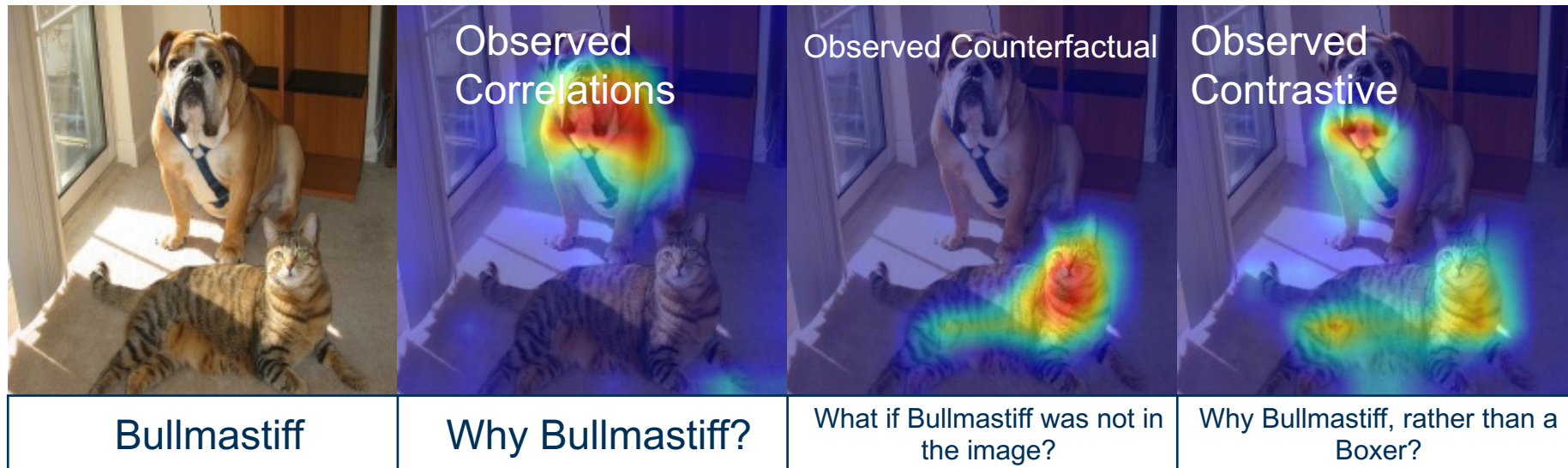
# Introspection

## 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  $x$ .*

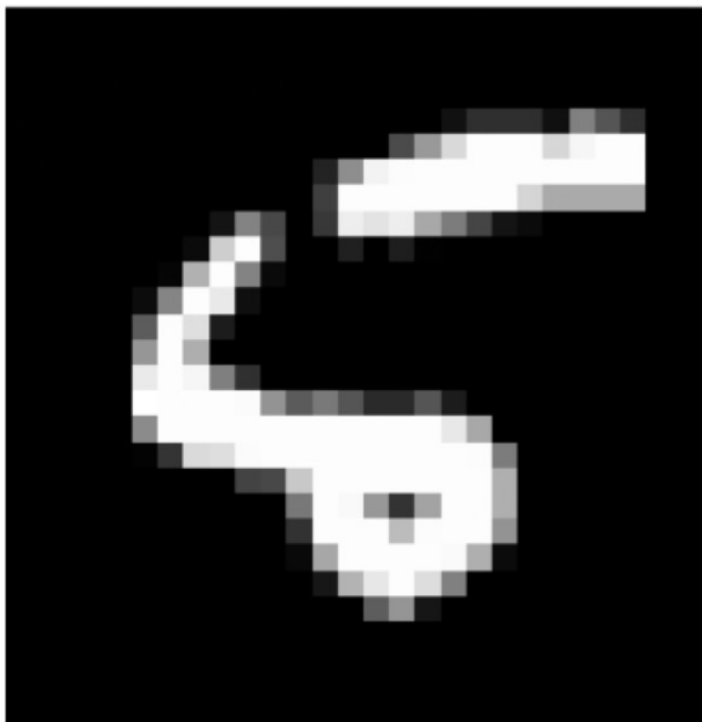
# Introspection

## 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  $x$



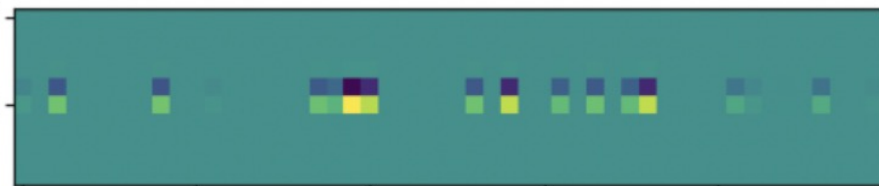
Why 5, rather than 0?



