

**E9: 309 ADL 23-12-2020**

<http://leap.ee.iisc.ac.in/sriram/teaching/ADL2020/>



# Housekeeping

## \* Midterm project II presentations ←

→ Done during Dec. 29,30th

→ Same format as previous evaluation

## \* Midterm project III ←

→ Abstract submission deadline (Jan 10th) ←

✓ Evaluation after final exam (1st week of Feb) ←

## \* Final Exam (as per IISc schedule) ←

✓ Jan ~~2~~nd afternoon!

(Jan 23<sup>rd</sup> final exam)



# Topics Discussed thus far

Explainable/Interpretable  
Deep Learning

Visualizing deep layer  
activations using tSNE

Interepoch evolution of  
activations

Transferability





# Topics Discussed thus far

Explainable/Interpretable  
Deep Learning

Visualizing deep layer  
activations using tSNE

Backpropagation based  
approach - Deconv net

Interepoch evolution of  
acvitations

Establishing hierarchal  
representation learning

Transferability

*Visualization  
in input*





# Topics Discussed thus far

Explainable/Interpretable  
Deep Learning

*Saliency*

Visualizing deep layer  
activations using tSNE

Backpropagation based  
approach - Deconv net

Using Global average  
pooling and  
interpolation

CAM

Interepoch evolution of  
acvitations

Establishing hierarchal  
representation learning

Using gradients from  
last layer

Grad-  
CAM

Transferability

Attention





# Topics Discussed thus far

Explainable/Interpretable  
Deep Learning

Visualizing deep layer  
activations using tSNE

Backpropagation based  
approach - Deconv net

Using Global average  
pooling and  
interpolation

CAM

Interepoch evolution of  
acvitations

Establishing hierarchal  
representation learning

Using gradients from  
last layer

Grad-  
CAM

Transferability

Attention

Causal  
Inference





# Causal inference

## \* Causal inference

→ Deriving the causal connection between conditions that cause an effect.

## \* Three levels of causation

→ **association** Seeing and observing the environment. ✓  
Is the incidence of lung cancer higher among smokers?

→ **intervention** Doing and intervening in the environment. ✓  
How do we reduce lung cancer? What is the effect if we ban cigarettes?

→ **counterfactuals** Imagining, retrospection, understanding the environment. ✓  
What if I had not smoked for the last two years?





# Pruning based approach to analyzing/compressing

Published as a conference paper at ICLR 2017

---

## PRUNING CONVOLUTIONAL NEURAL NETWORKS FOR RESOURCE EFFICIENT INFERENCE

**Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, Jan Kautz**

NVIDIA

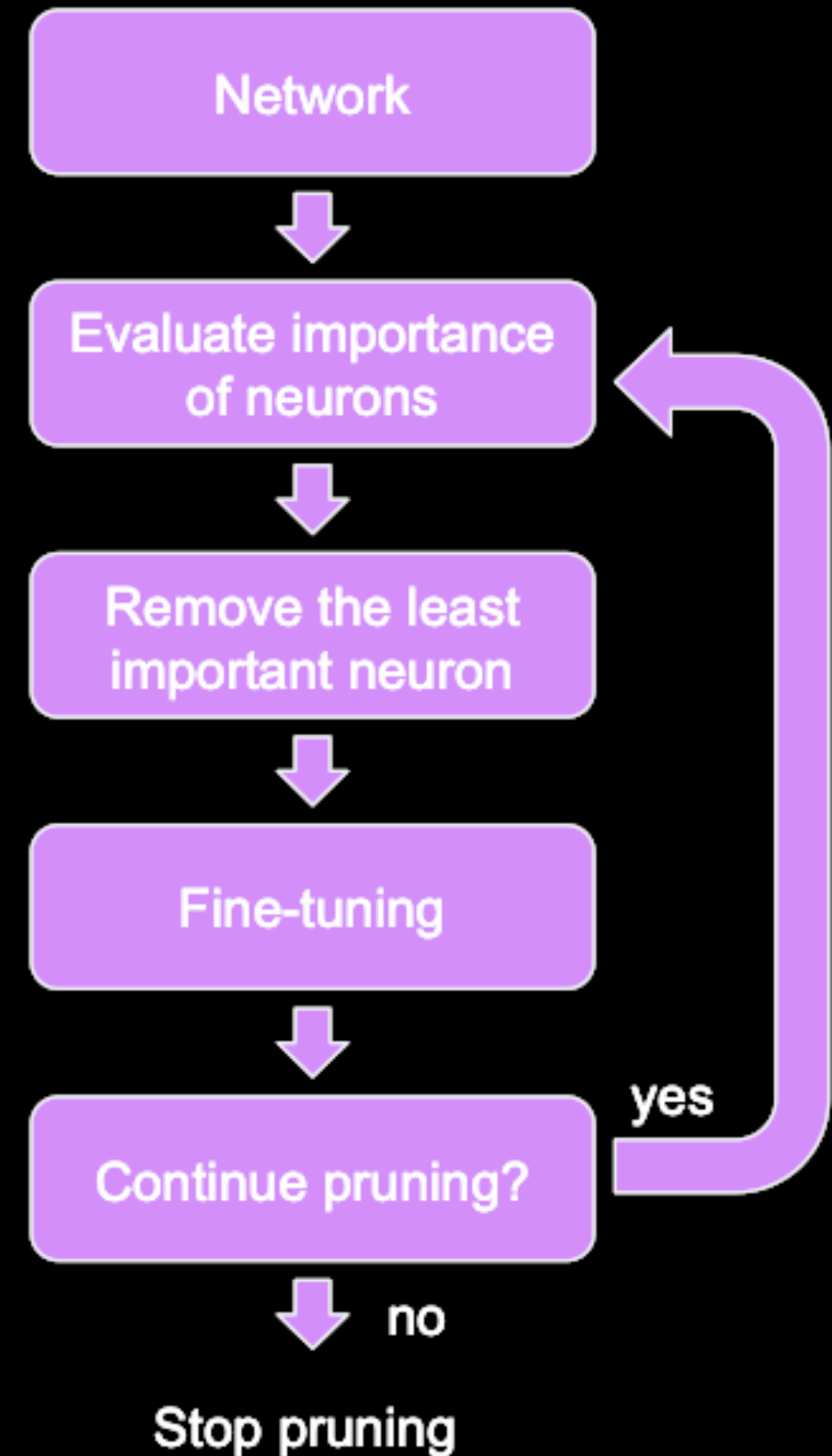
{`pmolchanov, styree, tkarras, taila, jkautz`}@nvidia.com





# Pruning based approach to analyzing/compressing

- \* Removing connections of a learned neural network
  - Analyzing the effect of this intervention on the output of the model.
    - ✓ Example of intervention based causal model analysis.
- \* Pruning is interleaved with fine-tuning





