

Visualising convolutional networks

Hakan Bilen

Machine Learning Practical - MLP Lecture 12

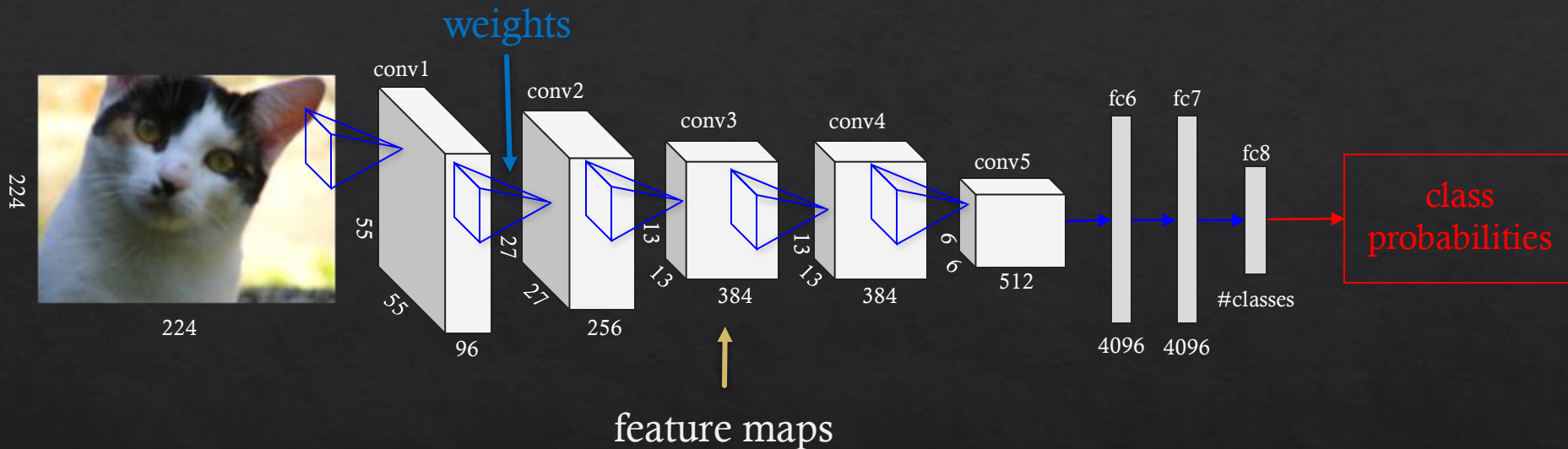
16 January 2019

<http://www.inf.ed.ac.uk/teaching/courses/mlp/>

Lectures in second semester

- ◇ Understanding convolutional networks
- ◇ Generative adversarial networks
- ◇ Domain adaptation and transfer learning
- ◇ Convolutional network design and compression (Dr Elliot Crowley)
- ◇ Object detection and semantic segmentation
- ◇ Language and vision models
- ◇ Video analytics

Recap: Convolutional Neural Networks (CNNs)



What is inside the black box (filters and feature maps)?

Why does it matter?

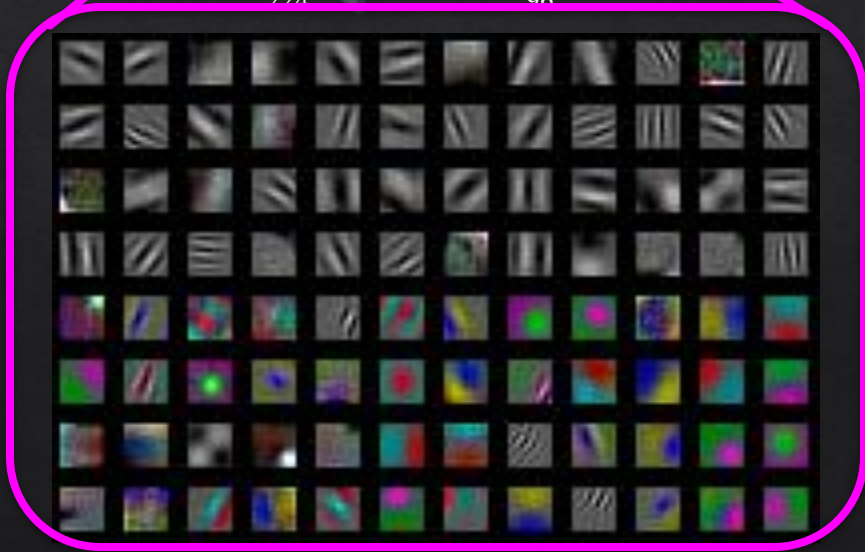
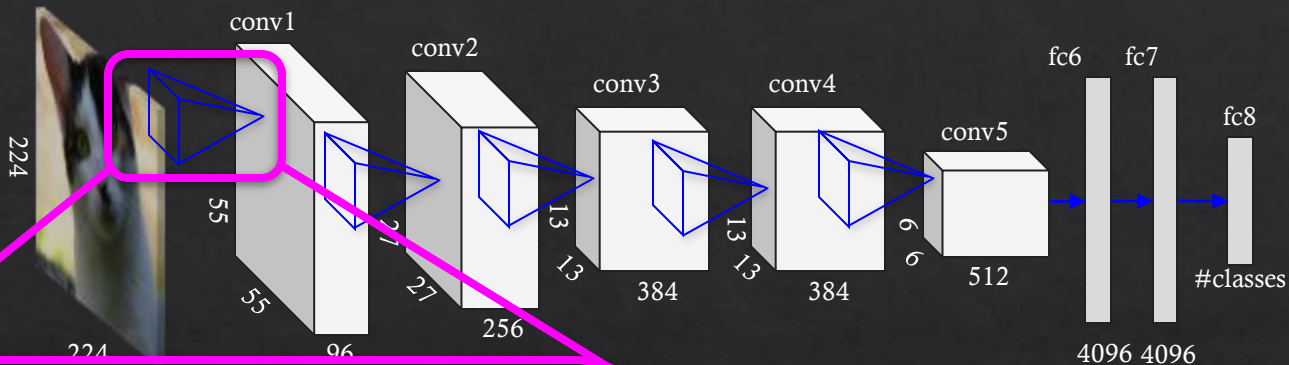
- Interpretability: understand what they learn and why they work
- Monitor training process (evolution of training)
- Gain intuitions to develop better models
- Diagnose potential problems

Today

1. Visualize filters / weights
2. Analyze activations
3. Deconvolutional networks
4. Saliency deconvolutional networks
5. Adversarial noise

