

# Explainability Evaluation and Global Explainability

CPSC680: Trustworthy Deep Learning

Rex Ying

# Readings

- Readings are updated on the website (syllabus page)
- **Lecture 6 readings:**
  - [1710.10547.pdf \(arxiv.org\)](#) Explanations can also be vulnerable to adversarial attacks
  - [2005.00631.pdf \(arxiv.org\)](#) Evaluating Explanations

# Content

- Evaluating Explainability Methods
- Model-level Explainability
- Intrinsic Explainability / Interpretability

# Content

- Evaluating Explainability Methods
- Model-level Explainability
- Intrinsic Explainability / Interpretability

# Criteria of Good Explanation

- Fidelity

The explanation maximally supports model's prediction

- Sensitivity

The explanation is stable for similar model input and output

- Conciseness

The explanation cannot be too large (Occam's razor)

- Interpretability

The explanation can be easily understood by human

***Dico ergo** ad qñem q  
qz pluralitas  
non est ponenda sine necessitate ⁊ non  
ē necessitas quare debeat poni tñus oī  
secretum mensurās motum angeli. naz*

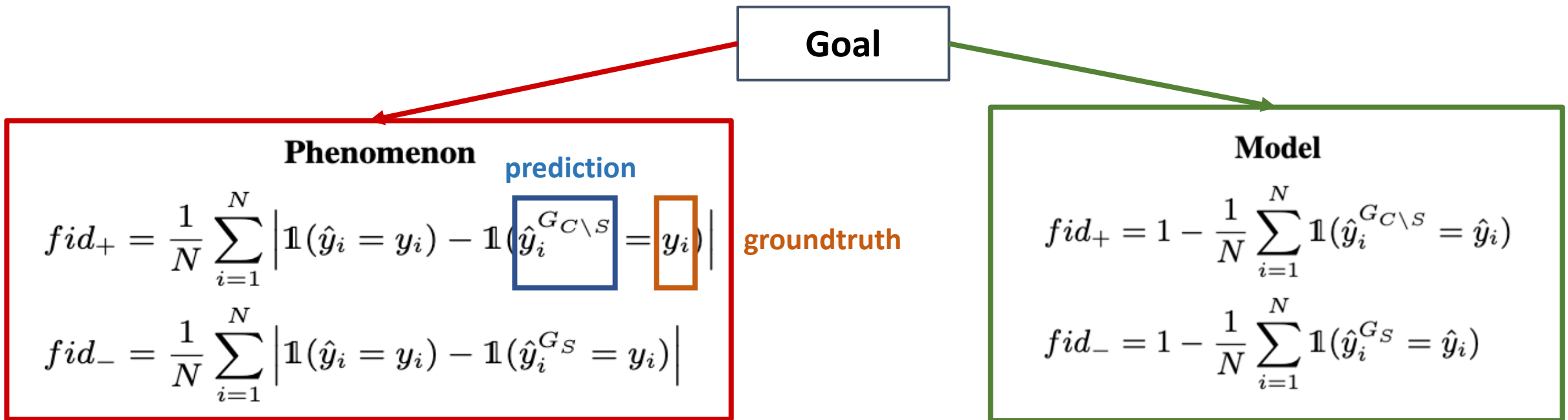
“pluarity is not to be  
posited without necessity”

# Explanation Goal

- **Phenomenon** Explanation
  - Explain the underlying reasons for the ground truth phenomenon
- **Model** Explanation
  - Explain why model makes a particular prediction
- We will explain the **fidelity** metric in both cases:

# Explanation Goal: Fidelity Metric

- Define 2 fidelity metrics:  $fid_+$  and  $fid_-$  to capture different aspects of **explanation quality**
- The formula of fidelity depends on the goal:
  - **Goal 1**: explain **phenomenon** of the data
  - **Goal 2**: explain what has the **model** learned



# Fidelity Metric Details

- **Characteristics of a good explanation**
- $fid_+$ : removal important subgraph will result in dramatic decrease of the confidence
- $fid_-$ : Using only the important subgraph will result in similar confidence

## Phenomenon

$$fid_+ = \frac{1}{N} \sum_{i=1}^N \left| \mathbb{1}(\hat{y}_i = y_i) - \mathbb{1}(\hat{y}_i^{G_{C \setminus S}} = y_i) \right|$$

Removal of important subgraph

$$fid_- = \frac{1}{N} \sum_{i=1}^N \left| \mathbb{1}(\hat{y}_i = y_i) - \mathbb{1}(\hat{y}_i^{G_S} = y_i) \right|$$

Keeping only the important subgraph

Original prediction  
probability / confidence



# Explanation Evaluation Criteria

- Notably, the explanation evaluation criteria are **multi-dimensional**
- **Explanation quality**
  - High fidelity / characterization scores
  - Sufficiency and necessity aspects (see the previous slide)
- **Explanation stability**
  - Explanations are consistent across random optimization seeds (measure variance)
- **Explanation complexity**
  - The explanation should be concise and easy to understand by human (measure size)

# Types of Explanations

- **Sufficiency**

- An explanation is sufficient if it leads by its own to the initial prediction of the model explanation. ( $fid_- \rightarrow 0$ )

- **Necessity**

- An explanation is necessary if the model prediction changes when removing it from the initial graph. ( $fid_+ \rightarrow 1$ )

- Use the **Characterization** score to summarize the explanation quality

$$character = \frac{w_+ + w_-}{\frac{w_+}{fid_+} + \frac{w_-}{1 - fid_-}} = \frac{(w_+ + w_-) \times fid_+ \times (1 - fid_-)}{w_+ \cdot (1 - fid_-) + w_- \cdot fid_+}$$

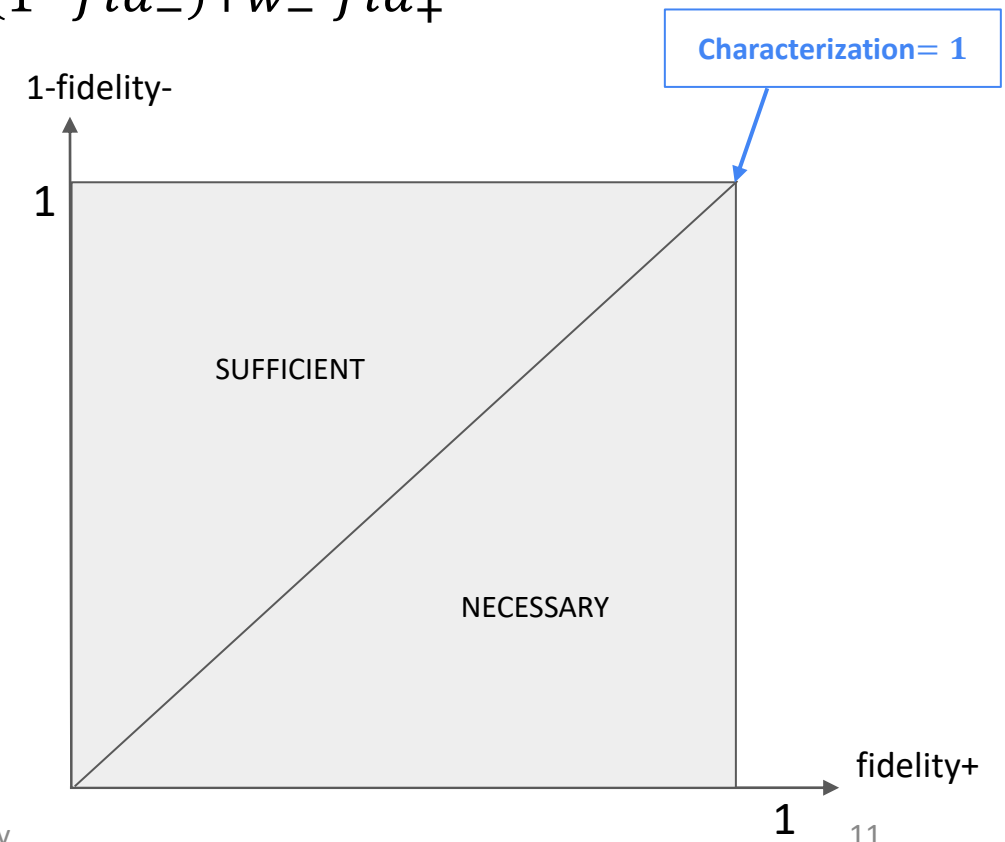
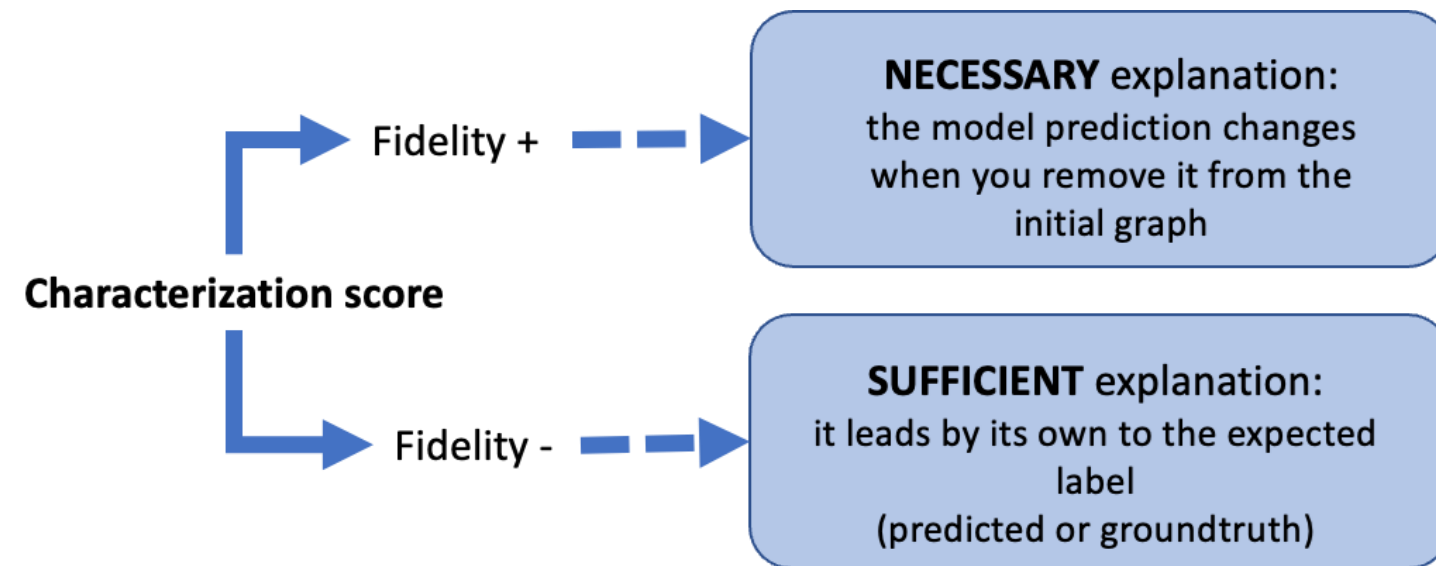
Where  $w_+$  and  $w_-$  are the weights of both fidelity metrics (commonly set  $w_+ = w_- = 1$ )