# Pruning based analysis of neural networks

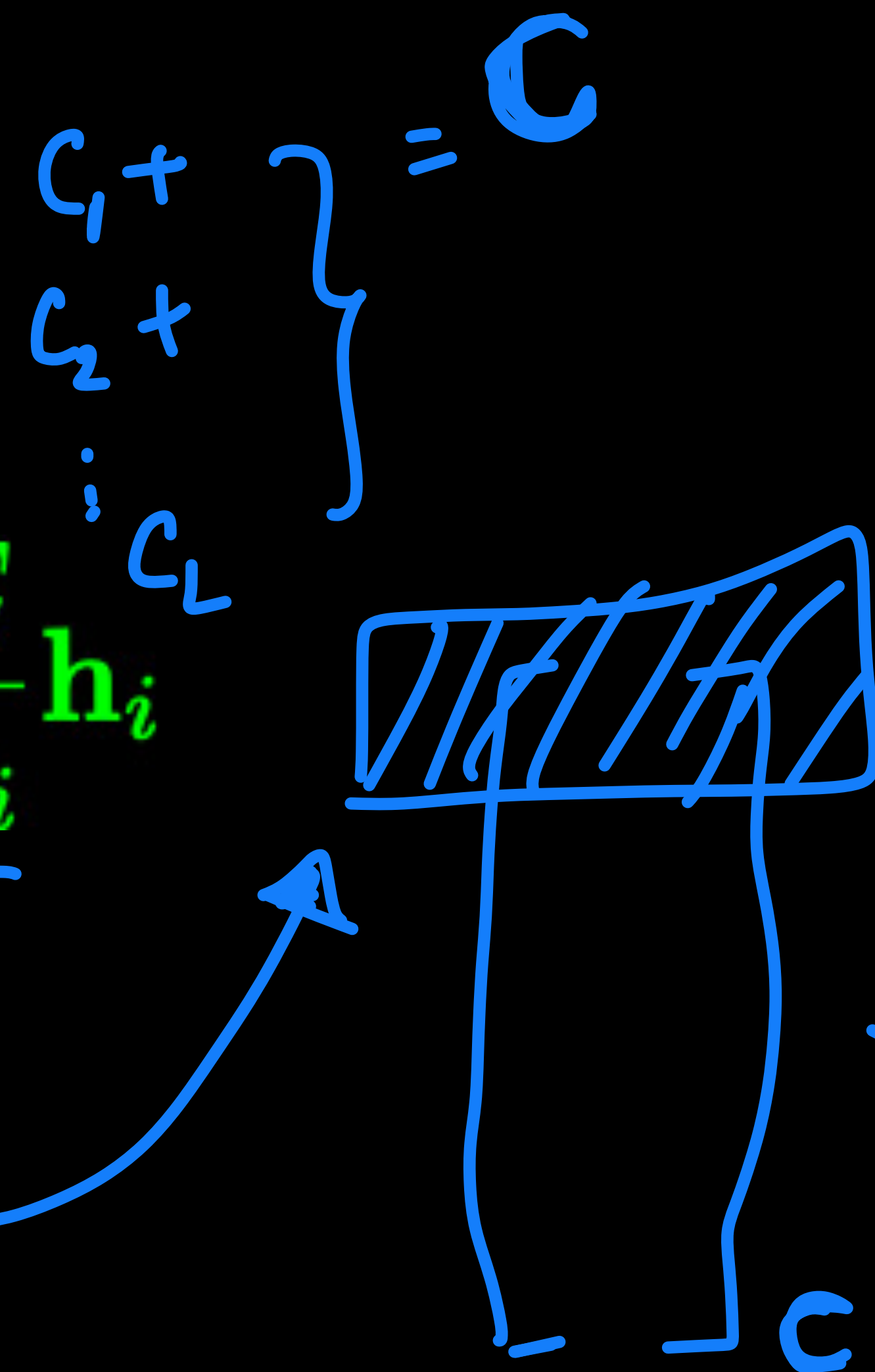
$$\mathbf{h} = \{ \mathbf{z}_1^{(1)}, \mathbf{z}_2^{(1)}, \dots, \mathbf{z}_{C_L}^{(L)} \}$$

\* Taylor series expansion based

$$E(\mathcal{D}|\mathbf{h}_i) \approx E(\mathcal{D}|\mathbf{h}_i = 0) + \frac{\partial E}{\partial \mathbf{h}_i} \mathbf{h}_i$$

\* Criterion for pruning feature maps

$$\left| \frac{\partial E^T}{\partial \mathbf{h}_i} \mathbf{h}_i \right|$$





# Criterion is involved in identifying importance

- \* With non-linearities like ReLU (deriving gradients can have effects of saturation)

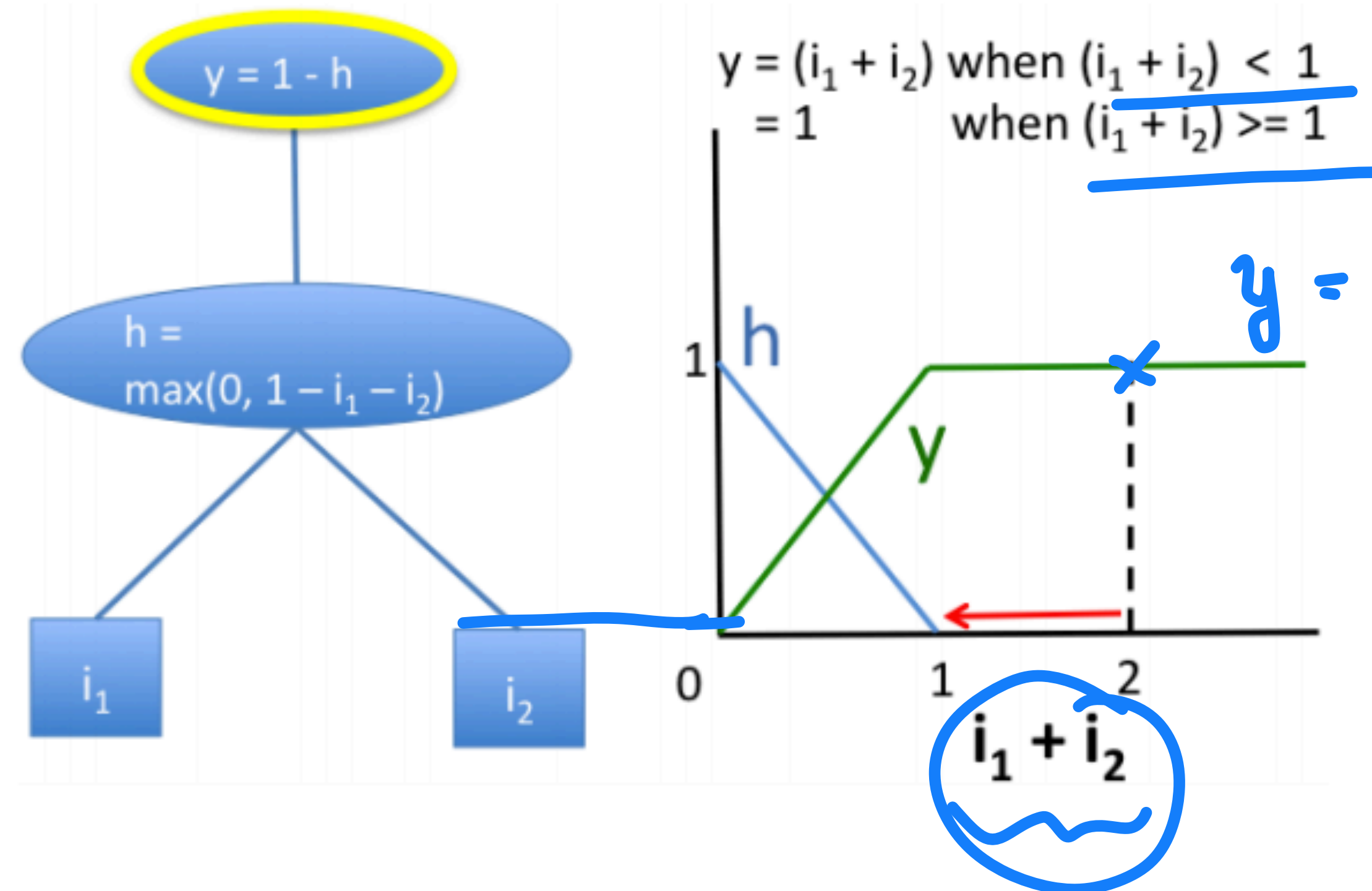


Figure 1. Perturbation-based approaches and gradient-based approaches fail to model saturation. Illustrated is a simple network exhibiting saturation in the signal from its inputs. At the point where  $i_1 = 1$  and  $i_2 = 1$ , perturbing either  $i_1$  or  $i_2$  to 0 will not produce a change in the output. Note that the gradient of the output w.r.t the inputs is also zero when  $i_1 + i_2 > 1$ .



