

# A Holistic View of Perception in Intelligent Vehicles



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

June 04, 2023 – Anchorage, Alaska

# Autonomous Vehicles

## Why Autonomous Vehicles?



### Safety in Mobility

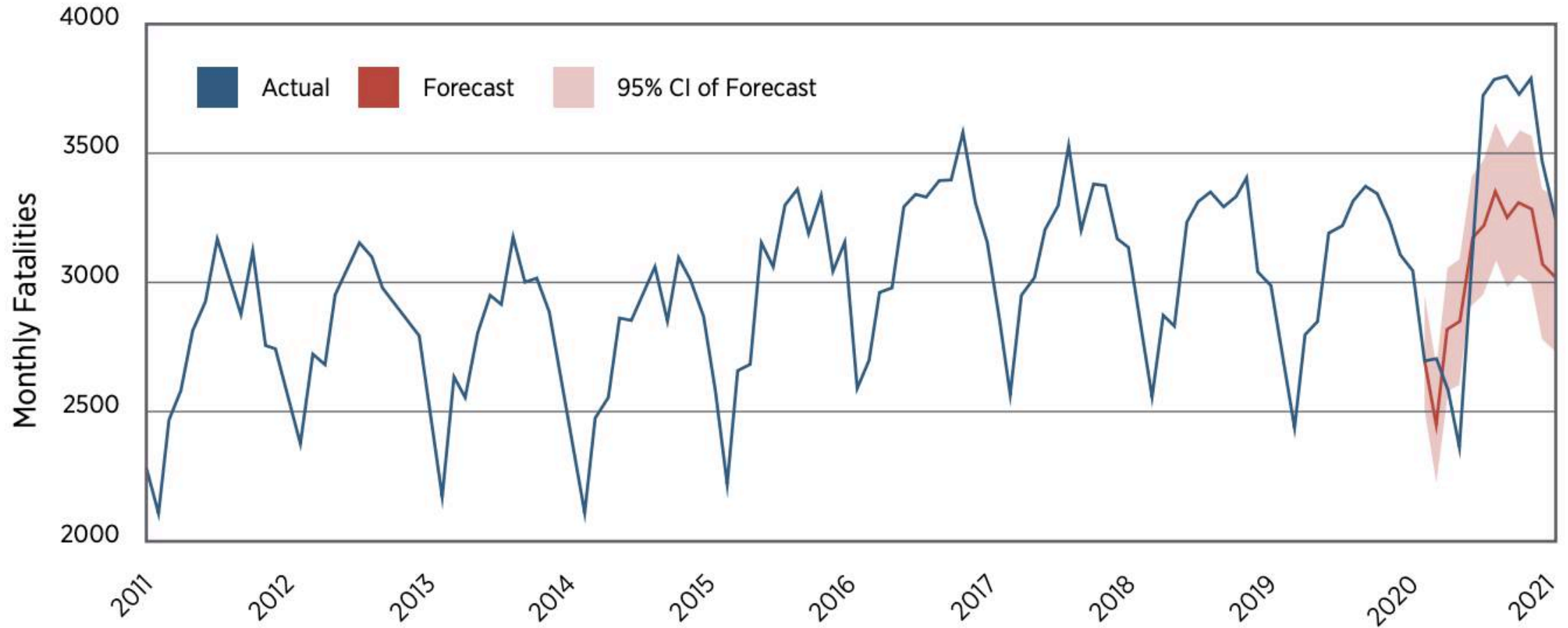


### Mobility Experience

# Autonomous Vehicles

## Why Autonomous Vehicles?

In 2020, despite COVID-19 restrictions, fatalities increased in the US



# Autonomous Vehicles

## Why Autonomous Vehicles?

**Next Revolution in Mobility Safety: AI**  
94% of all car accidents are due to human error



It is estimated that, globally, AVs can prevent 4.22 million accidents per year by 2050

# Autonomous Vehicles

## How will AI ensure Safety in Mobility?

**AI identifies and overcomes human limitations in sensing and simulates complex environments for testing**



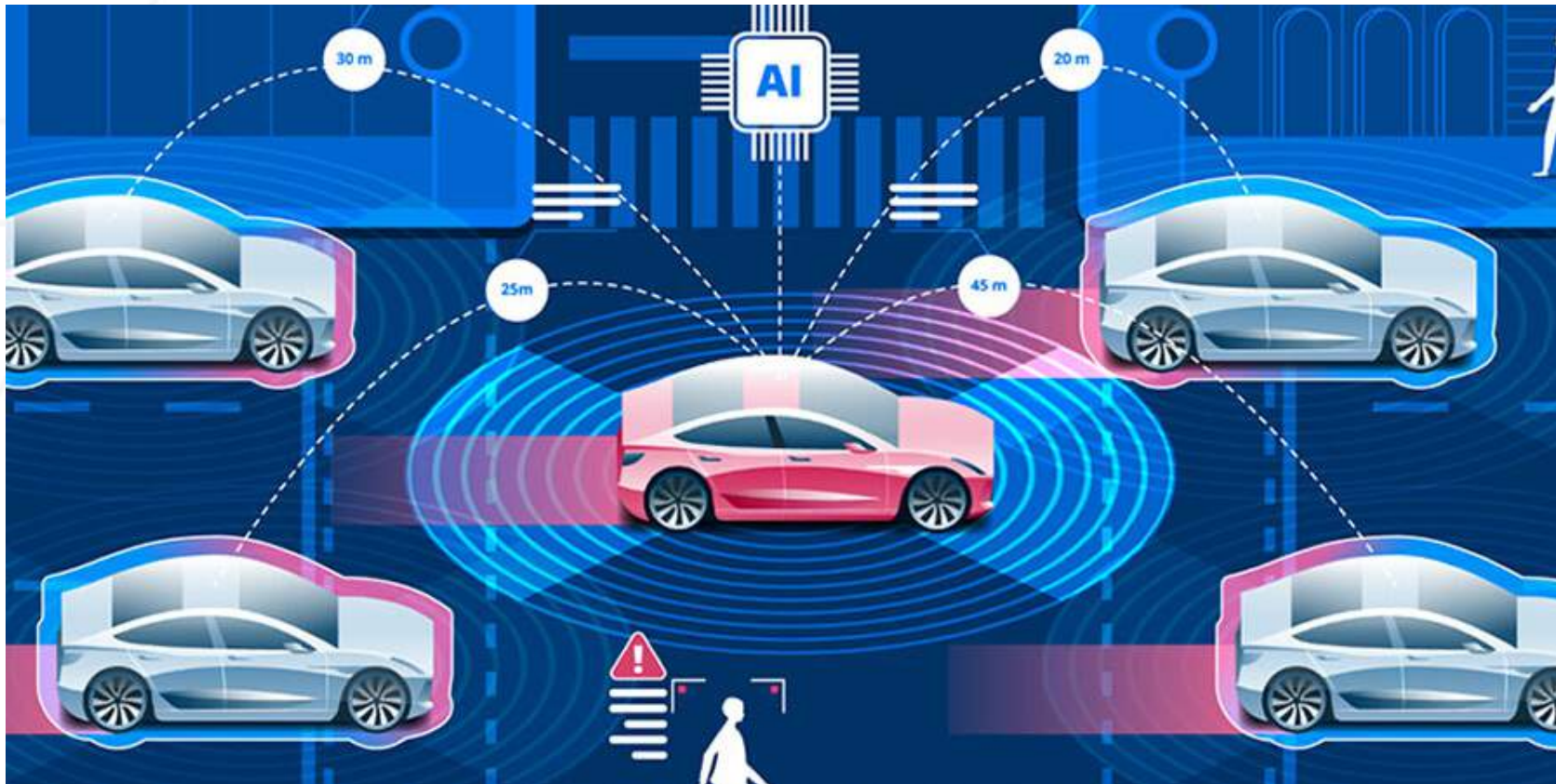
Active sensors like LIDAR overcome the limitations of passive vision sensing

Incredibly complex driving scenarios can be simulated using AI to test itself

# Autonomous Vehicles

## How will AI ensure Safety in Mobility?

AI provides technologies to handle large data modalities in real time environments



Real-time connection to other vehicles, pedestrians, infrastructure and networks is facilitated by AI

# Objectives

## Objectives of the Tutorial

- Part I: Challenges in Perception and Autonomy
- Part II: Deep Learning for Perception
- Part III: Existing Deep Learning solutions to Challenges in Perception
- Part IV: Remaining Challenges and Future Directions

# A Holistic View of Perception in Intel. Vehicles

## Part I: Perception and Autonomy



# Objectives

## Objectives in Part I

- Summarize the progress of AVs over the years
- Discuss the role of perception in AVs and where it fits within the AV workflow
- Review well-known failures of AVs in providing safety to drivers and to others
- Discuss major technical challenges currently facing AV
- Motivate deep learning as a holistic solution to perception challenges

# Perception

What is Perception?



What is perception?

See, process, understand.

# Perception

## Perception in AVs



+



=



# Perception in AVs

Tsubaka Mechanical Engineering Laboratory (1977)

**First standalone “autonomous” vehicle**



**Technology demonstrated:**

Two video cameras and an analog computer onboard for image processing, Detect street markings

# Perception in AVs

Eureka PROMETHEUS Project (1987 - 1995)



## New technologies demonstrated:

Vision enhancement, Lane keeping support, visibility range monitoring, Driver status monitoring, Collision avoidance, Cooperative driving, Autonomous intelligent cruise control

# Perception in AVs

DARPA Grand Challenge (2004 - 2005)



## New technologies demonstrated:

Wide sensor suite including stereo vision, LIDAR, radar, and ultrasound sensors, sensor fusion, obstacle detection, off-road path following, path finding

# Georgia Tech in DARPA Challenge

## Need for Failsafe in AVs

### Video/News Articles

# Remote Repositioning

A driver in the Cloud Remotely Drives a Completely Equipped Vehicle

**New technologies demonstrated:**

Low latency failsafe mechanisms in connected cars



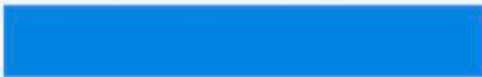

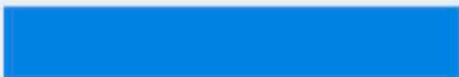


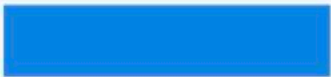






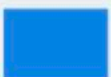









# Perception in AVs

## A Leap in Progress

### AV statistics in California (Dec 2019 – Nov 2020)

			Miles	Miles per disengagement
Waymo (Alphabet)			628,839	 29,945
Cruise (GM)			770,049	 28,520
AutoX			40,734	 20,367
Pony.AI			225,496	 10,738
Argo.AI (Ford, VW)			21,037	 10,519
WeRide			13,014	 6,507
DiDi Chuxing			10,401	 5,201
Nuro			55,370	 5,034

**Disengagement:** Cases where the car's software detects a failure or the driver perceived a failure, resulting in control being seized by the driver.



# Perception in AVs

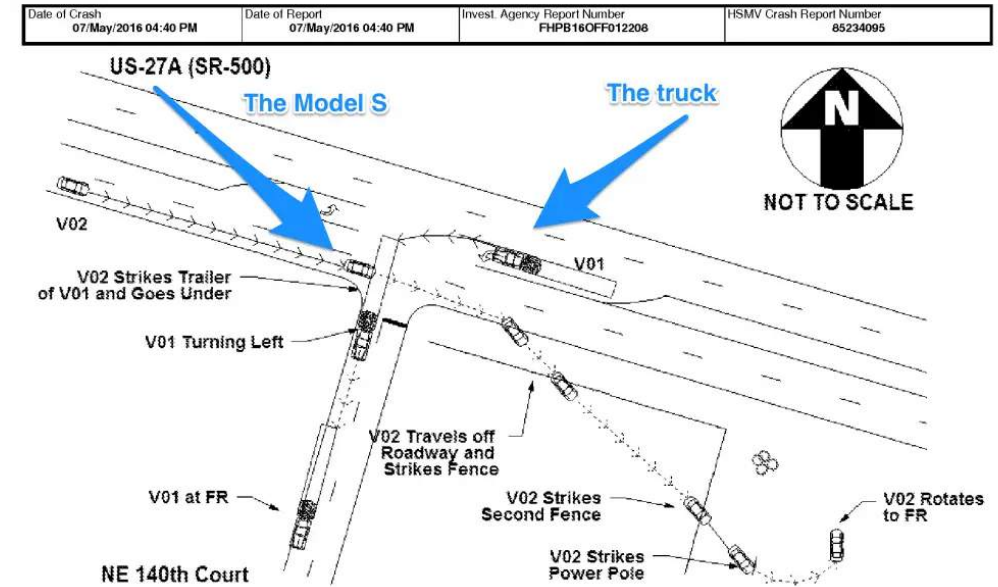
## Setbacks and Challenges

# Tesla driver dies in first fatal crash while using autopilot mode

The autopilot sensors on the Model S failed to distinguish a white tractor-trailer crossing the highway against a bright sky

**Autopilot didn't detect the trailer as an obstacle (NHTSA investigation and Tesla statements)**

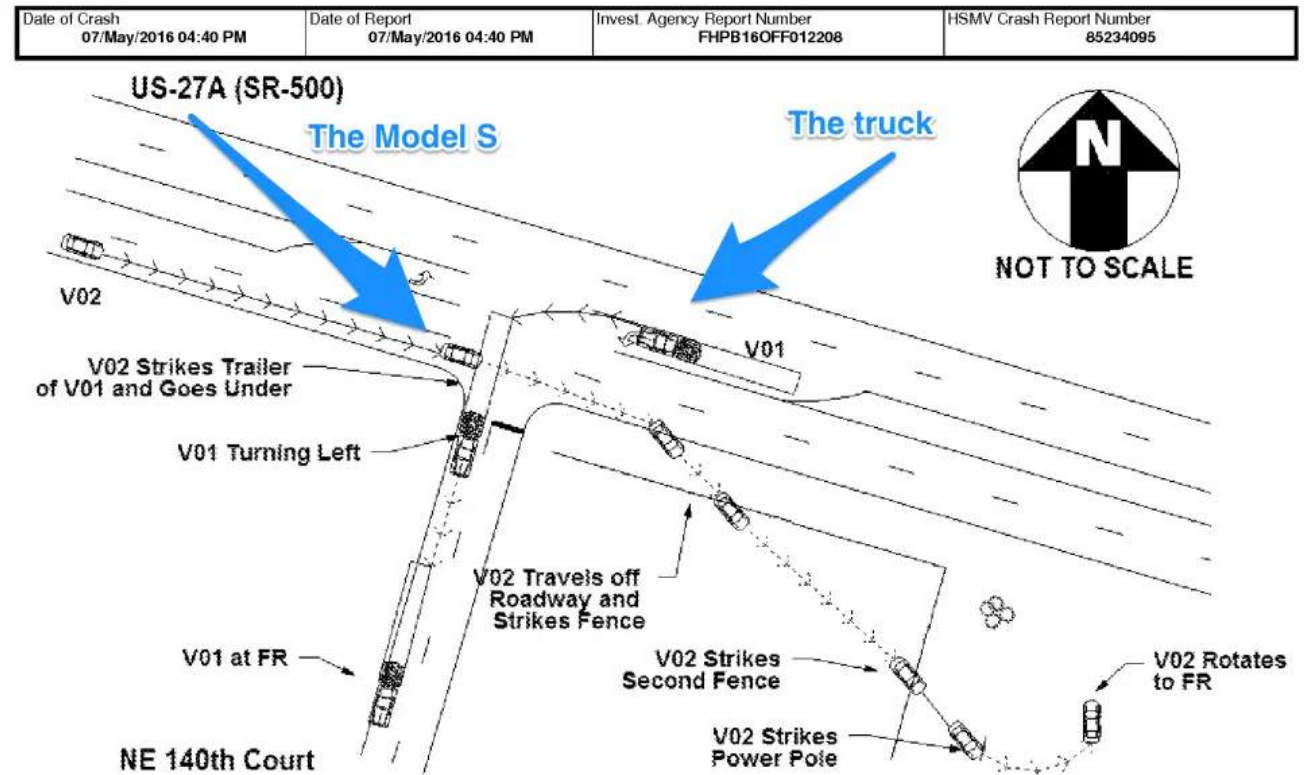
1. The National Highway Traffic Safety Administration (NHTSA) determined that a "lack of safeguards" contributed to the death
2. "Neither Autopilot nor the driver noticed the white side of the tractor trailer against a brightly lit sky, so the brake was not applied," Tesla said.



# Challenges in Perception in Autonomous Vehicles

## Tesla driver dies in first fatal crash while using autopilot mode

1. The National Highway Traffic Safety Administration (NHTSA) determined that a “lack of safeguards” contributed to the death
2. "Neither Autopilot nor the driver noticed the white side of the tractor trailer against a brightly lit sky, so the brake was not applied," Tesla said.



# Uber's self-driving SUV saw the pedestrian in fatal accident but didn't brake, officials say

PUBLISHED THU, MAY 24 2018•9:52 AM EDT | UPDATED THU, MAY 24 2018•10:43 AM EDT



Sensors on the fully autonomous Volvo XC-90 SUV spotted [REDACTED] while the car was traveling 43 miles per hour and determined that braking was needed 1.3 seconds before impact, according to the report.

# Perception in AVs

## Technical Challenges

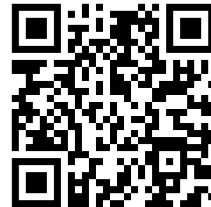
- Challenging weather
- Challenging sensing
- Challenging environments
- Context awareness
- Embedded perception
- V2X perception



# Technical Challenges in Perception for AVs

## Challenging Sensing and Weather

- Challenging weather
- Challenging sensing
- Challenging environments
- Context awareness
- Embedded perception
- V2X perception



# Technical Challenges in Perception for AVs

## Challenging Environments

- Challenging weather
- Challenging sensing
- Challenging environments
- Context awareness
- Embedded perception
- V2X perception



# Technical Challenges in Perception for AVs

## Context Awareness

### Does the fire impede mobility?

- Challenging weather
- Challenging sensing
- Challenging environments
- **Context awareness**
- Embedded perception
- V2X perception





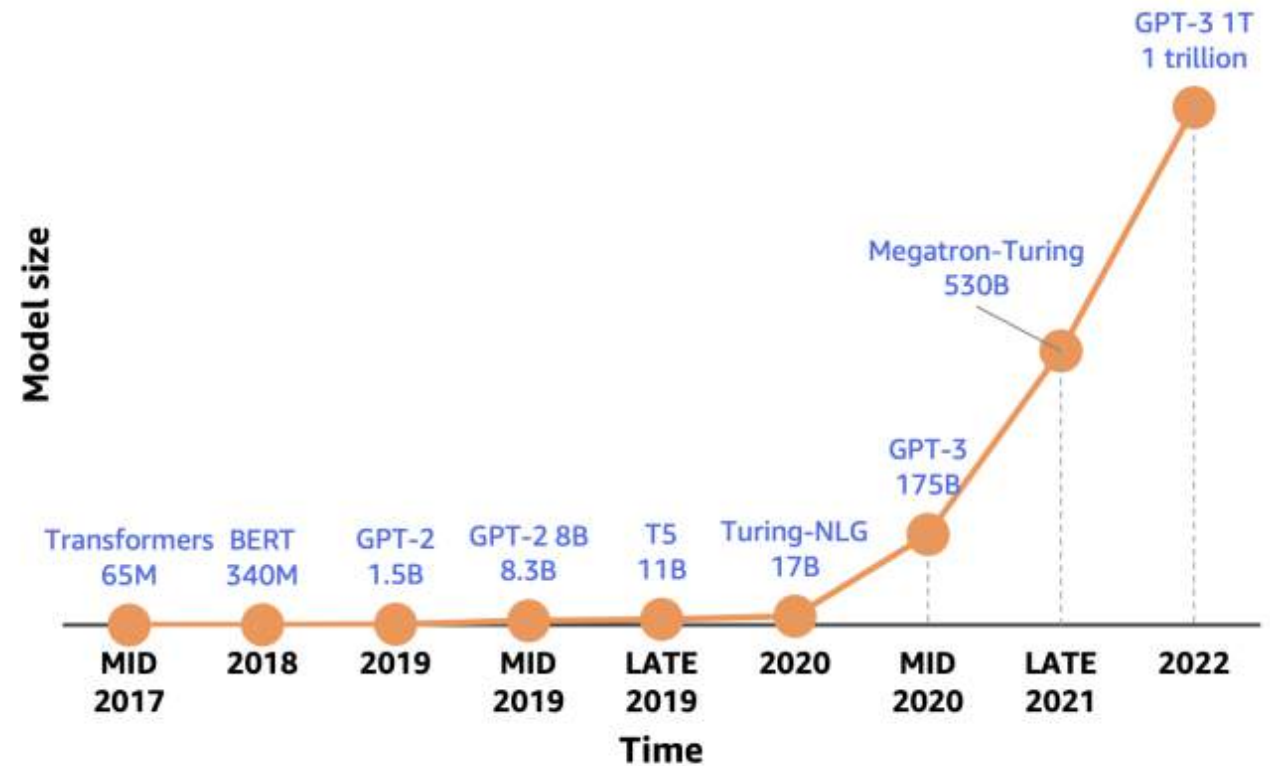
# Technical Challenges in Perception for AVs

## Embedded Perception

**On-board computational capabilities of modern deep learning algorithms is a challenge**

- Challenging weather
- Challenging sensing
- Challenging environments
- Context awareness
- Embedded perception
- V2X perception

**15,000x increase in 5 years**



# Technical Challenges in Perception for AVs

## V2X Perception

**Source: Fast and Furious 8!**

- Challenging weather
- Challenging sensing
- Challenging environments
- Context awareness
- Embedded perception
- V2X perception



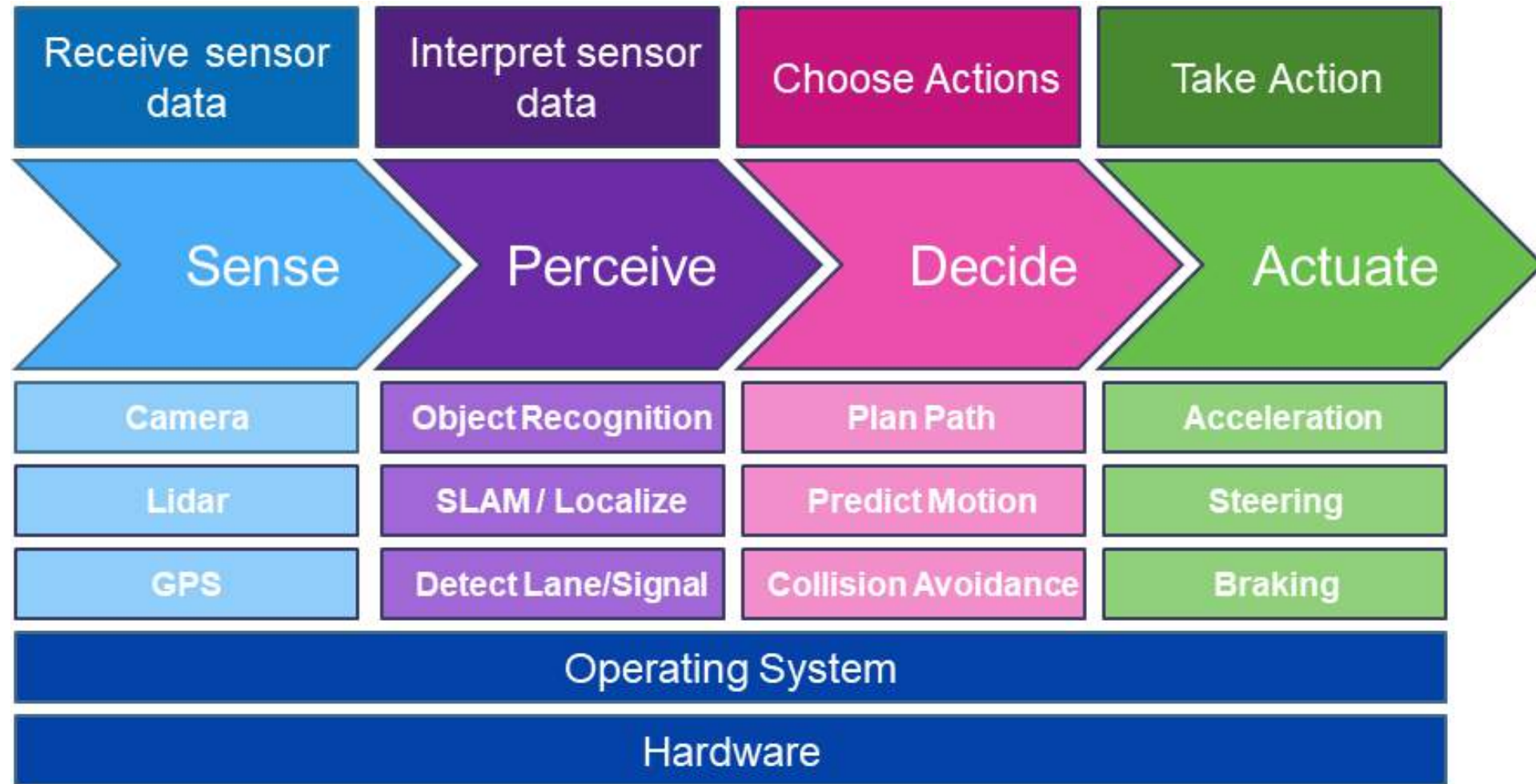
# Role of Perception

## Role of Perception within AVs

Role of Perception:

- Filter,
- process, and
- understand

sensor data



# Sensors

## Role of Sensors for Perception



Automatically Operated Car



Eureka PROMETHEUS Project (1987 - 1995)

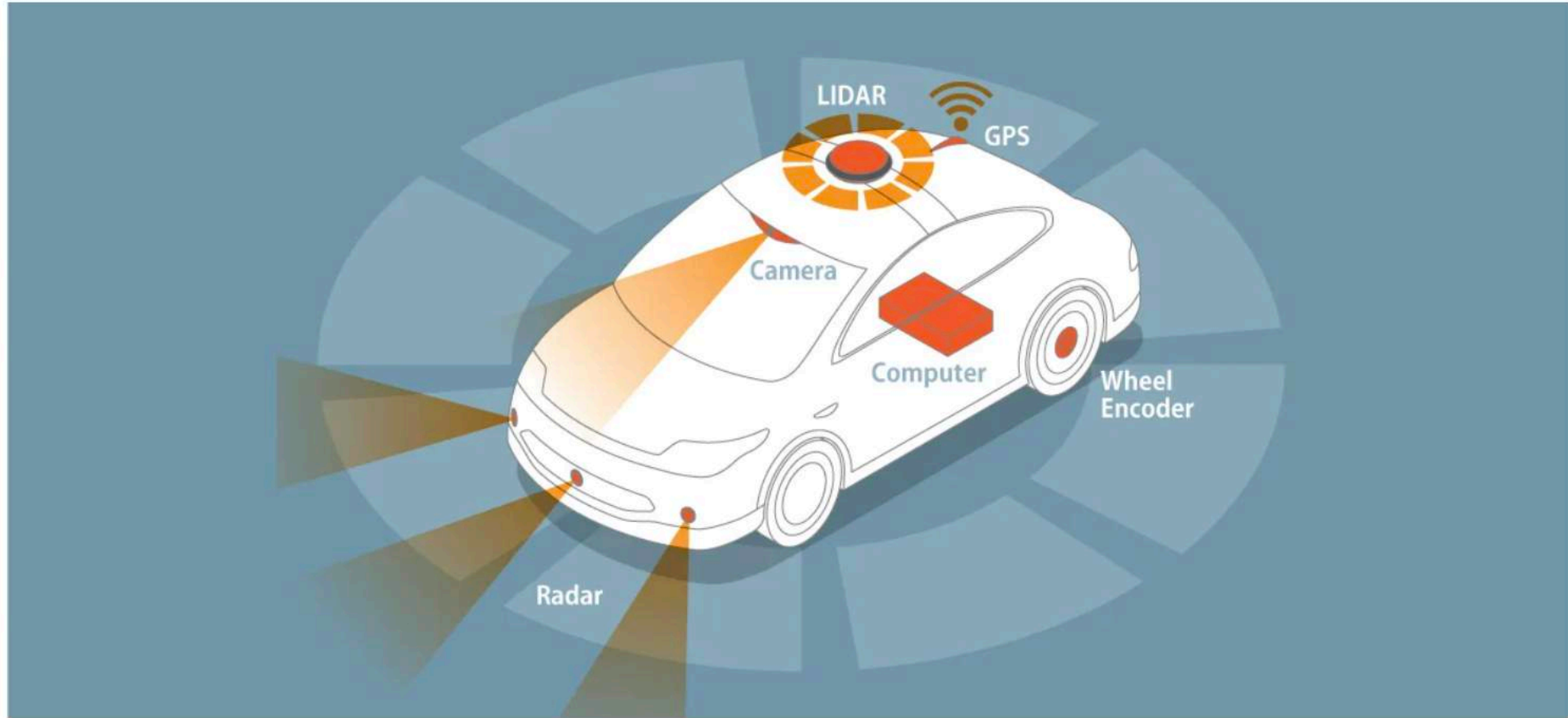


DARPA Grand Challenge (2004 - 2005)

More sensors and better fusion strategies!

# Sensors

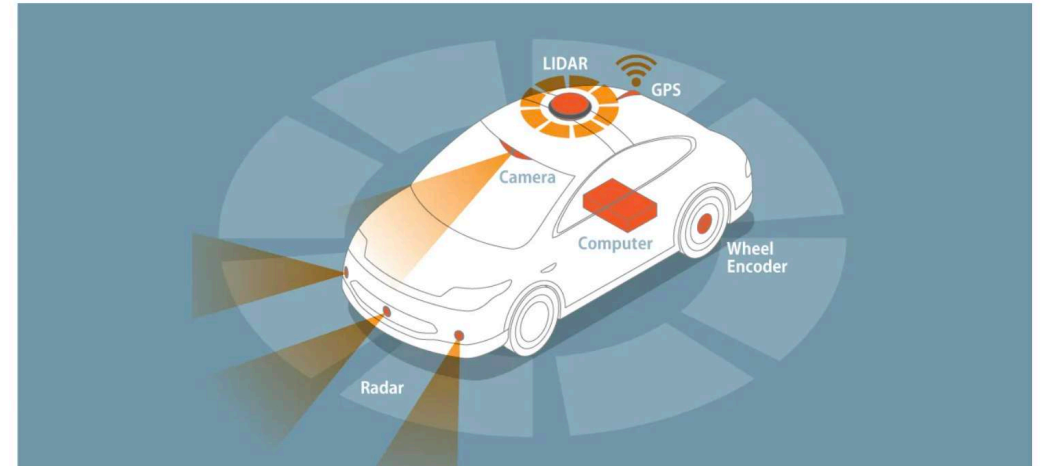
How can we choose the “appropriate” Sensors?



# Sensors

## Choosing the Appropriate Sensors

- Sensors need to work under **challenging weather conditions**
- Sensors need to have sensing capacity and resolution in meeting **challenging sensing environments**
- Sensors must be **cost effective**
- Sensor fusion and sensor registration must be **computationally effective**
- Sensors must output minimum **noise** or their working ranges must be known in advance
- Sensor data must be resistant to **cyber and adversarial attacks**



# Sensors

## Choosing the Appropriate Sensors

Factors	Camera	LiDAR	Radar	Fusion
Range	~	~	✓	✓
Resolution	✓	~	×	✓
Distance Accuracy	~	✓	✓	✓
Velocity	~	×	✓	✓
Color Perception, e.g., traffic lights	✓	×	×	✓
Object Detection	~	✓	✓	✓
Object Classification	✓	~	×	✓
Lane Detection	✓	×	×	✓
Obstacle Edge Detection	✓	✓	×	✓
Illumination Conditions	×	✓	✓	✓
Weather Conditions	×	~	✓	✓

# Sensors

## Choosing the Appropriate Sensors

TABLE I  
DIFFERENT SENSORS USED IN AV DEVELOPMENT

Vehicle	A <sup>#</sup>	B	C	D	E	F
Audi's Research Vehicle [48]	Y	Y	Y	Y	Y	Y
Ford: Hybrid Fusion [49]	Y			Y	Y	Y
Google: Toyota Prius [50]	Y	Y		Y	Y	
Nagoya and Nagasaki University's Open ZMP Robocar HV (Toyota Prius) [51]	Y			Y		
Volvo: (Stoklosa, Cars) [52]	Y		Y	Y	Y	Y
Apple: Lexus RX450h SUVs [53]	Y		Y	Y	Y	Y
DIDI's research vehicle [54]	Y		Y	Y	Y	Y
Infiniti Q50S [55]	Y				Y	Y
Lexus RX [56]	Y				Y	Y
Volvo XC90 [57]	Y				Y	Y
BMW750i xDrive [58]	Y	Y	Y		Y	Y
Mercedes-Benz E & S-Class [55]	Y	Y	Y		Y	Y
Otto Semi-Trucks [59]	Y			Y	Y	
Renault GT Nav [60]	Y				Y	Y
Tesla Model S [61]	Y				Y	Y
Baidu Apollo [62]	Y				Y	Y

#Note: A:Vision; B:Stereovision; C:IR Camera; D:LIDAR; E:Radar; and F:Sonar.



# Levels of Autonomy Taxonomy



## SAE J3016™ LEVELS OF DRIVING AUTOMATION™

Learn more here: [sae.org/standards/content/j3016\\_202104](https://www.sae.org/standards/content/j3016_202104)

Copyright © 2021 SAE International. The summary table may be freely copied and distributed AS-IS provided that SAE International is acknowledged as the source of the content.

	SAE LEVEL 0™	SAE LEVEL 1™	SAE LEVEL 2™	SAE LEVEL 3™	SAE LEVEL 4™	SAE LEVEL 5™
What does the human in the driver's seat have to do?	You <u>are</u> driving whenever these driver support features are engaged – even if your feet are off the pedals and you are not steering			You <u>are not</u> driving when these automated driving features are engaged – even if you are seated in “the driver’s seat”		
	You <u>must constantly supervise</u> these support features; you must steer, brake or accelerate as needed to maintain safety			When the feature requests, you must drive	These automated driving features will not require you to take over driving	

Copyright © 2021 SAE International.

Current technology:

- Levels 1 and 2 are in the market
- Extensive testing on Level 3

	These are driver support features			These are automated driving features		
What do these features do?	These features are limited to providing warnings and momentary assistance	These features provide steering <b>OR</b> brake/acceleration support to the driver	These features provide steering <b>AND</b> brake/acceleration support to the driver	These features can drive the vehicle under limited conditions and will not operate unless all required conditions are met	This feature can drive the vehicle under all conditions	
Example Features	<ul style="list-style-type: none"> <li>• automatic emergency braking</li> <li>• blind spot warning</li> <li>• lane departure warning</li> </ul>	<ul style="list-style-type: none"> <li>• lane centering <b>OR</b></li> <li>• adaptive cruise control</li> </ul>	<ul style="list-style-type: none"> <li>• lane centering <b>AND</b></li> <li>• adaptive cruise control at the same time</li> </ul>	<ul style="list-style-type: none"> <li>• traffic jam chauffeur</li> </ul>	<ul style="list-style-type: none"> <li>• local driverless taxi</li> <li>• pedals/steering wheel may or may not be installed</li> </ul>	<ul style="list-style-type: none"> <li>• same as level 4, but feature can drive everywhere in all conditions</li> </ul>

[Tutorial] | [Ghassan AlRegib and Mohit Prabhushankar] | [June 4, 2023]

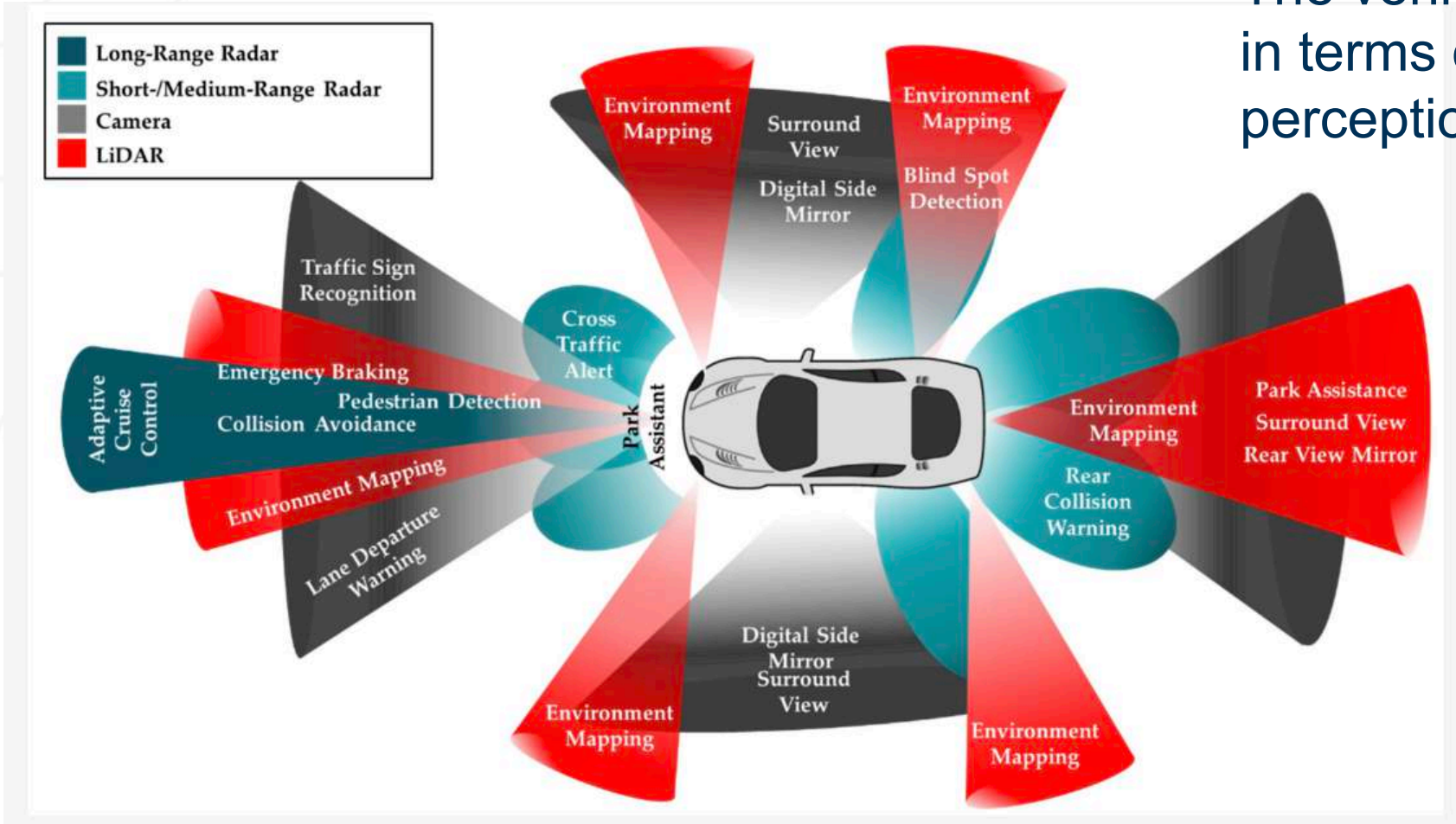
<https://www.sae.org/blog/sae-j3016-update>



# Levels of Autonomy

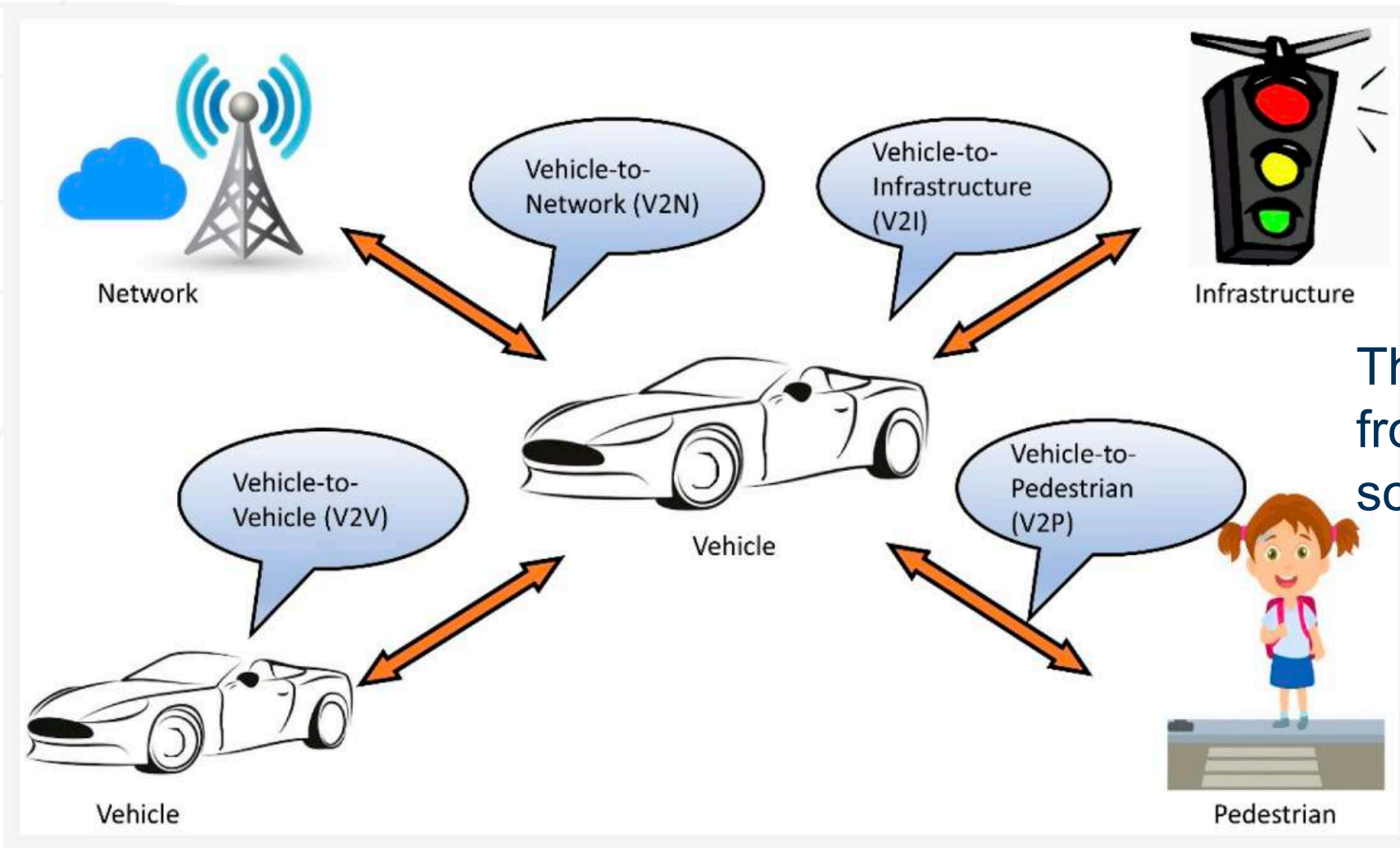
## Levels 1 and 2 Autonomy

The vehicle is self-sufficient in terms of onboard sensors and perception!