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?

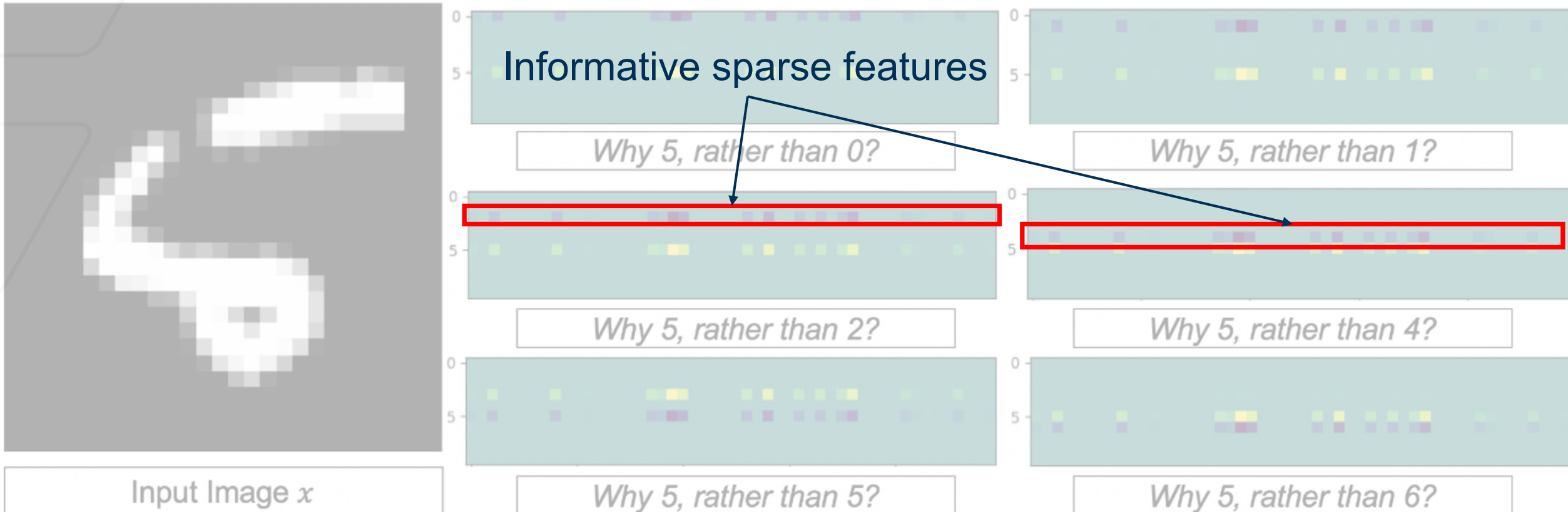
# Introspection

## 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



# Introspection

## Gradients as Features



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

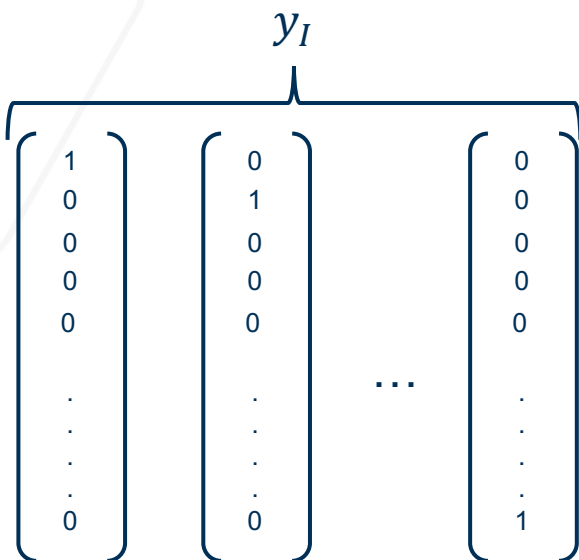
For a well-trained network, the gradients are robust

