

Robust Neural Networks

Part 3: Uncertainty at Inference

Objective

Objective of the Tutorial

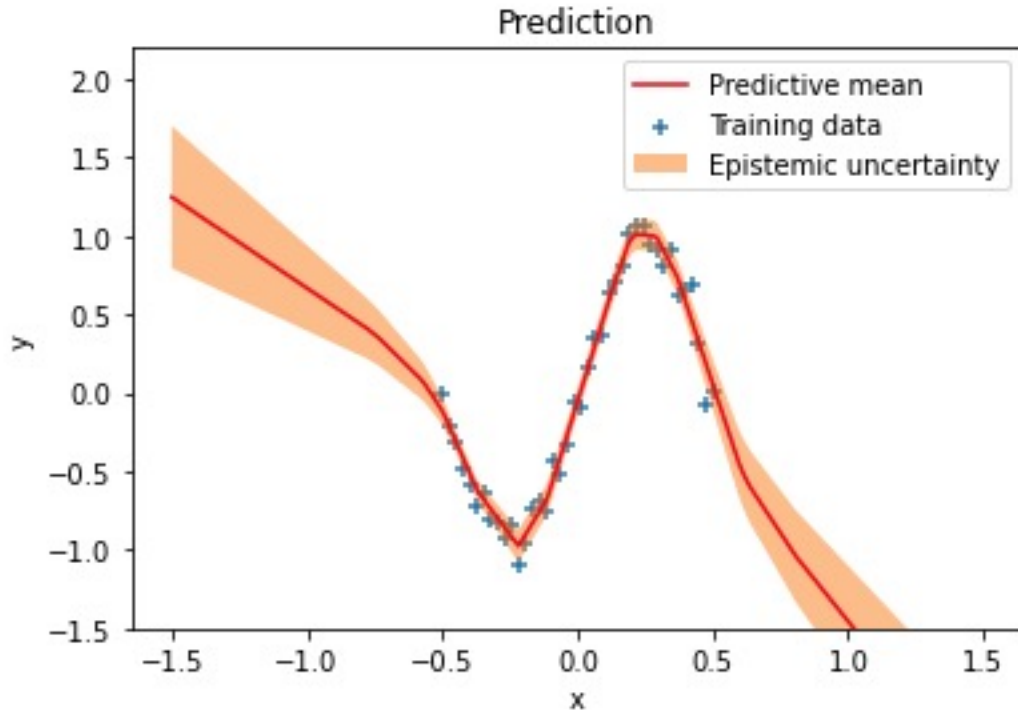
To discuss methodologies that promote robustness in neural networks at inference

- 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



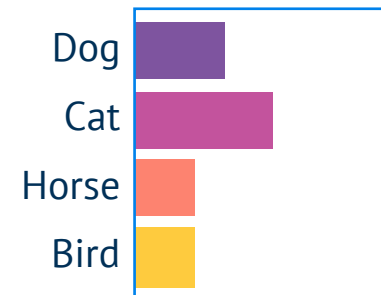
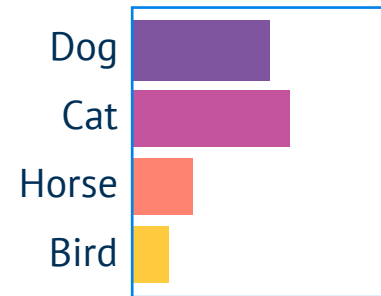
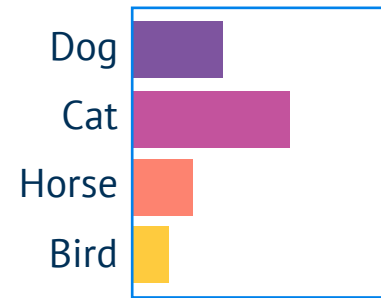
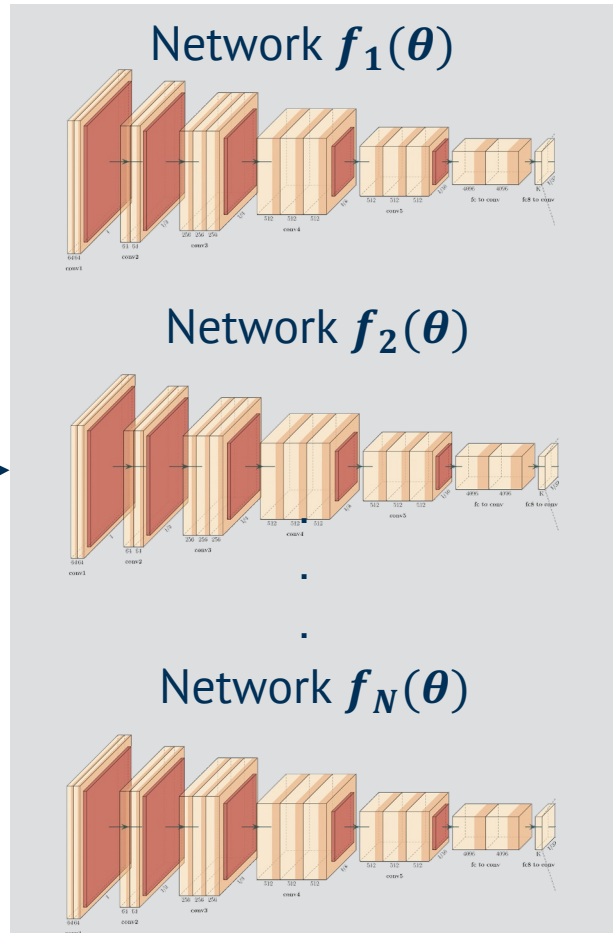
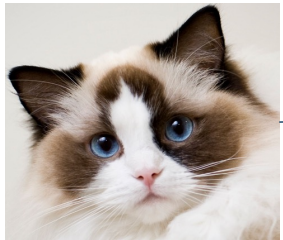
A simple example:

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

Uncertainty

Uncertainty Quantification in Neural Networks

Via Ensembles¹

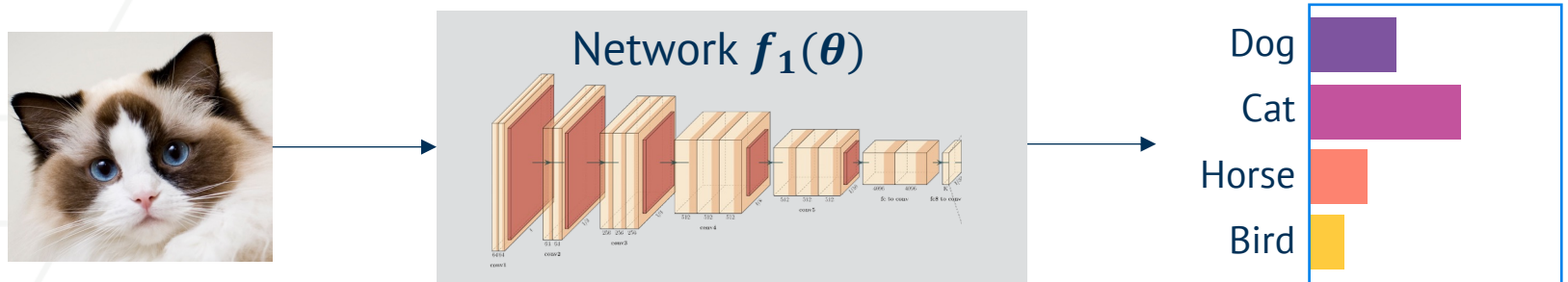


Variation within outputs $Var(y)$ is the uncertainty. Commonly referred to as **Prediction Uncertainty.**