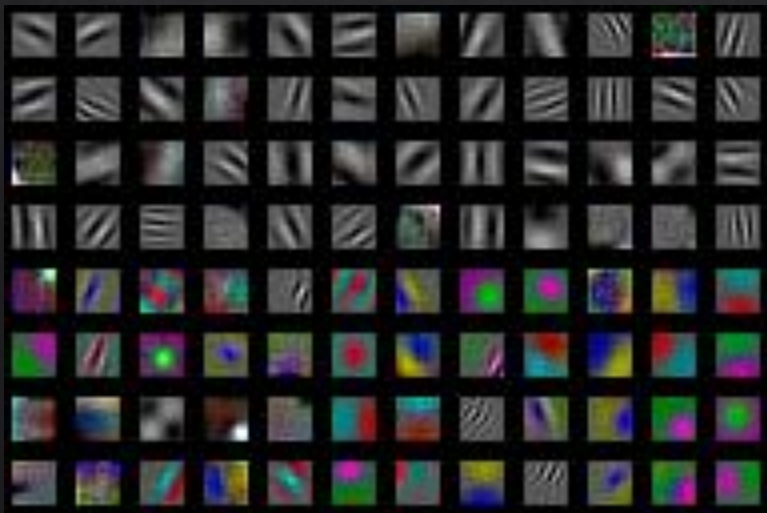
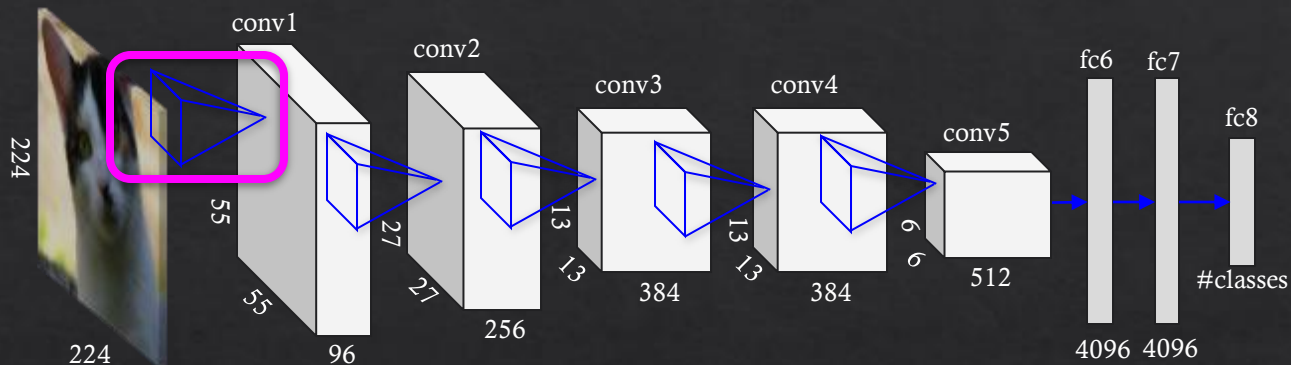
1. Visualize filters

conv1:
96 filters 11x11x3



1. Visualize filters

conv1:
96 filters 1x1 1x3

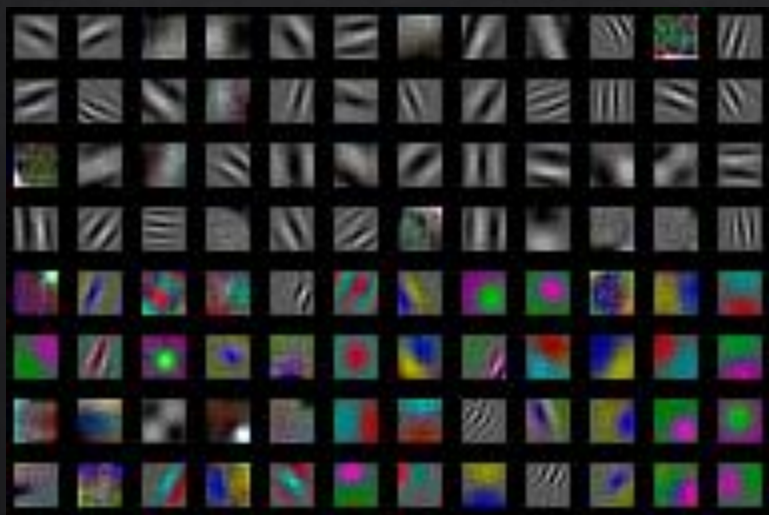
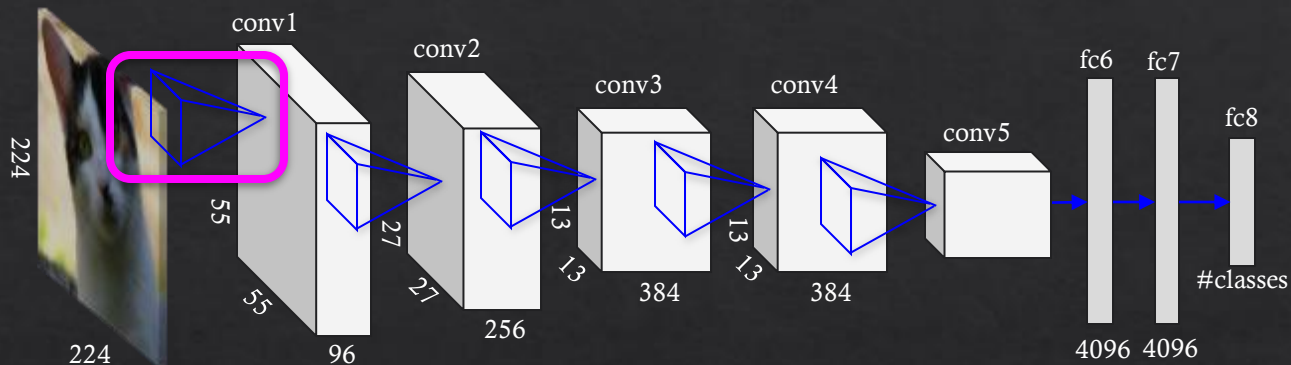


Question:

What do these filters detect?

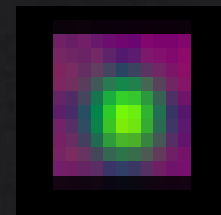
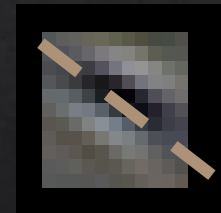
1. Visualize filters

conv1:
96 filters 1x1x3



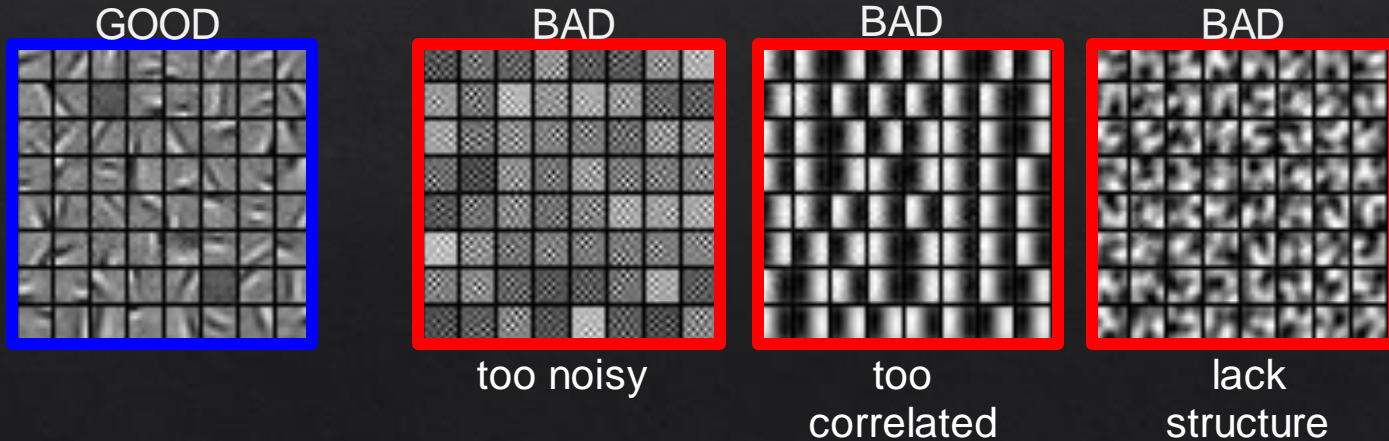
□ Oriented edge filters (similar to Gabor filters)

