

# Robustness at Inference: Towards 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  
Jan 07, 2024 – Waikaloa, HI, USA



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

## WACV 2024 Tutorial

# Robustness at Inference: Towards Explainability, Uncertainty, and Intervenability

*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/>

**Duration:** Half-Day event

## Expectation vs Reality of Deep Learning





# Deep Learning

## Expectation vs Reality

### LATEST TRICKS

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

Stop



Dumb-bell

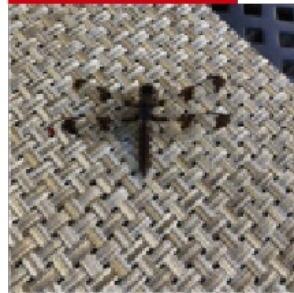


Racket

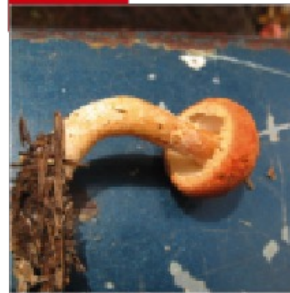


Even natural images can fool a DNN, because it might focus on the picture's colour, texture or background rather than picking out the salient features a human would recognize.

Manhole cover



Pretzel



©nature



# Deep Learning

## Expectation vs Reality

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

*- Engineers, probably*





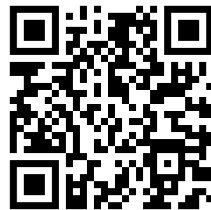
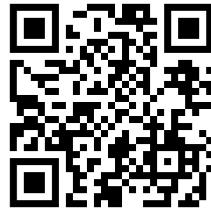
# Deep Learning

## Requirements and Challenges

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

**Novel data sources:**

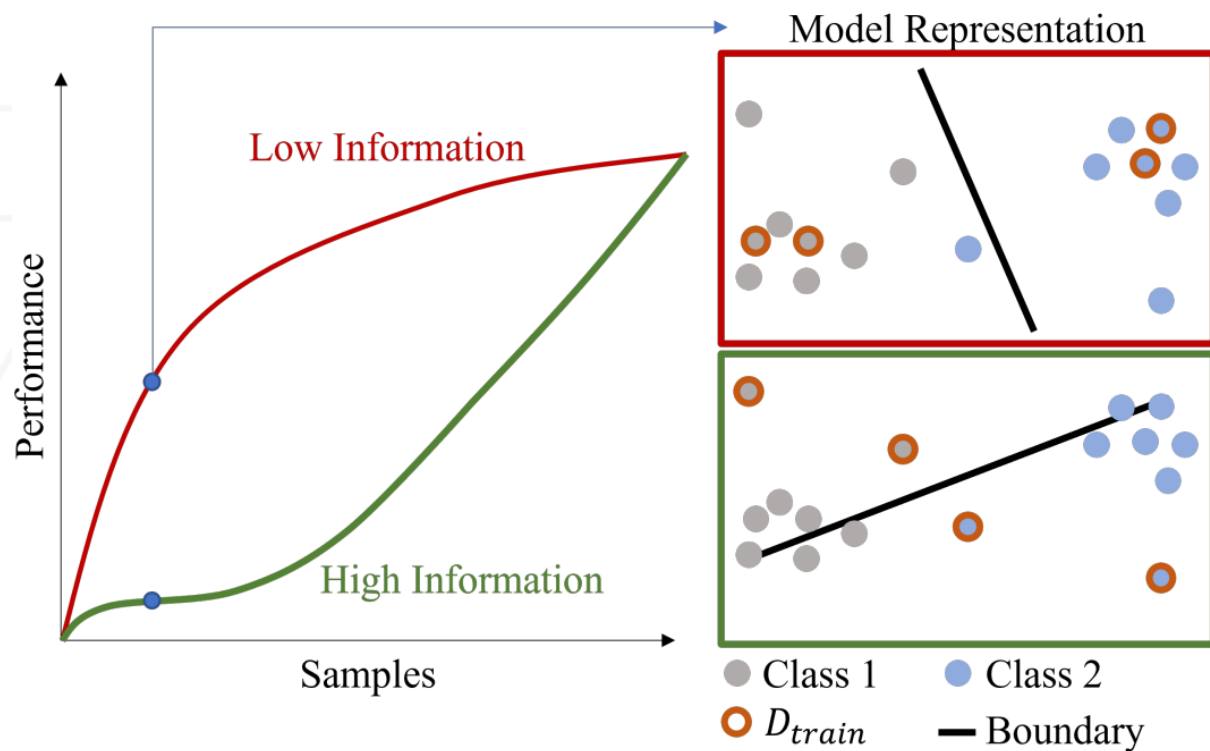
- Test distributions
- Anomalous data
- Out-Of-Distribution data
- Adversarial data
- Corrupted data
- Noisy data
- New classes
- ...



# Deep Learning at Training

## Overcoming Challenges at Training: Part 1

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



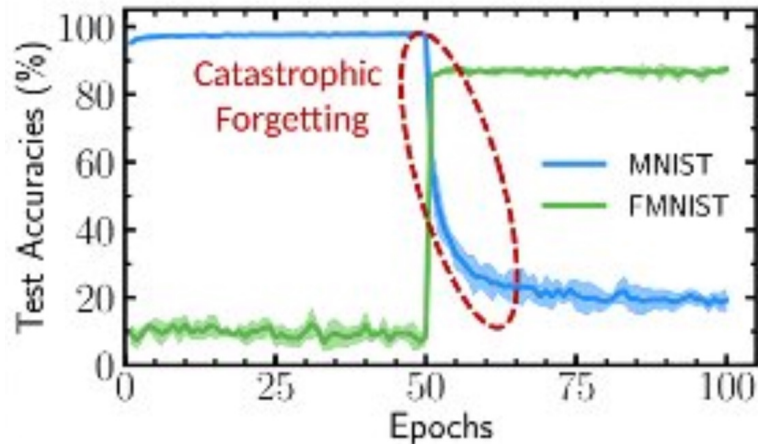
- The first instance of training must occur with less informative samples
- Ex: For autonomous vehicles, less informative means
  - Highway scenarios
  - Parking
  - No accidents
  - No aberrant events