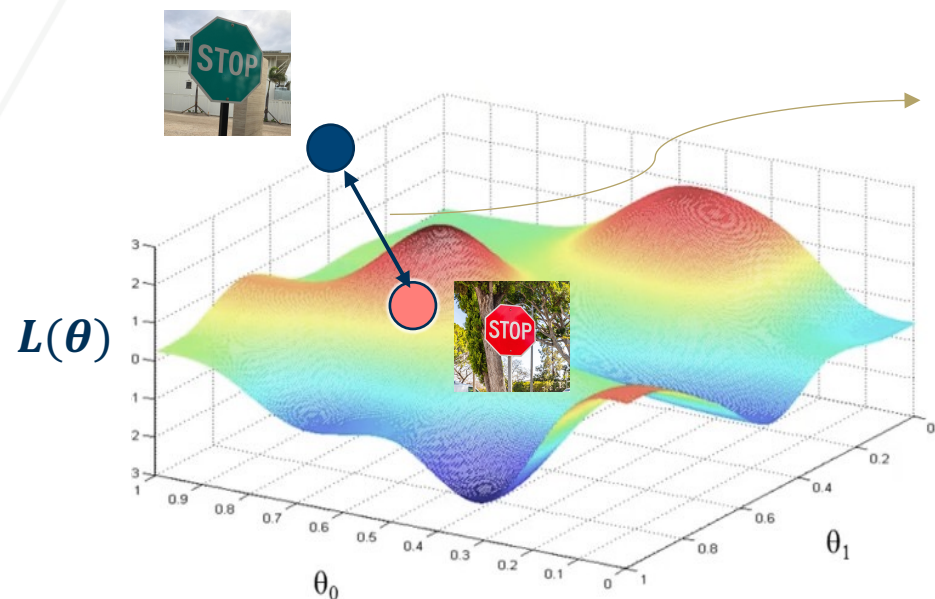
Uncertainty

Uncertainty Quantification in Neural Networks

Via Single pass methods¹



Uncertainty quantification using a single network and a single pass



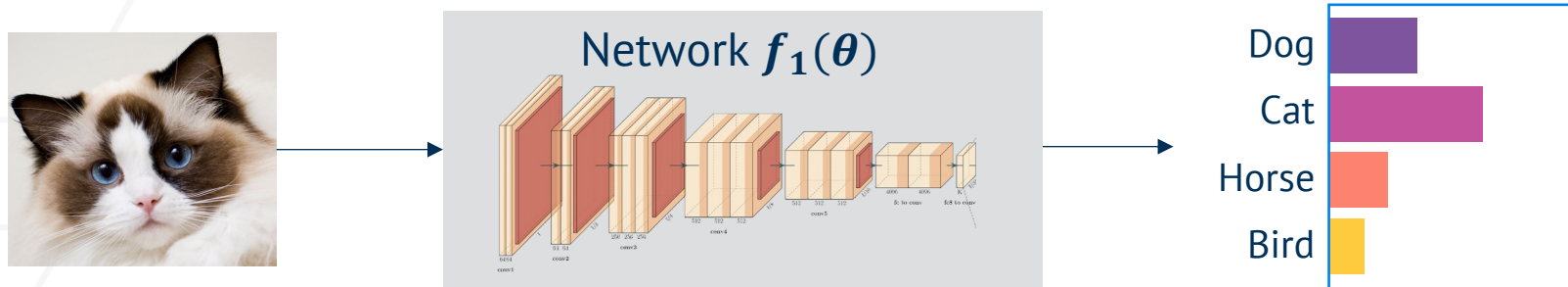
Calculate distance from some trained clusters

Does not require multiple networks!

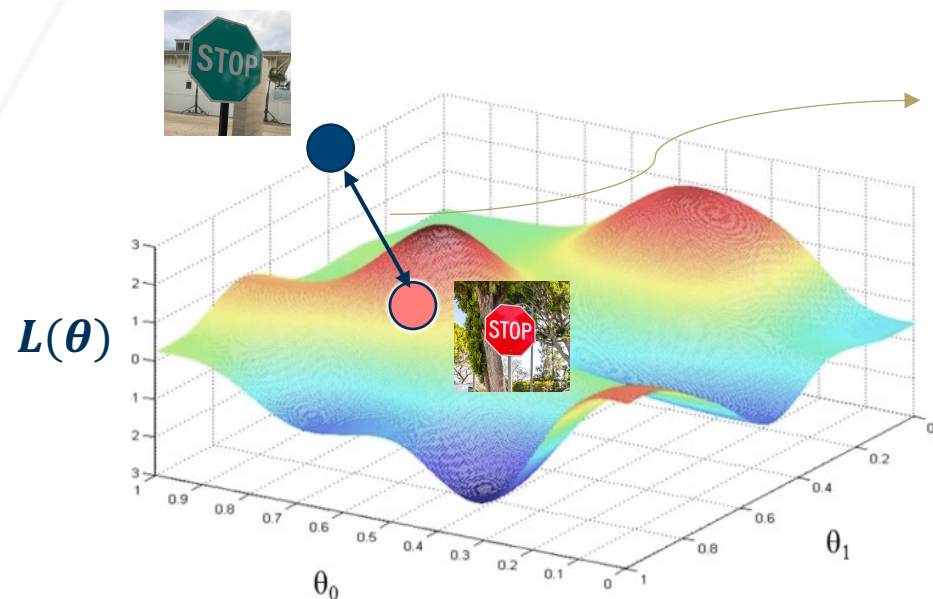
Uncertainty

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

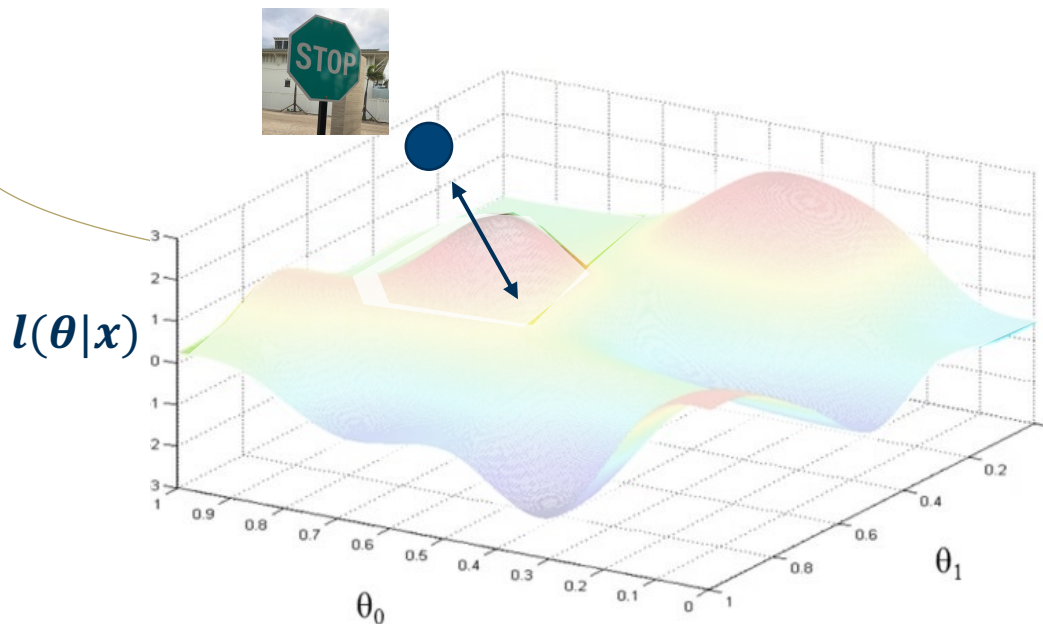
Gradients as Single pass Features

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

Distance from unknown cluster

Two techniques:

1. **Backpropagating Confounding labels for Adversarial Detection**
2. Backpropagating Confounding labels for Robust Prediction





