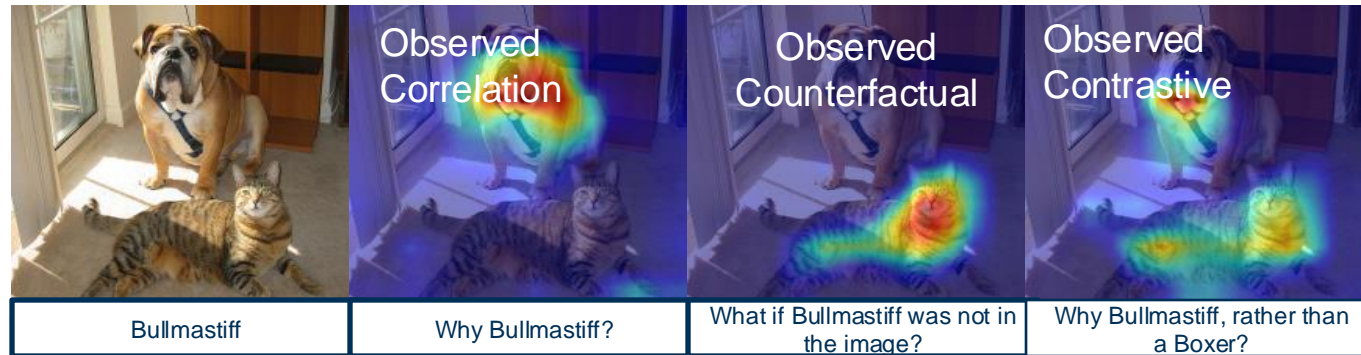


ECE 4252/8803: Fundamentals of Machine Learning (FunML)

Fall 2024

Lecture 22: Explainability in Neural Networks



Overview

In this Lecture..

Explainability

- Definition
- Why Explainability

Visualization of Convolutional Neural Networks

Types of Explanations

Explanatory Evaluation

What is Explainability?

The ability of an entity to explain or justify its decisions or predictions in human-understandable terms

Explainability

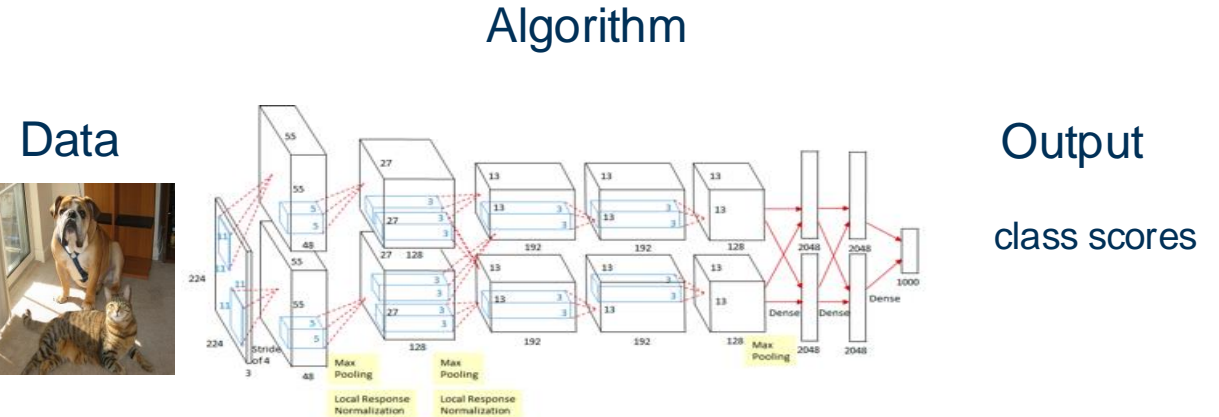
Why Explainability?

Explainability matters:

establish **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



Overview

In this Lecture..

Explainability

Visualization of Convolutional Neural Networks

- Explainability in CNNs
- Visualizing filters
- Dimensionality Reduction with Last layer Embeddings
- Visualizing activations
- Gradient-based visualizations
- Saliency Maps and Intermediate feature visualization
- Grad-CAM visualization and explanations

Types of Explanations

Explanatory Evaluation

Explainability

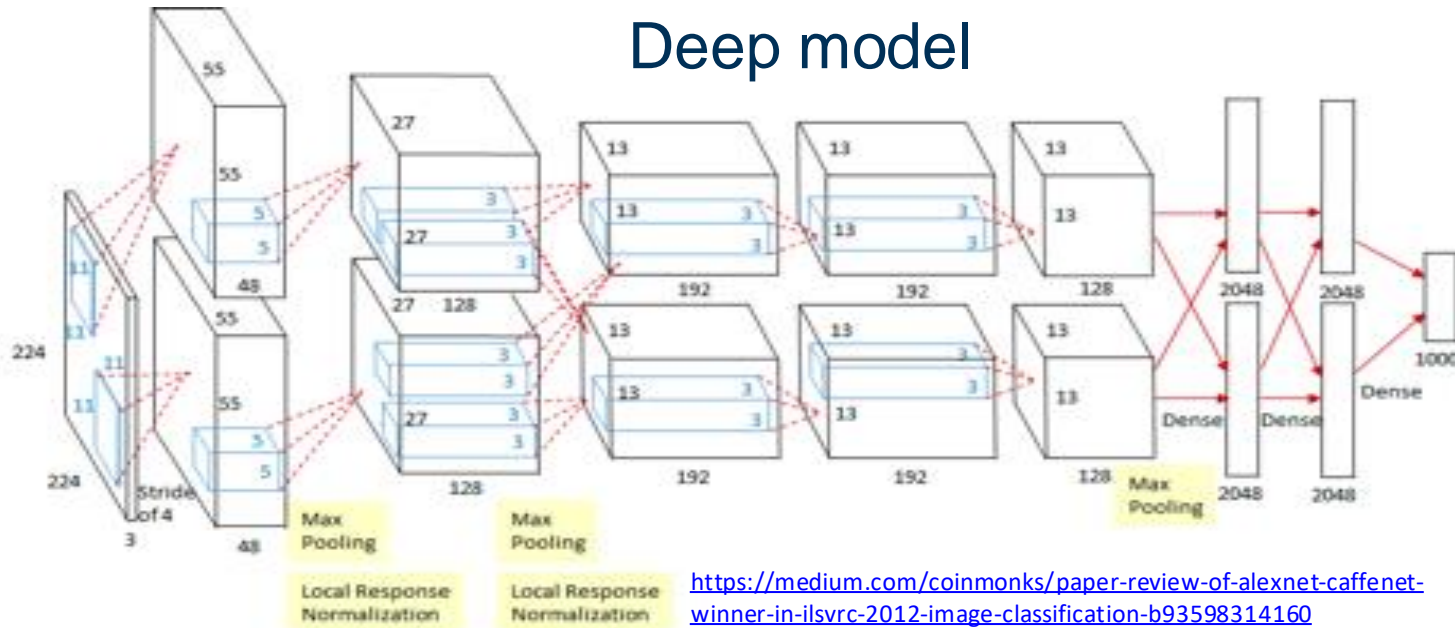
Explainability in CNNs

Data



Input Image:
3 x 224 x 224

Deep model



Output

Output class scores:
1000 numbers
(trained on ImageNet)

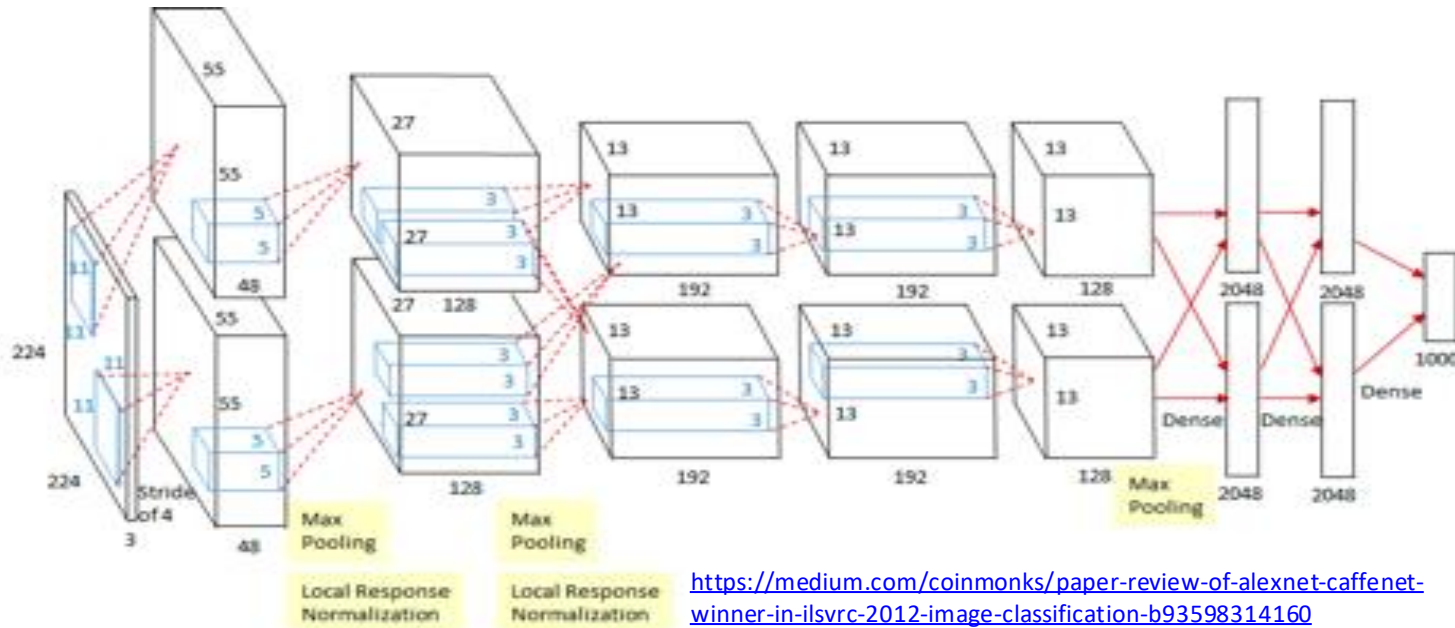
Not Explainable

Explainability

Explainability in CNNs



Input Image:
3 x 224 x 224



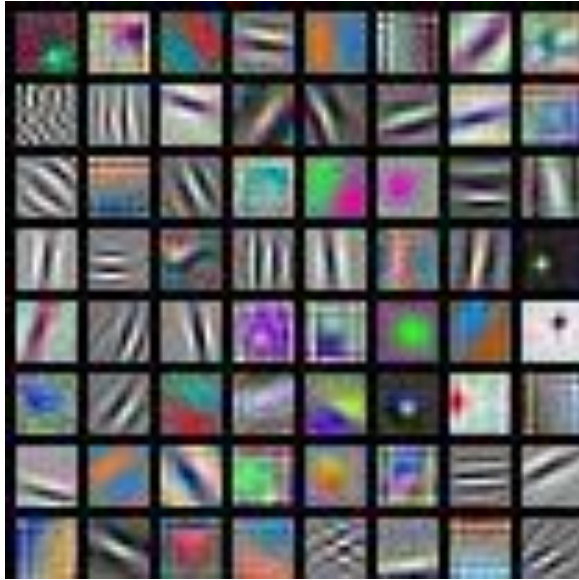
Output class scores:
1000 numbers
(trained on ImageNet)

<https://medium.com/coinmonks/paper-review-of-alexnet-caffenet-winner-in-ilsvrc-2012-image-classification-b93598314160>

What are these layers looking for?
Are they explainable?

Visualizing CNNs

Visualizing Filters in First Layers



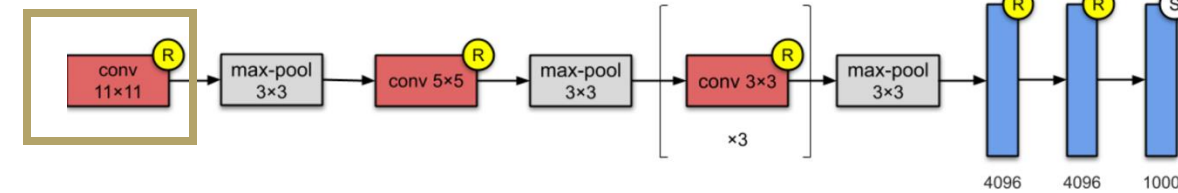
AlexNet:
64 x 3 x 11 x 11

Filters always extend the full depth of the input volume

- 64 filters in the first convolutional layer
- Filter size: 11 x 11 x 3 (visualize as RGB images)
- Filters are looking for **low-level oriented edges, color blobs, textures, background etc.**



