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

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

Feb 21, 2024 – Vancouver, BC, Canada

OLIVES
@GeorgiaTech

Georgia Tech

# Tutorial Materials
## Accessible Online



https://alregib.ece.gatech.edu/aaai-2024-tutorial/

{alregib, mohit.p}@gatech.edu

# AAAI 2024 Tutorial



The 38th Annual AAAI Conference on Artificial Intelligence

FEBRUARY 20-27, 2024 | VANCOUVER, CANADA
VANCOUVER CONVENTION CENTRE – WEST BUILDING

**Presented by:** *Ghassan AlRegib, and Mohit Prabhushankar*

Georgia Institute of Technology

www.ghassanalregib.info

**Duration:** Half Day (3 hours, 30 mins)

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

OLIVES @GeorgiaTech

Georgia Tech

**Expectation vs Reality of Deep Learning**

[Tutorial@AAAI'24] | [Ghassan AlRegib and Mohit Prabhushankar] | [Feb 21, 2024]

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

*- Engineers, probably*

OLIVES
@GeorgiaTech

Georgia Tech

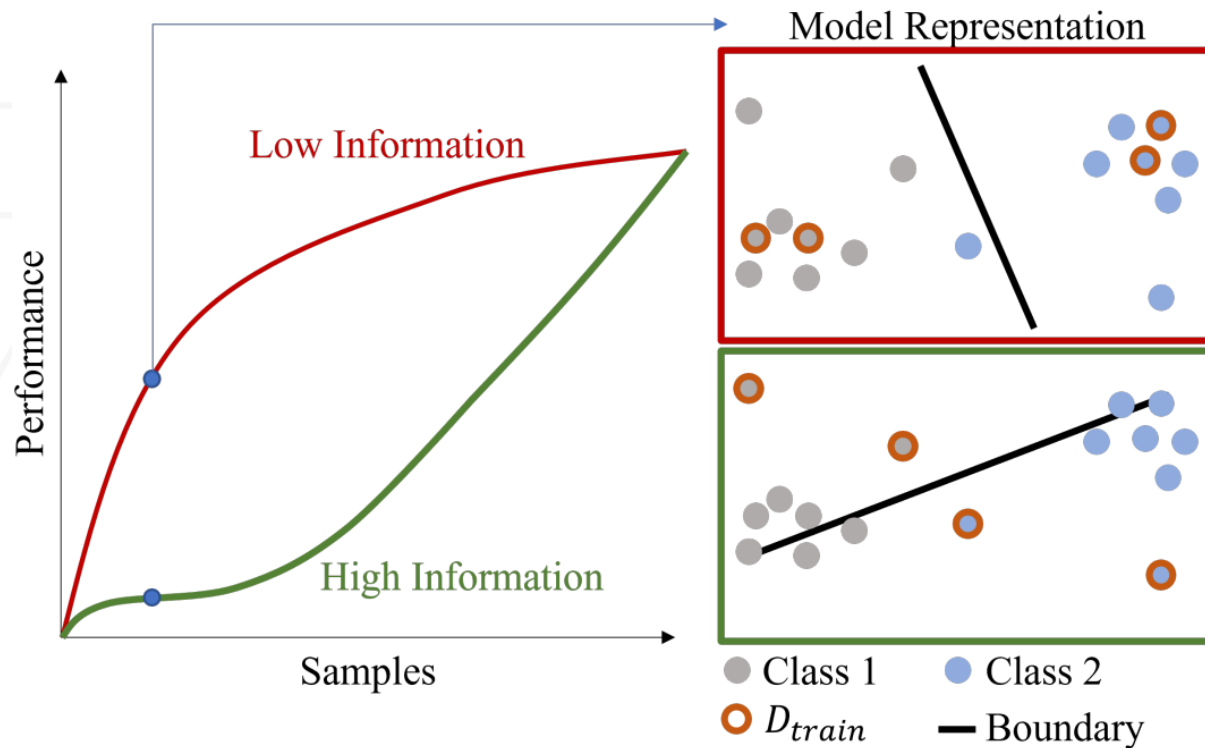## Requirements: Deep Learning-enabled systems must predict correctly on novel data

**Novel** data sources:

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

OLIVES
@GeorgiaTech

Georgia Tech

**The most novel/aberrant samples should <u>not</u> be used in early training**
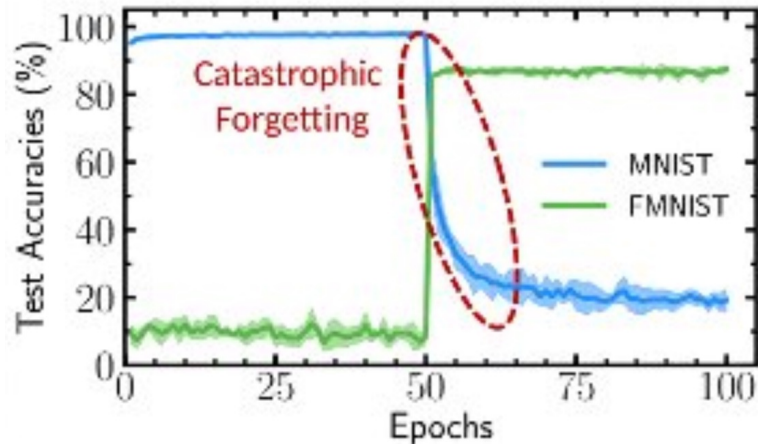


- The first instance of training must occur with less informative samples

- Ex: For autonomous vehicles, less informative means
  - Highway scenarios
  - Parking
  - No accidents
  - No aberrant events

Novel samples = Most Informative

OLIVES
@GeorgiaTech

Georgia Tech

## Subsequent training must <u>not</u> focus only on novel data



- The model performs well on the new scenarios, while forgetting the old scenarios