Probing the Purview of Neural Networks via Gradient Analysis



Jinsol Lee,
PhD Candidate



Mohit Prabhushankar, PhD
Postdoc



Ghassan AlRegib, PhD
Professor



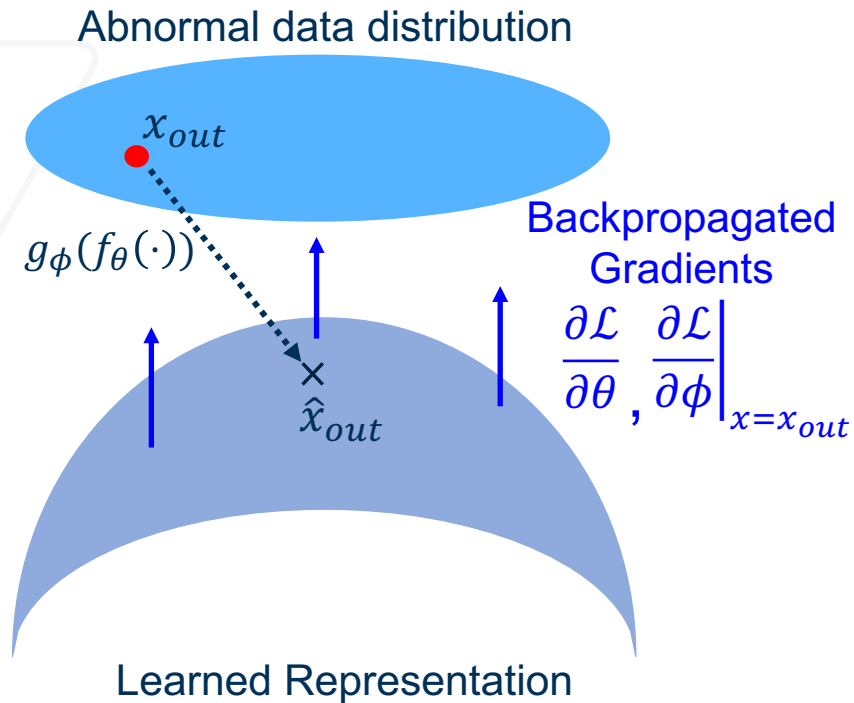
Uncertainty in Neural Networks

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



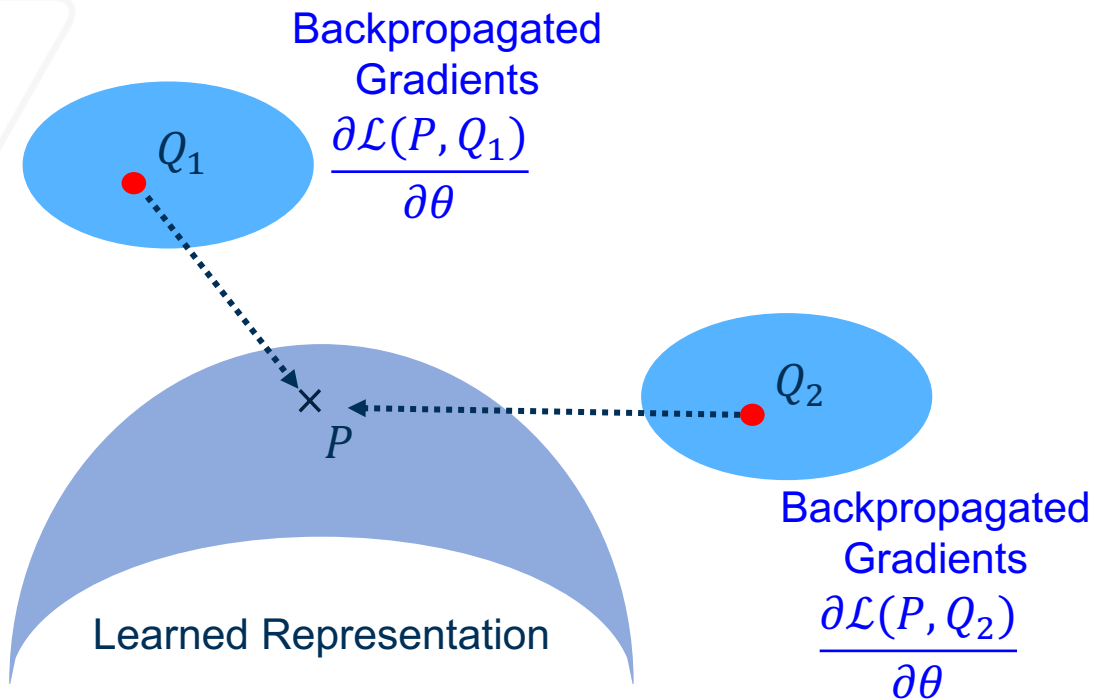
Probing the Purview of Neural Networks via Gradient Analysis

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

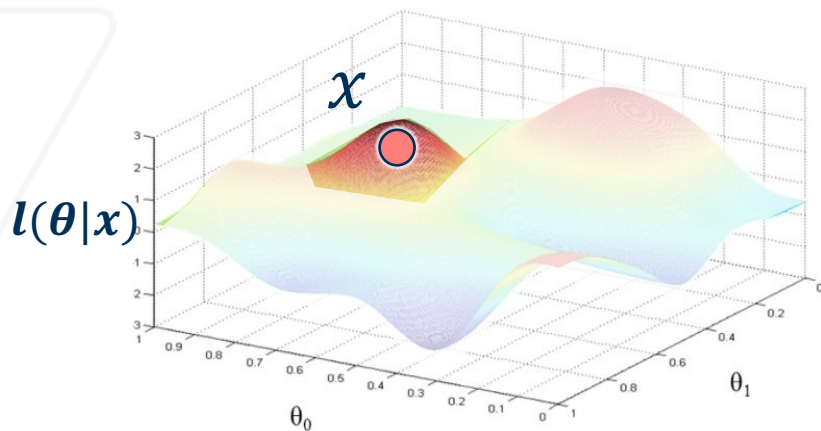
What is uncertainty?



Probing the Purview of Neural Networks via Gradient Analysis

Gradients represent the local required change in manifold

Similar to introspective learning!



Contrast class 1



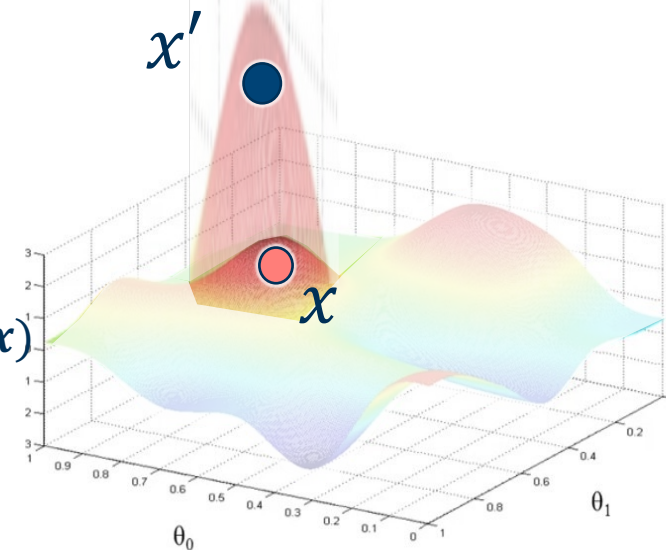
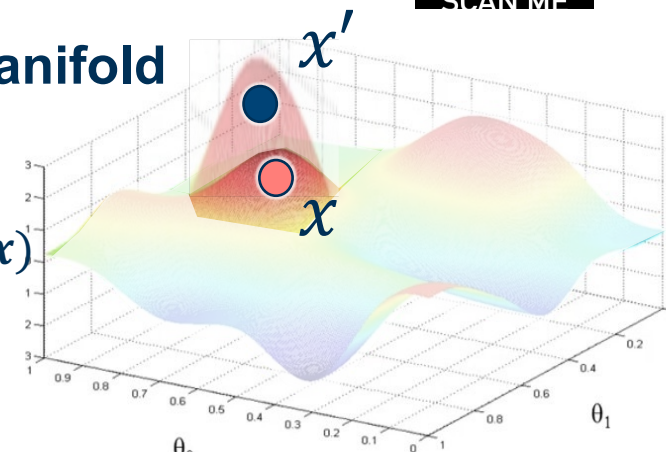
$l(\theta|x)$

·
·
·

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)!

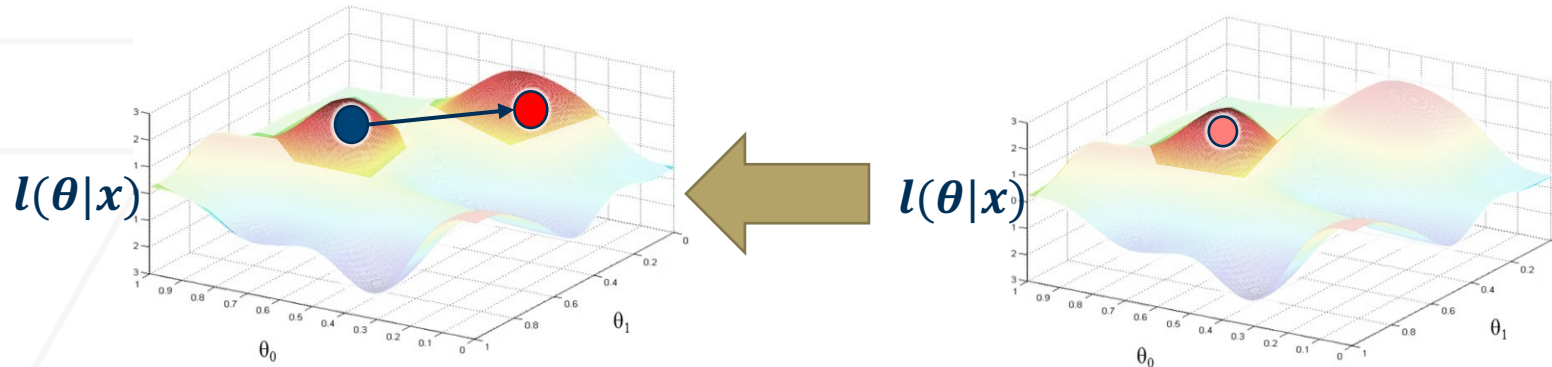
Toy Manifold Example

How is this different from Explainability?