□ Coloured blob detectors

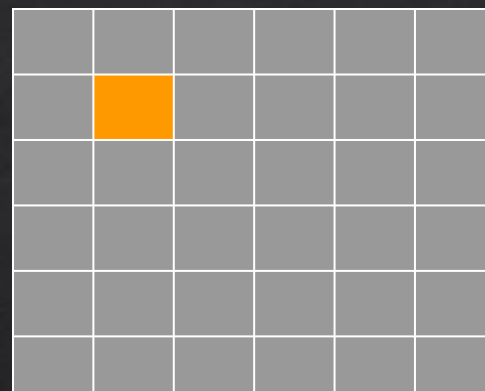


Monitoring filters during training

Good training: learned filters should exhibit structure and are uncorrelated



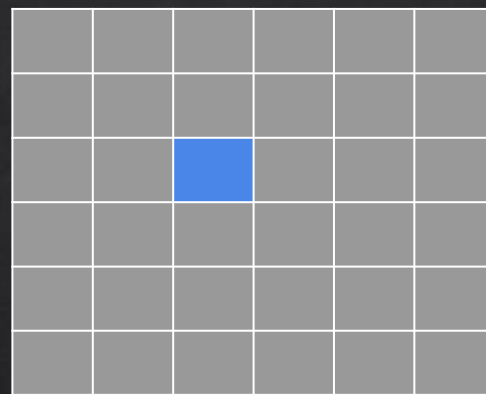
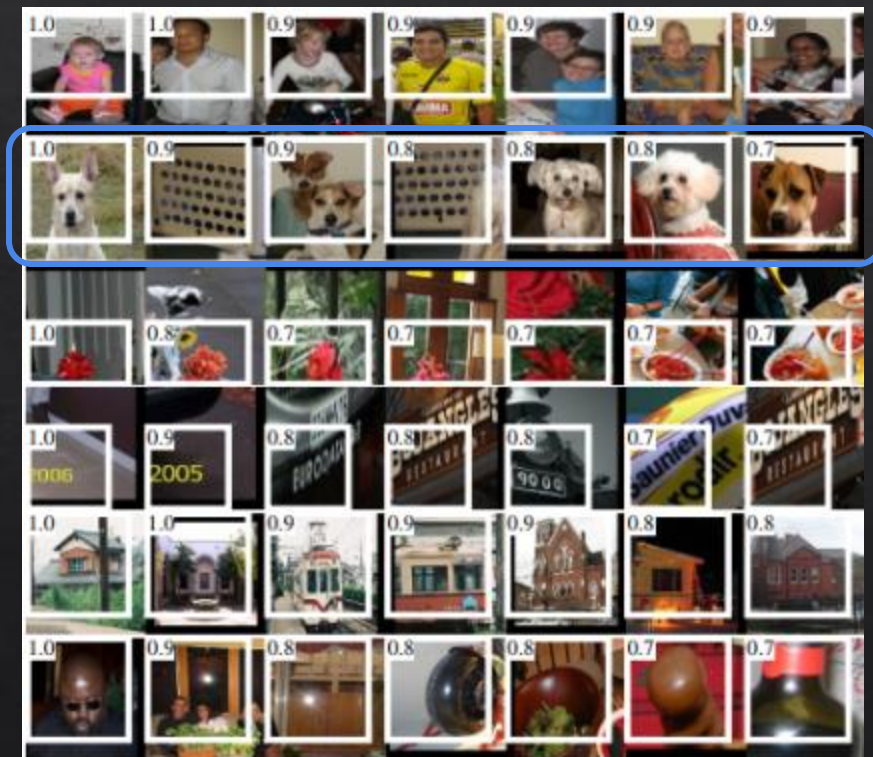
2. Analyze activations



Pool5

1. Pick a neuron at a layer
2. Record it for multiple images
3. Show the images with highest activation value
4. See whether the images correspond to a common concept

2. Analyze activations



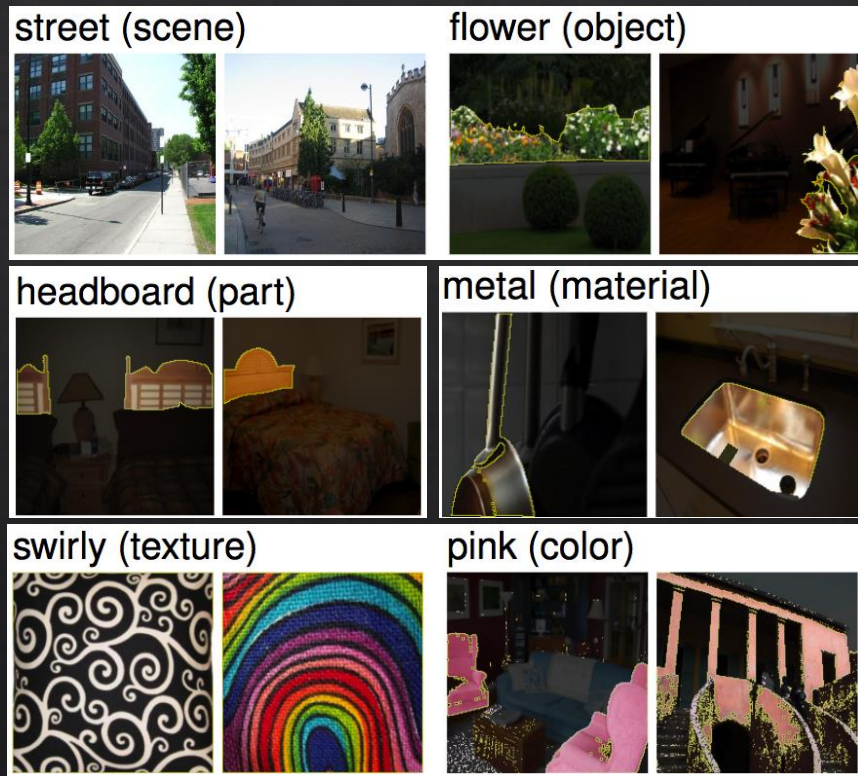
- 😊 Easy to implement
- ☹️ Only qualitative analysis

2. Analyze activations quantitatively

Q. How can we quantify alignment with visual concepts?

I. Collect images and label all the pixels with various concepts

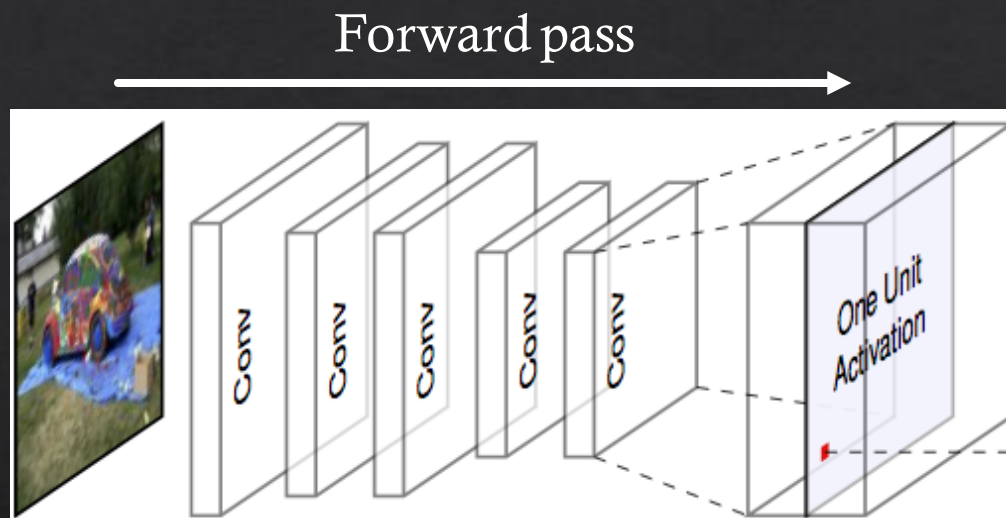
- objects, parts, scenes, textures, colours and materials



2. Analyze activations quantitatively

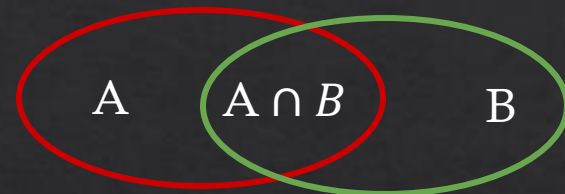
II. Gather responses of neurons to known concepts