- A number of techniques exist to overcome this trend

- However, they affect the overall performance in large-scale settings

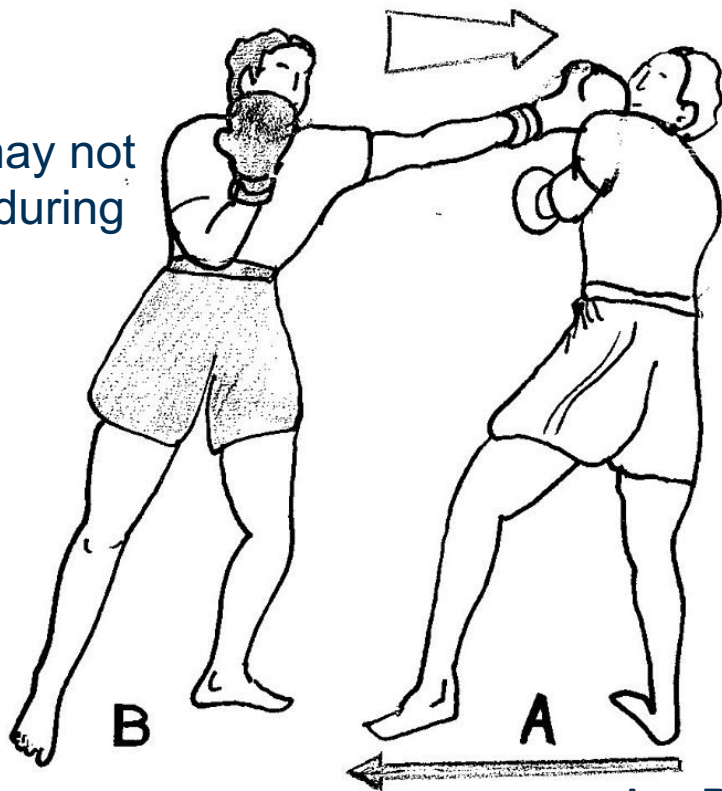- It is not always clear **if and when** to incorporate novel scenarios in training

[Tutorial@AAAI'24] | [Ghassan AlRegib and Mohit Prabhushankar] | [Feb 21, 2024]

Laborieux, Axel, et al. "Synaptic metaplasticity in binarized neural networks." *Nature communications* 12.1 (2021): 2549.

**Novel data packs a 1-2 punch!**



Novel data may not be available during training

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

A = Deep Neural Networks
B = Novel data

OLIVES
@GeorgiaTech

Georgia Tech

**We must handle novel data at Inference!!**

**Novel** data sources:

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

## Model Train



## At Inference

OLIVES
@GeorgiaTech

Georgia Tech

# Objective
## Objective of the Tutorial

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

- Part 1: Inference in Neural Networks

- Part 2: Explainability at Inference

- Part 3: Uncertainty at Inference

- Part 4: Intervenability at Inference

- Part 5: Conclusions and Future Directions

[Tutorial@AAAI'24] | [Ghassan AlRegib and Mohit Prabhushankar] | [Feb 21, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

# Robust Neural Networks
# Part I: Inference in Neural Networks
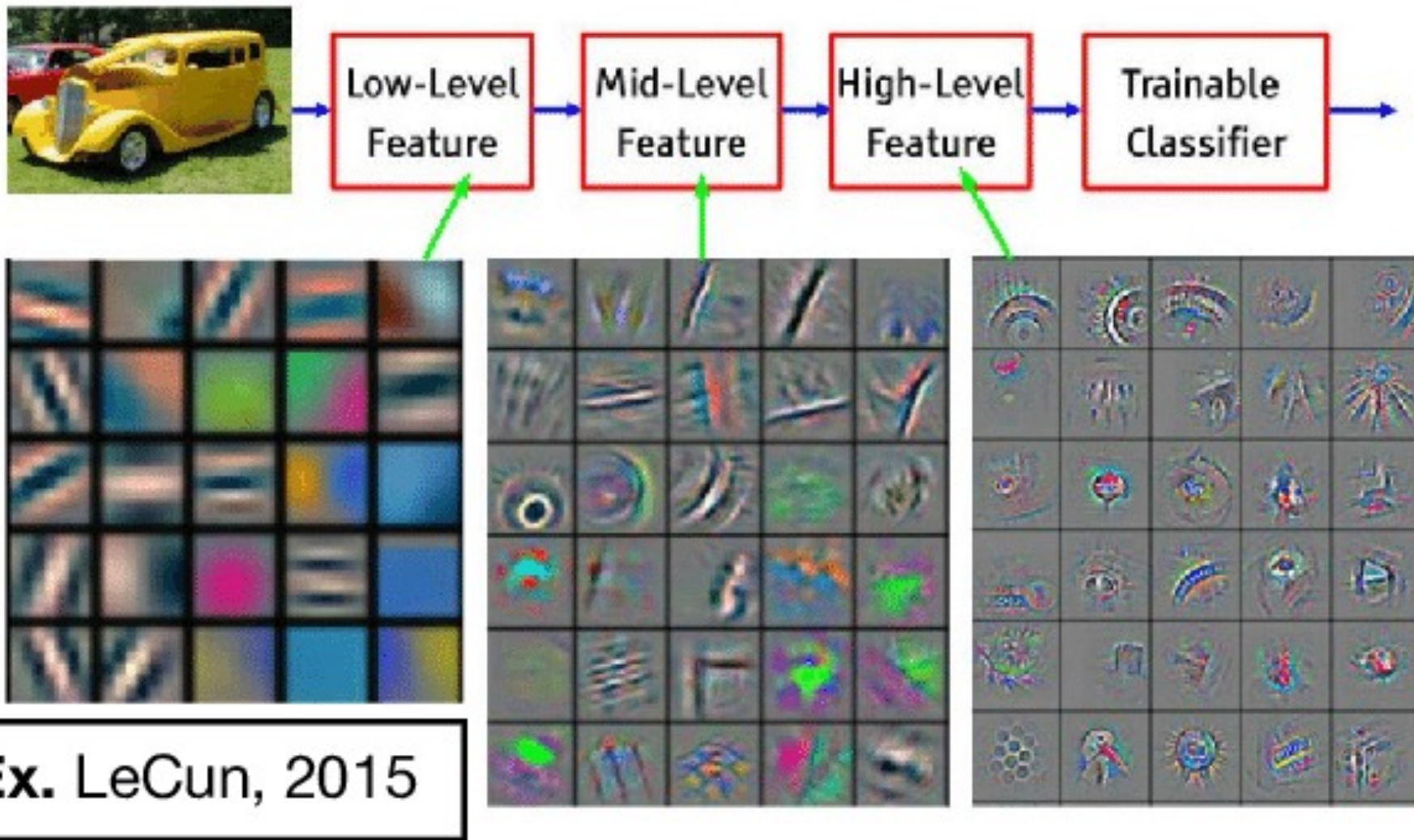
# Objective
## Objective of the Tutorial