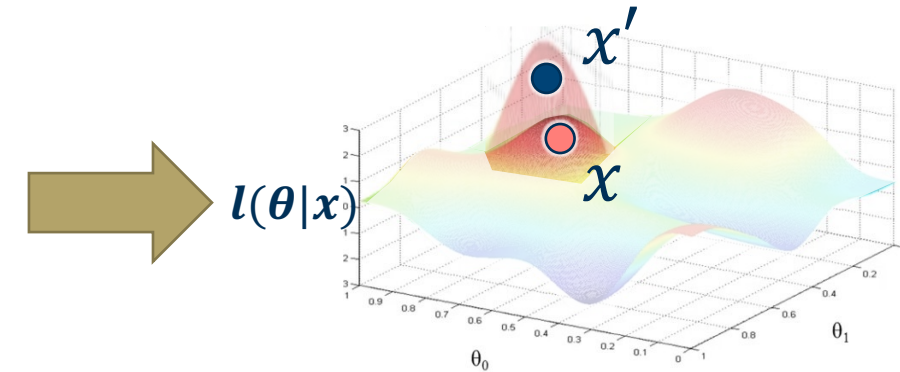
Probing the Purview of Neural Networks via Gradient Analysis

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

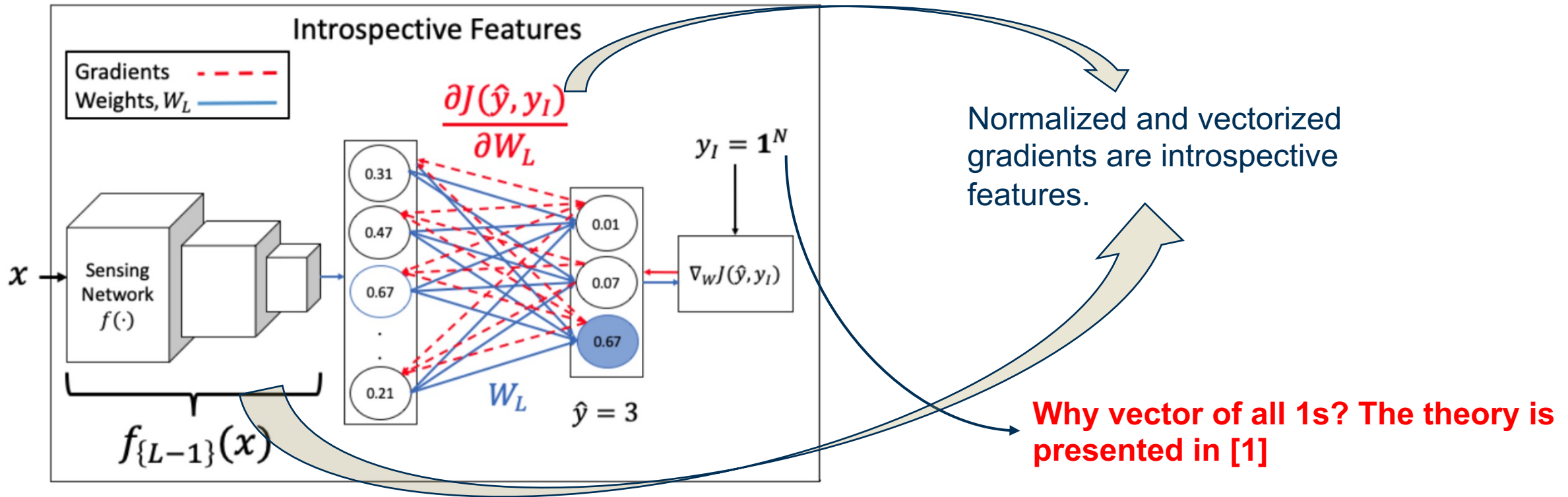
Uncertainty in Neural Networks

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



Uncertainty in Neural Networks

Utilizing Gradient Features



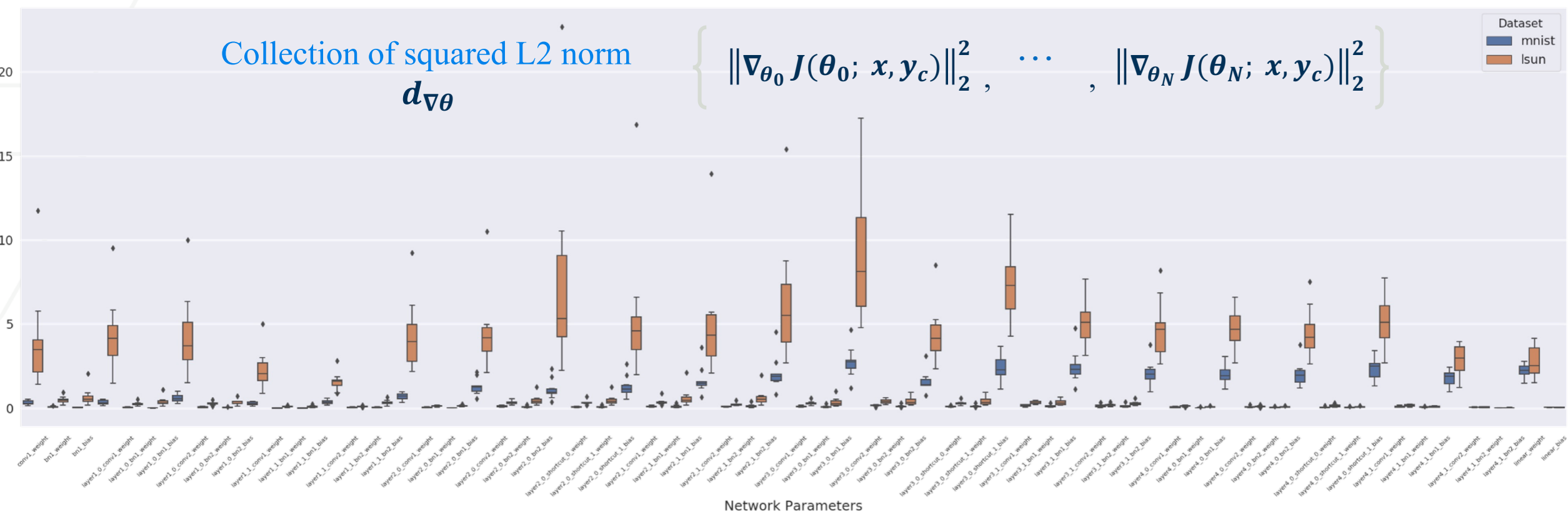
Probing the Purview of Neural Networks via Gradient Analysis