# Criterion involved in identifying importance

- \* With non-linearities like ReLU (deriving gradients can have effects of saturation)

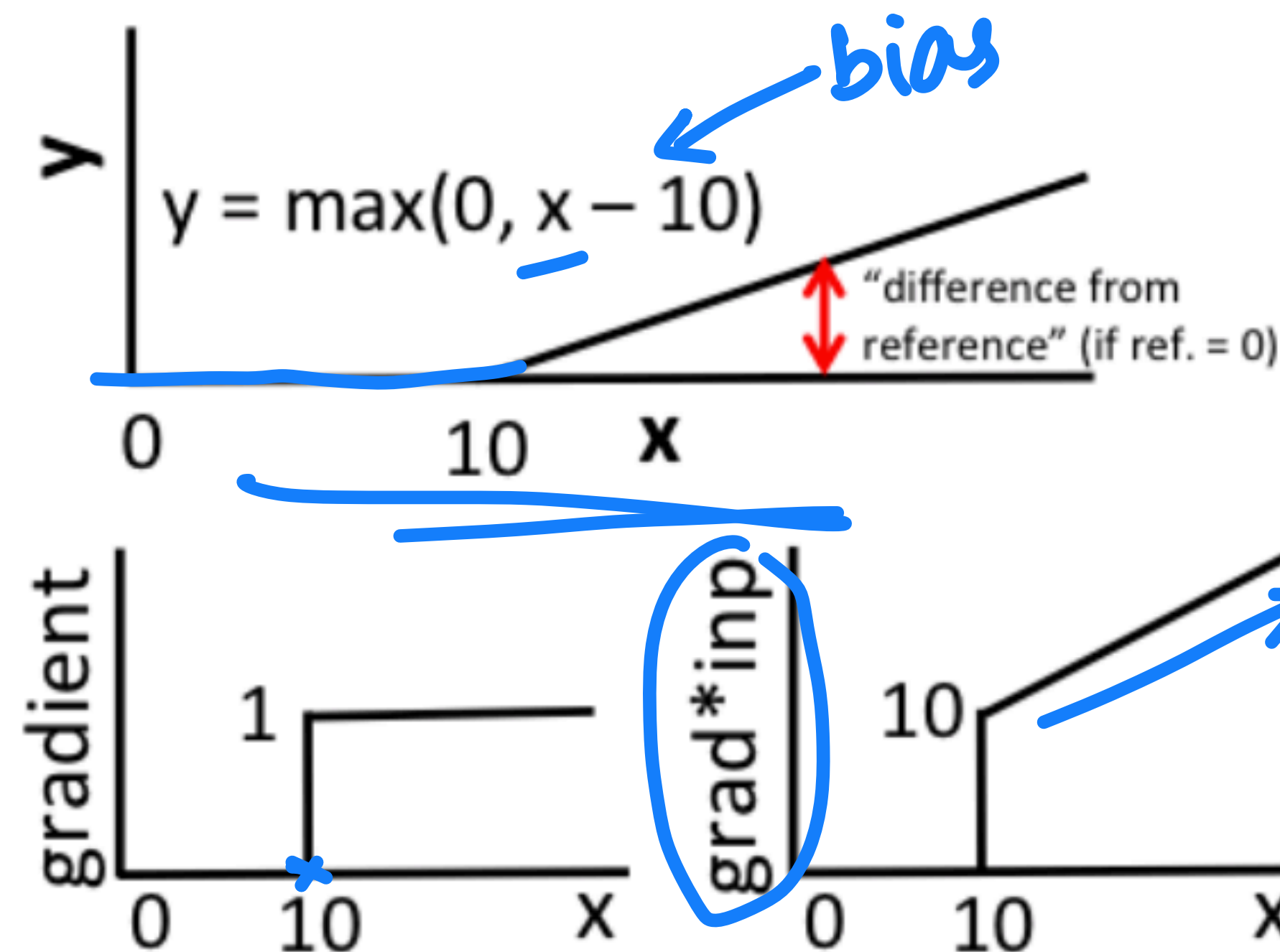


Figure 2. **Discontinuous gradients can produce misleading importance scores.** Response of a single rectified linear unit with a bias of  $-10$ . Both gradient and gradient  $\times$  input have a discontinuity at  $x = 10$ ; at  $x = 10 + \epsilon$ , gradient  $\times$  input assigns a contribution of  $10 + \epsilon$  to  $x$  and  $-10$  to the bias term ( $\epsilon$  is a small positive number). When  $x < 10$ , contributions on  $x$  and the bias term are both 0. By contrast, the difference-from-reference (red arrow, top figure) gives a continuous increase in the contribution score.



# Approximating gradients

- \* Using finite differences w.r.t. reference.

---

**Learning Important Features Through Propagating Activation Differences**

---

Avanti Shrikumar<sup>1</sup> Peyton Greenside<sup>1</sup> Anshul Kundaje<sup>1</sup>





# Approximating gradients

\* Let input neurons be defined as (for a given input)

$$\{x_1, x_2, \dots, x_D\}$$

images like MNIST

\* let reference input be defined as

$$\{x_1^0, x_2^0, \dots, x_D^0\}$$

References background of 0's for all pixels.

\* Let  $\underline{y_c}, \underline{y_c^0}$  denote the output of the model (at some dimension) for the given input and the reference. *input*





# Approximating gradients

- \* Use finite differences for gradients

$$m_{\Delta y_c, \Delta x_j} = \frac{\Delta y_c}{\Delta x_j}$$

$$\Delta y_c = y_c - y_c^0$$
$$\Delta x_j = x_j - x_j^0$$

- ✓ Called as multipliers
- ✓ Chain rule of partial derivatives can also be extended to multipliers.
- ✓ Can be used instead of actual gradients to compute the importance of a feature/hidden layer output.