# Characterization Score

- **Characterization** score to summarize the explanation quality

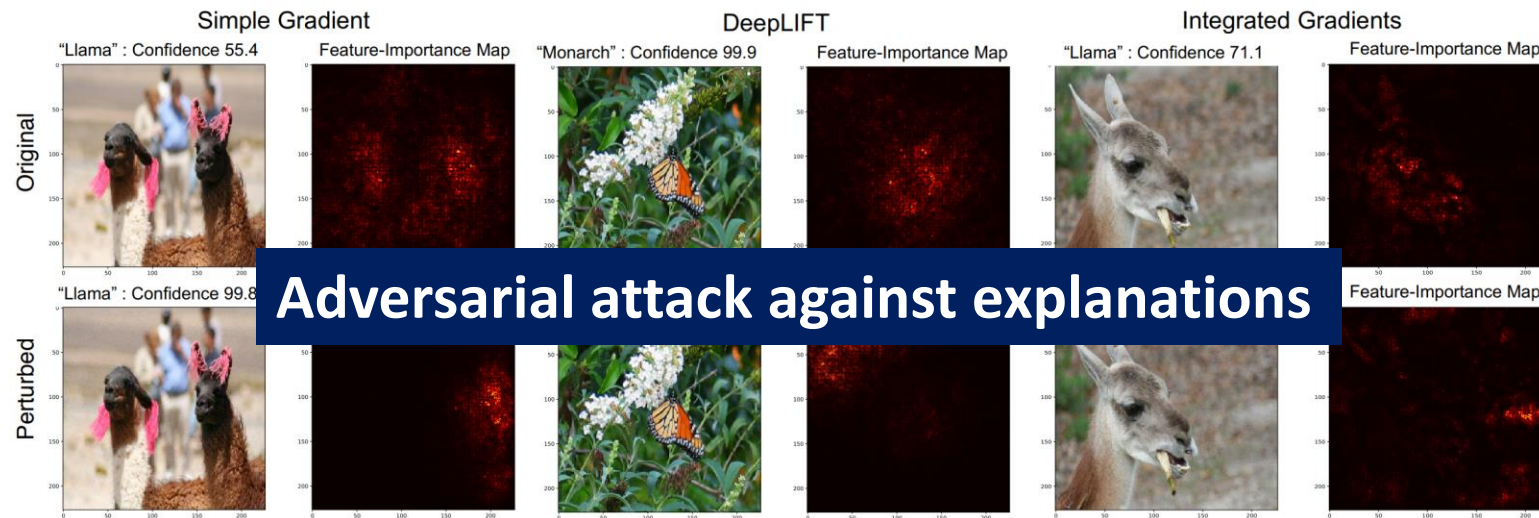
$$charact = \frac{w_+ + w_-}{\frac{w_+}{fid_+} + \frac{w_-}{1 - fid_-}} = \frac{(w_+ + w_-) \times fid_+ \times (1 - fid_-)}{w_+ \cdot (1 - fid_-) + w_- \cdot fid_+}$$

- Necessary AND sufficient



# Sensitivity Desiderata

- Similar input & output  $\rightarrow$  explanations should be similar
- Also called “stability”
- Local smoothness is not usually true for deep neural networks, but is a very common assumption in human cognition



# Sensitivity Definition

- Define the neighborhood of a point of interest  $x$

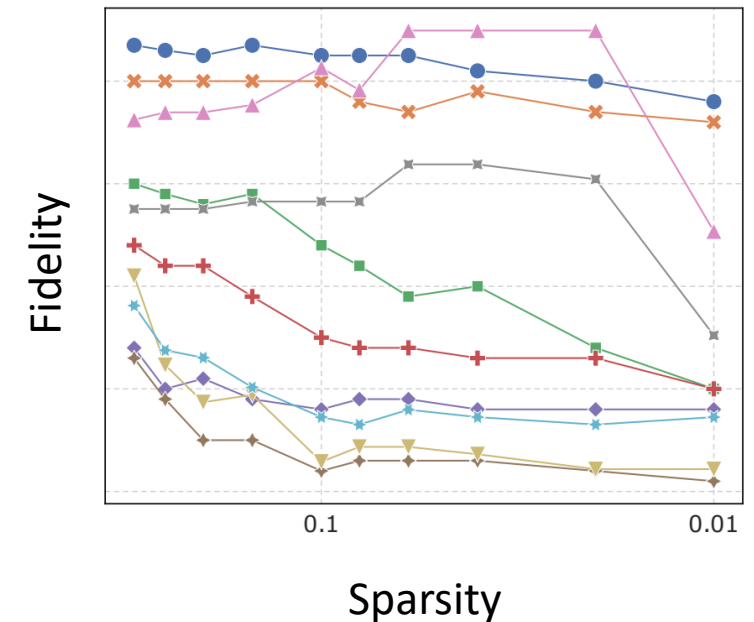
$$N_r = \{z \in D_x | \rho(x, z) \leq r, f(x) = f(z)\}$$

- The local region around prediction of  $x$  that's stable
- Max Sensitivity  $\mu_M$ 
$$\mu_M(f, g, r; x) = \max_{z \in N_r} D(g(f, x), g(f, z))$$
- Average Sensitivity

$$\mu_A(f, g, r; x) = \int_{z \in N_r} D(g(f, x), g(f, z)) dz$$

# Conciseness – Low Complexity

- **Sparsity** is important in order to ensure that explanation highlights the most relevant part of the input
- Sparsity can be measured by the **size** of the explanation
  - Often controlled in the experiments
- Sparsity can also be measured by **entropy**
  - For explanations with importance scores
  - Attribution methods, Mask-based methods etc.



Allows us to investigate the tradeoff

# Global Explainability Evaluation

- Relative performance loss

$$RPF = \frac{(\log \mathcal{L}(M_{-F}) - \log \mathcal{L}(M))}{\log \mathcal{L}(M)}$$

- $\log \mathcal{L}(M)$  : loss function value on all test data
- $\mathcal{L}(M_{-F})$  : loss function value after feature pruning (on all test data)
- Analogous to instance-level fidelity

# Human Evaluation



(a) Raw input image. Note that this is not a part of the tasks (b) and (c)

What do you see?



Your options:

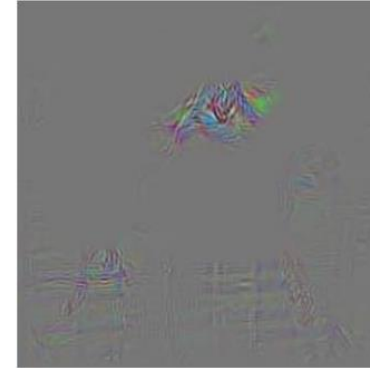
- ☐ Horse
- ☐ Person

(b) AMT interface for evaluating the class-discriminative property

Both robots predicted: Person

Robot A based its decision on

Robot B based its decision on



Which robot is more reasonable?

- ☐ Robot A seems clearly more reasonable than robot B
- ☐ Robot A seems slightly more reasonable than robot B
- ☐ Both robots seem equally reasonable
- ☐ Robot B seems slightly more reasonable than robot A
- ☐ Robot B seems clearly more reasonable than robot A

(c) AMT interface for evaluating if our visualizations instill trust in an end user

Utilizes human evaluation platform such as Amazon Mechanical Turk (AMT)



# Content

- Evaluating Explainability Methods
- **Global-level Explainability**
- Intrinsic Explainability / Interpretability

# Model-level Explanation

- **Model-level explanations** aim to shed light on a model's overall decision-making process on *a set of inputs*, instead of a specific instance.
- provides a **bird-eye-view** of the model behavior, analyzing potential bias affecting a group/subgroup of instances.
- Examples:
  - Concept-based explanations: provide importance measurement for high-level concepts, instead of individual features or pixels.
  - Influence functions: measure the impact of each data point in the training set on the model's predictions

# Global Explanation via Model Distillation

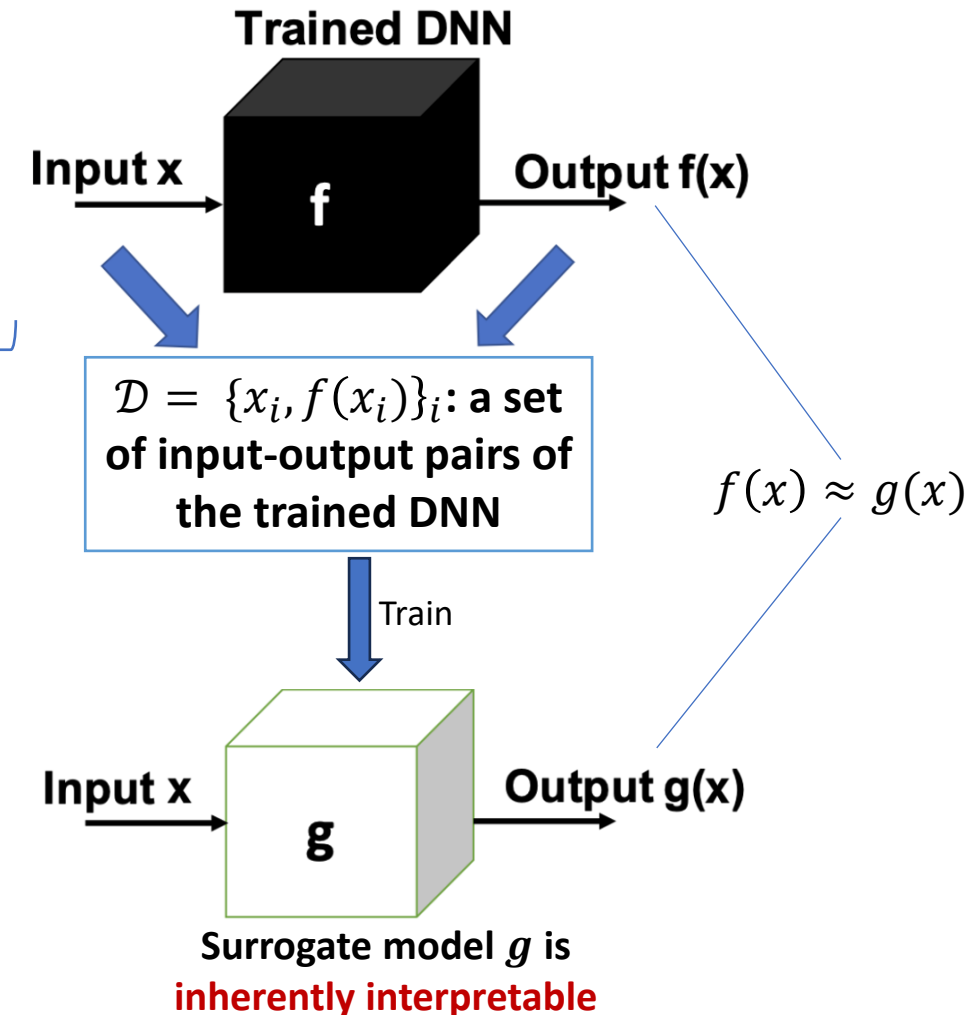
- Generalized Additive Model (GAM)

$$g(x) = h_0 + \underbrace{\sum_i h_i(x_i)}_{\text{Shape functions of individual features}} + \underbrace{\sum_{i \neq j} h_{ij}(x_i, x_j) + \sum_{i \neq j} \sum_{i \neq k} h_{ijk}(x_i, x_j, x_k) + \dots}_{\text{Higher-order feature interaction terms}}$$

Shape functions of individual features

Higher-order feature interaction terms

What are the potential issues?



# Concept Definition

- **Concept**: high-level units that are more understandable to human than individual features, pixels, etc.
- For example, **the wheel** and the **police logo** are important concepts for police vans.