# Levels of Autonomy

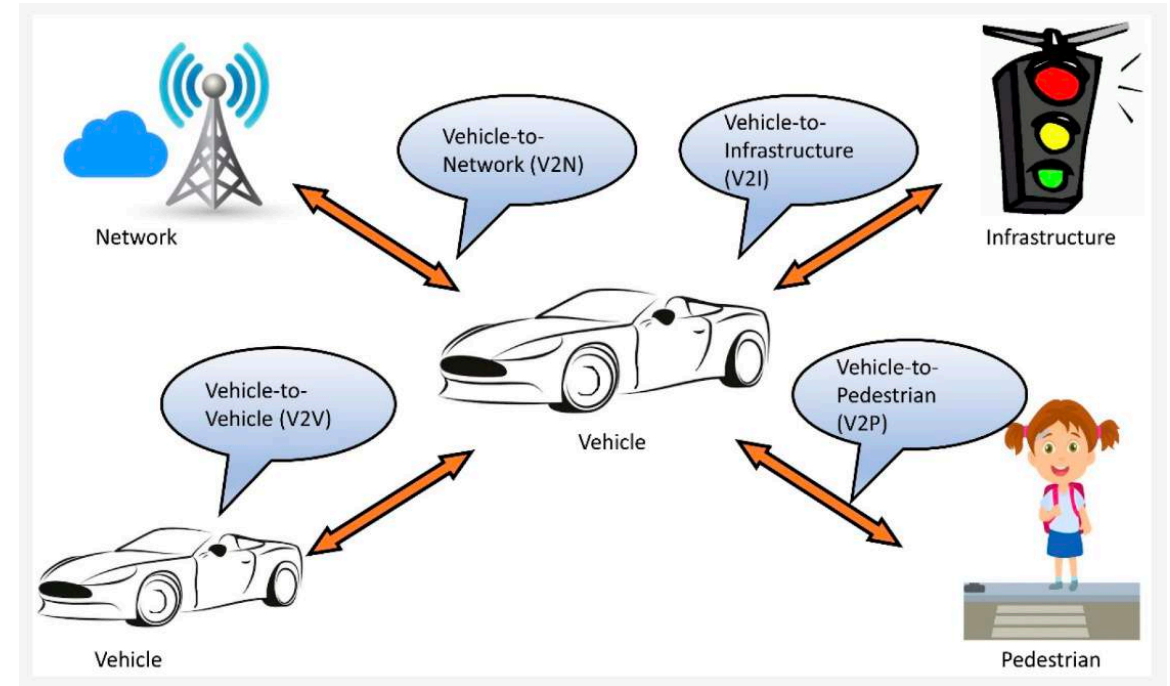
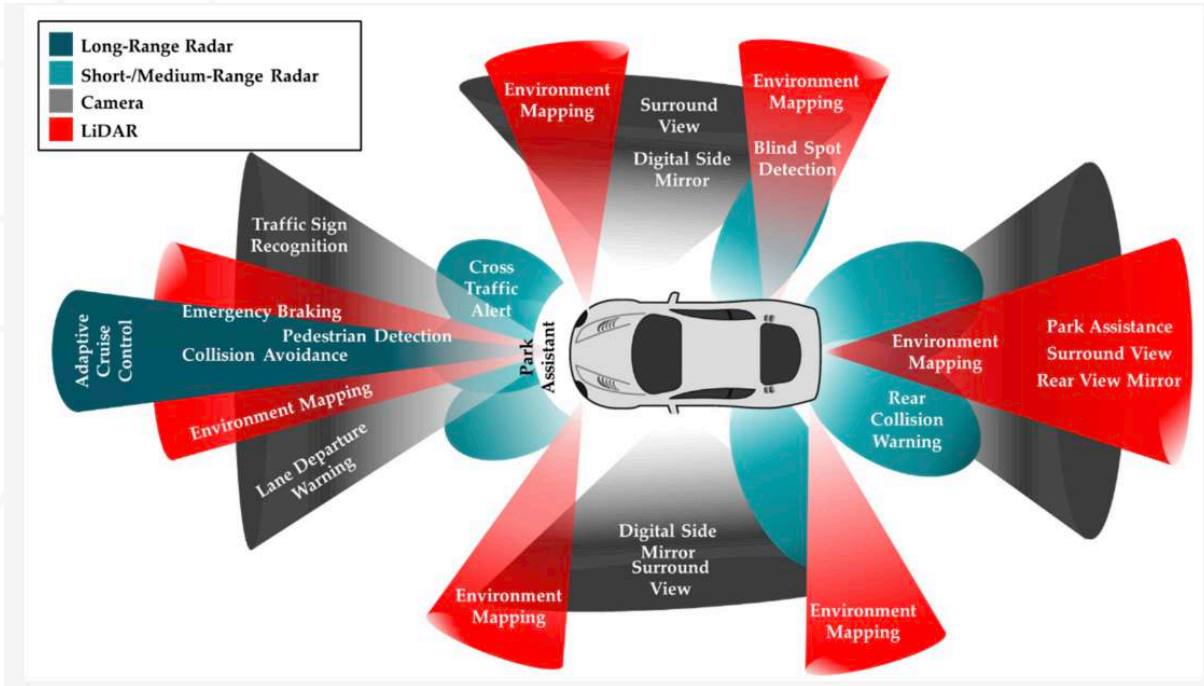
## Levels 3 and Beyond



The vehicle needs help from other sensors, sources, and processors!

# Levels of Autonomy

## Achieving Perception

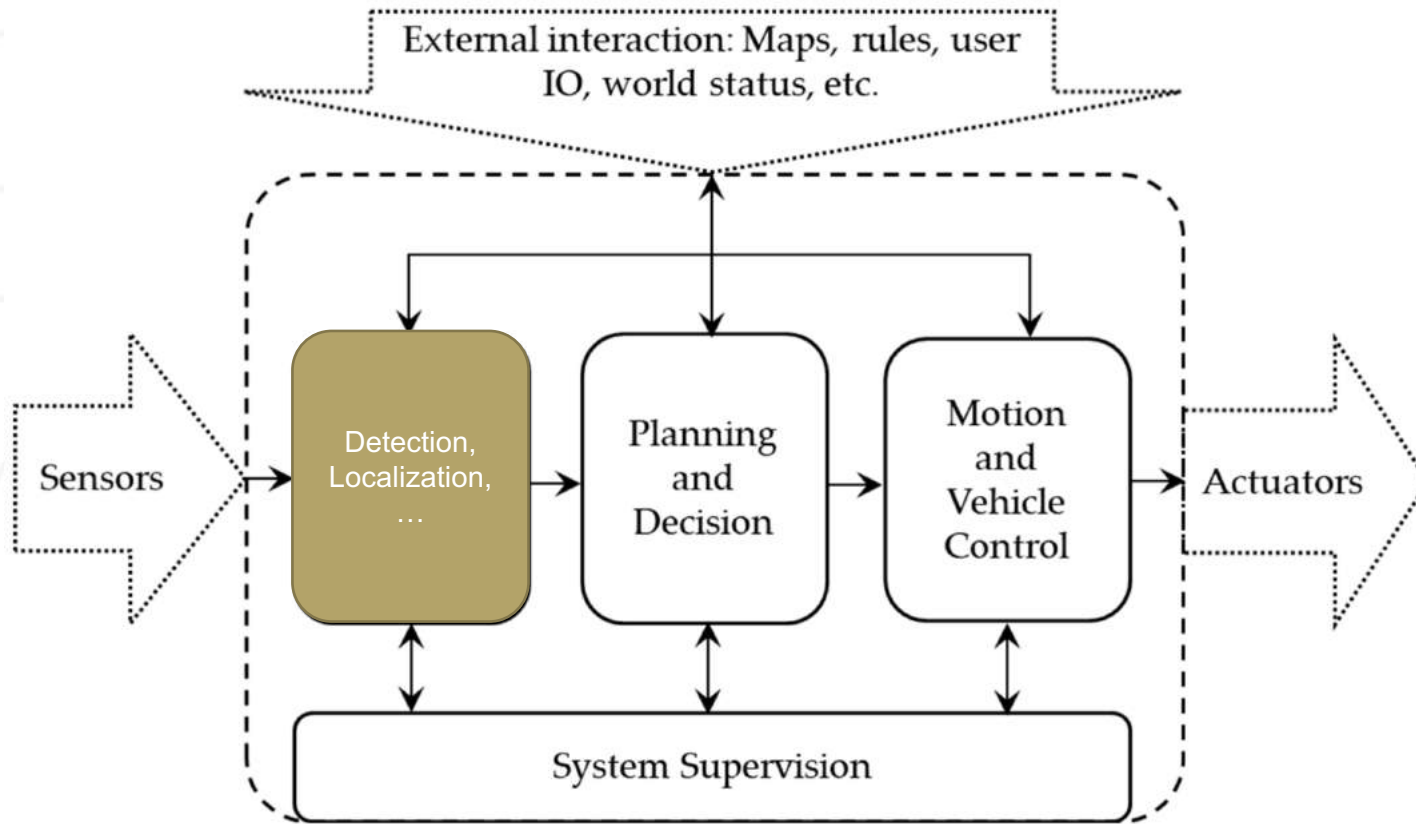


How to filter, process, and understand sensor data?

# Levels of Autonomy

## Achieving Perception

**Before: Perception is decomposed into a number of manageable applications**

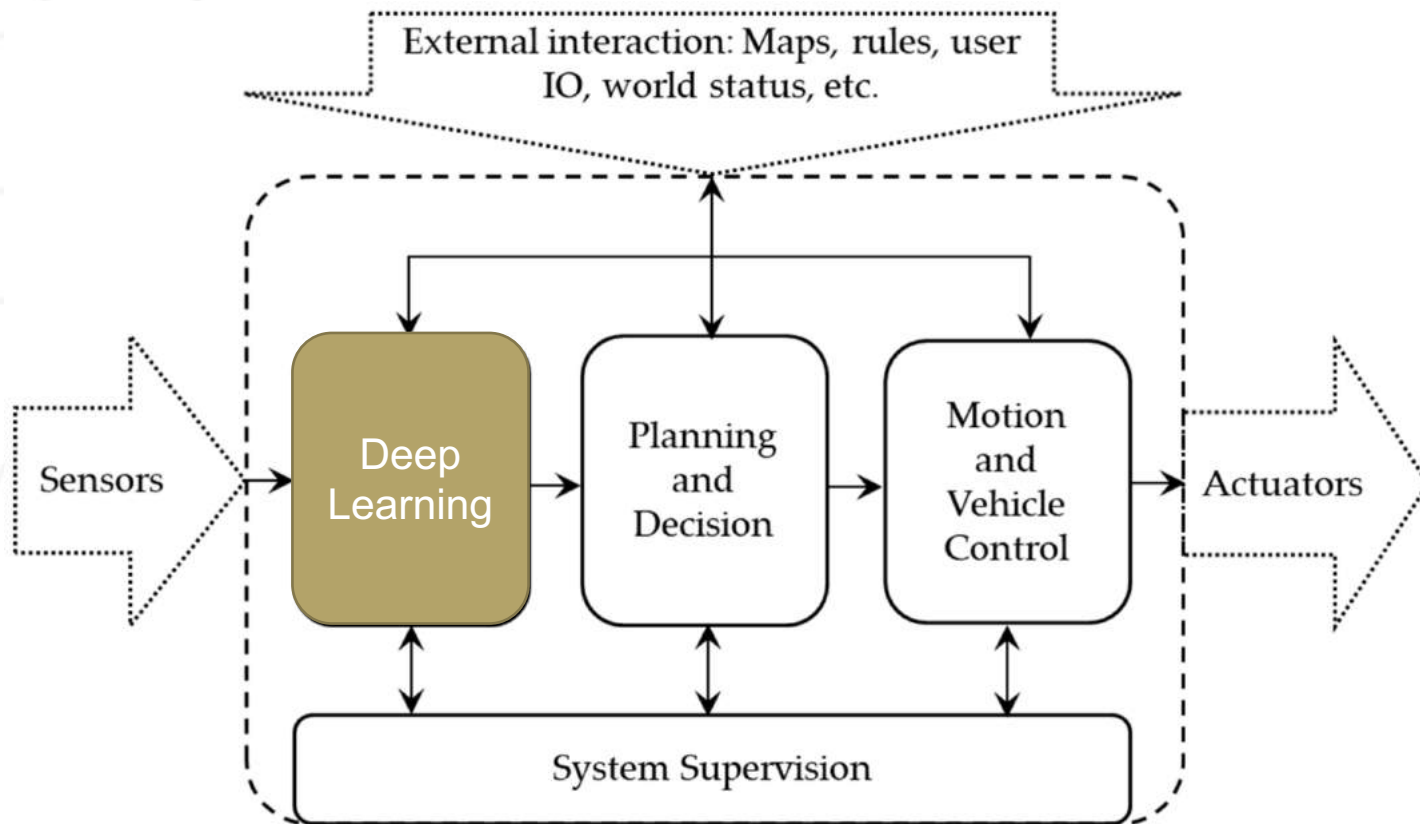


How to filter, process, and understand sensor data?

# Levels of Autonomy

## Goal of the Tutorial

**Deep Learning: Provides a holistic solution to perception**



How to filter, process, and understand sensor data?

# Objectives

## Takeaways from Part I

- **Part I: Challenges in Perception and Autonomy**
  - Robustness under challenging conditions, environments, context and surroundings-awareness are challenges in AV perception
  - Deep Learning promises a holistic solution to a number of the above challenges
- Part II: Deep Learning for Perception
- Part III: Existing Deep Learning solutions to Challenges in Perception
- Part IV: Remaining Challenges and Future Directions

# A Holistic View of Perception in Intel. Vehicles

## Part II: Deep Learning for Perception



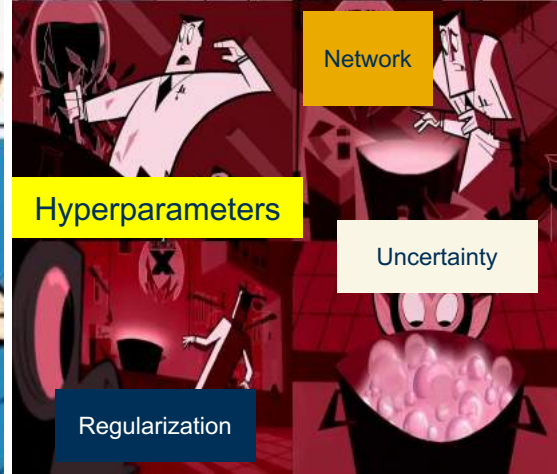
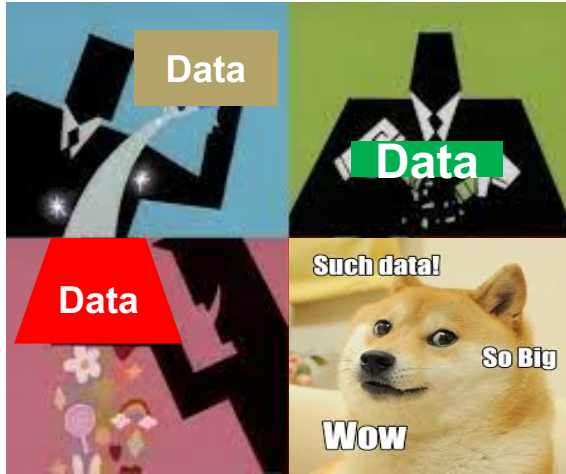
# Objectives

## Objectives in Part II

- Discuss myths surrounding deep learning
- Brief history of deep learning
- Review deep learning models for vision
- Deep learning extensions into sensor domain
- Transfer Learning and foundation models
- Self-supervised learning
- Case study: Self-supervised learning for fisheye images

# Deep Learning

Meme to start off with



Expectation

Reality

Generalizable

Explainable

Robust



# Deep Learning

Meme to start off with

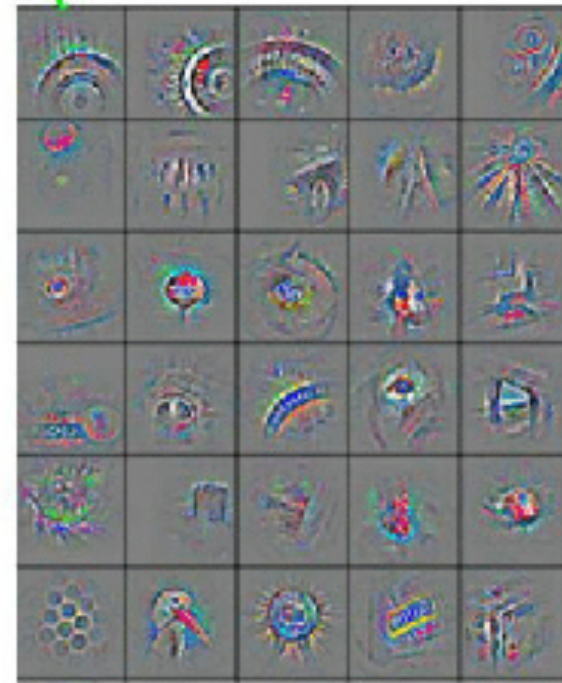
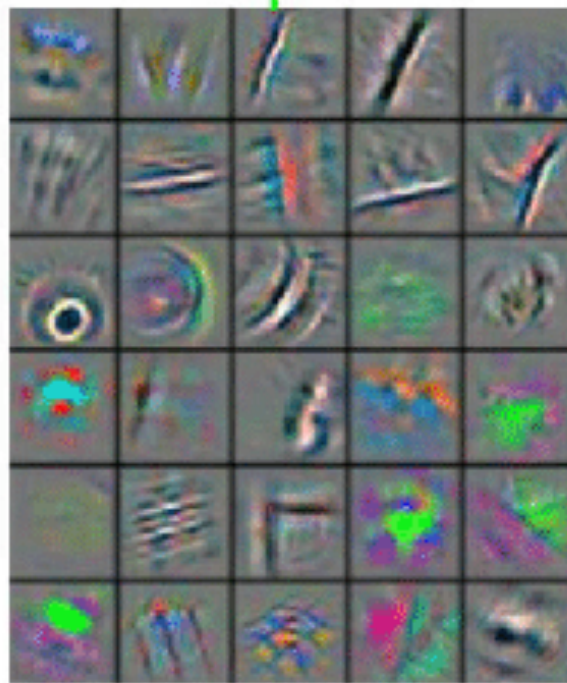
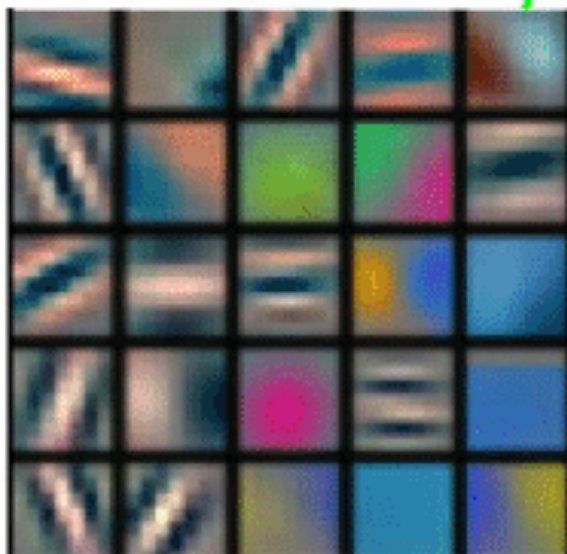
People's expectation of AI and Deep Learning



[Tutorial] | [Ghassan AlRegib and Mohit Prabhushankar] | [June 4, 2023]

# Deep Learning

## Model Decomposition



**Ex. LeCun, 2015**

# Deep Learning

## Some Common Myths about Deep Learning

*“Deep learning is hard to train”*

# PyTorch 2.0

pytorch

PyTorch Conference 2023

October 16 - 17 | San Francisco, CA | #pytorchconf

### Convolution Layers

109,392 repository results

`nn.Conv1d`

Applies a 1D convolution over an input signal composed of several input planes.

`nn.Conv2d`

Applies a 2D convolution over an input signal composed of several input planes.

`nn.Conv3d`

Applies a 3D convolution over an input signal composed of several input planes.

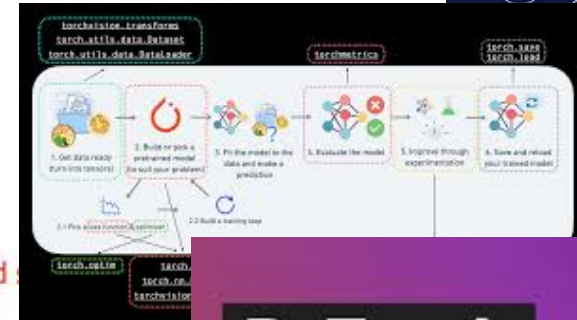
`nn.ConvTranspose1d`

Applies a 1D transposed convolution operator over an input image composed of several input planes.

`nn.ConvTranspose2d`

Applies a 2D transposed convolution operator over an

- Containers
- Convolution Layers
- Pooling layers
- Padding Layers
- Non-linear Activations (weighted)
- Non-linear Activations (other)
- Normalization Layers
- Recurrent Layers
- Transformer Layers
- Linear Layers



# PyTorch



## CRASH COURSE

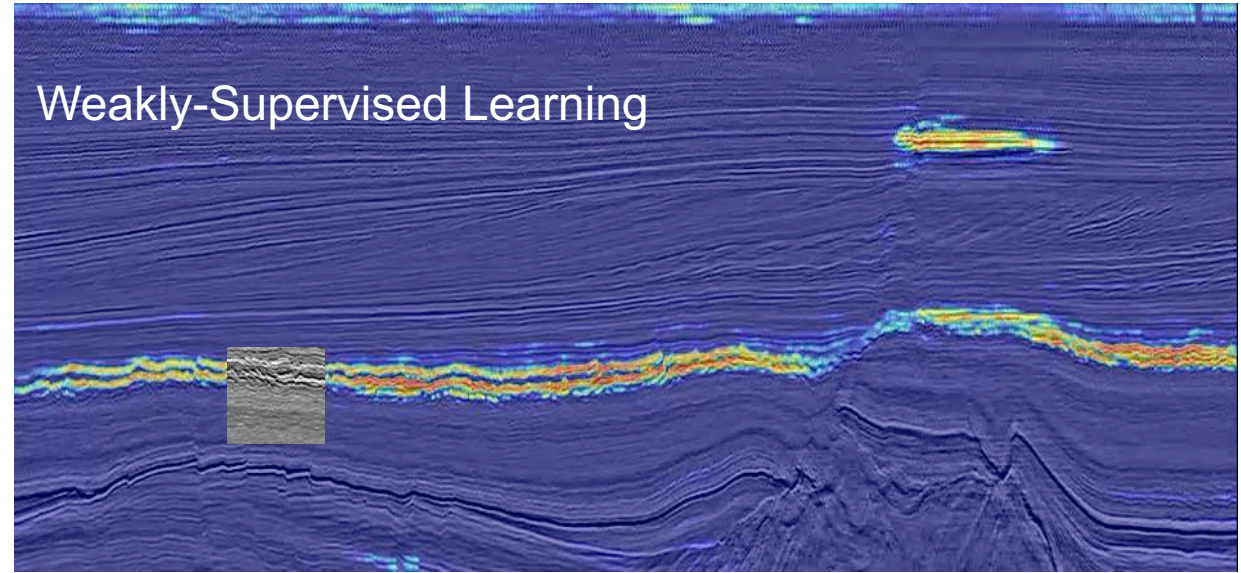
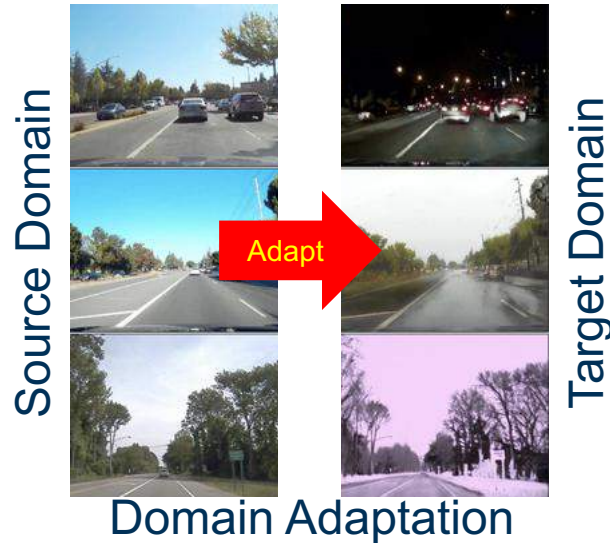
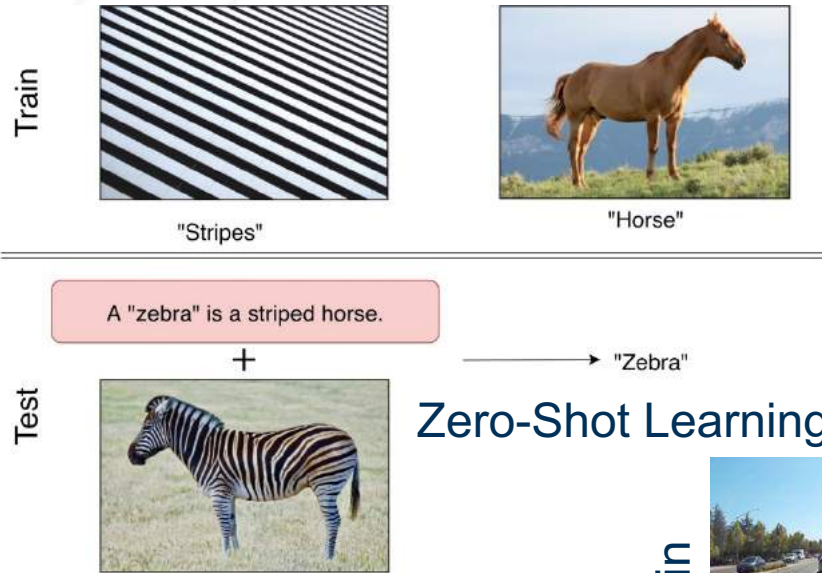
ZERO TO HERO IN 50 MINUTES



# Deep Learning

## Some Common Myths about Deep Learning

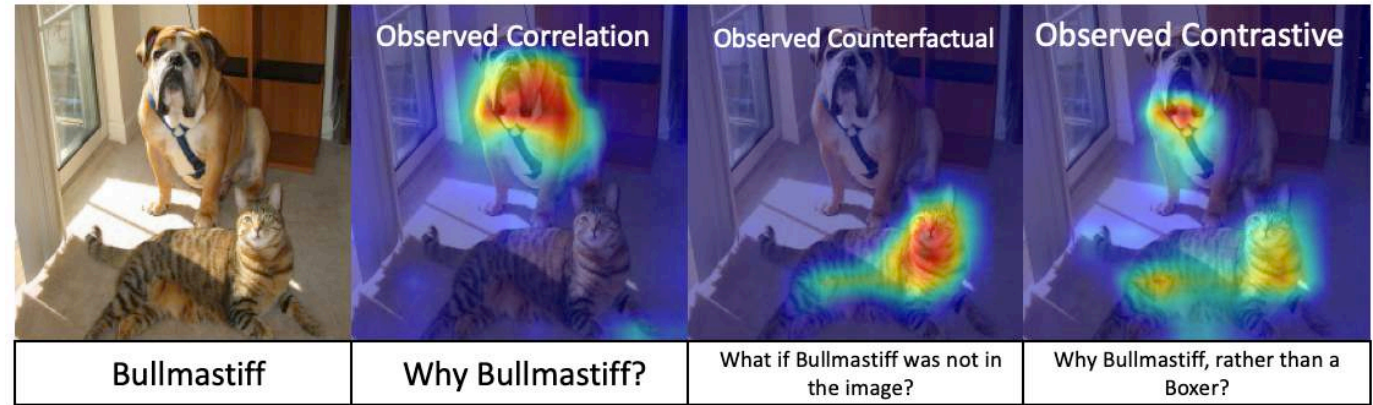
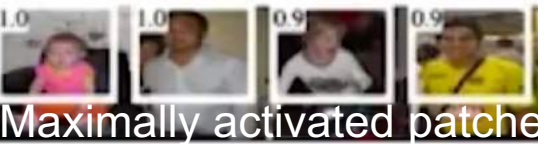
*“Deep learning requires lots of data”*



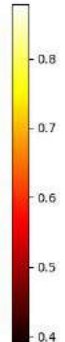
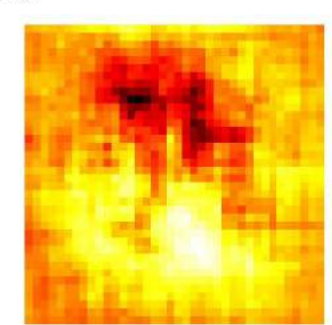
# Deep Learning

## Some Common Myths about Deep Learning

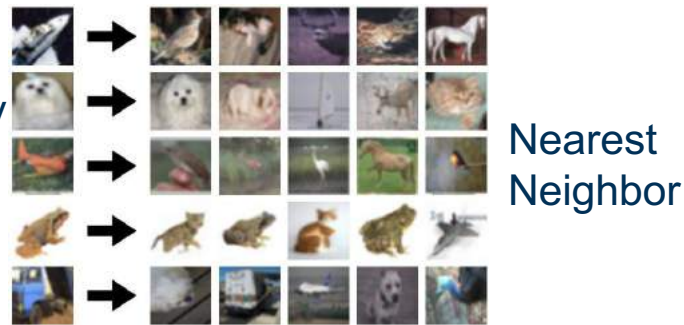
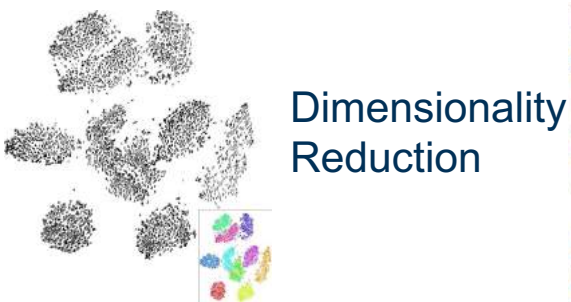
*“Deep learning has poor interpretability”*



African elephant, *Loxodonta africana*



Saliency via occlusion

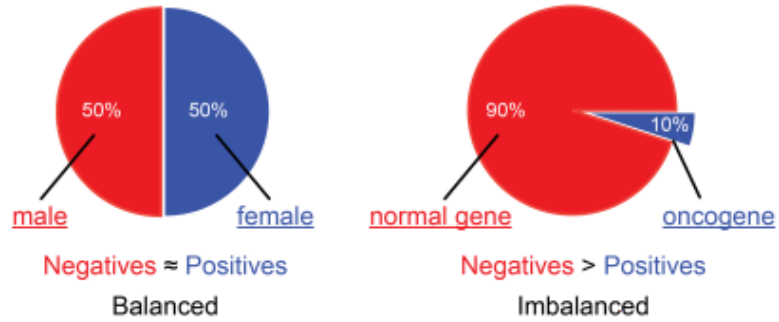


# Deep Learning

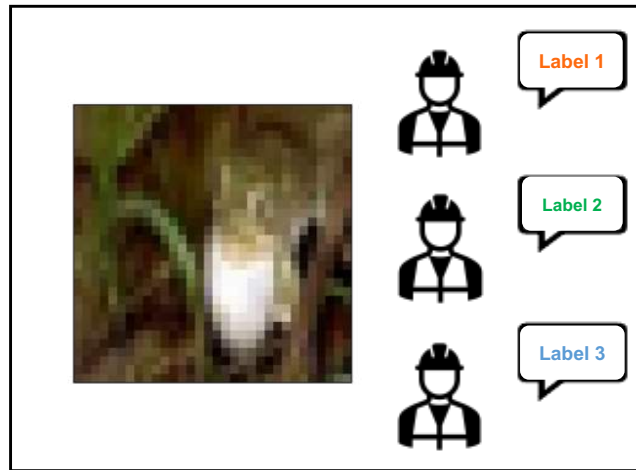
## Some Common Myths about Deep Learning

*“More the data, better the model”*

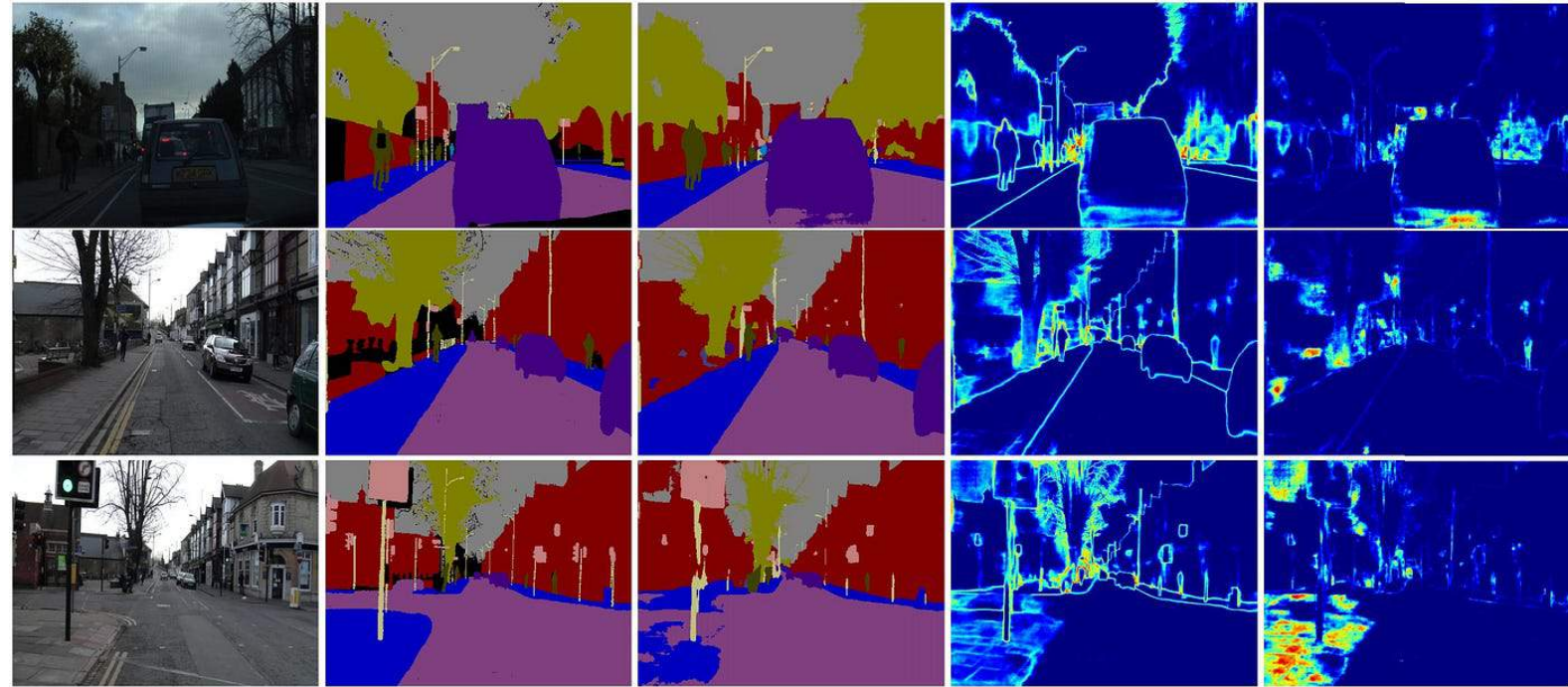
Example of balanced and imbalanced data



### Data imbalance issues



### Human labeling issues



(a) Input Image

(b) Ground Truth

(c) Semantic Segmentation

(d) Aleatoric Uncertainty

(e) Epistemic Uncertainty

### Dataset uncertainties



# Deep Learning

## Some Common Myths about Deep Learning

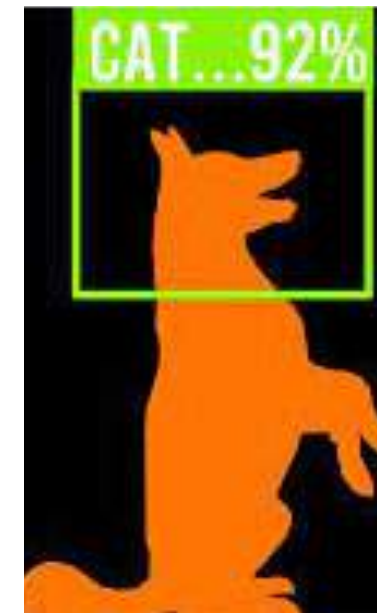
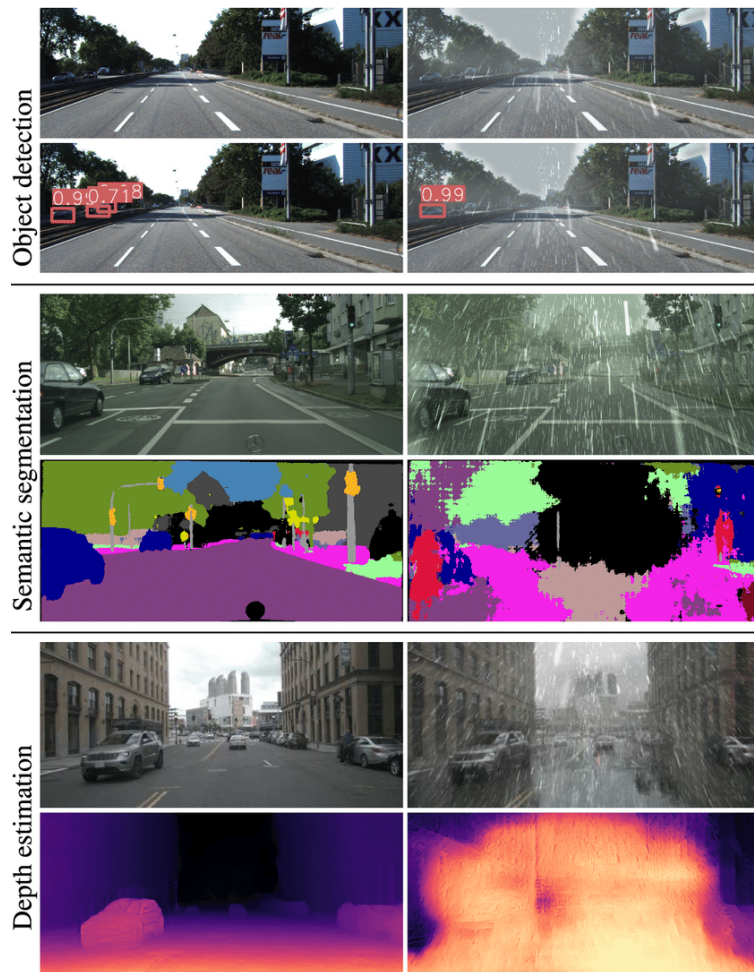
*“Deep learning is State-of-the-Art in every field”*



241 - (-241) + 1



241 - (-241) + 1 is equivalent to 241 + 241 + 1, which simplifies to 483 + 1. So 241 - (-241) + 1 is equal to 484.



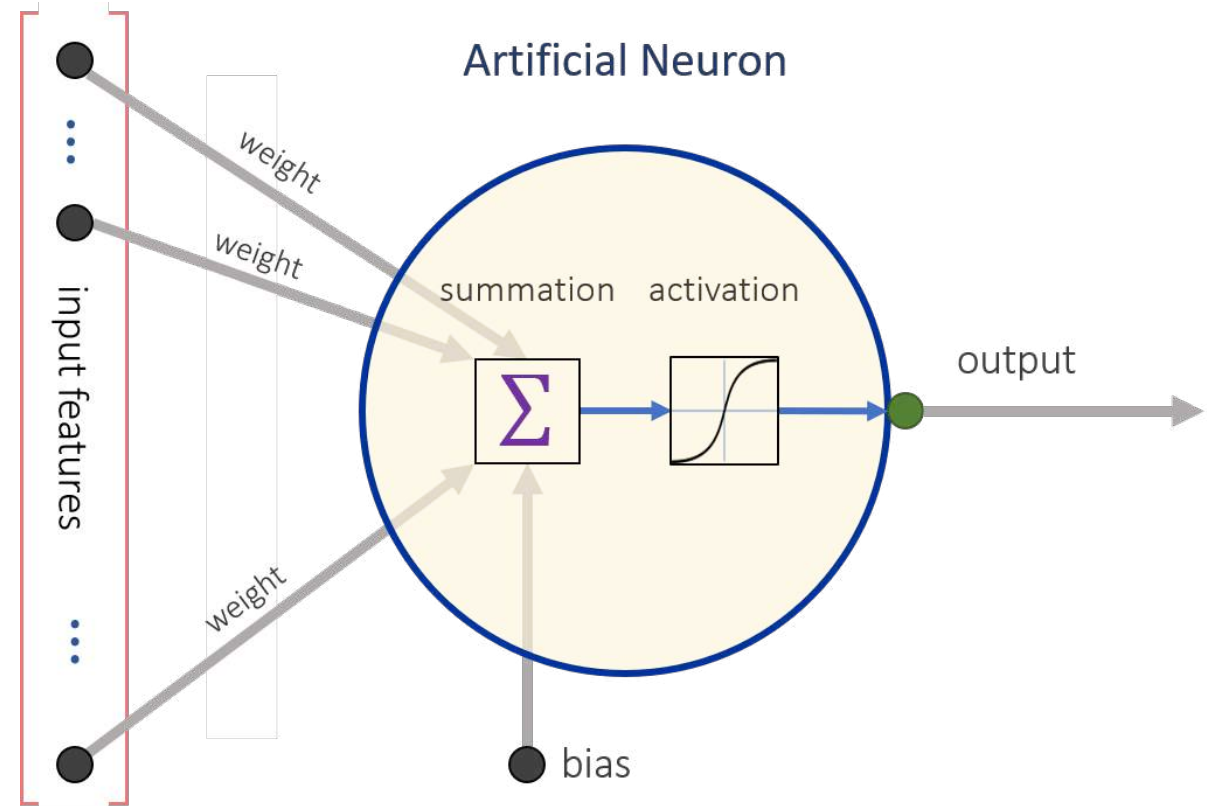
# Deep Learning

## The Building Block

The underlying computational unit is the artificial neuron

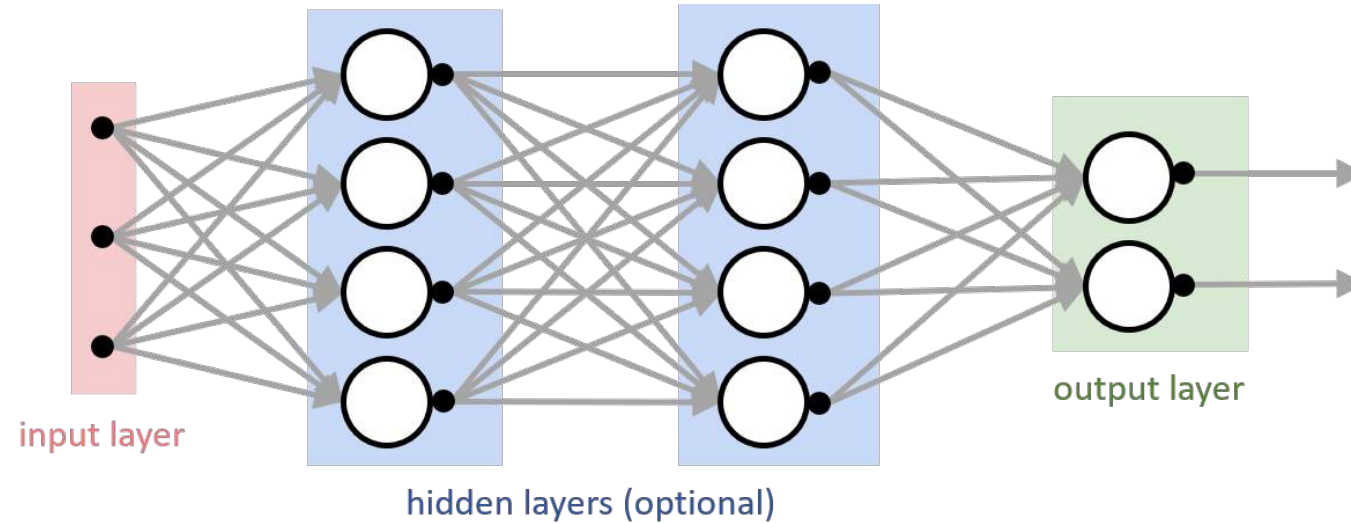
Artificial neurons consist of:

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



# Deep Learning

## Artificial Neural Networks

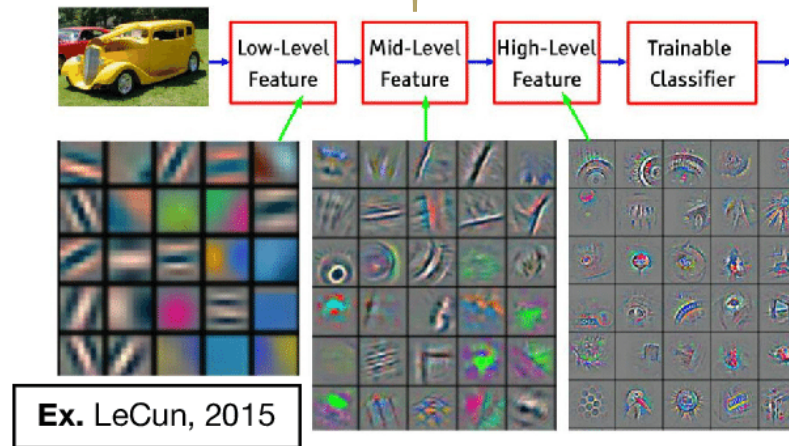
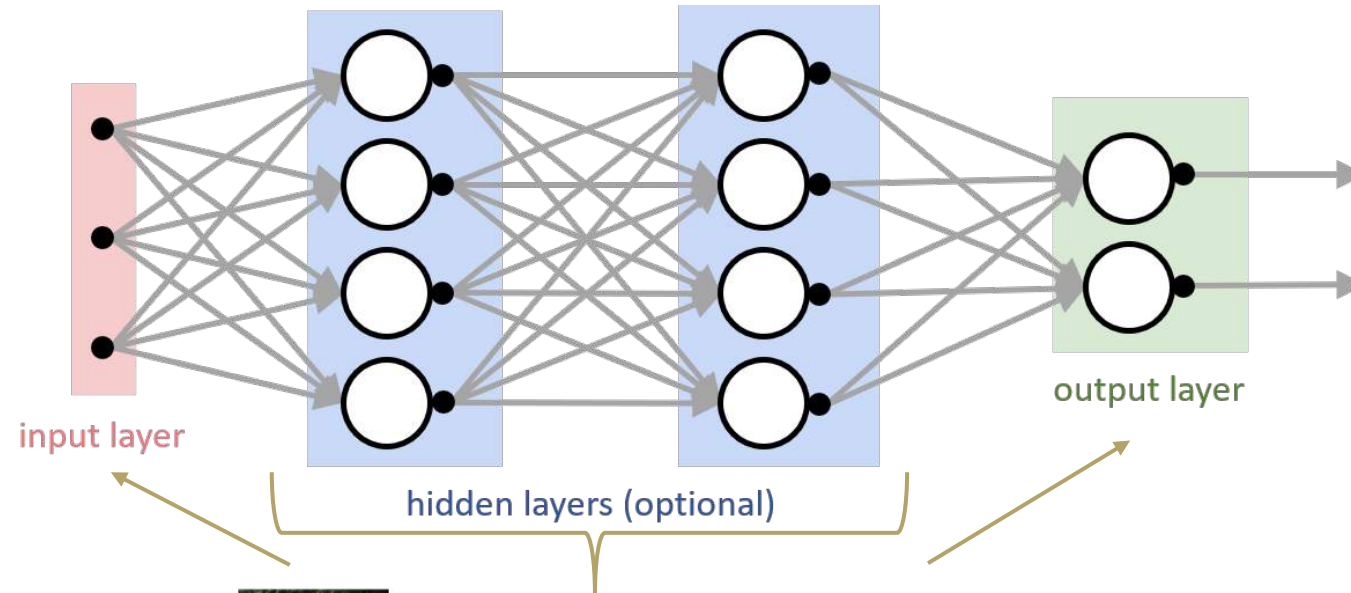