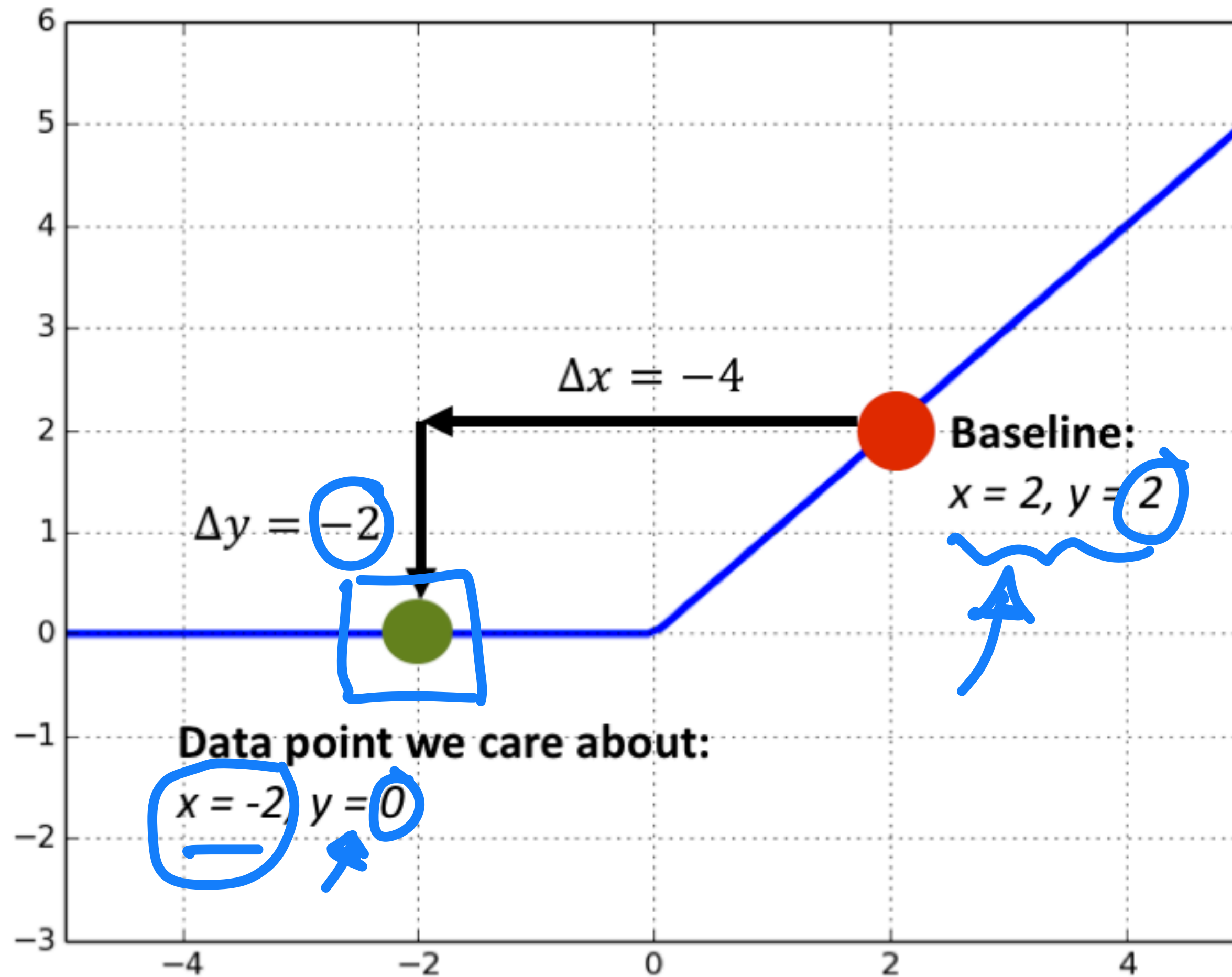
Interesting math





# Illustration of finite gradients for ReLU

$$y = \text{ReLU}(x) = \max(0, x)$$



$\begin{cases} x_i & \text{— internal neuron} \\ x_i^0 & \text{— neuron (ref)} \end{cases}$

1. Calculating the slope

$$\frac{\Delta y}{\Delta x} = \frac{-2}{-4} = 0.5 \quad \text{[approx]}$$

2. Finding the feature importance

$$x_i \times \frac{\Delta y}{\Delta x} = -4 \times 0.5 = -2$$

$$x_i \frac{\partial y_i}{\partial x_i} \quad \leftarrow \text{approx.}$$





# Adversarial Examples and Explainability





# Adversarial examples and learning

## Adversarial Examples: Attacks and Defenses for Deep Learning

Xiaoyong Yuan, Pan He, Qile Zhu, Xiaolin Li\*

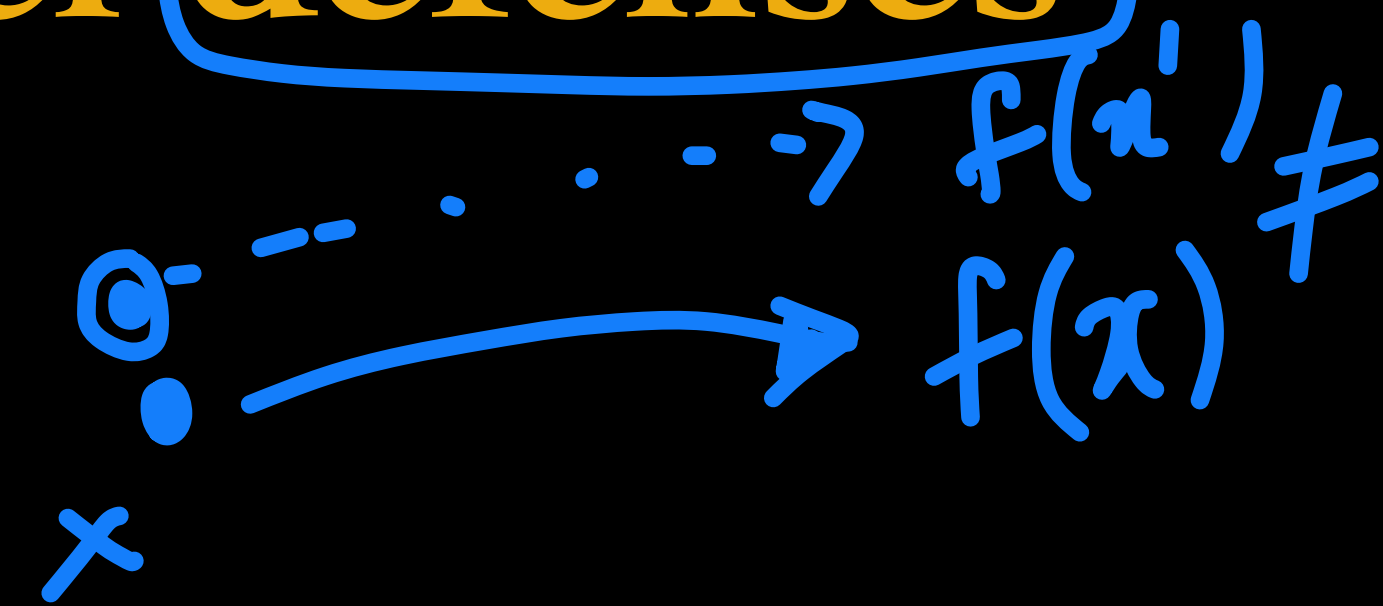
National Science Foundation Center for Big Learning, University of Florida

{chbrian, pan.he, valder}@ufl.edu, andyli@ece.ufl.edu





# Adversarial attacks and model defenses



## \* Examples that fool the model

→ Using a trained image classifier published by a third party, a user inputs one image to get the prediction of class label. Adversarial images are original clean images with small perturbations, often barely recognizable by humans. However, such perturbations misguide the image classifier

$$\begin{aligned} \min_{x'} \quad & \|x' - x\| \\ \text{s.t.} \quad & f(x') = l' \\ & f(x) = l \\ & l \neq l', \end{aligned}$$

$f$ : classifier





# Types of Adversarial attacks

- \* **False positive attacks** generate a negative sample which is misclassified as a positive one (Type I Error). In a malware detection task, a benign software being classified as malware is a false positive. In an image classification task, a false positive can be an adversarial image unrecognizable to human, while deep neural networks predict it to a class with a high confidence score.
- \* **False negative attacks** generate a positive sample which is misclassified as a negative one (Type II Error). In a malware detection task, a false negative can be the condition that a malware (usually considered as positive) cannot be identified by the trained model. False negative attack is also called machine learning evasion. This error is shown in most adversarial images, where human can recognize the image, but the neural networks cannot identify it.