Novel samples = Most Informative

# Deep Learning at Training

## Overcoming Challenges at Training: Part 2

Subsequent training must not focus only on novel data



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

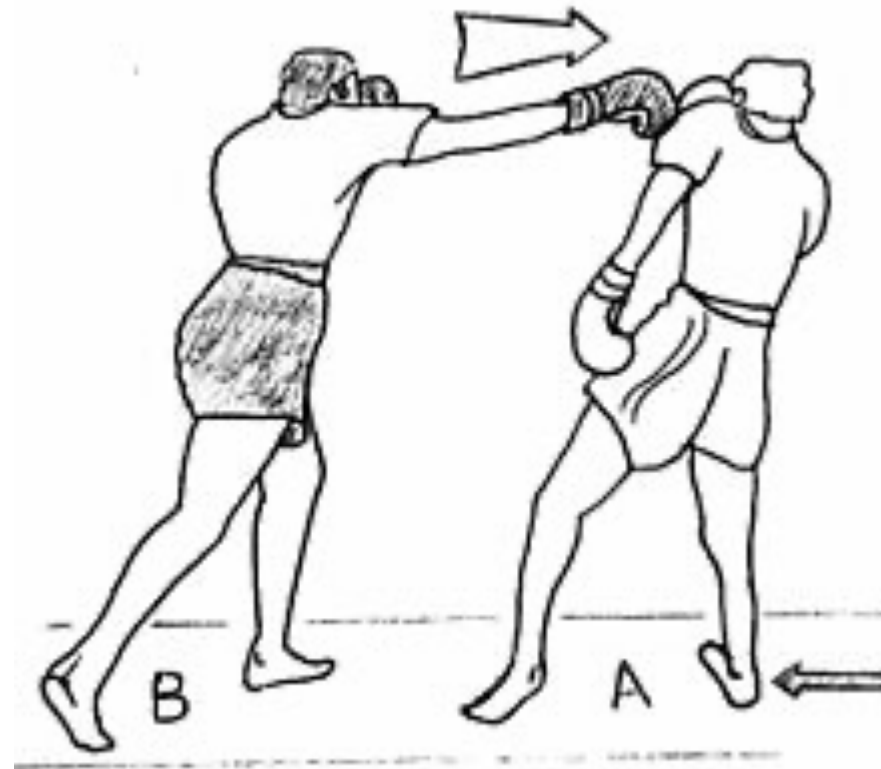
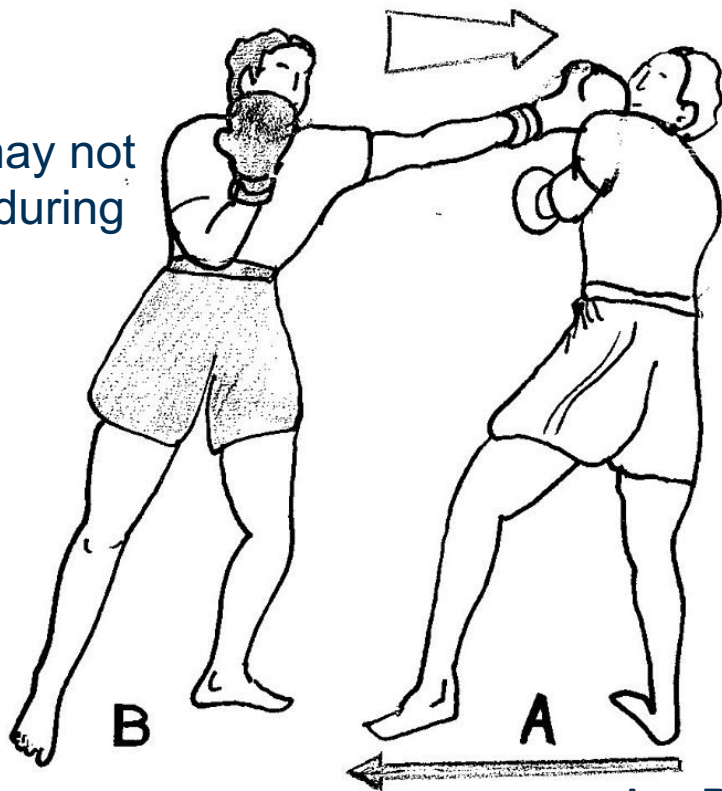