- Input image x to CNN
- Take an activation map $A_k(x)$ at layer l
- Threshold $P(A_k(x) > T)$
- Upscale to image size



2. Analyze activations quantitatively

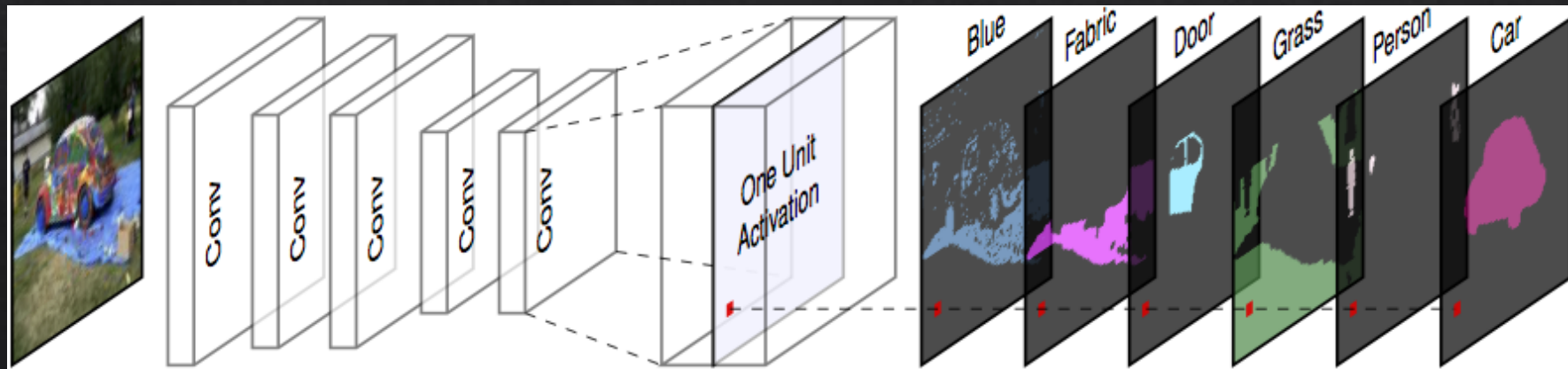
III. Measure overlap with human labelled concepts



Intersection over union

$$IoU = \frac{A \cap B}{A \cup B}$$

Forward pass



conv5 unit 79

car (object)

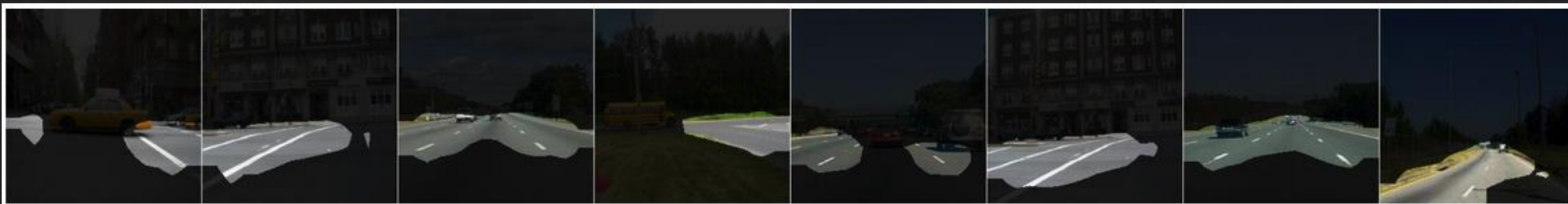
IoU=0.13

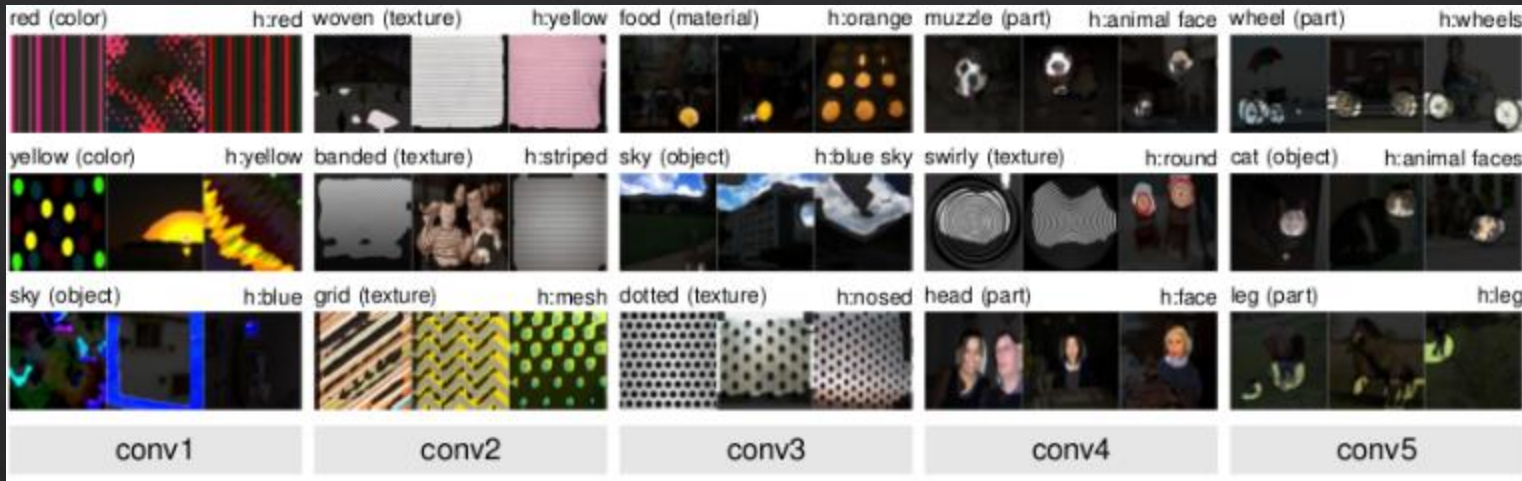
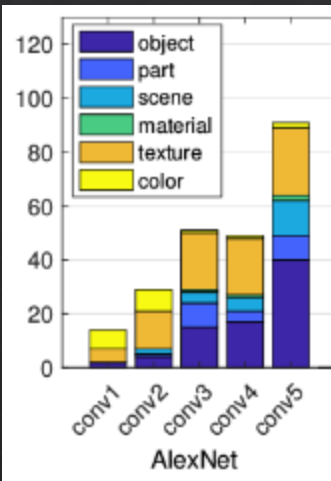


conv5 unit 107

road (object)