# Types of Adversarial attacks

\* White box versus Black box

↑ easier

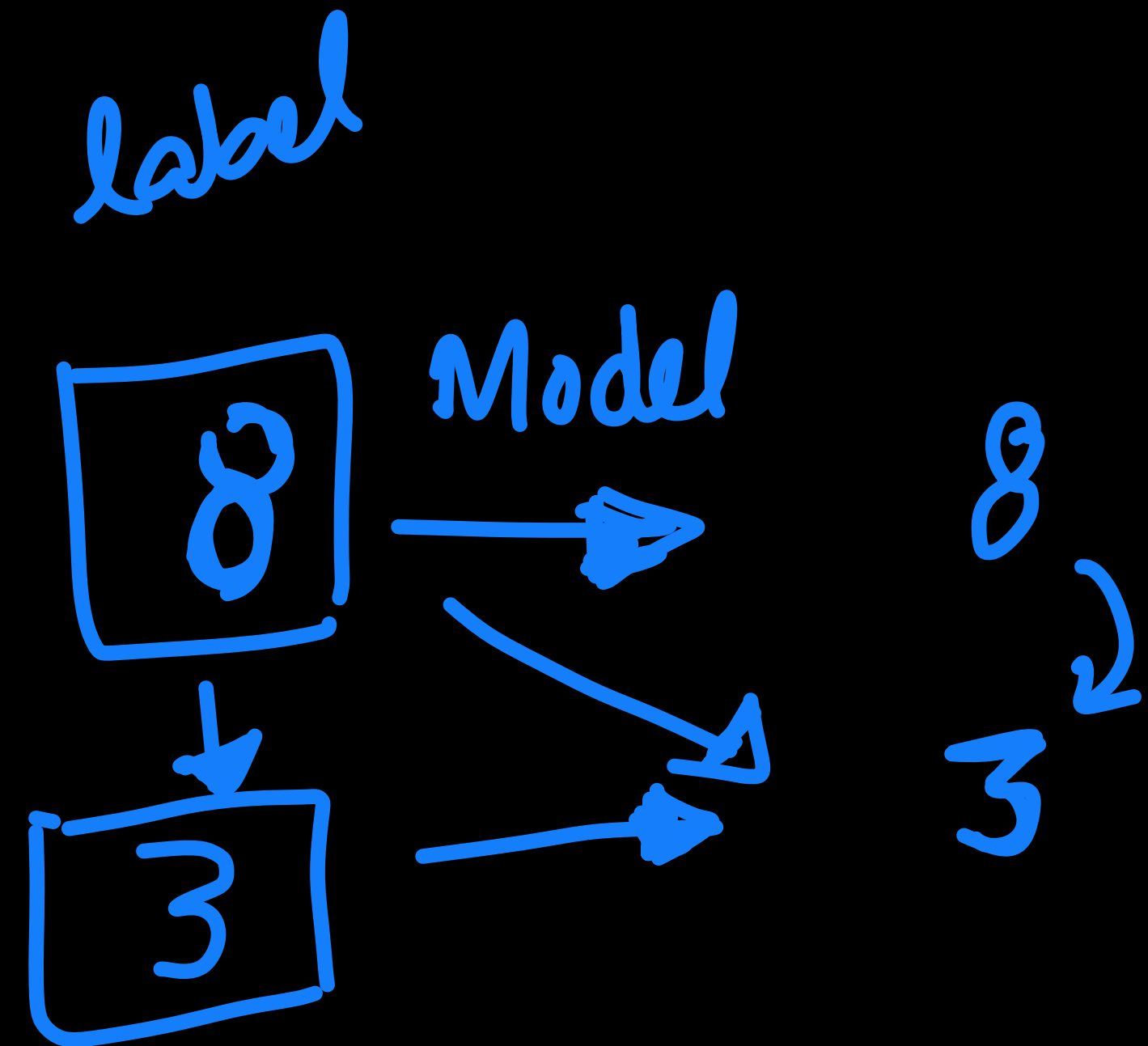
for more consequential

\* Targeted versus Non-targeted

↑ make it all to one class

\* One-time versus many time

do you one-shot at fooling the system





# Simple adversarial attack

## \* Fast Gradient Sign Method

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i + \epsilon \text{sign}\left(\frac{\partial E(\mathbf{x}_i, l)}{\partial \mathbf{x}_i}\right)$$

non-targeted  
adversarial

\* Move in the direction of the gradient ascent *on inputs*

\* Other similar rules - gradient value based adversarial learning

$$\mathbf{x}_i = \mathbf{x}_i + \epsilon \left(\frac{\partial E(\mathbf{x}_i, l)}{\partial \mathbf{x}_i}\right)$$



# Simple adversarial attack

## \* Fast Gradient Sign Method

$x$   
“panda” /  $x = 0.57$   
57.7% confidence

$+ 0.007 \times$

$\text{sign}(\nabla_x J(\theta, x, y))$   
“nematode”  
8.2% confidence

$=$

$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$   
“gibbon”  
99.3% confidence