# Deep Learning at Training

## Overcoming Challenges at Training

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

Novel data may not be available during training



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

# Deep Learning at Inference

## Overcoming Challenges at Inference

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

**Novel** data sources:

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

Model Train



At Inference



# Objective

## Objective of the Tutorial

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

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



# Robust Neural Networks

## Part I: Inference in Neural Networks

# Objective

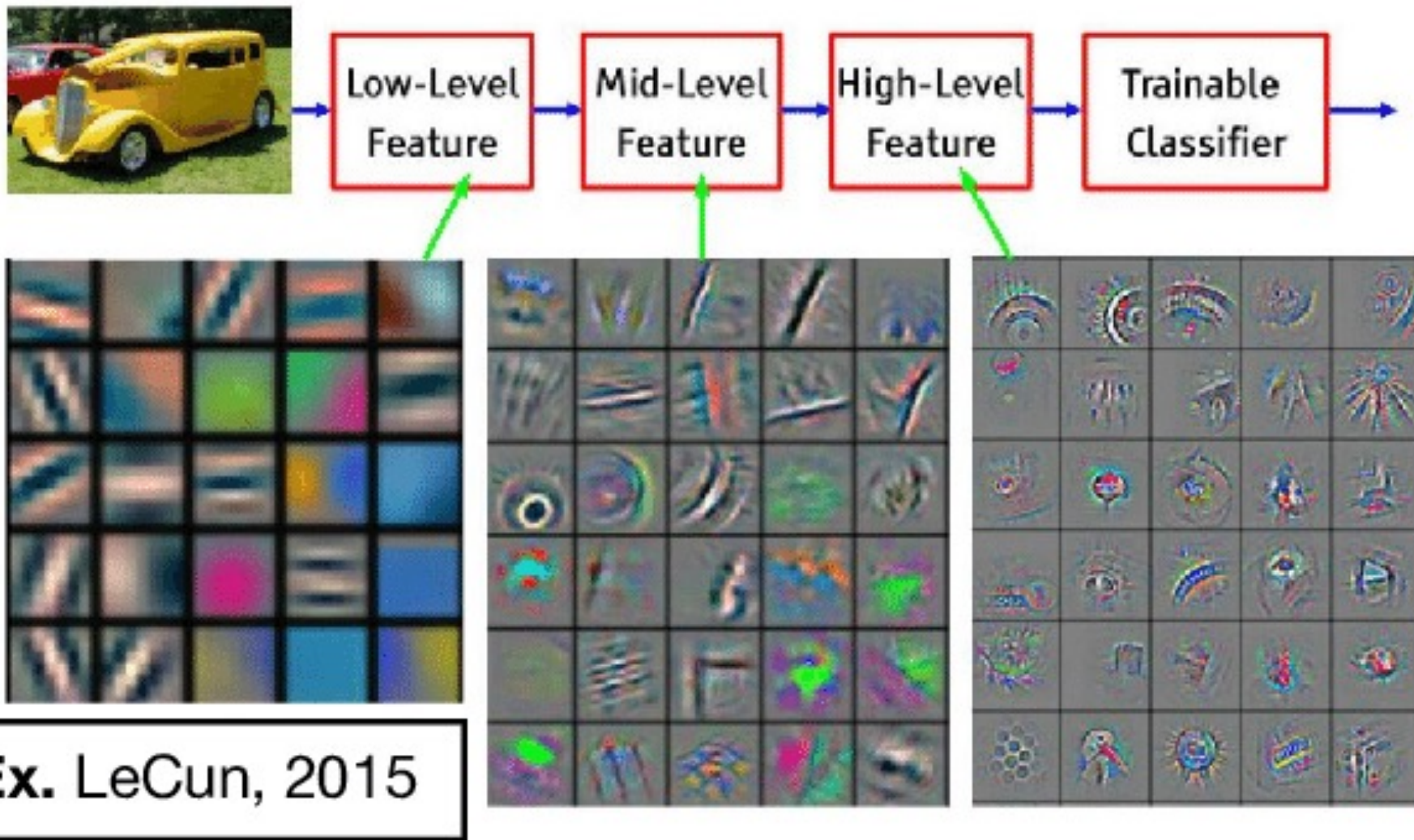
## Objective of the Tutorial

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

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

# Deep Learning

## Overview





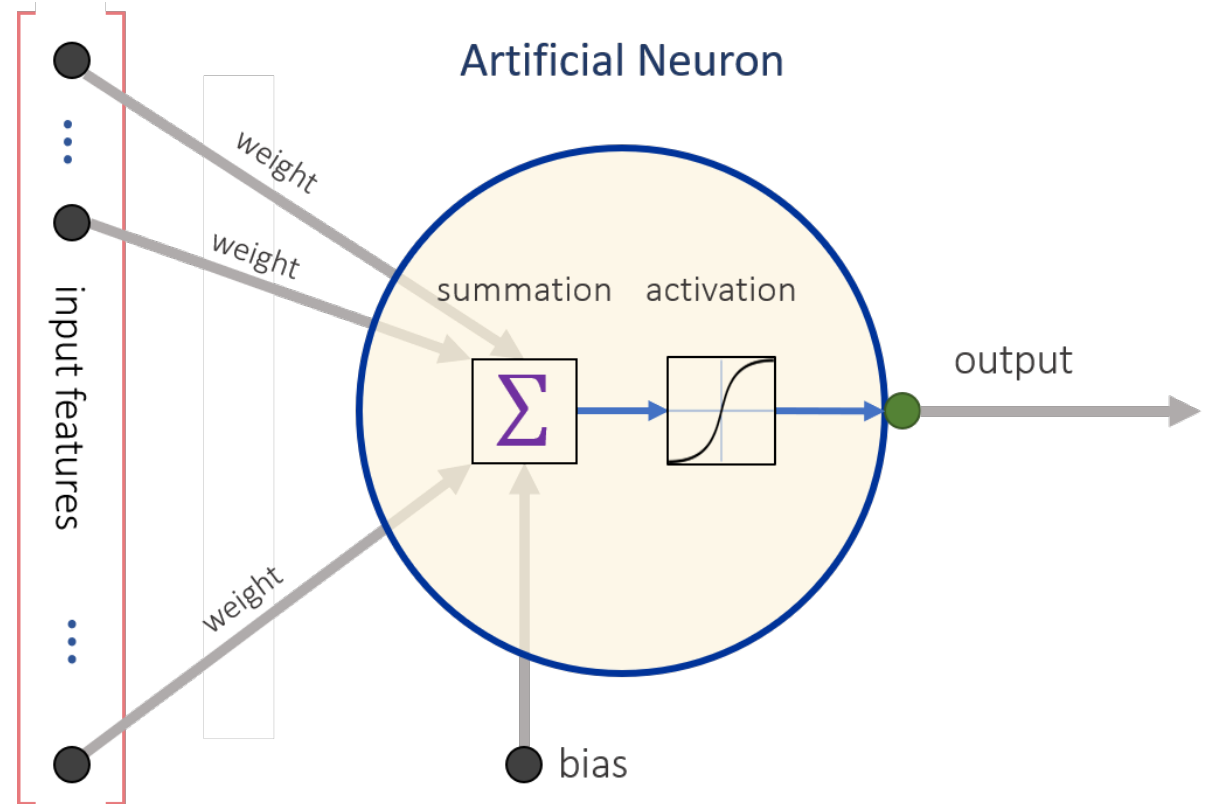
# Deep Learning

## Neurons

### The underlying computation unit is the Neuron

Artificial neurons consist of:

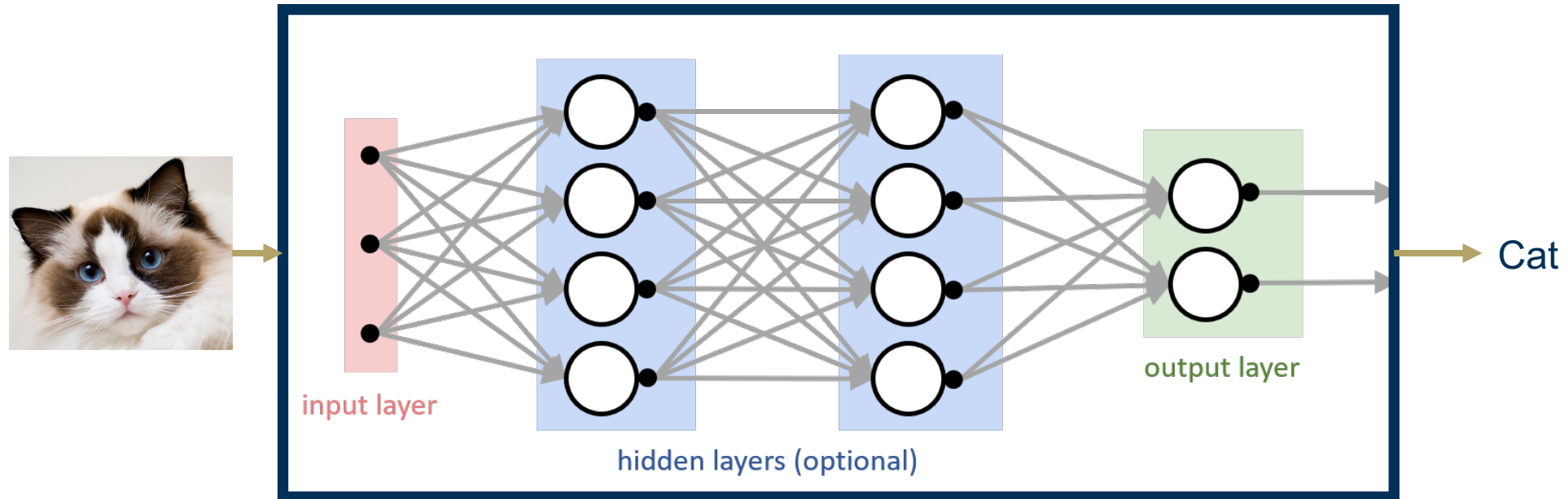
- A single output
- Multiple inputs
- Input weights
- A bias input
- An activation function