IoU=0.15



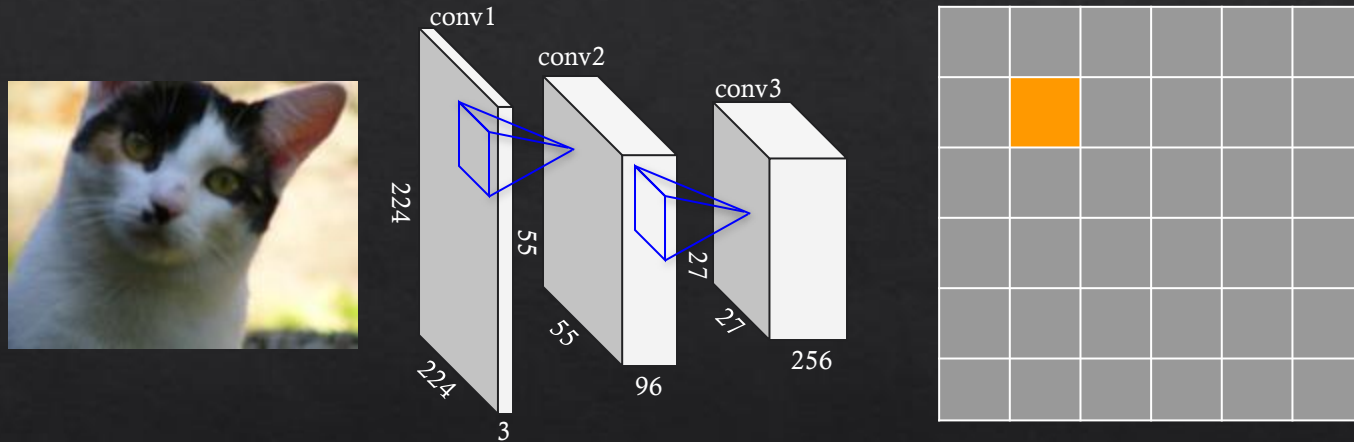


□ More complex concepts emerge at the later layers

□ Some low level concepts at the later layer are still useful for classification

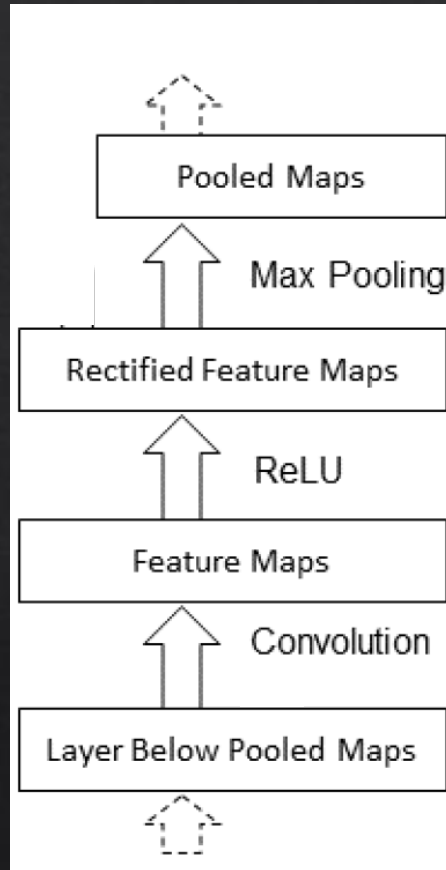
3. Deconvolutional networks

So far, finding correlations between a set of images and activations



What input pattern originally caused a given activation in the feature maps?

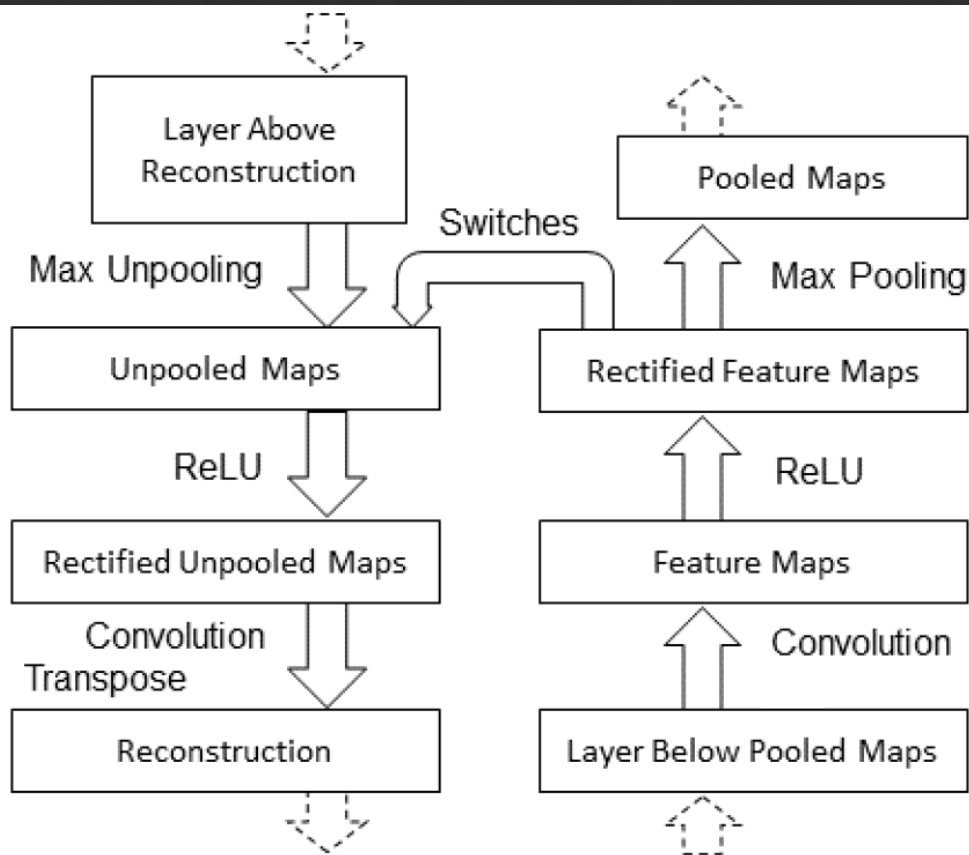
Convnet



How to project the activations back to the input pixel space?

Deconvnet

Convnet



- Deconvnet aims to project the activations back to the input pixel space
- Invert convnet by
 - Unpooling
 - (Un)rectification
 - Convolution transpose

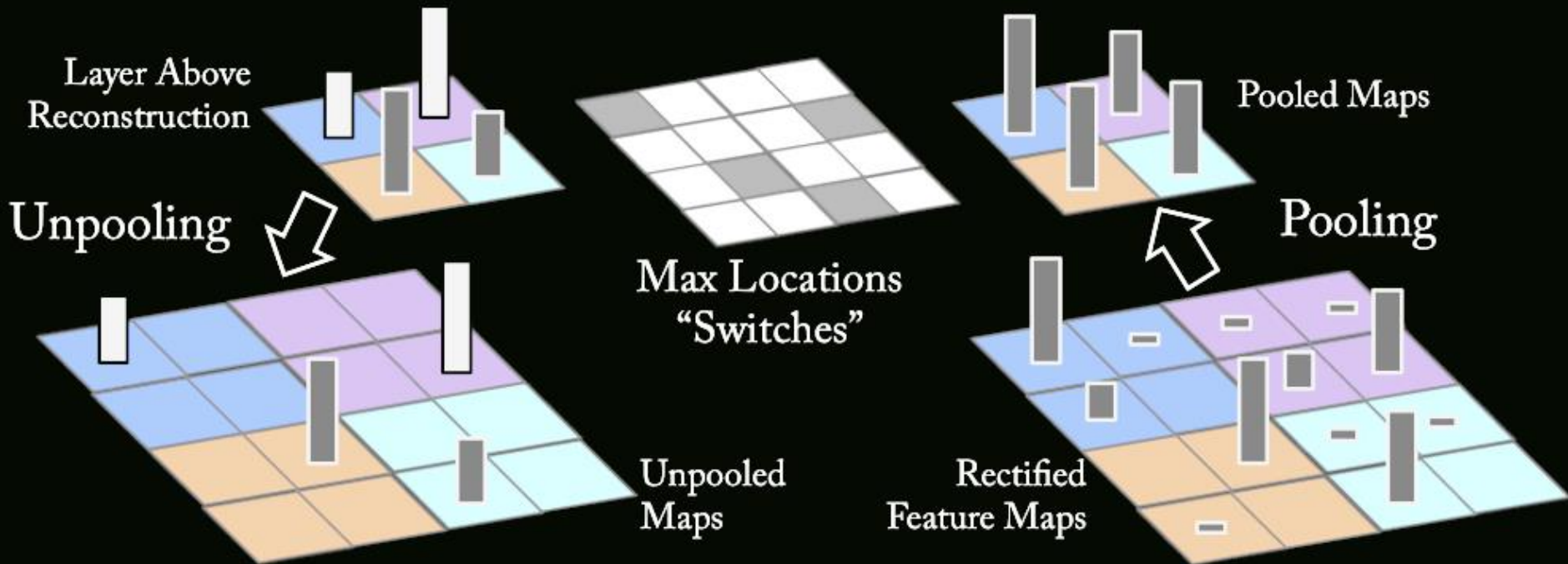
Question

Is max pool operation invertible?