Input Image:
3 x 224 x 224

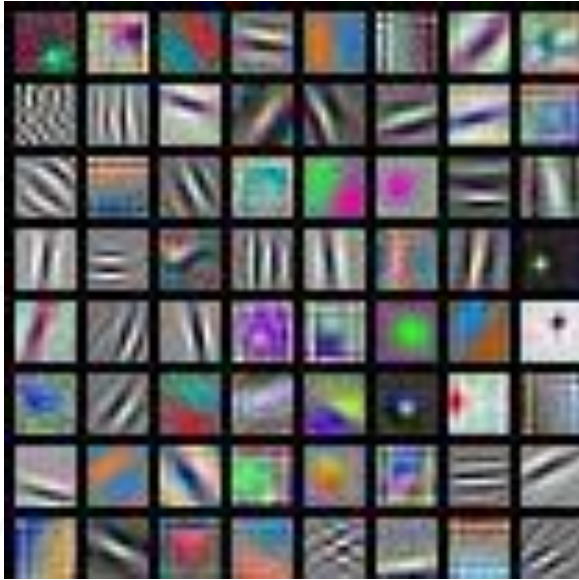


AlexNet

Filters = Weights

Visualizing CNNs

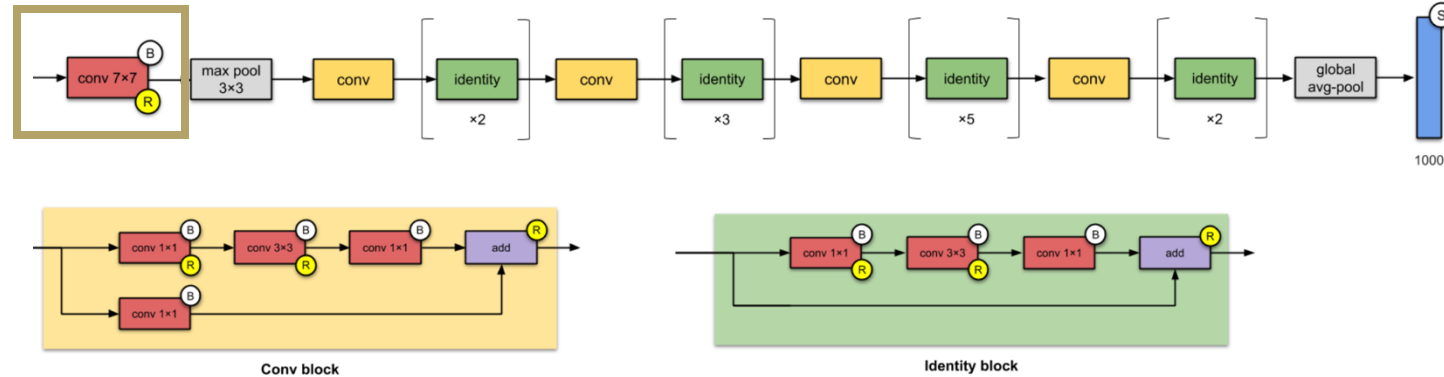
Visualizing Filters in First Layers



AlexNet:
 $64 \times 3 \times 11 \times 11$



ResNet-18:
 $64 \times 3 \times 7 \times 7$



ResNet

- Filters in the first convolutional layers **across different architectures learn similar patterns**

Visualizing CNNs

Visualizing Filters in Intermediate Layers

Visualize the filters
(raw weights)

Filters in **higher**
convolutional **layers** are
not as interpretable as
filters in **the first layer**

Weights:



Weights:



Weights:



Conv layer 1 weights
 $16 \times 3 \times 7 \times 7$
(visualize as RGB images)

Conv layer 2 weights
 $20 \times 16 \times 7 \times 7$
(visualize as 16
grayscale images)

Conv layer 3 weights
 $20 \times 20 \times 7 \times 7$
(visualize as 20
grayscale images)

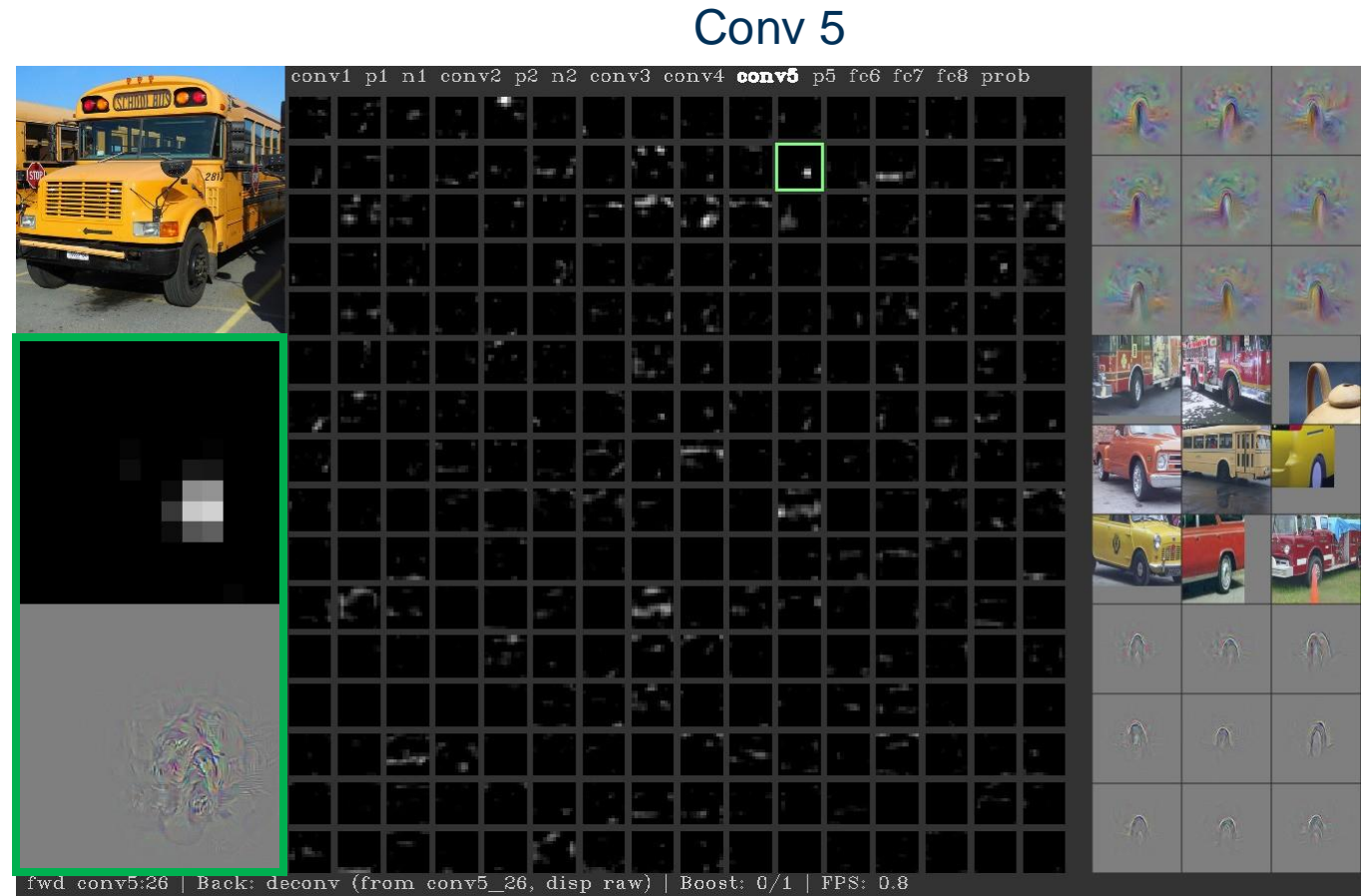
Visualizing CNNs

Visualizing Filters in Intermediate Layers

Intermediate layers:

- Weights: not very interpretable
- Activations: interpretable

The filter in conv 5 layer is activated when it sees a wheel

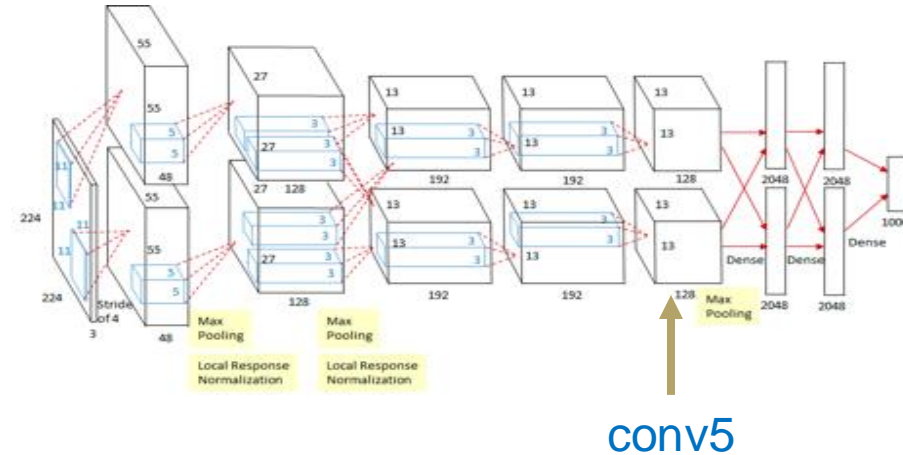


AlexNet

Visualizing CNNs

Visualizing Filters in Intermediate Layers

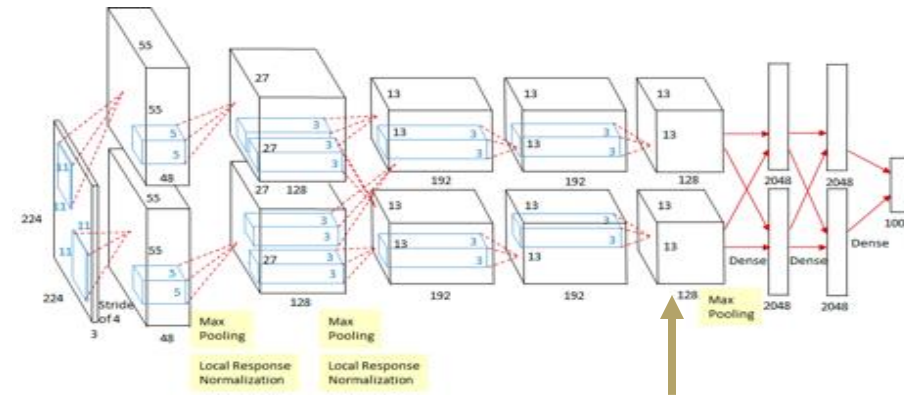
- Apart from visualizing the intermediate activations, we can also visualize what **similar visual patterns** in images that cause the maximum activations of certain neurons
- Maximally Activating Patches:
 - Image patches in the input that cause the maximum activations of certain filters



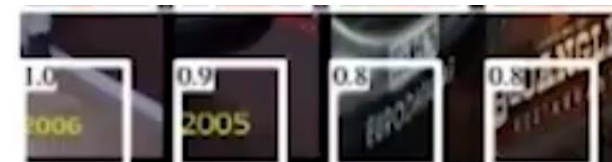
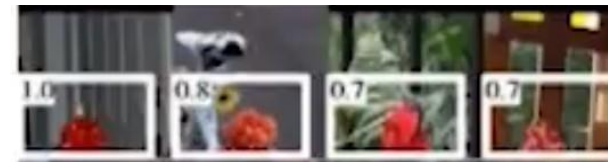
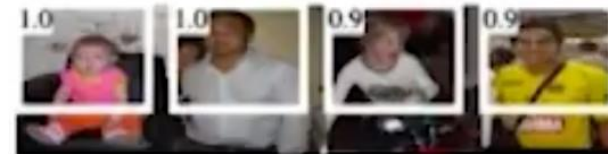
Visualizing CNNs

Visualizing Maximally Activating Patches

- Maximally Activating Patches:
 - Image patches in the input that cause the maximum activations of certain filters
- Obtaining Maximally Activating Patches:
 - Pick activations in a layer, e.g., conv5 (128 x 13 x 13), pick one of the channels 17/128
 - Feed forward many images through the network, record values of the chosen channel
 - Visualize image patches that correspond to **maximal activation**
- Maximally activating patches share **similar visual patterns**



conv5

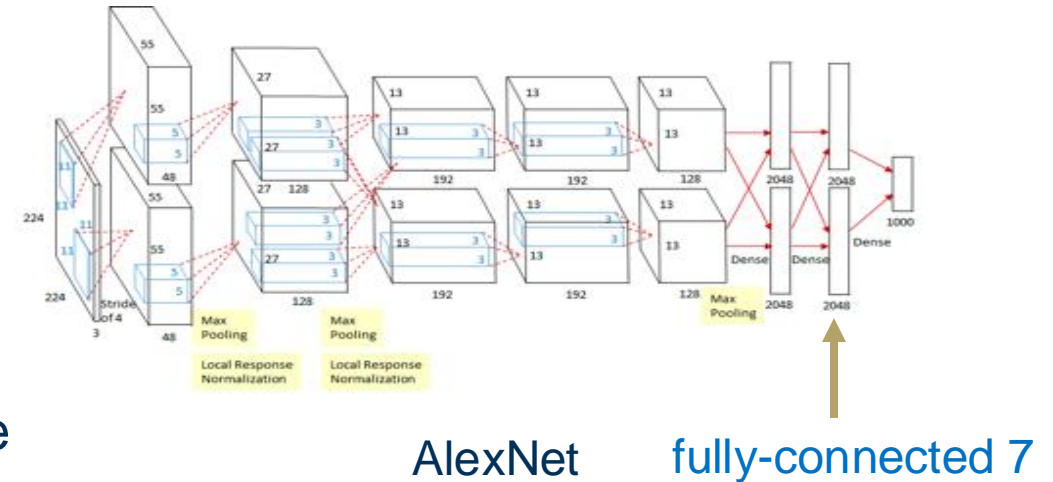


Maximally activating patches
Each row corresponds to a particular neuron in conv5