# Deep Learning

## Artificial Neural Networks

Neurons are stacked and densely connected to construct ANNs



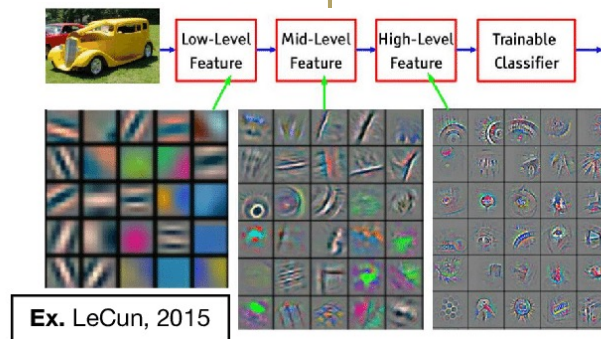
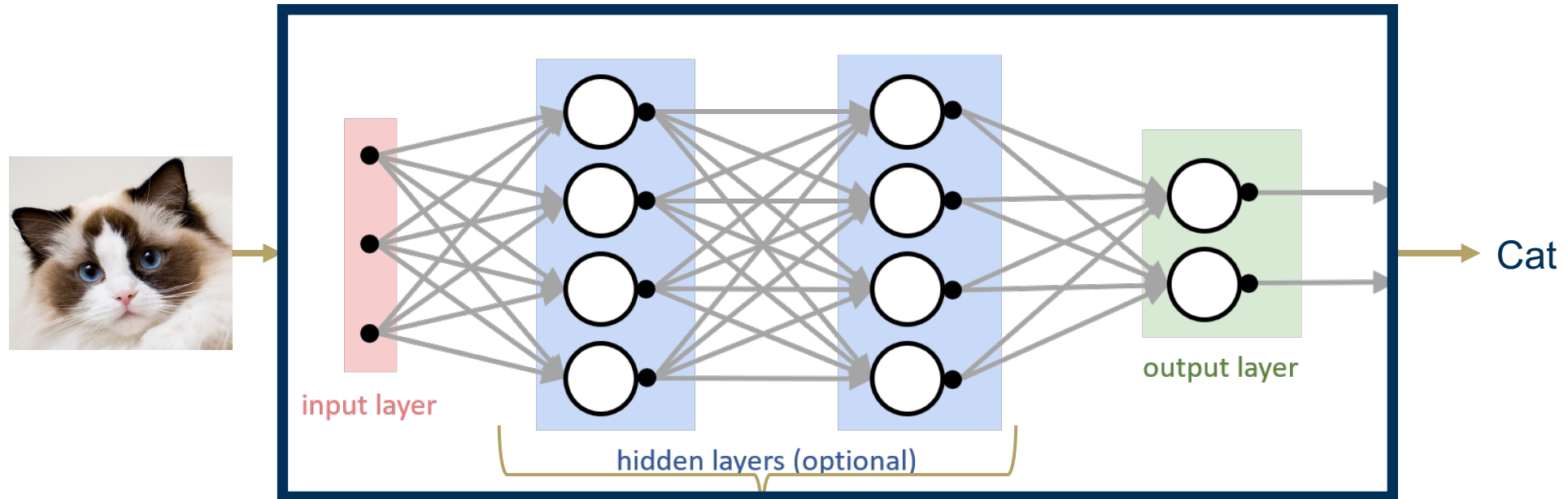
Typically, a neuron is part of a network organized in layers:

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

# Deep Learning

## Convolutional Neural Networks

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

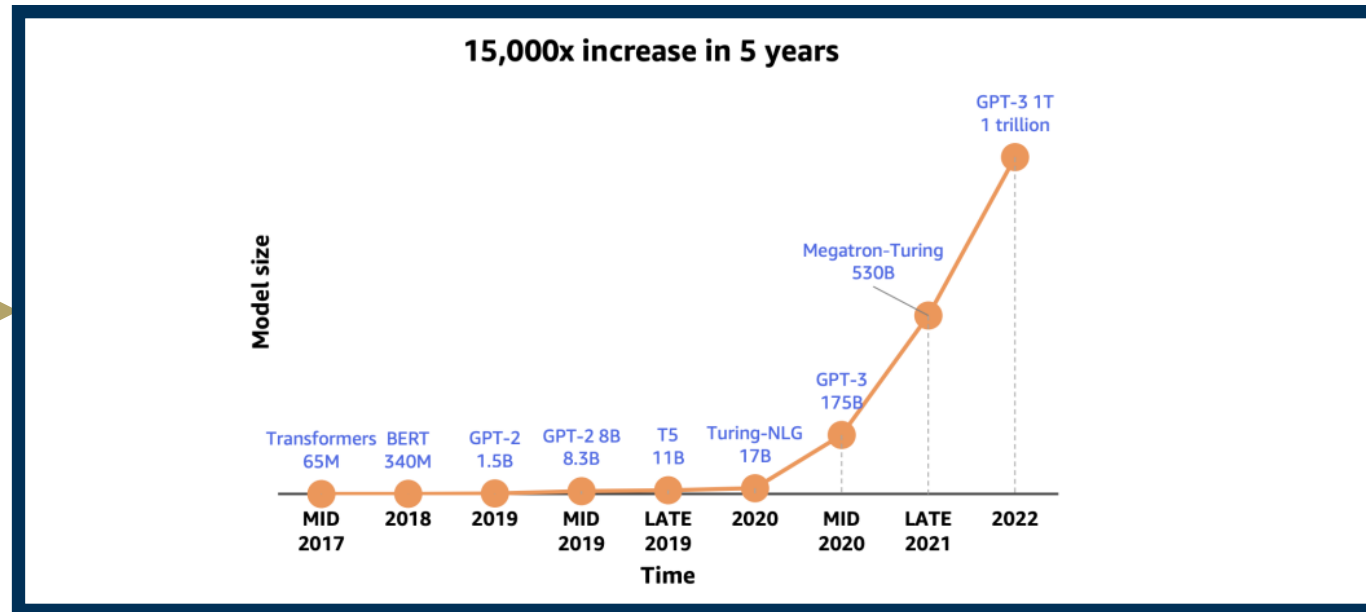
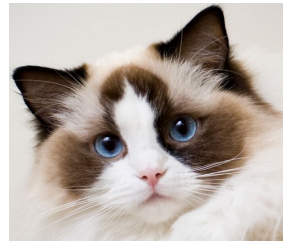




# Deep Deep Deep Deep Deep ... Learning

## Recent Advancements

### Transformers, Large Language Models and Foundation Models



Cat

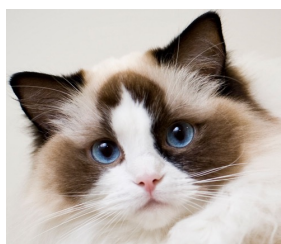
Primary reasons for advancements:

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

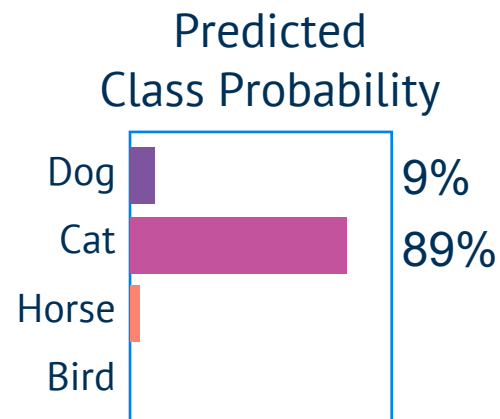
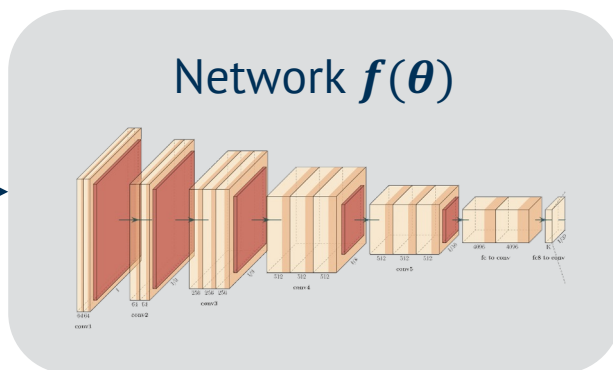
# Deep Learning at Inference