Typically, a neuron is part of a network organized in layers:

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

# Deep Learning

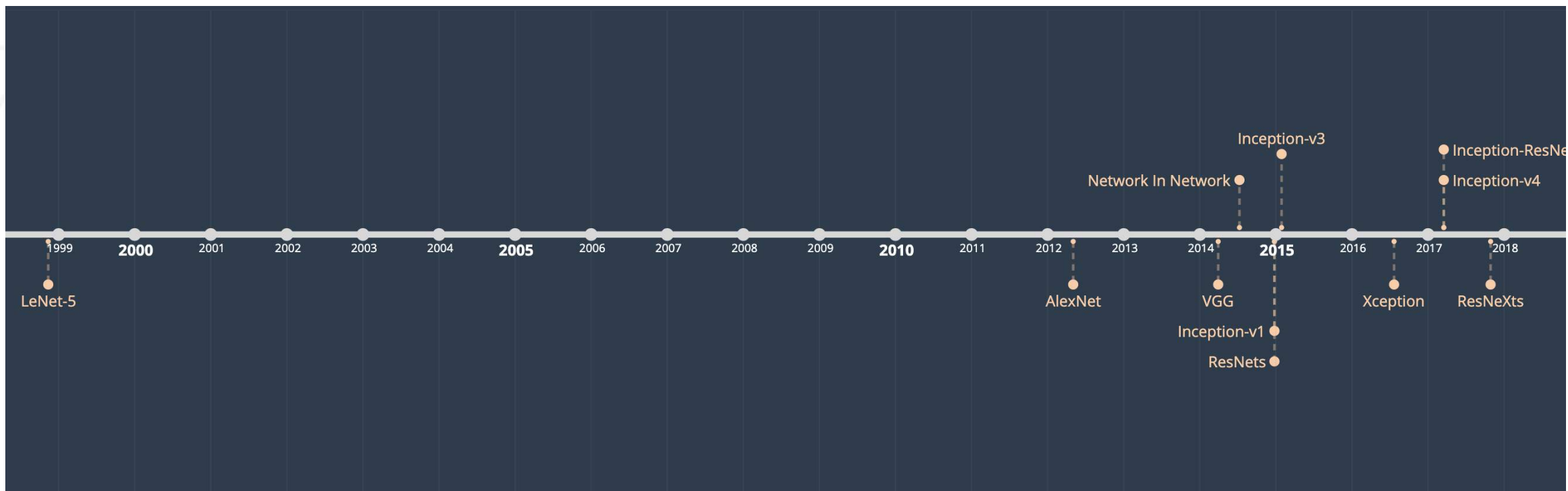
## Convolutional Neural Networks



# Deep Learning

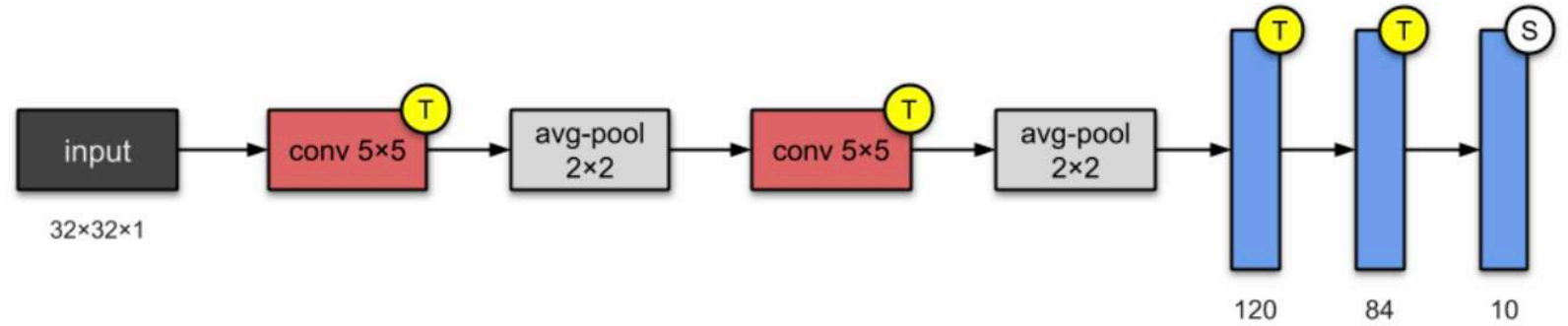
## Evolution of CNN Architectures

- LeNet
- AlexNet
- VGG
- GoogLeNet (Inception-V1)
- ResNet



# CNN Architectures

## LeNet5 (1998)



### Novelty:

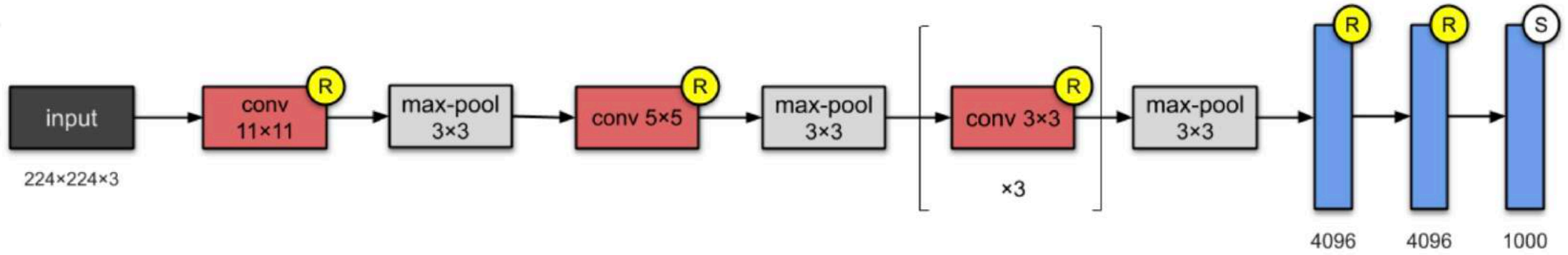
- Reduced number of learnable parameters and learned from raw pixels automatically
- The 1<sup>st</sup> popular CNN that became the “standard” template of CNNs
  - Stacking convolutional, activation, pooling layers
  - Ending with fully connected layers
- Good results on small datasets
  - Top-5 error rate on MNIST is 0.95%

# Long Gap (1998 – 2012)

- Working to improve **computational power**
  - Existing accelerators were not yet sufficiently powerful to make deep multichannel, multilayer CNNs with a large number of parameters.
- Existing **datasets** were relatively **small**
  - Limited storage capacity of computers
- **Tricks for neural network training** were not established yet
  - Parameter initialization
  - Variants of stochastic gradient descent
  - Non-squashing activation functions
  - Effective regularization techniques

# CNN Architectures

## AlexNet (2011)



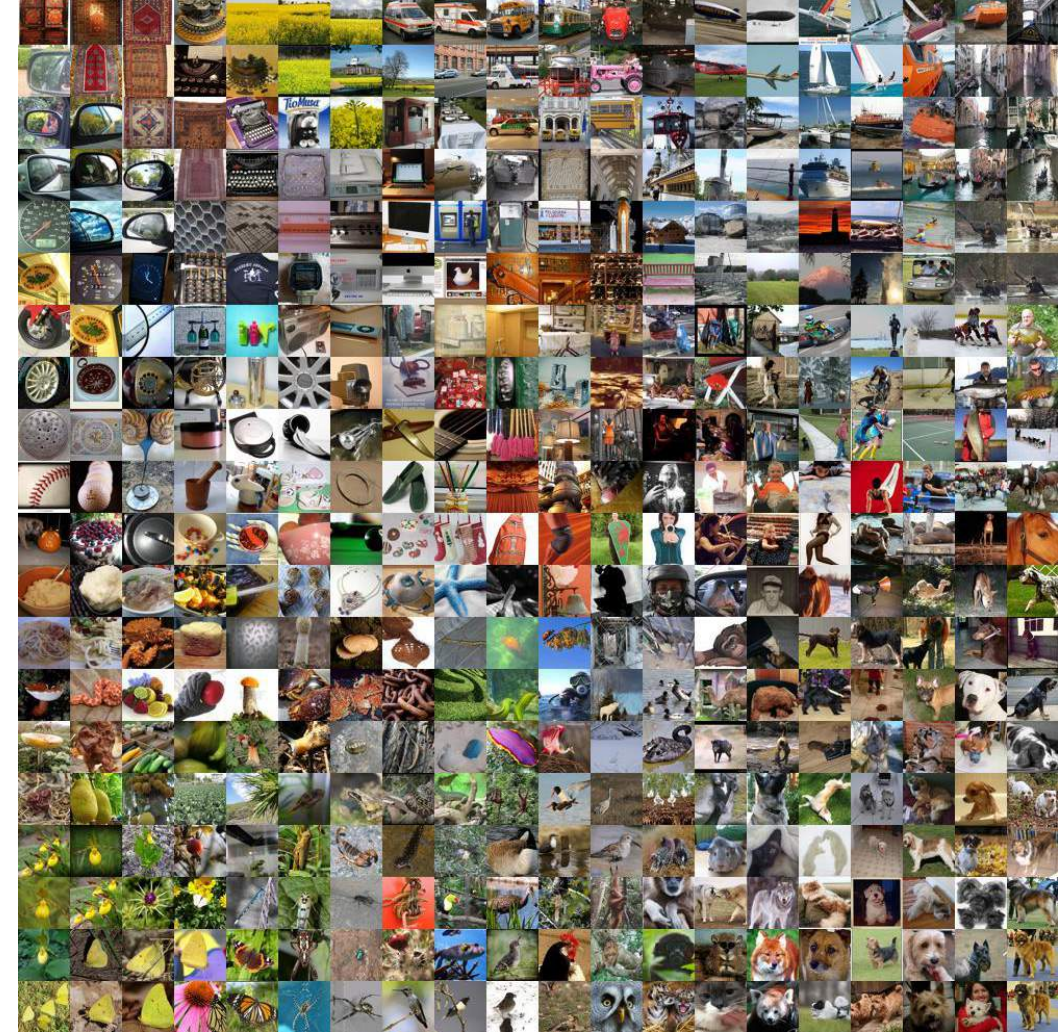
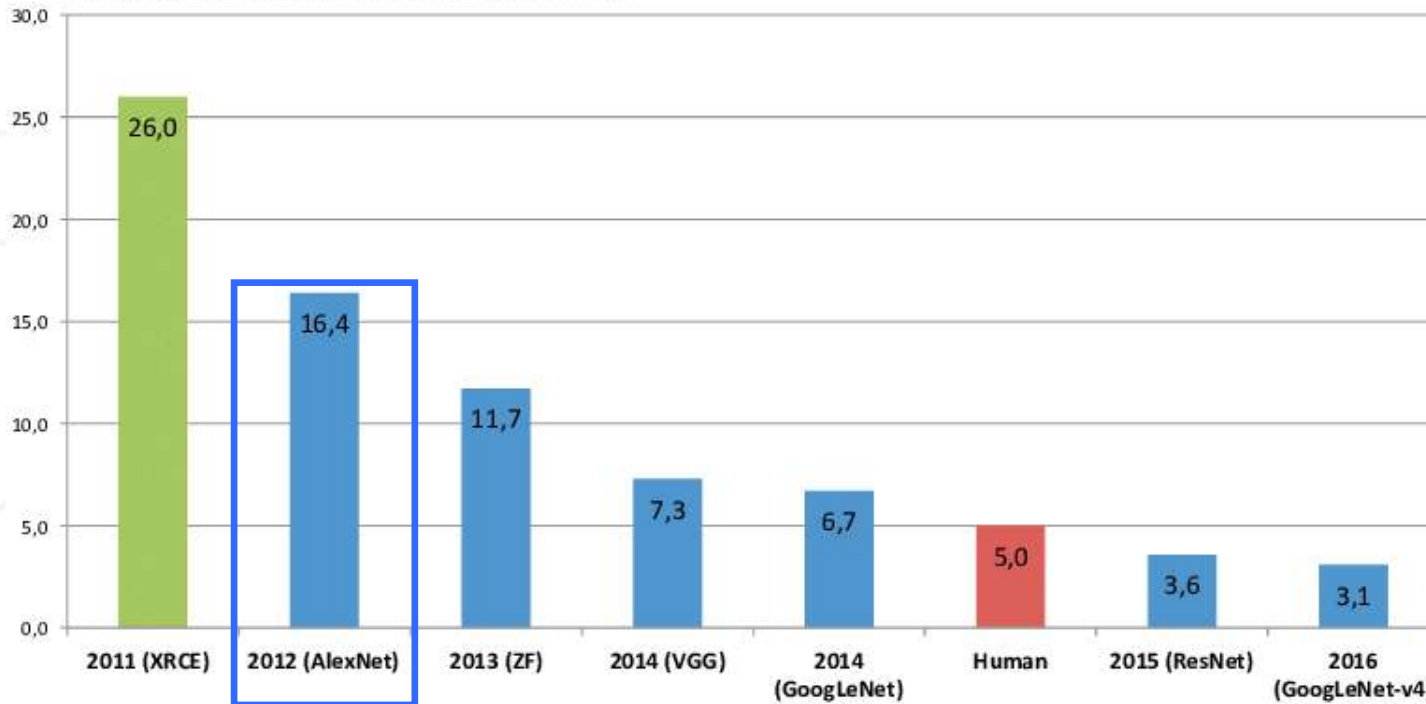
### Novelty:

- First to implement Rectified Linear Units (ReLUs) as activation, solving the vanishing gradient problem
- Applied dropout regularization to fully connected layer to control complexity
- Deep CNN that runs on GPU hardware
- Deeper and wider than LeNet
- More robust than LeNet (data augmentation)
- **Won ImageNet Challenge and significantly outperformed traditional methods**



# AlexNet (2012)

ImageNet Classification Error (Top 5)



Imagenet:  
1000 classes, 1.2M training images, 150K for testing

## 16.4% top 5 error in ILSVRC 2012

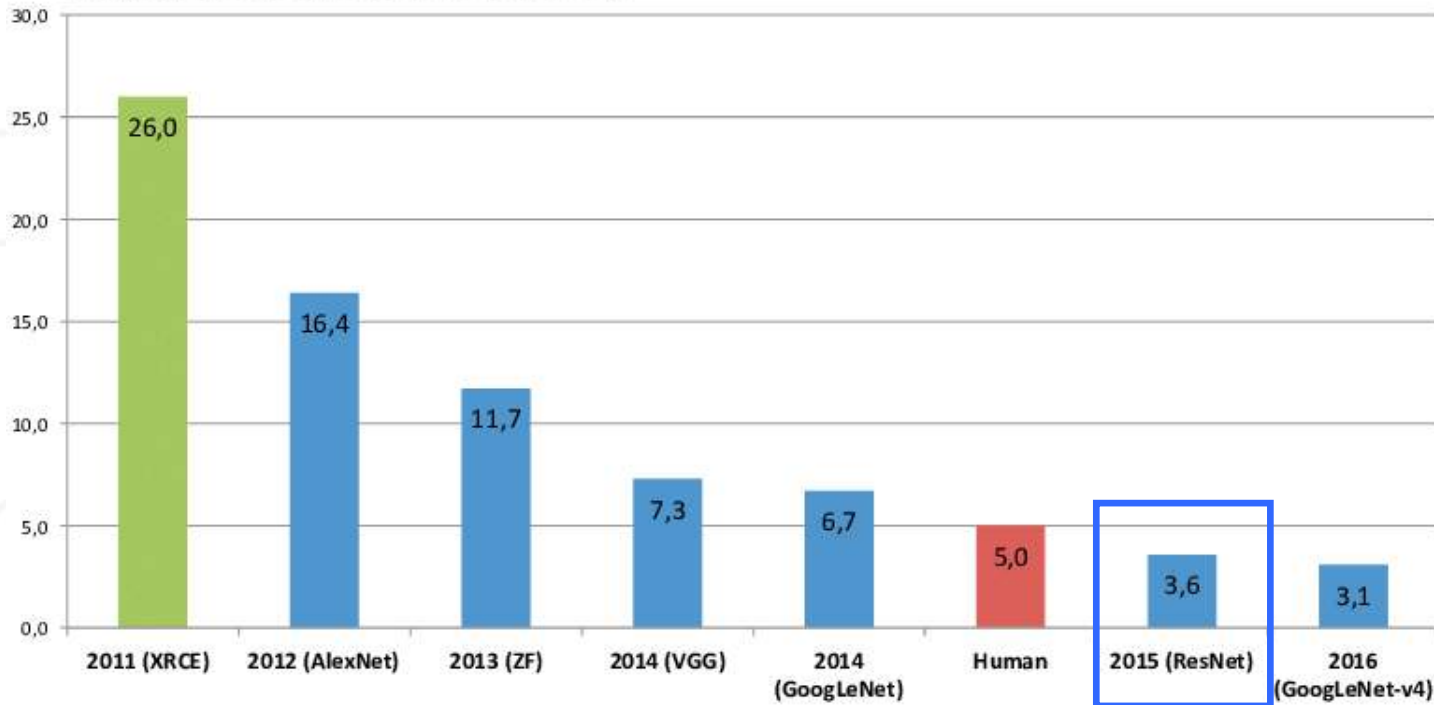
Figure Credit: Zitzewitz, Gustav. "Survey of neural networks in autonomous driving." (2017)

Krizhevsky et al. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems* 25 (2012)

[Tutorial] | Ghassan AlRegib and Mohit Prabhushankar | June 4, 2023

# ResNet (2015)

ImageNet Classification Error (Top 5)



~3.6% top 5 error in ILSVRC 2015,  
lower than human recognition error!

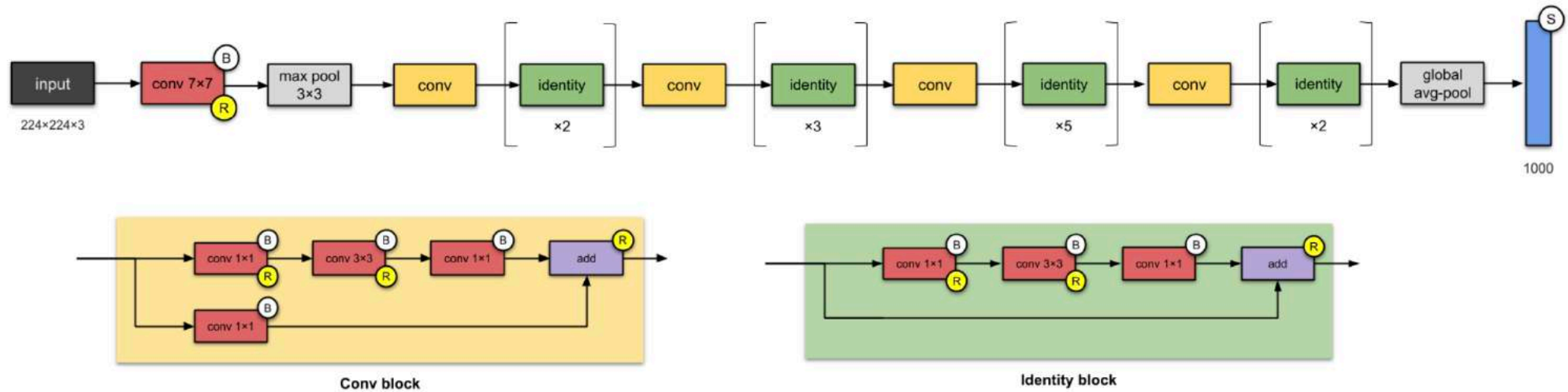
Figure Credit: Zitzewitz, Gustav. "Survey of neural networks in autonomous driving." (2017)



Imagenet:  
1000 classes, 1.2M training images, 150K for testing

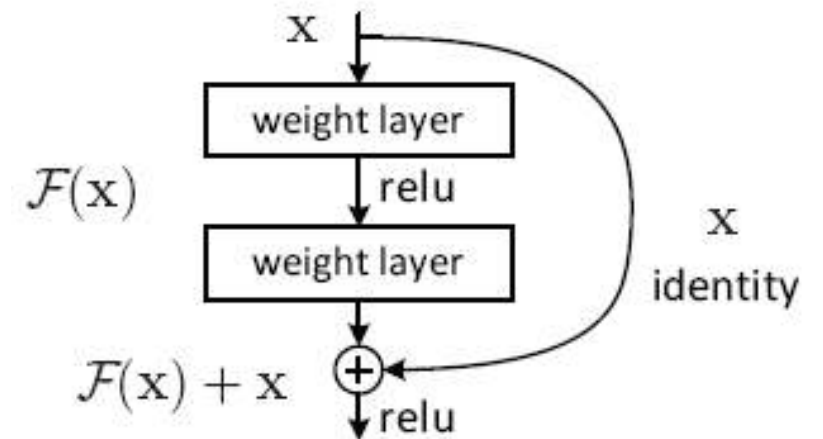
# CNN Architectures

## ResNet (2015)



### Novelty:

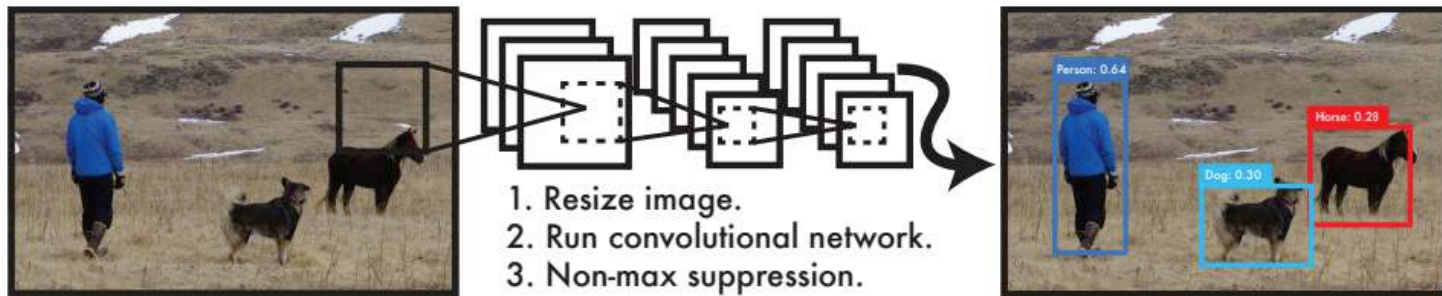
- Introduced residual learning (Residual blocks)
  - Shortcut connections with identity mapping
- Popularized skip connections
- 20 and 8 times deeper than AlexNet and VGG, respectively with less computational complexity and without compromising generalization power



# Object Detection Architectures

## YOLO (2016 - Ongoing)

All previous object detection techniques required multiple stages of detection



**Figure 1: The YOLO Detection System.** Processing images with YOLO is simple and straightforward. Our system (1) resizes the input image to  $448 \times 448$ , (2) runs a single convolutional network on the image, and (3) thresholds the resulting detections by the model's confidence.

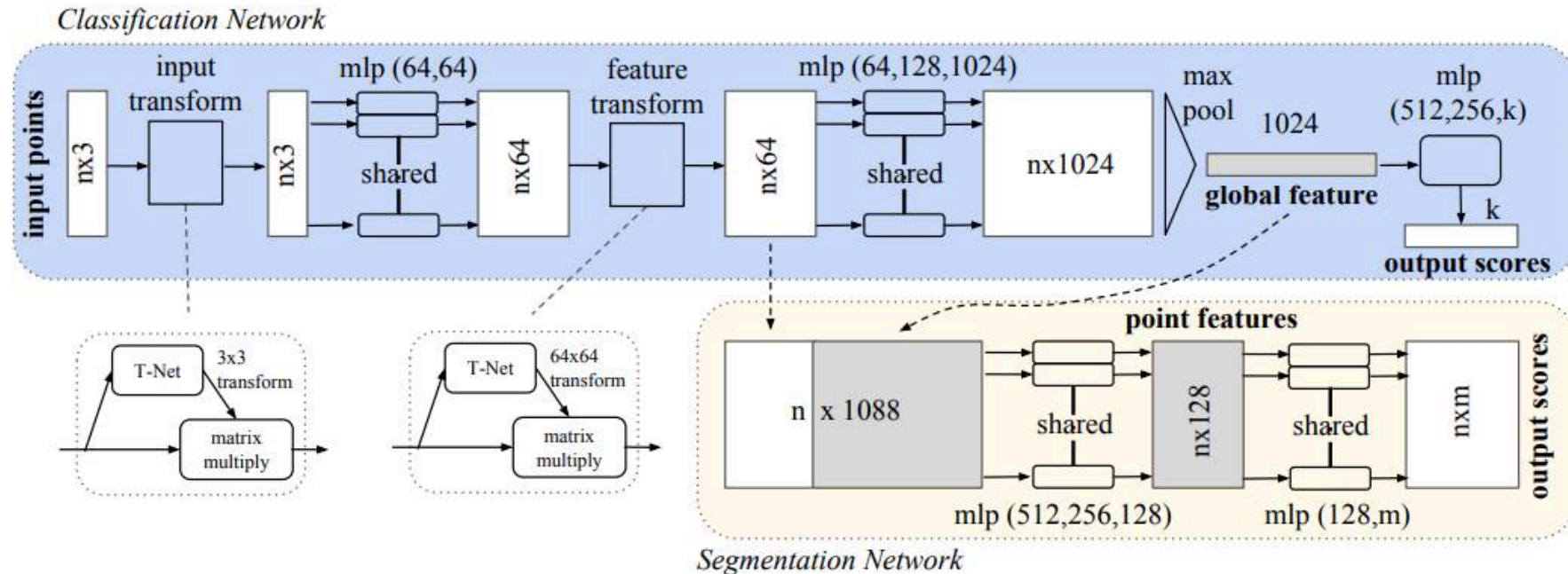
Novelty:

- Object detection is reformulated as a regression problem from image space to bounding-box coordinate space
- Single stage object detectors
  - Feature extraction, detection, classification performed in one go
- Contextual information is encoded within each prediction

# Deep Learning for LIDAR data

## PointNet (2017)

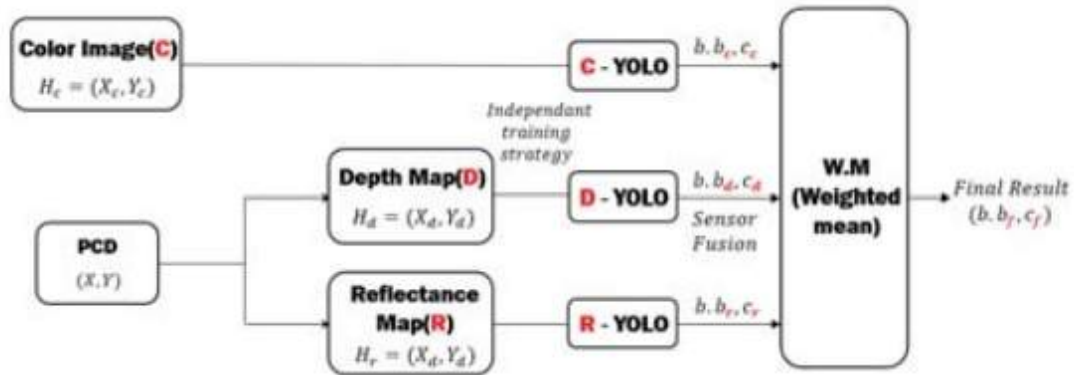
The challenge in utilizing LIDAR data is the volume of point cloud data and the permutation of their processing



- Performed classification and segmentation on  $n$  points of LIDAR data. Input  $n \times 3$  refers to  $n$  points with  $\{x, y, z\}$  coordinate dimensions
- Used RNNs to overcome the permutation issues within LIDAR data

# Deep Learning for Sensor Fusion

## Vision and LIDAR



YOLO Framework is used to independently extract features from cameras and LIDAR sensors and fused to detect missed boxes

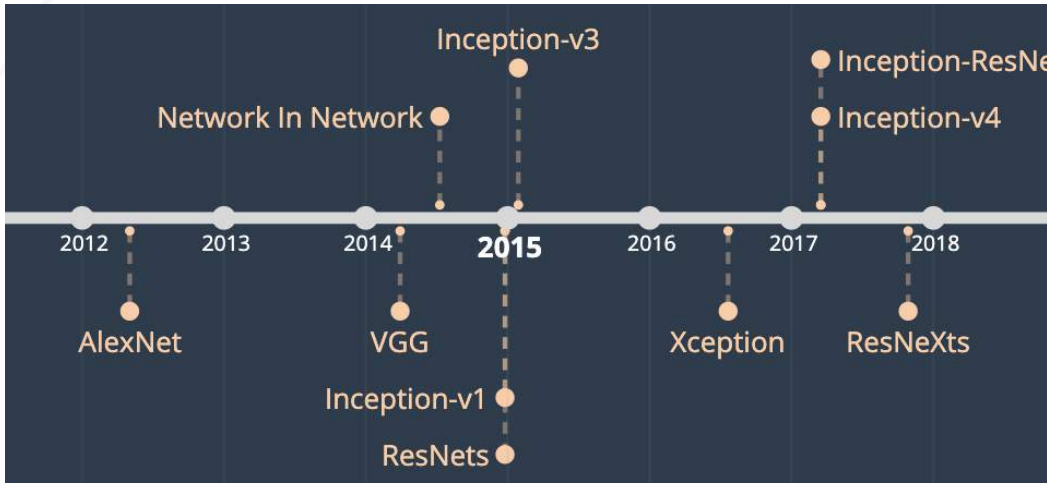
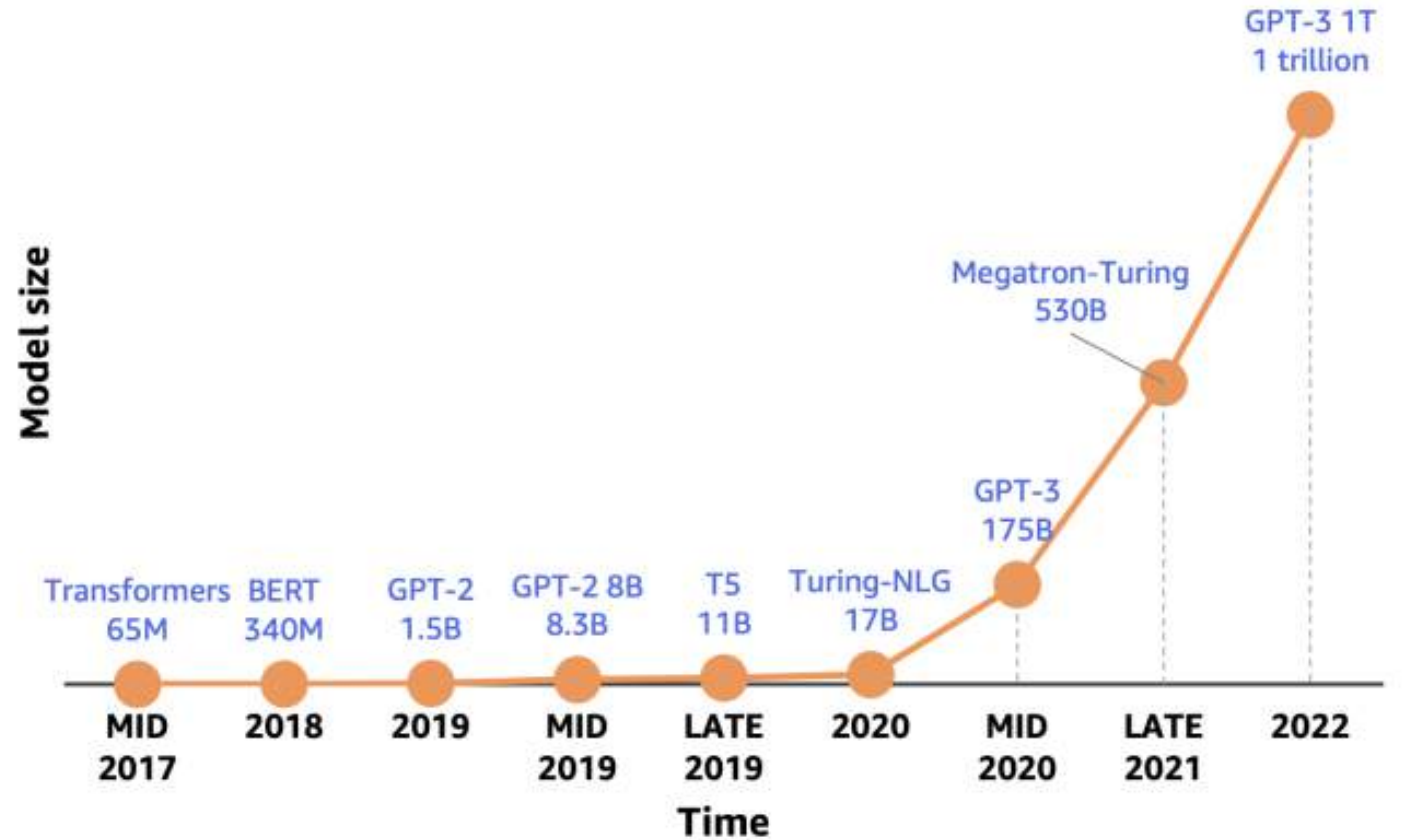
This is *'late fusion'*, in the sense that each sensor modality is independently evaluated

# Deep Deep Deep ... Deep Deep Learning

## Recent Advancements

The number of parameters in models has increased exponentially

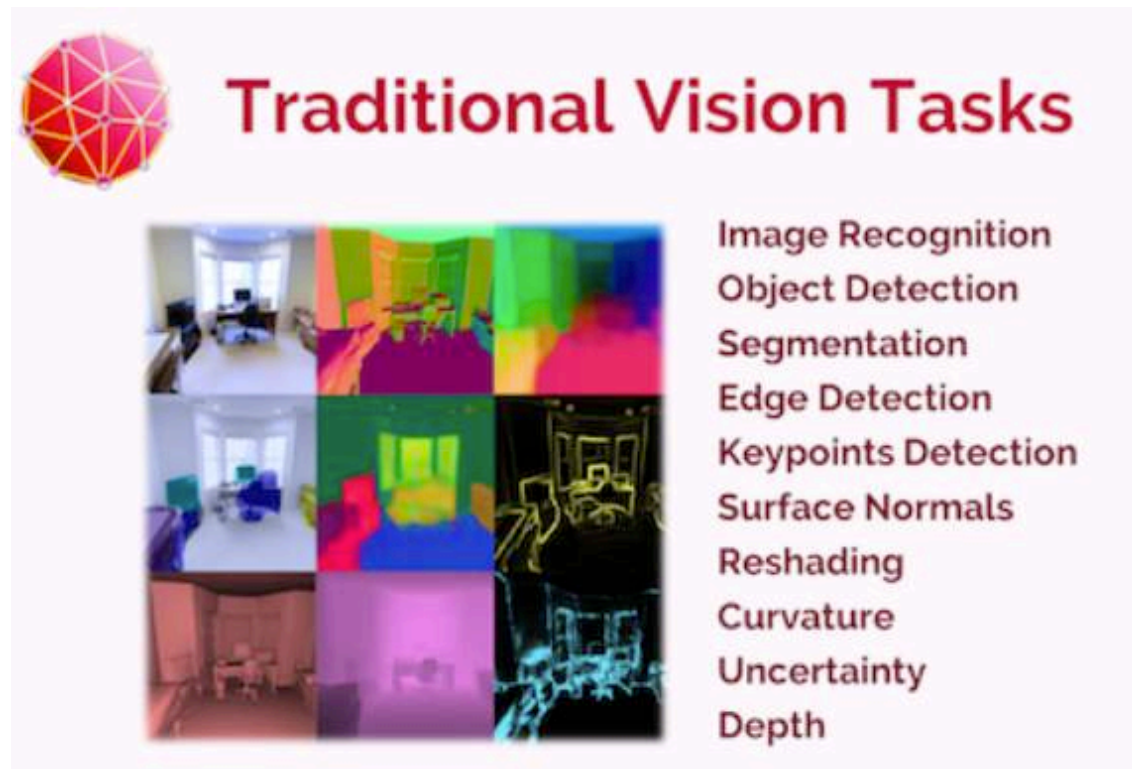
15,000x increase in 5 years



# Deep Deep Deep ... Deep Deep Learning

## Motivation

**Underlying features among different vision tasks are similar**



**Traditional Vision Tasks**

- Image Recognition
- Object Detection
- Segmentation
- Edge Detection
- Keypoints Detection
- Surface Normals
- Reshading
- Curvature
- Uncertainty
- Depth

**This similarity leads to Transfer Learning**

[Tutorial] | [Ghassan AlRegib and Mohit Prabhushankar] | [June 4, 2023]

Bommasani, Rishi, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein et al. "On the opportunities and risks of foundation models." *arXiv preprint arXiv:2108.07258* (2021).





# Transfer Learning

## What is Transfer Learning?

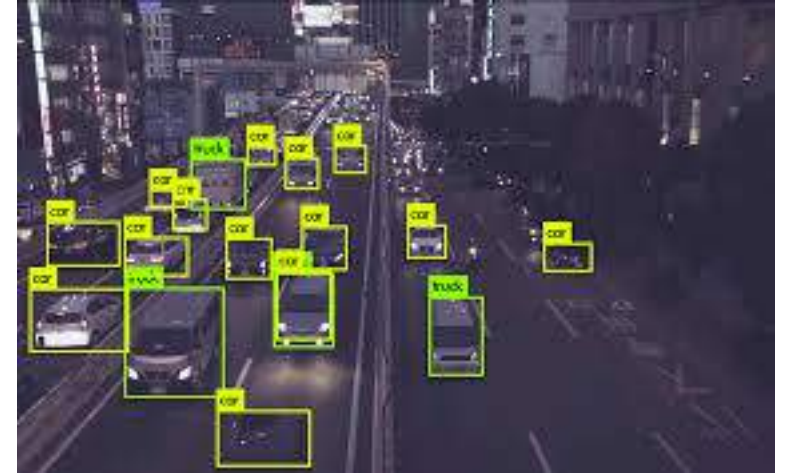
- Deep networks tend to **learn common representations** for various tasks in their earlier layers
- Can be exploited **to transfer representations from networks trained on large datasets** on one task (i.e., Image Classification on ImageNet) called the *source* to a different task called the *target* task
- Usually done by **taking large pretrained network** and then **finetuning last layer** (with all other layers frozen) on target dataset
- **Pre-trained frozen backbone** acts as a **feature extractor** while **finetuned last layer** acts to project the representations into the **decision boundary for the target task**
- Utility depends on how closely related the source and target datasets and/or tasks are

# Transfer Learning

## Foundation Models



Source: <https://gluon-cv.mxnet.io/>



Source: <https://www.move-lab.com/blog/tracking-things-in-object-detection-videos>



Pretraining



Foundation Model



Finetuning

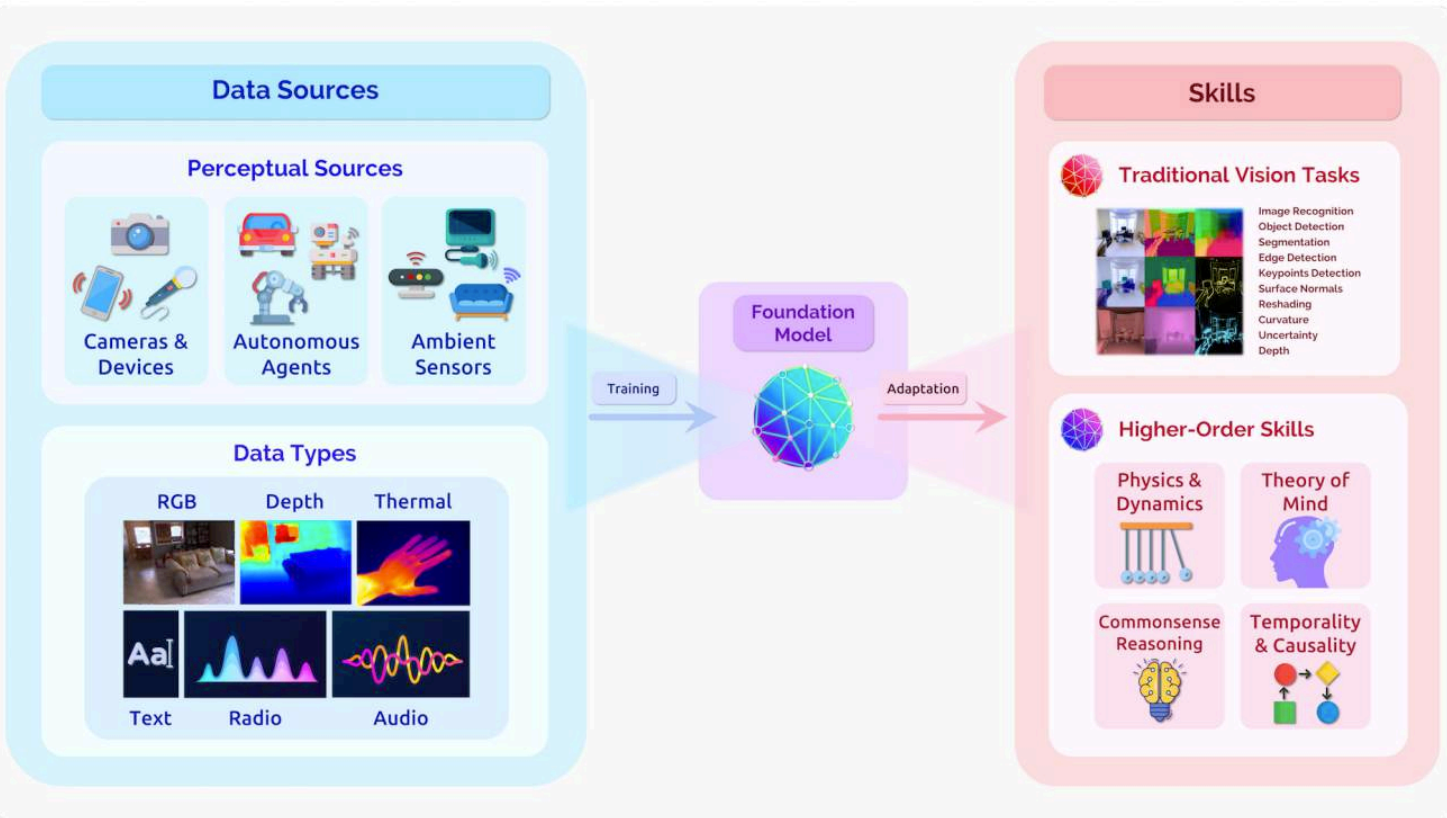
# Foundation Models

## Origin of the term Foundation Models

- **Foundation models** are like any other deep network that have employed **transfer learning**, except at **scale**
- **Scale** brings about **emergent properties** that are common between tasks
- **Before 2019**: Base architectures that powered multiple neural networks were **ResNets, VGG** etc.
- **Since 2019**: **BERT, DALL-E, GPT, Flamingo**
- Changes since 2019: **Transformer architectures and Self-Supervision**

# Foundation Models

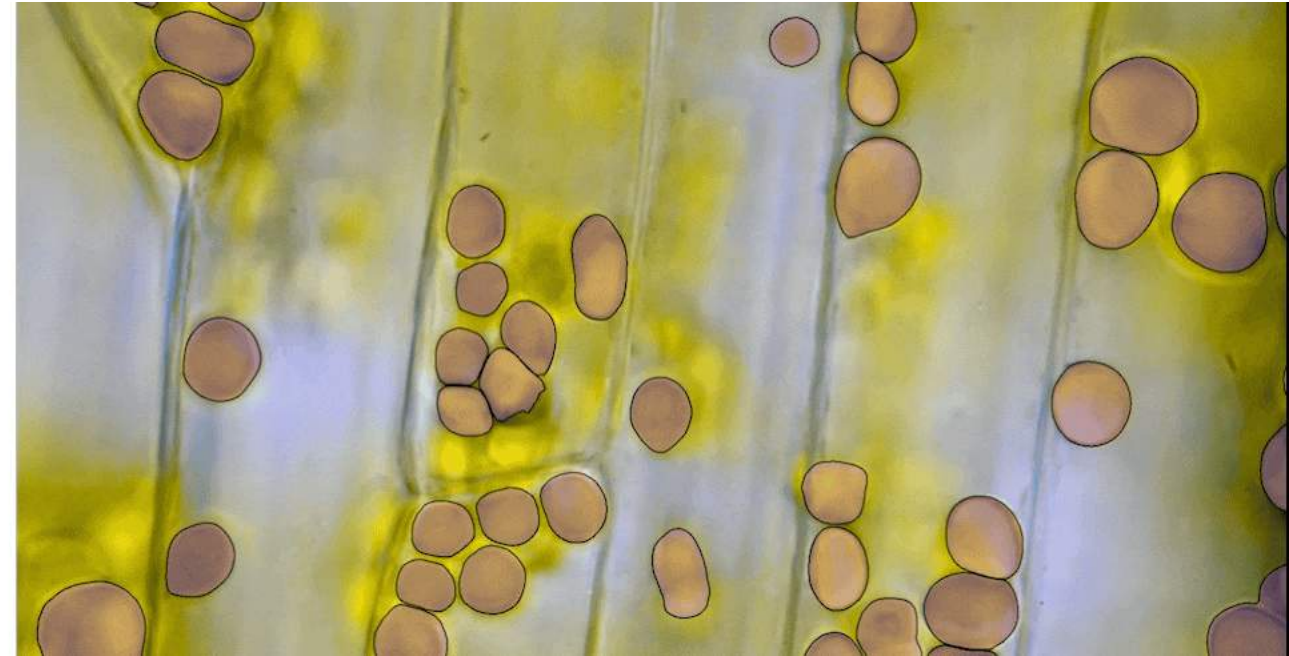
## Origin of the term Foundation Models



*'By harnessing self-supervision at scale, foundation models for vision have the potential to distill raw, multimodal sensory information into visual knowledge, which may effectively support traditional perception tasks and possibly enable new progress on challenging higher-order skills like temporal and commonsense reasoning. These inputs can come from a diverse range of data sources and application domains, suggesting promise for applications in healthcare and embodied, interactive perception settings.'*

# Foundation Models

## Segment Anything Model



Segment Anything Model (SAM) released by Meta on April 5, 2023 was trained on Segment Anything 1 Billion dataset with 1.1 billion high-quality segmentation masks from 11 million images

# Foundation Models

## Segment Anything Model



Cityscapes dataset  
semantic segmentation  
annotation took ~90  
mins per image

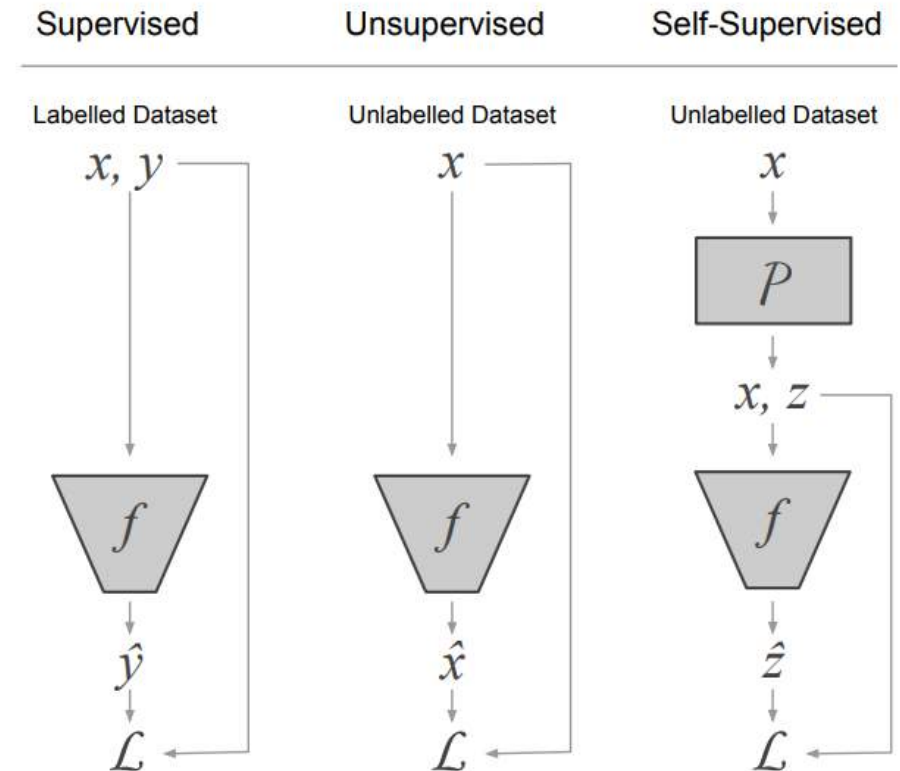
# Foundation Models

## Training Foundation Models

### Foundation models are trained via Self-Supervision

#### Self-Supervision:

- Type of unsupervised learning
- Primary difference is the introduction of a **“pre-text task.”**
- The pre-text task generates pseudo-labels that are used to train a network.