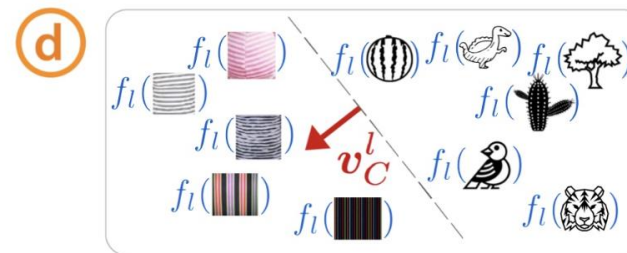
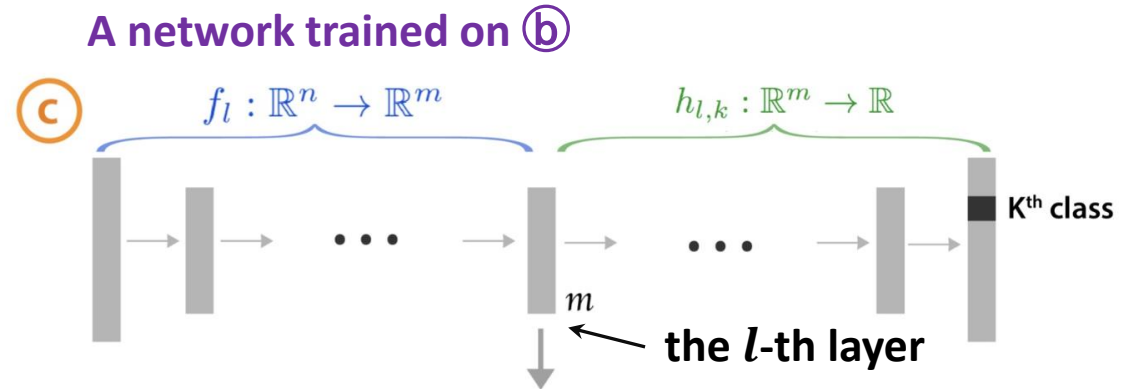
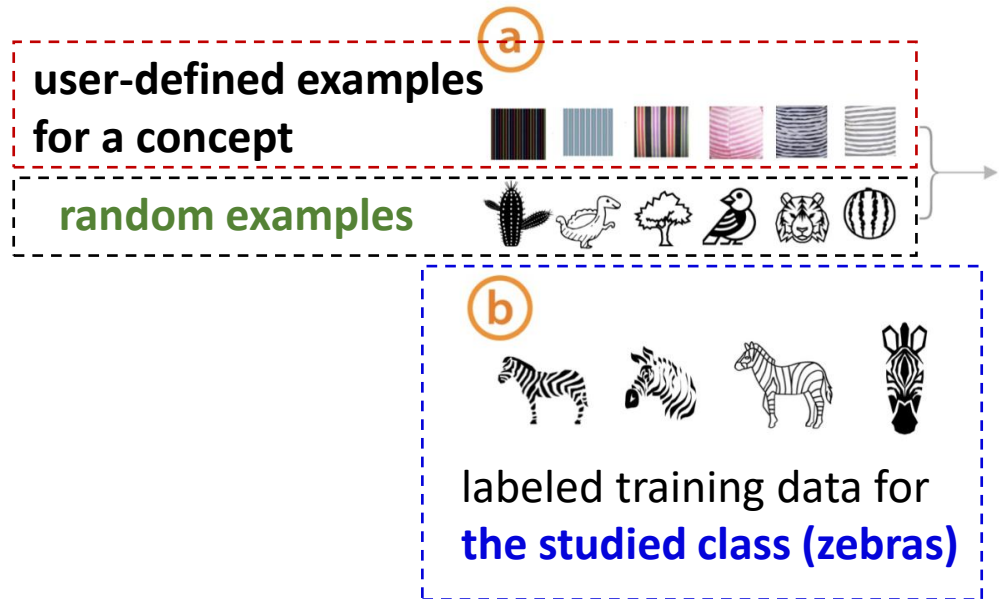
concept 1: **wheel**



concept 2: **police logo**

# TCAV Pipeline

- Testing with Concept Activation Vectors (CAV) [[paper](#)]



train a linear classifier in the activation space of the  $l$ -th layer

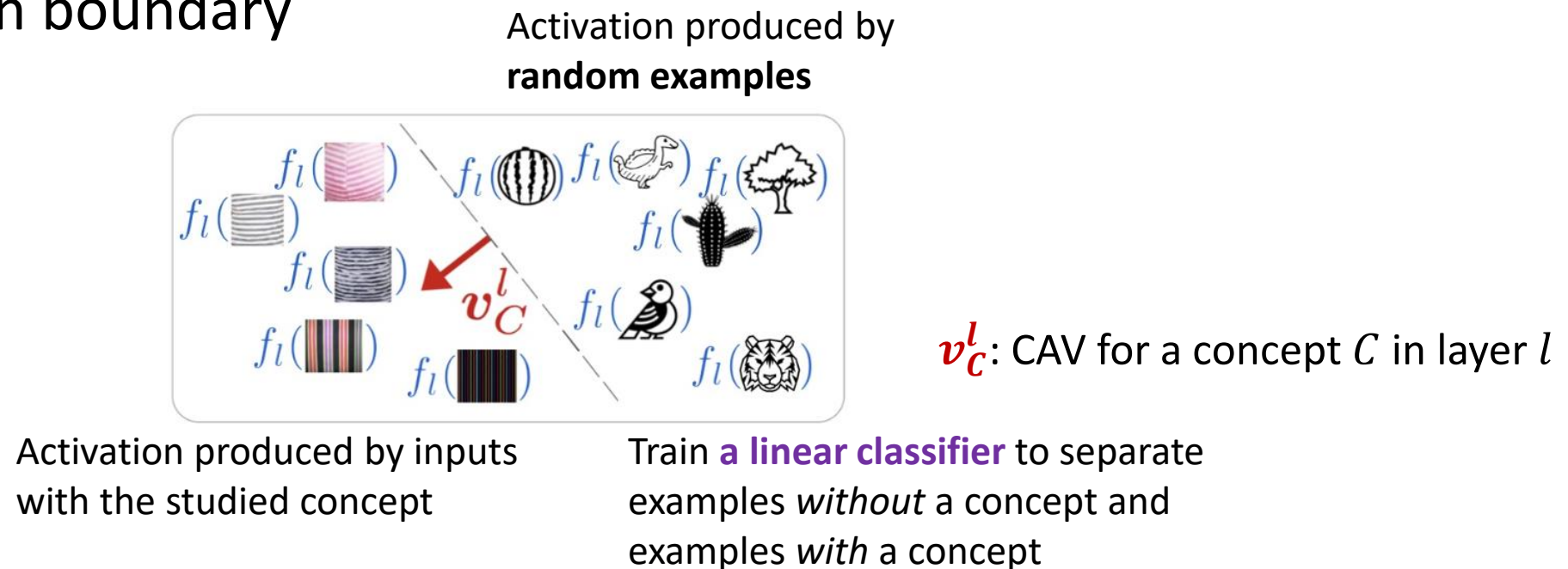
**(e)**

$$S_{C,k,l}(\text{zebra}) = \nabla h_{l,k}(f_l(\text{zebra})) \cdot v_C^l$$

Conceptual Sensitivity to “**striped**” for the class of “**zebras**”

# CAV Definition

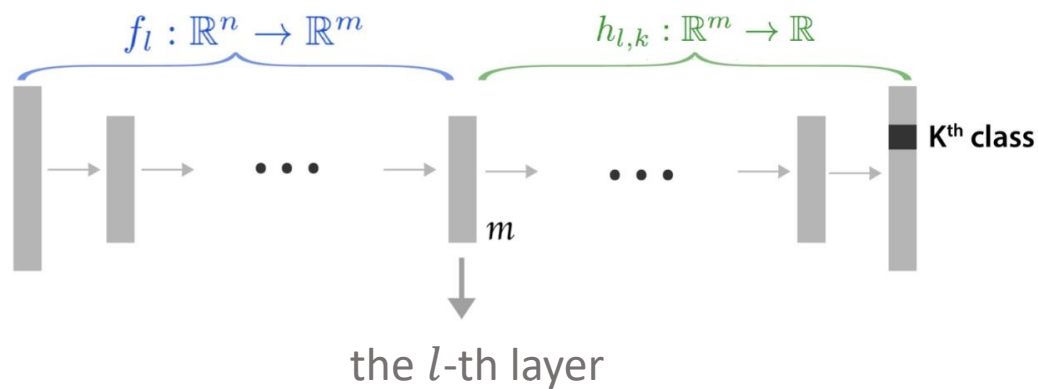
- For a user-defined concept, we seek a **vector in the embedding space** of the  $l$ -th layer that represents this concept
- Concept Activation Vector (CAV): a unit vector orthogonal to the classification boundary



# Conceptual Sensitivity

- $f_l(\mathbf{x})$ : the activations for input  $\mathbf{x}$  at layer  $l$  ;  $h_{l,k}(f_l(\mathbf{x}))$ : the logit for class  $k$
- **Sensitivity of class  $k$  to concept  $C \in \mathbb{R}$ :**

$$S_{C,k,l}(\mathbf{x}) = \lim_{\epsilon \rightarrow 0} \frac{h_{l,k}(f_l(\mathbf{x}) + \epsilon \mathbf{v}_C^l) - h_{l,k}(f_l(\mathbf{x}))}{\epsilon}$$
$$= \nabla h_{l,k}(f_l(\mathbf{x})) \cdot \mathbf{v}_C^l,$$



$\mathbf{v}_C^l \in \mathbb{R}^m$ : CAV for a concept  $C$  in layer  $l$

# Testing with CAVs

- TCAV score is defined as:

$$\text{TCAVQ}_{C,k,l} = \frac{|\{\mathbf{x} \in X_k : S_{C,k,l}(\mathbf{x}) > 0\}|}{|X_k|} \in [0,1]$$

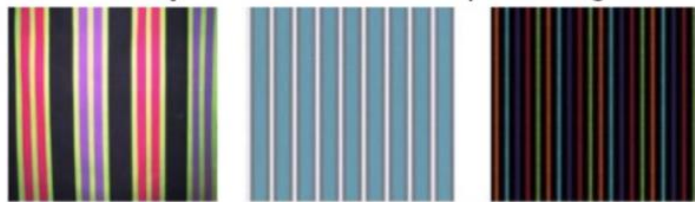
- $X_k$ : all inputs with the class  $k$
- TCAV measures the fraction of inputs with the class  $k$  whose  $l$ -th layer activation vector was **positively sensitive** to concept  $C$  (i.e.,  $S_{C,k,l}(\mathbf{x}) > 0$ )
- Note: TCAV only depends on the sign of  $S_{C,k,l}$ 
  - could be further improved to consider the magnitude



# Application: Sorting Images with CAVs

- CAV essentially encodes the direction of a concept.
- The **cosine similarity** between the picture of interest to the CAV reflects the relation between the picture and the concept.
  - First learn a CAV from CEO / Model Women class (collected from ImageNet)
  - Sort similar/dissimilar images with respect to the learned CAVs

*CEO concept: most similar striped images*



*Model Women concept: most similar necktie images*



*CEO concept: least similar striped images*

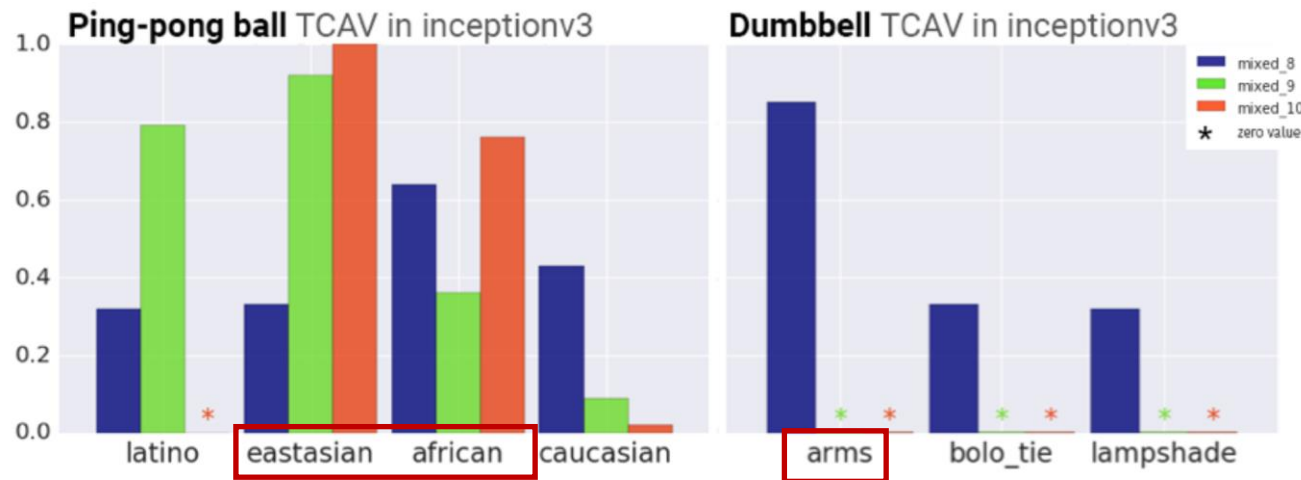
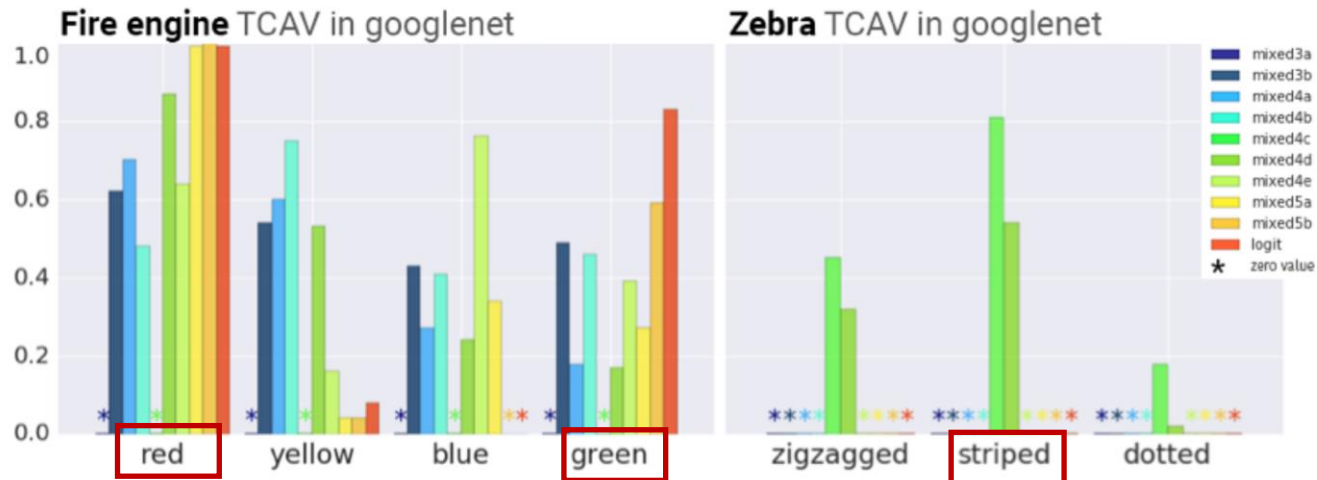