## Classification

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



$x$



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

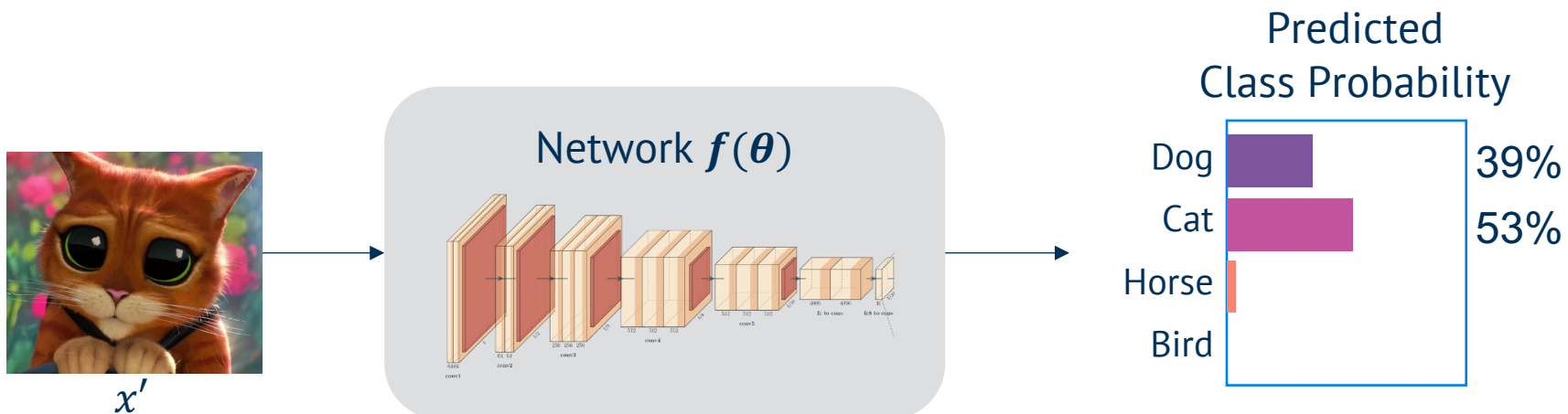
$$\hat{y} = f(x)$$
$$y = \operatorname{argmax}_i \hat{y}$$
$$p(\hat{y}) = T(f(x))$$

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

# Deep Learning at Inference

## Robust Classification in Deep Networks

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



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

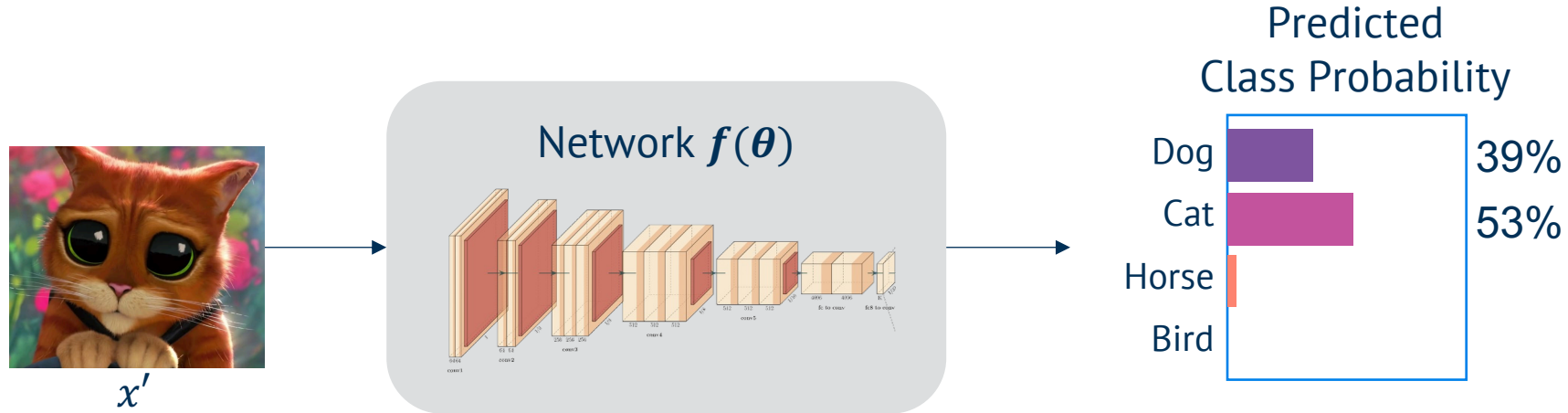
$$\begin{aligned} \hat{y} &= f(x' + \epsilon) & \hat{y} &= \text{Logits} \\ y &= \operatorname{argmax}_i \hat{y} & y &= \text{Predicted Class} \\ p(\hat{y}) &= T(f(x' + \epsilon)) & p(\hat{y}) &= \text{Probabilities} \\ & & f(\cdot) &= \text{Trained Network} \\ & & \chi &= \text{Training data} \\ & &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



# Deep Learning at Inference

## Robust Classification in Deep Networks

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



To achieve robustness at Inference, we need the following:

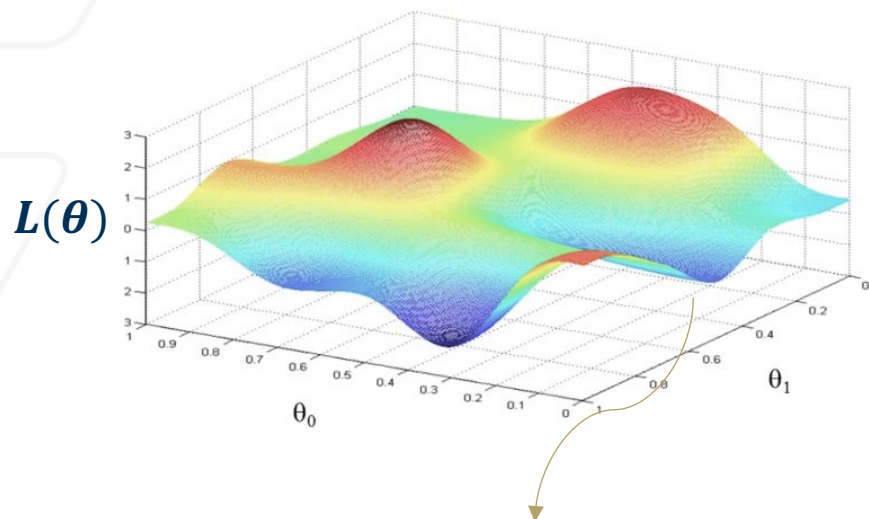
- **Information** provided by the novel data as a **function of training distribution**
- Methodology to **extract information** from novel data
- **Techniques** that utilize the information from novel data

**Why is this Challenging?**

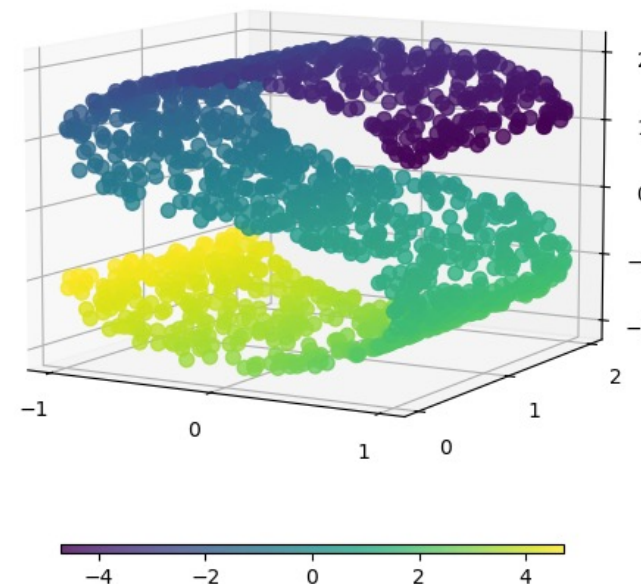
# Challenges at Inference

A Quick note on Manifolds..