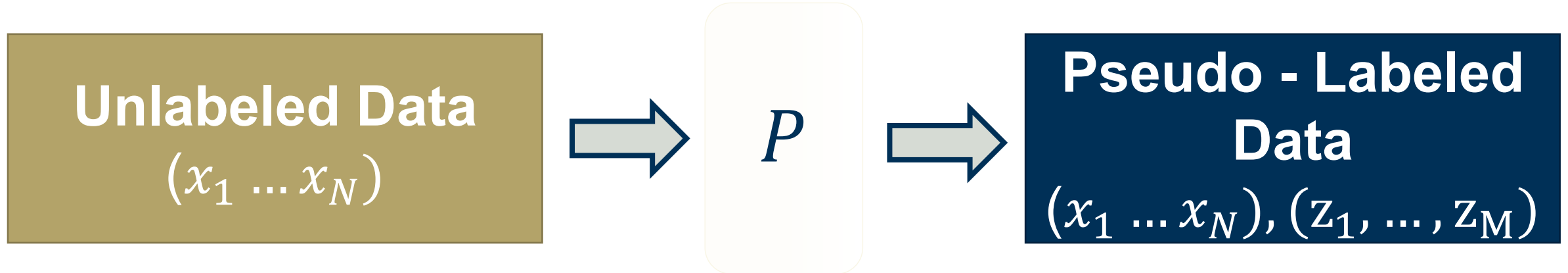
# Self-Supervision

## Overall Training Process

### 1. Identify Labeled and Unlabeled Data



### 2. Generate pseudo-labels with some pre-text task $P$

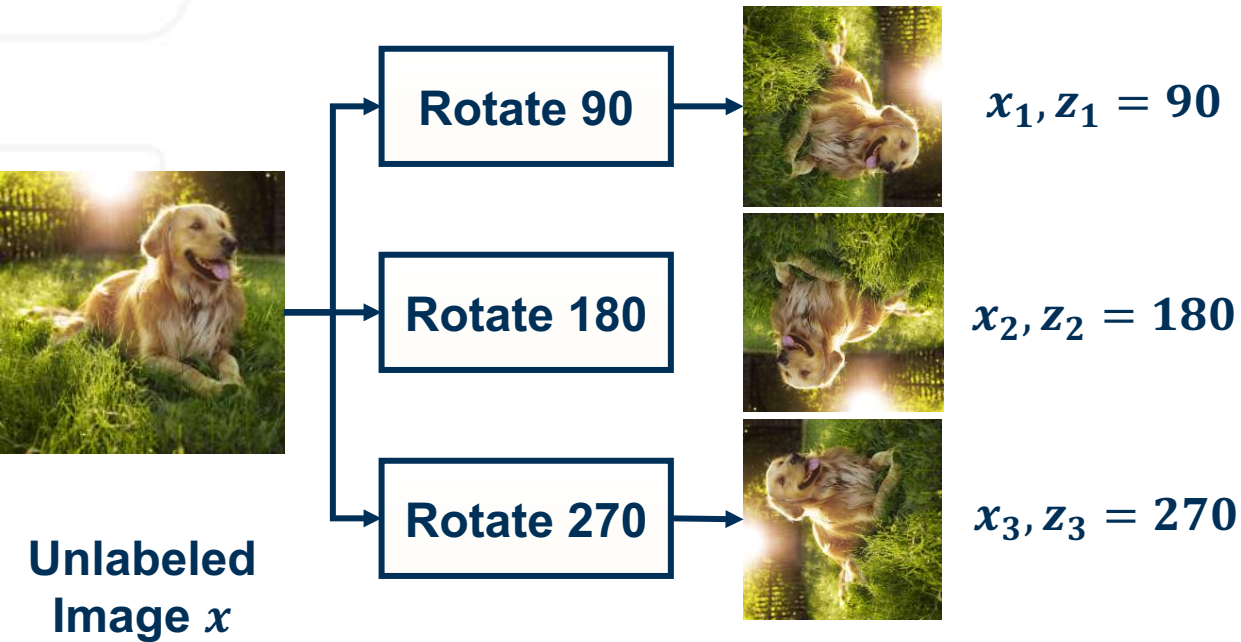




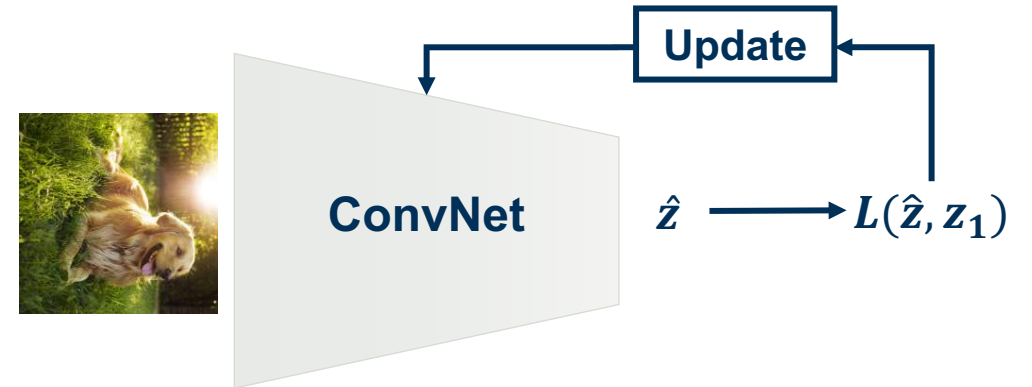
# Self-Supervision

## Example Training Process

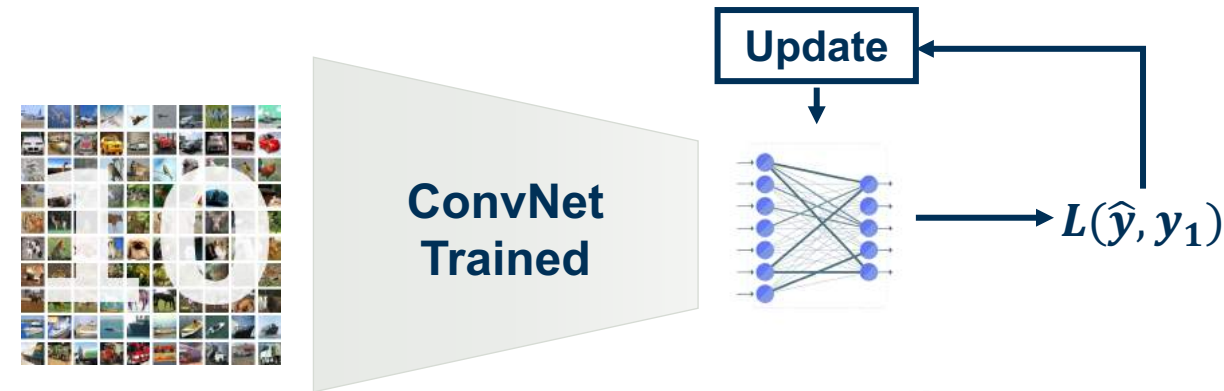
### Step 1: Generate pseudo-labels via image rotations



### Step 2: Network learns to predict angle image is rotated



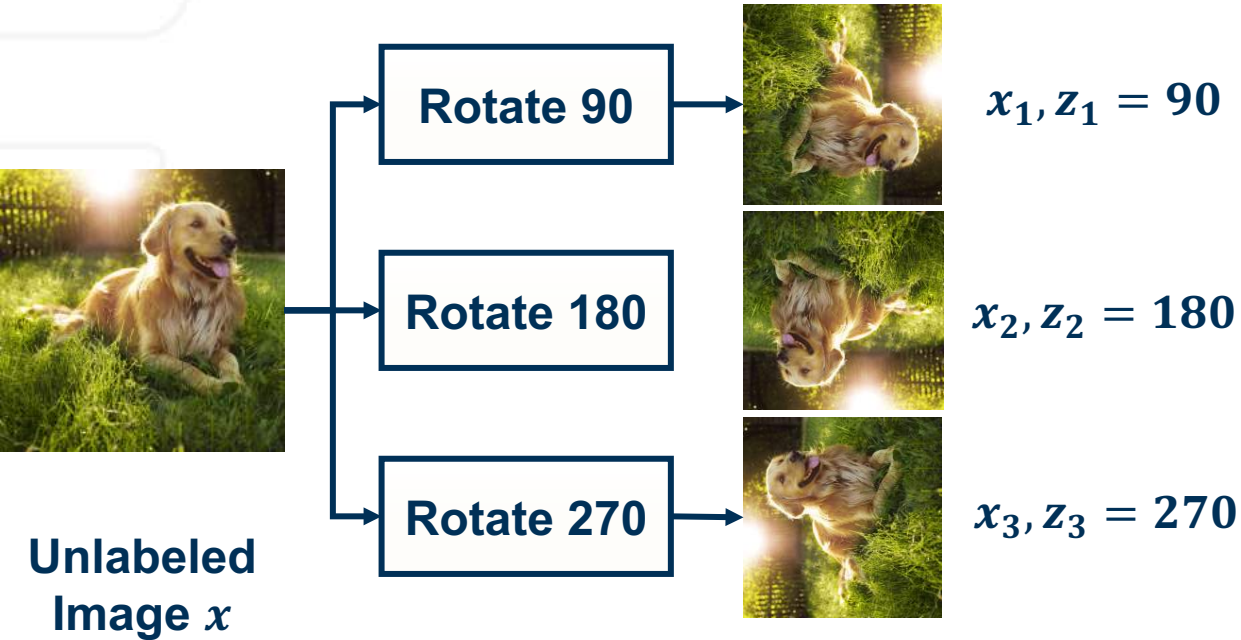
### Step 3: Attach linear layer and train to classify labels ( $y$ ) on labeled dataset



# Self-Supervision

## Motivation

### Step 1: Generate pseudo-labels via image rotations



Learning pre-text task will allow network to learn relevant features without needing explicit labels!

# Self-Supervision

## Types of Pre-text Tasks

**Differences in self-supervision are based on the type of pre-text task that is defined**

### Transformation Prediction

- Pre-text task performs some transformation on data and tasks model with trying to learn nature of transformation.

### Masked Prediction

- Pre-text task removes some part of the data and the model is tasked with trying to predict what was removed.

### Deep Clustering

- Identify clusters of features and iteratively assign pseudo-labels to train model.

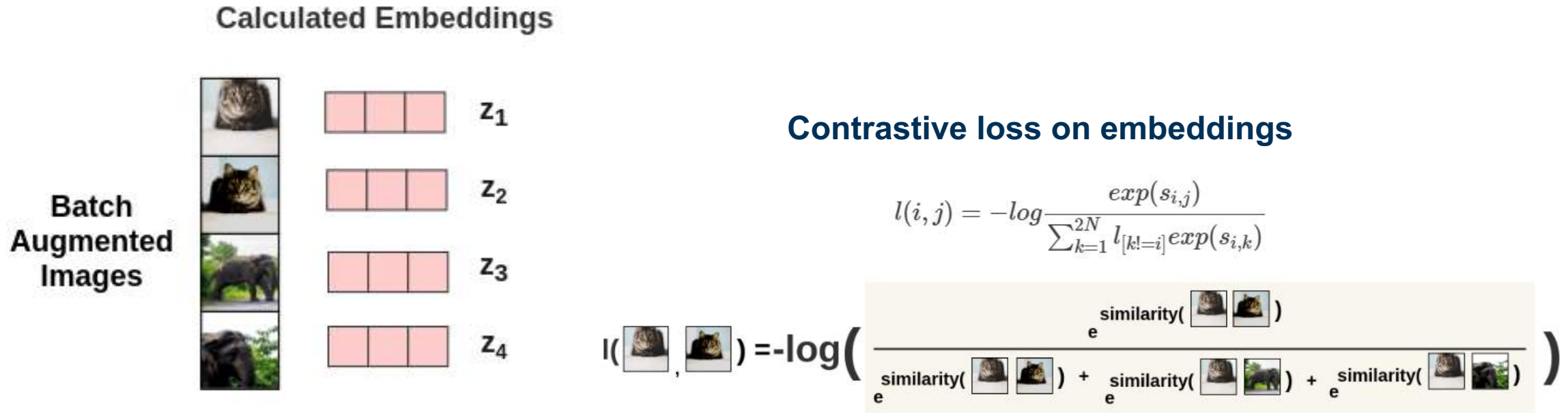
### Contrastive Learning

- Pre-text task identifies positive and negative pairs of data and the model is tasked with learning similarities to discriminate between positive and negatives.

# Contrastive Learning

## Sim-CLR Framework

The Pseudo-labels are used to create positive-negative pairs within each batch

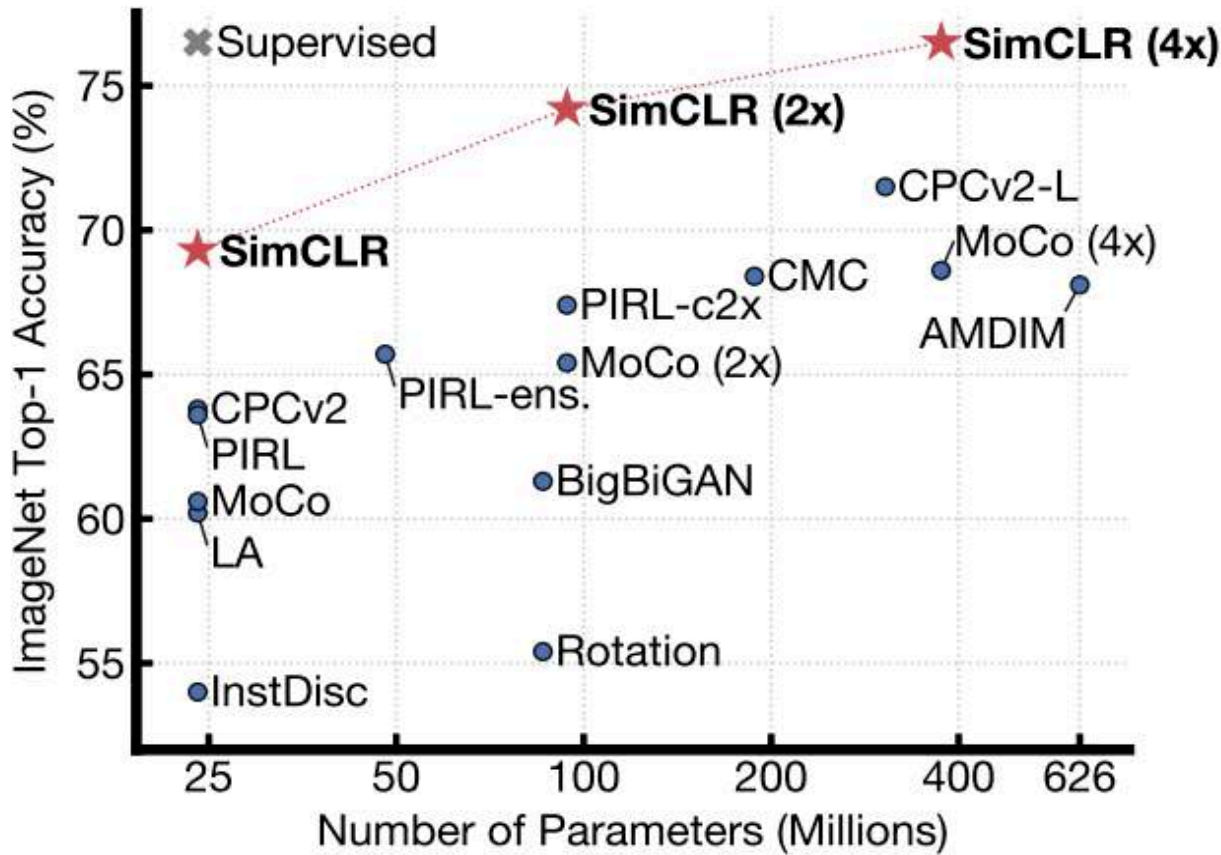


**Note: The positive pairs are only the augmentations and negative pairs are all other images in the batch**

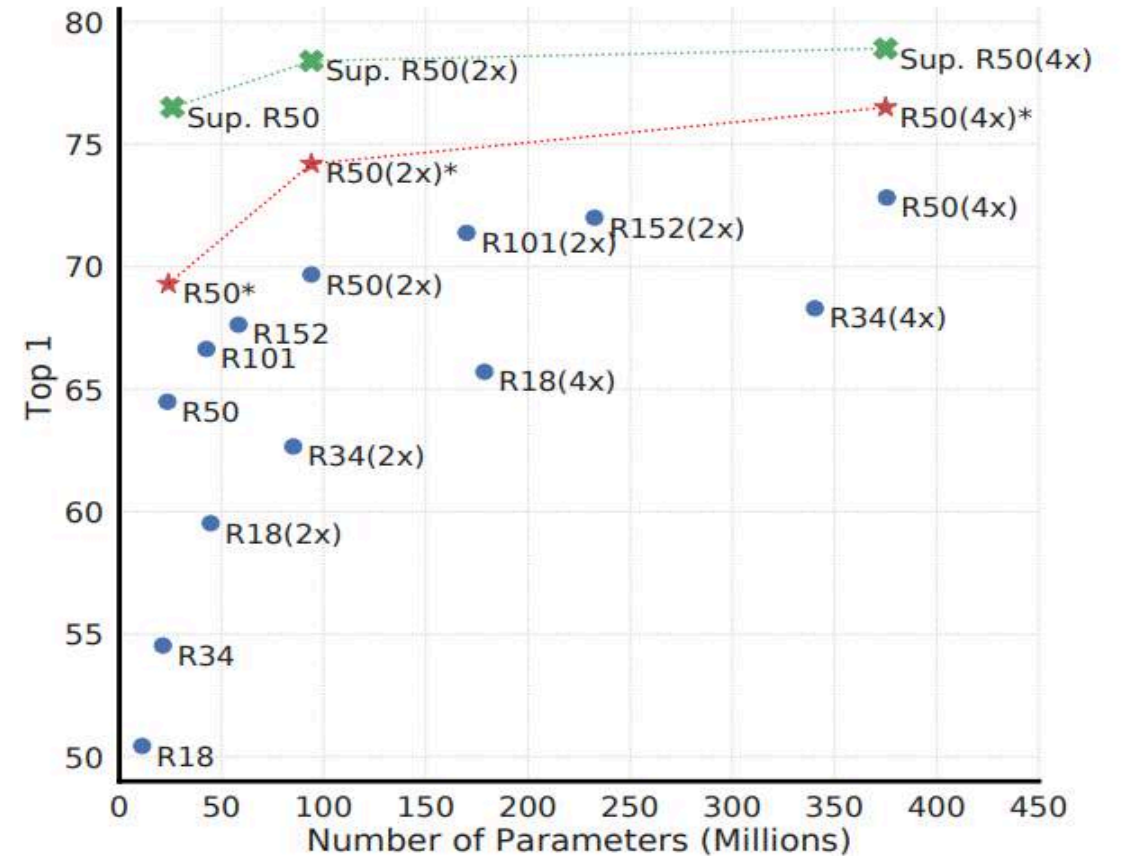
# Contrastive Learning

## Contrastive Learning vs Supervised Learning

### Performance vs Models



### Performance vs Parameters



# Contrastive Learning

## Contrastive Learning other than SIM-CLR

### What differentiates other Contrastive Learning methods from Sim-CLR?

Paper	Short description	Topics of contribution
Becker and Hinton [8]	Maximise MI between two views	Foundation
Bromley et al. [11]	Siamese network in metric learning setting	Foundation
Chopra, Hadsell, and LeCun [20]	Learn similarity metric with contrastive pair loss	Energy-based loss, Application
Hadsell, Chopra, and LeCun [39]	Learn invariant representation from pair loss	Energy-based loss, Application
Weinberger, Blitzer, and Saul [108]	Learn distance metric with triplet loss	Energy-based loss
Collobert and Weston [21]	Learn language model with triplet loss	Application
Chechik et al. [15]	Learn image retrieval model with triplet loss	Application
Noise Contrastive Estimation [38]	Introduce NCE, a general methods to learn unnormalised probabilistic model	Probabilistic loss
Mnih and Teh [71]	Learn language model with NCE-based loss	Application
Mikolov et al. [68]	Learn word embedding with Negative Sampling (NEG), a modified version of NCE	Probabilistic loss, Application
Wang et al. [105]	Learn fine-grained image similarity using deep network and triplet loss	Application
Wang and Gupta [107]	Use video's sequential coherence to learn unsupervised video representation	Similarity, Application
Lifted-structure loss [75]	Extend triplet loss to multiple positive and negative pairs per query	Energy-based loss
N-pair loss [92]	Proposed non-parametric classification loss with multiple negative pairs per query	Probabilistic loss
Wu et al. [109]	Focus on the quality of negative samples through a distance-weighted margin loss	Similarity, Energy-based loss
Hermans, Beyer, and Leibe [45]	State the important of mining hard samples in triplet loss	Similarity
Wu et al. [110]	Self-supervised representation with instance discrimination Memory bank to holds keys for next epoch	Application Encoder
CPC [77]	Mutual Information with the contrastive loss Define similarity with past-future context-instance relationship	Mutual Information loss Similarity
DIM [46]	Evaluate multiple mutual information bound for the contrastive loss Global-local context-instance relationship	Mutual Information Loss Similarity
MoCo [43]	Use momentum encoder to store features to memory queue	Encoder
SimCLR [16]	Simplify and demonstrate large empirical improvement in instance discrimination task Focus on the use of separate heads	Application Transform heads
BYOL [34]	Learning similarity without negative samples	Loss

The way that similar pairs (positives) and dissimilar pairs (negatives) are generated.

IEEE Open Journal of  
Signal Processing

## Exploiting the Distortion-Semantic Interaction in Fisheye Data



Kiran Kokilepersaud,  
PhD Student



Mohit Prabhushankar, PhD  
Postdoc



Ghassan AlRegib, PhD  
Professor



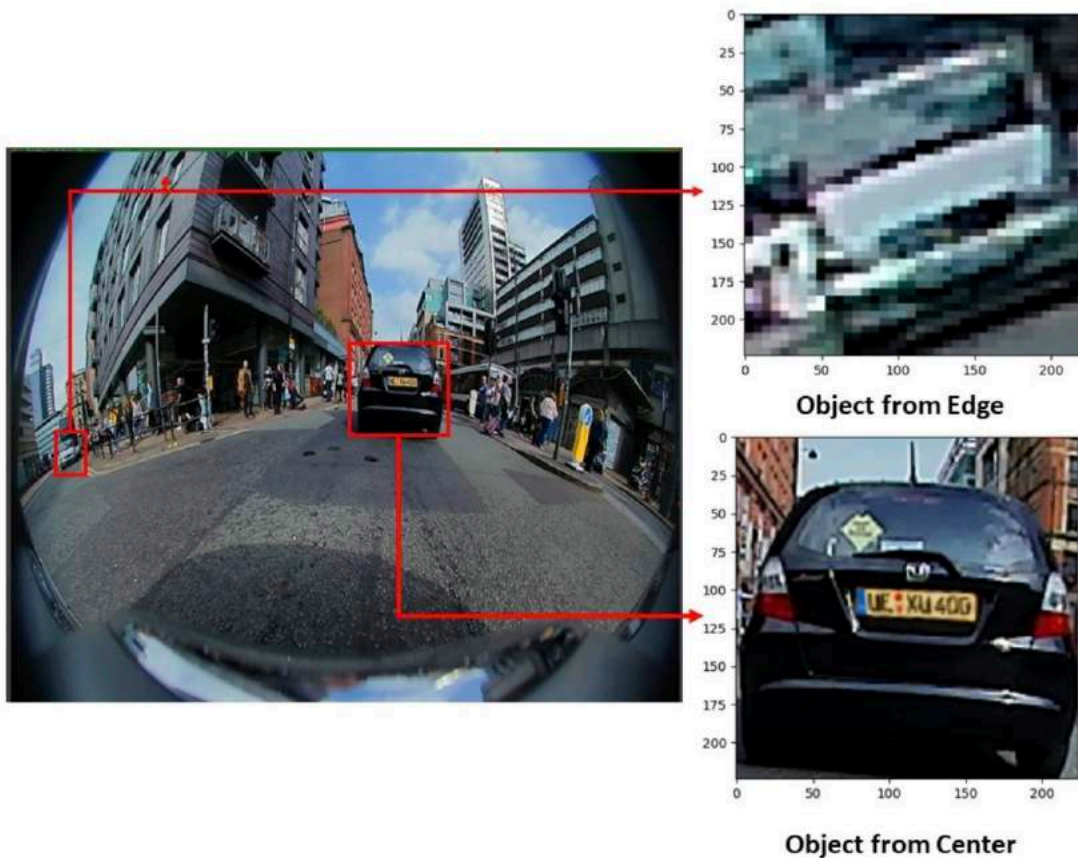
# Contrastive Learning for Fisheye Images

## Positive-negative pairs in Fisheye Images



Exploiting the Distortion-Semantic Interaction in Fisheye Data

**Intuition: Regions within a fisheye image are their own class. Hence, any object within them are positives**



Intuition for Loss 1:

All objects from the edge (be it a car, bike, pedestrian) are positives and objects from the centre (be it a car, bike, pedestrian) are negatives

Intuition for Loss 1:

All objects from labeled car (be it in the center or the edge) are positives and all other objects (be it in the center or the edge) are negatives



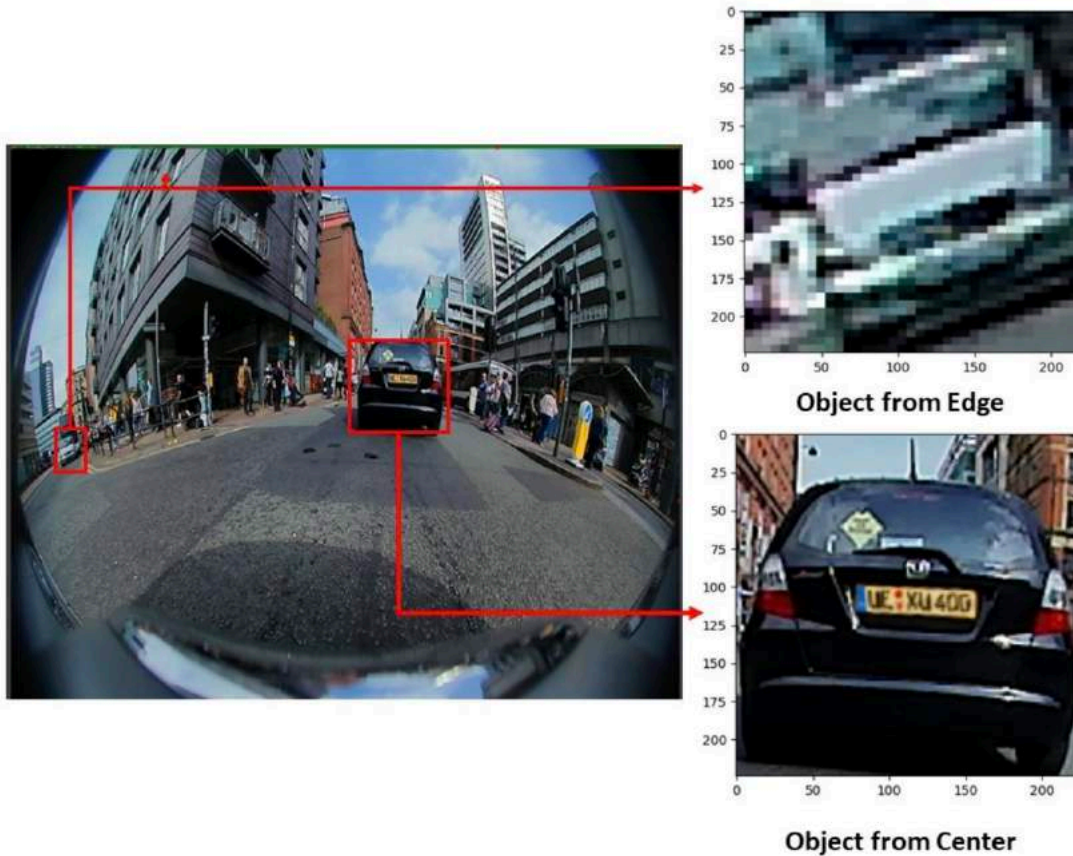
# Contrastive Learning for Fisheye Images

## Positive-negative pairs in Fisheye Images



Exploiting the Distortion-Semantic Interaction in Fisheye Data

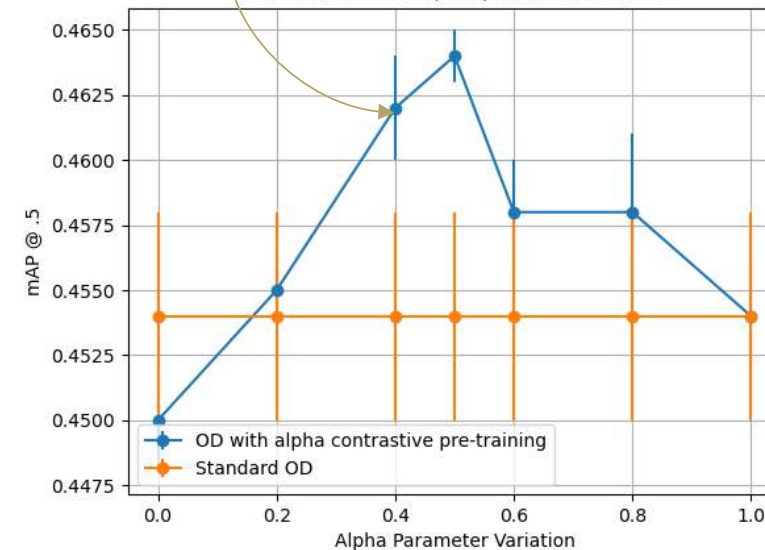
**Intuition: Regions within a fisheye image are their own class. Hence, any object within them are positives**



$$\alpha L_{class} + (1 - \alpha) L_{RegionClass}$$

$\alpha$  controls the level of unsupervised contrastive learning

Performance as alpha parameter varies



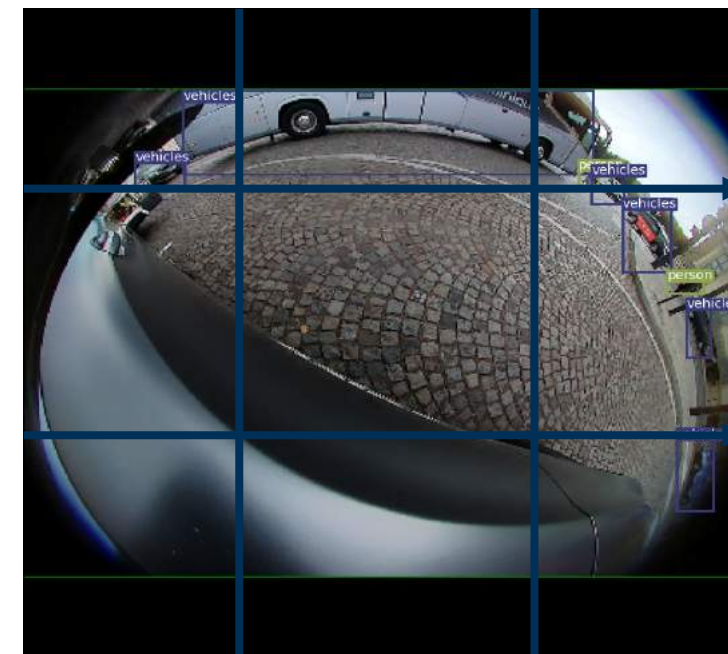
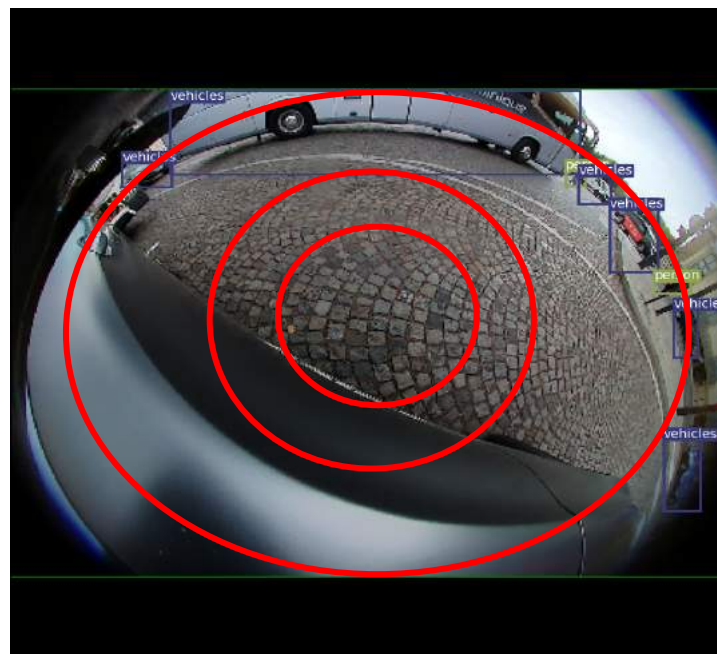
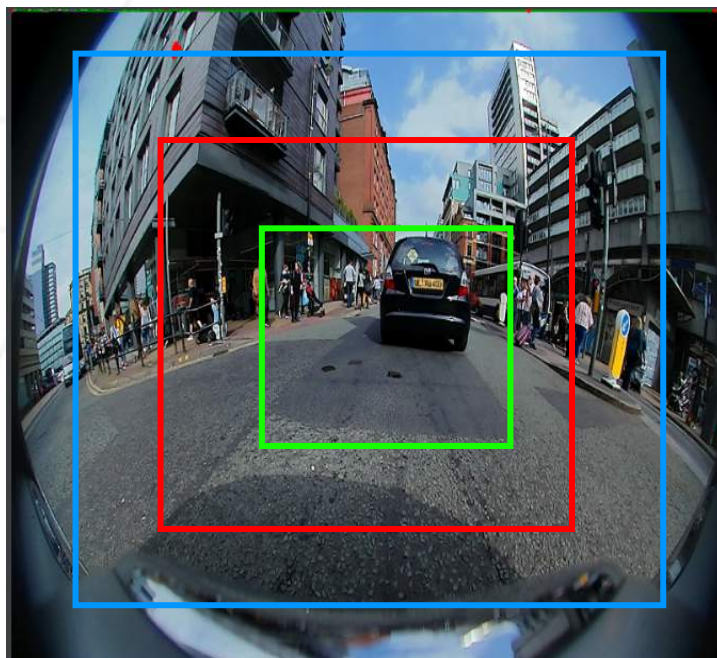
# Contrastive Learning for Fisheye Images

## Positive-negative pairs in Fisheye Images



Exploiting the Distortion-Semantic Interaction in Fisheye Data

Are there alternative ways of partitioning the regions?



Defining the positive-negative pairs is application dependent

# Objectives

## Takeaways from Part II

- Part I: Challenges in Perception and Autonomy
- **Part II: Deep Learning for Perception**
  - Transfer Learning and training at scale are essential for foundation model development
  - Self-supervised Learning provides a framework for large scale learning on unannotated data
- Part III: Existing Deep Learning solutions to Challenges in Perception
- Part IV: Remaining Challenges and Future Directions

# A Holistic View of Perception in Intel. Vehicles

## Part III: Deep Learning at Inference

# Objectives

## Objectives in Part III

- Challenging conditions at training
- Inference
  - Deficiencies at Inference
- Overcoming deficiencies at Inference
  - Anomaly Detection
  - Uncertainty
  - Explainability
- Case study 1: Robustness to challenging conditions
- Case study 2: Aberrant Object Detection

# Perception in AVs

## Technical Challenges

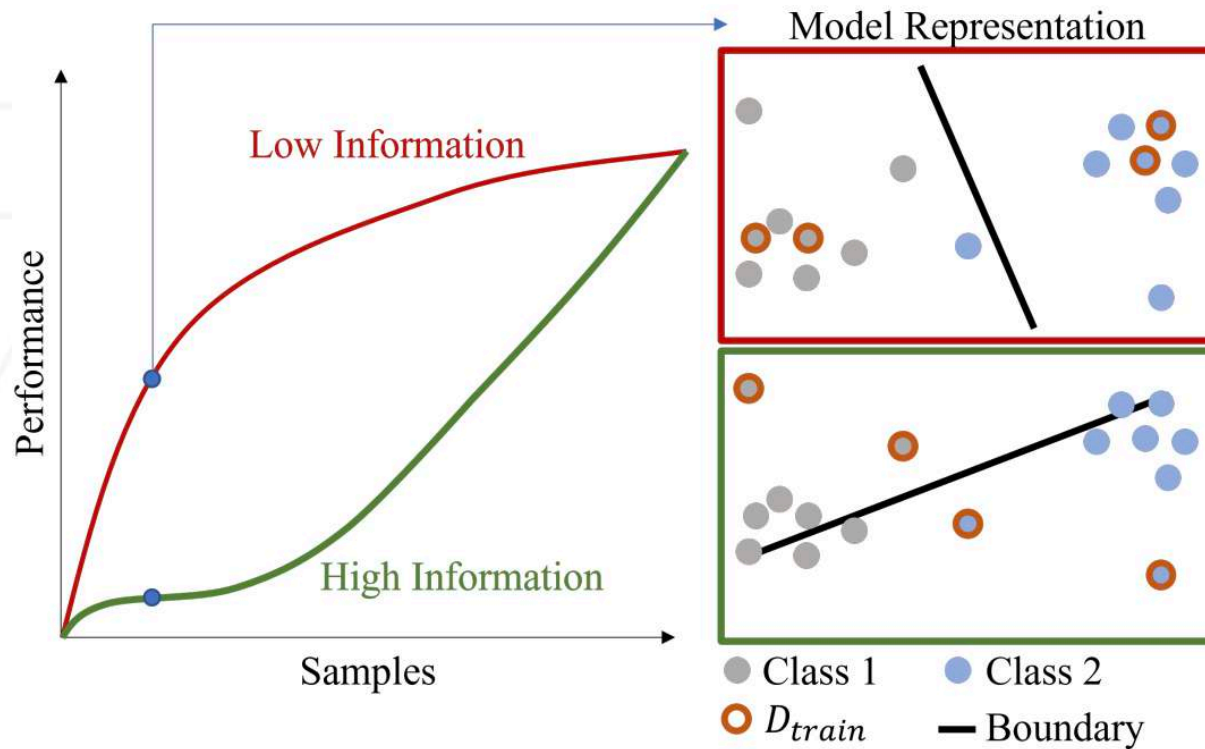
- Challenging weather
- Challenging sensing
- Challenging environments
- Context awareness
- Embedded perception
- V2X perception



# Challenging Conditions in Deep Learning

## Integrating Challenging Conditions in Training

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



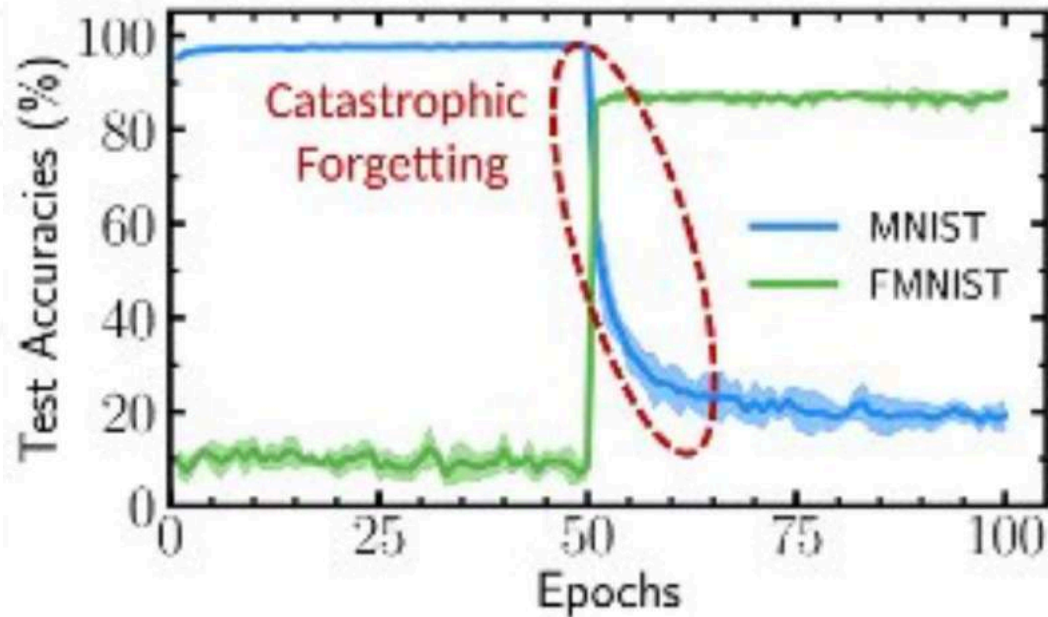
- The first instance of training must occur with less informative samples
- Less informative:
  - Highway scenarios
  - Parking
  - No accidents
  - No aberrant events

Novel samples = Most Informative

# Challenging Conditions in Deep Learning

## Integrating Challenging Conditions in Training

Subsequent training must not focus only on novel data



Catastrophic Forgetting

- 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

**Handle challenging conditions at Inference!**



# Inference

## What is Inference?

**Ability of a system to predict correctly on novel data**

**Novel data sources:**

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

**Model Train**



**At Deployment**



# Inference

## What is Inference?

**Ability of a system to predict correctly on novel data**

### Novel data sources

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



Trained Model

Cat

# Inference

## Deficiencies at Inference



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

*- Engineers, probably*

# Inference

## Overcoming Deficiencies at Inference

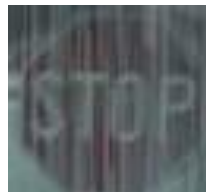
**What is required when networks are met with challenging data at inference?**

To overcome deficiencies, predictions from neural networks must be equipped with:

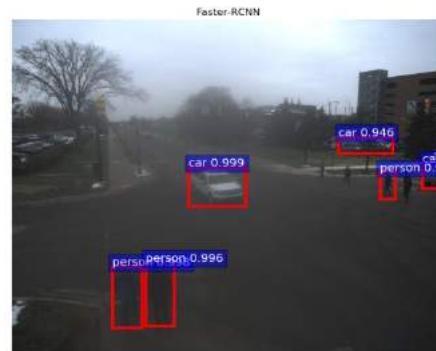
- Anomaly scores: How *close* to the training data is the novel data at inference?
- Uncertainty scores: How close to the *best* possible network is the trained network?
- Contextual Explainability: How *relevant* are the network explanations for its prediction?