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

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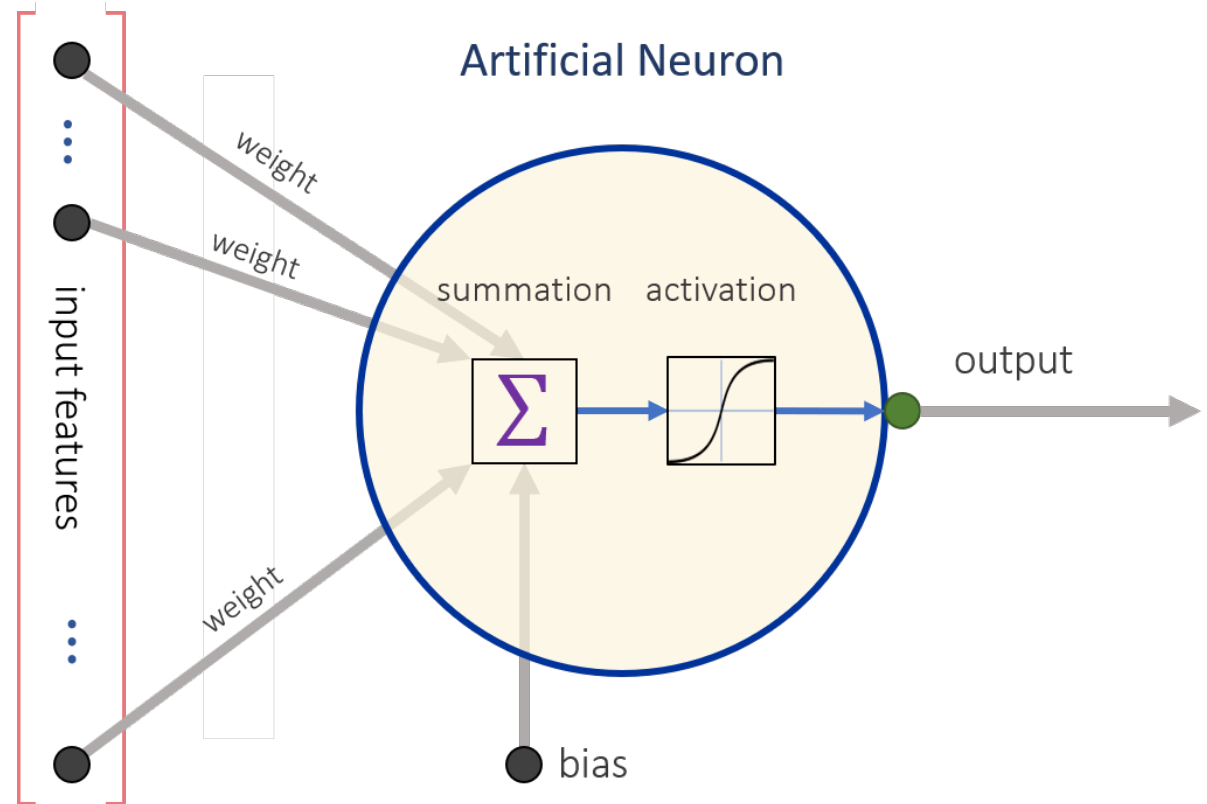- Part 5: Conclusions and Future Directions

OLIVES
@GeorgiaTech

Georgia Tech

Ex. LeCun, 2015

**The underlying computation unit is the Neuron**

Artificial neurons consist of:

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



Artificial Neuron

input features

weight

weight

weight

summation    activation

$\Sigma$

output

bias

OLIVES
@GeorgiaTech

Georgia Tech

**Neurons are stacked and densely connected to construct ANNs**



input layer

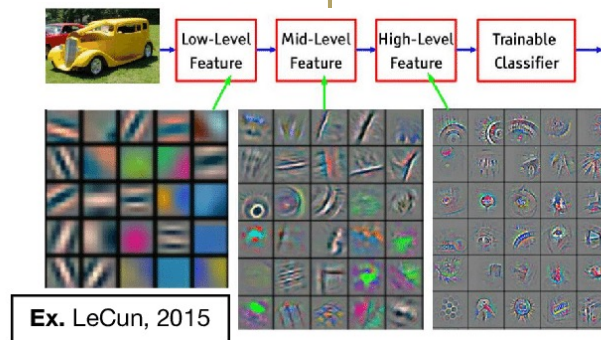hidden layers (optional)
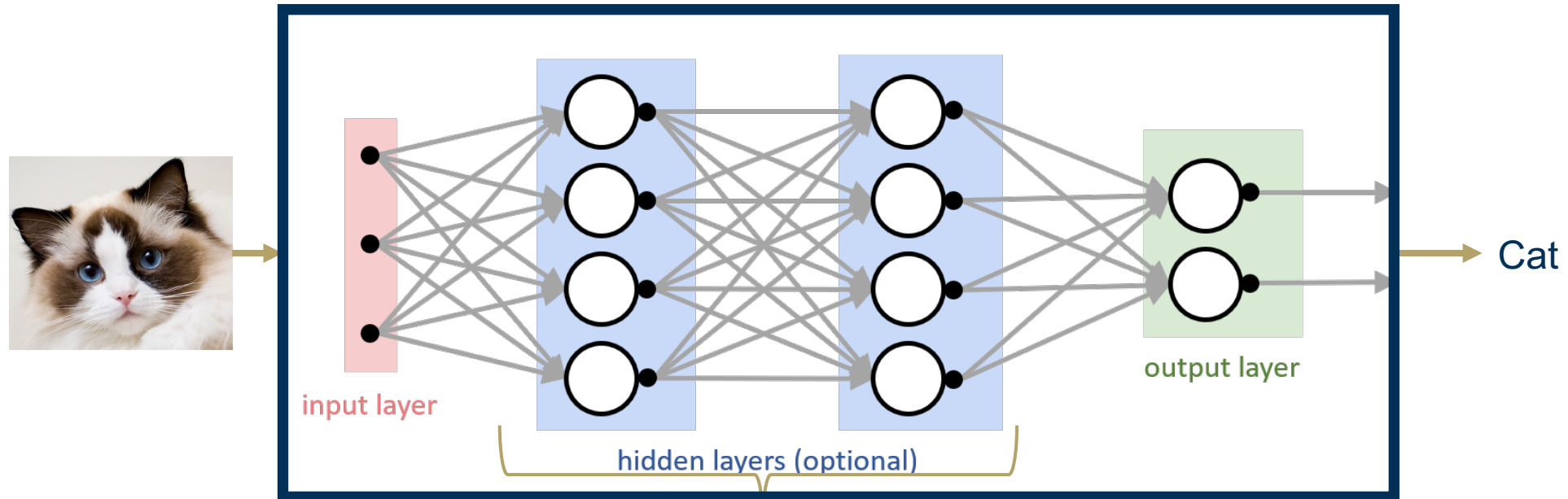
output layer

Cat

Typically, a neuron is part of a network organized in layers:
- An input layer (Layer $0$)
- An output layer (Layer $K$)
- Zero or more hidden (middle) layers (Layers $1 \ldots K-1$)

OLIVES
@GeorgiaTech

Georgia Tech

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



input layer

hidden layers (optional)

output layer

Cat

**Ex.** LeCun, 2015

Low-Level Feature — Mid-Level Feature — High-Level Feature — Trainable Classifier

[Tutorial@AAAI'24] | [Ghassan AlRegib and Mohit Prabhushankar] | [Feb 21, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

**Transformers, Large Language Models and Foundation Models**



Primary reasons for advancements:
1. Expanded interests from the research community
2. Computational resources availability
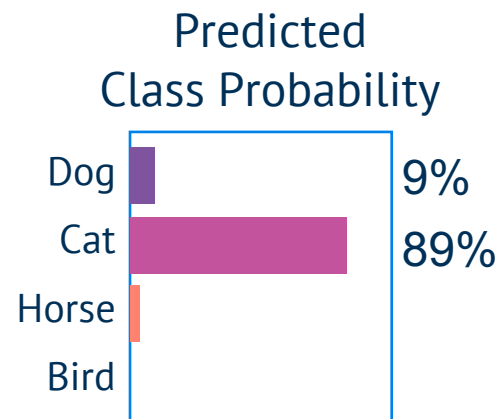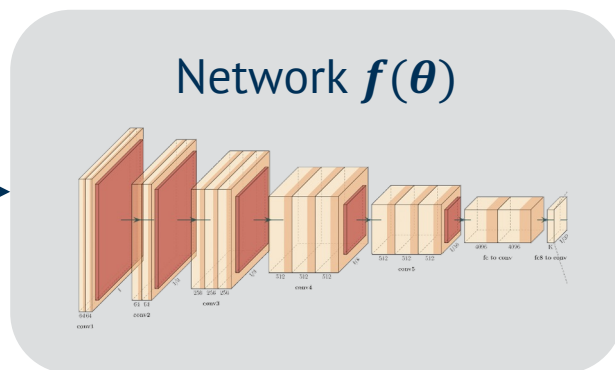3. **Big data availability**