Visualizing CNNs

Visualizing Last Layer Activations

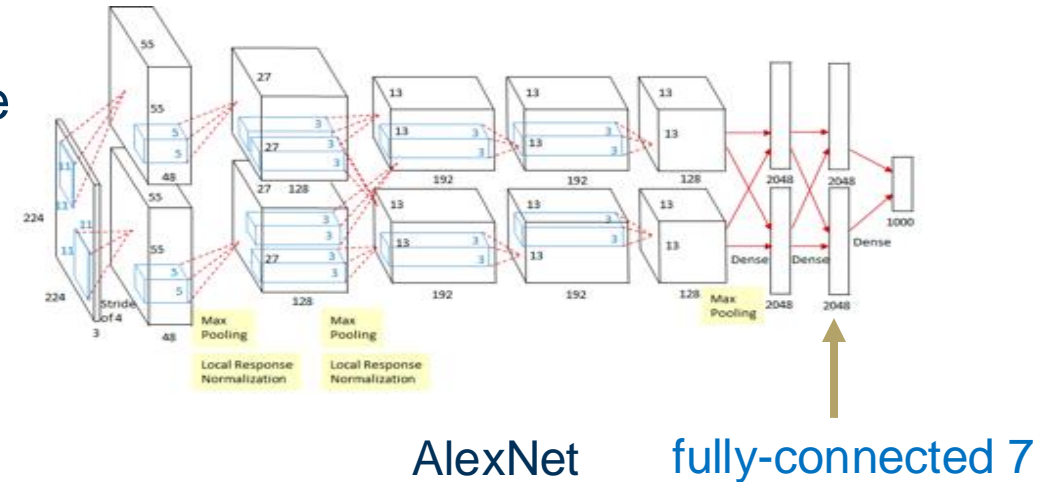
- We can group the images that have similar class-specific information by exploring last layer activations
- Last layer activations (embedding):
 - 4096-dimensional feature vector for an image (layer immediately before the classifier)
 - Representations of *entire input images* instead of specific patches
 - Similar embeddings correspond to same classes of input images



Visualizing CNNs

Visualizing Last Layer Activations

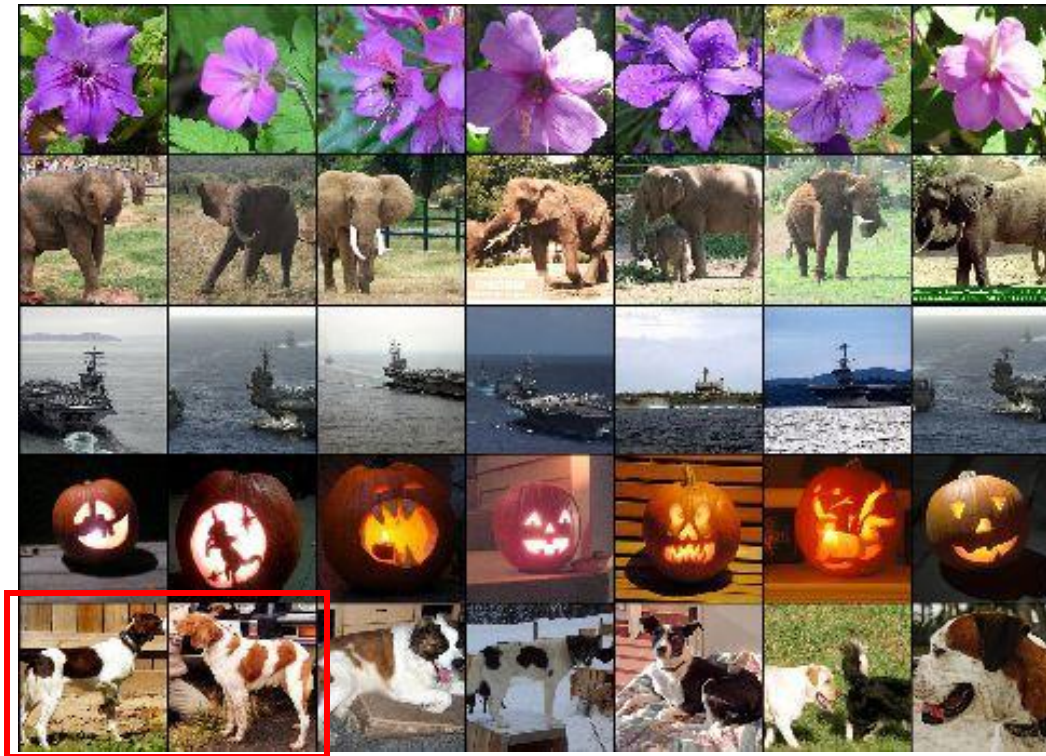
- Last layer activations (embedding):
 - 4096-dimensional feature vector for an image (layer immediately before the classifier)
 - Representations of *entire input images* instead of specific patches
 - Similar embeddings correspond to same classes of input images
- Feed forward many images through the network, collect the final layer feature vectors
- Visualize input images that have similar last layer embeddings



Visualizing Last Layer Activations

-

Test image L2 Nearest neighbors in feature space



- 16 of 53

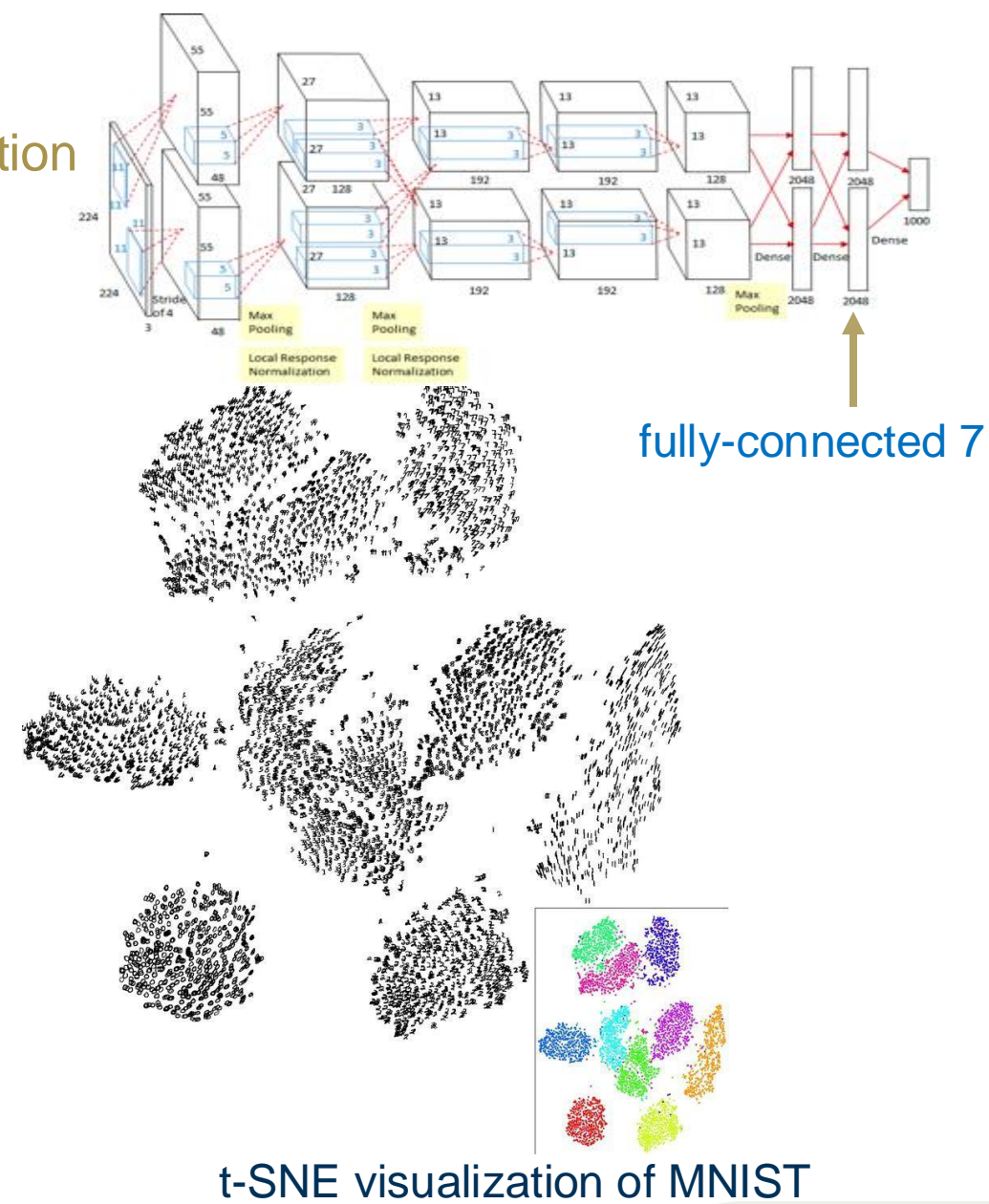
Visualizing Last Layer Activations via Dimensionality Reduction

-
- The diagram illustrates a deep convolutional neural network (CNN) architecture for handwritten digit recognition. The input is a 28x28x3 volume. The first layer is a convolution with 11x11 kernels, stride 4, resulting in 128 volumes of size 5x5x3. This is followed by two layers of max pooling and local response normalization, each resulting in 128 volumes of size 3x3x3. The next two layers are convolution with 3x3 kernels, resulting in 192 volumes of size 3x3x3. This is followed by two more layers of convolution with 3x3 kernels, resulting in 128 volumes of size 3x3x3. The final layer is a fully connected layer with 1000 units, indicated by an upward arrow.

Visualizing CNNs

Visualizing Last Layer Activations via Dimensionality Reduction

- Last layer embedding:
 - 4096-dimensional feature vector for an image
- Dimensionality reduction using t-SNE (t-distributed stochastic neighbor embedding):
- Embed **high-dimensional data** points so that **locally, pairwise distances are conserved** i.e., similar things end up in clusters, while dissimilar things end up wherever



Visualizing CNNs

Summary

We have been visualizing:

- Weights (filters) in conv layers
- Maximally activating patches
- Nearest neighbor images in features space
- Last layer embeddings

Next, we will look into how pixels affect model decisions



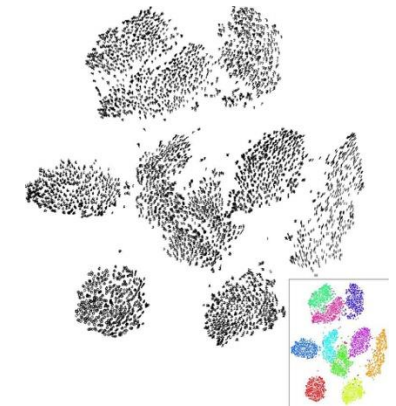
filters



Maximally activating patches



Nearest neighbor images

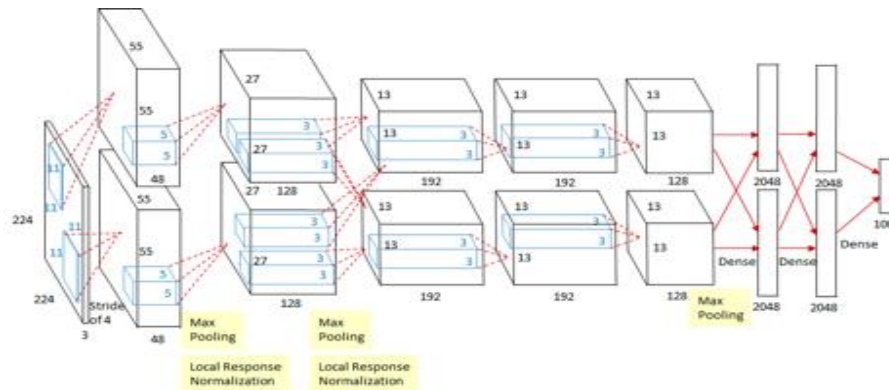
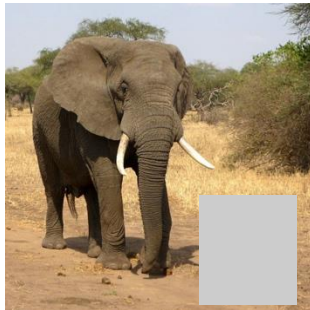


Last layer embeddings

Saliency Methods

Saliency via Occlusion

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



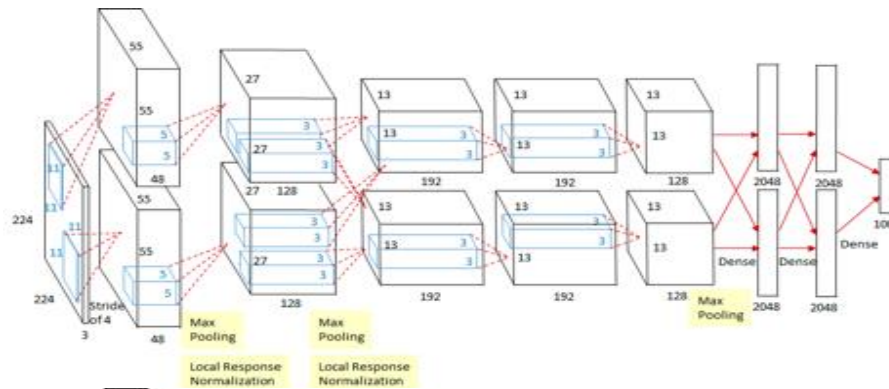
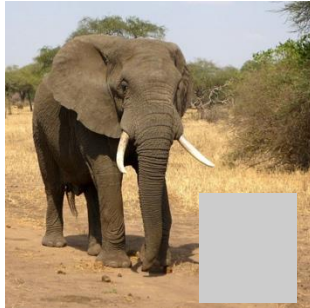
$$P(\text{elephant}) = 0.95$$

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.