$$y = \text{maxpool}(x)$$

$$x? = (\text{maxpool})^{-1}(y)$$

Unpooling



Unpooling

Relation to backprop (see lecture 8)

		0		0		
		0		1		
		0		0		
		0		0		
		0		1		
		0		0		
		0		0		

$\frac{\partial E}{\partial H^l} \quad \times \quad \frac{\partial H^l}{\partial H^{l-1}} \quad = \quad \frac{\partial E}{\partial H^{l-1}}$

- E is loss function
- $\frac{\partial E}{\partial H_{l-1}} = \frac{\partial E}{\partial H_l} \frac{\partial H_l}{\partial H_{l-1}}$
- Unpooling corresponds to backprop of maxpooling

Unrectification (UnReLU)



$$H_l = \max(H_{l-1}, 0)$$

Relation to backpropagation

$$\frac{\partial E}{\partial H_{l-1}} = \frac{\partial E}{\partial H_l} \cdot \mathbf{1}(H_l > 0)$$



$$R_{l-1} = \max(R_l, 0)$$

UnReLU does not utilise $R_l \cdot \mathbf{1}(R_l > 0)$ but $\max(R_l, 0)$

Transpose convolution (deconvolution?)



Convolution

$$H_l = \text{conv}(H_{l-1}, W_l)$$

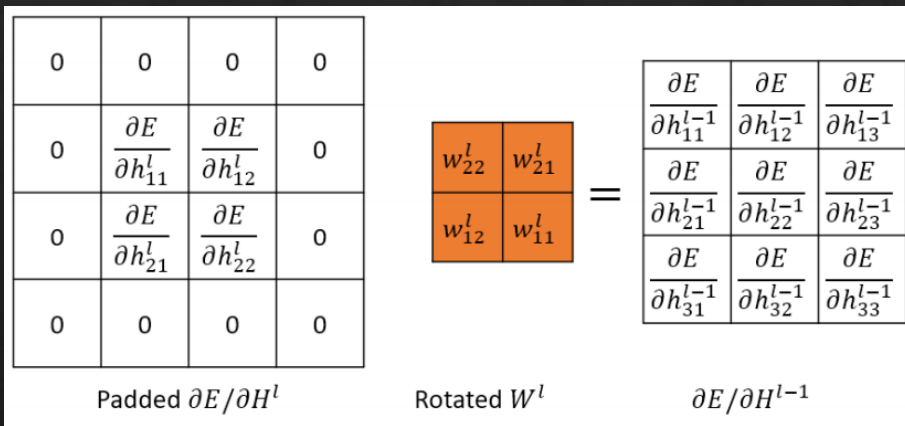
- It is not inverse convolution!
- Usually $\text{conv}(H_{l-1}, W_l) \neq \text{conv}(H_{l-1}, W_l^T)$



Transpose convolution

$$R_{l-1} = \text{conv}(R_l, W_l^T)$$

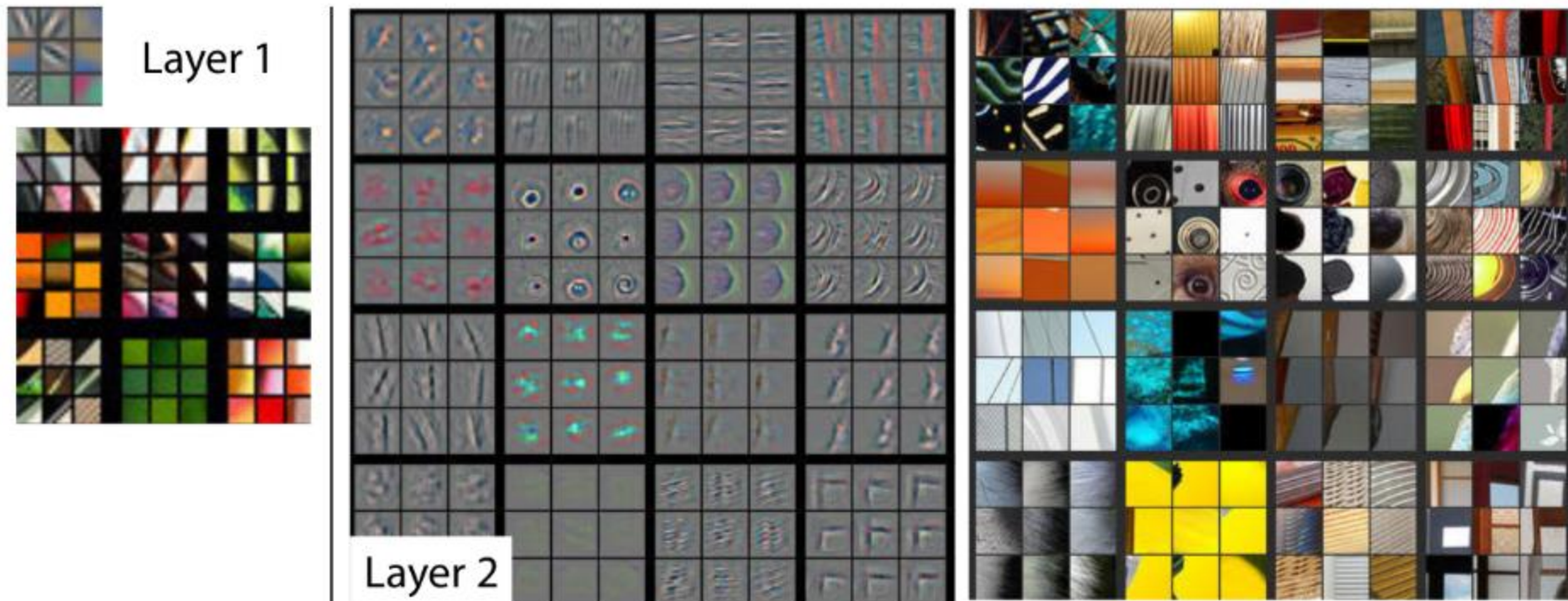
Relation to backprop (see lecture 8)



Layer 1-2: Top-9 Patches

Top 9 activations are projected down to pixel space using deconvolutional net

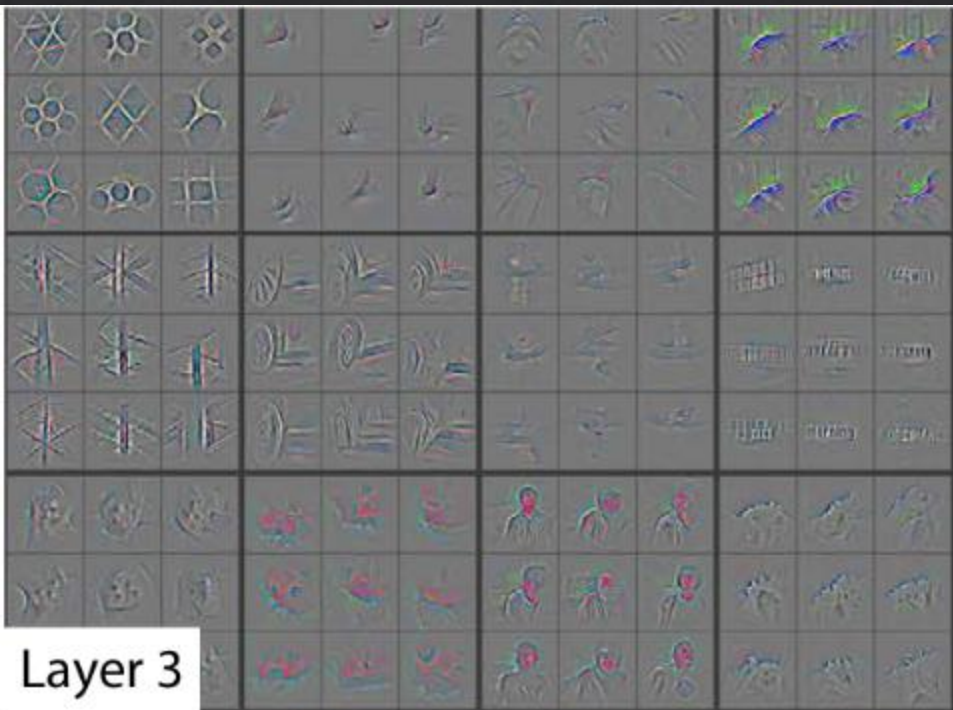
Patches from validation images that give maximal activation of a given feature map