# Adversarial examples

## Adversarial Examples for Evaluating Reading Comprehension Systems

**Robin Jia**

Computer Science Department

Stanford University

`robinjia@cs.stanford.edu`

**Percy Liang**

Computer Science Department

Stanford University

`pliang@cs.stanford.edu`

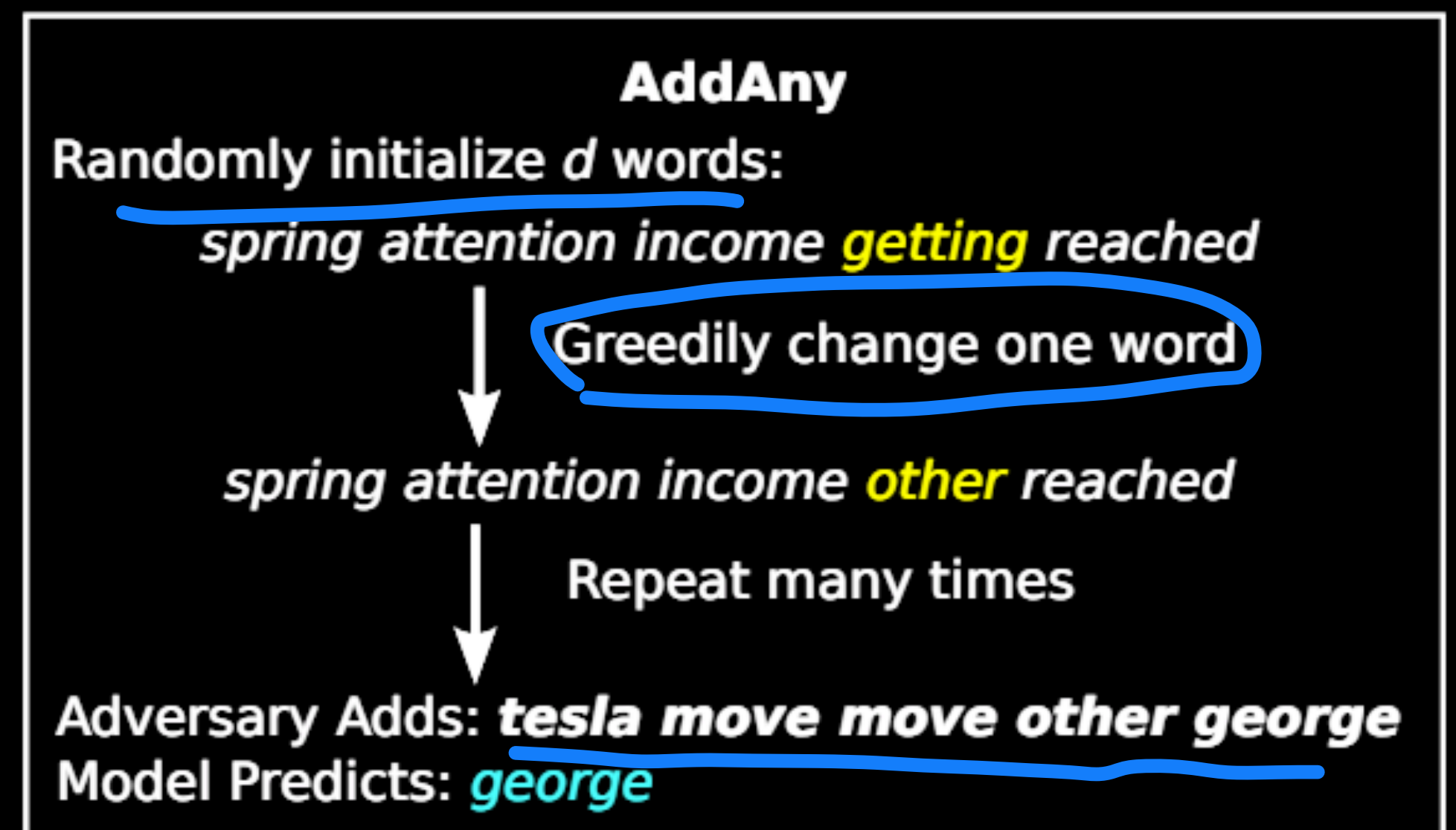


# Adversarial examples in text

- ✳ Using similar methods to gradient based update.
- ✳ Adding sentences confuses models which will typically not confuse humans

≠ Locally very smooth

Article: **Nikola Tesla**  
Paragraph: "In January 1880, two of Tesla's uncles put together enough money to help him leave Gospić for **Prague** where he was to study. Unfortunately, he arrived too late to enroll at Charles-Ferdinand University; he never studied Greek, a required subject; and he was illiterate in Czech, another required subject. Tesla did, however, attend lectures at the university, although, as an auditor, he did not receive grades for the courses."  
Question: "What city did Tesla move to in 1880?"  
Answer: **Prague**  
Model Predicts: Prague





# Adversarial attacks

## \* Understanding adversarial attacks

- ✓ Allows explainability
- ✓ Build defenses to these attacks

Adversarial example

Defense (data augmentation)

regularization  
(locally smooth)

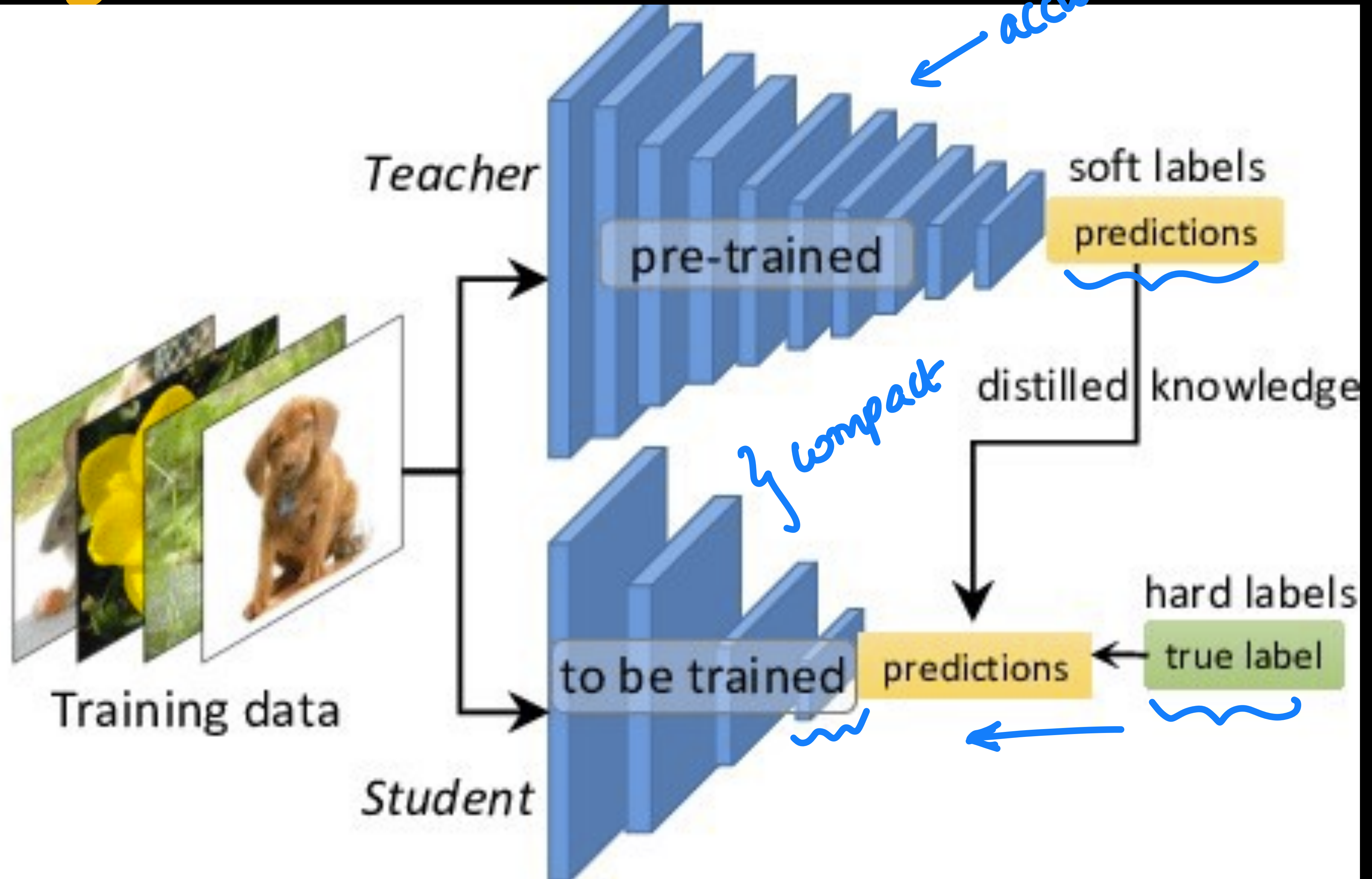


# Explainability with distillation





# Knowledge distillation



# Knowledge distillation

- \* Teacher models are complex large neural networks

  - Student models are typically lighter models.

- \* Useful in semi-supervised learning

  - Student model has to approximate outputs from a teacher model.

    - ✓ Also needs to learn from small amounts of labelled data.