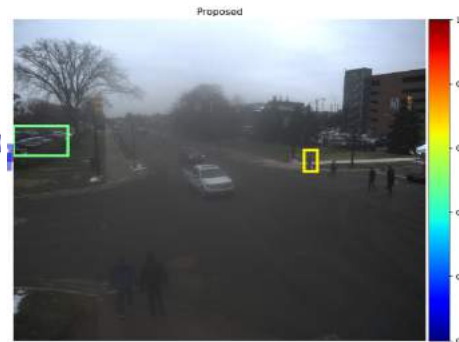
Training data



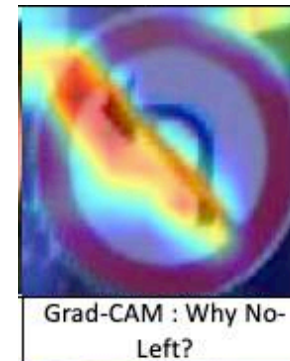
Anomalous data



Certain objects



Uncertain objects



Grad-CAM : Why No-Left?

'Why P'



Why No-Left, rather than Stop?

'Why P, rather than Q?'

# Inference

## Overcoming Deficiencies at Inference

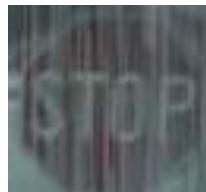
What is required when networks are met with challenging data at inference?

To overcome deficiencies, predictions from neural networks must be equipped with:

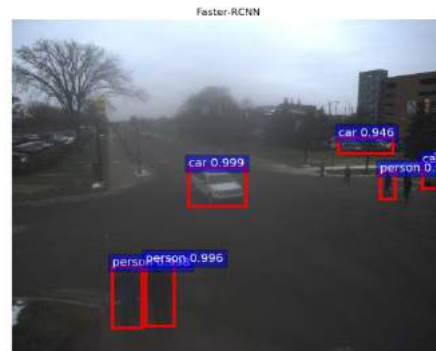
- **Anomaly scores:** How *close* to the training data is the novel data at inference?
- Uncertainty scores: How close to the *best* possible network is the trained network?
- Contextual Explainability: How *relevant* are the network explanations for its prediction?



Training data



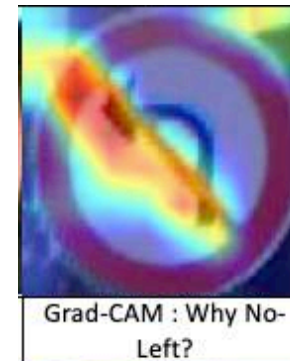
Anomalous data



Certain objects



Uncertain objects



'Why P'



'Why P, rather than Q?'



# Backpropagated Gradient Representations for Anomaly Detection



Gukyeong Kwon, PhD  
Amazon AWS



Mohit Prabhushankar, PhD  
Postdoc, Georgia Tech



Ghassan AlRegib, PhD  
Professor, Georgia Tech

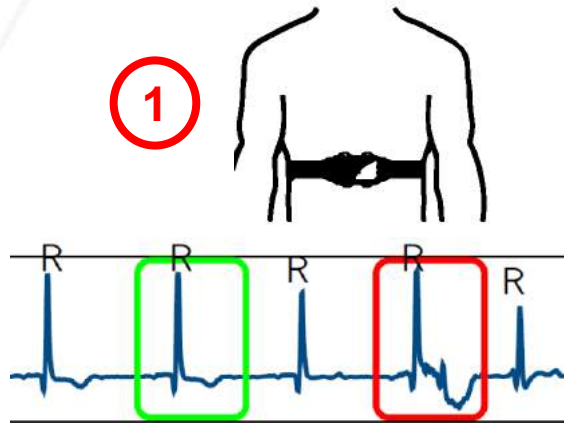


# Anomalies

## Finding Rare Events in Normal Patterns



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

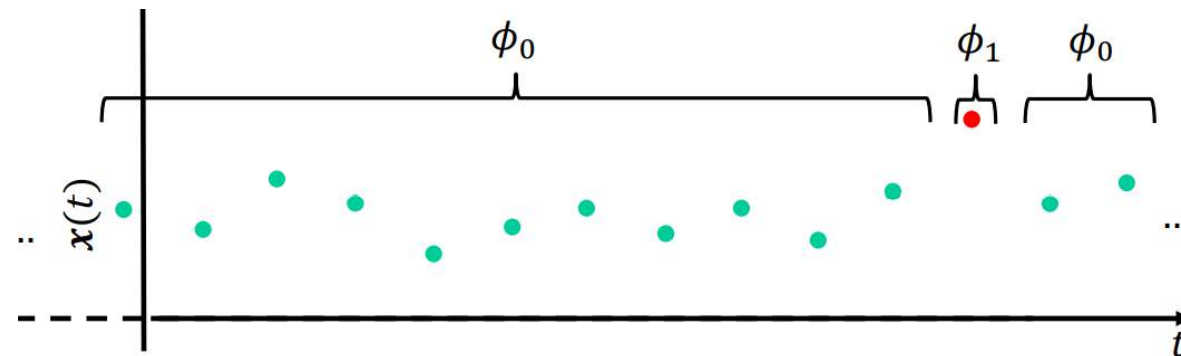


Statistical Definition:

- 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

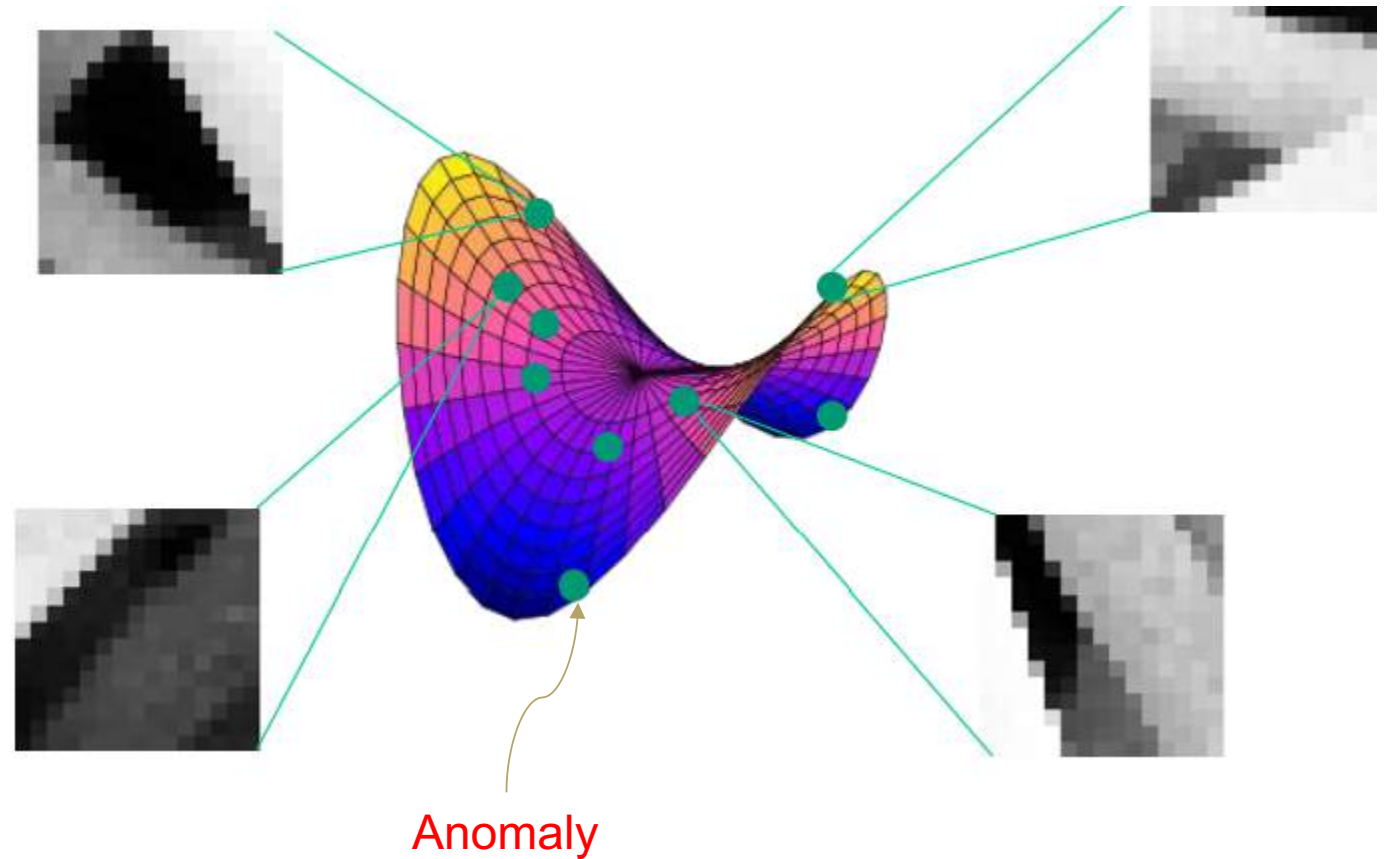
## Steps for Anomaly Detection



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





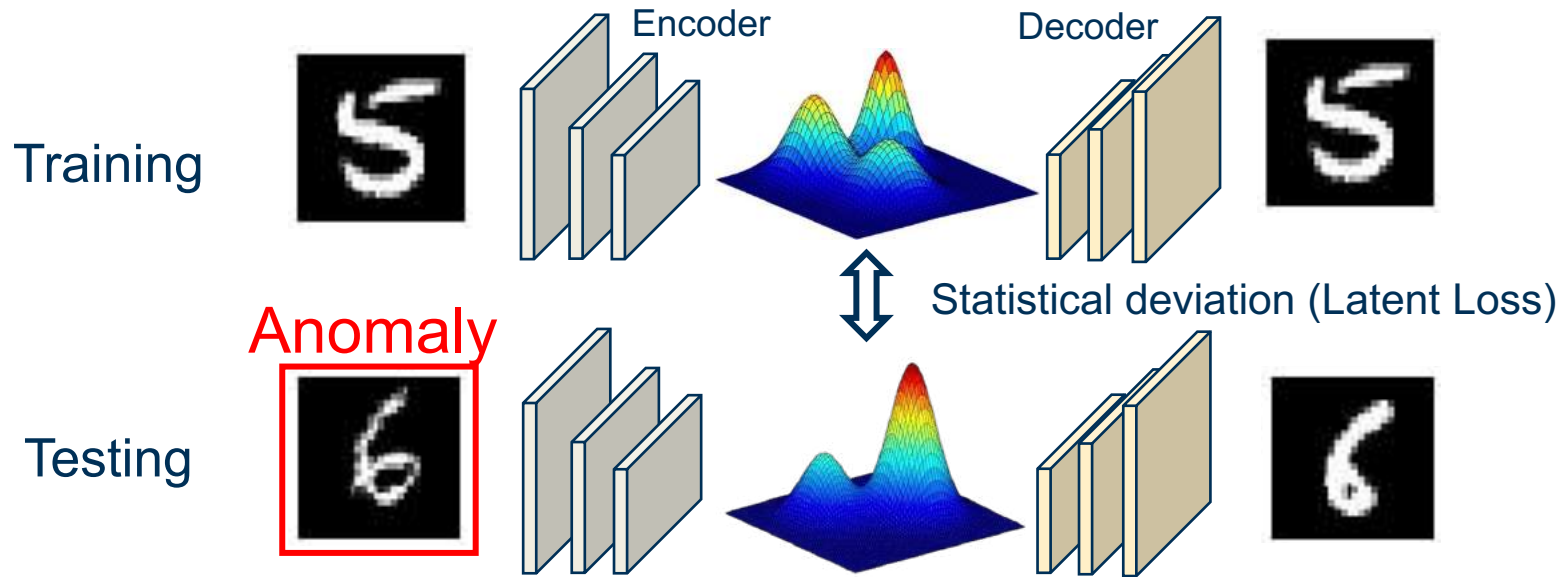
# Constraining Manifolds

## 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. Almohsen, 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.

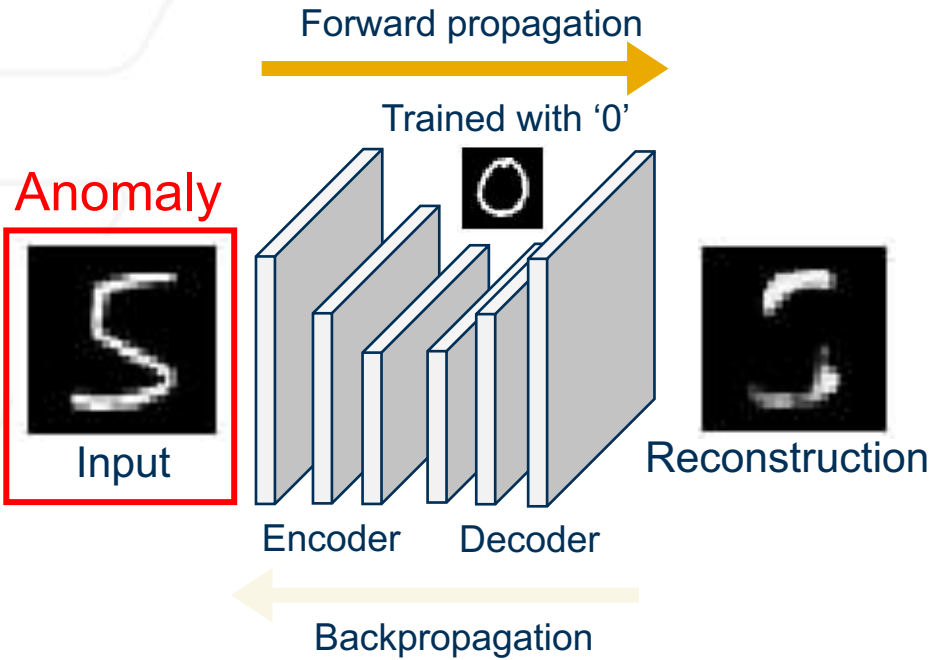
# Constraining Manifolds

## Gradient-based Constraints



Backpropagated Gradient Representations for Anomaly Detection

### Activation Constraints



Activation-based representation  
(Data perspective)

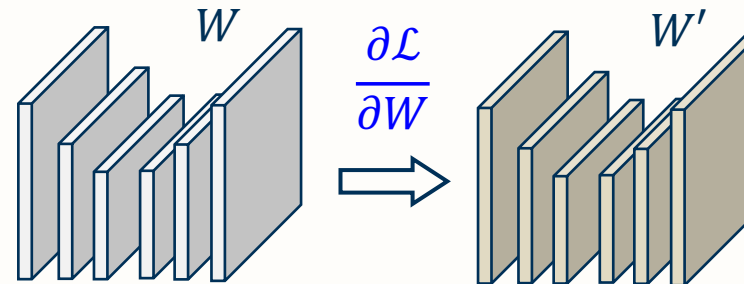
e.g. Reconstruction error ( $\mathcal{L}$ )



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

### Gradient Constraints

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



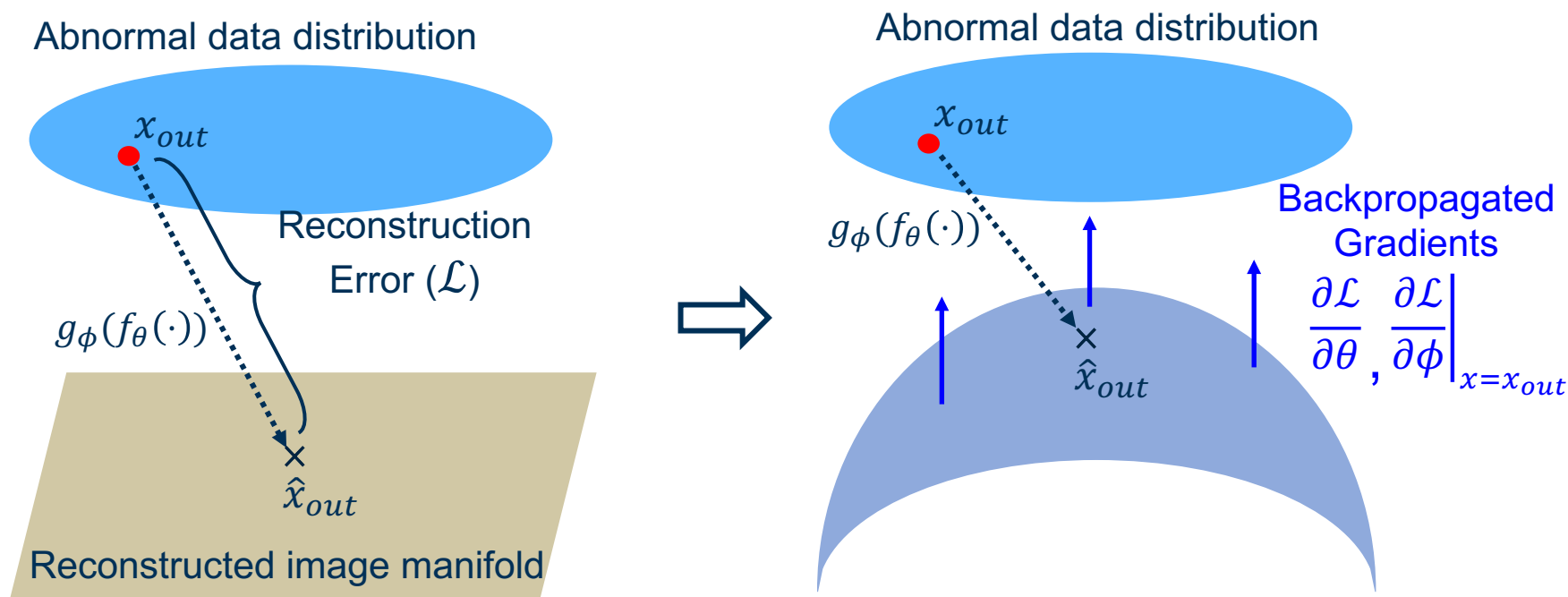
How much **model update** is required by the input?

# 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**



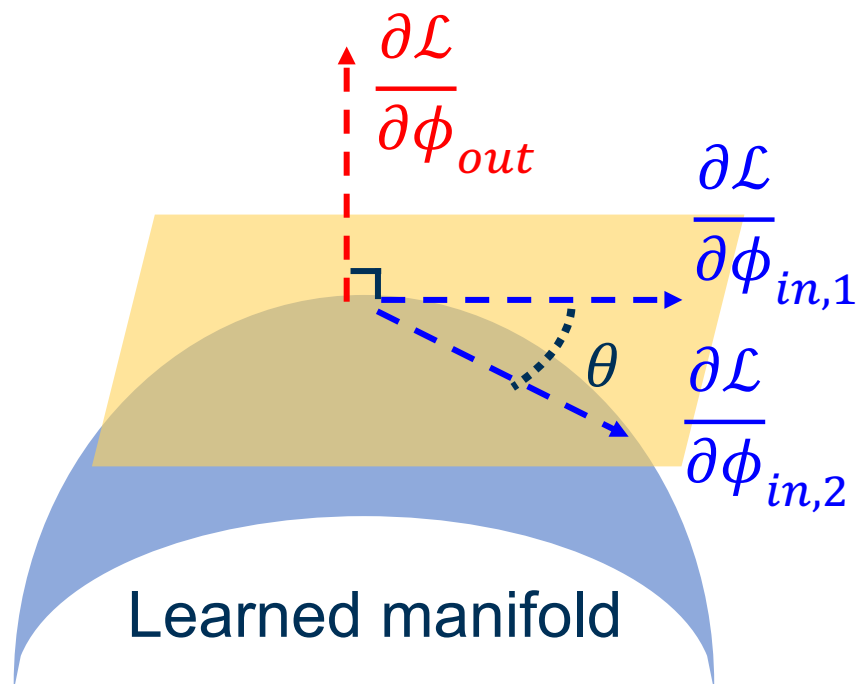
# GradCON: Gradient Constraint

## 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



$\phi$ : 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

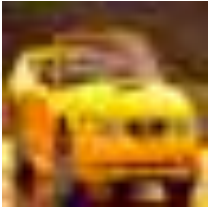


Backpropagated Gradient Representations for Anomaly Detection

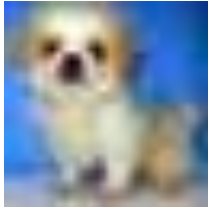
### AUROC Results

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
VAE	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
+ Grad	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)

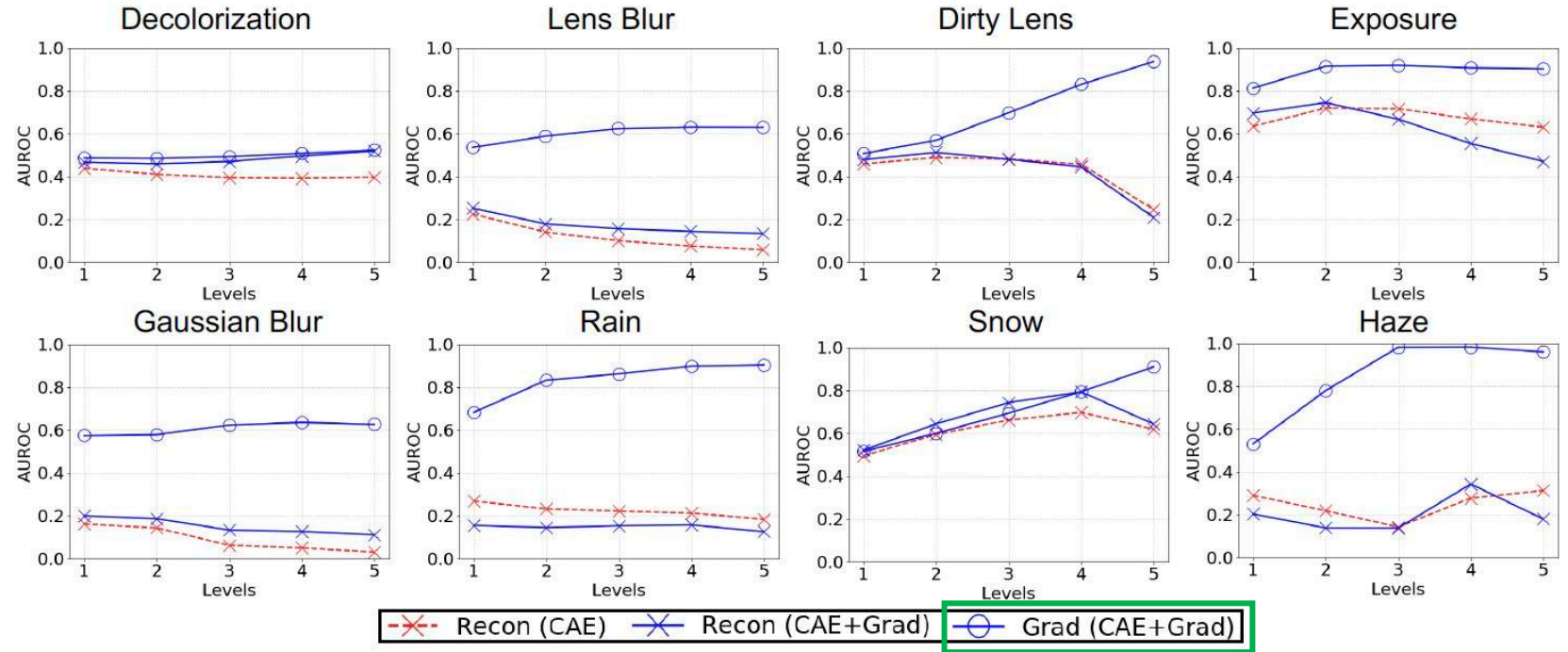


Normal



Abnormal

## AUROC Results



Recon: Reconstruction error, Grad: Gradient loss

# Inference

## Overcoming Deficiencies at Inference

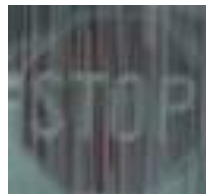
What is required when networks are met with challenging data at inference?

To overcome deficiencies, predictions from neural networks must be equipped with:

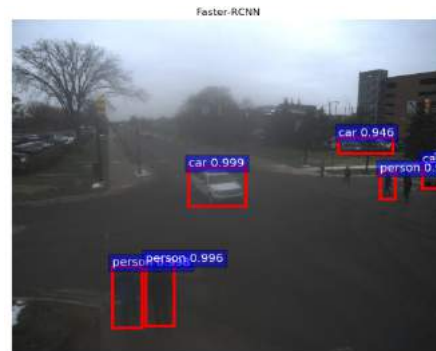
- Anomaly scores: How *close* to the training data is the novel data at inference?
- **Uncertainty scores**: How close to the *best* possible network is the trained network?
- Contextual Explainability: How *relevant* are the network explanations for its prediction?



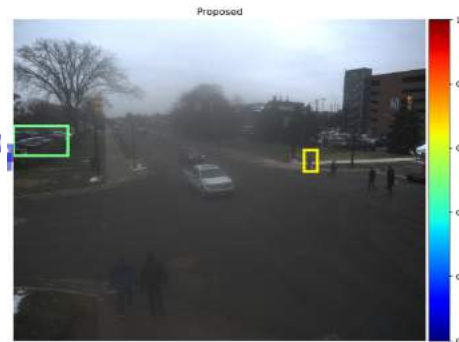
Training data



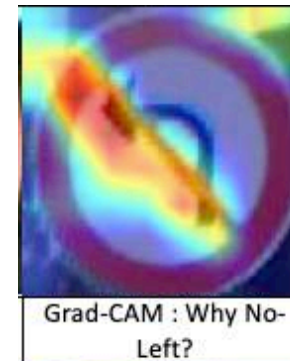
Anomalous data



Certain objects



Uncertain objects



'Why P'



'Why P, rather than Q?'



# Probing the Purview of Neural Networks via Gradient Analysis



Jinsol Lee,  
PhD Candidate



Mohit Prabhushankar, PhD  
Postdoc



Ghassan AlRegib, PhD  
Professor



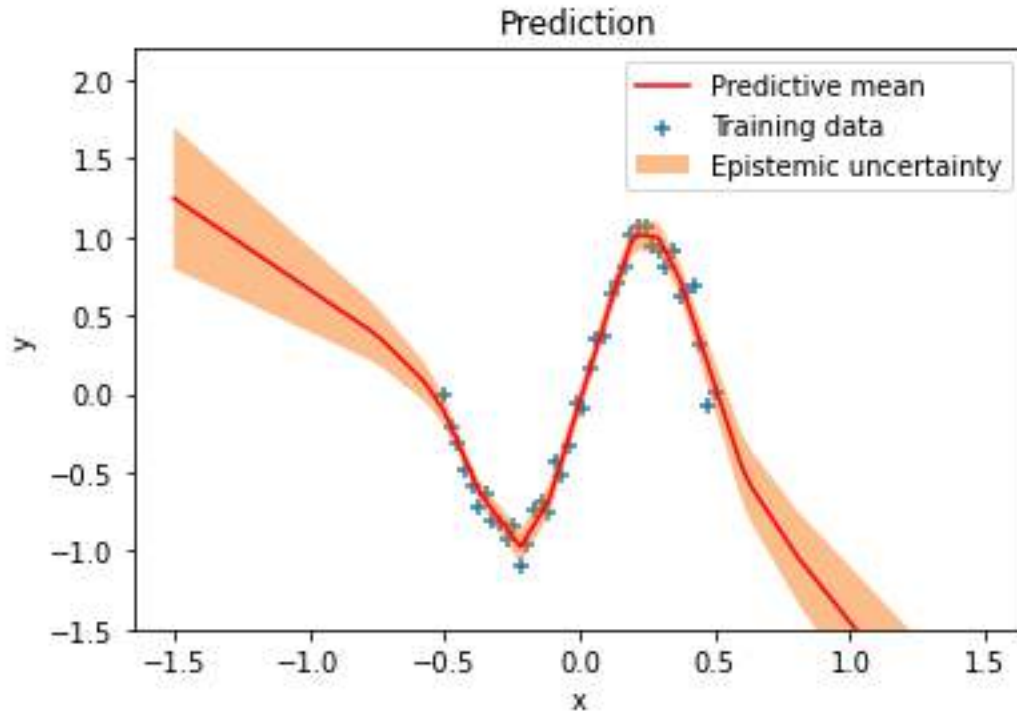


# Uncertainty

## What is Uncertainty?



**Uncertainty is a model knowing that it does not know**



A simple example: More the training data, lesser the uncertainty

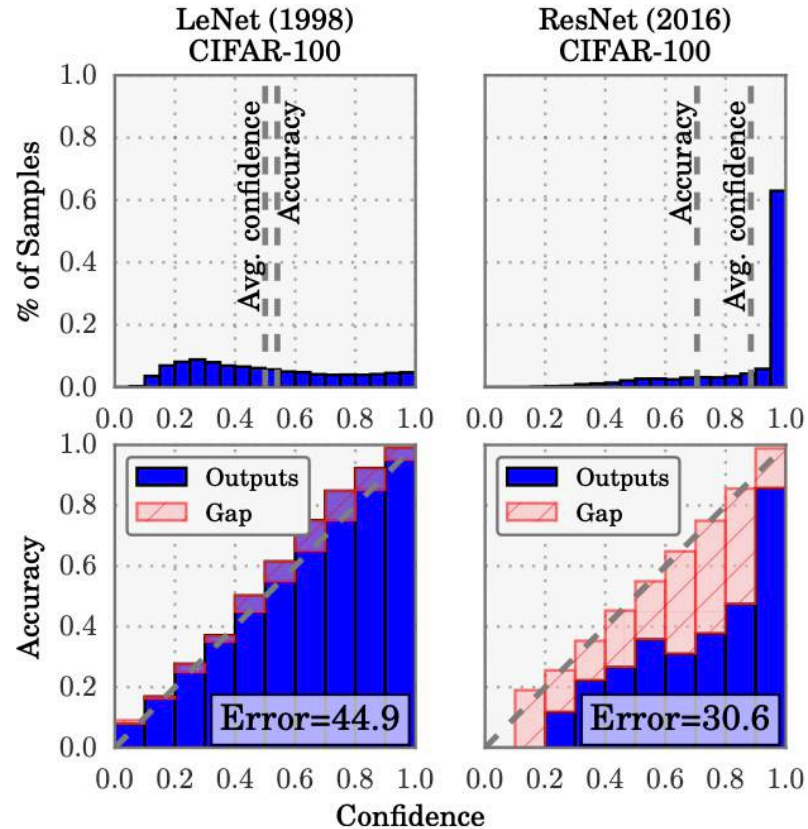
# Uncertainty

When is uncertainty an issue?



Probing the Purview of Neural Networks via Gradient Analysis

**Uncertainty is a model knowing that it does not know**



- Larger the model, more misplaced is a network's confidence
- On ResNet, the gap between prediction accuracy and its corresponding confidence is significantly high
- On OOD data, uncertainty is not easy to quantify

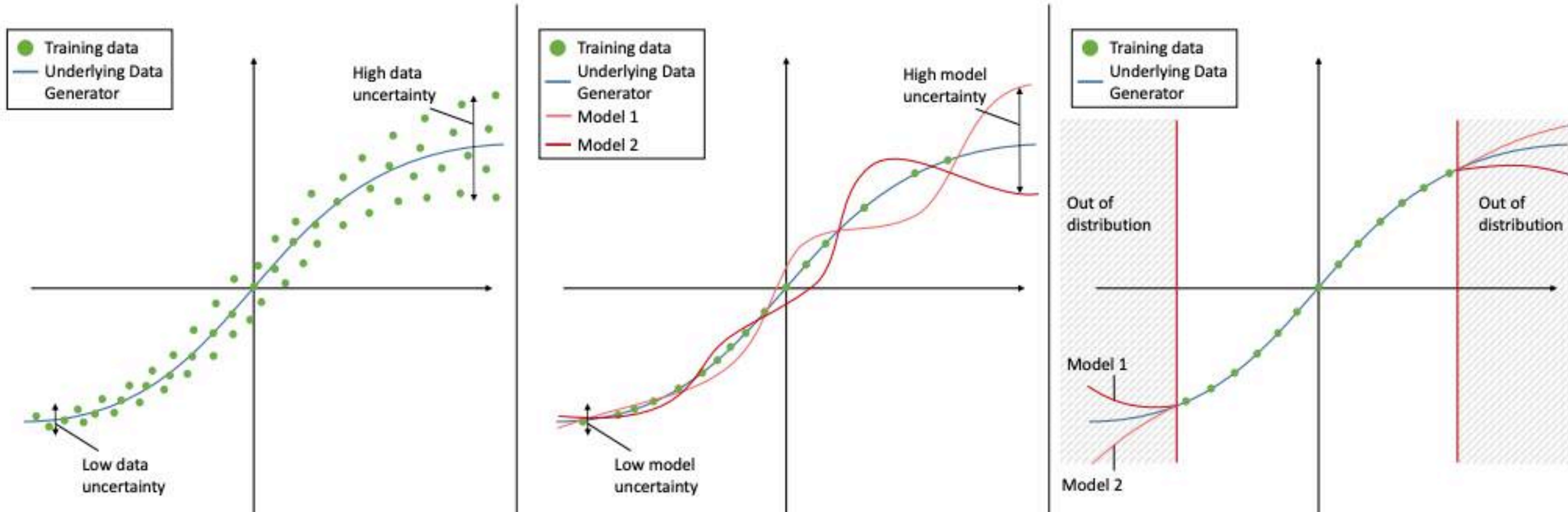
# Uncertainty

## Types of Uncertainty



Probing the Purview of Neural Networks via Gradient Analysis

Two major types of uncertainty: Uncertainty in data and uncertainty in model

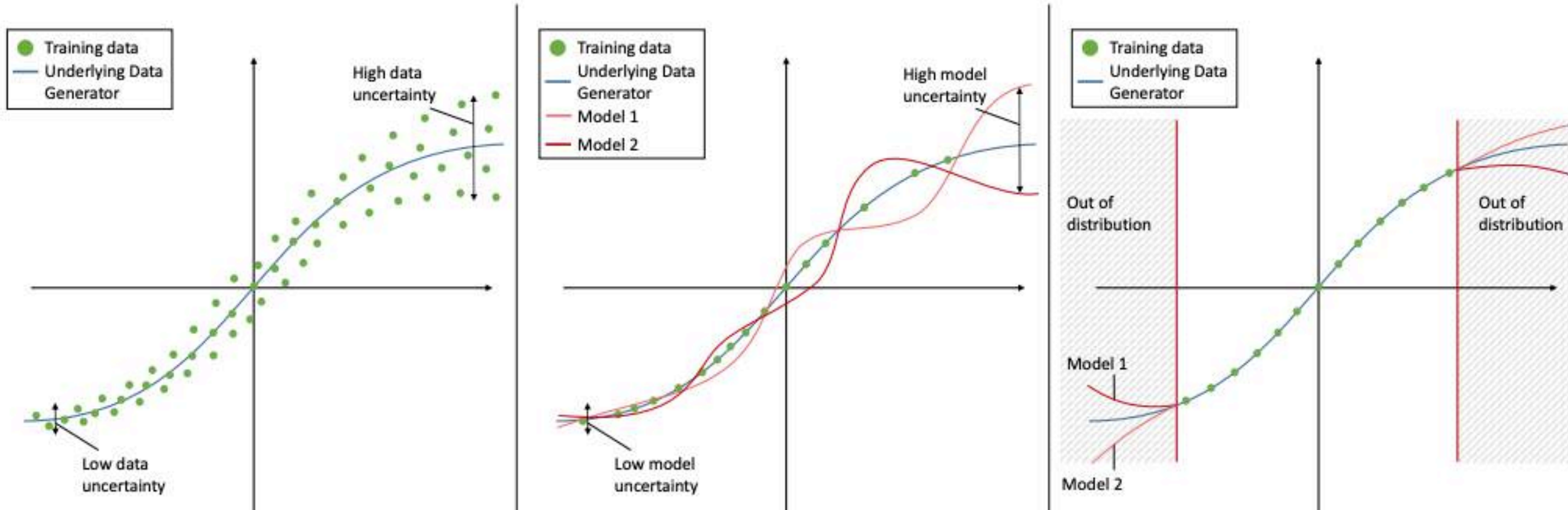


# Uncertainty

## Types of Uncertainty



For the purpose of predictions: Both uncertainties are combined as Predictive Uncertainty



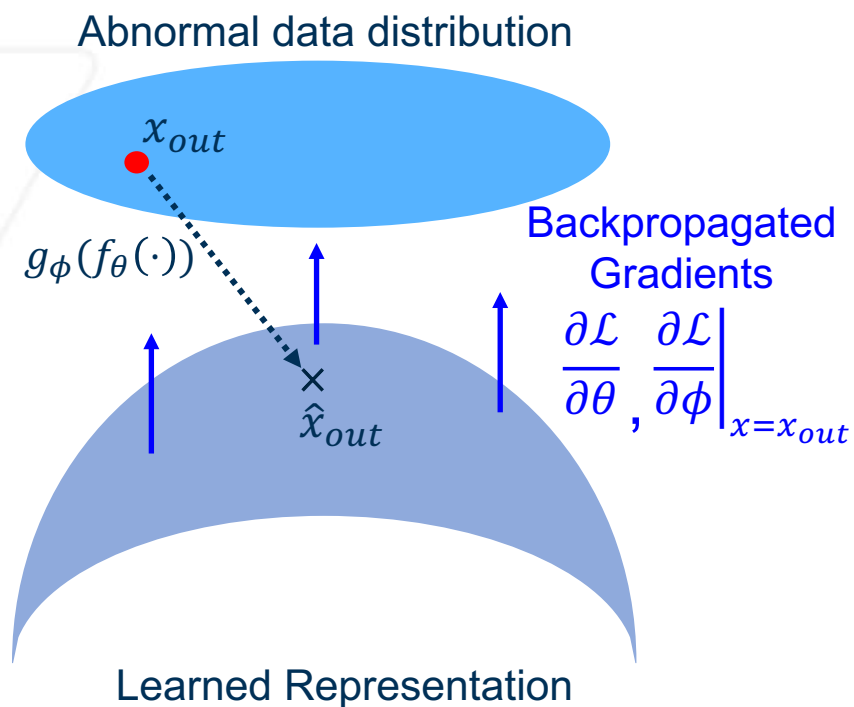
# Uncertainty in Neural Networks

## 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 or ground truth

# Uncertainty in Neural Networks

## Principle



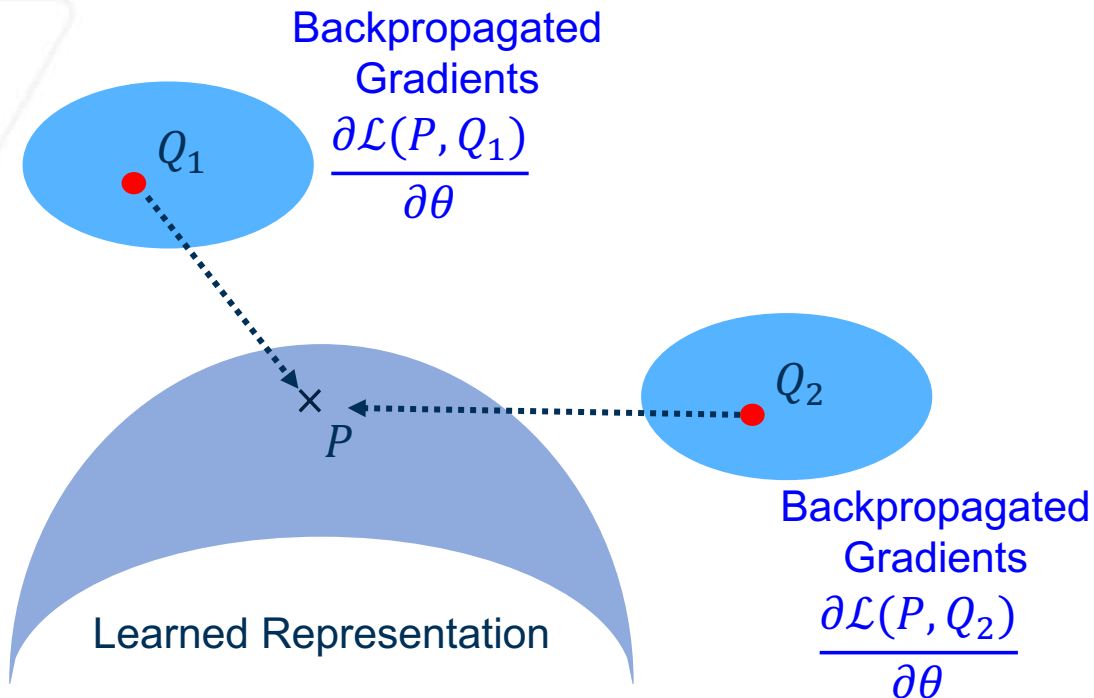
Probing the Purview of Neural Networks via Gradient Analysis

**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 or ground truth
- **We backpropagate all possible classes -  $Q_1, Q_2 \dots Q_N$  by backpropagating N one-hot vectors**
- Higher the distance to all classes, higher the uncertainty score