*Model Women concept: least similar necktie images*



**the CAVs correctly reflect  
the concept of interest**

# TCAV Results



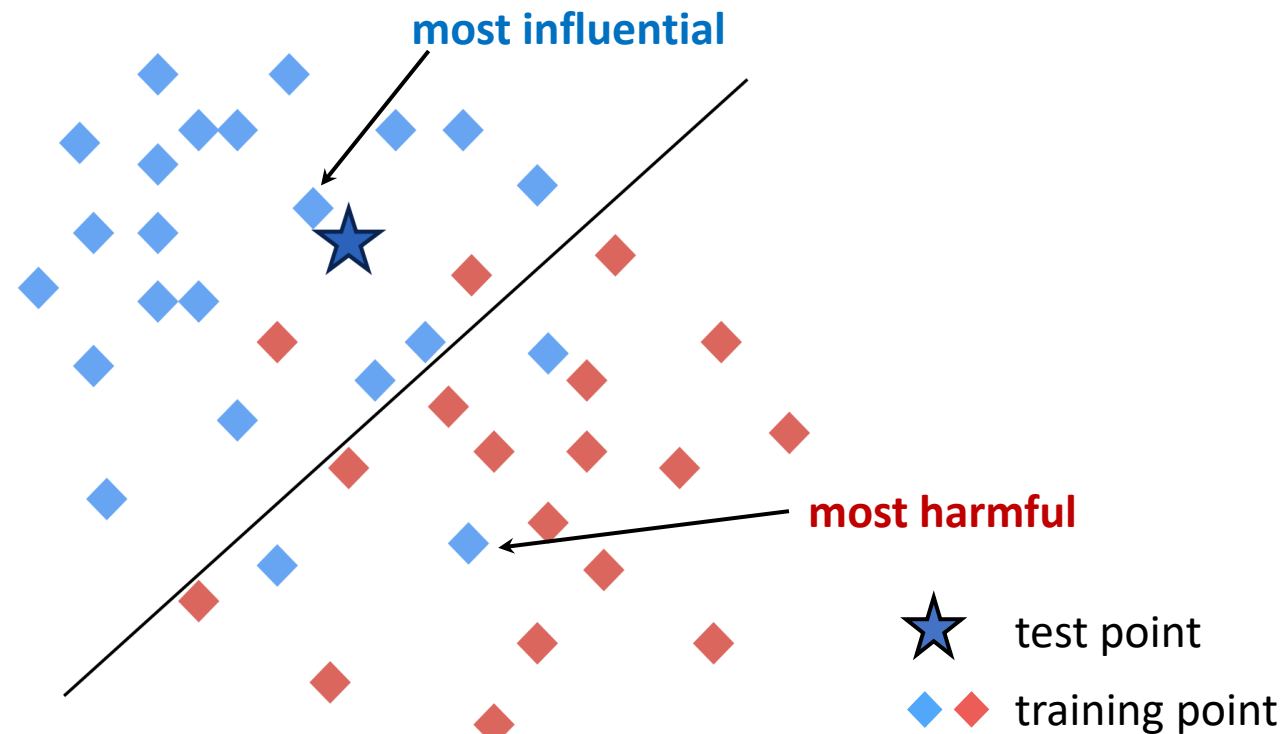
Last 3 layers in [Inception v3](#)

TCAVQs in layers close to the logit layer (red) represent more direct influence on the prediction than lower layers in general.

# Influence Functions: Motivation

Given a well-trained deep learning model, we are interested in

- Which training points were most **influential** for this prediction?
- Which training points were most **harmful** for the prediction?



# Influence Functions: Setting (1)

- **Question:** How to measure the impact of a training point on a prediction?
- We are given training points  $z_1, \dots, z_n$ . How to measure the impact of a training point  $z_{train}$  on the prediction of  $z_{test}$
- Instead of retraining the model on  $\hat{Z} = \{z_i\}_{i=1}^n \cup z_{train}$ , we use influence functions to measure the model changes as we upweight  $z_{train}$  by an infinitesimal amount

# Influence Functions: Setting (2)

- Let  $L(z_i, \theta)$  be the loss, where  $\theta \in \Theta$  represents model parameters.
- $\hat{\theta} := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n L(z_i, \theta)$  is the original optimal parameters
- Given  $\hat{Z} = \{z_i\}_{i=1}^n \cup \mathbf{z}_{train}$ , the optimal parameters become:

$$\hat{\theta}_{\varepsilon, \mathbf{z}_{train}} := \operatorname{argmin}_{\theta \in \Theta} \left[ \frac{1}{n} \sum_{i=1}^n L(z_i, \theta) \right] + \varepsilon L(\mathbf{z}_{train}, \theta)$$

**Assumption:** the empirical risk is twice-differentiable and strictly convex in  $\theta$ .

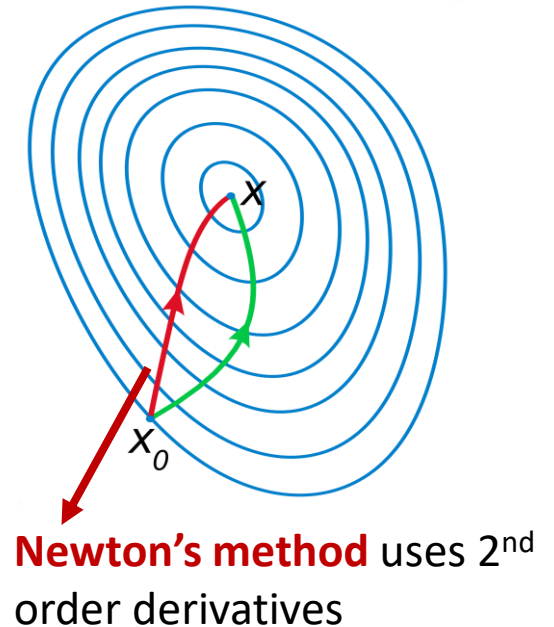
- Goal: approximate the change in  $L(\mathbf{z}_{test}, \hat{\theta}_{\varepsilon, \mathbf{z}_{train}})$  as we increase  $\varepsilon$

# Influence Functions: Definition

- Under smoothness assumptions:

$$\begin{aligned} \mathcal{I}_{up,loss}(\mathbf{z}_{train}, \mathbf{z}_{test}) &\stackrel{\text{def}}{=} \left. \frac{dL(\mathbf{z}_{test}, \hat{\theta}_{\epsilon, \mathbf{z}_{train}})}{d\epsilon} \right|_{\epsilon=0} \\ &= \nabla_{\theta} L(\mathbf{z}_{test}, \hat{\theta})^{\top} \left. \frac{d\hat{\theta}_{\epsilon, \mathbf{z}_{train}}}{d\epsilon} \right|_{\epsilon=0} \\ &= -\nabla_{\theta} L(\mathbf{z}_{test}, \hat{\theta})^{\top} H_{\hat{\theta}}^{-1} \nabla_{\theta} L(\mathbf{z}_{train}, \hat{\theta}) \end{aligned}$$

- where  $H_{\hat{\theta}} = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta}^2 L(z_i, \hat{\theta})$
- In essence, influence functions form a **quadratic approximation** to the empirical risk around  $\hat{\theta}$  and take a single Newton step.



# Use case: Understand Model Behavior (1)

Influence functions reveal insights about how models rely on and extrapolate from the training data.

- Dataset: Dog & Fish image classification from ImageNet dataset
- Two well trained models: (1) [Inception v3 network](#) and (2) an SVM with an RBF kernel
- Investigate the impact of training points on a test image (**fish**) that both models got correct prediction

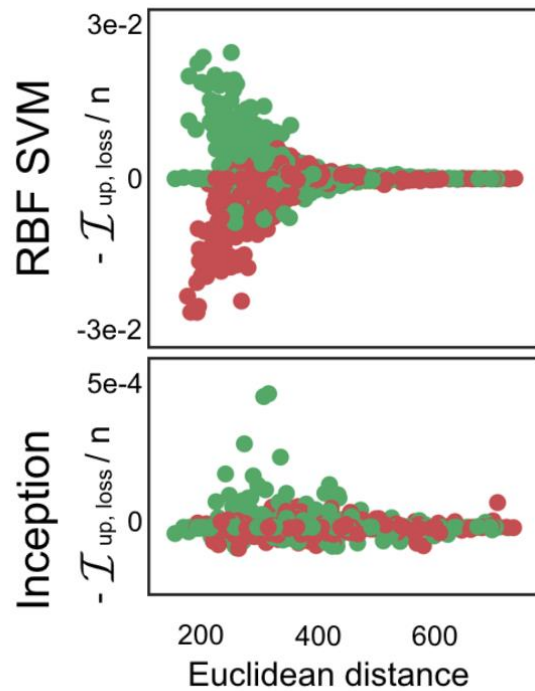
Test image





# Use case: Understand Model Behavior (2)

$\mathcal{I}_{up,loss}(Z_{train}, Z_{test})$  V.S. Euclidean distance  $\|Z_{train} - Z_{test}\|$



**In RBF-SVM:** training images far from the test image in pixel space having almost no influence;  
(emphasizing nearby samples)  
**Fish** images are mostly helpful, while  
**dog** images are mostly harmful

**In Inception network,** **fish** and **dogs** both could be helpful or harmful for correctly classifying the test image. The influence is not related to the distance.