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

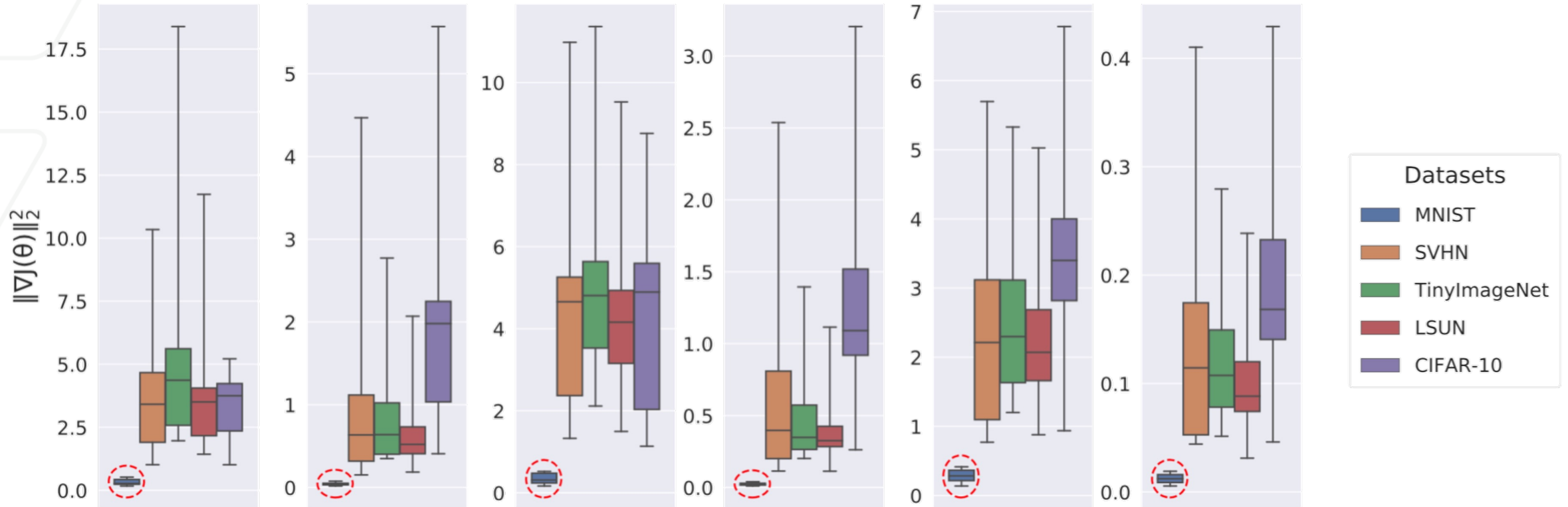
Gradient-based Uncertainty

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

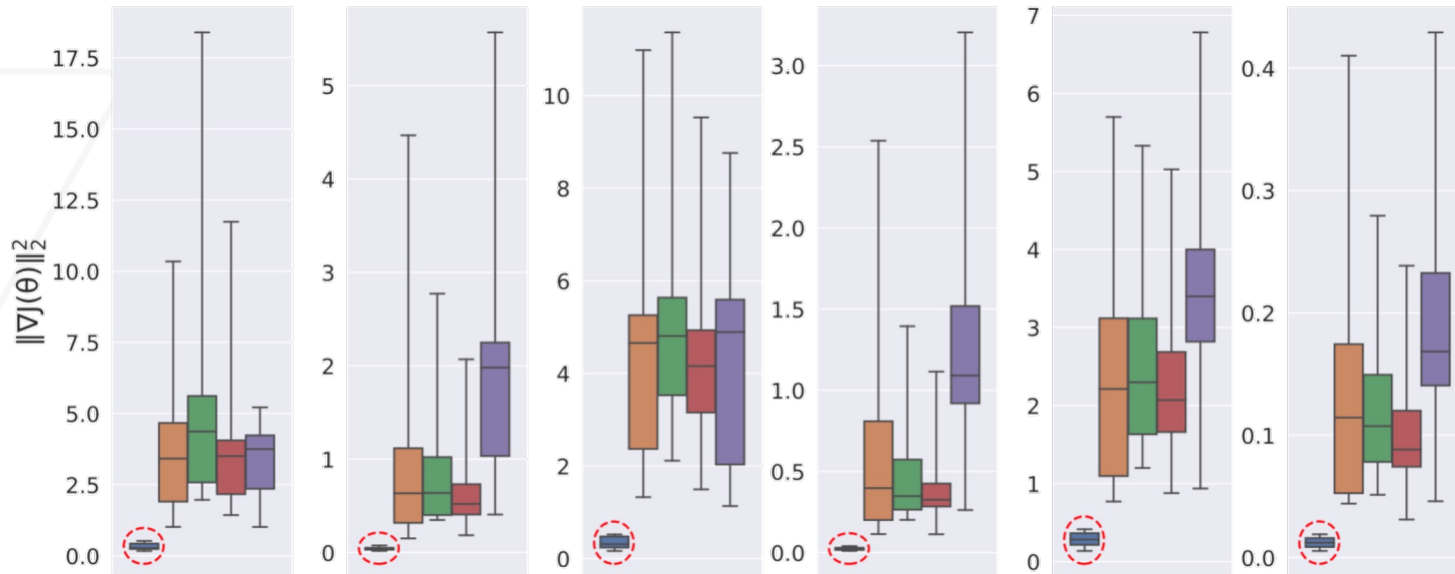
Gradient-based Uncertainty

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**

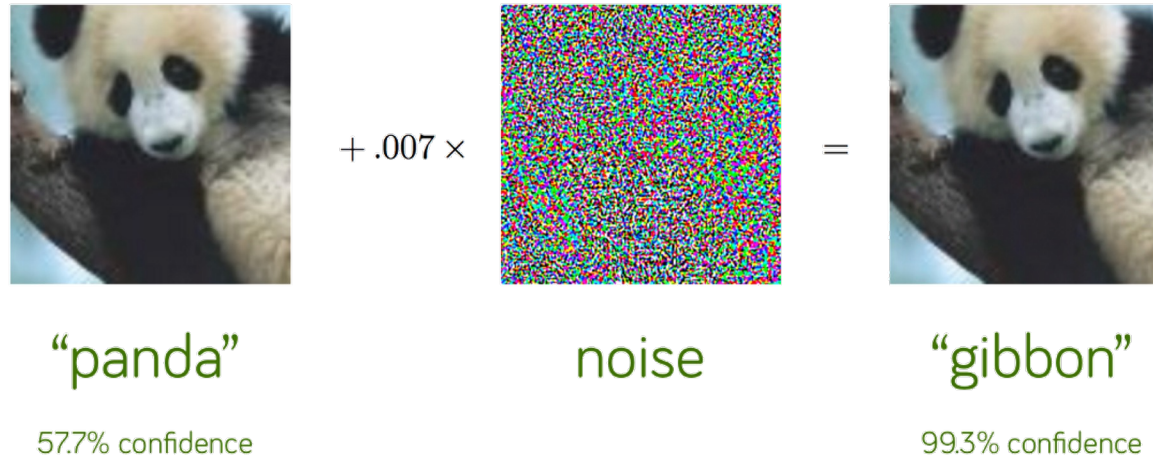
Gradient-based Uncertainty

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



SCAN ME

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

Uncertainty

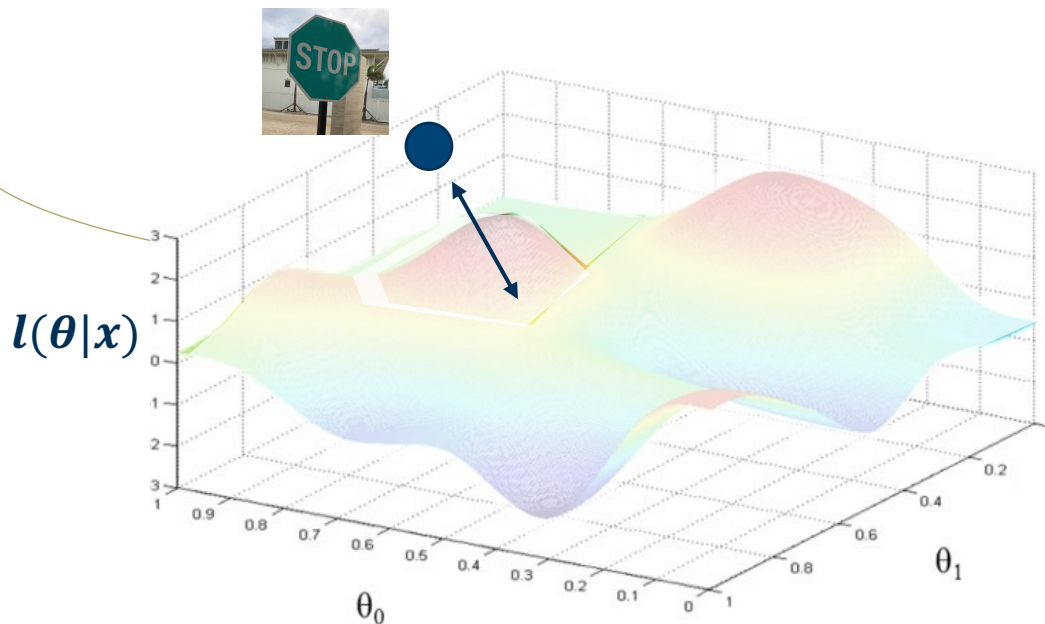
Gradients as Single pass Features

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

Distance from unknown cluster

Two techniques:

1. Backpropagating Confounding labels for Adversarial Detection
2. **Backpropagating Confounding labels for Robust Prediction**





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?



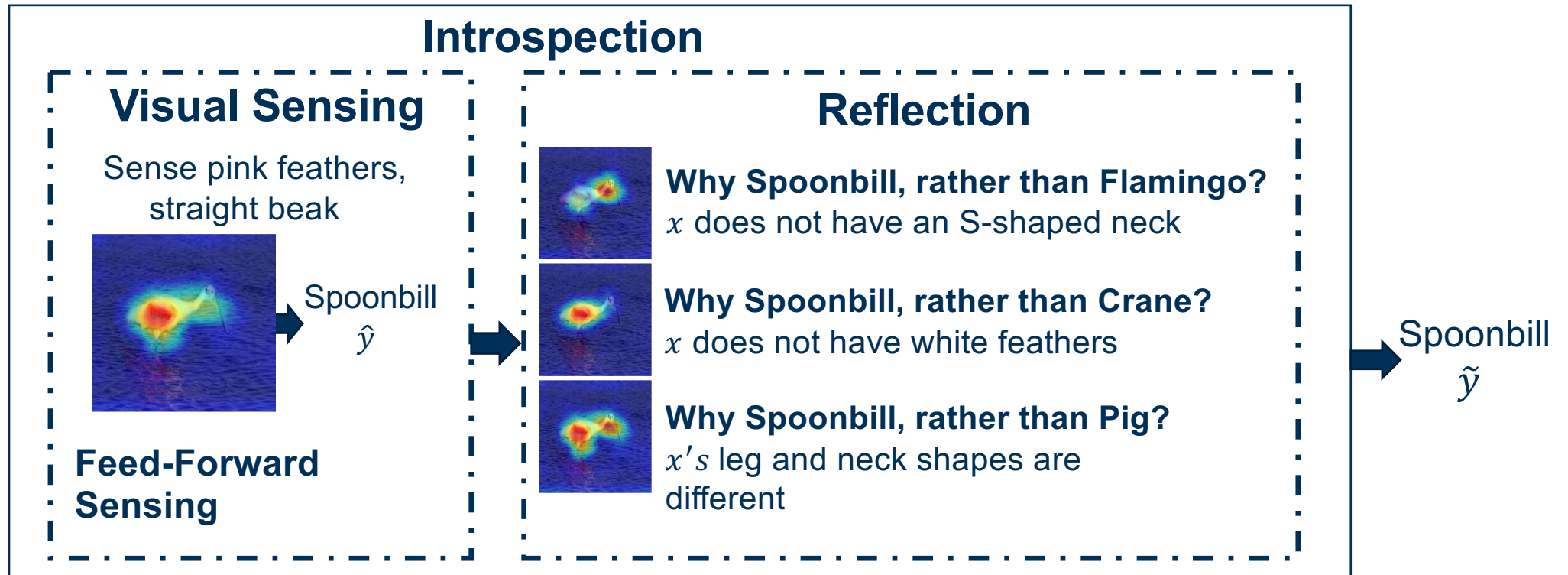
Introspection

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?