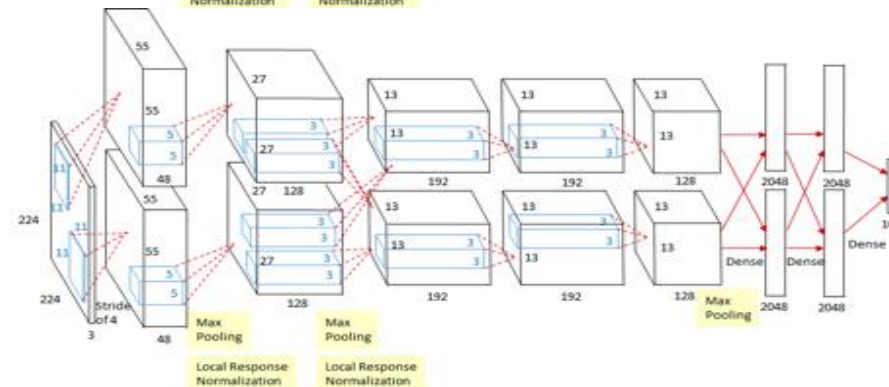
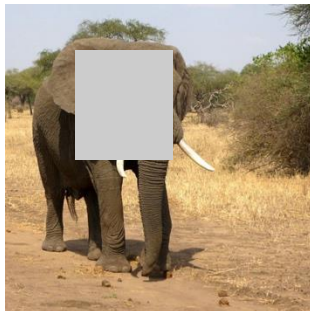
Saliency Methods

Saliency via Occlusion

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



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



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

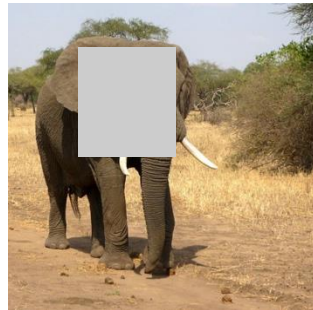
These pixels
affect decisions
more

Saliency Methods

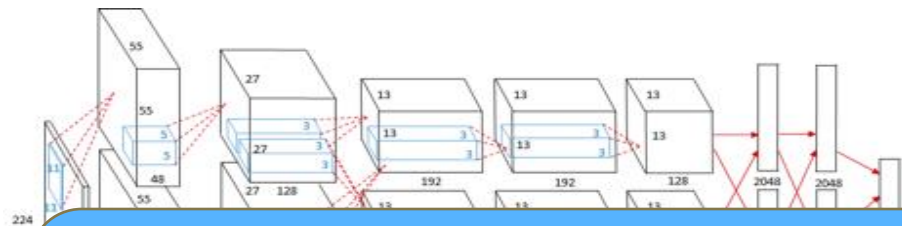
Saliency via Occlusion

- Visualizing the heatmap of pixels that cause decrease in probabilities when blocked

These pixels affect decisions more



go-kart



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



African elephant, *Loxodonta africana*



go-kart



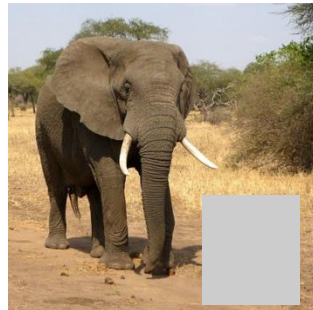
Saliency Methods

Saliency via Occlusion

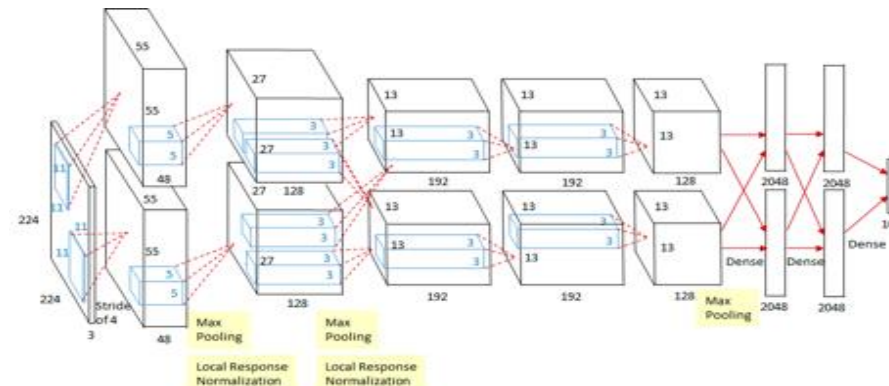
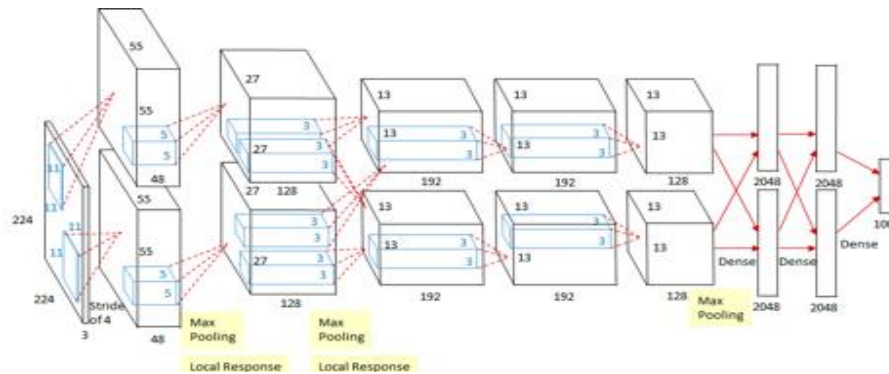
- Visualizing the heatmap of pixels that cause decrease in probabilities when blocked

Very expensive

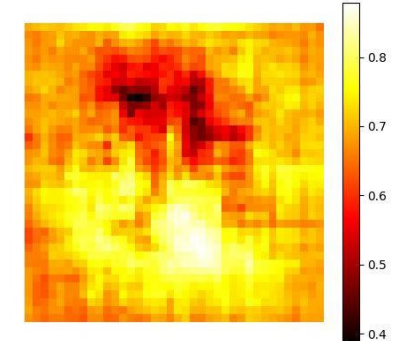
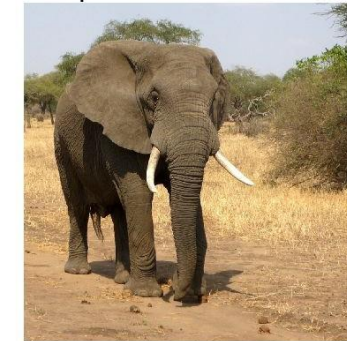
These pixels affect decisions more



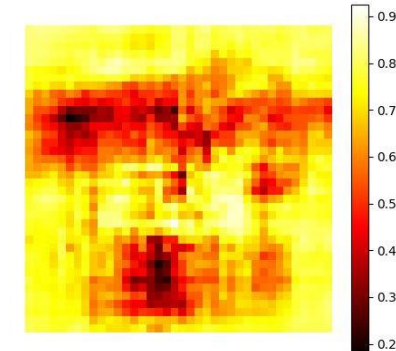
go-kart



Faithful
African elephant, *Loxodonta africana*



go-kart



Saliency Methods

Saliency via Occlusion

- Saliency via Occlusion:
 - a **very expensive** approach especially for high resolution input images
- We want to find a less expensive approach
- Recall **feature importance** in logistic regression:

$$P(y = 1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

$$P(y = 0|\mathbf{x}) = 1 - \sigma(\mathbf{w}^T \mathbf{x} + b)$$

- for each feature x_i , a **weight** w_i represents its **importance**
- We want to generate pixel saliency maps by deep models as feature importance maps

Saliency Methods

Saliency via Gradients

- Generate pixel saliency maps by deep models as feature importance maps

- *Highly non-linear mapping function*
 $f_{\theta}: \mathcal{X} \rightarrow \mathcal{Y}$:

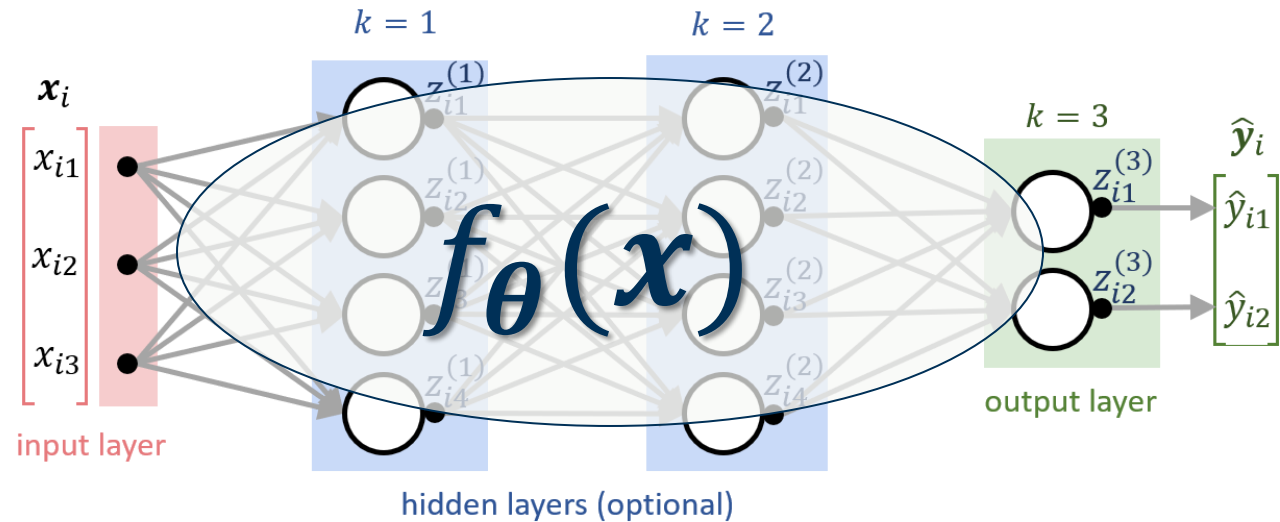
$$\hat{Y} = \varphi \left(\varphi \left(\varphi \left(X(W^{(1)})^T + o(b^{(1)})^T \right) (W^{(2)})^T + o(b^{(2)})^T \right) (W^{(3)})^T + o(b^{(3)})^T \right)$$

- Assume that we can ‘linearize’ the model using Taylor series

$$\hat{Y} \approx X(W)^T + o(b)^T$$

$$W \approx \frac{\partial \hat{Y}}{\partial X}$$

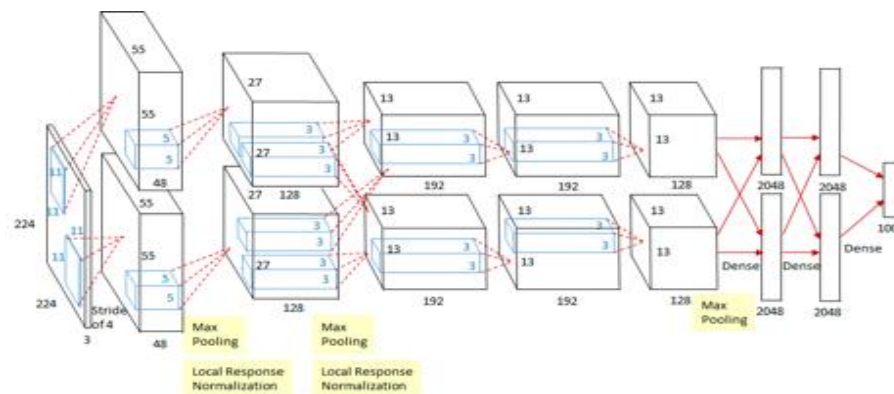
saliency approximated by **gradients w.r.t. input**, which can be **obtained via backpropagation**



Saliency Methods

Saliency via Gradients

Forward pass: Compute probabilities

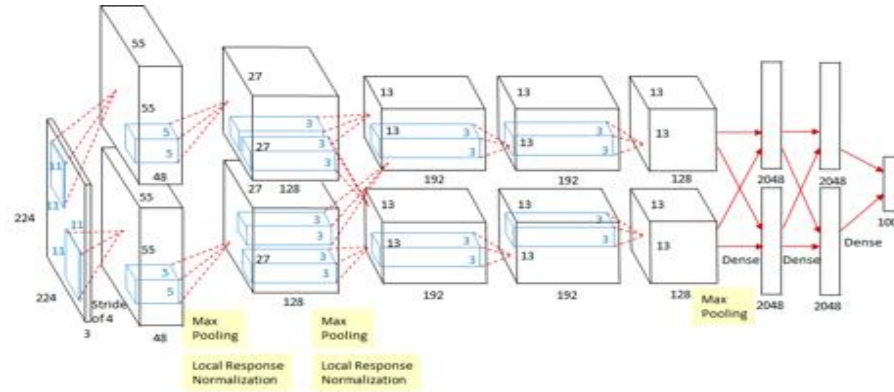


$$\hat{y} = \begin{bmatrix} 0.05 \\ 0.10 \\ 0.85 \end{bmatrix} \text{ Dog}$$

Saliency Methods

Saliency via Gradients

Forward pass: Compute probabilities

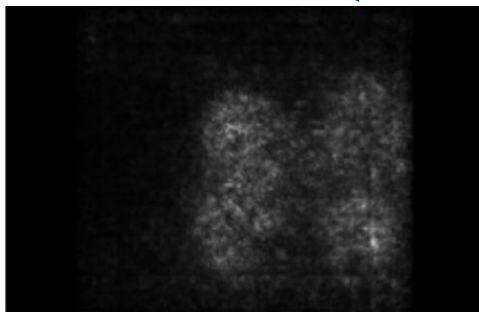


$$\hat{y} = \begin{bmatrix} 0.05 \\ 0.10 \\ 0.85 \end{bmatrix}$$

$\hat{y}_c = 0.85$ (dog)