# Knowledge distillation for explainability

- \* Use a simpler explainable model for student model to approximate the deeper model

## “Why Should I Trust You?” Explaining the Predictions of Any Classifier

Marco Tulio Ribeiro  
University of Washington  
Seattle, WA 98105, USA  
marcotcr@cs.uw.edu

Sameer Singh  
University of Washington  
Seattle, WA 98105, USA  
sameer@cs.uw.edu

Carlos Guestrin  
University of Washington  
Seattle, WA 98105, USA  
guestrin@cs.uw.edu



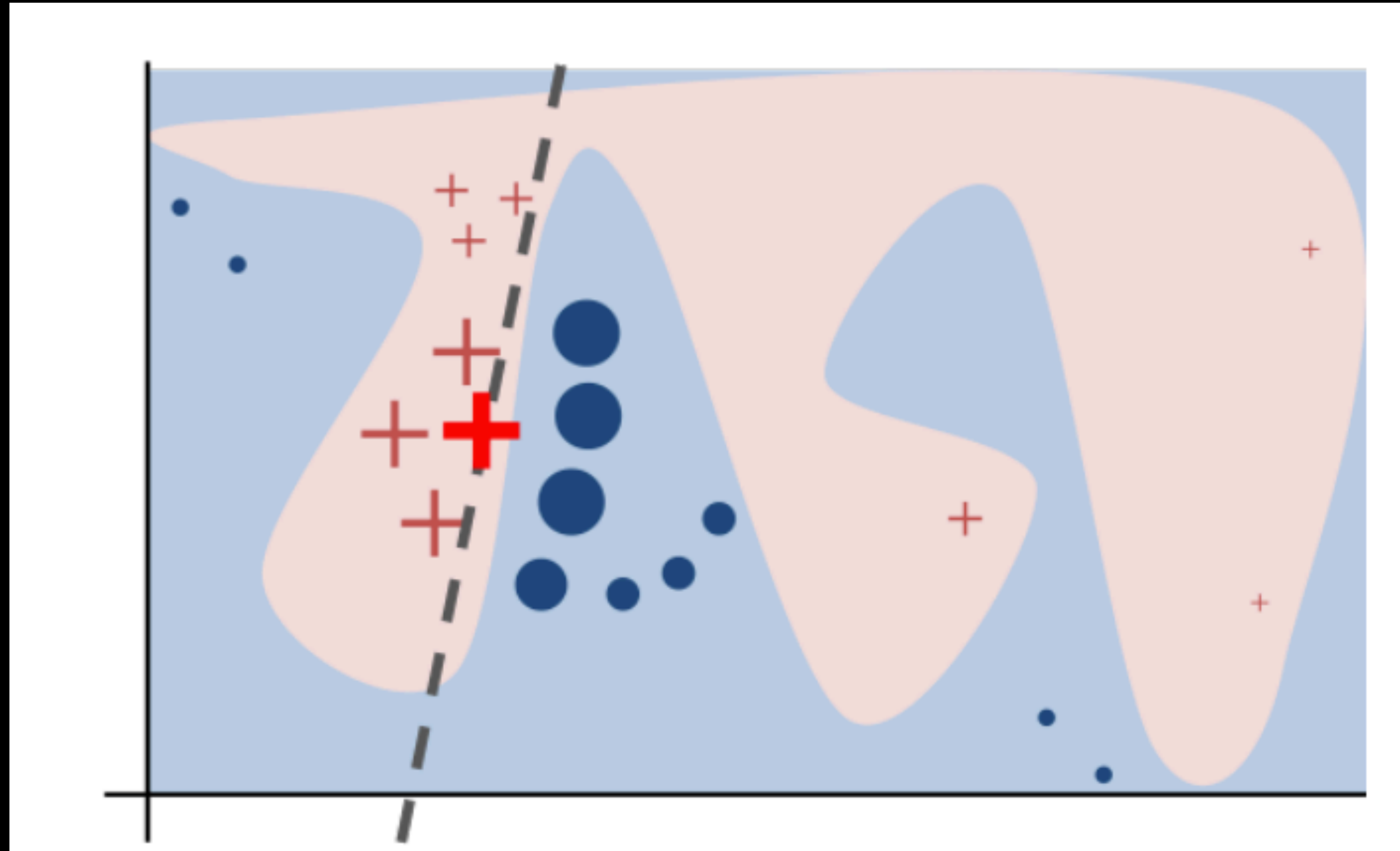
# Knowledge distillation for explainability

- ✳ Use a simpler explainable model for student model to approximate the deeper model.
- ✳ Use locality preservation as a criterion for sampling
  - ✓ Method - Local Interpretable Model Agnostic Representations
  - ✓ Explainability for each sample under consideration





# Knowledge distillation for explainability



# Knowledge distillation for explainability

- \* Let  $f(\mathbf{x})$  denote the original neural network
- \* Let  $\mathbf{x}'$  denote the interpretable version of input
- \* Let  $\mathbf{z}, \mathbf{z}'$  denote samples drawn around input and its interpretable version.

$$\mathcal{L}(f, g, \pi_x) = \sum_{z, z' \in \mathcal{Z}} \pi_x(z) (f(z) - g(z'))^2$$

- \* The function  $g$  can be sparse linear regression





# LIME model - text example

\* Building sparse linear regression for each output class

Example #3 of 6 True Class: ● Atheism Instructions Previous Next

Algorithm 1	Algorithm 2																								
<p><b>Words that A1 considers important:</b></p> <table style="width: 100%;"><tr><td>GOD</td><td style="background-color: red; width: 80%;"></td></tr><tr><td>mean</td><td style="background-color: red; width: 70%;"></td></tr><tr><td>anyone</td><td style="background-color: green; width: 60%;"></td></tr><tr><td>this</td><td style="background-color: green; width: 55%;"></td></tr><tr><td>Koresh</td><td style="background-color: red; width: 30%;"></td></tr><tr><td>through</td><td style="background-color: green; width: 20%;"></td></tr></table>	GOD		mean		anyone		this		Koresh		through		<p><b>Words that A2 considers important:</b></p> <table style="width: 100%;"><tr><td>Posting</td><td style="background-color: red; width: 80%;"></td></tr><tr><td>Host</td><td style="background-color: red; width: 80%;"></td></tr><tr><td>Re</td><td style="background-color: red; width: 50%;"></td></tr><tr><td>by</td><td style="background-color: green; width: 50%;"></td></tr><tr><td>in</td><td style="background-color: green; width: 50%;"></td></tr><tr><td>Nntp</td><td style="background-color: red; width: 20%;"></td></tr></table>	Posting		Host		Re		by		in		Nntp	
GOD																									
mean																									
anyone																									
this																									
Koresh																									
through																									
Posting																									
Host																									
Re																									
by																									
in																									
Nntp																									
<p><b>Predicted:</b></p> <p><span style="color: red;">●</span> Atheism</p> <p><b>Prediction correct:</b></p> <p style="font-size: 2em; color: green;">✓</p>	<p><b>Predicted:</b></p> <p><span style="color: red;">●</span> Atheism</p> <p><b>Prediction correct:</b></p> <p style="font-size: 2em; color: green;">✓</p>																								
<p><b>Document</b></p> <p>From: pauld@verdix.com (Paul Durbin) Subject: Re: DAVID CORESH IS! <b>GOD!</b> Nntp-Posting-Host: sarge.hq.verdix.com Organization: Verdix Corp Lines: 8</p>	<p><b>Document</b></p> <p>From: pauld@verdix.com (Paul Durbin) Subject: <b>Re:</b> DAVID CORESH IS! GOD! <b>Nntp-Posting-Host:</b> sarge.hq.verdix.com Organization: Verdix Corp Lines: 8</p>																								