Classification

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



$x$

Predicted
Class Probability

| | |
|---|---|
| Dog | 9% |
| Cat | 89% |
| Horse | |
| Bird | |

If $x \in \chi$, the data is **not novel**
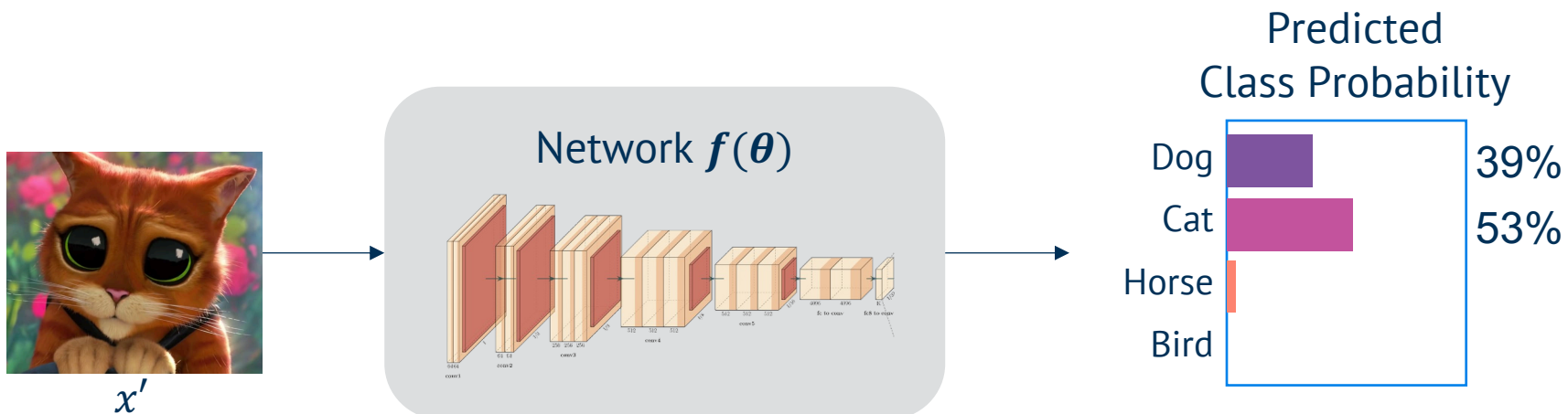
$\hat{y} = f(x)$
$y = argmax_i \ \hat{y}$
$p(\hat{y}) = T(f(x))$

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

OLIVES
@GeorgiaTech

Georgia
Tech

**Deep learning robustness: Correctly predict class even when data is <u>novel</u>**

Network $f(\boldsymbol{\theta})$

$x'$

Predicted
Class Probability

Dog — 39%
Cat — 53%
Horse
Bird

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

$$\hat{y} = f(x' + \epsilon)$$
$$y = argmax_i \, \hat{y}$$
$$p(\hat{y}) = T(f(x' + \epsilon))$$

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

OLIVES
@GeorgiaTech

Georgia Tech

**Deep learning robustness: Correctly predict class even when data is <u>novel</u>**



Predicted
Class Probability

| | |
|---|---|
| Dog | 39% |
| Cat | 53% |
| Horse | |
| Bird | |

Network $f(\theta)$

$x'$

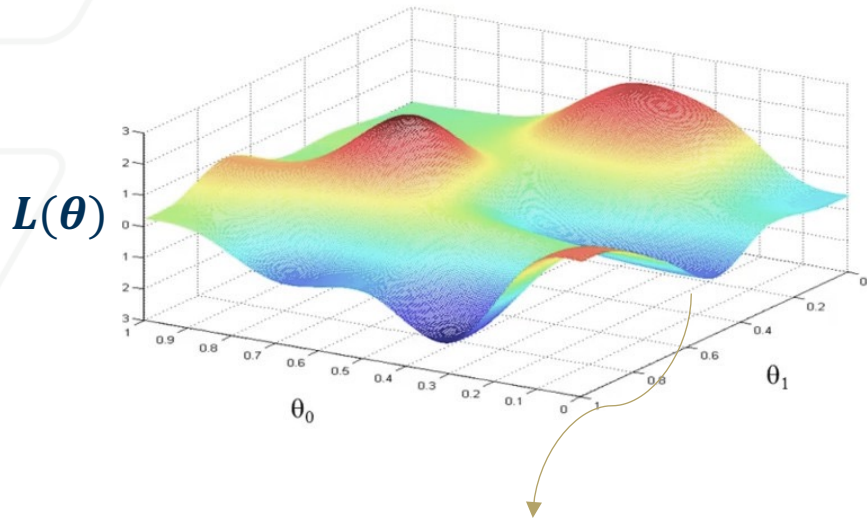To achieve robustness at Inference, we need the following:

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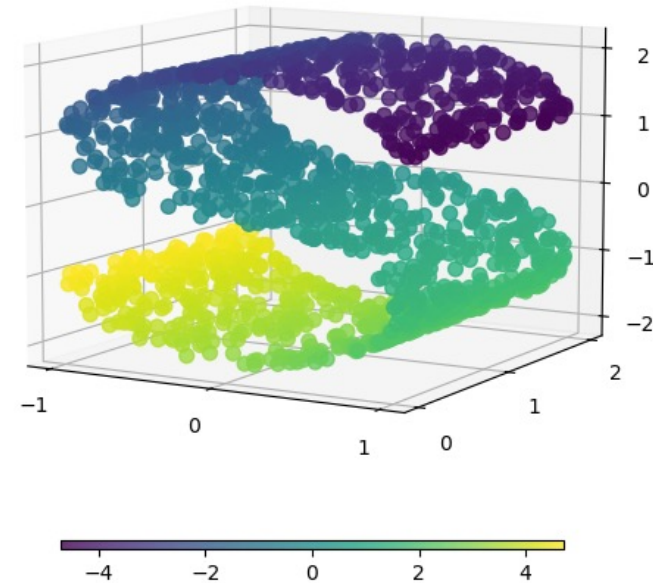- **Techniques** that utilize the information from novel data