**Green dot: fish**

**Red dot: dog**

Note: the test image is a **fish** image

Most helpful training images for **RBF-SVM**



Most helpful training images for **Inception**



the 5th most helpful training image for **Inception** is a dog image



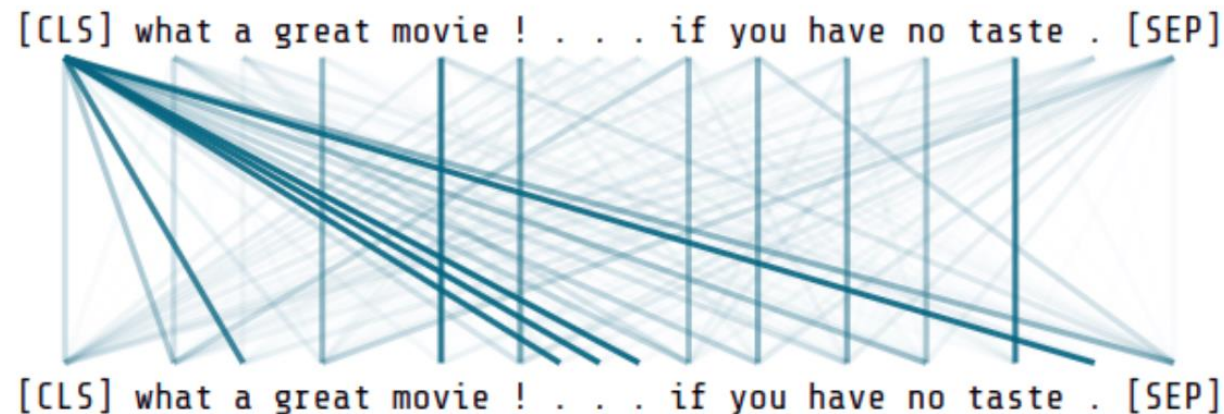
# Content

- Evaluating Explainability Methods
- Model-level Explainability
- Intrinsic Explainability / Interpretability

# Explainability via Attention

- **Attention Mechanisms**

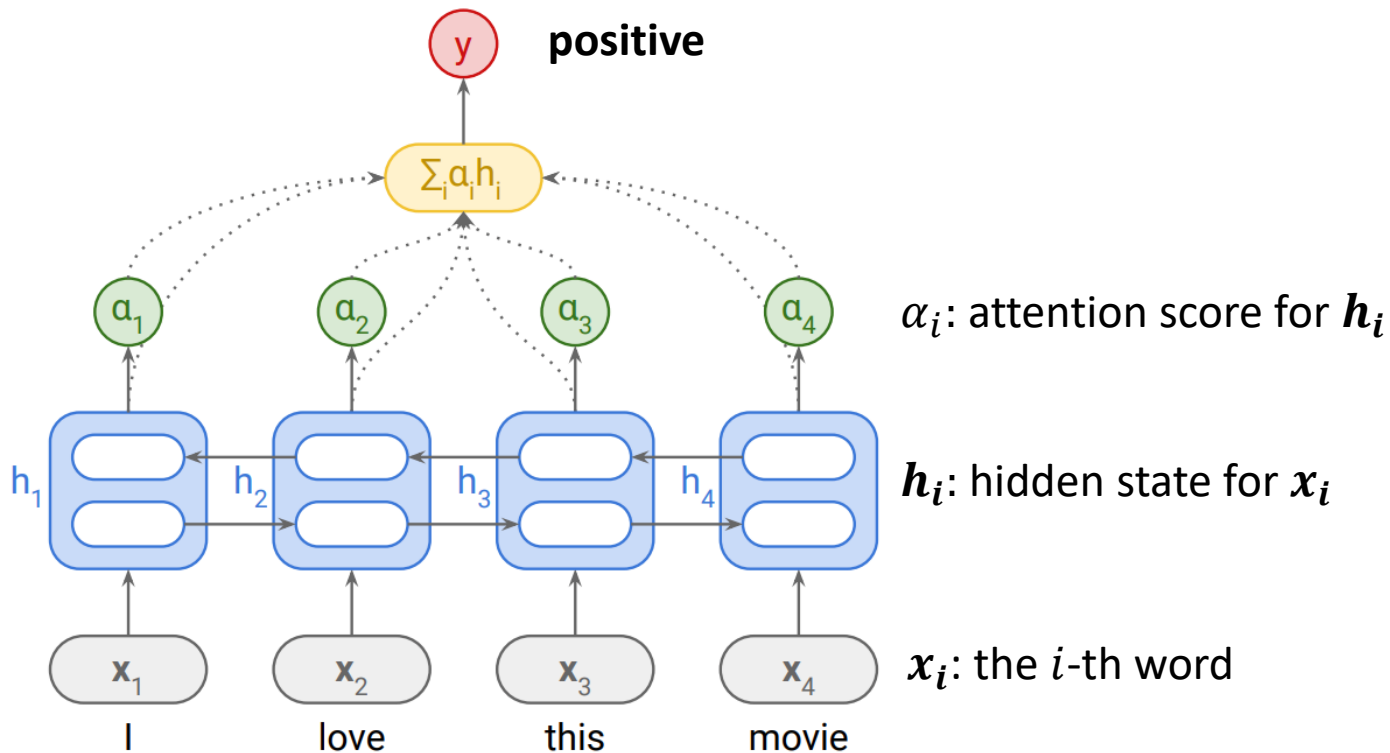
- DNNs can be endowed with attention mechanisms that simultaneously
  - preserve or even **improve their performance**
  - obtain **explainable outputs**
- Visualize attention weights in a attention models:



Color represents the value of attention weight  
**darker blue ⇔ larger attention weight**

# Attention in Text Classification

## BiLSTM for sentiment analysis (text classification)



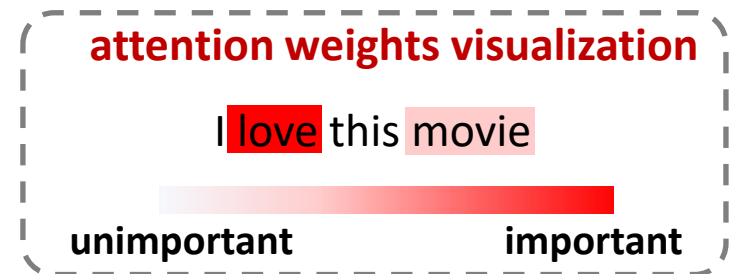
## Attention score:

$$\alpha_i = \frac{\exp e_i}{\sum_k \exp e_k}$$

where  $e_i = \mathbf{v}^\top \tanh(W_h \mathbf{h}_i + W_q \mathbf{q})$  and  $\mathbf{v}, \mathbf{q}, W_h, W_q$  are learnable parameters

(example of **MLP-based attention**)

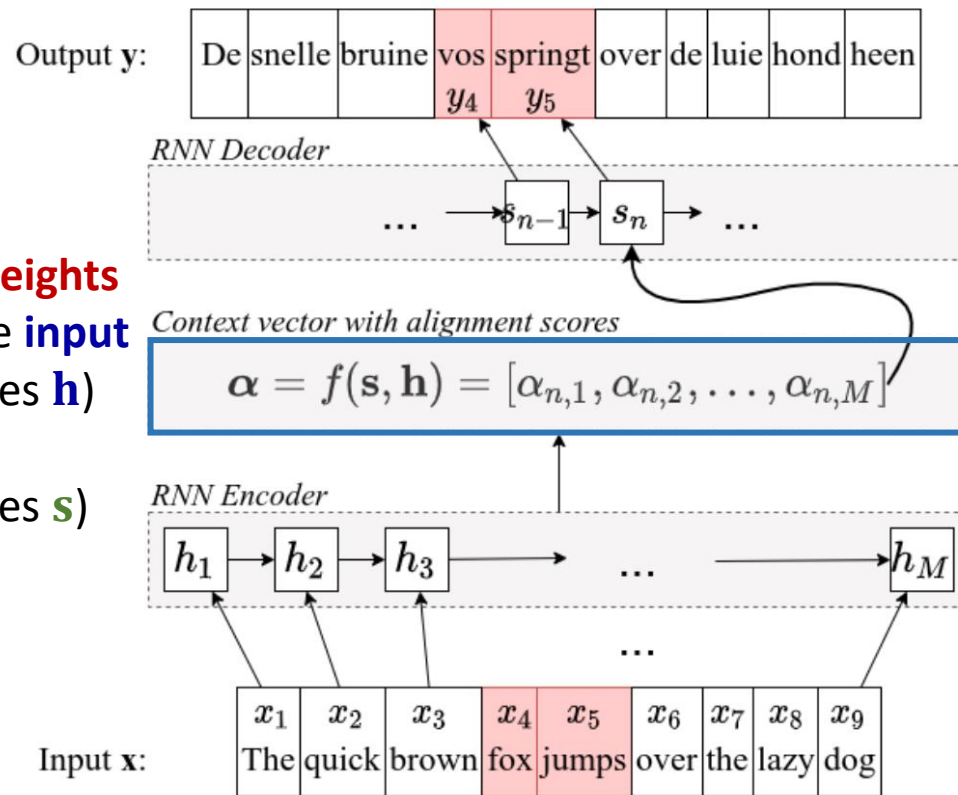
$\alpha_i$  shows the **importance** of the word  $x_i$  to the prediction " $y = \text{positive}$ "



Bastings, Jasmijn, and Katja Filippova. "The elephant in the interpretability room: Why use attention as explanation when we have saliency methods?."

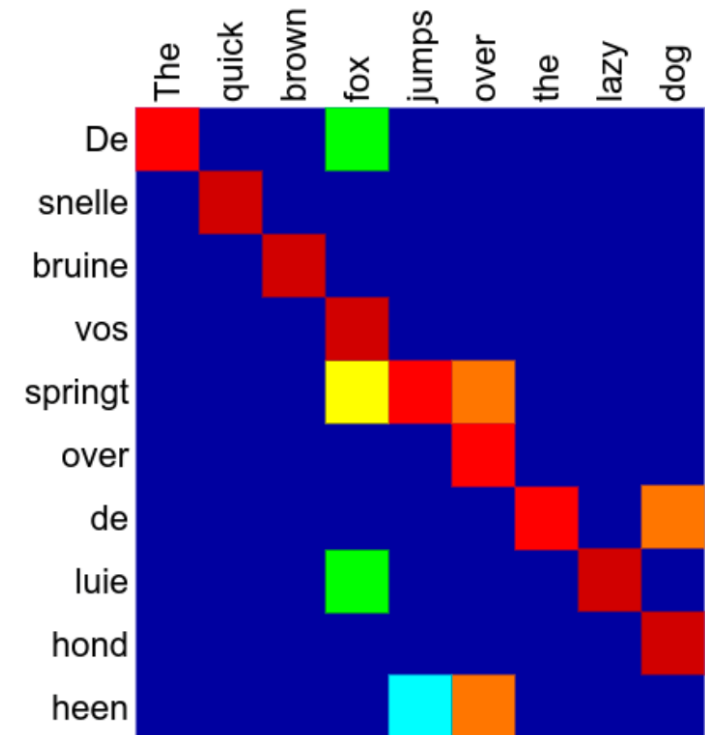
# Attention in Sequential Generation

## RNN model for Natural Language Translation



Attention weights between the **input** (hidden states  $h$ ) and **output** (hidden states  $s$ )

Attention weights  $\alpha_{ij}$  shows the importance of input  $x_j$  to the output  $y_i$



unimportant

important

Ras, Gabrielle, et al. "Explainable deep learning: A field guide for the uninitiated."

# Signed Attention (1)

- **Attention weight is always positive**
- Ideal explanation should discriminate between **positive** and **negative** contributions towards a prediction
- Solution: **signed attention**

$$A_i = - \boxed{\frac{\partial \mathcal{L}}{\partial \alpha_i}} \times \alpha_i$$

Indicates the **positive or negative contribution**

Recap: for hidden state  $\mathbf{h}_i$

$$\text{attention weight: } \alpha_i = \frac{\exp e_i}{\sum_k \exp e_k} \geq 0$$

where  $e_i = \mathbf{v}^\top \tanh(W_h \mathbf{h}_i + W_q \mathbf{q})$

$\mathbf{v}, \mathbf{q}, W_h, W_q$ : learnable parameters

- $\mathcal{L}$ : loss function
- $\alpha_i$ : original attention weight
- value of  $\alpha_i$  measures the strength of the contribution