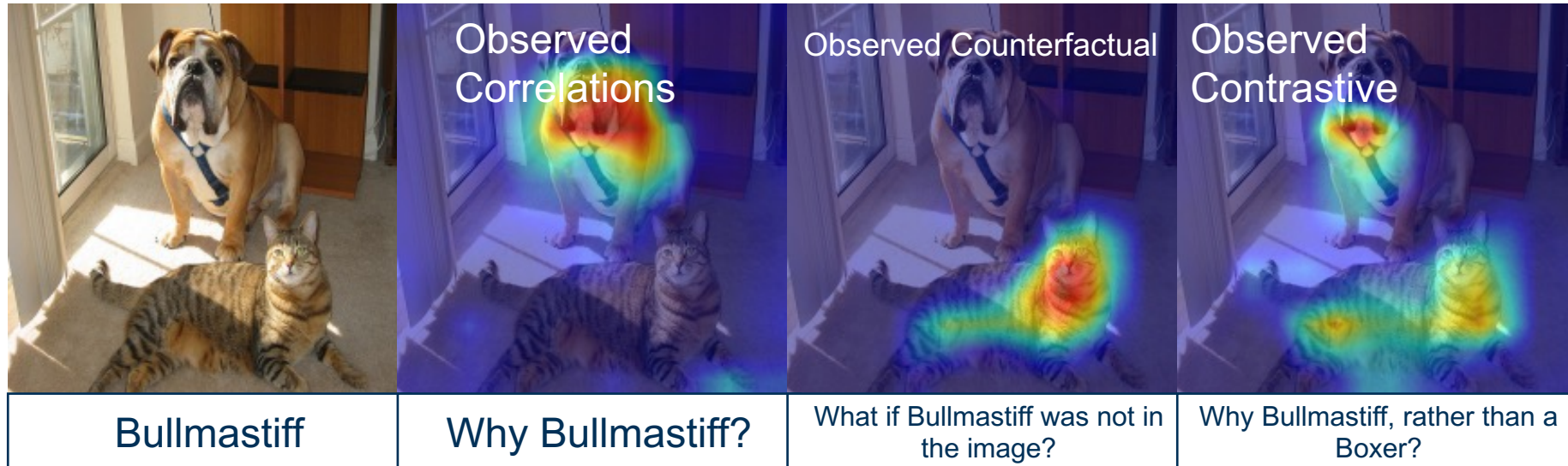
Introspection

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 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