Backward pass: Compute gradients

$$\frac{\partial \hat{y}_c}{\partial X}$$



Saliency map

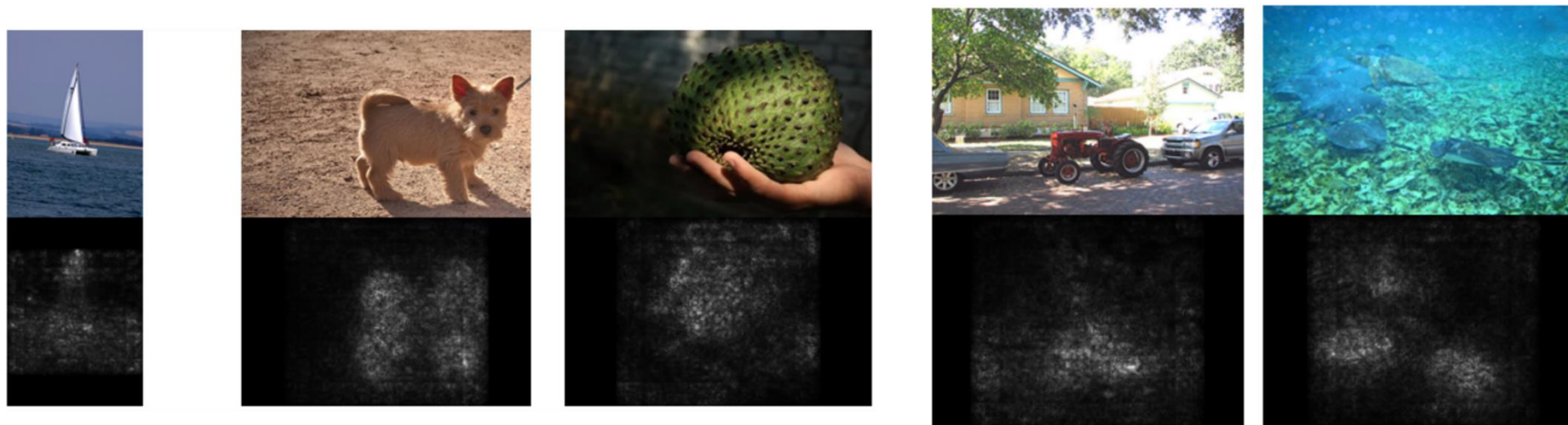
Compute **gradient** of (unnormalized) class score
with respect to image pixels

Then visualize the max of absolute value over RGB
channels

Saliency Methods

Saliency via Gradients

- Examples of saliency maps via backpropagation

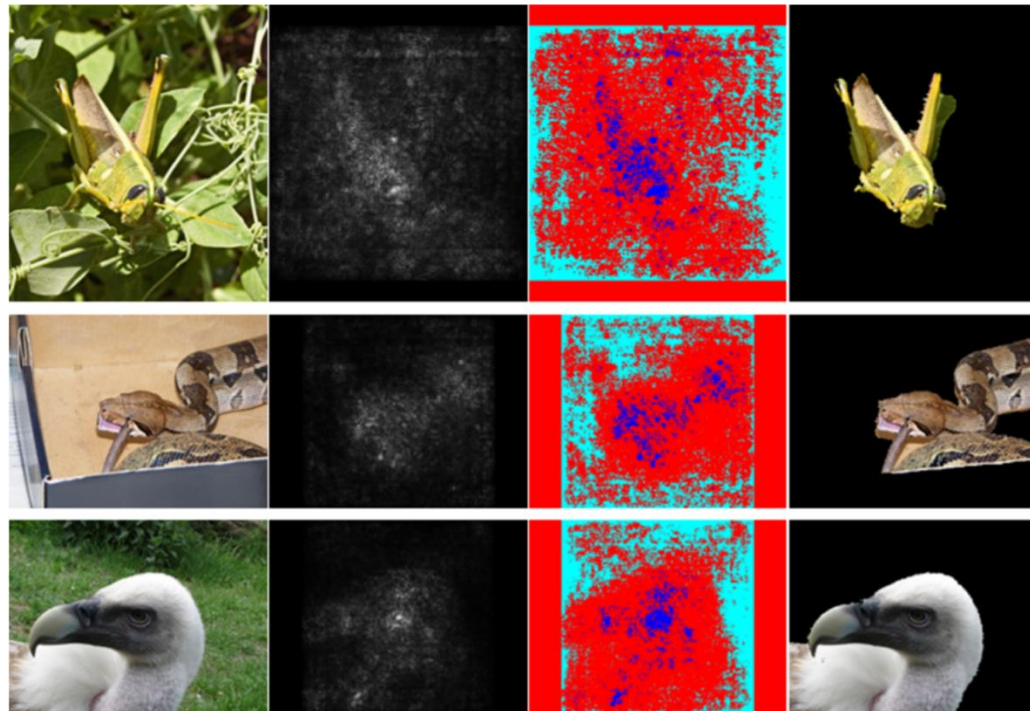


Saliency Methods

Weakly-supervised Segmentation

- Saliency maps can be used to help unsupervised semantic segmentation

Gradients w.r.t. pixels



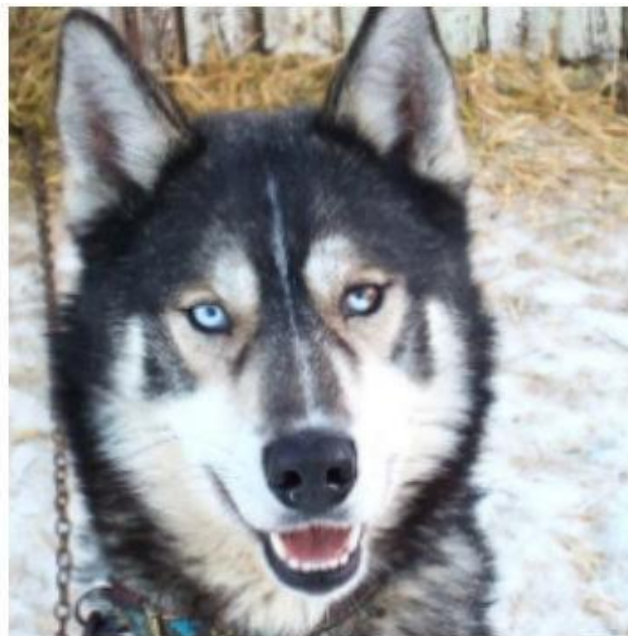
Saliency Methods

Uncovering Biases

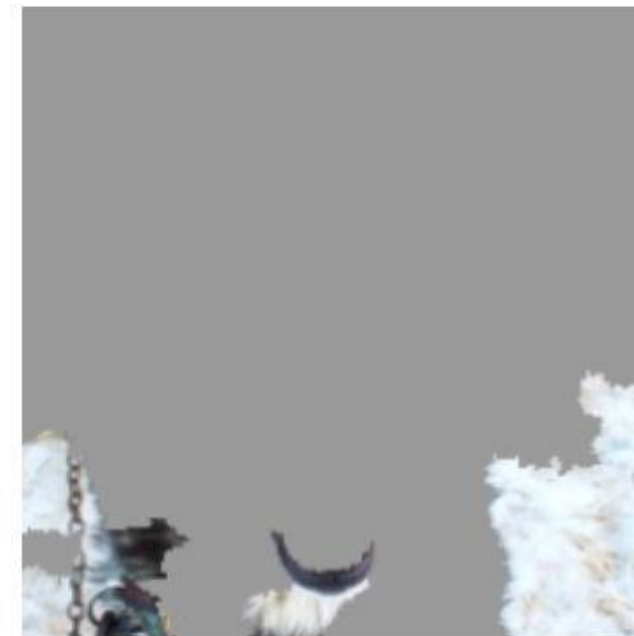
- Saliency Maps also find biases

When snow is presented in most of wolf images, network may use these snow pixels as salient regions for prediction

Wolf vs. dog classifier is actually a snow vs. no-snow classifier



(a) Husky classified as wolf

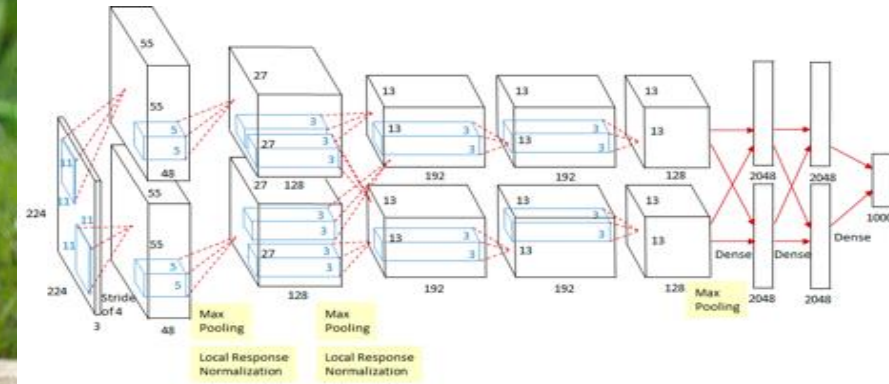


(b) Explanation
snow pixels as salient regions

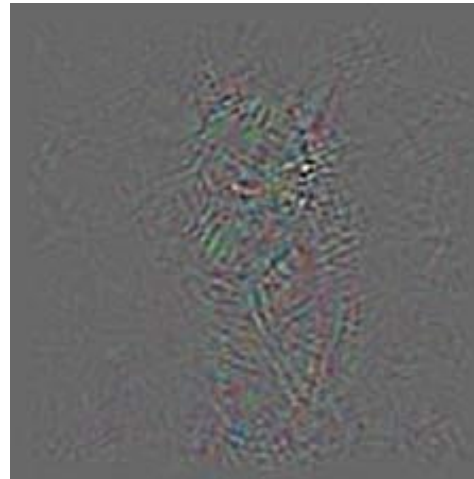
Gradient-based Explanations

Vanilla Backpropagation

- Saliency map by vanilla backprop can be noisy
- We want cleaner saliency maps:
 - *Modifying (rectifying) the backpropagation* empirically produce better visualizations



Backward pass: Compute gradients



Saliency map by vanilla backprop

Gradient-based Explanations

Vanilla Backpropagation

Going backward performing all the operations of the network (Unpooling, Filtering...), and for ReLU non-linearities, **only pass gradients to regions of positive activations**

$$h^{l+1} = \max\{0, h^l\}$$

Forward pass

1	-1	5
2	-5	-7
-3	2	4

h^l

→

1	0	5
2	0	0
0	2	4

h^{l+1}

$$\frac{\partial L}{\partial h^l} = \llbracket h^l > 0 \rrbracket \frac{\partial L}{\partial h^{l+1}}$$

Backward pass: backpropagation

positive activations in the previous layer

-2	0	-1
6	0	0
0	-1	3

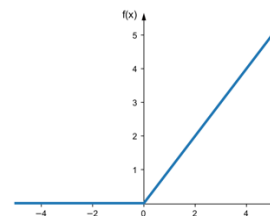
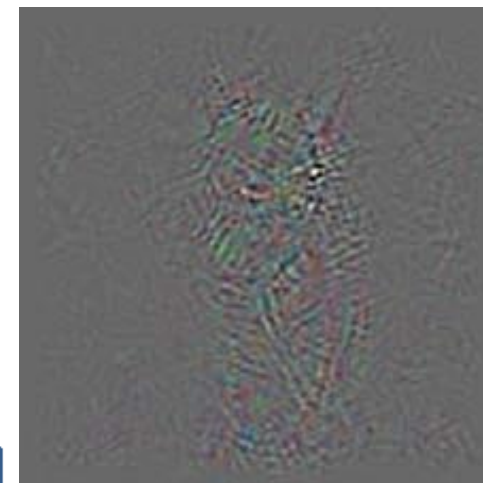
←

-2	3	-1
6	-3	1
2	-1	3

$\frac{\partial L}{\partial h^{l+1}}$

Gradients from the later layer

Saliency map by vanilla backprop

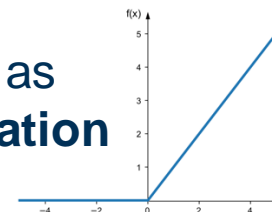


Gradient-based Explanations

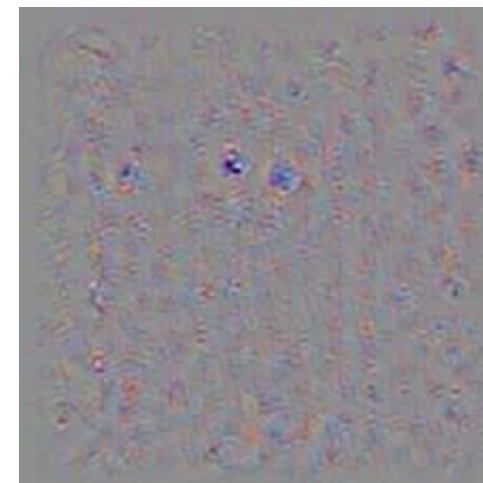
Deconvnet Backpropagation

- The way DeconvNet Backpropagation handle the ReLU non-linearities is different as they propose to **only propagate positive gradient**
- *Rectifying the backpropagation* empirically produce better saliency visualizations

We can think of **Deconvnet** as **rectified gradients propagation**



Saliency map by
Deconv backprop



Cleaner saliency map

$$\frac{\partial L}{\partial h^l} = \left[\left[\frac{\partial L}{\partial h^{l+1}} > 0 \right] \right] \frac{\partial L}{\partial h^{l+1}}$$

Backward pass:
"deconvnet"

positive gradient in the **later layer**

0	3	0
6	0	1
2	0	3



-2	3	-1
6	-3	1
2	-1	3

$\frac{\partial L}{\partial h^{l+1}}$

Gradients from the
later layer

Gradient-based Explanations

Guided Backpropagation

- Guided backpropagation propose to **propagate positive gradient and rectified by positive activations**
- **Non-intuitive** approach but **empirically** produce better saliency visualizations

$$h^{l+1} = \max\{0, h^l\} \quad \text{Forward pass}$$

$$\frac{\partial L}{\partial h^l} = \mathbb{I}[h^l > 0] \mathbb{I}\left[\frac{\partial L}{\partial h^{l+1}} > 0\right] \frac{\partial L}{\partial h^{l+1}} \quad \text{Backward pass: guided backpropagation}$$

positive activations in the previous layer **positive gradient** in the later layer