$\nabla_W$  = 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  $\hat{y}$

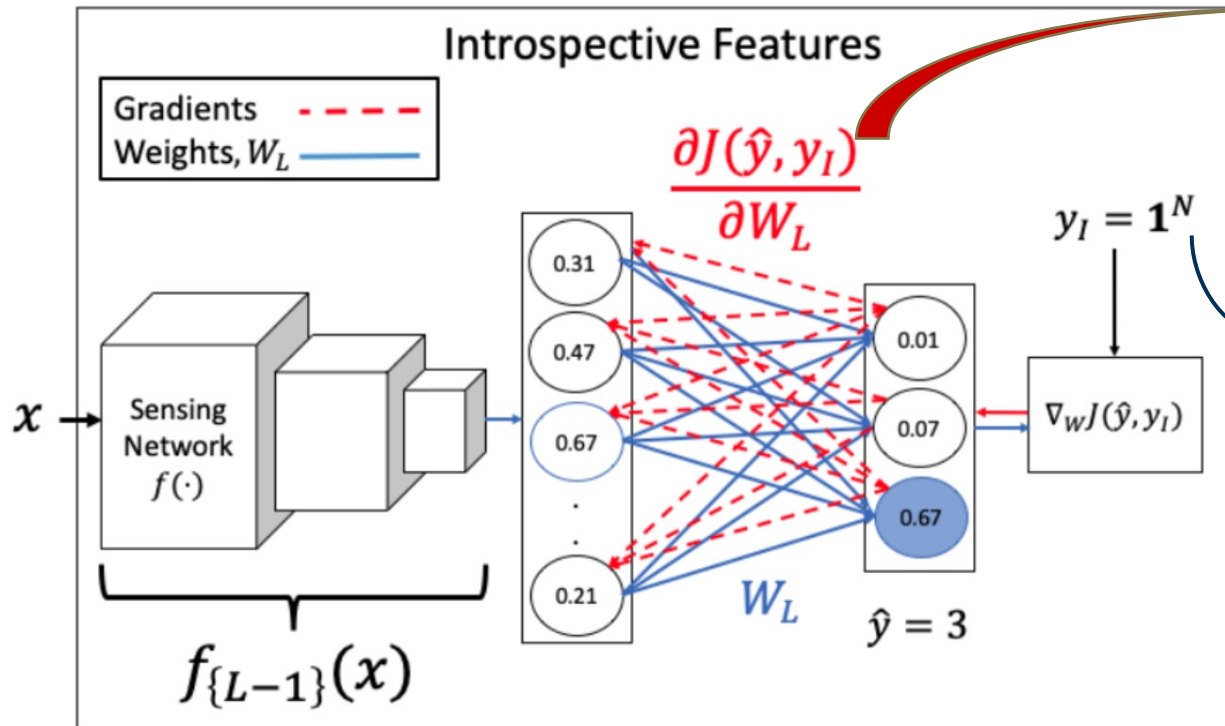
# Introspection

## 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



Normalized and vectorized gradients are introspective features

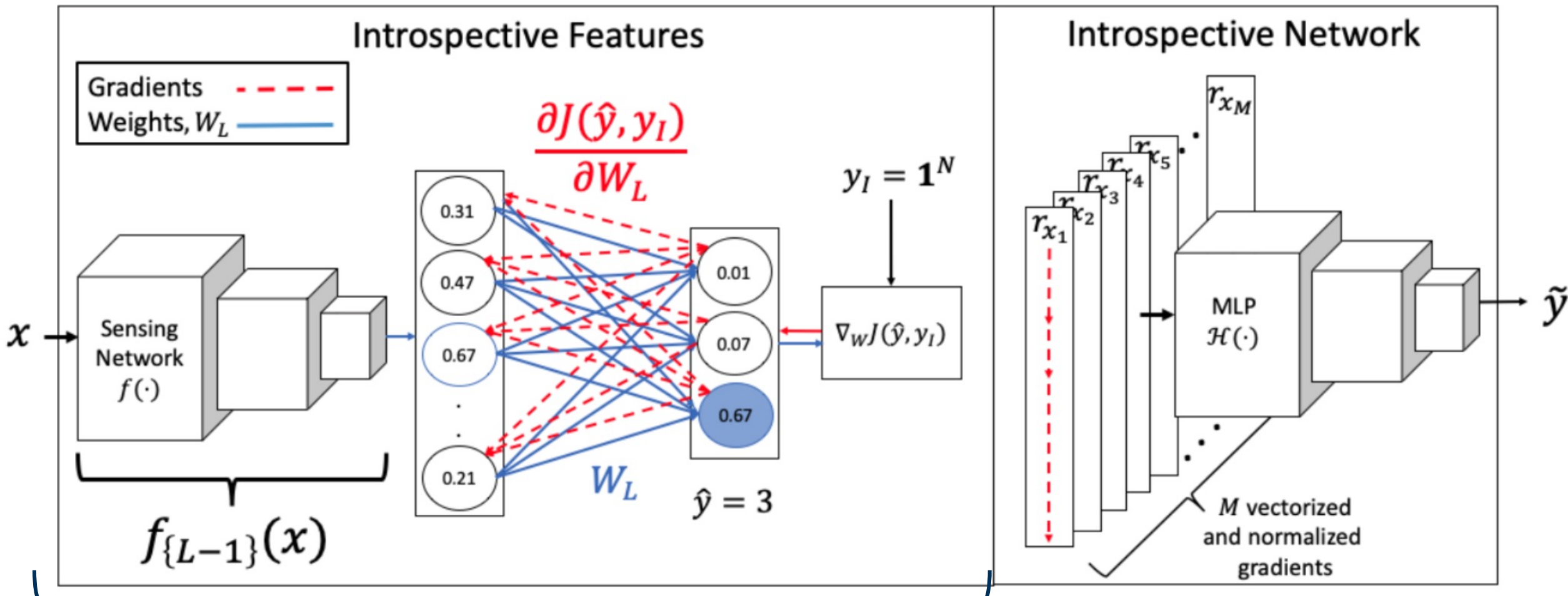
**Vector of all ones: A confounding label!**

# Introspection

## Utilizing Gradient Features



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



## Introspective Features

# Introspection

When is Introspection Useful?



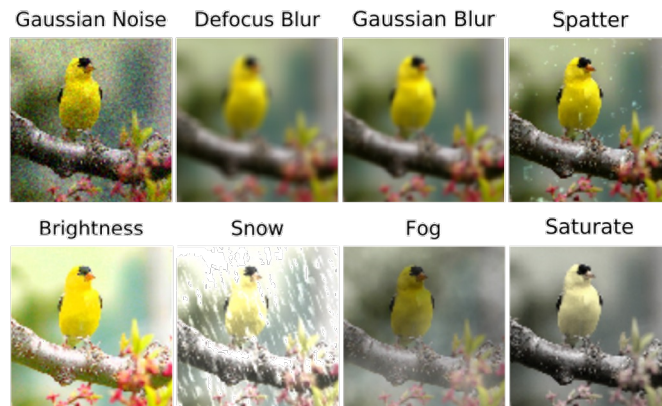
Introspective Learning: A Two-stage Approach for Inference in Neural Networks