**Why is this Challenging?**

OLIVES
@GeorgiaTech

Georgia Tech

**Manifolds are compact topological spaces that allow exact mathematical functions**



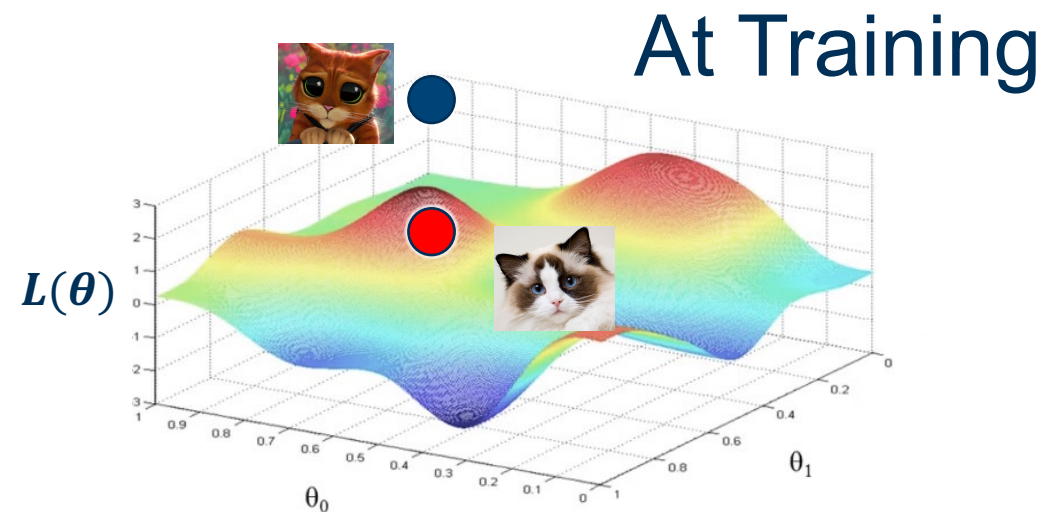Toy visualizations generated using functions (and thousands of generated data points)

Real data visualizations generated using dimensionality reduction algorithms (Isomap)

[Tutorial@AAAI'24] | [Ghassan AlRegib and Mohit Prabhushankar] | [Feb 21, 2024]

OLIVES
@GeorgiaTech

Georgia Tech

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

## At Inference



$L(\theta)$

Trained network knowledge is not easily accessible

## At Training



$L(\theta)$

At training, we have access to all training data.

OLIVES
@GeorgiaTech

Georgia Tech

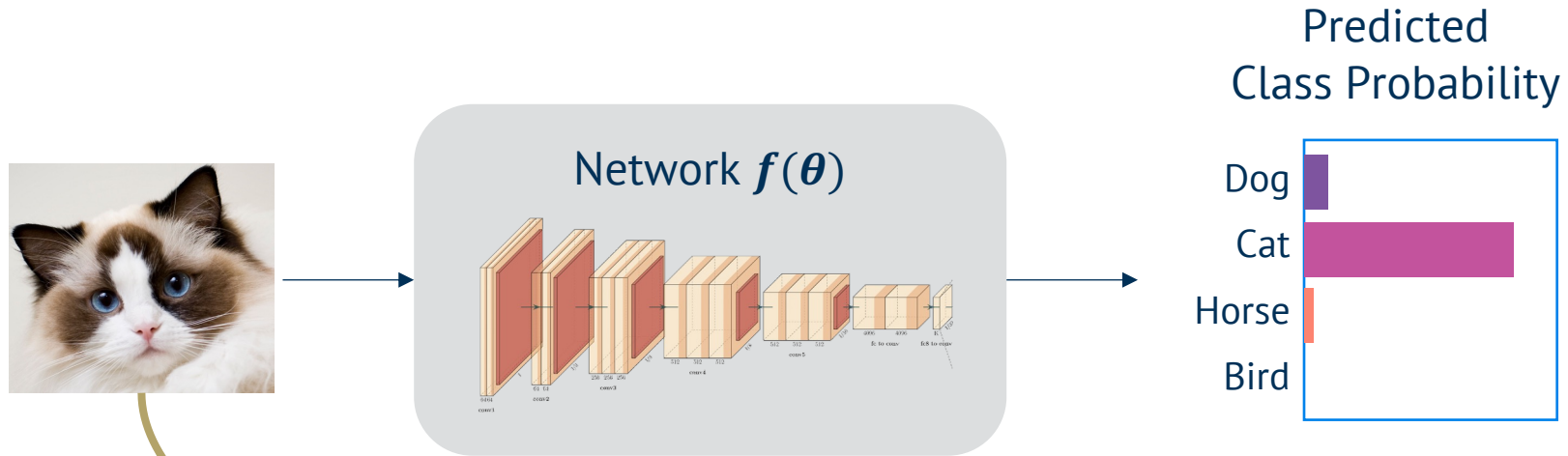**Colloquially, Fisher Information is the "surprise" in a system that observes an event**

Predicted
Class Probability

Network $f(\boldsymbol{\theta})$



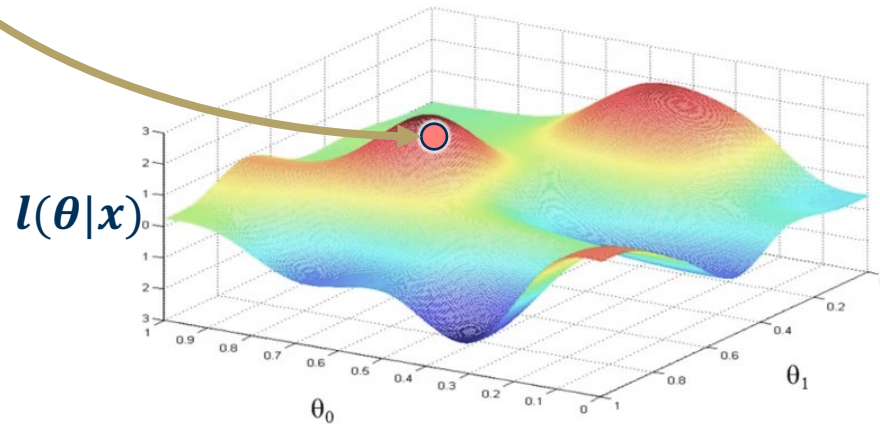$l(\boldsymbol{\theta}|\boldsymbol{x})$

Likelihood function

**Fisher Information**

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

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

[Tutorial@AAAI'24] | [Ghassan AlRegib and Mohit Prabhushankar] | [Feb 21, 2024]

[1] A good blogpost about Fisher Information: https://towardsdatascience.com/an-intuitive-look-at-fisher-information-2720c40867d8