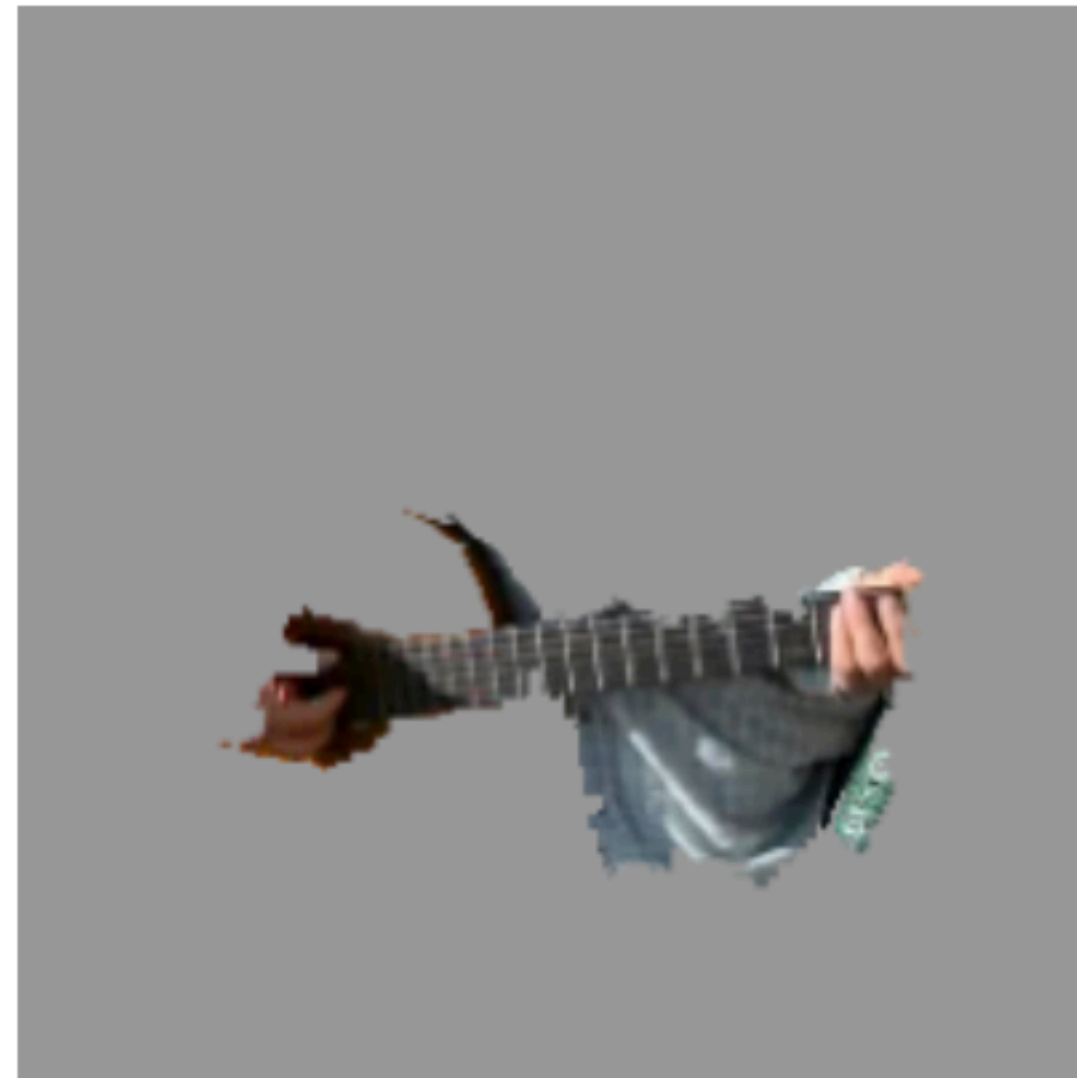


# LIME model - Image example

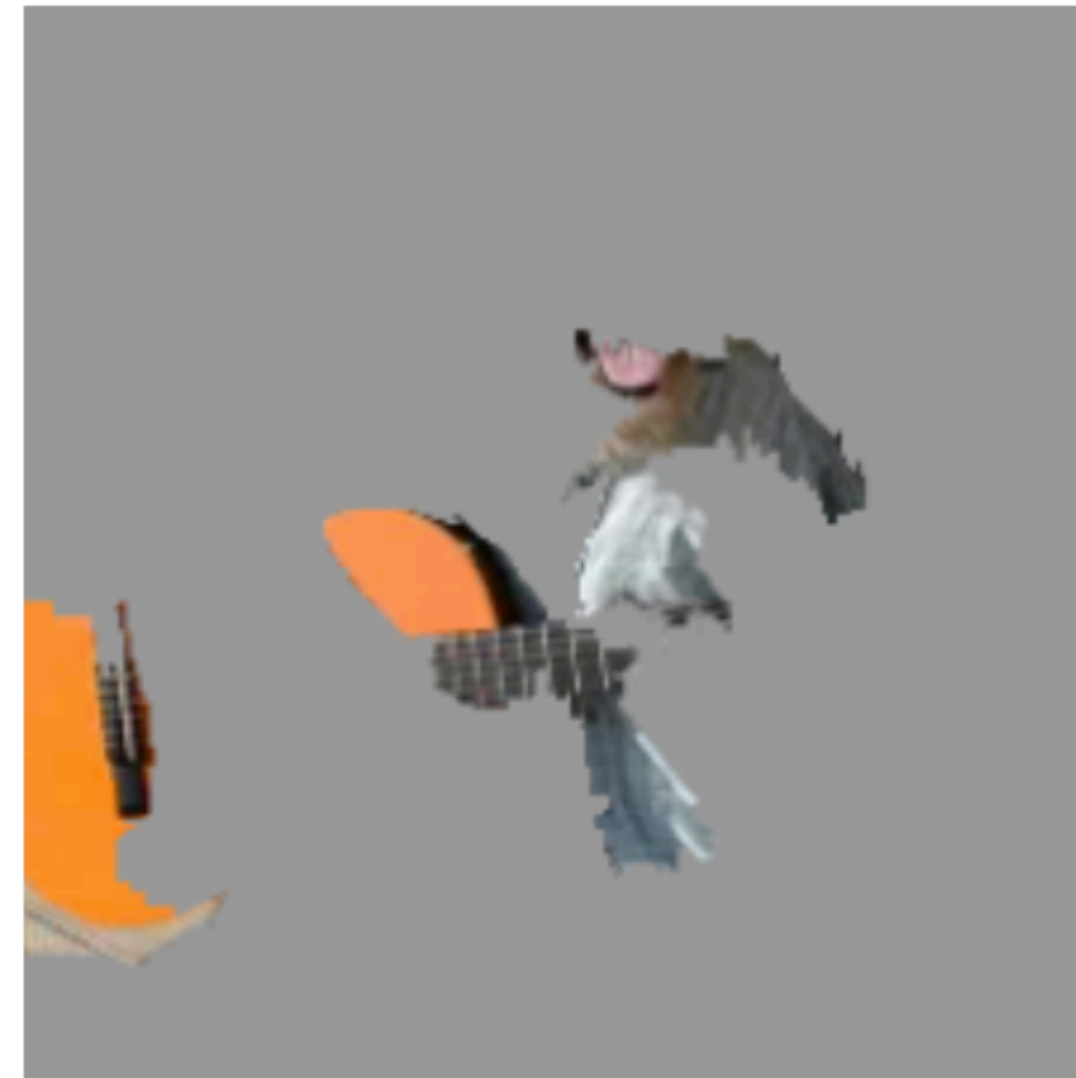
- \* Building sparse linear regression for each output class



(a) Original Image



(b) Explaining *Electric guitar*



(c) Explaining *Acoustic guitar*



(d) Explaining *Labrador*