1	-1	5
2	-5	-7
-3	2	4

0	0	0
6	0	0
0	0	3



1	0	5
2	0	0
0	2	4



-2	3	-1
6	-3	1
2	-1	3

Saliency map by Guided Backpropagation



Cleaner saliency map

h^{l+1}

$\frac{\partial L}{\partial h^{l+1}}$ Gradients from the later layer

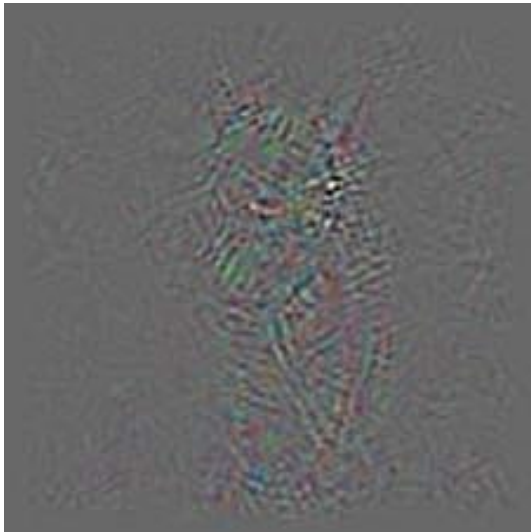
Gradient-based Explanations

Vanilla vs Deconvnet vs Guided Backpropagation

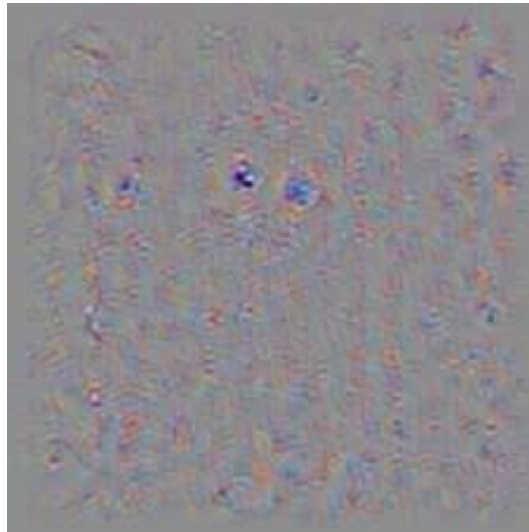
- Guided Backprop tends to be “cleanest”



Backprop



Deconv



Guided Backprop



Gradient-based Explanations

Problems with Guided Backpropagation

- Not very “class-discriminative”

GB for “airliner”



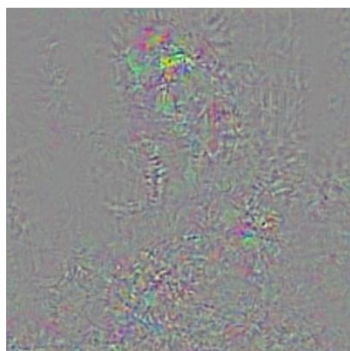
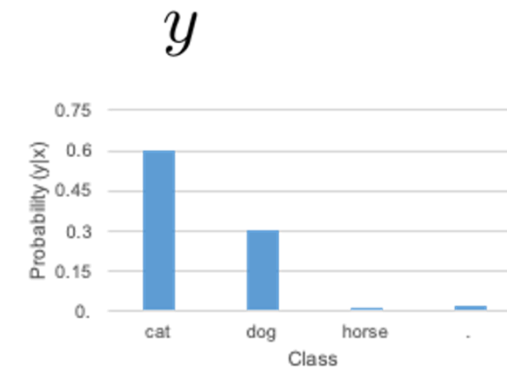
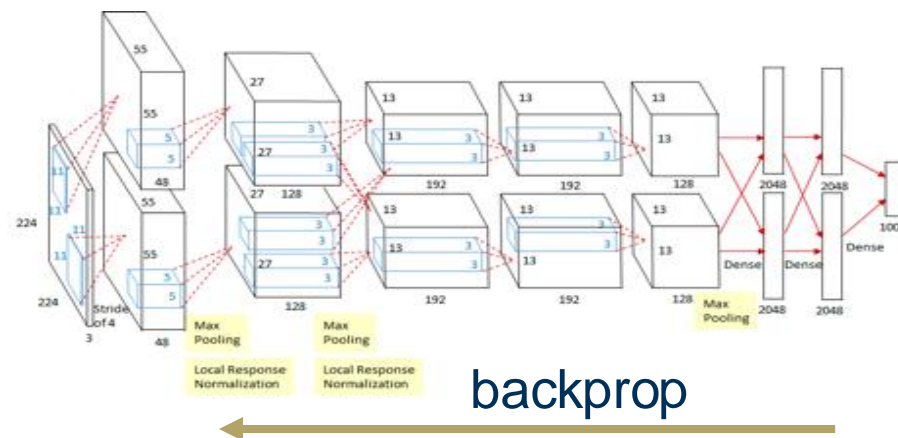
GB for “bus”



- Less related to the decision-making of neural networks

Gradient-based Explanations

Problems with Guided Backpropagation



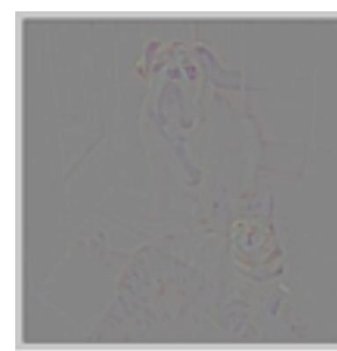
Backprop for
'cat'



Backprop for
'dog'

Noisy

$$w_c = \frac{\partial y_c}{\partial x} \Big|_{x=x_0}$$



Guided Backprop for
'cat'



Guided Backprop for
'dog'

Not Class-Discriminative

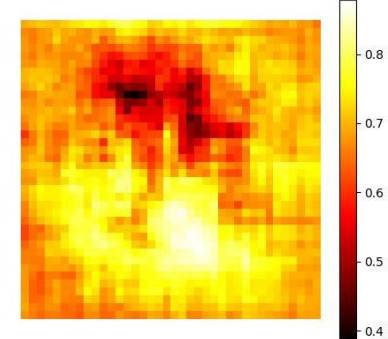
Gradient-based Explanations

Summary

Saliency maps:

- Intervention-based:
 - perturbing pixels and see how the decision change
 - expensive
- Gradient-based:
 - approximates feature importance by backpropagation
 - computational efficient
 - not reflect decision making

African elephant, *Loxodonta africana*



Expensive



Guided Backprop for
'cat'



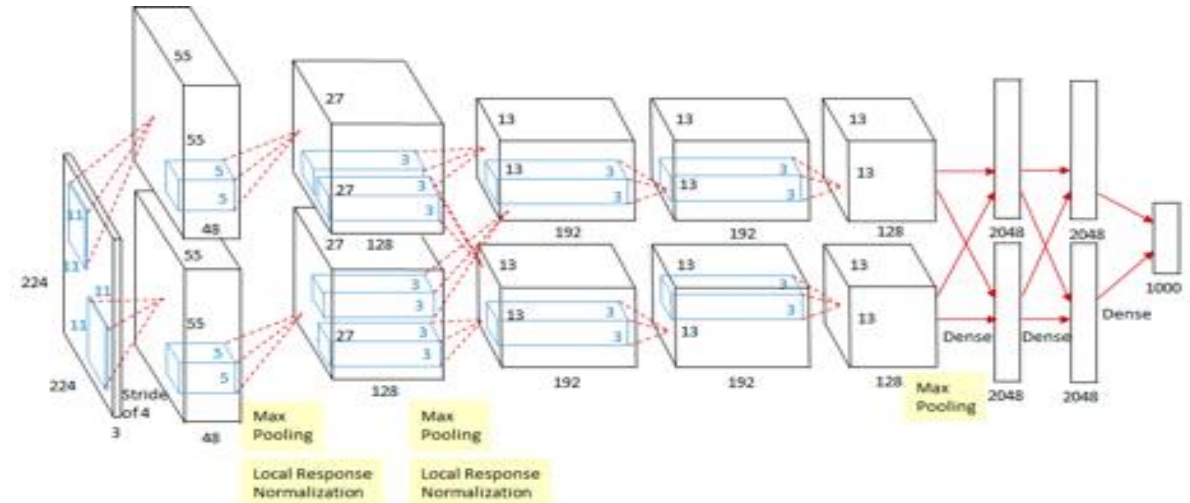
Guided Backprop for
'dog'

Not Class-Discriminative

Gradient-based Class Activation Map (Grad-CAM)

Motivation

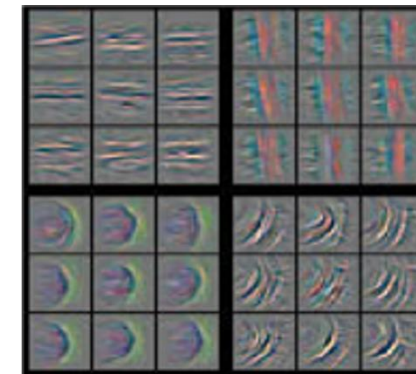
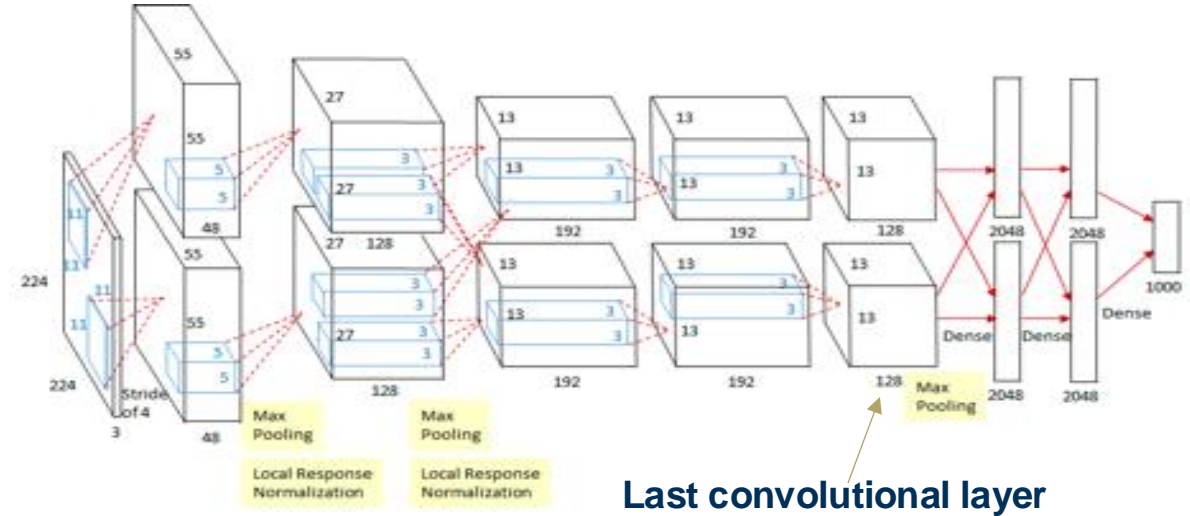
- To find the important activations that are responsible for a particular class
- We want the activations:
 - **Class-discriminative** to reflect decision-making
 - **Preserve spatial information** to ensure spatial coverage of important regions



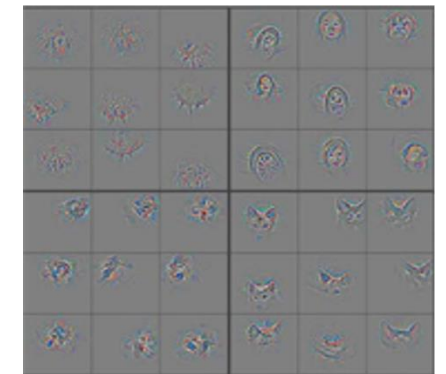
Gradient-based Class Activation Map (Grad-CAM)

Motivation

- To find the important activations that are responsible for a particular class
- We want the activations:
 - Class-discriminative
 - Preserve spatial information
- Which layer to perturb:
 - **Higher layer** capture **class specific information**
 - **Spatial information is lost** in **fully-connected** layers
 - **Last convolutional layer** forms a best compromise between high-level semantics and detailed spatial resolution



Low-level features captured by lower-layers



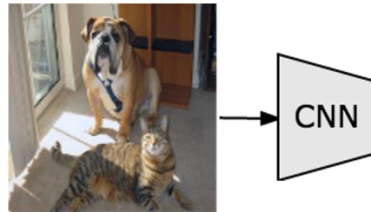
High-level features captured by higher-layers

Gradient-based Class Activation Map (Grad-CAM)

Method

- Given an image, feed forward through CNN

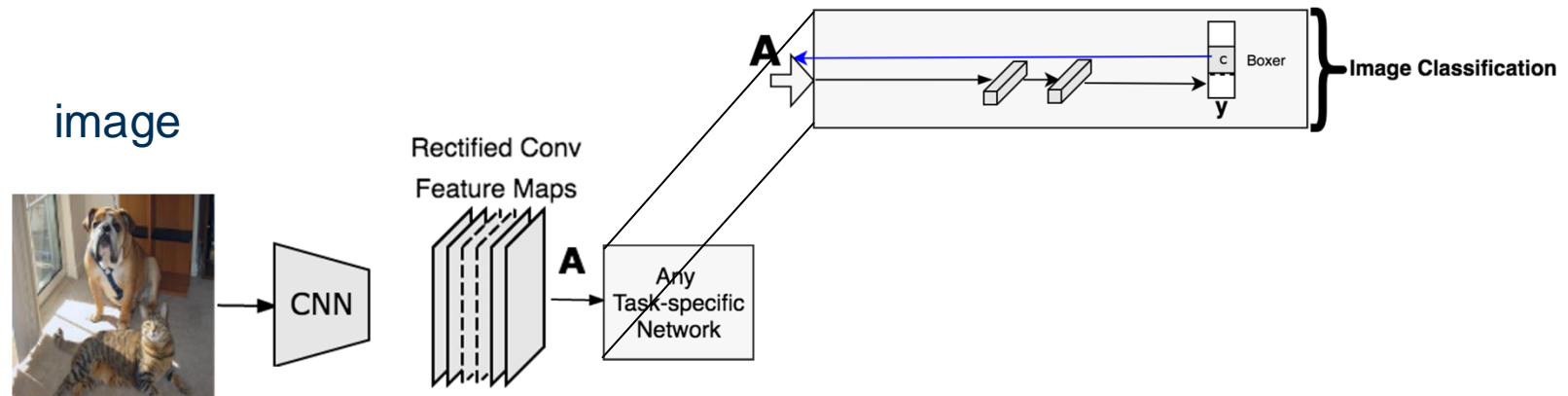
image



Gradient-based Class Activation Map (Grad-CAM)