# Signed Attention (2)

## Explanation for sentiment analysis

**Input:** *"Though the price may be too expensive, I love its surprisingly high quality."*

**Output:**  $y$  = "Positive"

### original attention visualization

*Though the price may be too expensive, I love its surprisingly high quality.*

unimportant  important

with  
signed attention



### signed attention visualization

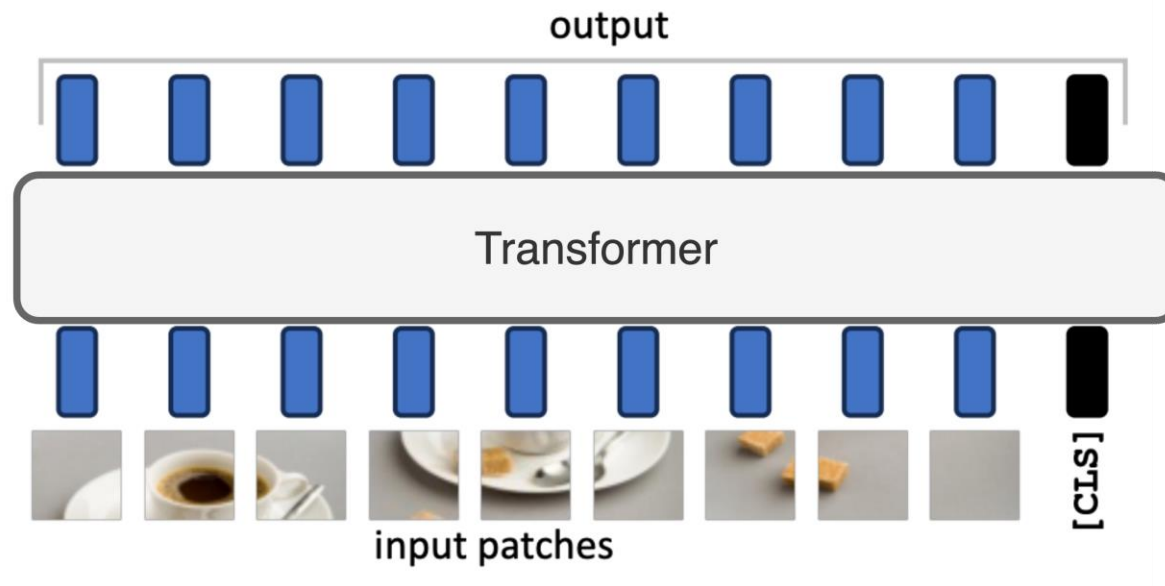
*Though the price may be too expensive, I love its surprisingly high quality.*

negative  positive

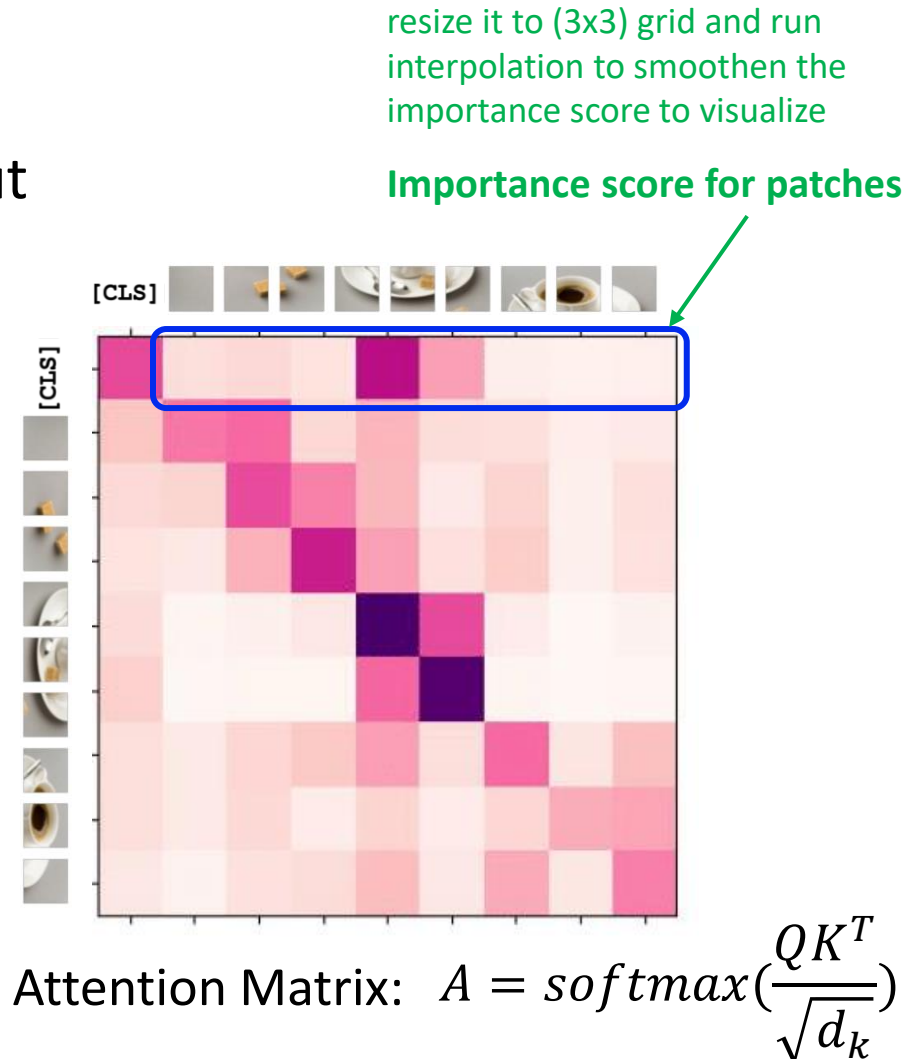
# Explainability in Vision Transformer (ViT)

- **Vision Transformer**

- Split an input image into NxN patches
- Add a [CLS] token as a global embedding of the input



Architecture



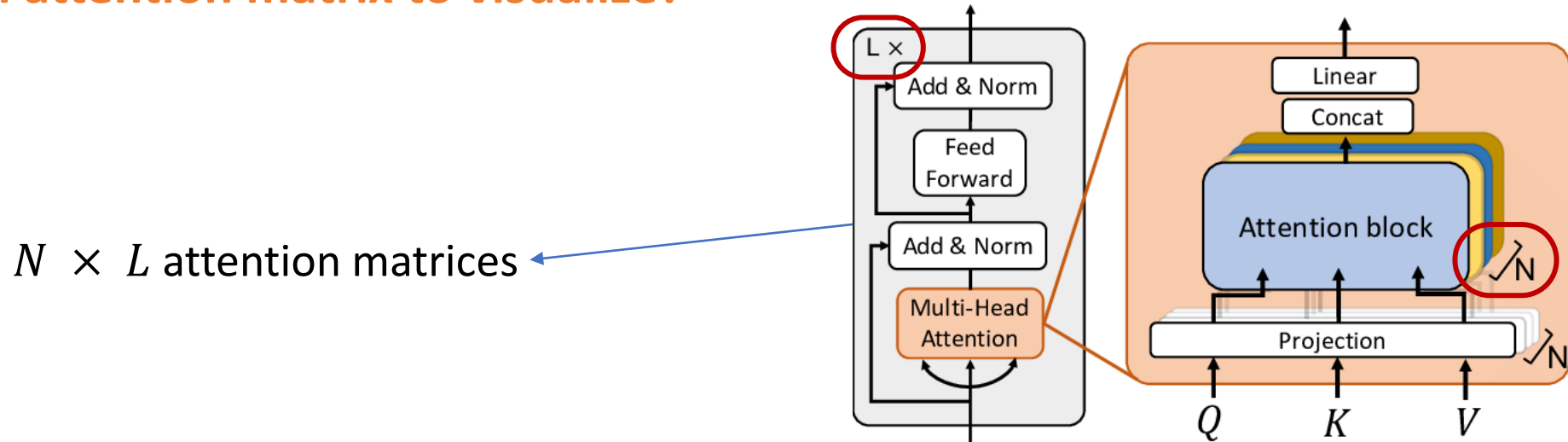
# Explainability in Vision Transformer (ViT)

- **Vision Transformer**

- Split an input image into  $N \times N$  patches
- Add a [CLS] token as a global embedding of the input

- **Challenge:**

- Transformer has multiple heads and layers, thus having multiple attention matrices  
→ **Which attention matrix to visualize?**





# Explainability in Vision Transformer (ViT)

- **Naïve approach (Rollout)**

- Aggregate attention matrices across multiple heads: **Mean averaging**

$$\bar{A}^{(l)} = \boxed{I +} \sum_i A^{(l,i)}$$

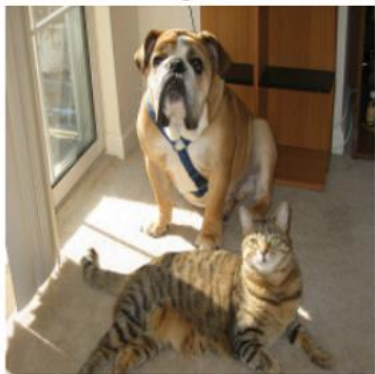
Why do we need to add I?

To account for skip connections in Transformer blocks

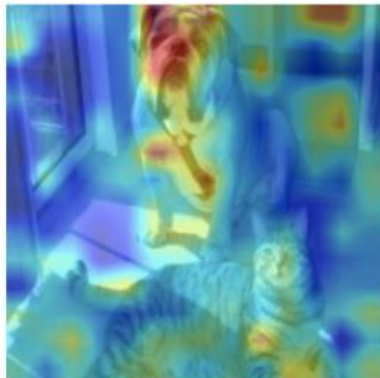
- Aggregate attention matrices across multiple heads: **Matrix multiplication**

$$C = \bar{A}^1 \bar{A}^2 \dots \bar{A}^N$$

Dog →

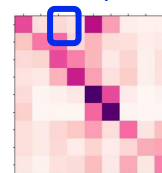


Attention visualization



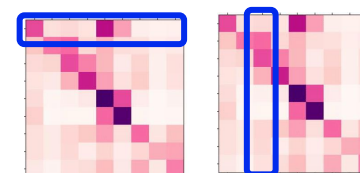
Why should it be matrix multiplication?  
Weighted aggregation importance score of nodes

Importance of patch 3



Importance of each patch

Attention weights of patch 3 to others



# Explainability in Vision Transformer (ViT)

- **Naïve approach (Rollout)**

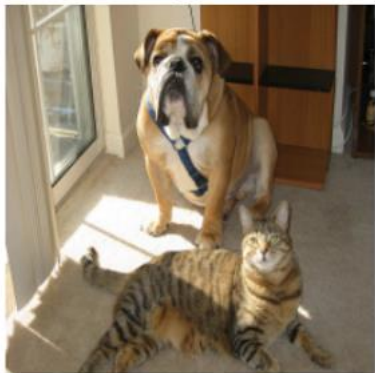
- Aggregate attention matrices across multiple heads: **Mean averaging**

$$\bar{A}^{(l)} = I + \sum_i A^{(l,i)}$$

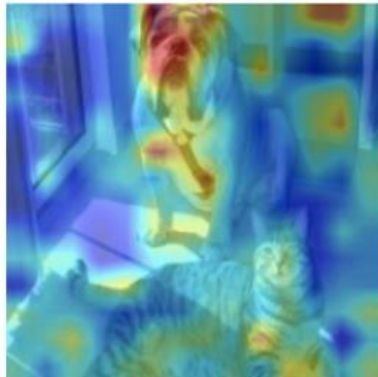
- Aggregate attention matrices across multiple heads: **Matrix multiplication**

$$C = \bar{A}^1 \bar{A}^2 \dots \bar{A}^N$$

Dog →

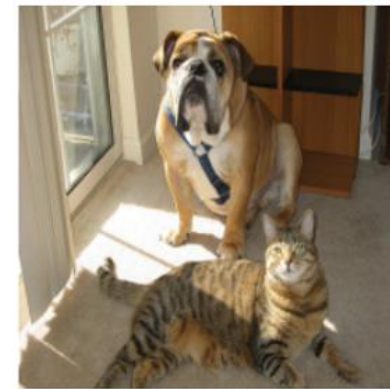


Attention visualization

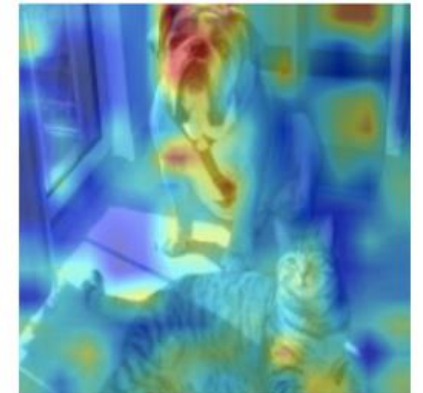


**However, the explanation  
for cat prediction is the  
same as for dog**

Cat →



Attention visualization



# Explainability in Vision Transformer (ViT)

- **Targeted Explanation**

- Aggregate attention matrices across multiple heads: **Relevance** and **gradient** diffusion