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



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

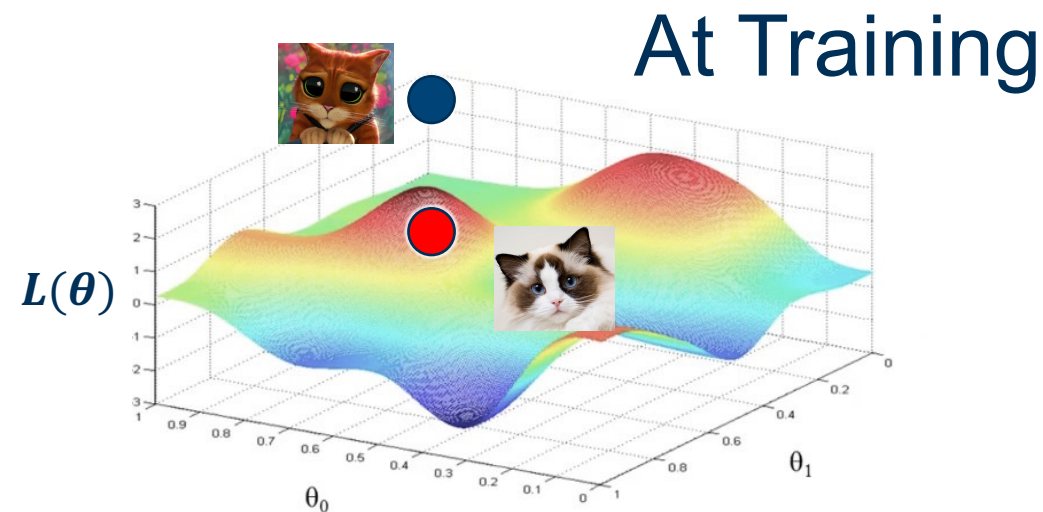
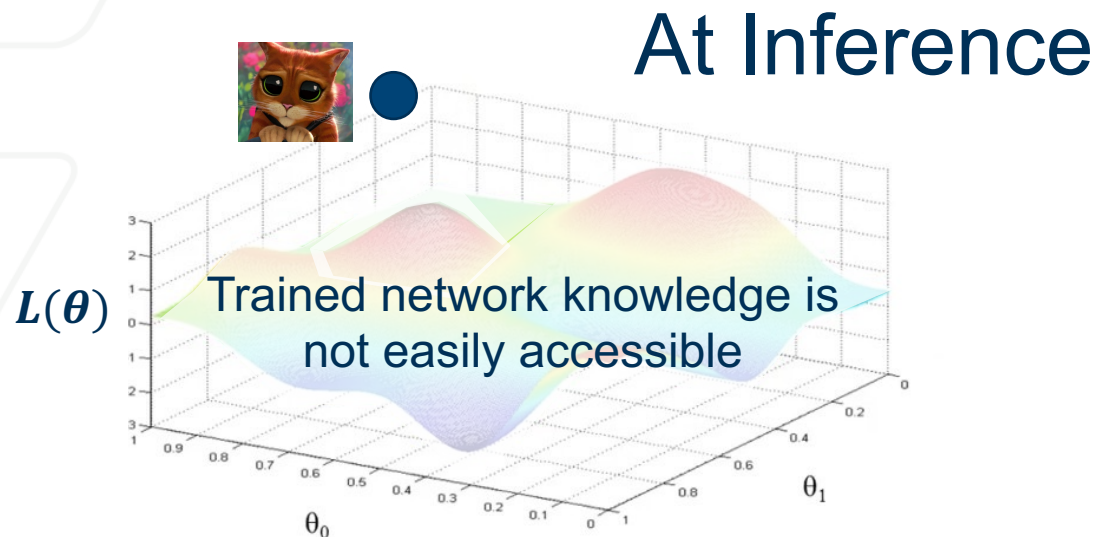


Real data visualizations generated using  
dimensionality reduction algorithms (Isomap)

# Challenges at Inference

## Inference

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

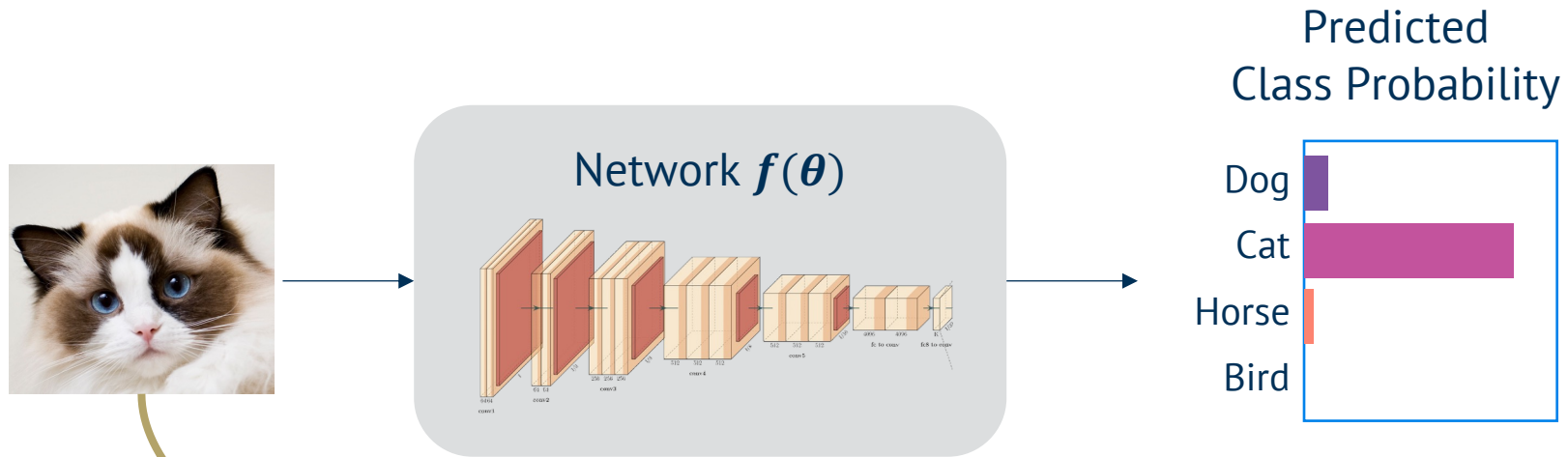


At training, we have access to all training data.