Method

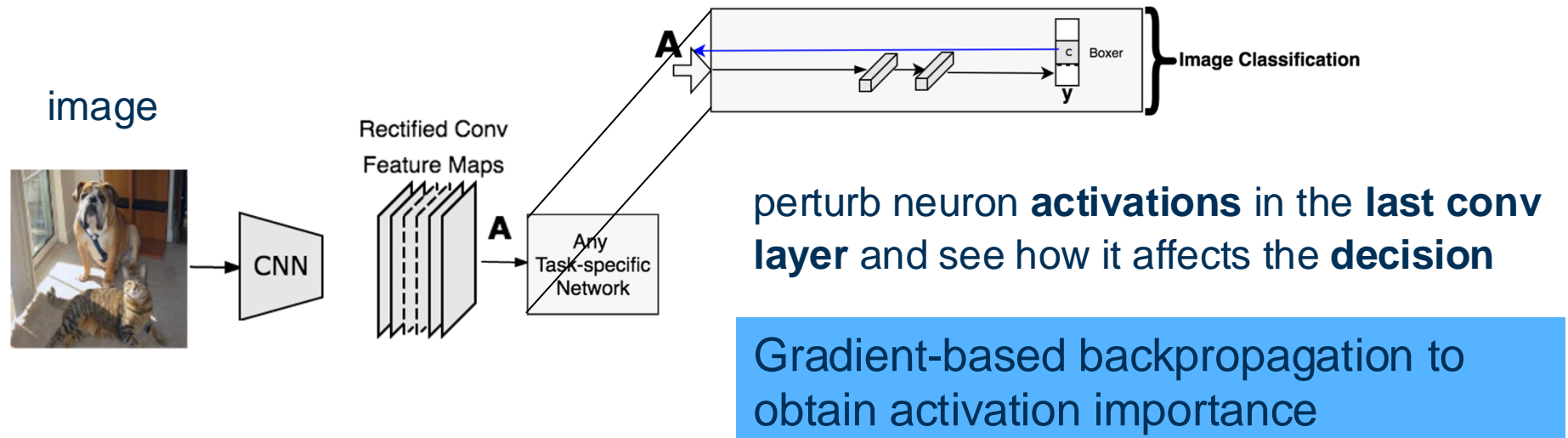
- Given an image, feed forward through CNN
- Final convolutional layer output feature maps for later task-specific layers, i.e., fc layer for classification



Gradient-based Class Activation Map (Grad-CAM)

Method

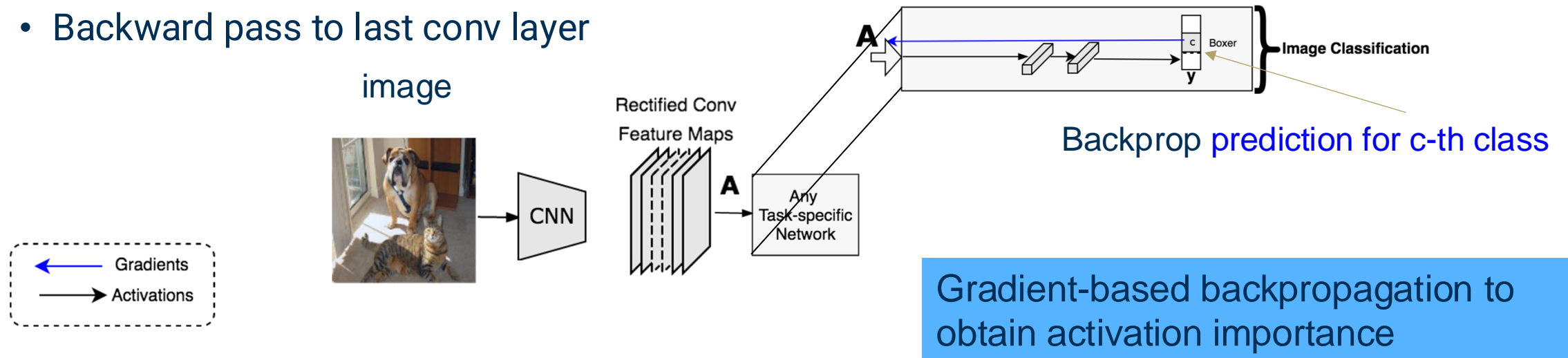
- Given an image, feed forward through CNN
- Final convolutional layer output feature maps for later task-specific layers, i.e., fc layer for classification



Gradient-based Class Activation Map (Grad-CAM)

Method

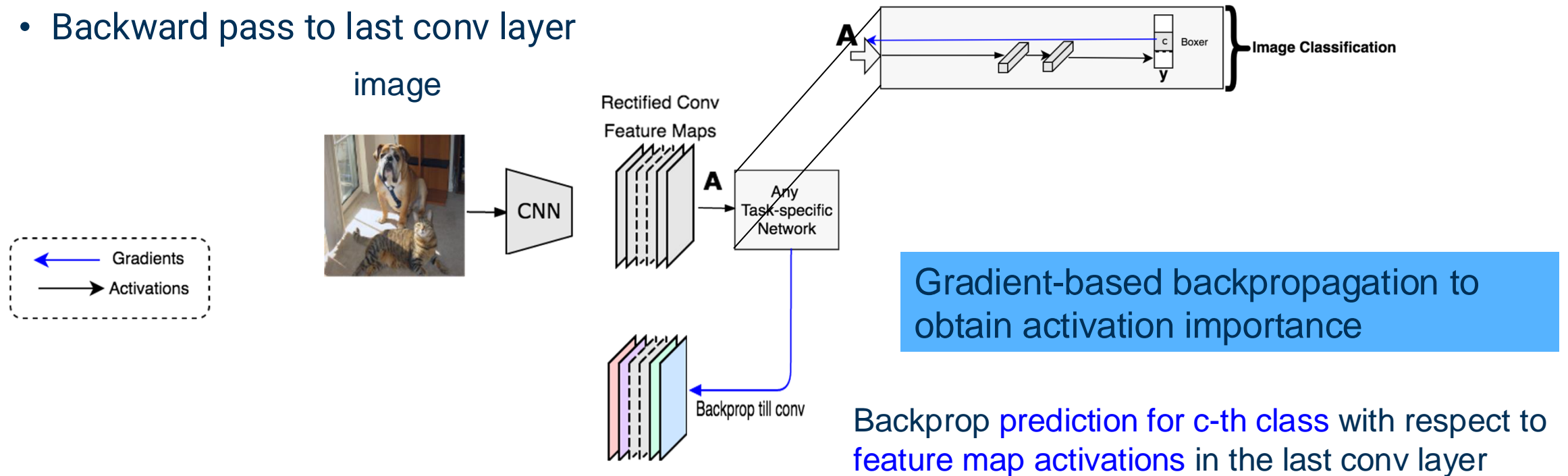
- Given an image, feed forward through CNN
- Final convolutional layer output feature maps for later task-specific layers, i.e., fc layer for classification
- Backward pass to last conv layer



Gradient-based Class Activation Map (Grad-CAM)

Method

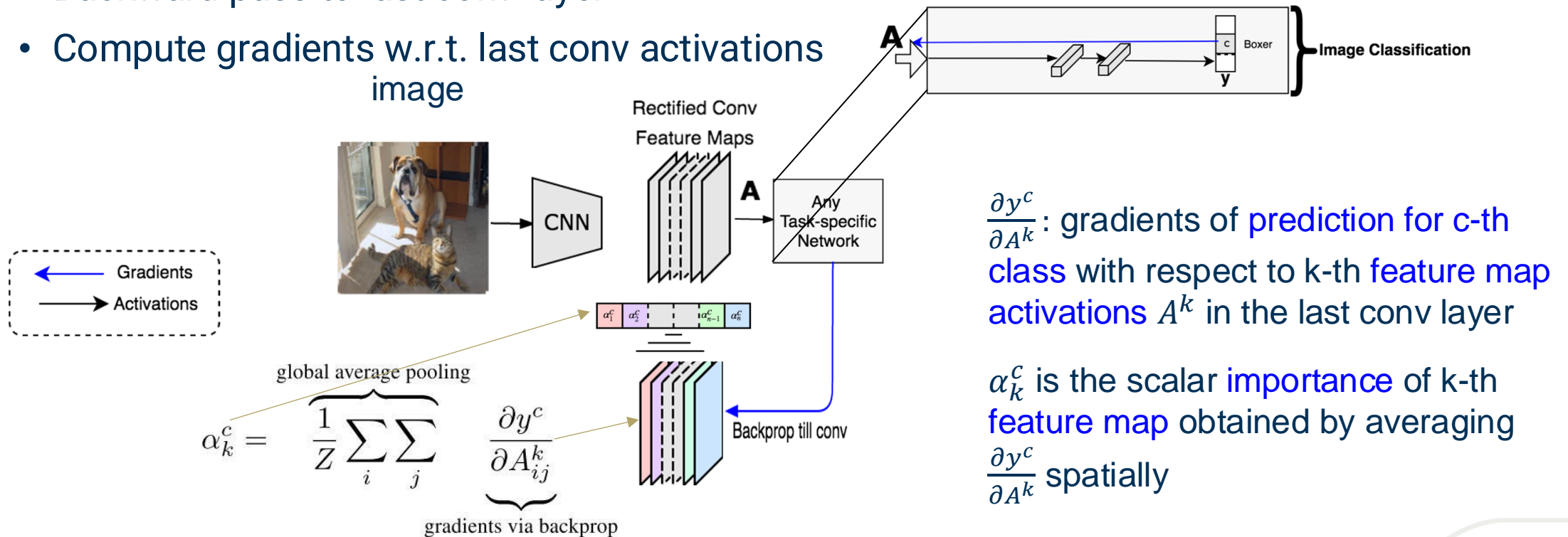
- Given an image, feed forward through CNN
- Final convolutional layer output feature maps for later task-specific layers, i.e., fc layer for classification
- Backward pass to last conv layer



Gradient-based Class Activation Map (Grad-CAM)

Method

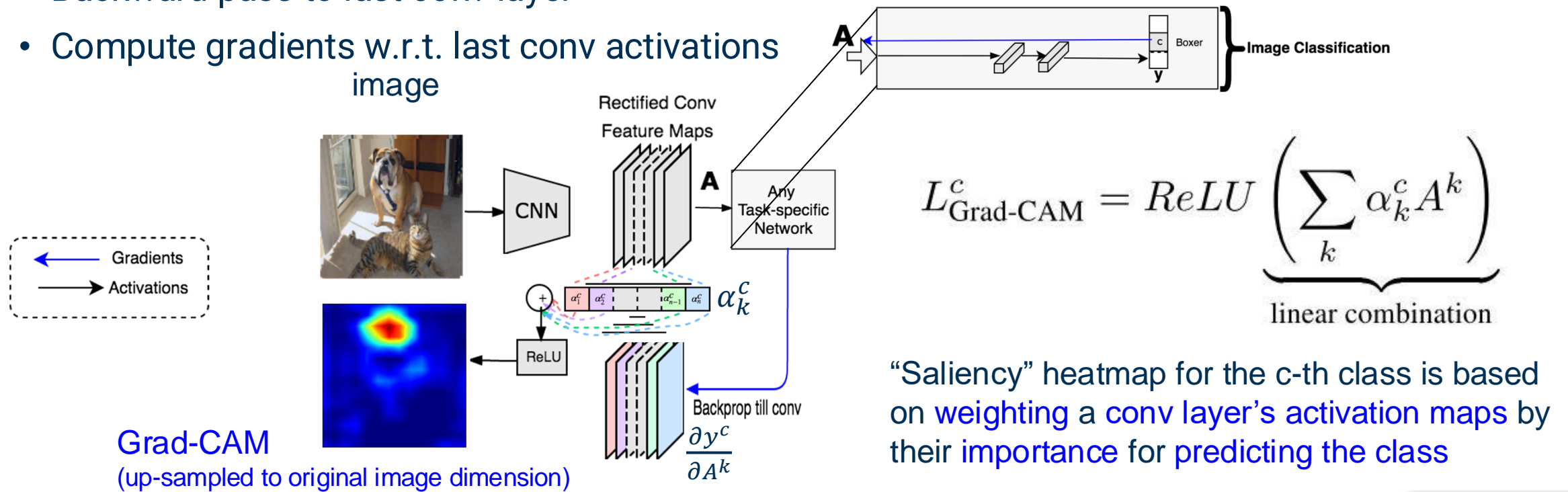
- Given an image, feed forward through CNN
- Final convolutional layer output feature maps for later task-specific layers
- Backward pass to last conv layer
- Compute gradients w.r.t. last conv activations image



Gradient-based Class Activation Map (Grad-CAM)