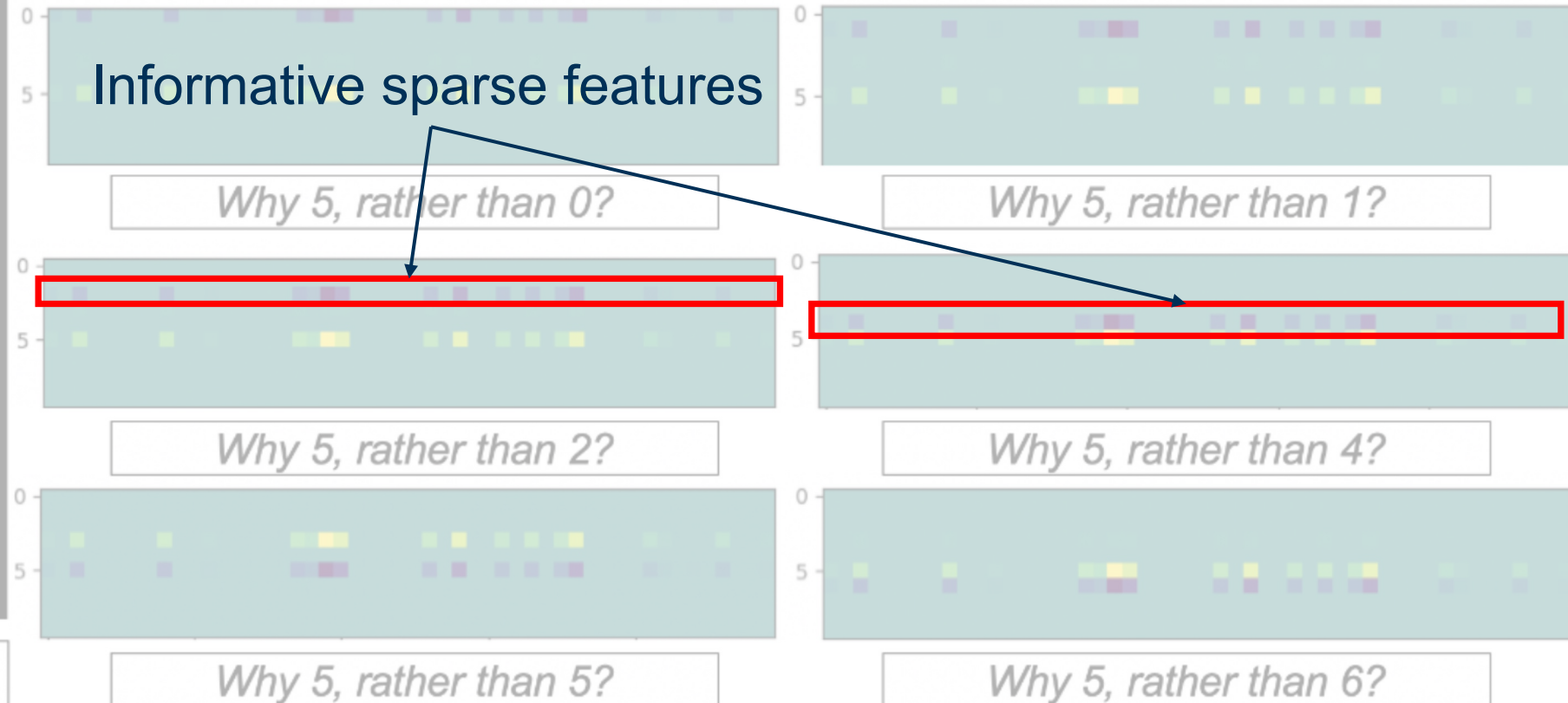


Introspective Learning: A Two-stage Approach for Inference in Neural Networks

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



Input Image x



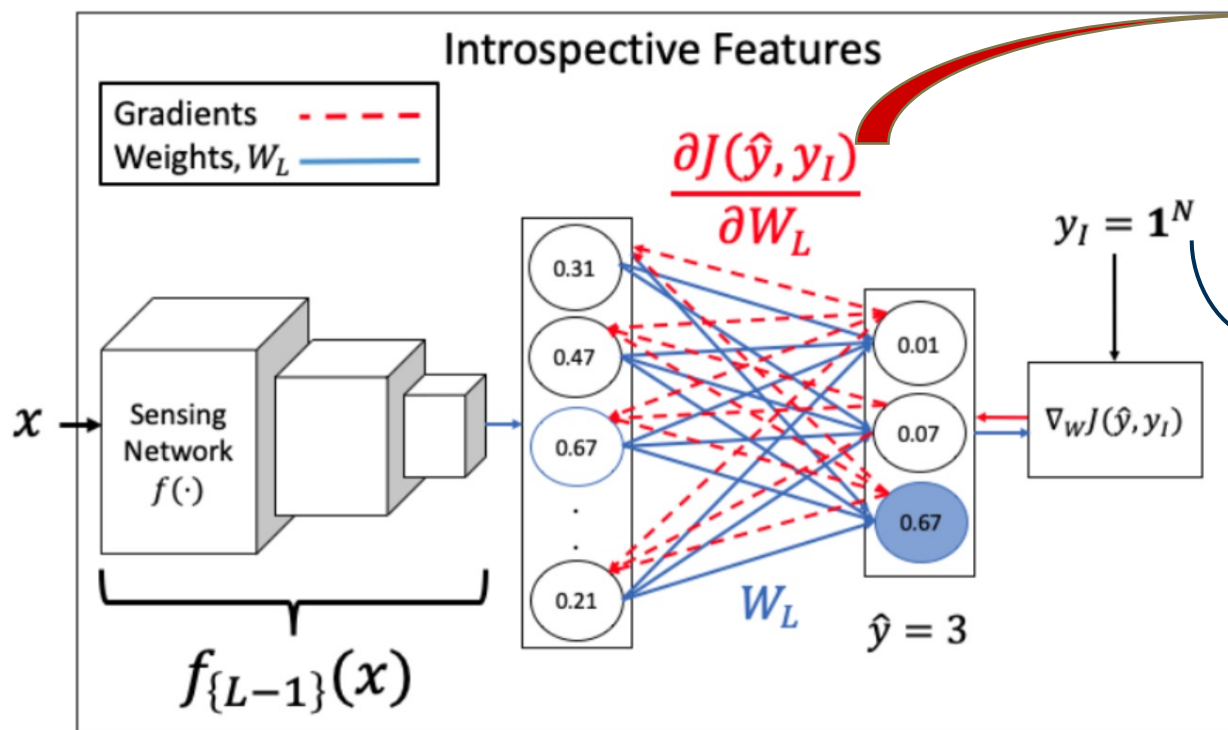
Introspection

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

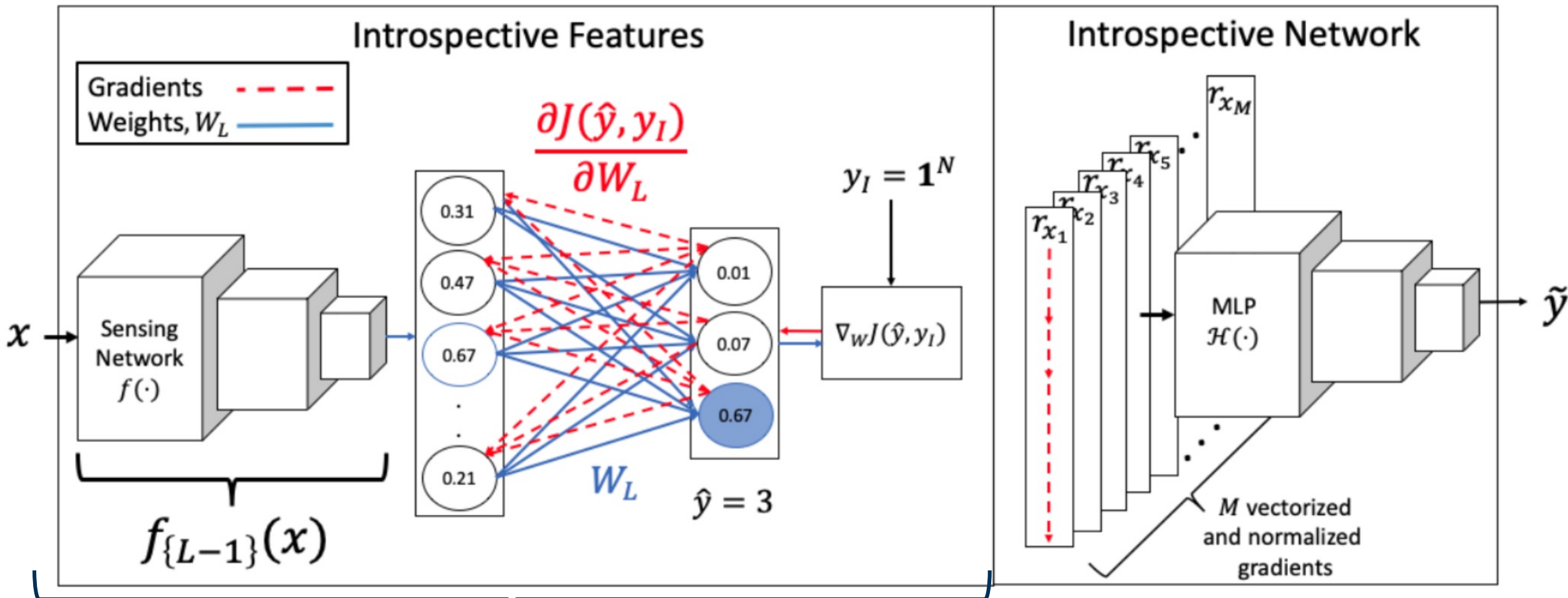
Vector of all ones: A confounding label!

Introspection

Utilizing Gradient Features



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



Introspective Features

Introspection

When is Introspection Useful?



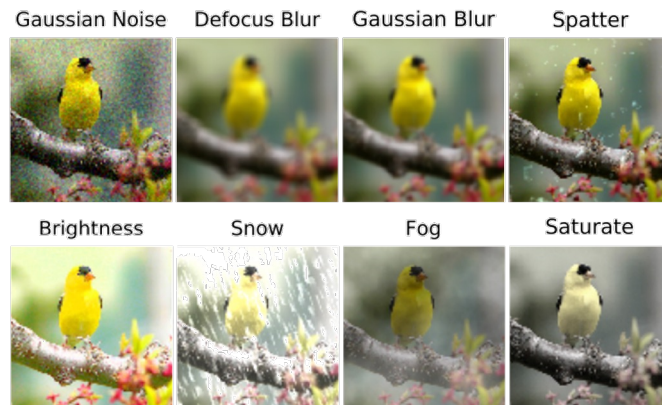
Introspective Learning: A Two-stage Approach for Inference in Neural Networks