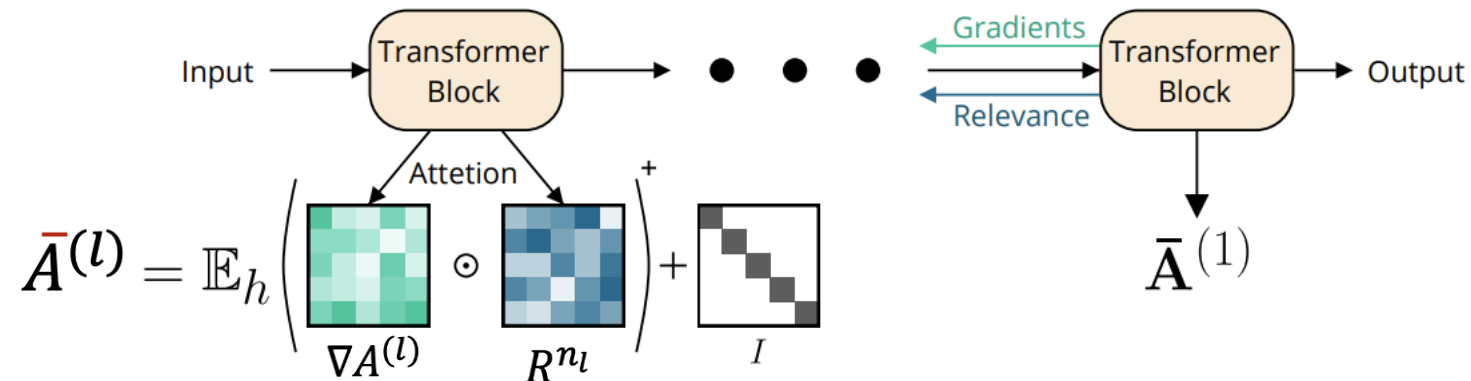
$$\bar{A}^{(l)} = I + \sum_h \boxed{\nabla A^{(l)}} \odot \boxed{R^{n_l}}$$

The gradient of the targeted output to the attention matrix

Layer-wise Relevance Propagation matrix

- Aggregate attention matrices across multiple heads: **Matrix multiplication**

$$C = \bar{A}^1 \bar{A}^2 \dots \bar{A}^N$$



# Layer-wise Relevance Propagation (LRP)

- **Layer-wise Relevance Propagation** is a method to propagate the relevance of input to the target output through layers.

- Let  $L^{(l)}(\mathbf{X}, \mathbf{Y})$  be the layer's operation on two tensors  $\mathbf{X}$  and  $\mathbf{Y}$
- Relevance score is computed by Deep Taylor Decomposition:

$$\begin{aligned} R_j^{(n)} &= \mathcal{G}(\mathbf{X}, \mathbf{Y}, R^{(n-1)}) \\ &= \sum_i \mathbf{X}_j \frac{\partial L_i^{(n)}(\mathbf{X}, \mathbf{Y})}{\partial \mathbf{X}_j} \frac{R_i^{(n-1)}}{L_i^{(n)}(\mathbf{X}, \mathbf{Y})}, \end{aligned}$$

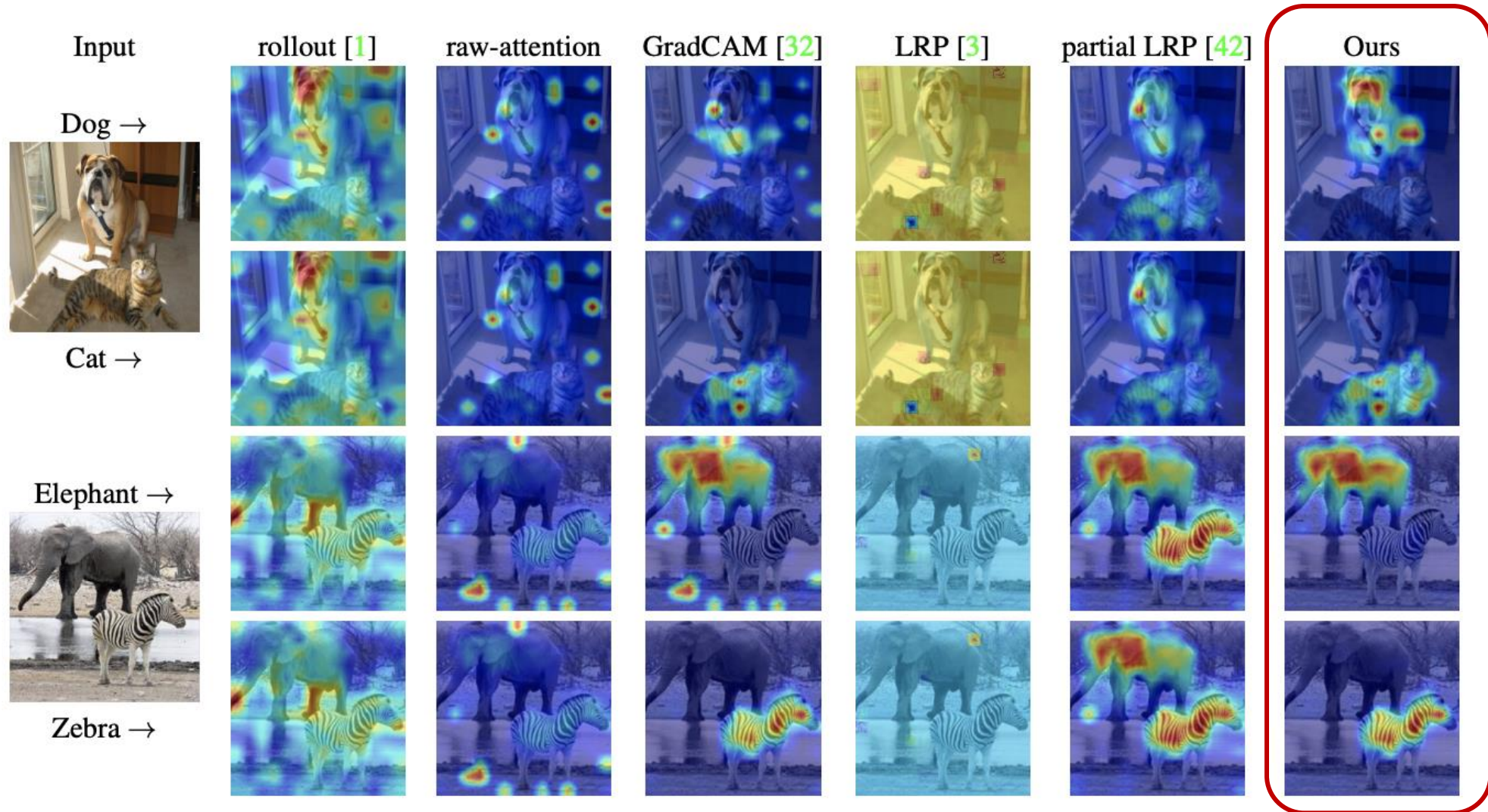
- And then propagate relevance through layers as follows

$$R_j^{(n)} = \mathcal{G}(x^+, w^+, R^{(n-1)}) = \sum_i \frac{x_j^+ w_{ji}^+}{\sum_{j'} x_{j'}^+ w_{j'i}^+} R_i^{(n-1)}$$

where  $\mathbf{X} = x$  and  $\mathbf{Y} = w$  are the layer's input and weights.  
The superscript denotes the operation  $\max(0, v)$  as  $v^+$ .



# Transformer Interpretability



# Vision Transformer Need Registers

- **Artifacts in Vision Transformer**

- Most of the existing transformer-based models exhibit artifacts on their attention matrix.

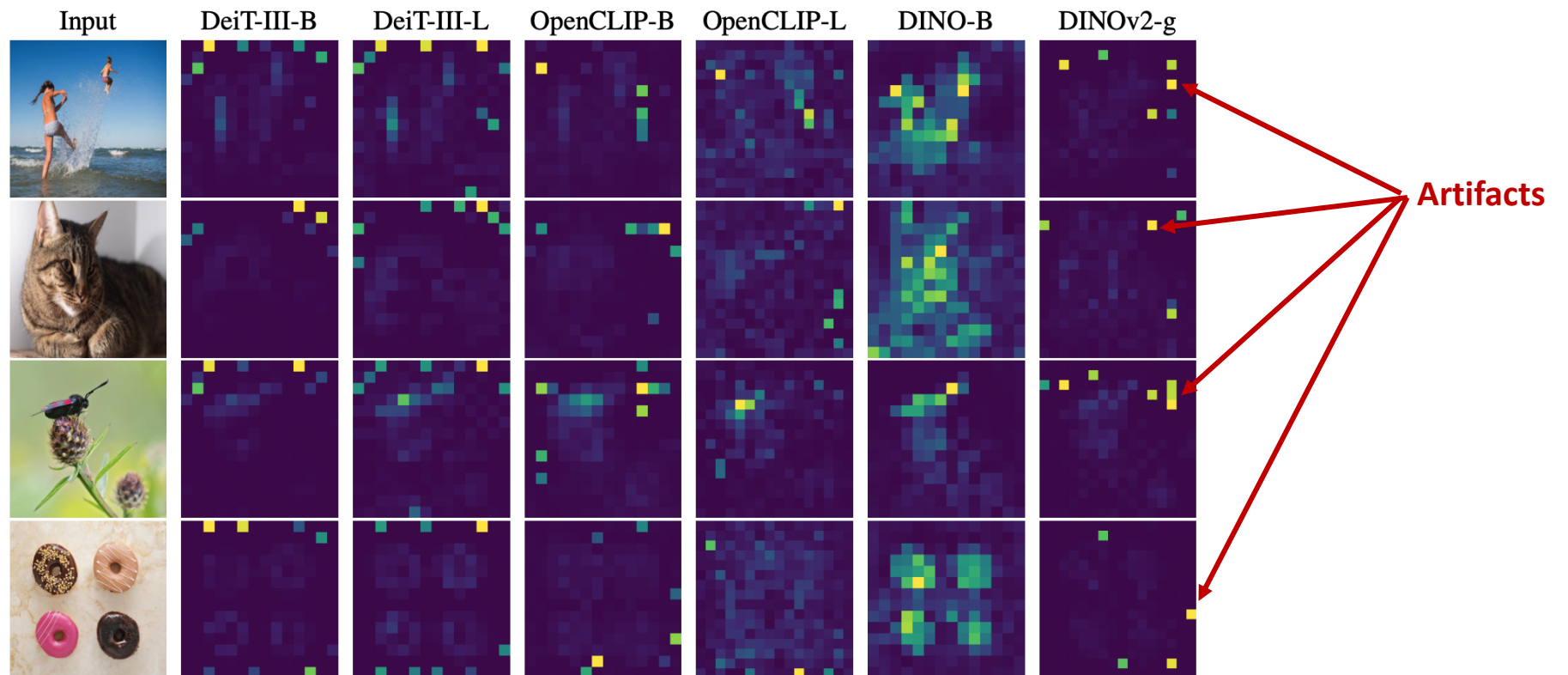
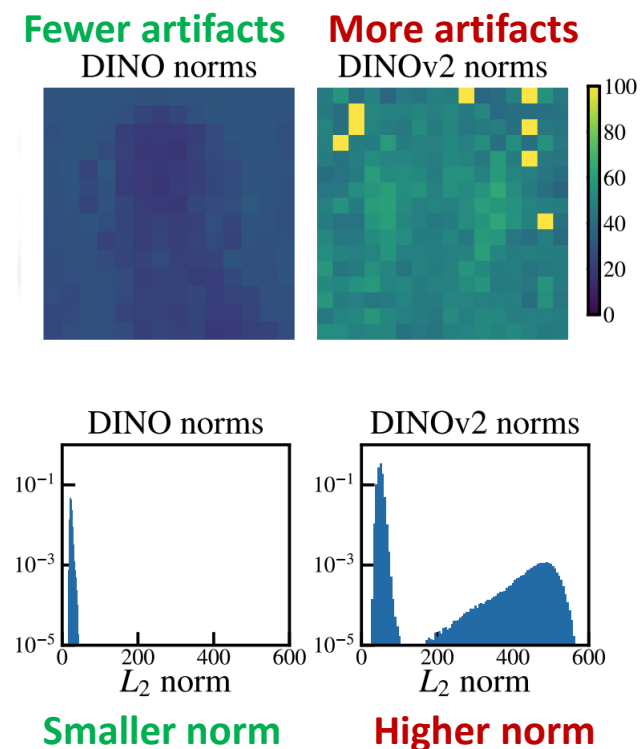
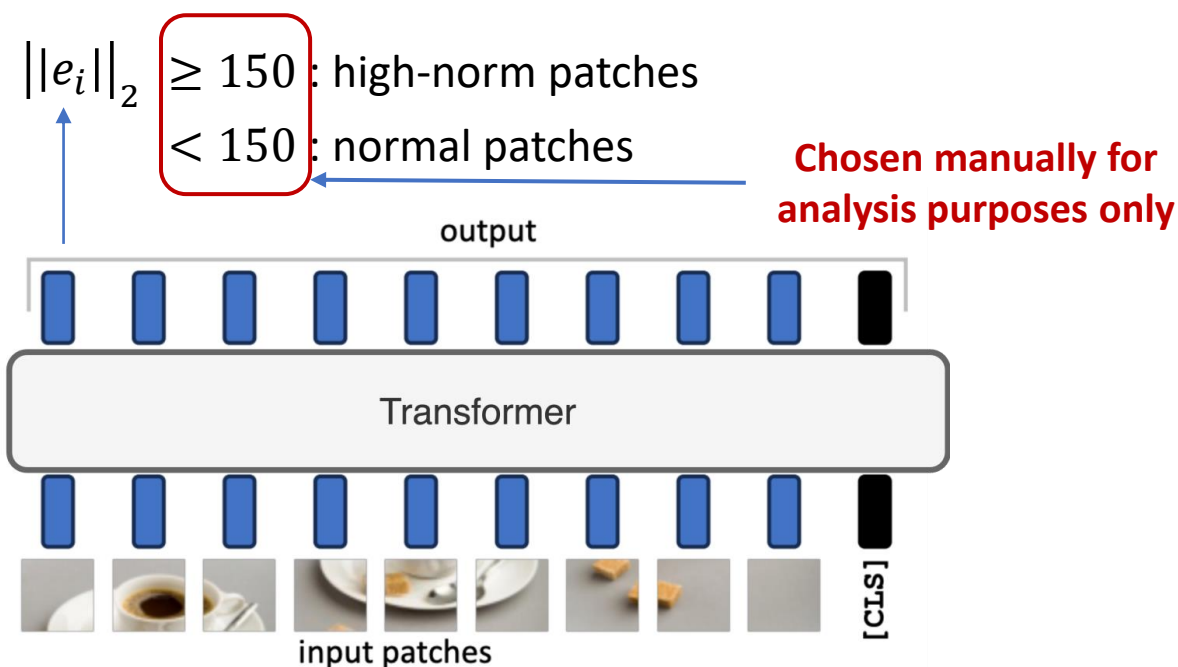