Method

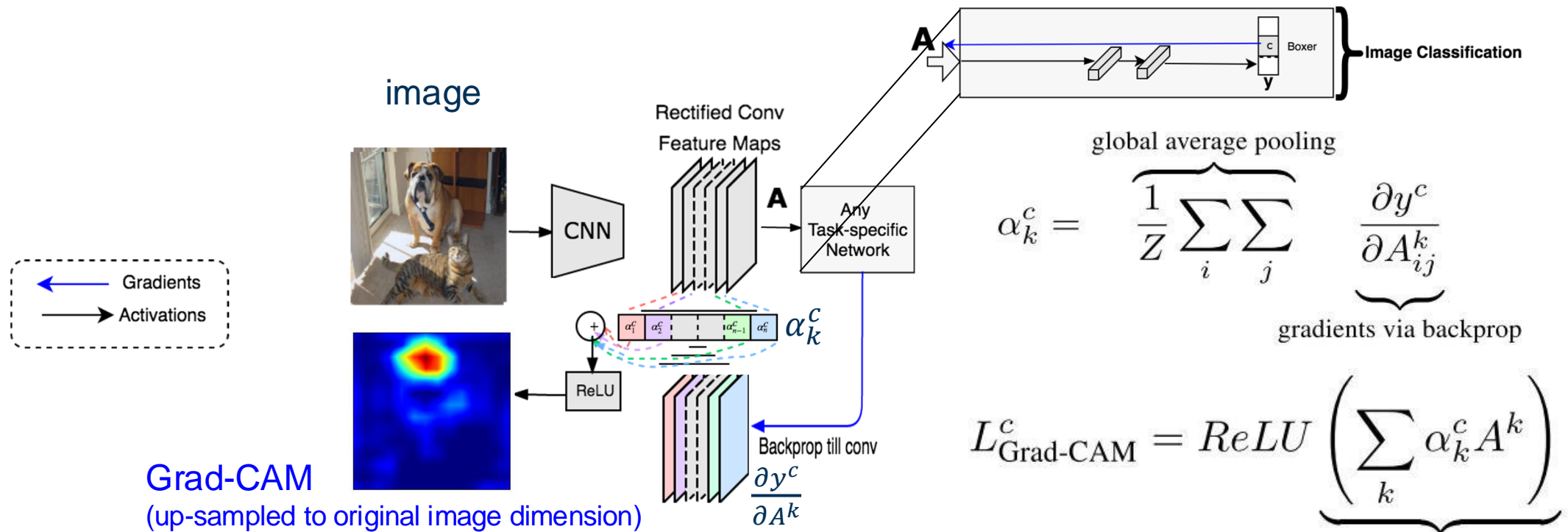
- Given an image, feed forward through CNN
- Final convolutional layer output feature maps for later task-specific layers
- Backward pass to last conv layer
- Compute gradients w.r.t. last conv activations image



Gradient-based Class Activation Map (Grad-CAM)

Method

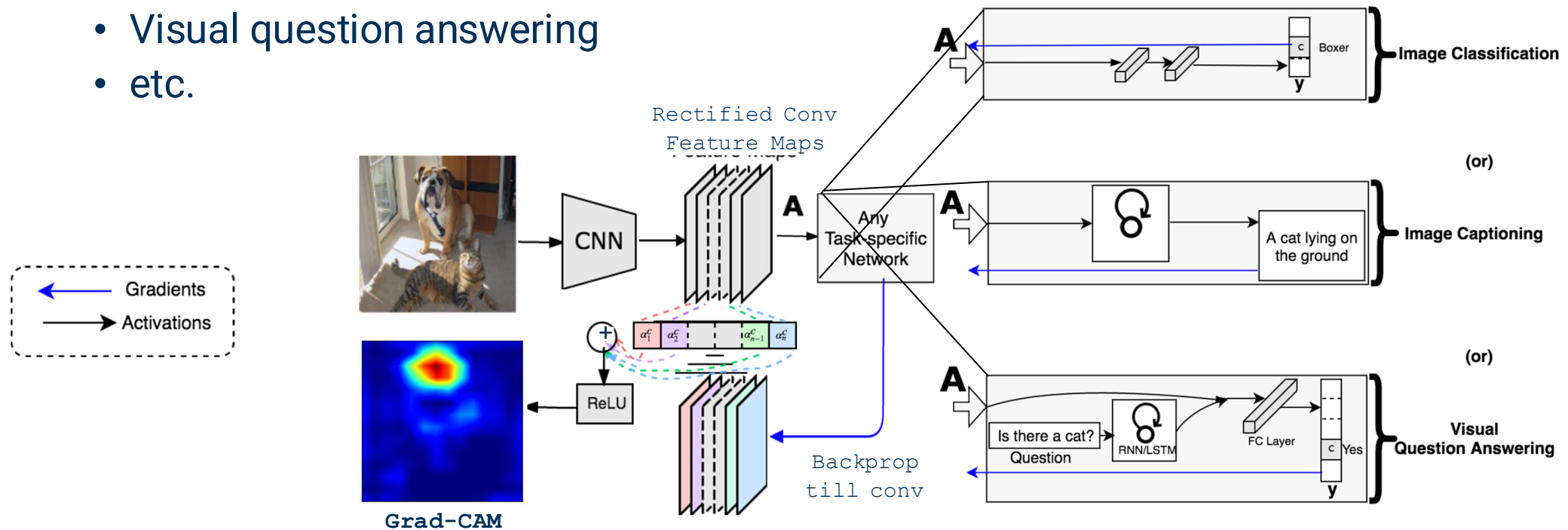
Grad-CAM uses the **gradient** information flowing into the **last convolutional layer** of the CNN to assign **importance values** to each activation for a **particular decision of interest**.



Gradient-based Class Activation Map (Grad-CAM)

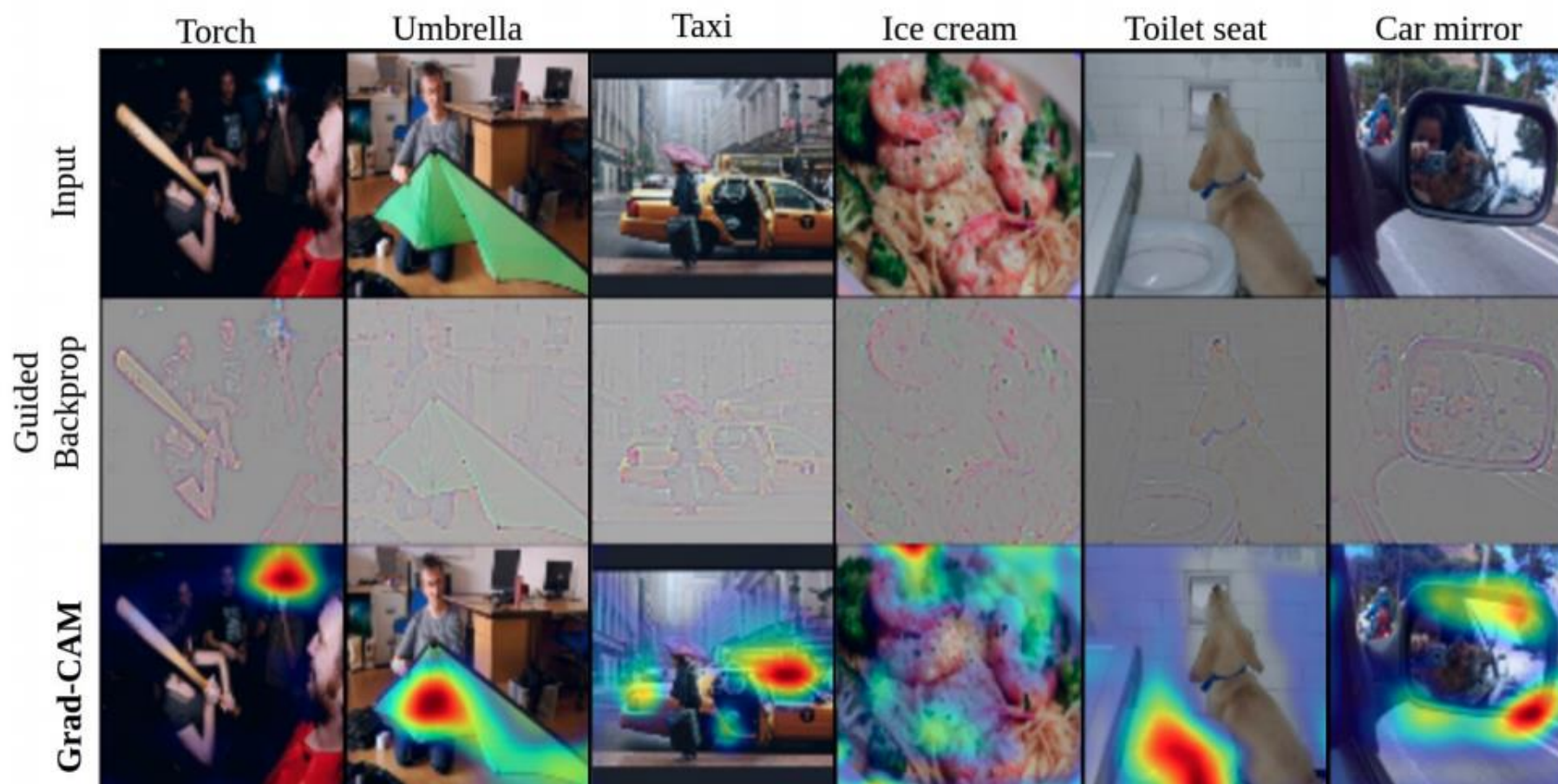
Method

- Grad-CAM generalizes to different tasks:
 - Image classification
 - Image captioning
 - Visual question answering
 - etc.



Gradient-based Class Activation Map (Grad-CAM)

Method



Gradient-based Class Activation Map (Grad-CAM)

Method

Guided Backprop

Grad-CAM



A bathroom with a toilet and a sink

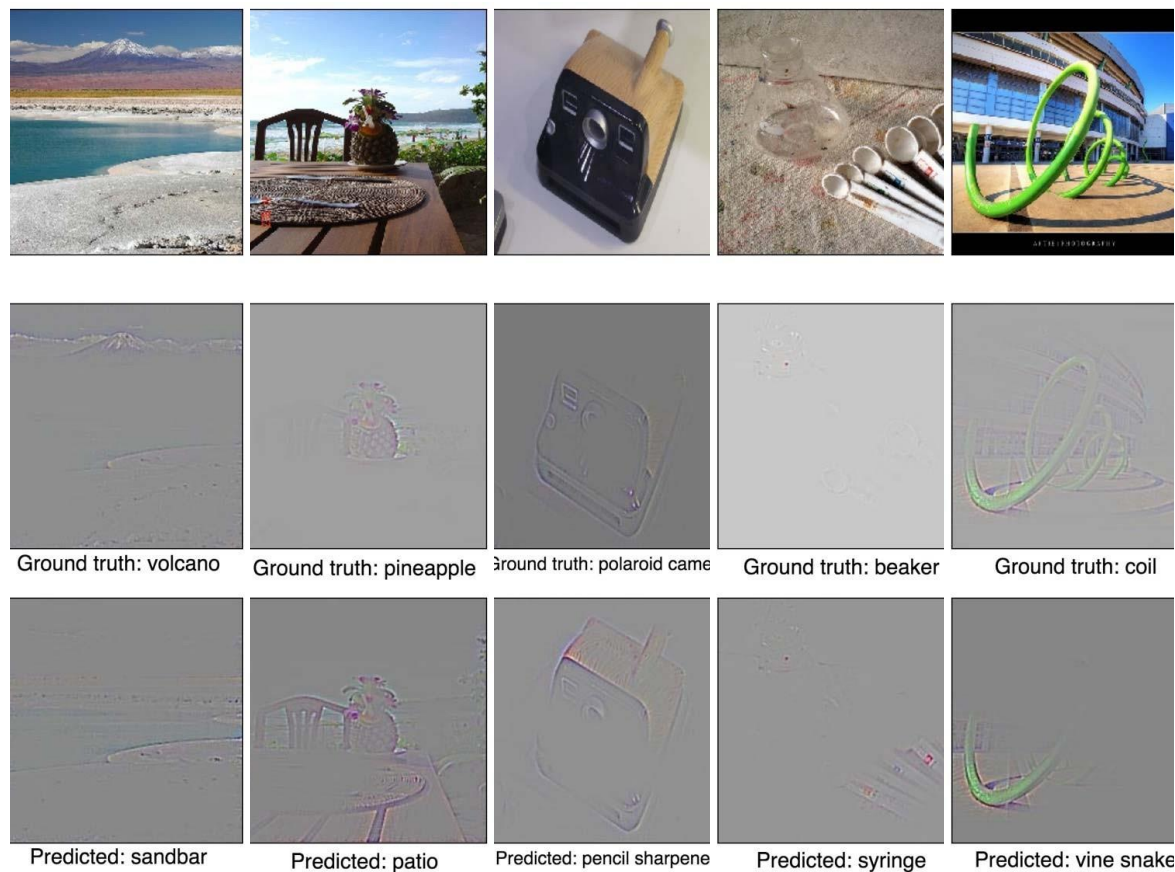


A horse is standing in a field with a fence in the background

Gradient-based Class Activation Map (Grad-CAM)

Method

- Even unreasonable predictions seem justifiable



Appendix

Notations

- x_i : a single feature
- \mathbf{x}_i : feature vector (a data sample)
- $\mathbf{x}_{:,i}$: feature vector of all data samples
- \mathbf{X} : matrix of feature vectors (dataset)
- N : number of data samples
- \mathbf{W} : weight matrix
- \mathbf{b} : bias vector
- $\mathbf{v}(t)$: first moment at time t
- $\mathbf{G}(t)$: second moment at time t
- $\mathbf{H}(\boldsymbol{\theta})$: Hessian matrix
- P : number of features in a feature vector
- α : learning rate
- Bold letter/symbol: vector
- Bold capital letters/symbol: matrix