OLIVES
@GeorgiaTech

Georgia
Tech

Predicted Class Probability

Dog
Cat
Horse
Bird

Network $f(\theta)$

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

$l(\theta|x)$

Likelihood function

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

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

OLIVES
@GeorgiaTech
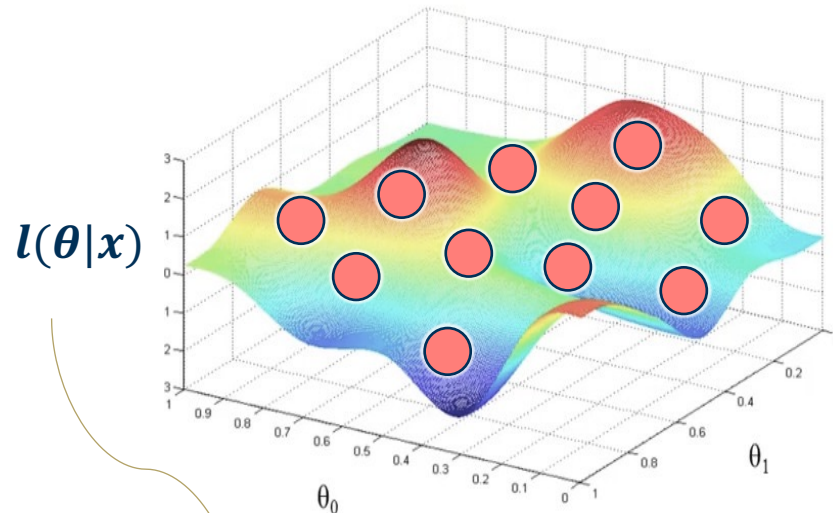
Georgia Tech

**Gradients infer information about the statistics of underlying manifolds**

$l(\theta|x)$



Likelihood function instead of loss manifold

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

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

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

$E[\cdot] = $ Expectation
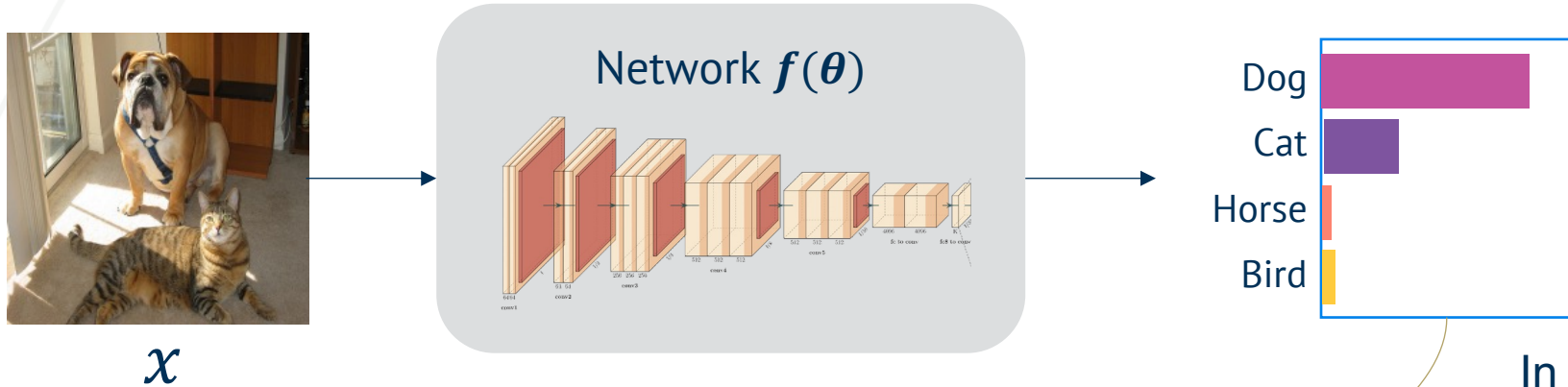$U_\theta = \nabla_\theta l(\theta|x)$, Gradients w.r.t. the sample

**Hence, gradients draw information from the underlying distribution as learned by the network weights!**
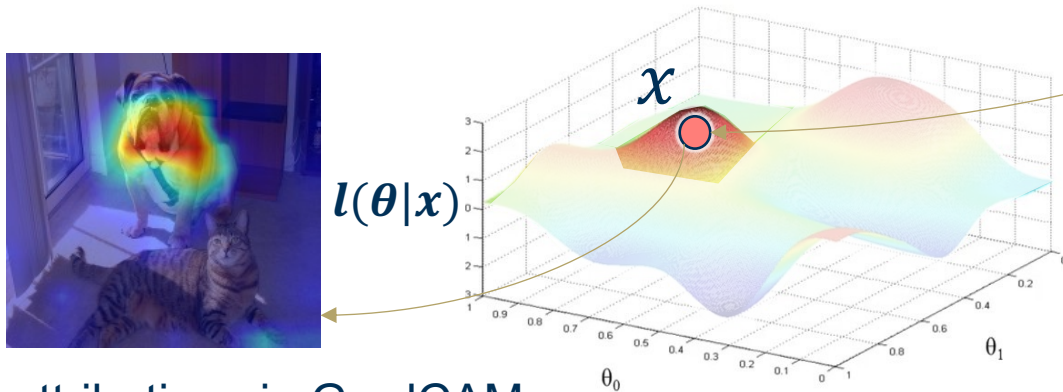
OLIVES
@GeorgiaTech

Georgia Tech

Case Study: Gradients as Fisher Information in Explainability

**Gradients infer information about the statistics of underlying manifolds**



Network $f(\boldsymbol{\theta})$

Dog
Cat
Horse
Bird

$x$

Local information (specific to $x$) is sufficient!