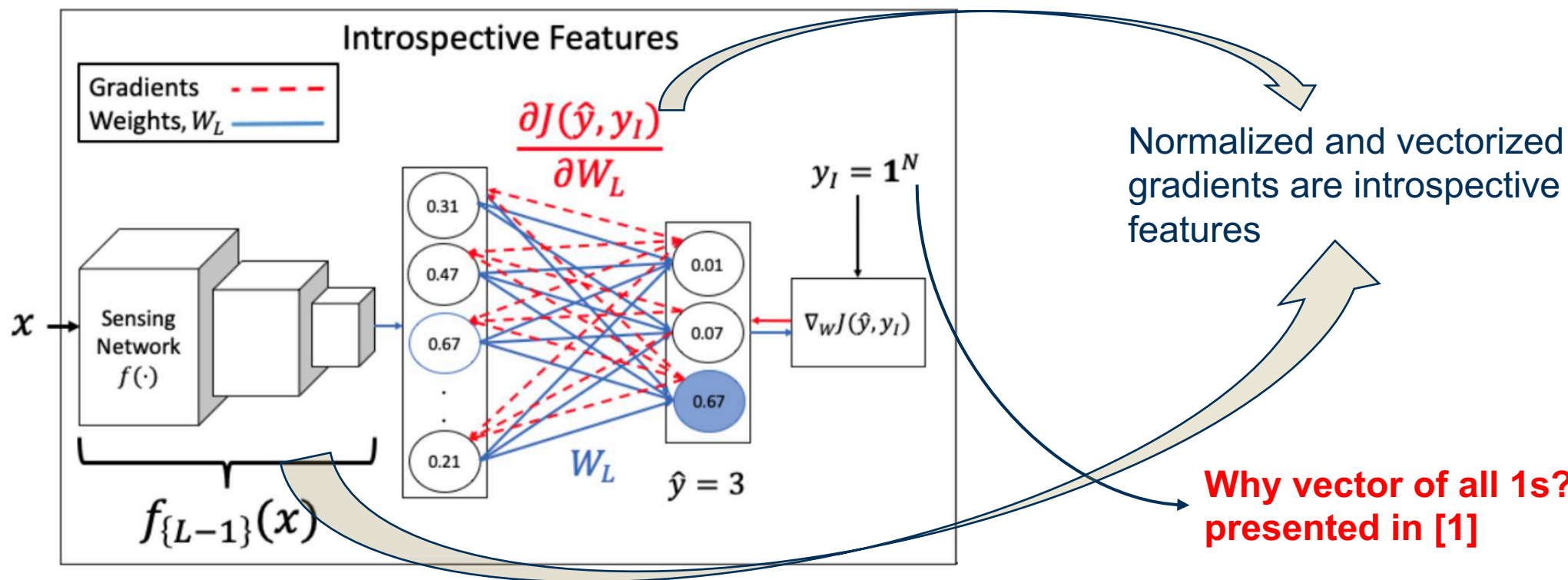
# Uncertainty in Neural Networks

## 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

## Deriving Gradient Features

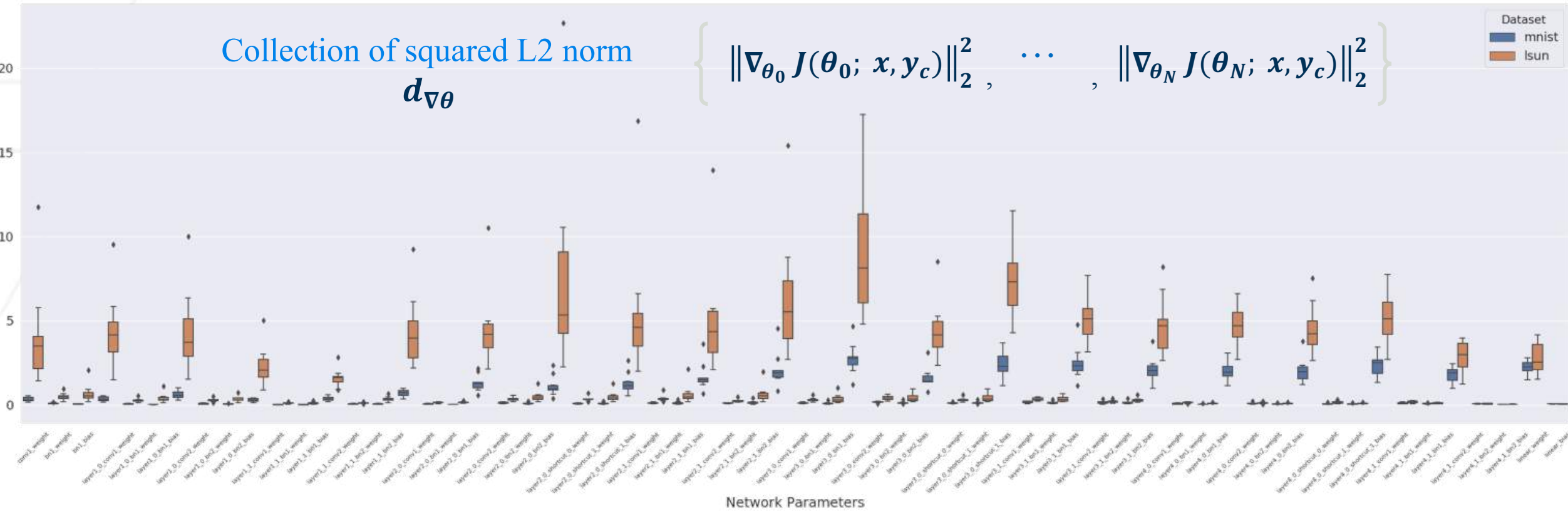


Probing the Purview of Neural Networks via Gradient Analysis

### Step 2: Take L2 norm of all generated gradients

Collection of squared L2 norm  $d_{\nabla\theta}$

$$\left\{ \|\nabla_{\theta_0} J(\theta_0; x, y_c)\|_2^2, \dots, \|\nabla_{\theta_N} J(\theta_N; x, y_c)\|_2^2 \right\}$$



### MNIST: In-distribution, SUN: Out-of-Distribution



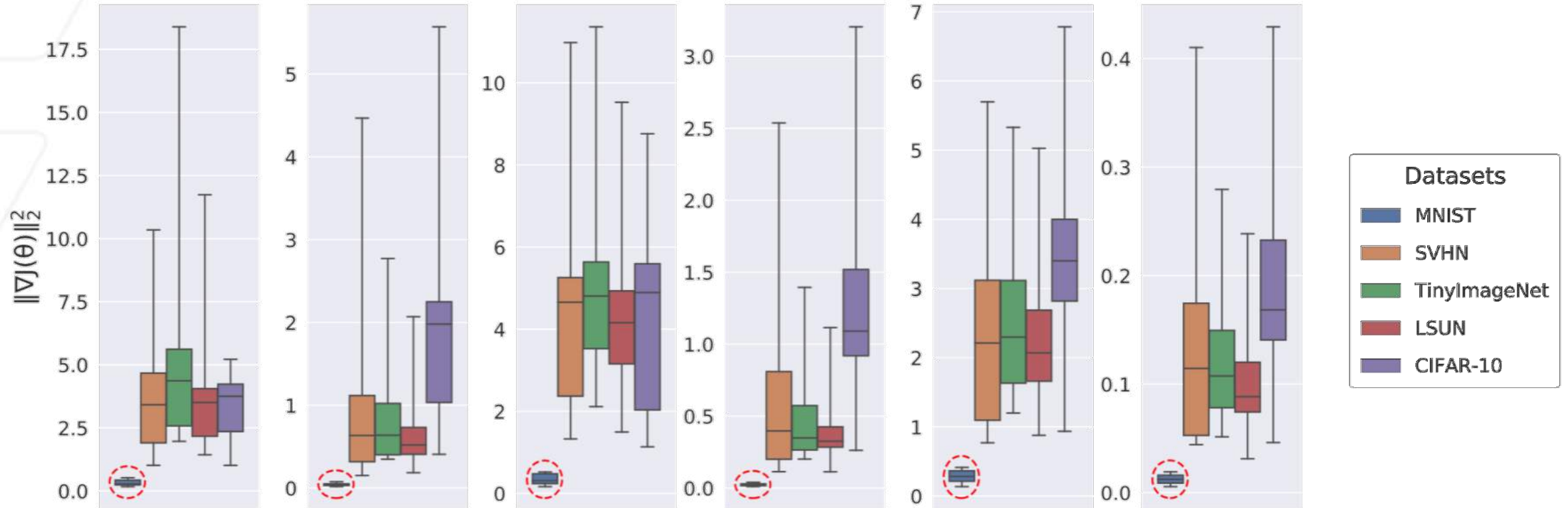
# Gradient-based Uncertainty

## Uncertainty Results 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**

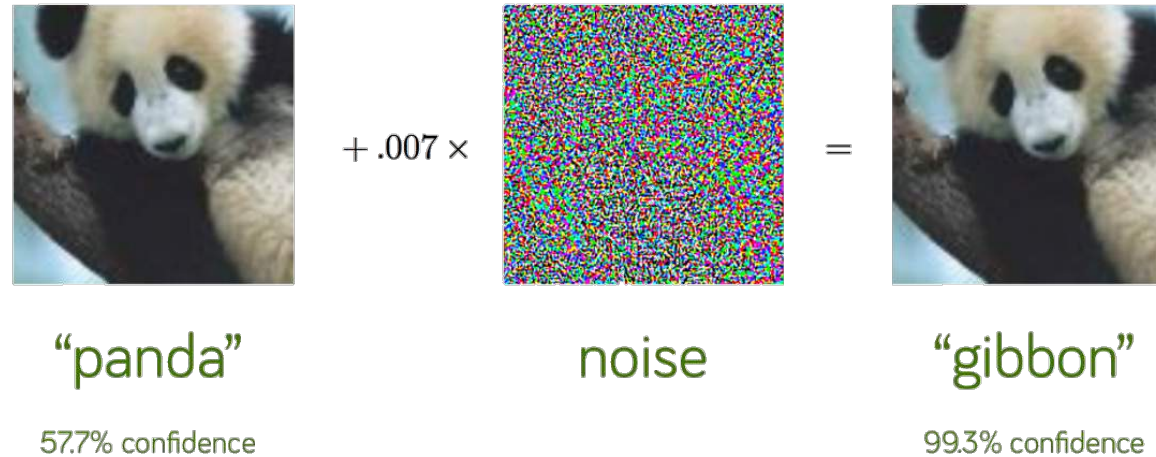
# Gradient-based Uncertainty

## Uncertainty Results 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 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	<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

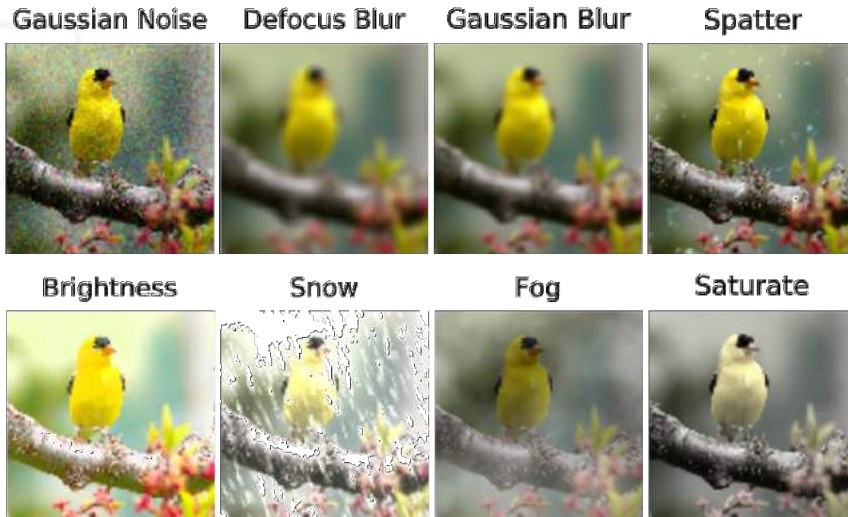
## Uncertainty Results to Detect Challenging Conditions



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



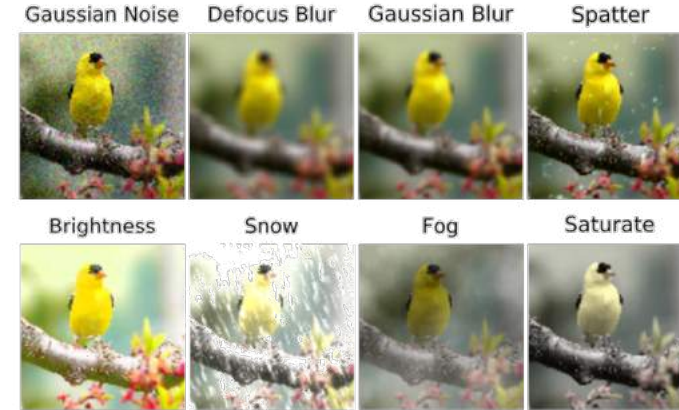
# Gradient-based Uncertainty

## Uncertainty Results to Detect 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>



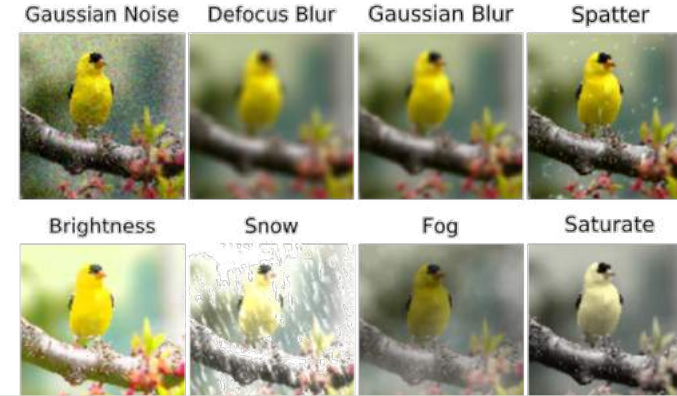
# Gradient-based Uncertainty

## Uncertainty Results to Detect Challenging Conditions



Probing the Purview of Neural Networks via Gradient Analysis

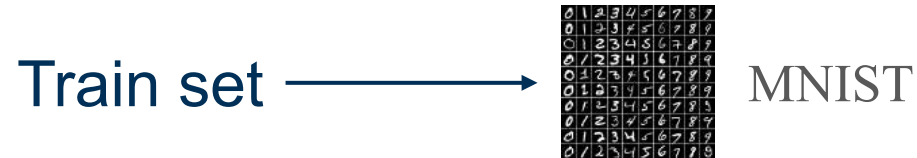
Dataset	Method	Mahalanobis [12] / Ours				
		Corruption	Level 1	Level 2	Level 3	Level 4
CIFAR-10-C	Noise	96.63 / <b>99.95</b>	98.73 / <b>99.97</b>	99.46 / <b>99.99</b>	99.62 / <b>99.97</b>	99.71 / <b>99.99</b>
	LensBlur	94.22 / <b>99.95</b>	97.51 / <b>99.99</b>	99.26 / <b>100.0</b>	99.78 / <b>100.0</b>	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



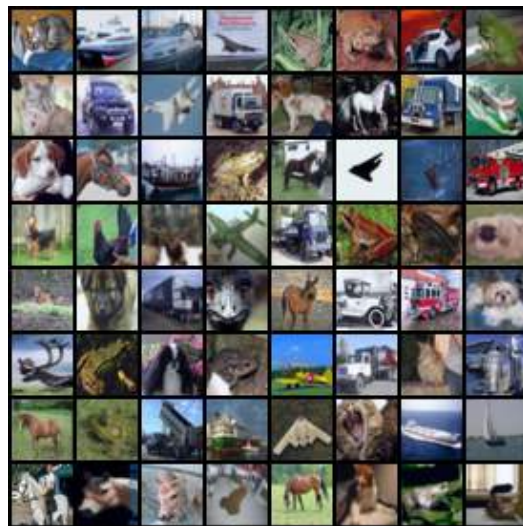
Probing the Purview of Neural Networks via Gradient Analysis



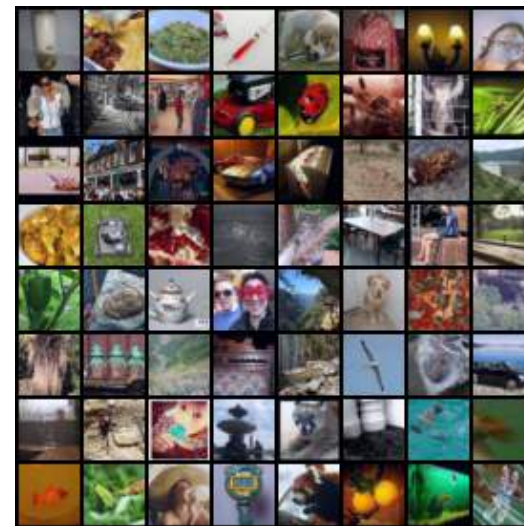
**Goal:** To detect that these datasets are not part of training



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>

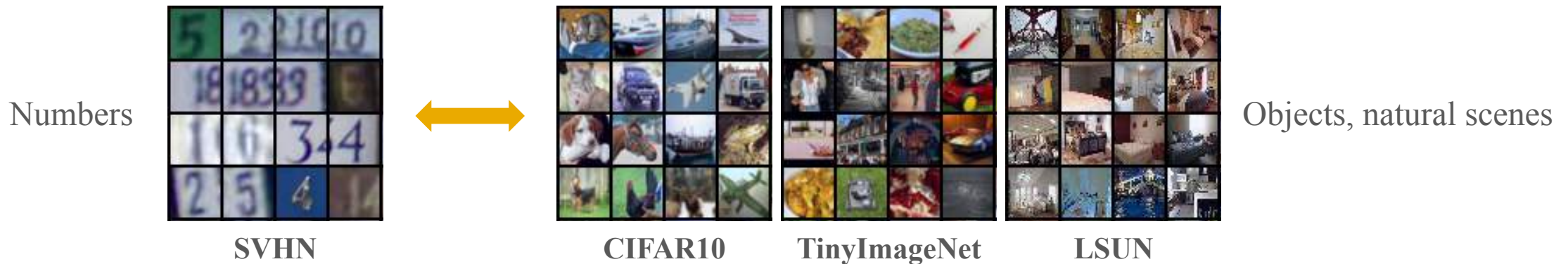


# Out-of-Distribution Detection



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>



# Inference

## Overcoming Deficiencies at Inference

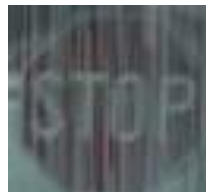
What is required when networks are met with challenging data at inference?

To overcome deficiencies, predictions from neural networks must be equipped with:

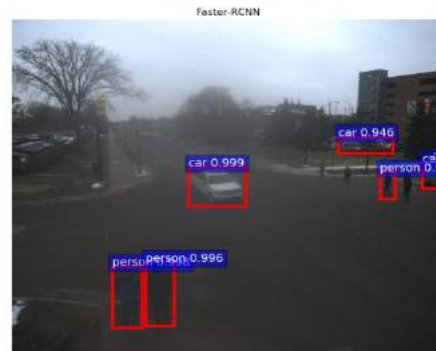
- Anomaly scores: How *close* to the training data is the novel data at inference?
- Uncertainty scores: How close to the *best* possible network is the trained network?
- **Contextual Explainability**: How *relevant* are the network explanations for its prediction?



Training data



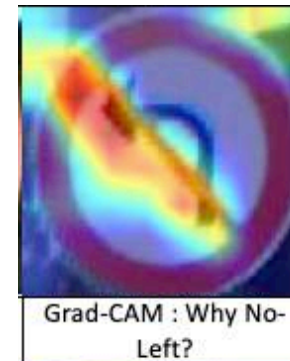
Anomalous data



Certain objects



Uncertain objects



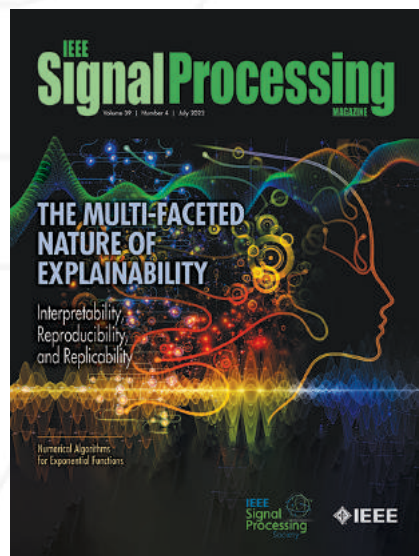
Grad-CAM : Why No-Left?

'Why P'



Why No-Left, rather than Stop?

'Why P, rather than Q?'



# Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations



Mohit Prabhushankar, PhD  
Postdoc



Ghassan AlRegib, PhD  
Professor



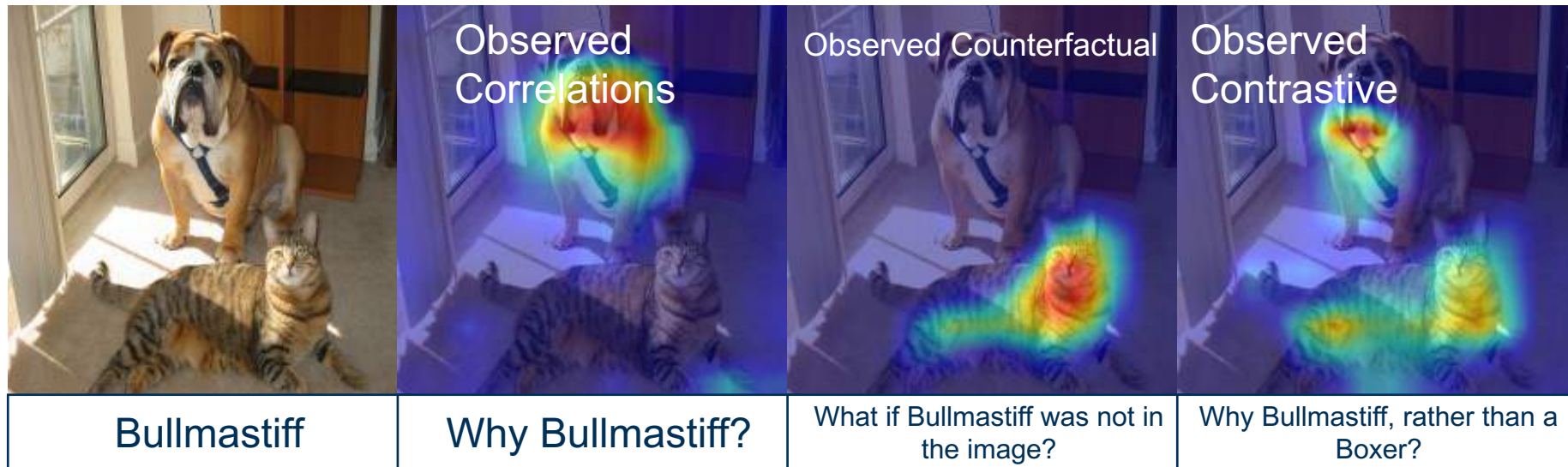
# Explanations

## What are Visual Explanations?



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

- Explanations are defined as a set of rationales used to understand the reasons behind a decision
- If the decision is based on visual characteristics within the data, the decision-making reasons are visual explanations



# Explanations

## Why Explainability?

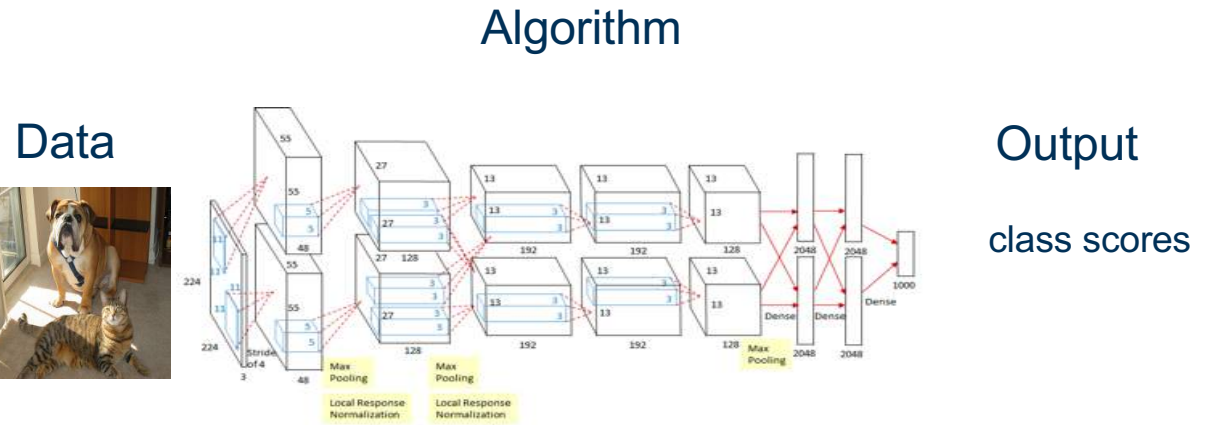


Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Explainability matters establishes trust in deep learning systems by developing *transparent* models that can explain *why they predict what they predict* to humans

Explainability is useful in:

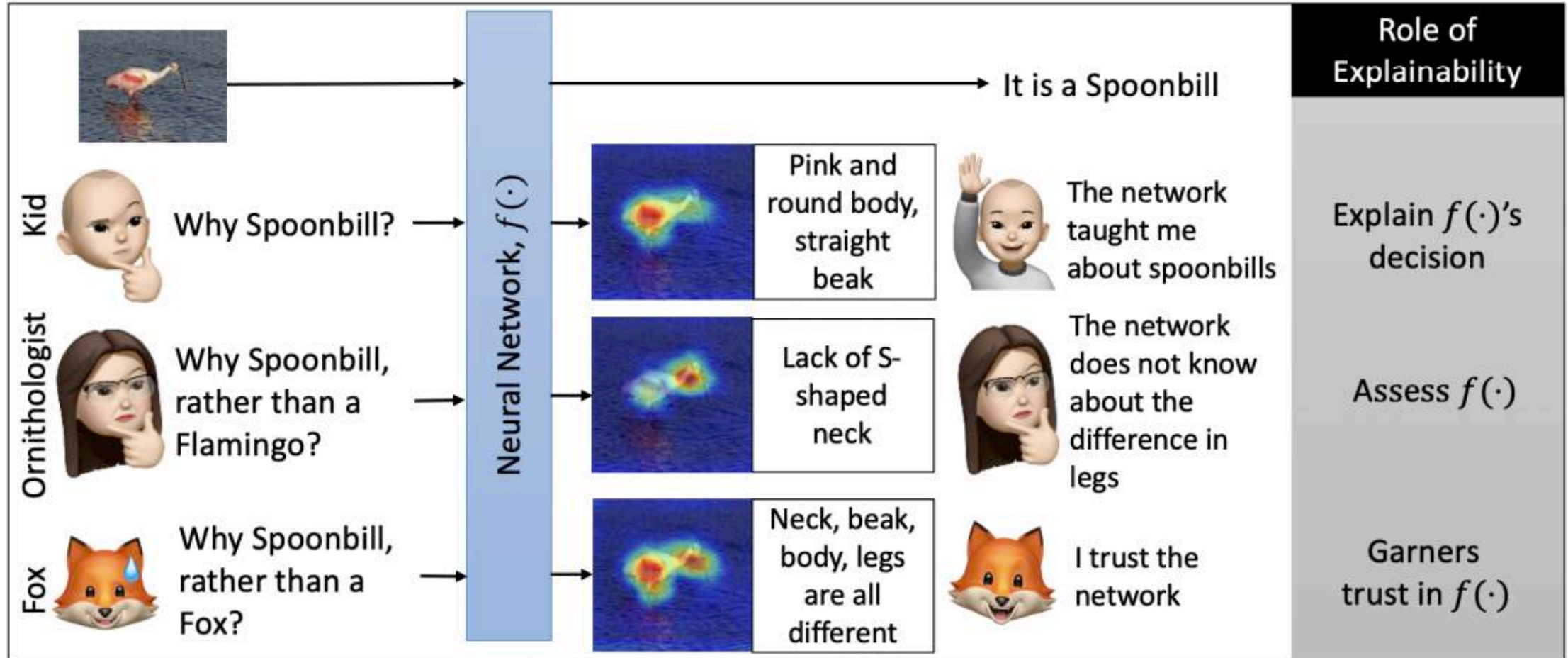
- Medical: help doctors diagnose
- Seismic: help interpreters label seismic data
- Autonomous Systems: build appropriate trust and confidence



Deep models act as algorithms that take data and output something **without** being able to **explain** their methodology

# Explanations

## Role of Visual Explanations



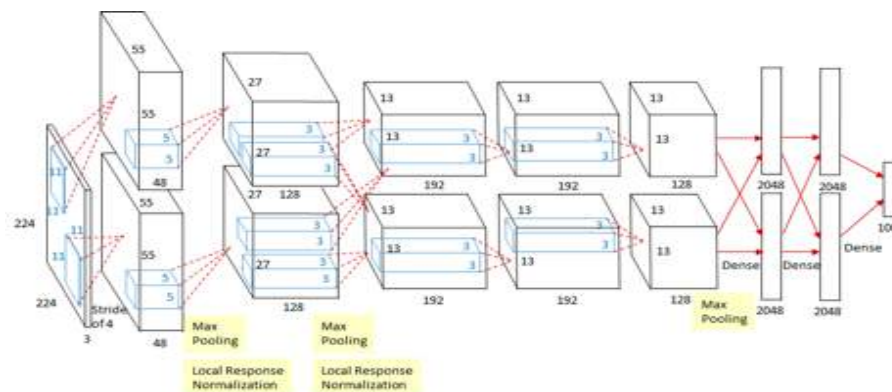
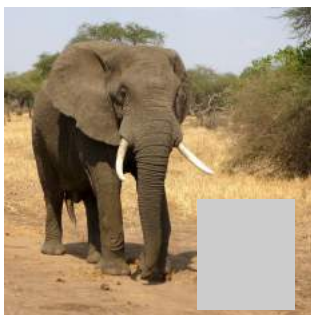
# Explanations

## Input Saliency via Occlusion



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

**Intervention: Mask part of the image before feeding to CNN, check how much predicted probabilities change**



$P(\text{elephant}) = 0.95$

A gray patch or patch of average pixel value of the dataset  
Note: not a black patch because the input images are centered to zero in the preprocessing.



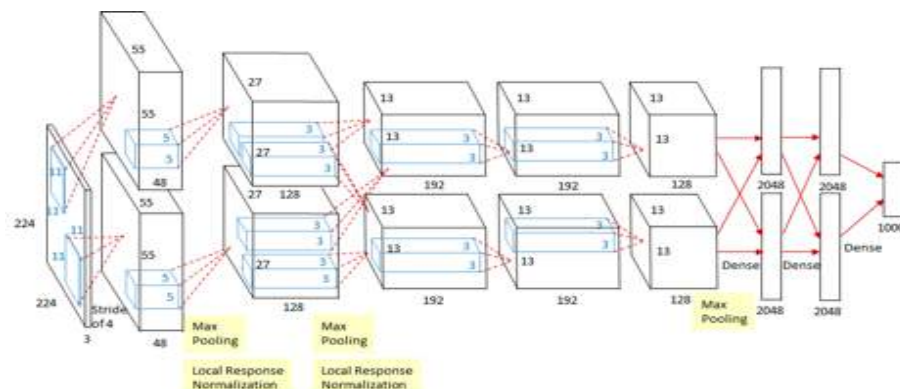
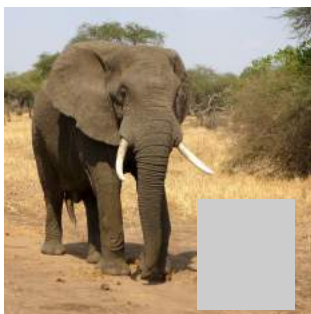
# Explanations

## Input Saliency via Occlusion



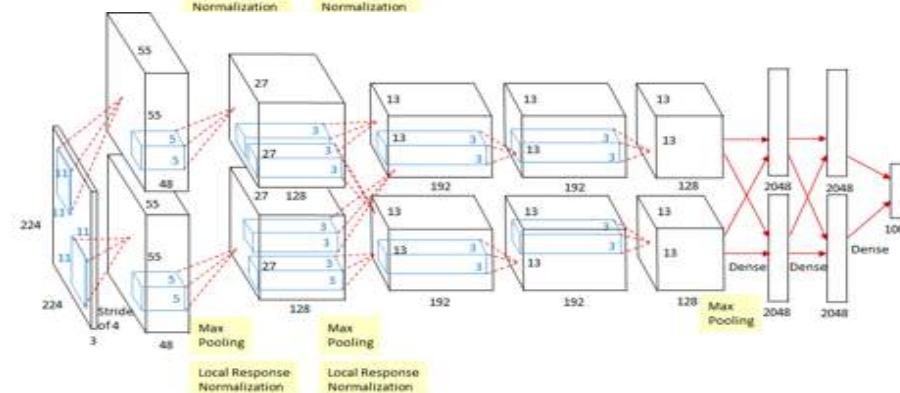
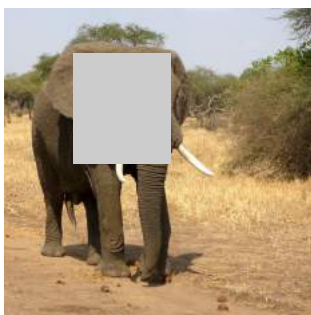
Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

**Intervention: Mask part of the image before feeding to CNN, check how much predicted probabilities change**



$P(\text{elephant}) = 0.95$

These pixels affect decisions more



$P(\text{elephant}) = 0.75$

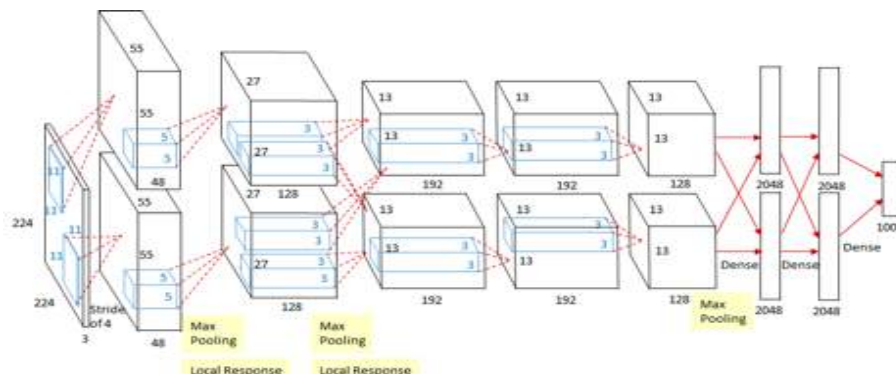
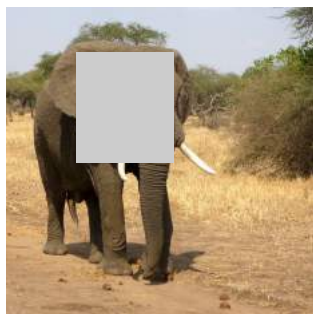
# Explanations

## Input Saliency via Occlusion

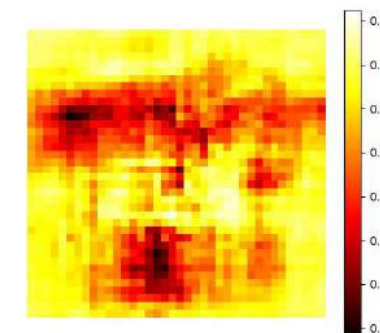
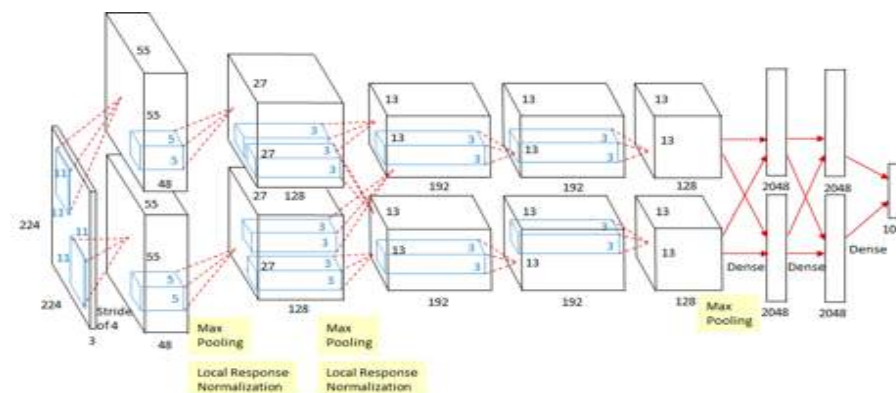
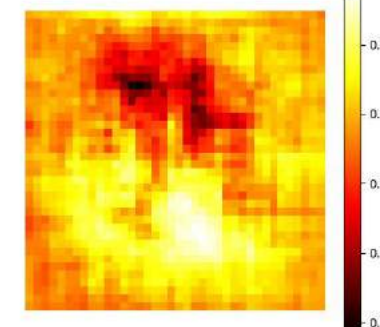


Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

The network is trained with image- labels, but it is sensitive to the common visual regions in images



African elephant, *Loxodonta africana*



# Explanations

## Input Saliency via Gradients



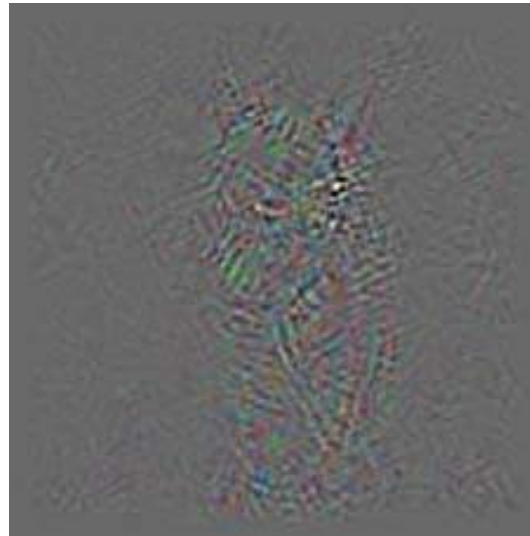
Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

**Gradients provide a one-shot means of perturbing the input that changes the output**

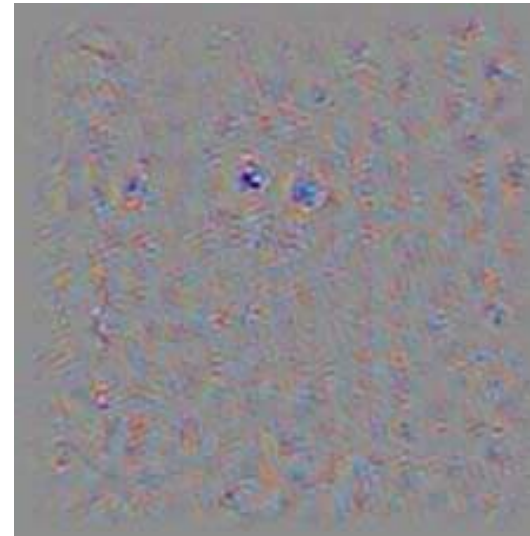
Input



Vanilla Gradients



Deconvolution Gradients



Guided Backpropagation



**However, localization remains an issue**

# Gradient and Activation-based Explanations

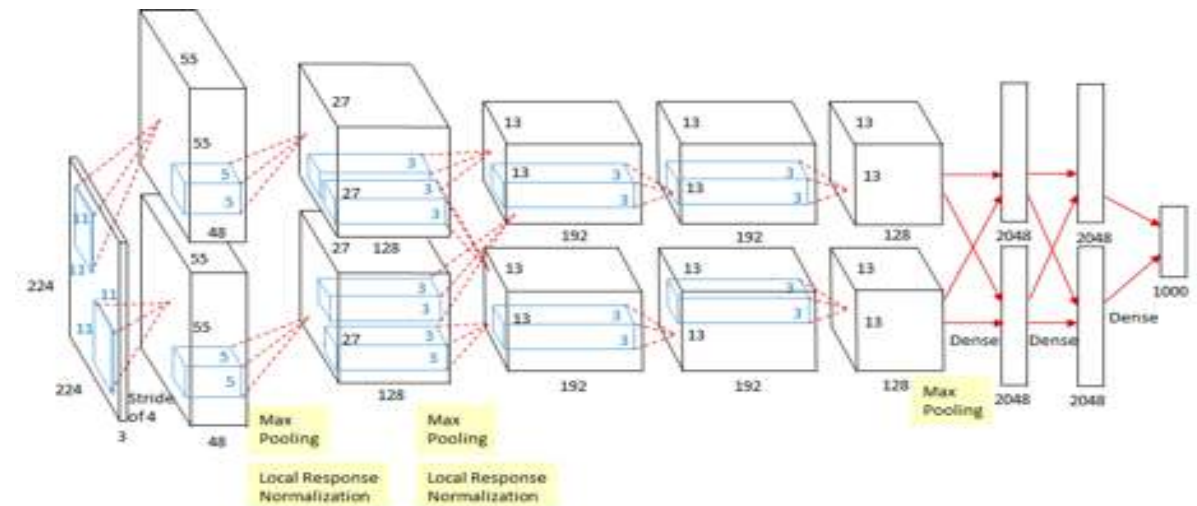
## GradCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

**Gradients provide a one-shot means of perturbing the input that changes the output.  
Activations provide the localization.**

- To find the important activations that are responsible for a particular class
- We want the activations:
  - **Class-discriminative** to reflect decision-making
  - **Preserve spatial information** to ensure spatial coverage of important regions



# Gradient and Activation-based Explanations

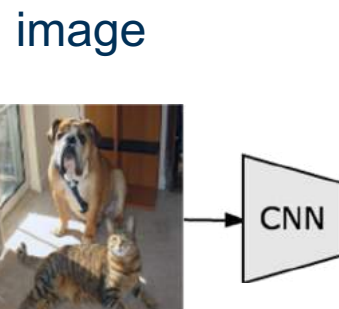
## GradCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

**Gradients provide a one-shot means of perturbing the input that changes the output.  
Activations provide the localization.**

- Given an image, feed forward through CNN



# Gradient and Activation-based Explanations

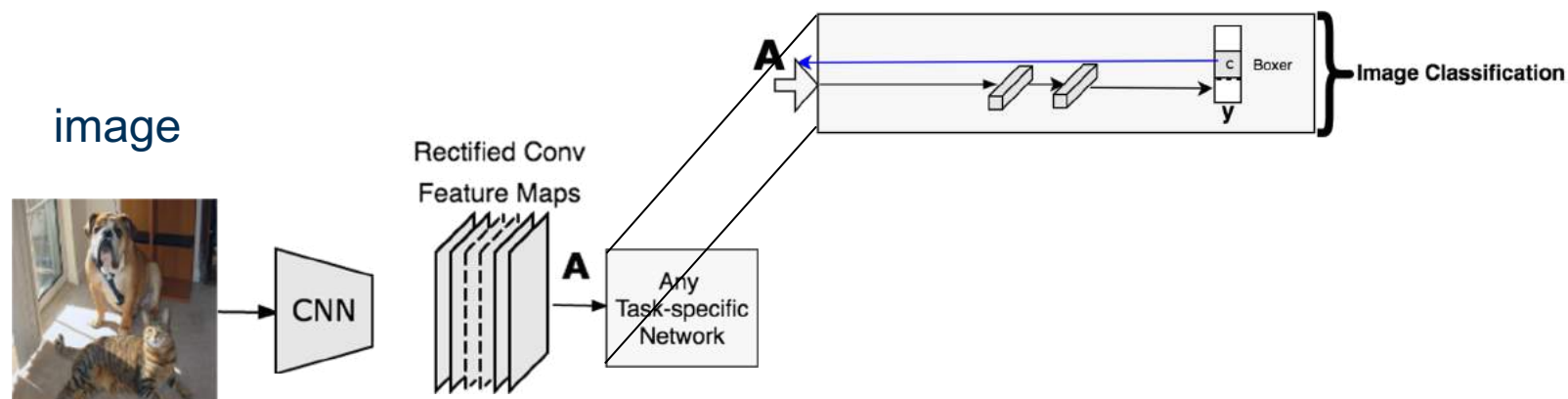
## GradCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

**Gradients provide a one-shot means of perturbing the input that changes the output.  
Activations provide the localization.**

- Given an image, feed forward through CNN
- Final convolutional layer output feature maps for later task-specific layers, i.e., fc layer for classification



# Gradient and Activation-based Explanations

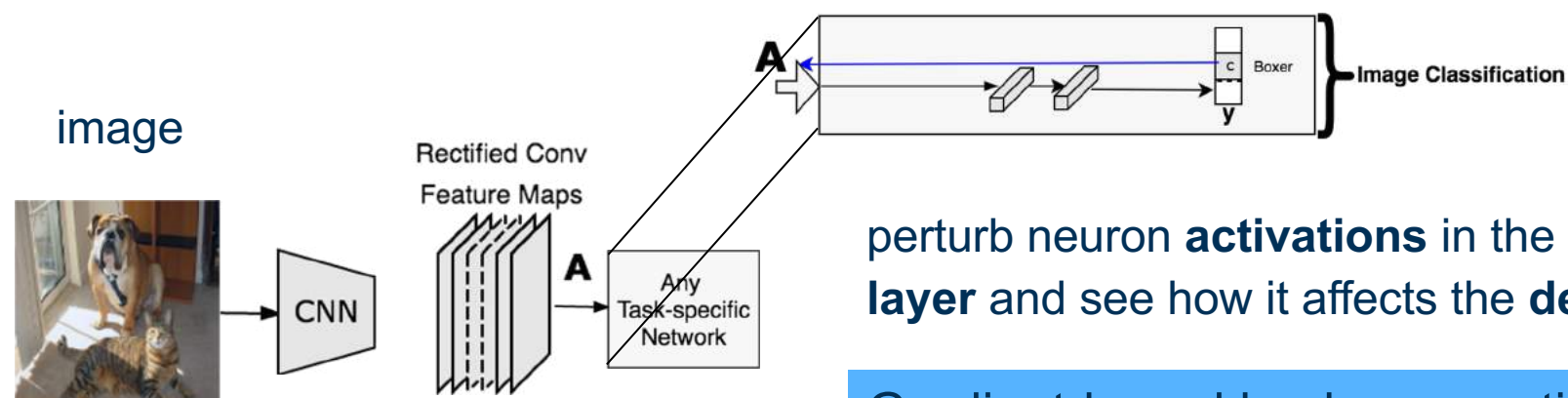
## GradCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

**Gradients provide a one-shot means of perturbing the input that changes the output. Activations provide the localization.**

- Given an image, feed forward through CNN
- Final convolutional layer output feature maps for later task-specific layers, i.e., fc layer for classification



**perturb neuron activations in the last conv layer and see how it affects the decision**

**Gradient-based backpropagation to obtain activation importance**

# Gradient and Activation-based Explanations

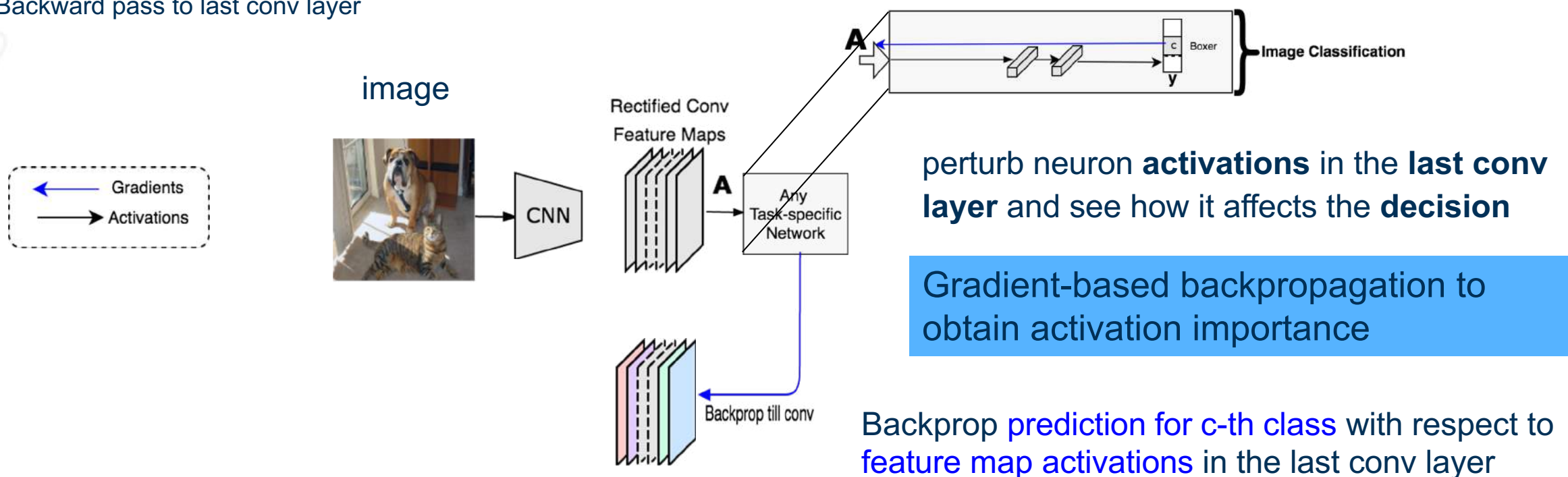
## GradCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

**Gradients provide a one-shot means of perturbing the input that changes the output.  
Activations provide the localization.**

- Given an image, feed forward through CNN
- Final convolutional layer output feature maps for later task-specific layers, i.e., fc layer for classification
- Backward pass to last conv layer





# Gradient and Activation-based Explanations

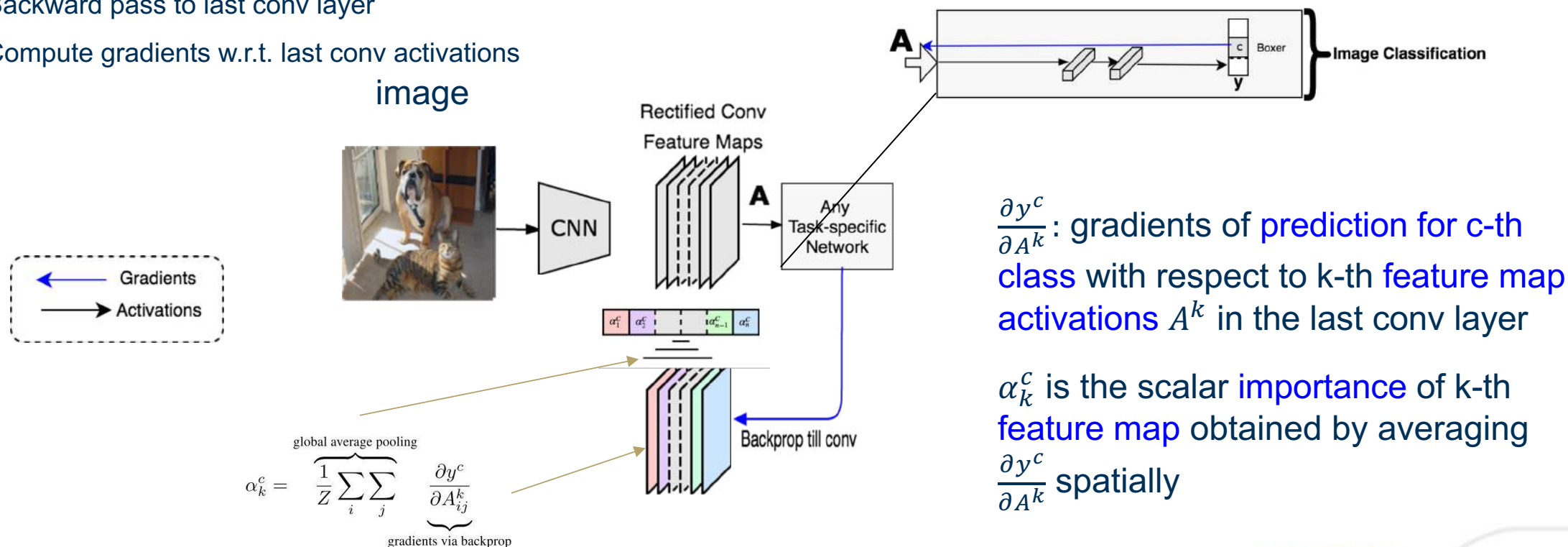
## GradCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

**Gradients provide a one-shot means of perturbing the input that changes the output. Activations provide the localization.**

- Given an image, feed forward through CNN
- Final convolutional layer output feature maps for later task-specific layers, i.e., fc layer for classification
- Backward pass to last conv layer
- Compute gradients w.r.t. last conv activations



[Tutorial] | [Ghassan AlRegib and Mohit Prabhushankar] | [June 4, 2023]

Selvaraju, Ramprasaath R., et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization." *Proceedings of the IEEE international conference on computer vision*. 2017.