Layer 3: Top-9 Patches

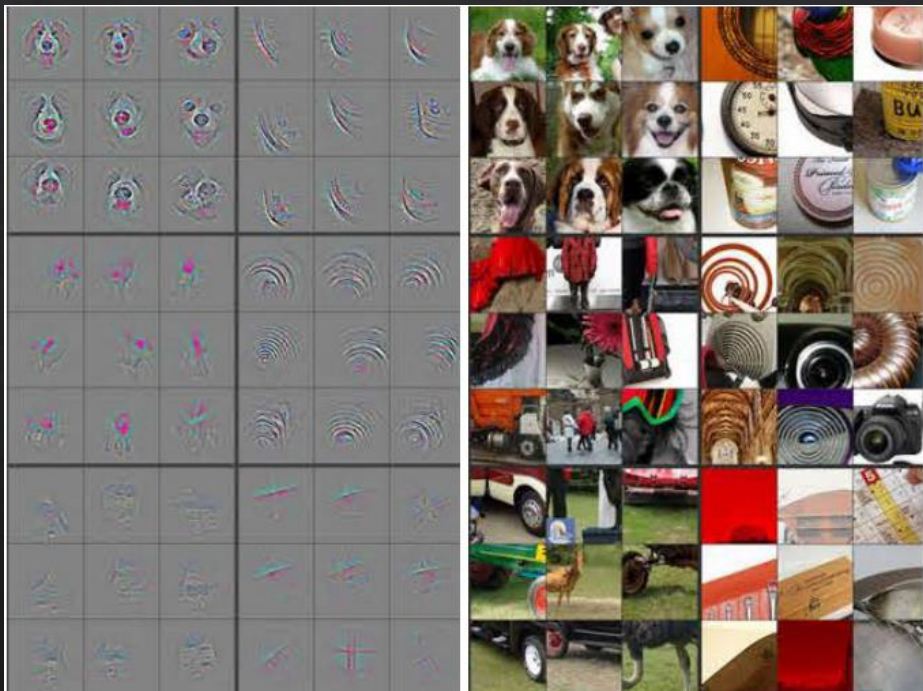
Top 9 activations are projected down to pixel space using deconvolutional net

Patches from validation images that give maximal activation of a given feature map