In this case, the image and its prediction extracts nose, mouth and jowl features.
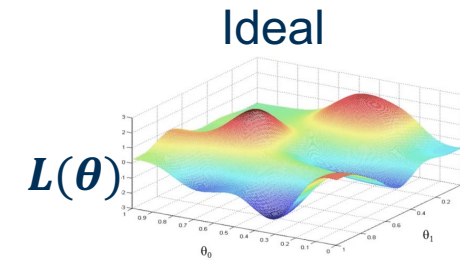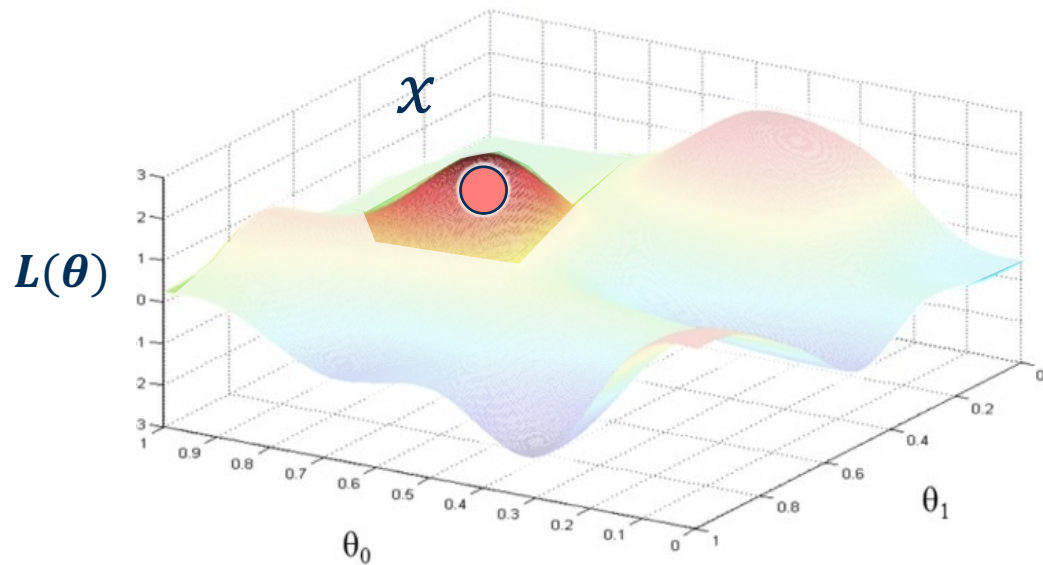
$l(\boldsymbol{\theta}|x)$

$x$

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

Feature attribution via GradCAM

OLIVES @GeorgiaTech

Georgia Tech

[1] A good blogpost about Fisher Information: https://towardsdatascience.com/an-intuitive-look-at-fisher-information-2720c40867d8

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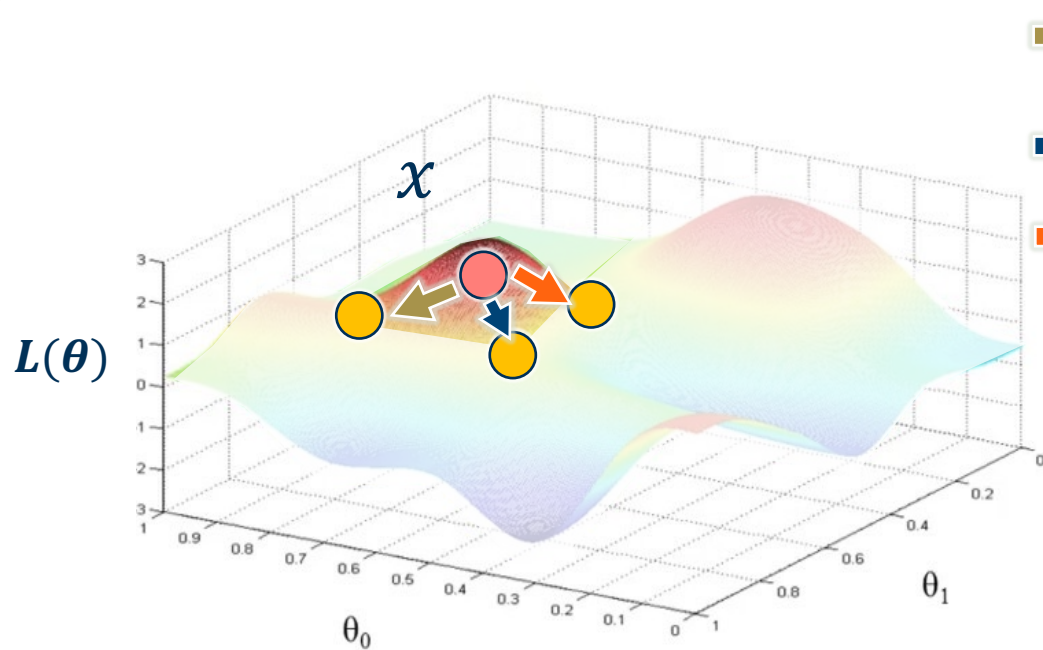


$$x$$

$$L(\theta)$$

Ideal

$$L(\theta)$$

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

OLIVES
@GeorgiaTech

Georgia Tech

# Gradients at Inference
## Direction of Steepest Descent

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



Path 1?

Path 2?

Path 3?

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

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

OLIVES
@GeorgiaTech

Georgia Tech

## To Characterize the Novel Data at Inference



At Inference

Trained network knowledge is not easily accessible

$L(\theta)$

Part 2, 3

Part 4

Counterfactual and Contrastive Representations using Gradients

$x'$

$x$

$L(\theta)$

Representation Traversal using Interventions

$L(\theta)$

[Tutorial@AAAI'24] | [Ghassan AlRegib and Mohit Prabhushankar] | [Feb 21, 2024]

OLIVES
@GeorgiaTech

Georgia Tech