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



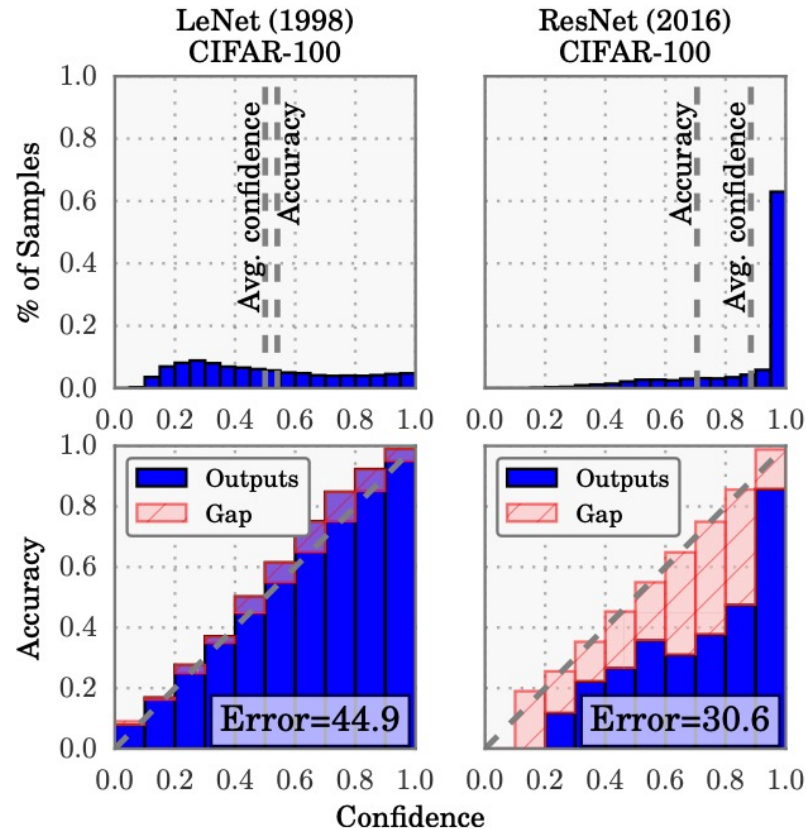
Calibration

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

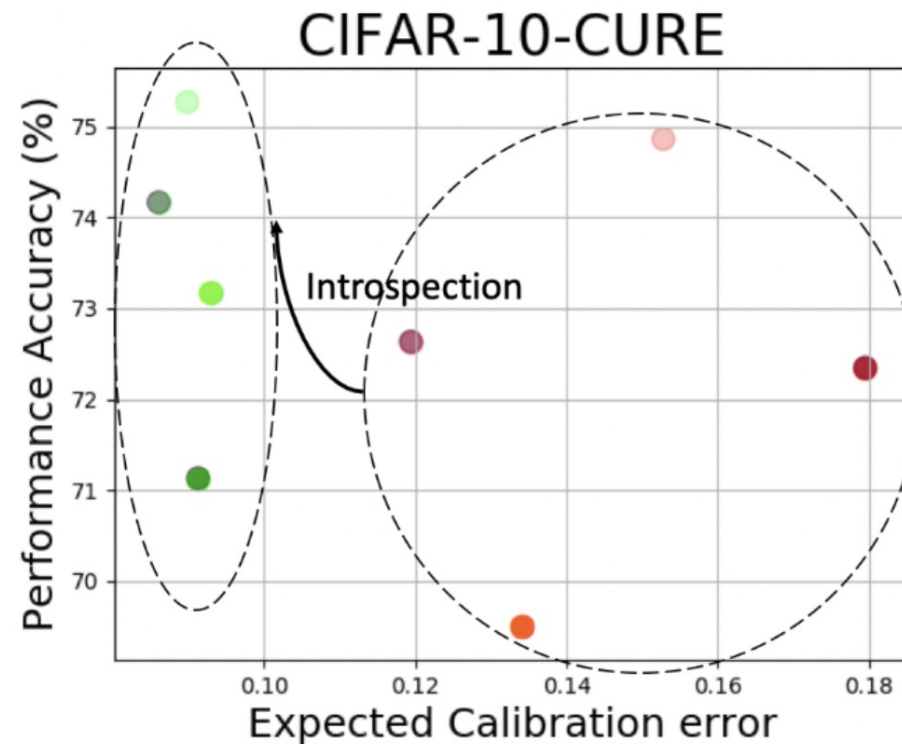
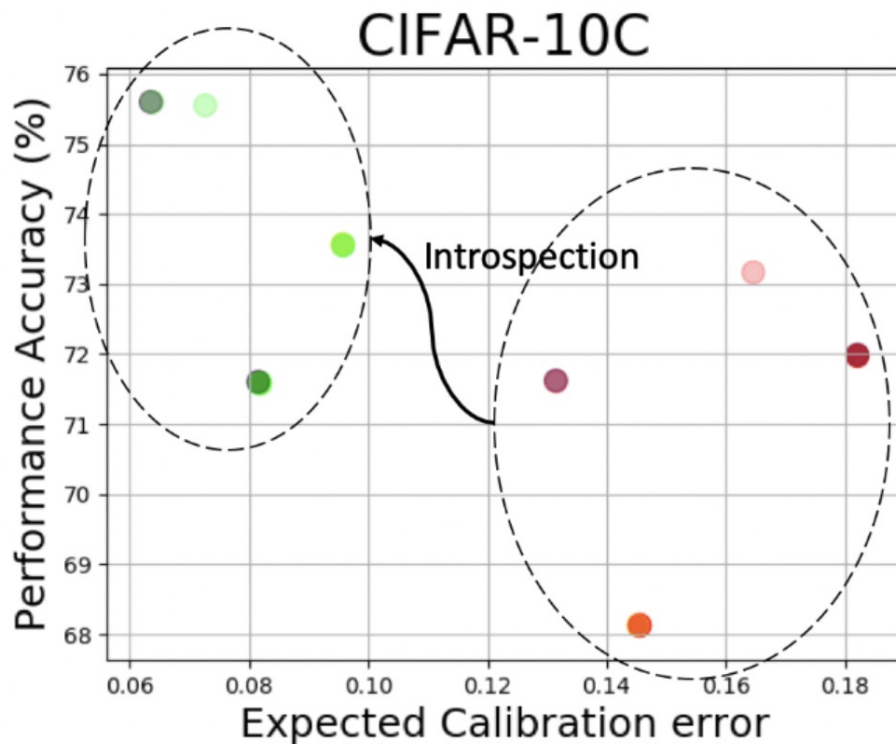


Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Ideal: Top-left corner

Y-Axis: Generalization

X-Axis: Calibration



Legend

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

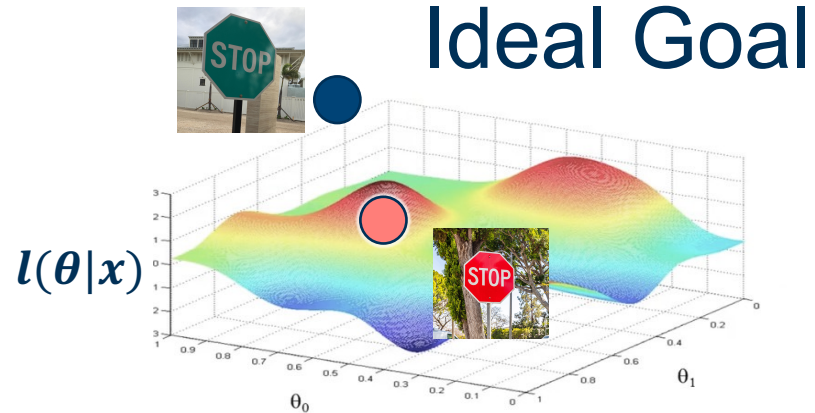
Table 1: Introspecting on top of existing robustness techniques.

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%

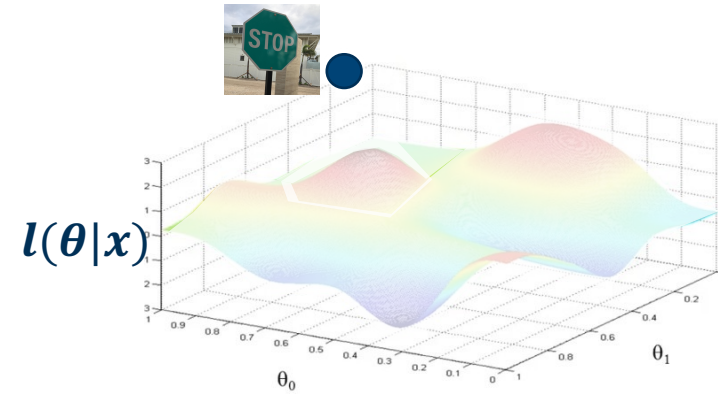
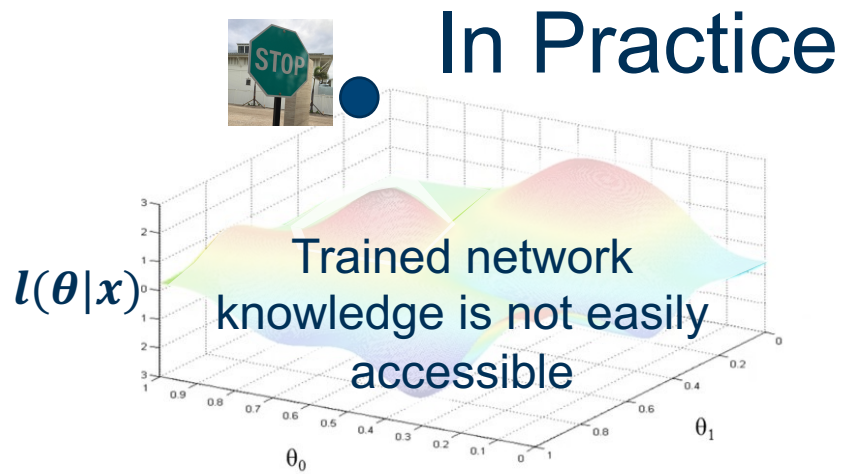
Introspection is a **plug-in approach** that works on all networks and on any downstream task!

Part I, II and III

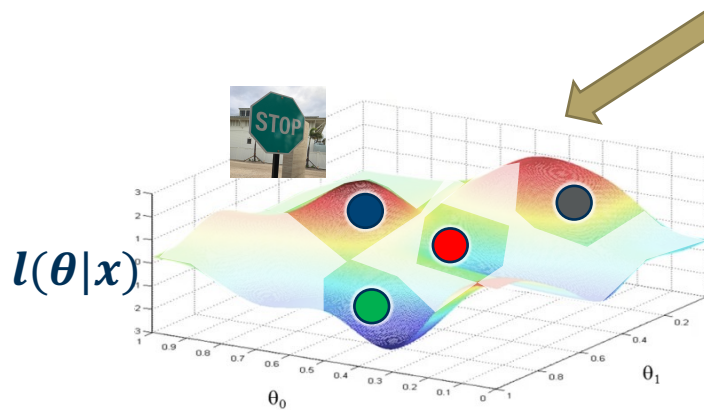
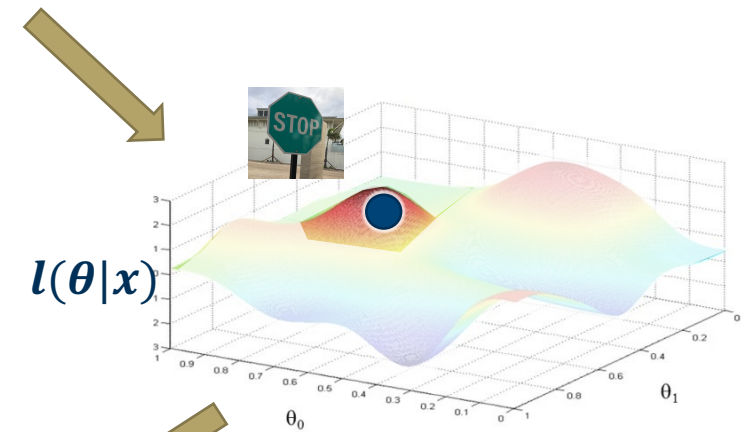
Tying it Back



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