# Information at Inference

## Fisher Information

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

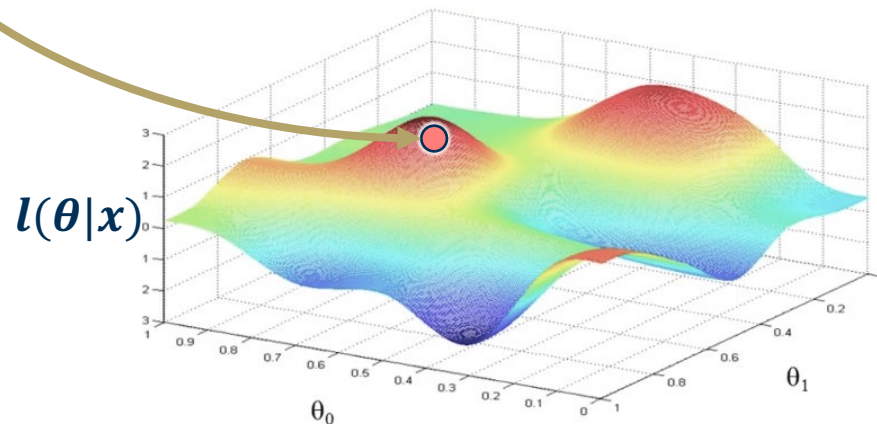


Fisher Information

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

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

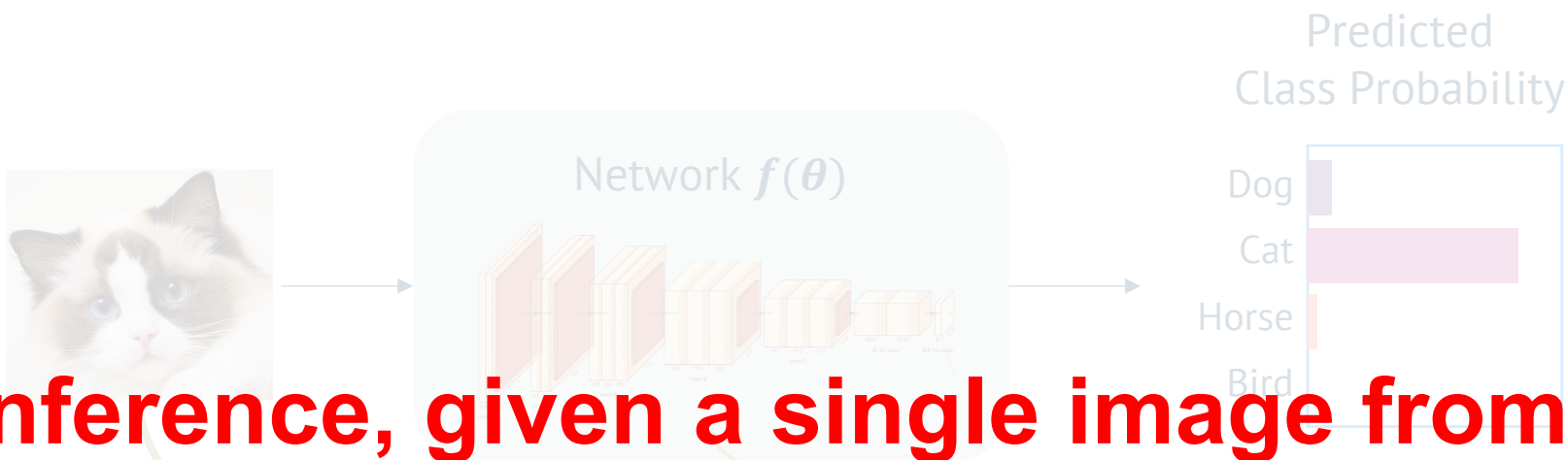
Likelihood function



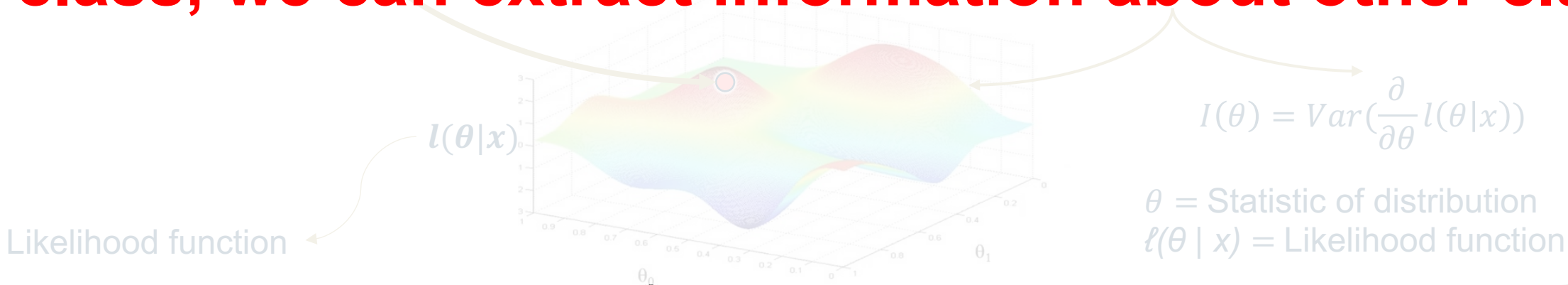


# Information at Inference

Information at Inference



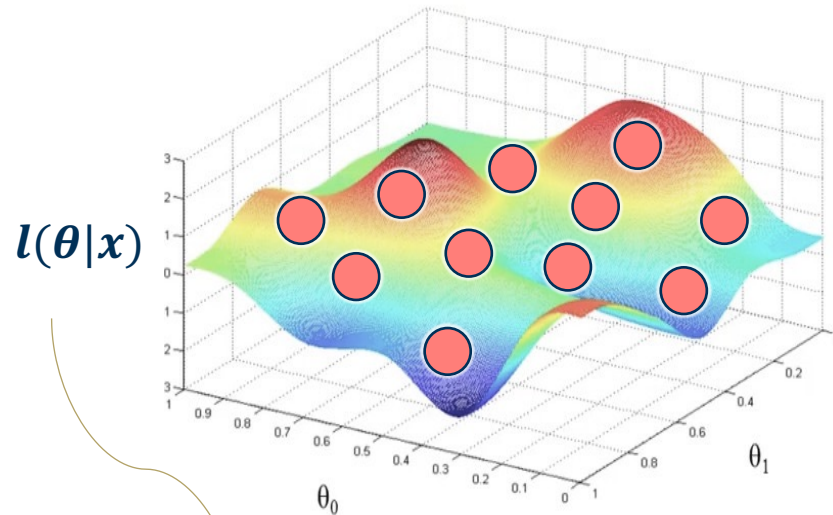
**At inference, given a single image from a single class, we can extract information about other classes**



# Information at Inference

## Gradients as Fisher Information

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



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

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

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

$E[\cdot]$  = Expectation

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

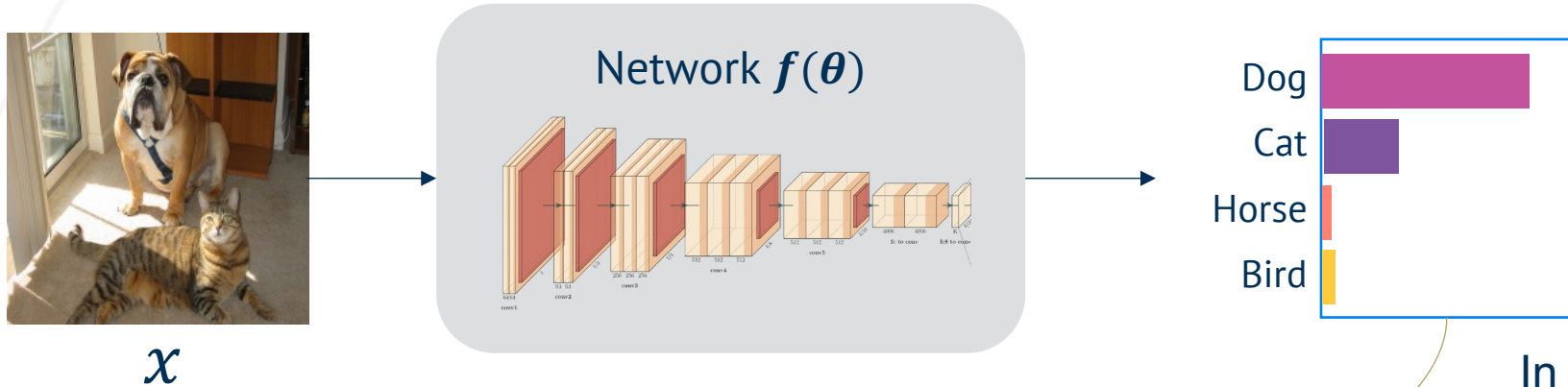
Likelihood function instead of loss manifold