Figure 2: Illustration of artifacts observed in the attention maps of modern vision transformers.

# Vision Transformer Need Registers

- **Why does it happen?**

- The artifacts have a connection to the norm of token embedding.  
→ We can analyze tokens with high-norm to analyze the artifacts created by the attention matrix.



# Vision Transformer Need Registers

## • Settings

- For each patch/token embedding, add simple linear layers to
  - Predict the position of the patch
  - Reconstruct pixel values on the patch
  - Predict image class from the patch embedding

## • Observation

- Position prediction & reconstruction: **normal patches give better results.**  
→ normal patches can maintain local information about patches
- Image class classification: **high-norm patches perform better**  
→ high-norm patches discard local information, having more global information (image class)

	position prediction		reconstruction
	top-1 acc	avg. distance ↓	L2 error ↓
normal	<b>41.7</b>	<b>0.79</b>	<b>18.38</b>
outlier	22.8	5.09	25.23

(b) Linear probing for local information.

	IN1k	P205	Airc.	CF10	CF100	CUB
[CLS]	<b>86.0</b>	<b>66.4</b>	<b>87.3</b>	<b>99.4</b>	<b>94.5</b>	<b>91.3</b>
normal	65.8	53.1	17.1	97.1	81.3	18.6
outlier	<u>69.0</u>	<u>55.1</u>	<u>79.1</u>	<u>99.3</u>	<u>93.7</u>	<u>84.9</u>

Image classification



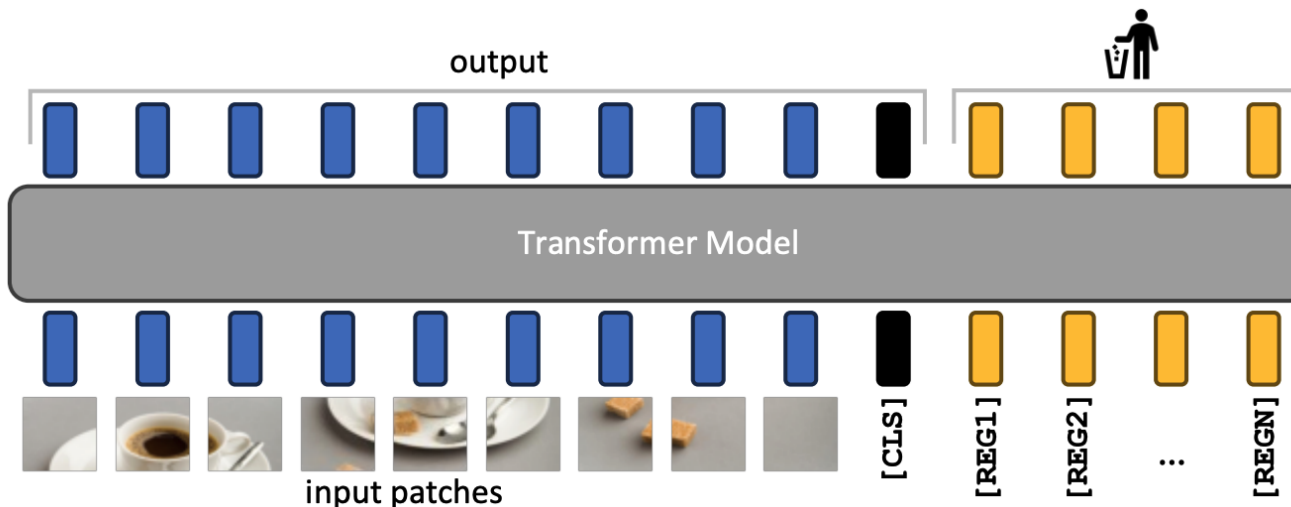
# Vision Transformer Need Registers

- **Hypothesis**

- Large Transformer models can *recognize redundant patches* (do not have much information) and leverage them to *store, process, and retrieve* global information.
- It may be undesirable as it may discard information from some patches.

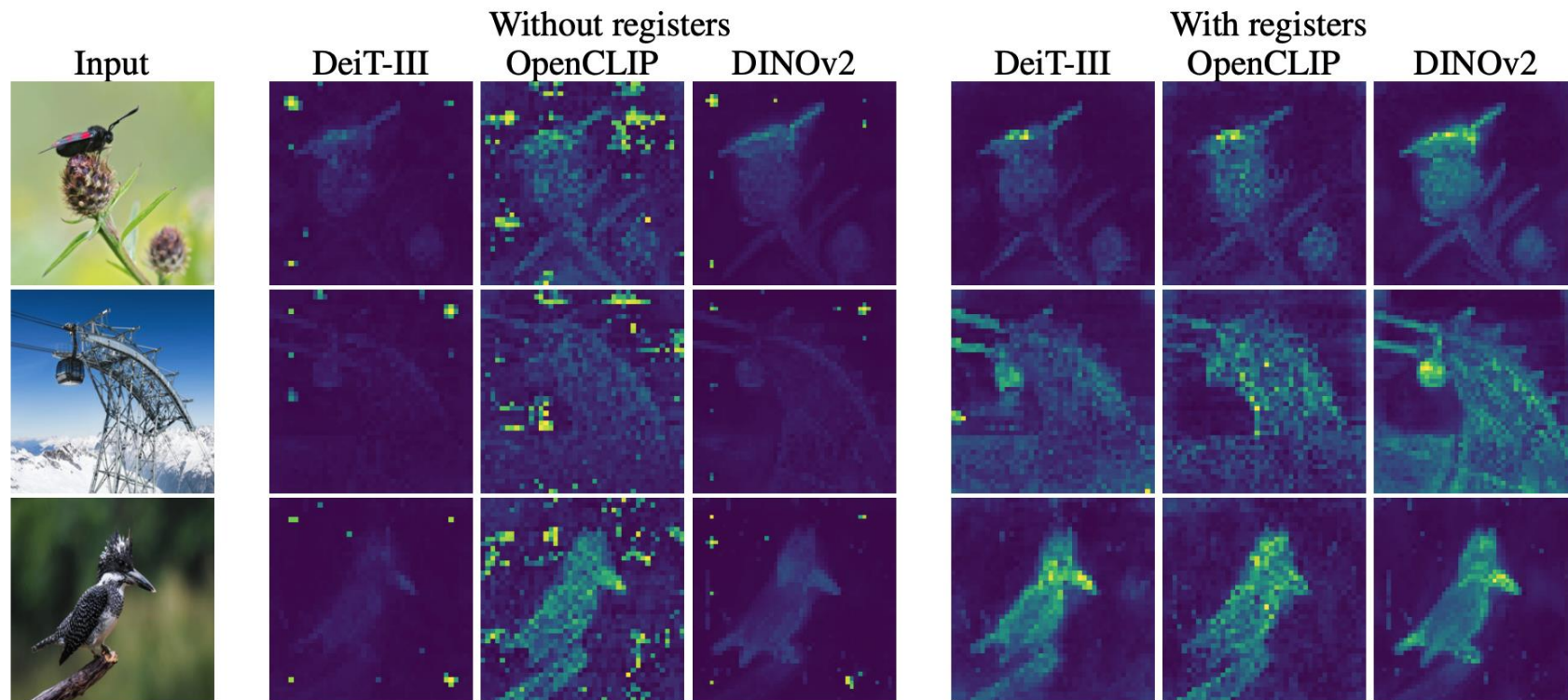
- **Solution**

- Use some additional tokens as registers (along with [CLS]) for saving global information purposes.



# Vision Transformer Need Registers

- Adding registers provides much better interpretability and reduces artifacts for the attention matrix.



# Vision Transformer Need Registers

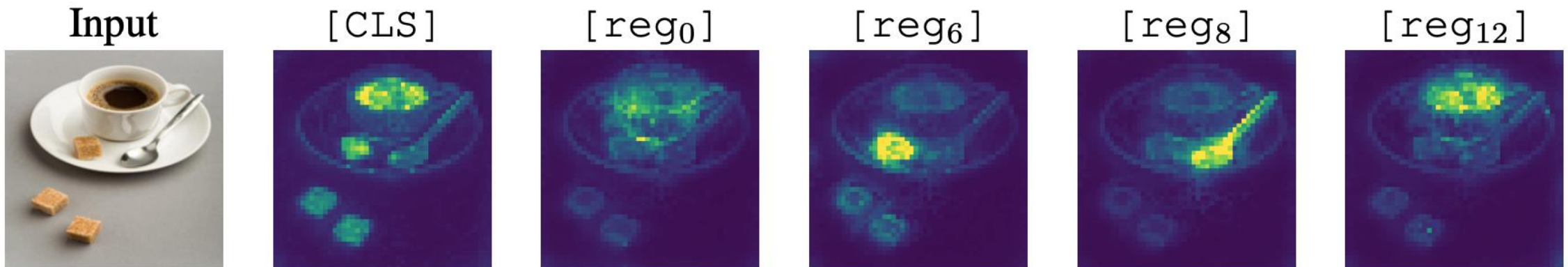
- Adding registers provides much better interpretability and reduces artifacts for the attention matrix.
- Performance slightly improved

	ImageNet Top-1	ADE20k mIoU	NYUd rmse ↓
DeiT-III	84.7	38.9	0.511
DeiT-III+reg	84.7	39.1	0.512
OpenCLIP	78.2	26.6	0.702
OpenCLIP+reg	78.1	26.7	0.661
DINOv2	84.3	46.6	0.378
DINOv2+reg	84.8	47.9	0.366

(a) Linear evaluation with frozen features.

# Vision Transformer Need Registers

- Adding registers provides much better interpretability and reduces artifacts for the attention matrix.
- Performance slightly improved.
- Each register pays attention to different regions (naturally emerged from training).



**Oral ICLR24 → A simple idea but good analyses/observation would also be appreciated**