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





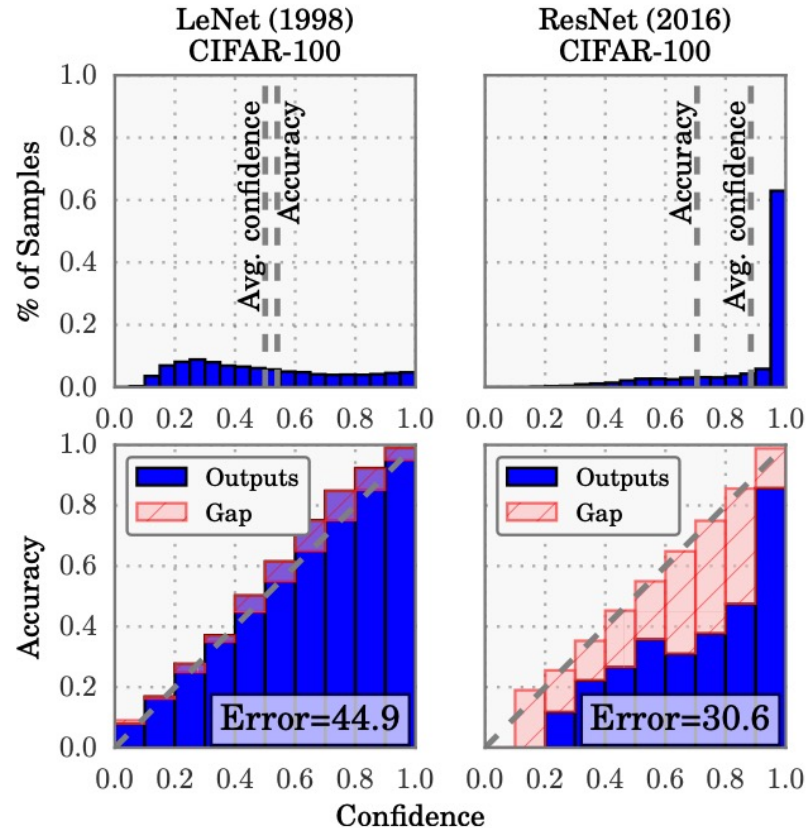
# Calibration

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

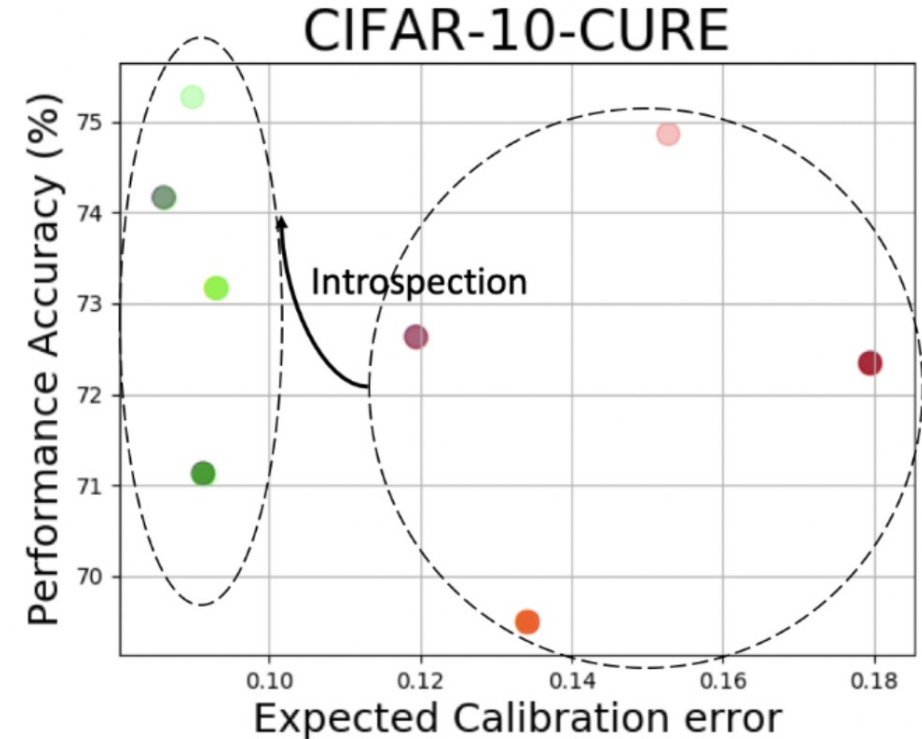
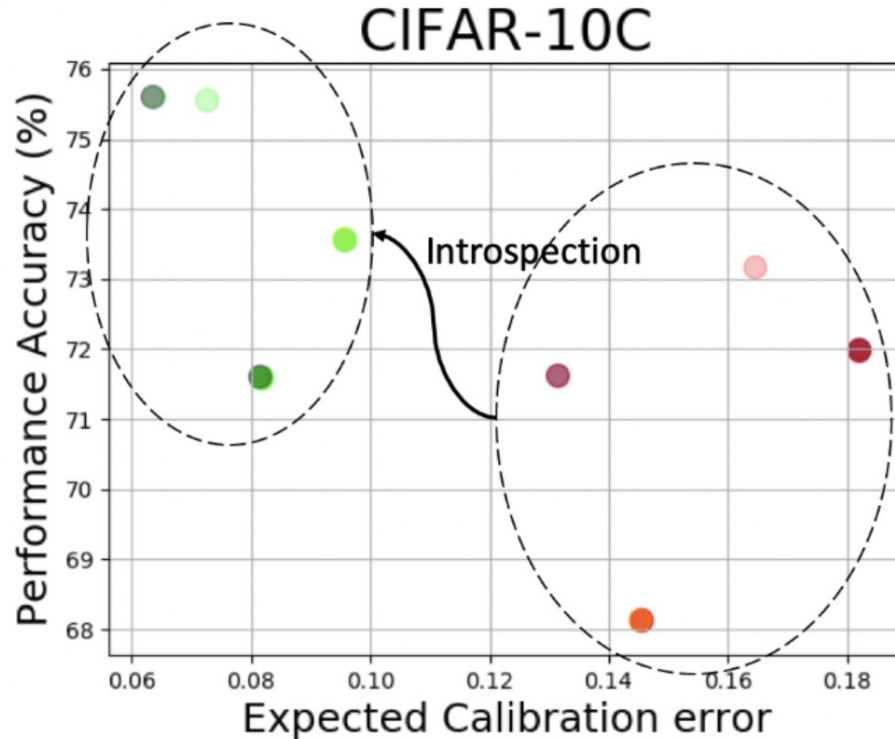


Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Ideal: Top-left corner

Y-Axis: Generalization

X-Axis: Calibration



**Legend**

<b>Feed-Forward Networks</b>	● ResNet-18	● ResNet-34	● ResNet-50	● ResNet-101
<b>After Introspection</b>	● 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	<b>71.4%</b>
DENOISING	FEED-FORWARD	65.02%
	INTROSPECTIVE	<b>68.86%</b>
ADVERSARIAL TRAIN (27)	FEED-FORWARD	68.02%
	INTROSPECTIVE	<b>70.86%</b>
SIMCLR (19)	FEED-FORWARD	70.28%
	INTROSPECTIVE	<b>73.32%</b>
AUGMENT NOISE (23)	FEED-FORWARD	76.86%
	INTROSPECTIVE	<b>77.98%</b>
AUGMIX (24)	FEED-FORWARD	89.85%
	INTROSPECTIVE	<b>89.89%</b>

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

Table 2: Recognition accuracy of Active Learning strategies.

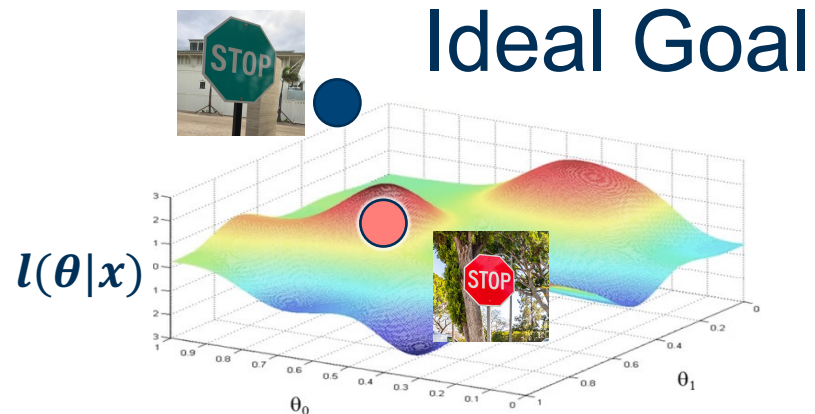
Methods	Architecture	Original Testset		Gaussian Noise	
		R-18	R-34	R-18	R-34
Entropy (E1)	Feed-Forward	0.365	0.358	0.244	0.249
	Introspective	0.365	0.359	<b>0.258</b>	<b>0.255</b>
Least (E1)	Feed-Forward	0.371	0.359	0.252	0.25
	Introspective	0.373	0.362	<b>0.264</b>	<b>0.26</b>
Margin (E2)	Feed-Forward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