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

# Information at Inference

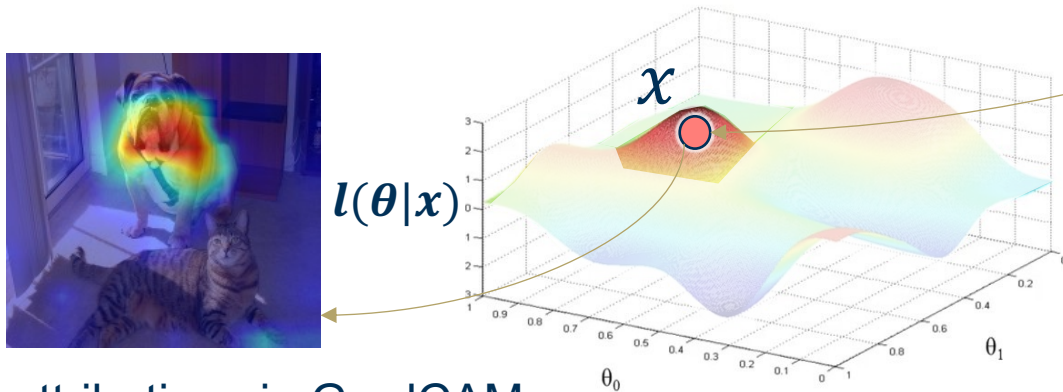
## Case Study: Gradients as Fisher Information in Explainability

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



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

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



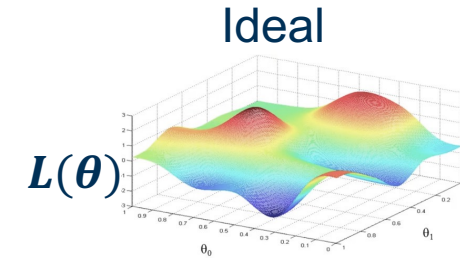
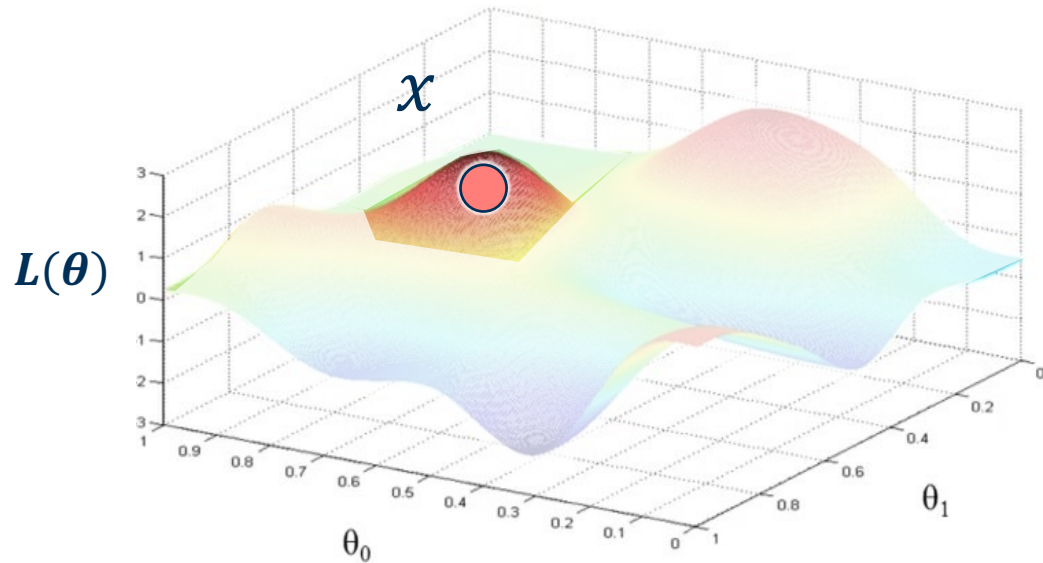
Hence, gradients draw information from the underlying distribution as learned by the network weights!

Feature attribution via GradCAM

# Gradients at Inference

## Local Information

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



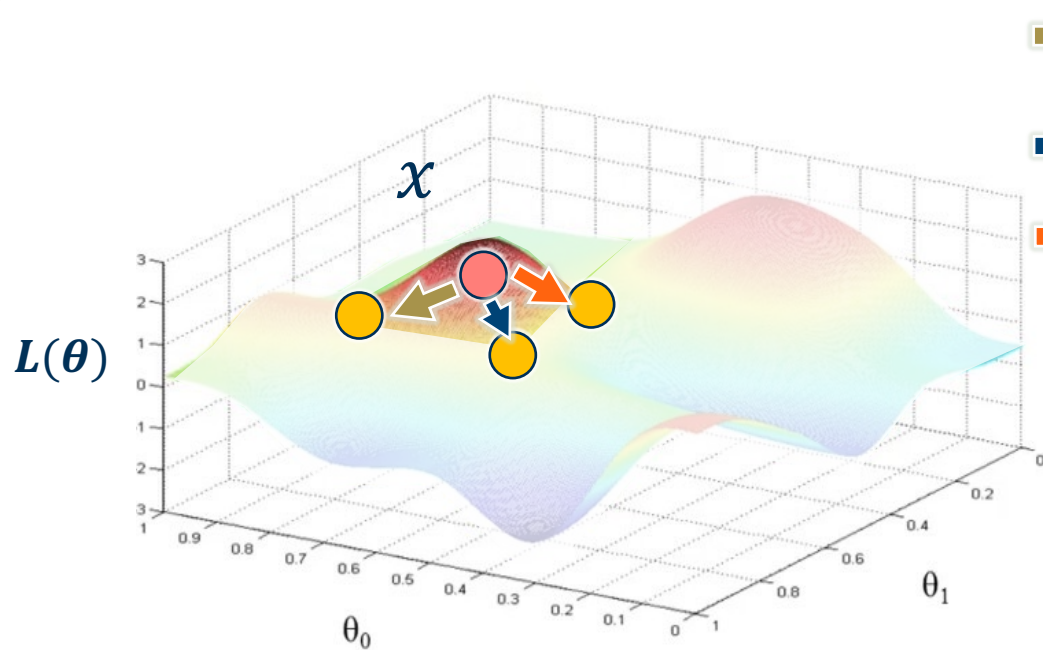
$\alpha \nabla_{\theta} L(\theta)$  provides local information up to a small distance  $\alpha$  away from  $x$



# Gradients at Inference

## Direction of Steepest Descent

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



Path 1?



Path 2?



Path 3?

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

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

# Gradients at Inference

To Characterize the Novel Data at Inference