# Gradient and Activation-based Explanations

## GradCAM

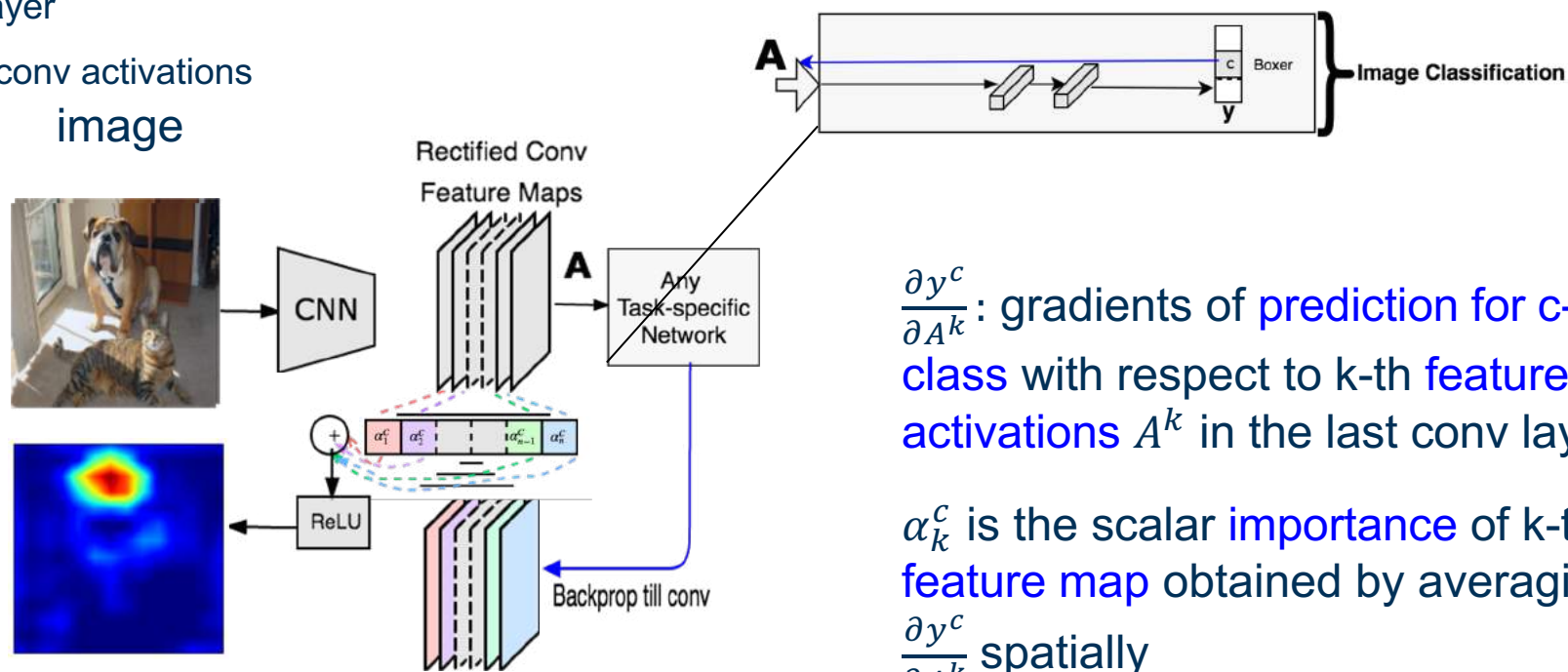
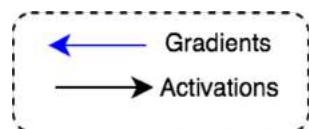


Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

**Gradients provide a one-shot means of perturbing the input that changes the output. Activations provide the localization.**

- Given an image, feed forward through CNN
- Final convolutional layer output feature maps for later task-specific layers, i.e., fc layer for classification
- Backward pass to last conv layer
- Compute gradients w.r.t. last conv activations

image



$\frac{\partial y^c}{\partial A^k}$  : gradients of prediction for c-th class with respect to k-th feature map activations  $A^k$  in the last conv layer

$\alpha_k^c$  is the scalar importance of k-th feature map obtained by averaging  $\frac{\partial y^c}{\partial A^k}$  spatially

Grad-CAM (up-sampled to original image dimension)

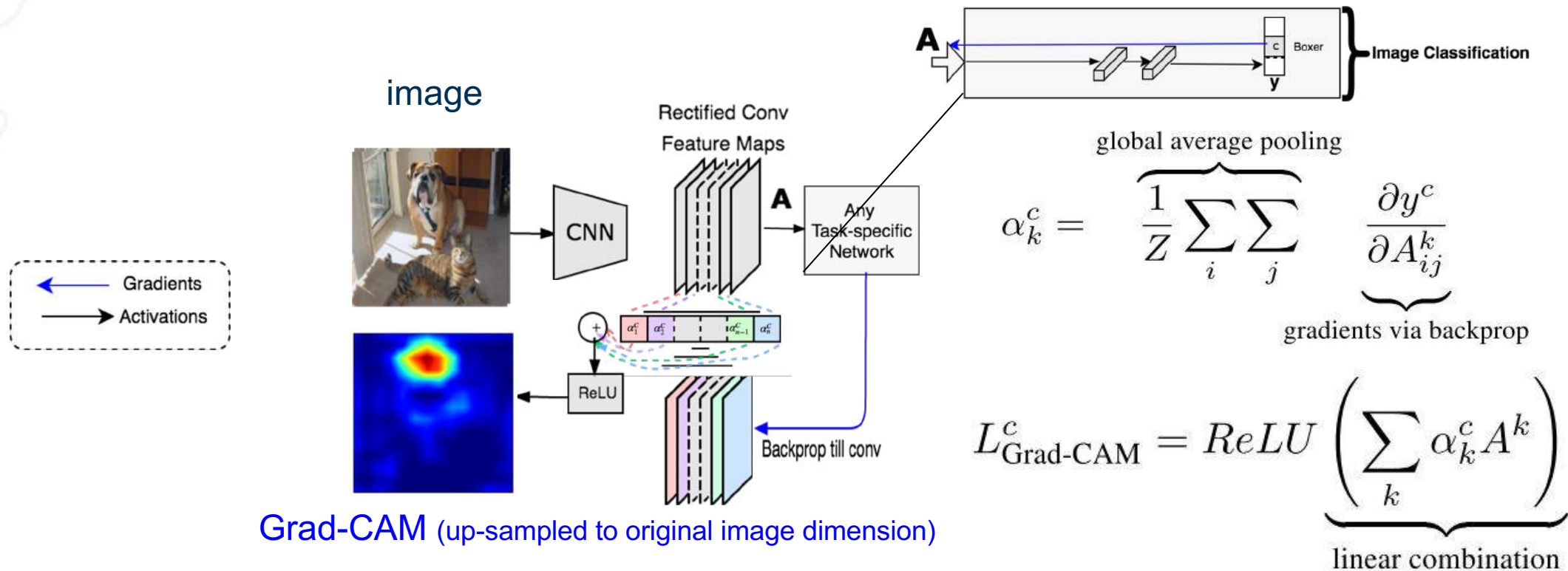
# Gradient and Activation-based Explanations

## GradCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Grad-CAM uses the gradient information flowing into the last convolutional layer of the CNN to assign importance values to each activation for a particular decision of interest.



# Gradient and Activation-based Explanations

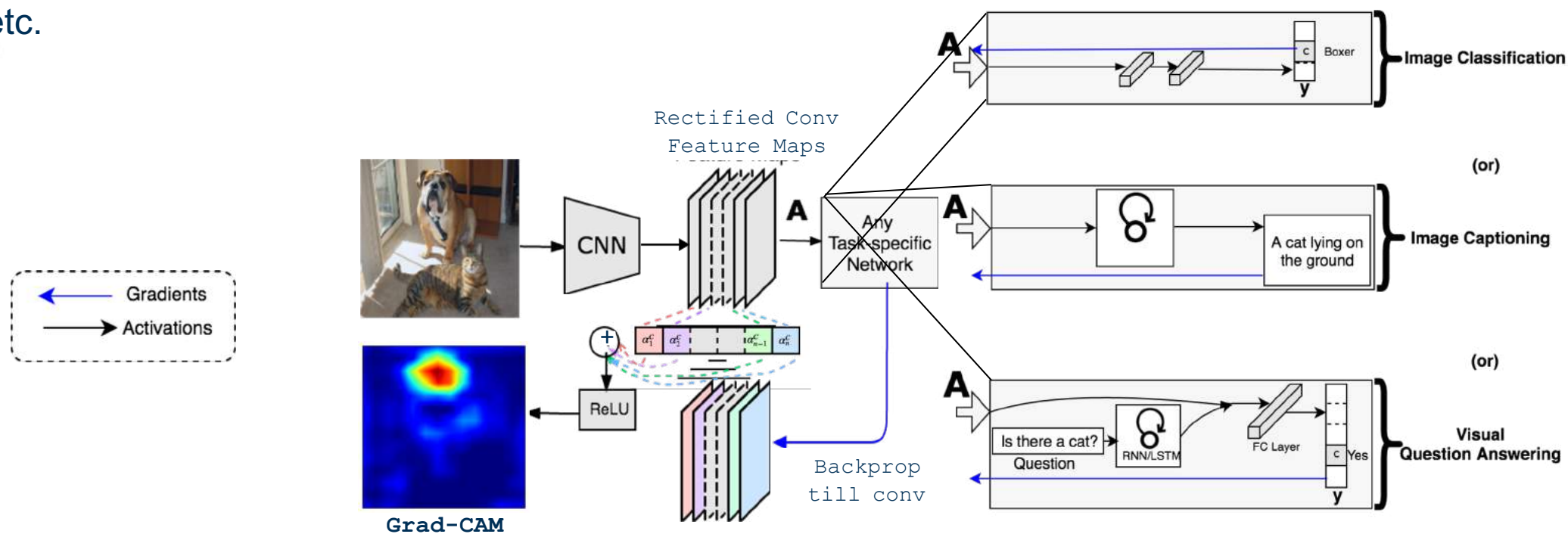
## GradCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Grad-CAM generalizes to any task:

- Image classification
- Image captioning
- Visual question answering
- etc.



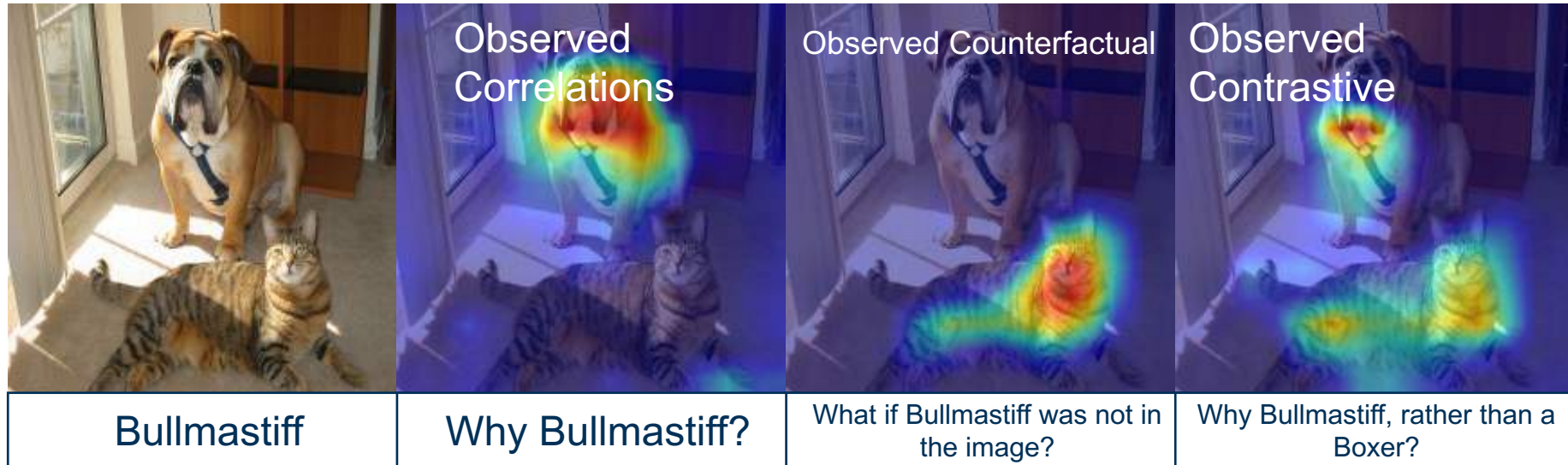
# Gradient and Activation-based Explanations

## Extensions of GradCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

**GradCAM provides answers to ‘Why P?’ questions. But different stakeholders require relevant and contextual explanations**



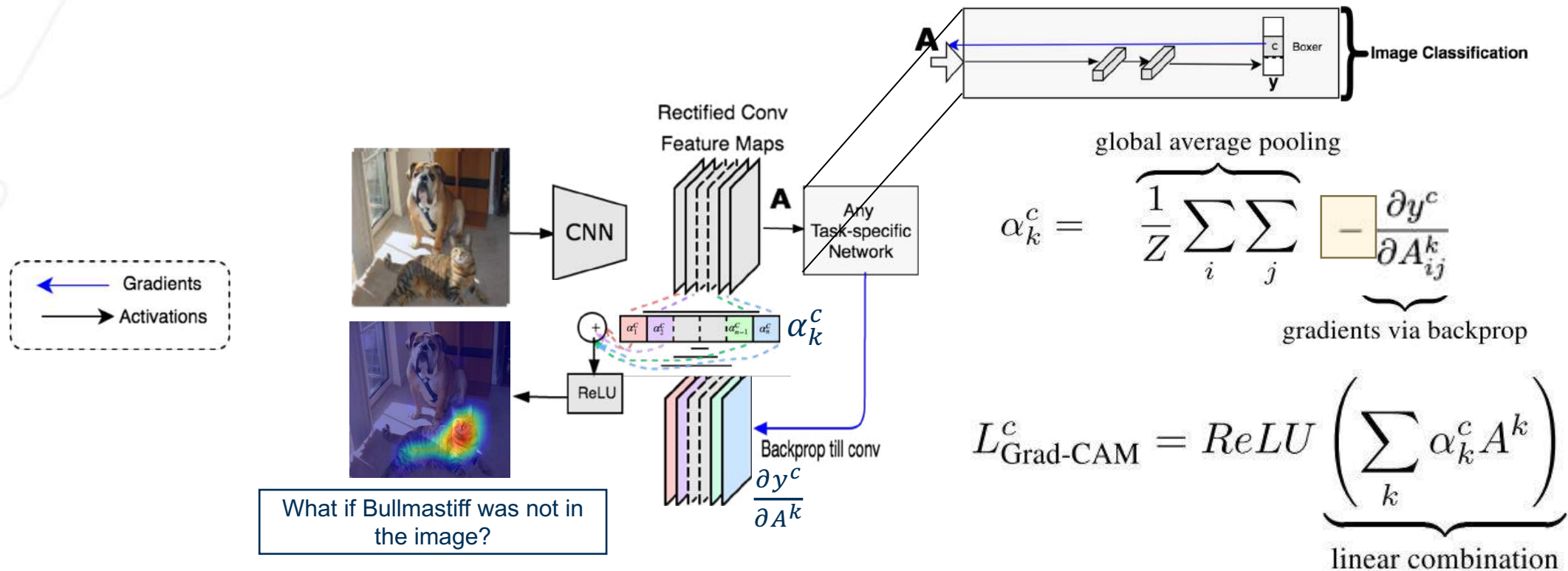
# Gradient and Activation-based Explanations

CounterfactualCAM: What if P is not there in the Image?



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

In GradCAM, global average pool the **negative of** gradients to obtain  $\alpha^c$  for each kernel  $k$



**Negating the gradients effectively removes these regions from analysis**

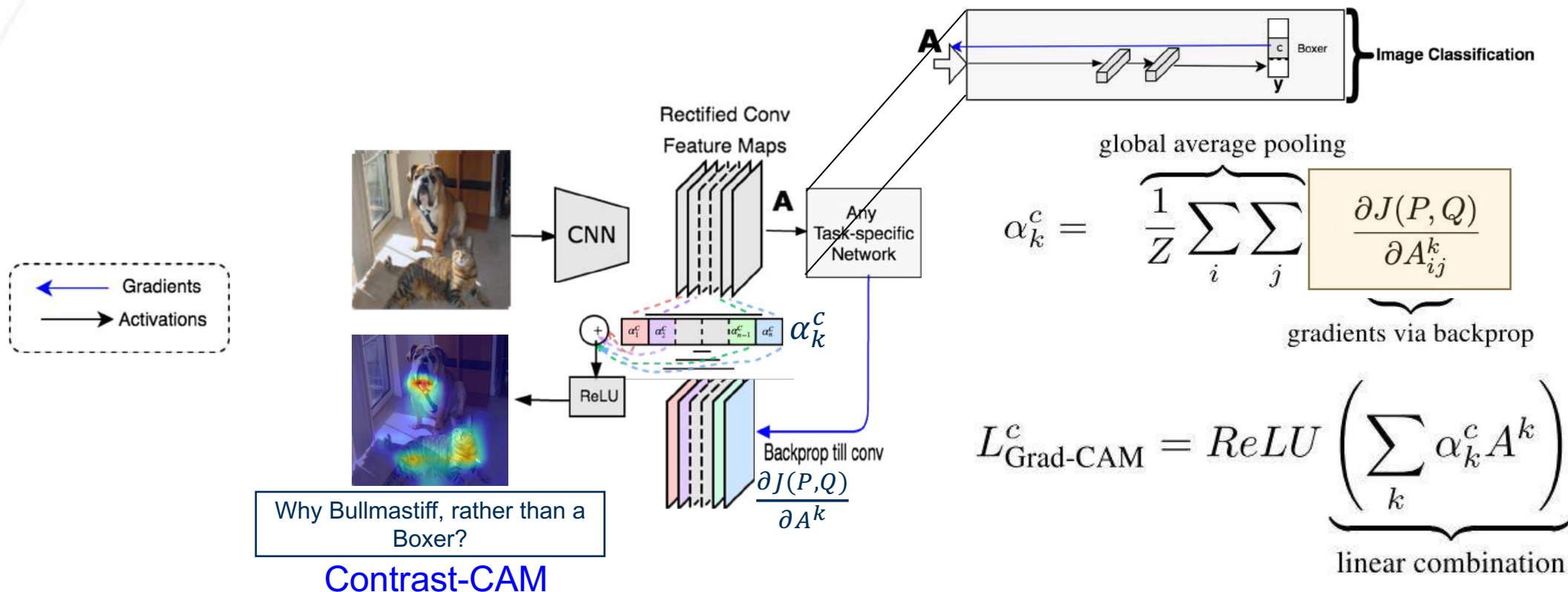
# Gradient and Activation-based Explanations

## ContrastCAM: Why P, rather than Q?



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

In GradCAM, backward pass the **loss between predicted class P and some contrast class Q** to last conv layer



Backpropagating the loss highlights the differences between classes P and Q.

# Gradient and Activation-based Explanations

## Results of GardCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? Why not Spoonbill, with 100% confidence?
ImageNet dataset : Bull Mastiff	Grad-CAM : Why : Bull Mastiff?	Representative Boxer image	Why Bull Mastiff, rather than Boxer?	Representative Blue jay image	Why Bull Mastiff, rather than Blue jay? Why not Bull Mastiff, with 100% confidence?
CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop? Why not No-Left with 100% confidence?
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6? Why not Bugatti with 100% confidence?

[Tutorial] | [Ghassan AlRegib and Mohit Prabhushankar] | [June 4, 2023]

Prabhushankar, M., Kwon, G., Temel, D., & AlRegib, G. (2020, October). Contrastive explanations in neural networks. In *2020 IEEE International Conference on Image Processing (ICIP)* (pp. 3289-3293). IEEE.





# Gradient and Activation-based Explanations

## Results of GardCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? Why not Spoonbill, with 100% confidence?
ImageNet dataset : Bull Mastiff	Grad-CAM : Why : Bull Mastiff?	Representative Boxer image	Why Bull Mastiff, rather than Boxer?	Representative Blue jay image	Why Bull Mastiff, rather than Blue jay? Why not Bull Mastiff, with 100% confidence?
CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop? Why not No-Left with 100% confidence?
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6? Why not Bugatti with 100% confidence?

Human Interpretable

# Gradient and Activation-based Explanations

## Results of GardCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? Why not Spoonbill, with 100% confidence?
ImageNet dataset : Bull Mastiff	Grad-CAM : Why : Bull Mastiff?	Representative Boxer image	Why Bull Mastiff, rather than Boxer	Representative Blue jay image	Why Bull Mastiff, rather than Blue jay? Why not Bull Mastiff, with 100% confidence?
CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop? Why not No-Left with 100% confidence?
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6? Why not Bugatti with 100% confidence?

Human Interpretable

Same as Grad-CAM

# Gradient and Activation-based Explanations

## Results of GardCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? Why not Spoonbill, with 100% confidence?
ImageNet dataset : Bull Mastiff	Grad-CAM : Why : Bull Mastiff?	Representative Boxer image	Why Bull Mastiff, rather than Boxer	Representative Blue jay image	Why Bull Mastiff, rather than Blue jay? Why not Bull Mastiff, with 100% confidence?
CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop? Why not No-Left with 100% confidence?
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6? Why not Bugatti with 100% confidence?

Human Interpretable

Same as Grad-CAM

Not Human Interpretable

# Gradient and Activation-based Explanations

## Results of GardCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? Why not Spoonbill, with 100% confidence?

Human Interpretable

Same as Grad-CAM



CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop?	Why not No-Left with 100% confidence?
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6?	Why not Bugatti with 100% confidence?



# Gradient and Activation-based Explanations

## Results of GardCAM, CounterfactualCAM, and ContrastCAM



Explanatory Paradigms in Neural Networks: Towards Relevant and Contextual Explanations

Input Image	Grad-CAM	Contrast 1	Contrastive Explanation 1	Contrast 2	Contrastive Explanation 2
ImageNet dataset : Spoonbill	Grad-CAM : Why Spoonbill?	Representative Flamingo image	Why Spoonbill, rather than Flamingo?	Representative Pig image	Why Spoonbill, rather than Pig? with 100% confidence?

Only traffic sign with a straight bottom-left edge – enough to say 'Not STOP Sign'

Human  
e.g. all  
Same as Grad-CAM



CURE-TSR dataset : No-Left Image	Grad-CAM : Why No-Left?	Representative No-Right image	Why No-Left, rather than No-Right?	Representative Stop Sign	Why No-Left, rather than Stop?	Why not No-Left with 100% confidence?	
Stanford Cars Dataset: Bugatti Convertible	Grad-CAM: Why Bugatti Convertible?	Representative Bugatti Coupe image	Why Convertible, rather than Coupe?	Representative Audi A6 image	Why Bugatti, rather than Audi A6?	Why not Bugatti with 100% confidence?	

# Case Study 1: Leveraging anomaly scores, uncertainty scores, and explanations for Robust Recognition



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



Mohit Prabhushankar, PhD  
Postdoc



Ghassan AlRegib, PhD  
Professor



# Robustness in Neural Networks

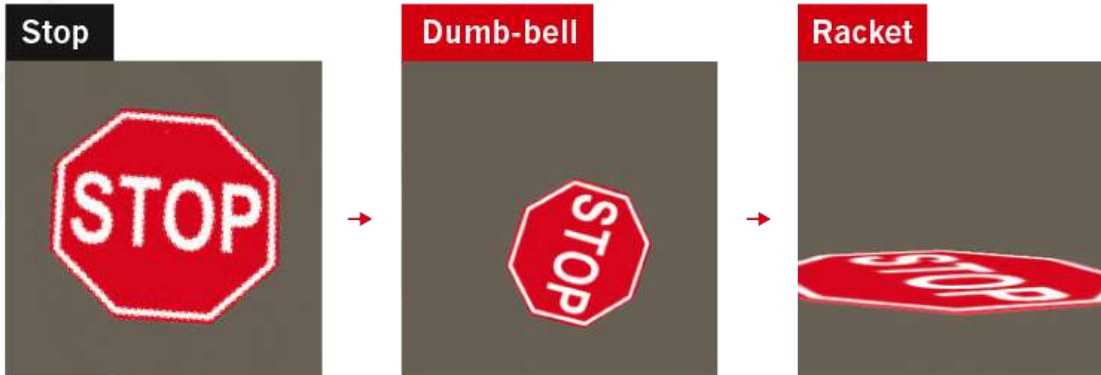
## Why Robustness?



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

### 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.



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.



©nature



# 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?





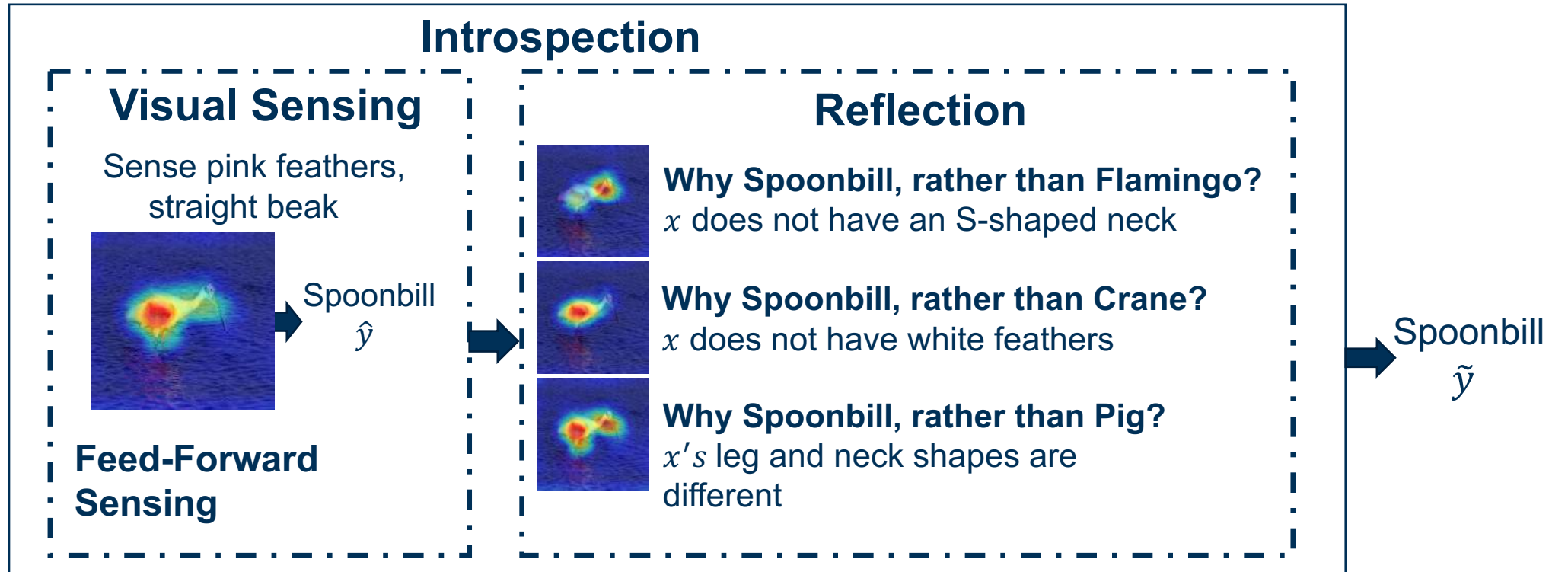
# Introspection

## 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?**

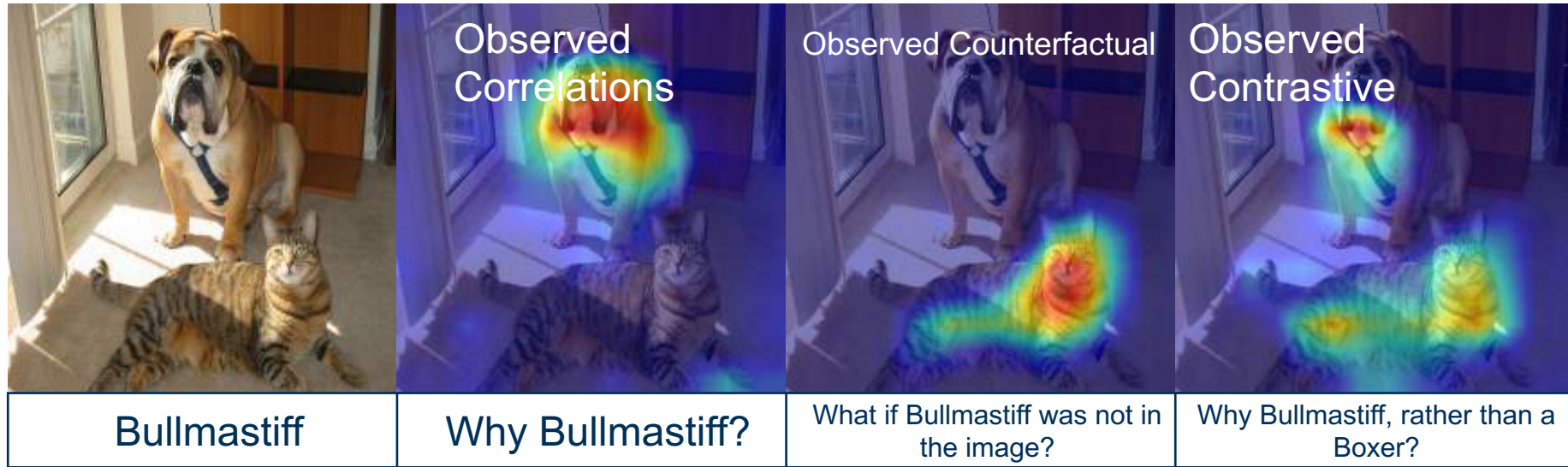
# Introspection

## 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  $x$ .*

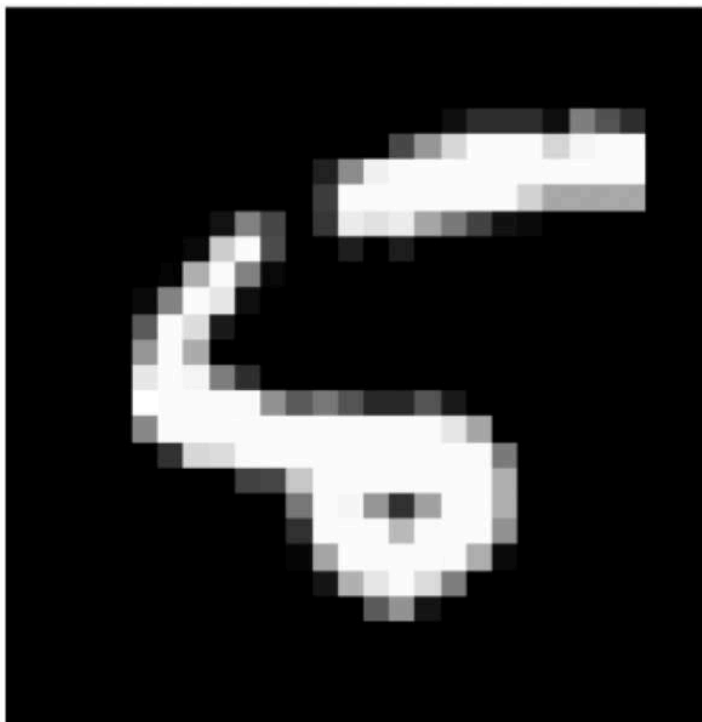
# Introspection in Neural Networks

## Gradients as Features

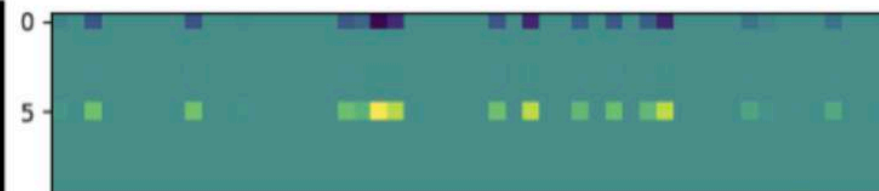


Introspective Learning: A Two-stage Approach for Inference in Neural Networks

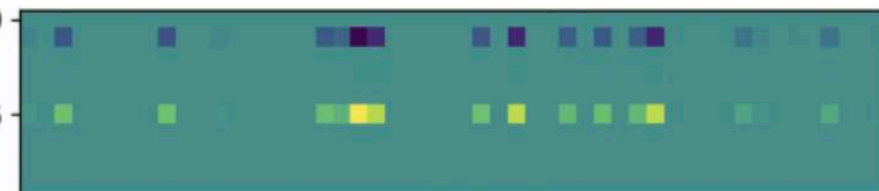
For a well-trained network, the gradients are sparse and informative



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?

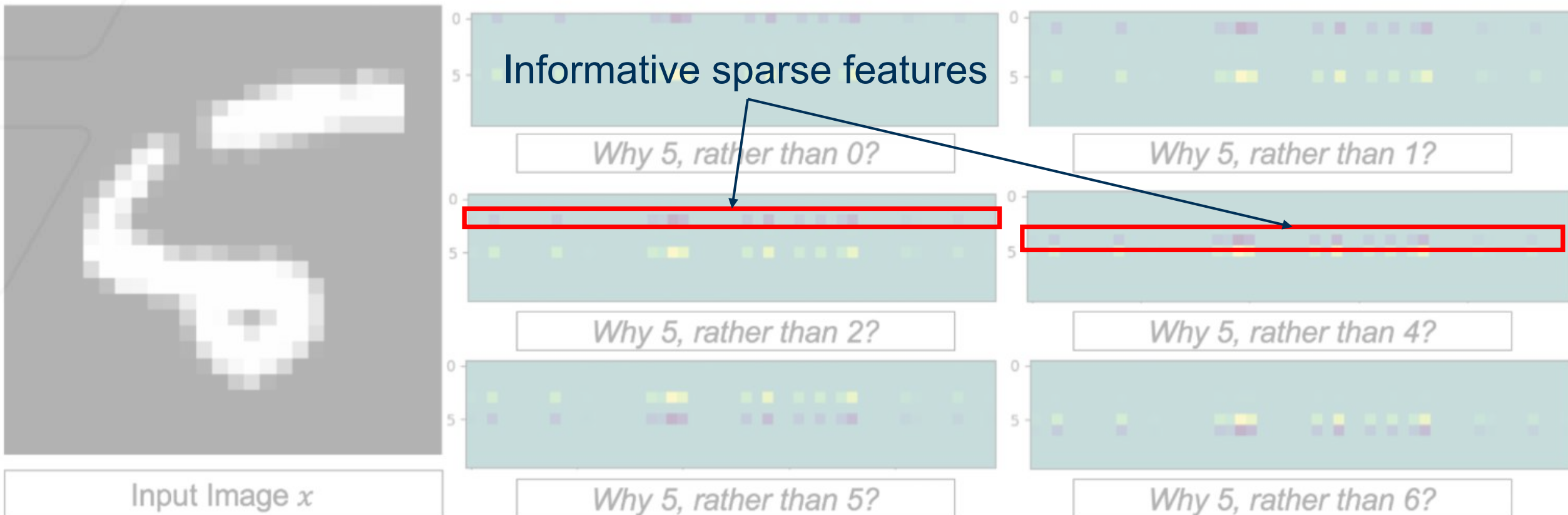
# Introspection in Neural Networks

## 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



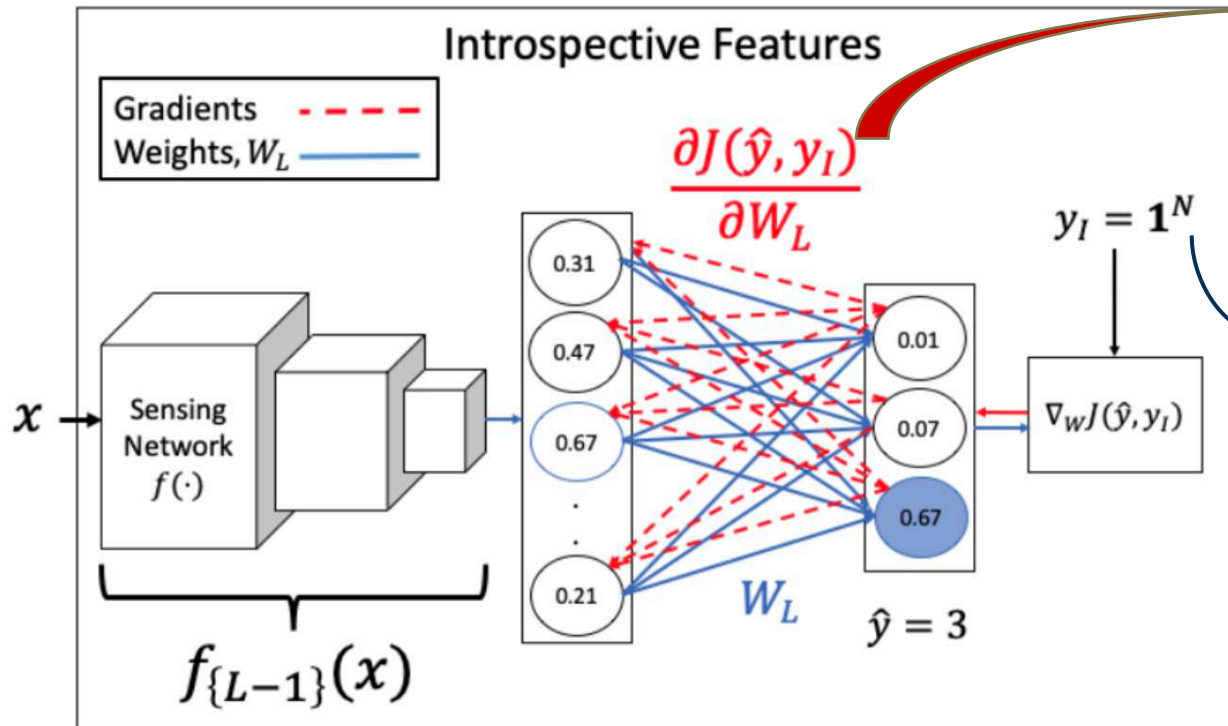
# Introspection in Neural Networks

## 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

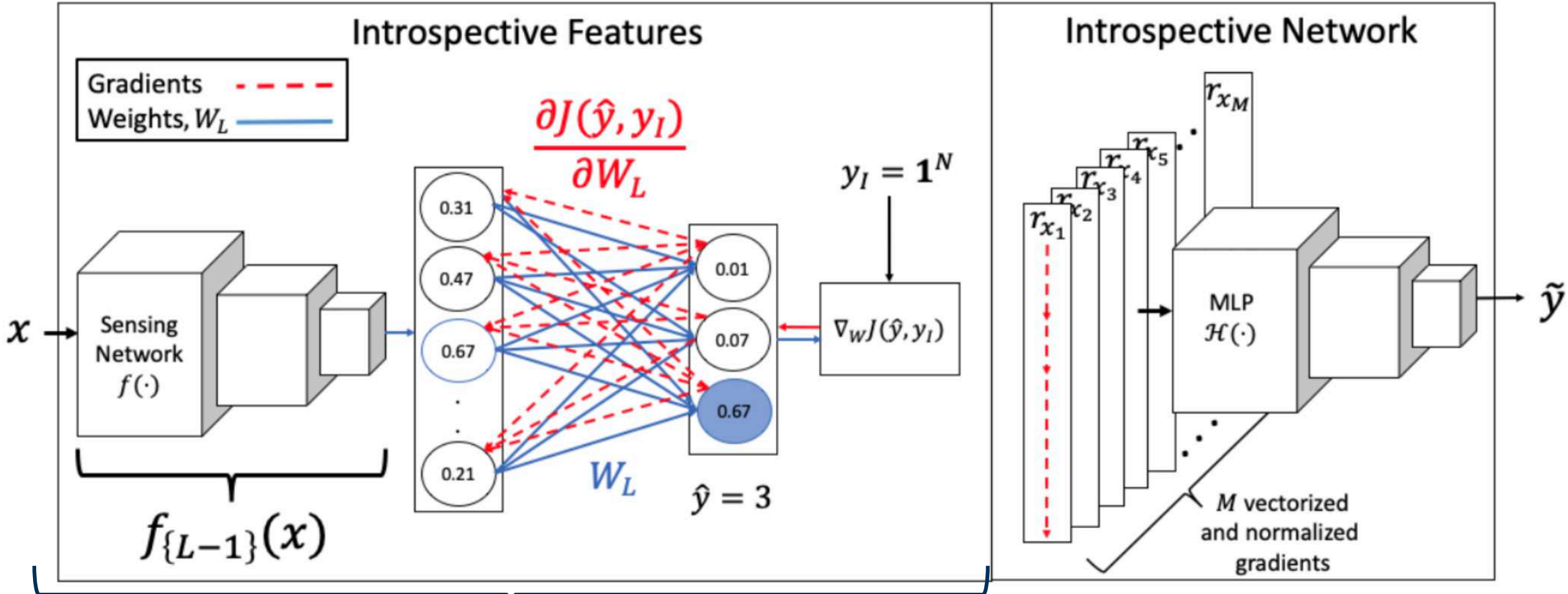
**Vector of all ones: A confounding label!**

# Introspection in Neural Networks

## Utilizing Gradient Features



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



## Introspective Features

[Tutorial] | [Ghassan AlRegib and Mohit Prabhushankar] | [June 4, 2023]

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.