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

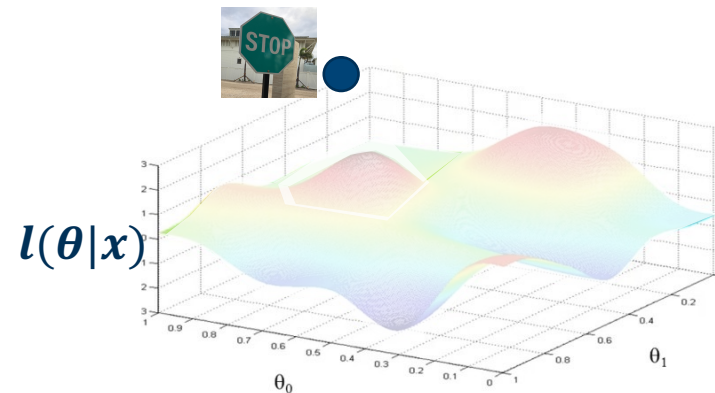
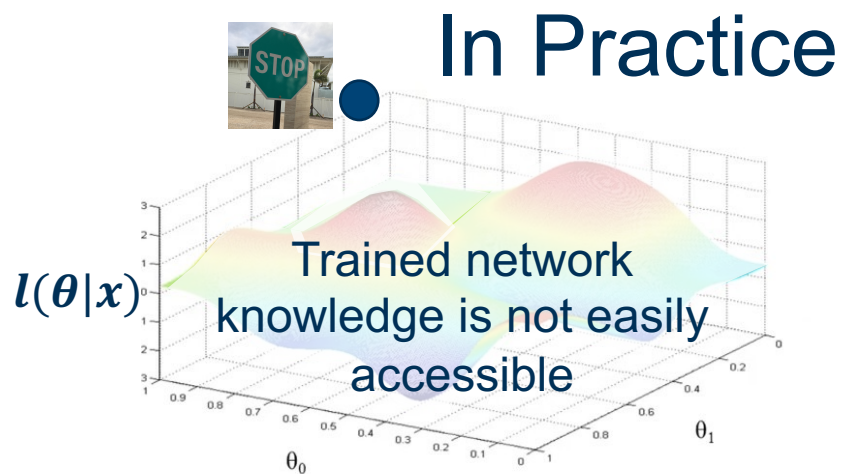
Methods	OOD Datasets	FPR (95% at TPR)	Detection Error	AUROC
		↓	↓	↑
Feed-Forward/Introspective				
MSP (E3)	Textures	58.74/19.66	18.04/7.49	88.56/97.79
	SVHN	61.41/51.27	16.92/15.67	89.39/91.2
	Places365	58.04/54.43	17.01/15.07	89.39/91.3
	LSUN-C	27.95/27.5	9.42/10.29	96.07/95.73
ODIN (E6)	Textures	52.3/9.31	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

# Part I, II and III

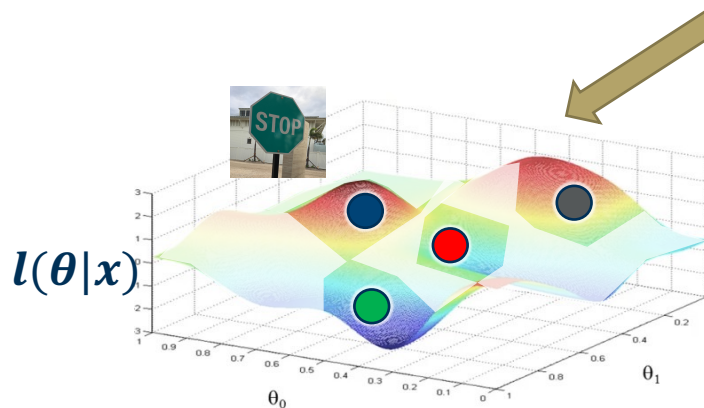
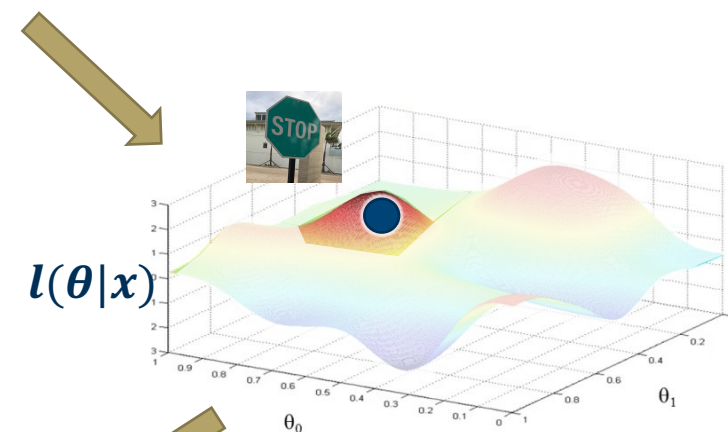
## Tying it Back



From Part I



Novel data projects onto the likelihood function (however incorrectly), and extracts fisher information around the projection



By backpropagating contrast classes (and not updating the network), the network finds the steepest descent towards other regions of likelihood function