# Introspection in Neural Networks

## When is Introspection Useful?



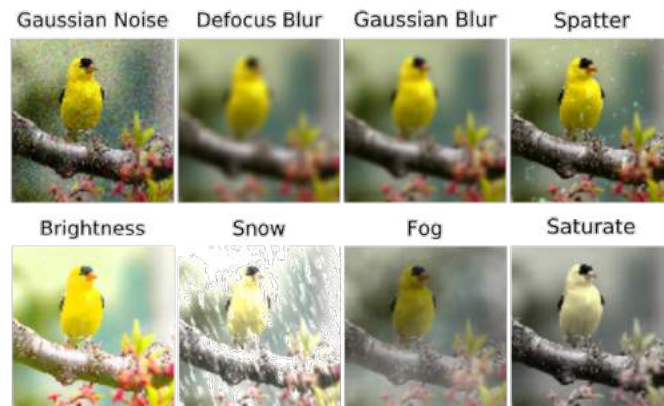
Introspective Learning: A Two-stage Approach for Inference in Neural Networks

**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**



# Introspection in Neural Networks

## Generalization and Calibration

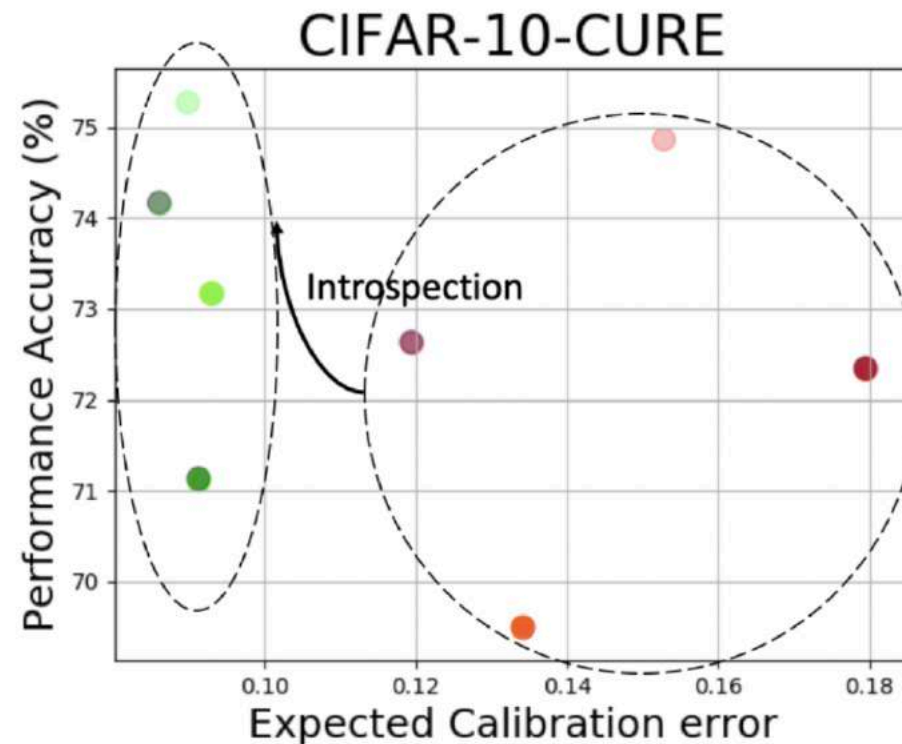
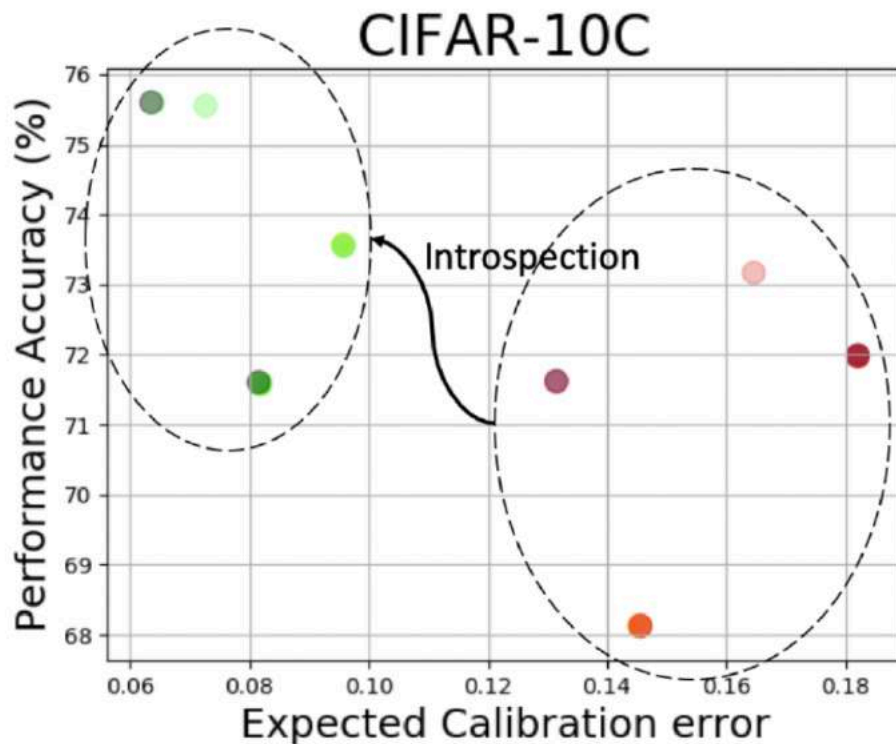


Introspective Learning: A Two-stage Approach for Inference in Neural Networks

Ideal: Top-left corner

Y-Axis: Generalization

X-Axis: Calibration



# 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

## Generalization and Calibration



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 (21)	Feed-Forward	0.365	0.358	0.244	0.249
	Introspective	0.365	0.359	<b>0.258</b>	<b>0.255</b>
Least (21)	Feed-Forward	0.371	0.359	0.252	0.25
	Introspective	0.373	0.362	<b>0.264</b>	<b>0.26</b>
Margin (22)	Feed-Forward	0.38	0.369	0.251	0.253
	Introspective	0.381	0.373	<b>0.265</b>	<b>0.263</b>
BALD (24)	Feed-Forward	0.393	0.368	0.26	0.253
	Introspective	0.396	0.375	<b>0.273</b>	<b>0.263</b>
BADGE (25)	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 (25)	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 (26)	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

[Tutorial] | [Ghassan AlRegib and Mohit Prabhushankar] | [June 4, 2023]

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.

## Case Study 2: Leveraging anomaly scores, uncertainty scores, and explanations for Anomalous object classification



# Detecting and Classifying Anomalies in Artificial Intelligence Systems



Gukyeong Kwon, PhD  
Amazon AWS



Mohit Prabhushankar, PhD  
Postdoc, Georgia Tech

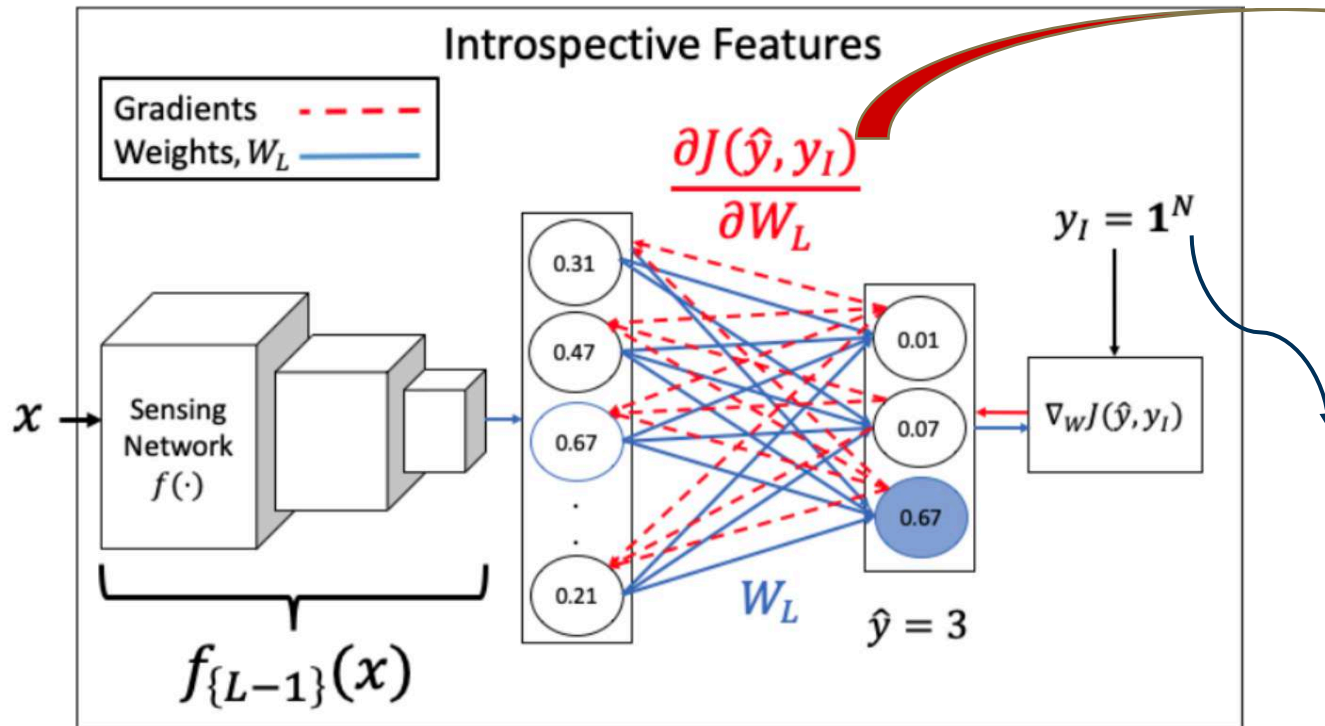


Ghassan AlRegib, PhD  
Professor, Georgia Tech

# Aberrant Object Detection

## Deriving Gradient Features

Measure the loss between the prediction  $\hat{y}$  and a vector of all ones and backpropagate to obtain the introspective features



**Uncertainty:** We took energy of all gradients  
**Robustness:** We trained a new network  
**Aberrant Objects:** We take variance across gradients from object detector

# Aberrant Object Detection

## Aberrance Detection

### Uncertainty using variance of introspective gradients rather than energy of gradients

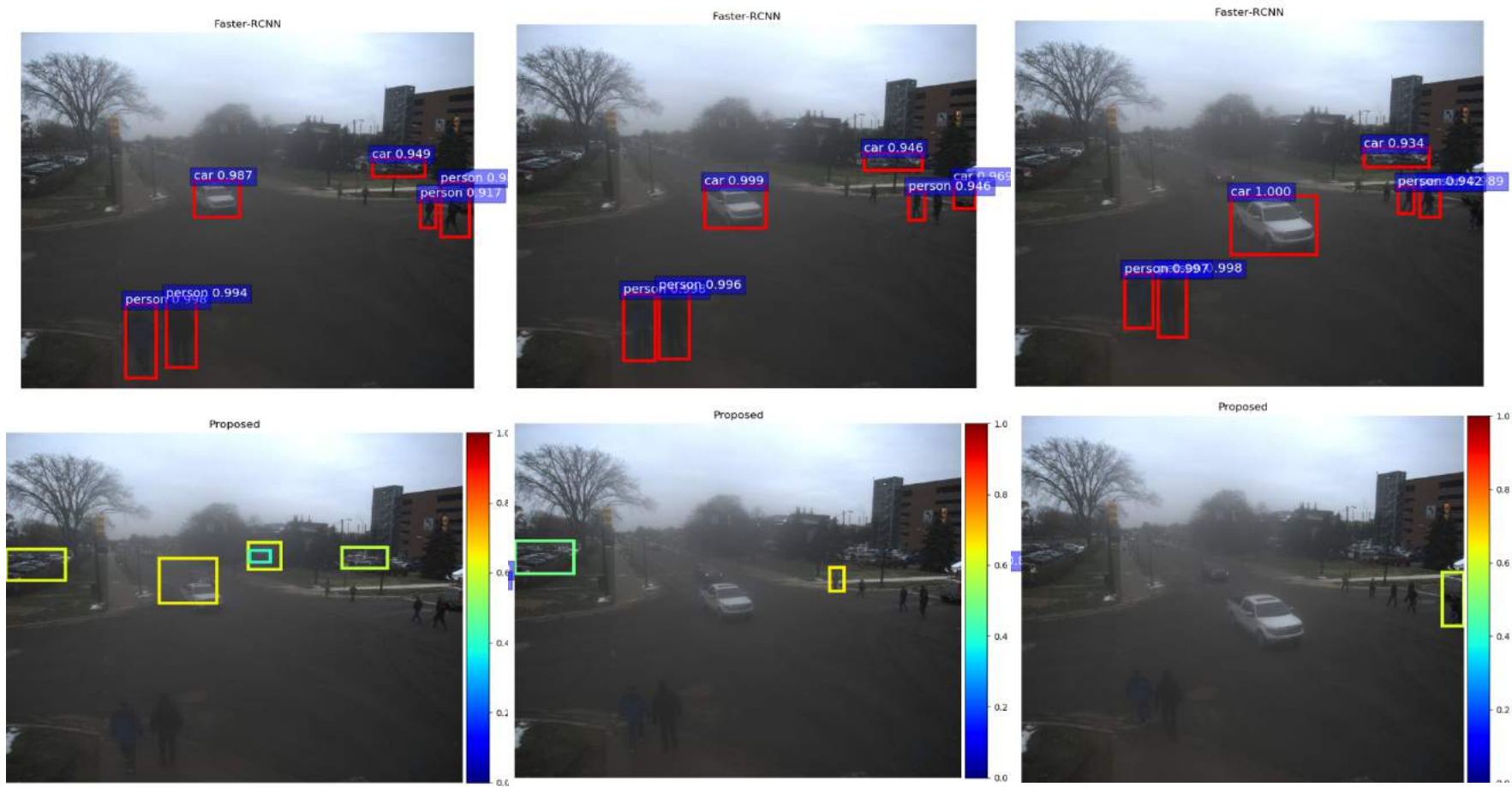


- Object detection algorithms would pick up on all the trained objects
- The gradient-based uncertainty approach picks up only the *aberrant* object – objects that bear a resemblance to novel classes

# Aberrant Object Detection

Complementary to object detectors

## Uncertainty using variance of introspective gradients rather than energy of gradients

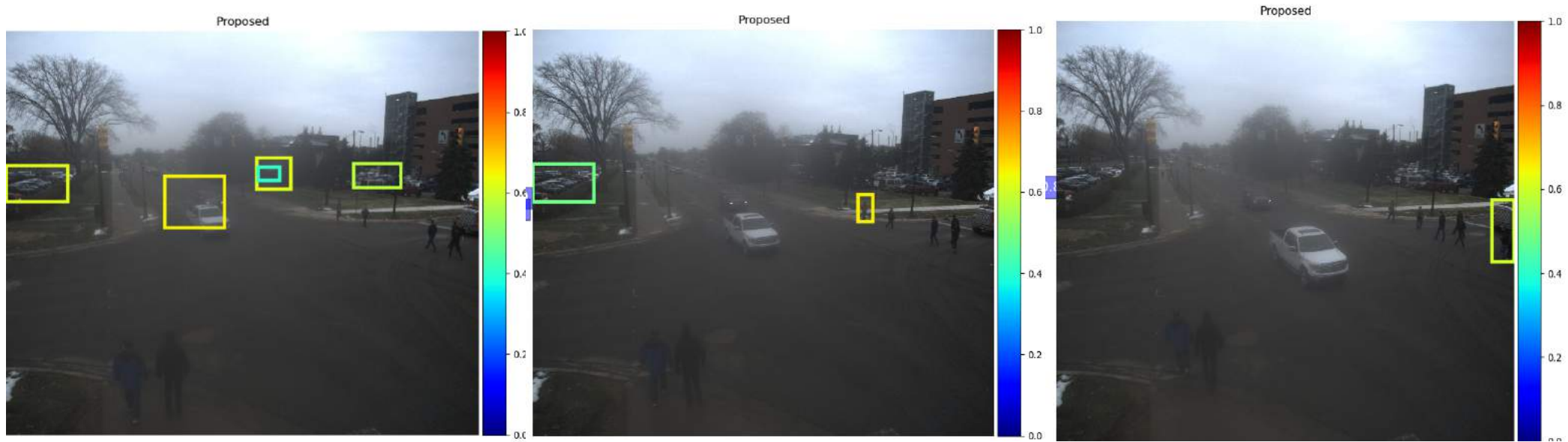




# Aberrant Object Detection

## Active Learning

Use the uncertain boxes for obtaining labels from annotators



Use new annotations for subsequent training in an active learning setting

# Objectives

## Takeaways from Part III

- Part I: Challenges in Perception and Autonomy
- Part II: Deep Learning for Perception
- **Part III: Existing Deep Learning solutions to Challenges in Perception**
  - It is not always clear if aberrant events and challenges must be incorporated in training
  - Instead, they can and should be equipped with diagnostic tools at predictions
  - These diagnostic tools are anomaly and uncertainty scores for decision making and contextual explainability for post-hoc stakeholders
  - Gradients provide the change induced by an aberrant event in the network and can be used to obtain the required prediction diagnosis
- Part IV: Key Takeaways and Future Directions

# A Holistic View of Perception in Intel. Vehicles

## Part IV: Key Takeaways and Future Directions

# Objectives

## Objectives in Part IV

- Takeaway Messages and Key Insights
- Unaddressed Challenges in Perception
  - Context Awareness
  - Embedded Perception
  - V2X Perception
- Future Research Directions
  - Temporal Processing
  - Sensor Processing Architectures
  - Sensors research
  - Infrastructure + AV Datasets

# Objectives

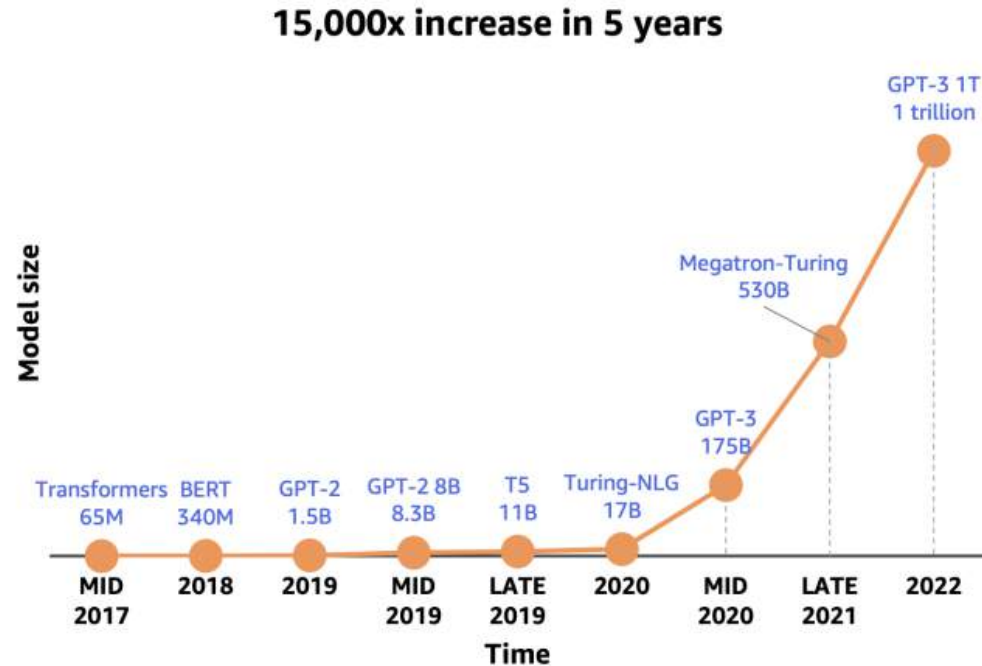
## Takeaway Messages and Key Insights

- **Robustness** under challenging conditions, environments, context and surroundings-awareness are **challenges** in AV perception
  - **Deep Learning** provides a **holistic solution** to a number of the above challenges
- **Transfer Learning** and **training at scale** help to create foundation models
  - **Self-supervised Learning** provides a framework for large scale learning on unannotated data
- It is not always clear if aberrant events and challenges must be incorporated in training
  - Instead, **model predictions** must be equipped with **diagnostic tools** at inference
  - These diagnostic tools are **anomaly and uncertainty scores** for decision making and **contextual explainability** for post-hoc stakeholders
  - **Gradients** provide the change induced by an aberrant event in the network and can be used to obtain the required **prediction diagnosis**

# Perception in AVs

## Unaddressed Technical Challenges for Level 3 Automation

- Challenging weather
- Challenging sensing
- Challenging environments
- **Context awareness**
- **Embedded perception**
- **V2X perception**



- Foundation models are great but the real-time feasibility is an issue
- The inaccuracies from model outputs is dangerous in urban settings

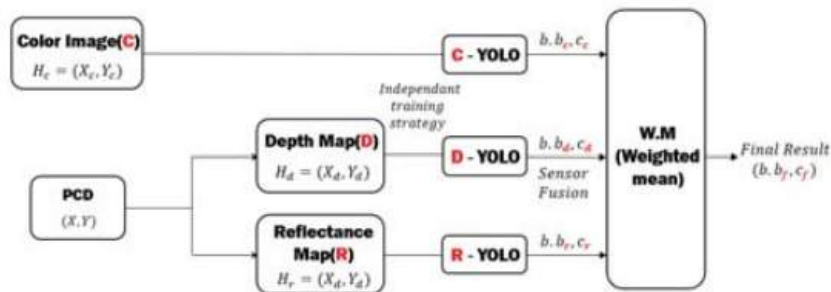
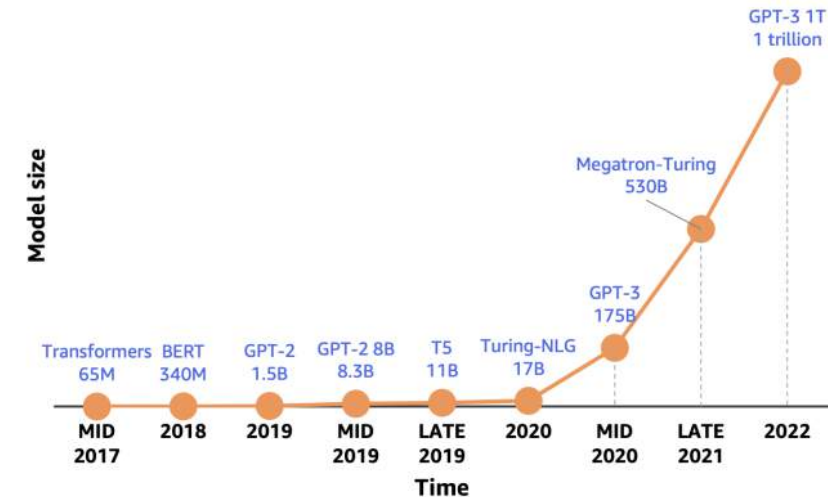
# Perception in AVs

## Unaddressed Technical Challenges for Levels 4 and 5

### Foundation models with multiple sensor modalities

- Challenging weather
- Challenging sensing
- Challenging environments
- **Context awareness**
- **Embedded perception**
- **V2X perception**

15,000x increase in 5 years



- Levels 4 and 5 automation relies on roadside infrastructure to obtain high-resolution predictions
- 10x is the rough estimate of the increase in processing power between levels of automation

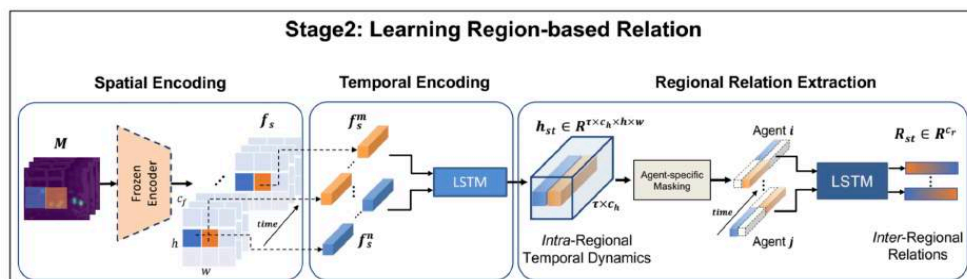
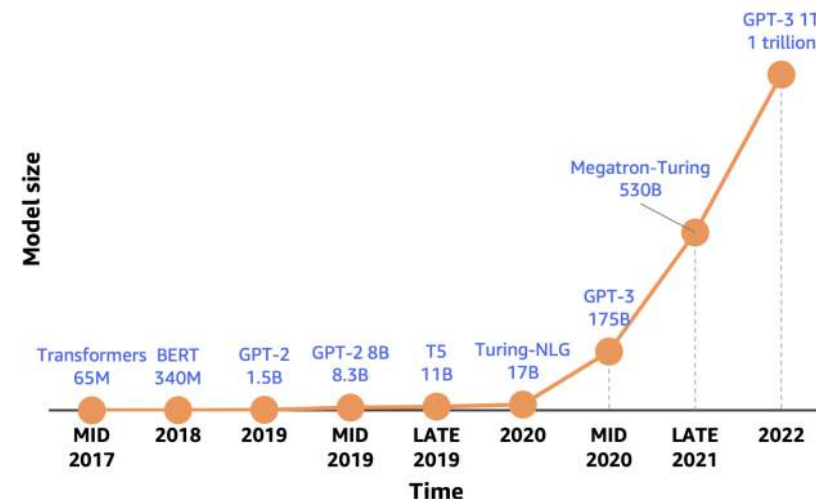
# Perception in AVs

## Unaddressed Technical Challenges for Levels 4 and 5

### Foundation models with multiple sensor modalities and on temporal data

- Challenging weather
- Challenging sensing
- Challenging environments
- **Context awareness**
- **Embedded perception**
- **V2X perception**

15,000x increase in 5 years



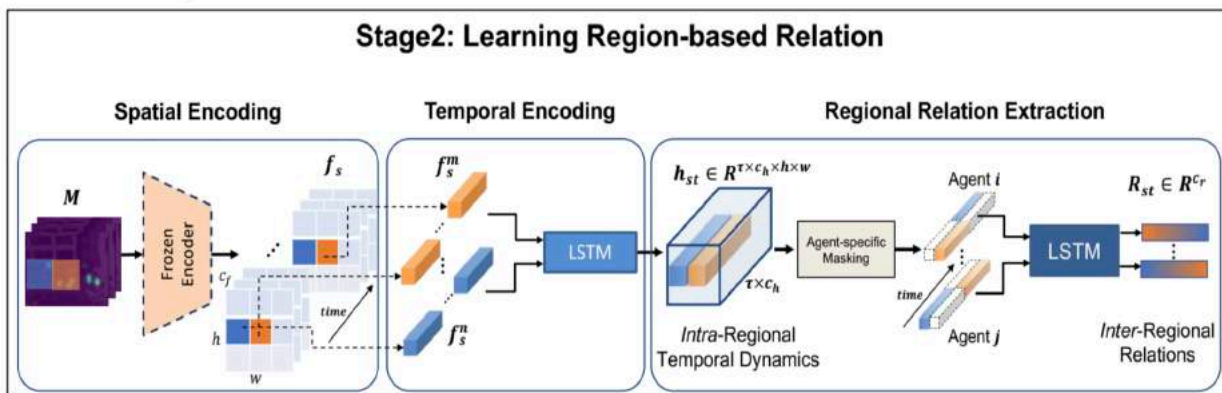
- Levels 4 and 5 automation relies on roadside infrastructure to obtain high-resolution predictions
- 10x is the rough estimate of the increase in processing power between levels of automation
- **Current temporal processing = linear spatial processing in time**



# Future Direction 1

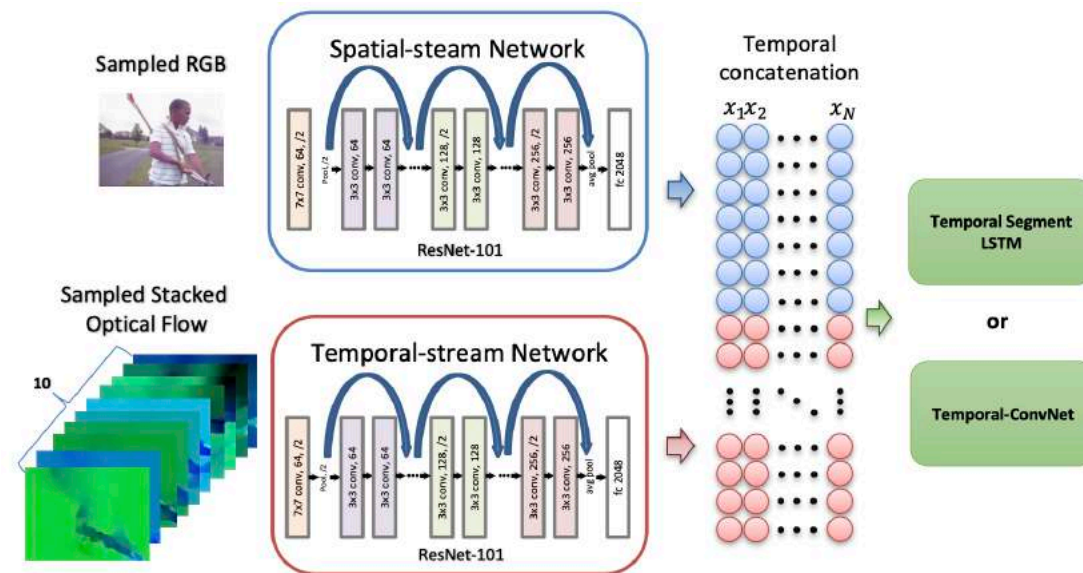
## Temporal processing of data

### Temporal processing $\neq$ Linear spatial processing



**Late temporal fusion:** Encode all spatial data in a time-wise fashion and determine temporal relationships

**Early temporal fusion:** Encode both spatial and temporal information together and fuse them within the network



# Future Direction 2

## Sensor processing architectures



Vision data processing was revolutionized by CNNs

Language data processing was revolutionized by Transformers

LIDAR data processing is revolutionized by ?

RADAR data processing is revolutionized by ?

...

# Future Direction 3

More data with less sensors!

## 4 Fisheye cameras provide a 360 degree surround view of the car

Results from Zero-shot (i.e. using the trained model out of the box) Segment Anything Model on Woodscape dataset



Important context and objects are not segmented

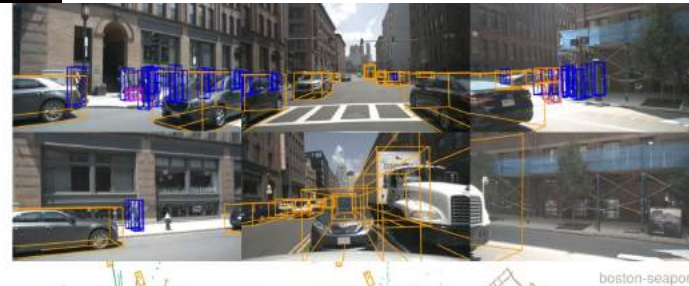
# Future Direction 4

## Infrastructure + AV Datasets

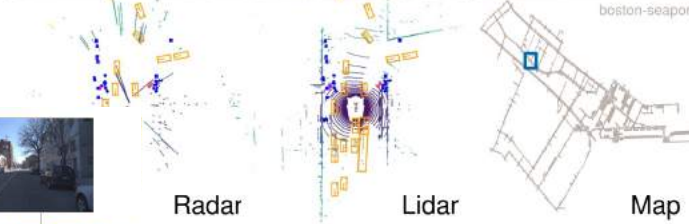
**Abundance of egocentric AV datasets! Dearth of Infrastructure + AV datasets**



Argoverse



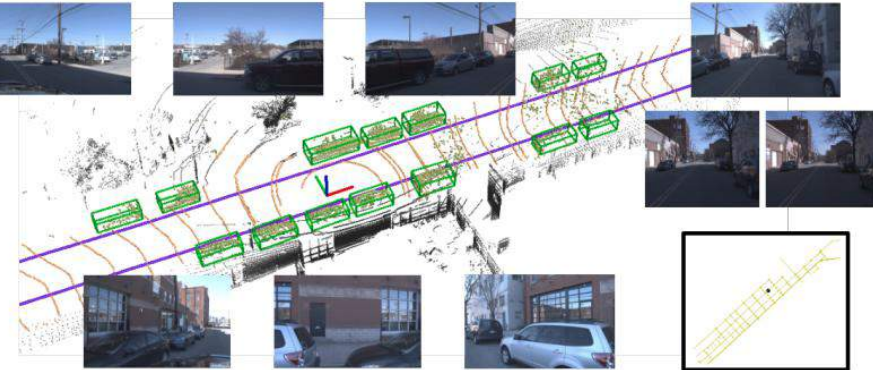
boston-seaport



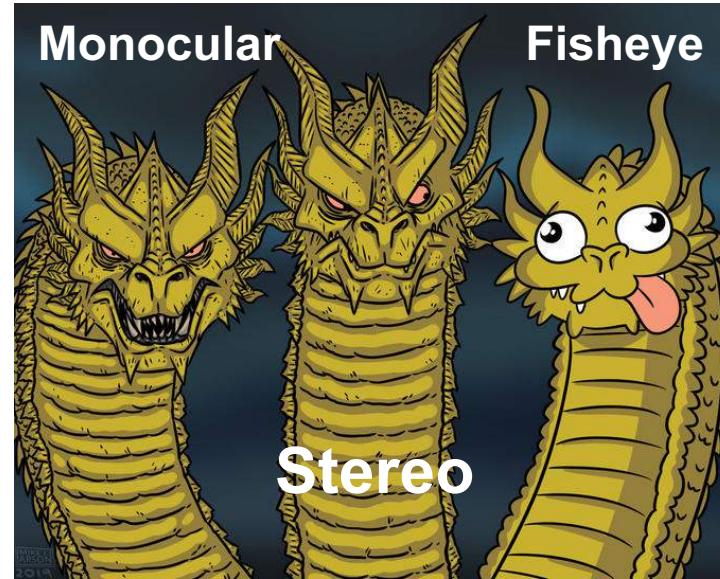
Radar Lidar Map  
"icycle, car makes a u-turn, lane change, peds crossing crosswalk"

NuScenes

- Infrastructure datasets: Stationary sensors at traffic junctures, streets, heavy pedestrian traffic areas etc.
- Infrastructure + AV datasets: Egocentric sensors on vehicles + stationary sensors for the same scenes



# Some Memes to Wrap it Up



# References

## Gradient representations for Robustness, OOD, Anomaly, Novelty, and Adversarial Detection

- **Gradients for robustness against noise:** 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
- **Gradients for adversarial, OOD, corruption detection:** J. Lee, M. Prabhushankar, and G. AlRegib, "Gradient-Based Adversarial and Out-of-Distribution Detection," in *International Conference on Machine Learning (ICML) Workshop on New Frontiers in Adversarial Machine Learning*, Baltimore, MD, Jul. 2022.
- **Gradients for Open set recognition:** Lee, Jinsol, and Ghassan AlRegib. "Open-Set Recognition With Gradient-Based Representations." *2021 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2021.
- **GradCon for Anomaly Detection:** Kwon, G., Prabhushankar, M., Temel, D., & AlRegib, G. (2020, August). Backpropagated gradient representations for anomaly detection. In *European Conference on Computer Vision* (pp. 206-226). Springer, Cham.
- **Gradients for adversarial, OOD, corruption detection :** J. Lee, C. Lehman, M. Prabhushankar, and G. AlRegib, "Probing the Purview of Neural Networks via Gradient Analysis," in IEEE Access, Mar. 21 2023.
- **Gradients for Novelty Detection:** Kwon, G., Prabhushankar, M., Temel, D., & AlRegib, G. (2020, October). Novelty detection through model-based characterization of neural networks. In *2020 IEEE International Conference on Image Processing (ICIP)* (pp. 3179-3183). IEEE.
- **Gradient-based Image Quality Assessment:** G. Kwon\*, M. Prabhushankar\*, D. Temel, and G. AlRegib, "Distorted Representation Space Characterization Through Backpropagated Gradients," in *IEEE International Conference on Image Processing (ICIP)*, Taipei, Taiwan, Sep. 2019.

## Explainability in Neural Networks

- **Explanatory paradigms:** AlRegib, G., & Prabhushankar, M. (2022). Explanatory Paradigms in Neural Networks: Towards relevant and contextual explanations. *IEEE Signal Processing Magazine*, 39(4), 59-72.
- **Contrastive Explanations:** Prabhushankar, M., Kwon, G., Temel, D., & AlRegib, G. (2020, October). Contrastive explanations in neural networks. In *2020 IEEE International Conference on Image Processing (ICIP)* (pp. 3289-3293). IEEE.
- **Explainability in Limited Label Settings:** M. Prabhushankar, and G. AlRegib, "Extracting Causal Visual Features for Limited Label Classification," in IEEE International Conference on Image Processing (ICIP), Sept. 2021.
- **Explainability through Expectancy-Mismatch:** M. Prabhushankar and G. AlRegib, "Stochastic Surprisal: An Inferential Measurement of Free Energy in Neural Networks," in *Frontiers in Neuroscience, Perception Science*, Volume 17, Feb. 09 2023.

# References

## Self Supervised Learning

- **Weakly supervised Contrastive Learning:** K. Kokilepersaud, S. Trejo Corona, M. Prabhushankar, G. AlRegib, C. Wykoff, "Clinically Labeled Contrastive Learning for OCT Biomarker Classification," in *IEEE Journal of Biomedical and Health Informatics*, 2023, May. 15 2023.
- **Contrastive Learning for Fisheye Images:** K. Kokilepersaud, M. Prabhushankar, Y. Yarici, G. AlRegib, and A. Parchami, "Exploiting the Distortion-Semantic Interaction in Fisheye Data," in *Open Journal of Signals Processing*, Apr. 28 2023.
- **Contrastive Learning for Severity Detection:** K. Kokilepersaud, M. Prabhushankar, G. AlRegib, S. Trejo Corona, C. Wykoff, "Gradient Based Labeling for Biomarker Classification in OCT," in *IEEE International Conference on Image Processing (ICIP)*, Bordeaux, France, Oct. 16-19 2022
- **Contrastive Learning for Seismic Images:** K. Kokilepersaud, M. Prabhushankar, and G. AlRegib, "Volumetric Supervised Contrastive Learning for Seismic Semantic Segmentation," in *International Meeting for Applied Geoscience & Energy (IMAGE)*, Houston, TX, , Aug. 28-Sept. 1 2022

## Human Vision and Behavior Prediction

- **Pedestrian Trajectory Prediction:** C. Zhou, G. AlRegib, A. Parchami, and K. Singh, "TrajPRed: Trajectory Prediction With Region-Based Relation Learning," *IEEE Transactions on Intelligent Transportation Systems*, submitted on Dec. 28 2022.
- **Human Visual Saliency in trained Neural Nets:** Y. Sun, M. Prabhushankar, and G. AlRegib, "Implicit Saliency in Deep Neural Networks," in *IEEE International Conference on Image Processing (ICIP)*, Abu Dhabi, United Arab Emirates, Oct. 2020.
- **Human Image Quality Assessment:** D. Temel, M. Prabhushankar and G. AlRegib, "UNIQUE: Unsupervised Image Quality Estimation," in *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1414-1418, Oct. 2016.

## Open-source Datasets to assess Robustness

- **CURE-TSD:** D. Temel, M-H. Chen, and G. AlRegib, "Traffic Sign Detection Under Challenging Conditions: A Deeper Look Into Performance Variations and Spectral Characteristics," in *IEEE Transactions on Intelligent Transportation Systems*, Jul. 2019
- **CURE-TSR:** D. Temel, G. Kwon\*, M. Prabhushankar\*, and G. AlRegib, "CURE-TSR: Challenging Unreal and Real Environments for Traffic Sign Recognition," in *Advances in Neural Information Processing Systems (NIPS) Workshop on Machine Learning for Intelligent Transportation Systems*, Long Beach, CA, Dec. 2017
- **CURE-OR:** D. Temel\*, J. Lee\*, and G. AlRegib, "CURE-OR: Challenging Unreal and Real Environments for Object Recognition," in *IEEE International Conference on Machine Learning and Applications (ICMLA)*, Orlando, FL, Dec. 2018

# References

## Active Learning

- **Active Learning and Training with High Information Content:** R. Benkert, M. Prabhushankar, G. AlRegib, A. Parchami, and E. Corona, "Gaussian Switch Sampling: A Second Order Approach to Active Learning," in *IEEE Transactions on Artificial Intelligence (TAI)*, Feb. 05 2023
- **Active Learning Dataset on vision and LIDAR data:** Y. Logan, R. Benkert, C. Zhou, K. Kokilepersaud, M. Prabhushankar, G. AlRegib, K. Singh, E. Corona and A. Parchami, "FOCAL: A Cost-Aware Video Dataset for Active Learning," *IEEE Transactions on Circuits and Systems for Video Technology*, submitted on Apr. 29 2023
- **Active Learning on OOD data:** R. Benkert, M. Prabhushankar, and G. AlRegib, "Forgetful Active Learning With Switch Events: Efficient Sampling for Out-of-Distribution Data," in *IEEE International Conference on Image Processing (ICIP)*, Bordeaux, France, Oct. 16-19 2022
- **Active Learning for Biomedical Images:** Y. Logan, R. Benkert, A. Mustafa, G. Kwon, G. AlRegib, "Patient Aware Active Learning for Fine-Grained OCT Classification," in *IEEE International Conference on Image Processing (ICIP)*, Bordeaux, France, Oct. 16-19 2022

## Uncertainty Estimation

- **Gradient-based Uncertainty:** J. Lee and G. AlRegib, "Gradients as a Measure of Uncertainty in Neural Networks," in *IEEE International Conference on Image Processing (ICIP)*, Abu Dhabi, United Arab Emirates, Oct. 2020
- **Uncertainty Visualization in Seismic Images:** R. Benkert, M. Prabhushankar, and G. AlRegib, "Reliable Uncertainty Estimation for Seismic Interpretation With Prediction Switches," in *International Meeting for Applied Geoscience & Energy (IMAGE)*, Houston, TX, , Aug. 28-Sept. 1 2022.
- **Uncertainty and Disagreements in Label Annotations:** C. Zhou, M. Prabhushankar, and G. AlRegib, "On the Ramifications of Human Label Uncertainty," in *NeurIPS 2022 Workshop on Human in the Loop Learning*, Oct. 27 2022
- **Uncertainty in Saliency Estimation:** T. Alshawi, Z. Long, and G. AlRegib, "Unsupervised Uncertainty Estimation Using Spatiotemporal Cues in Video Saliency Detection," in *IEEE Transactions on Image Processing*, vol. 27, pp. 2818-2827, Jun. 2018.