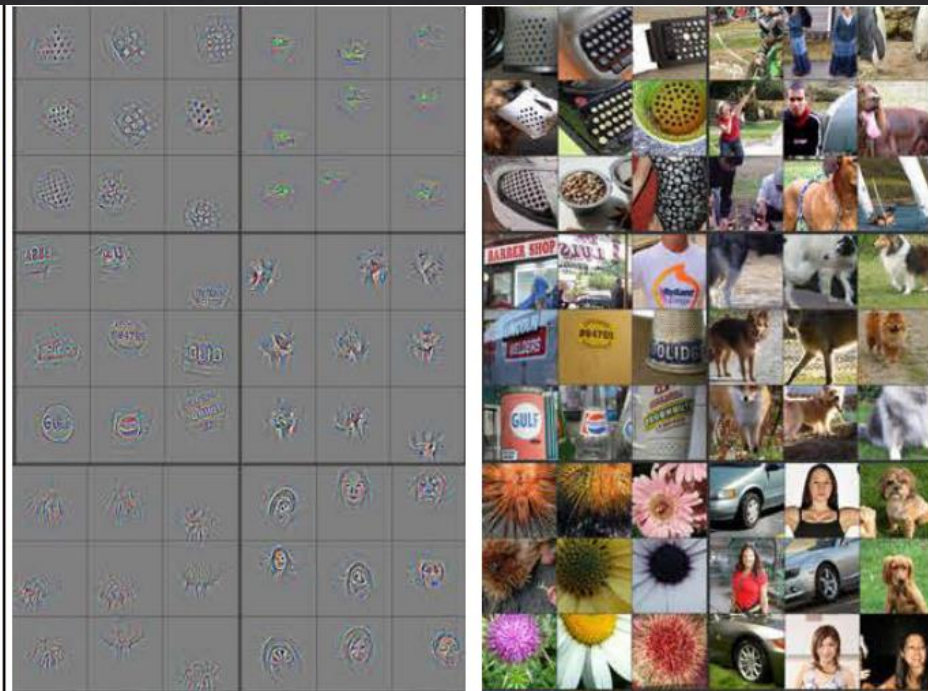


Layer 4-5: Top-9 Patches

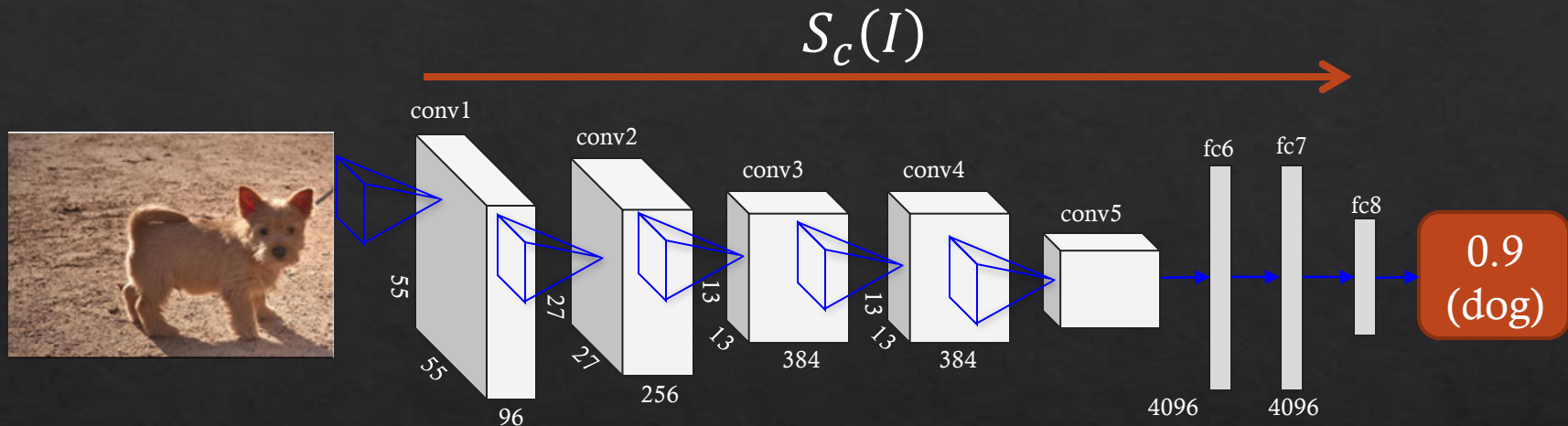
Layer 4



Layer 5

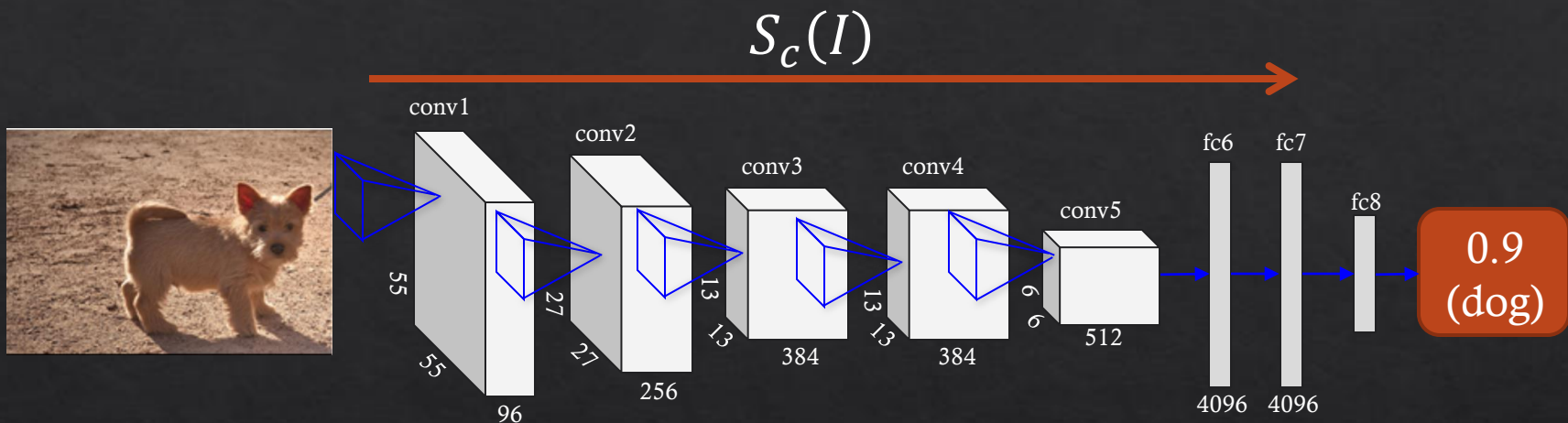


4. Image specific saliency



Which pixels matter most for the prediction?

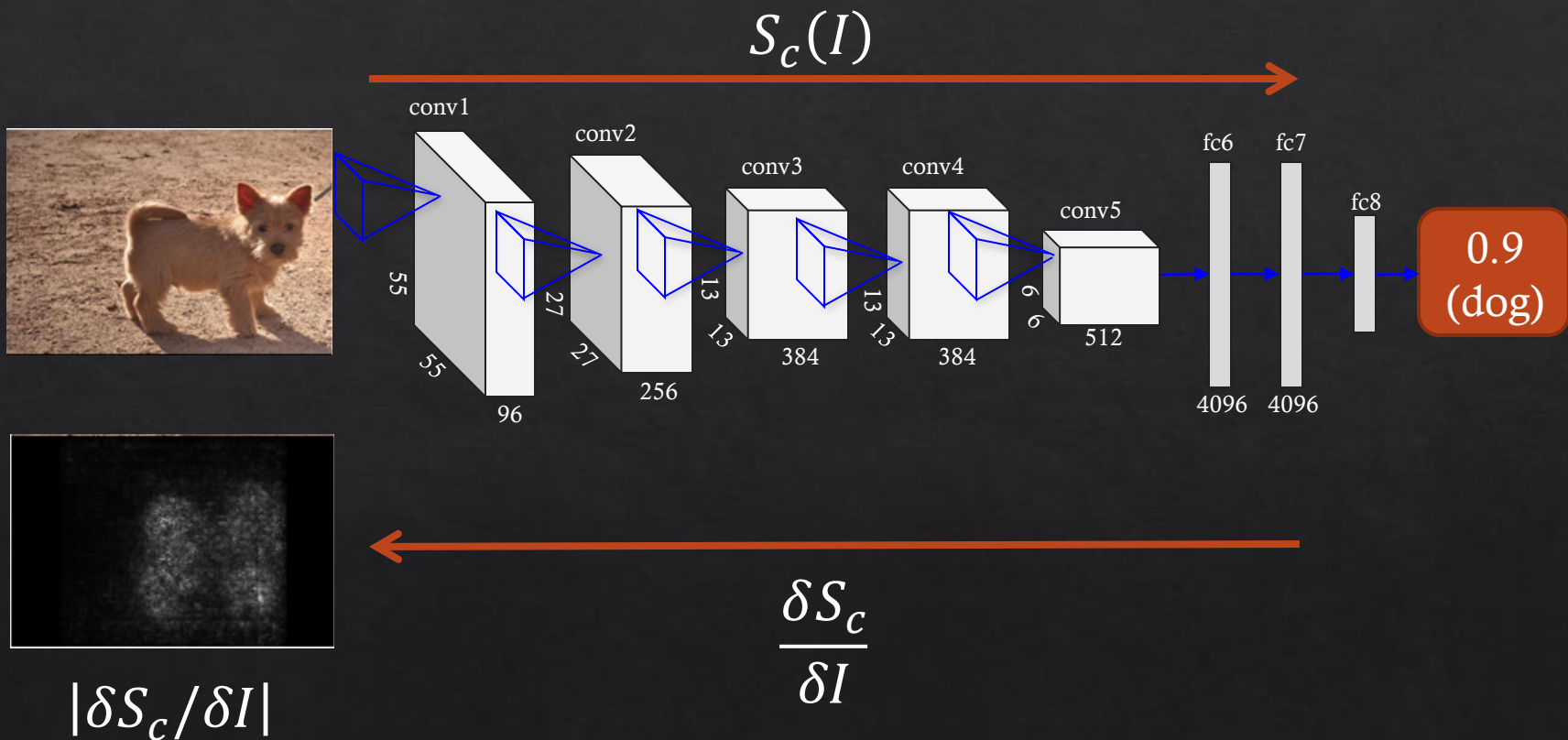
4. Image specific saliency



Question

Can we calculate influence of each pixel on the class probability?

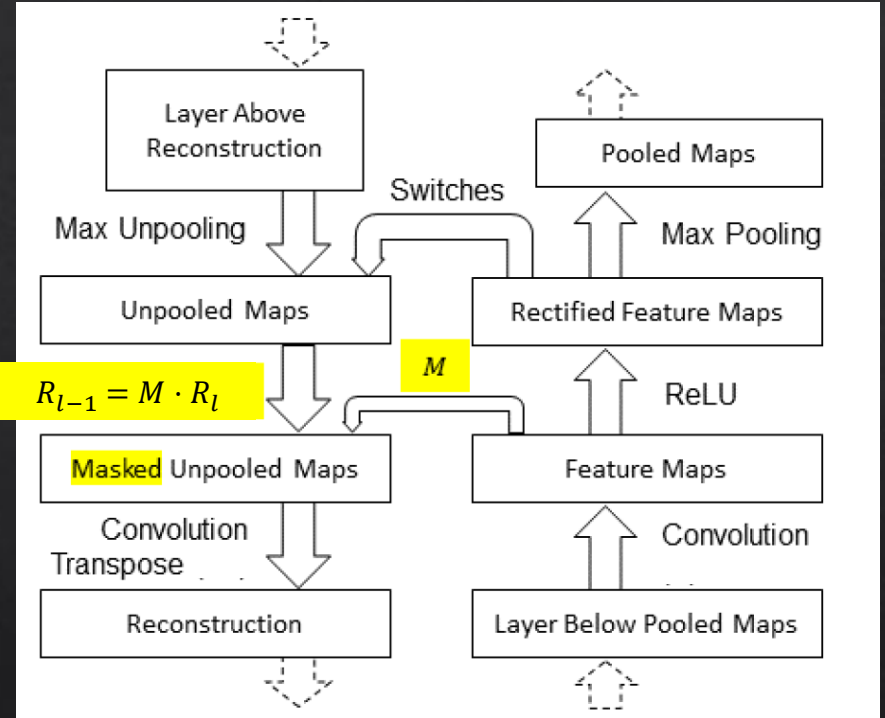
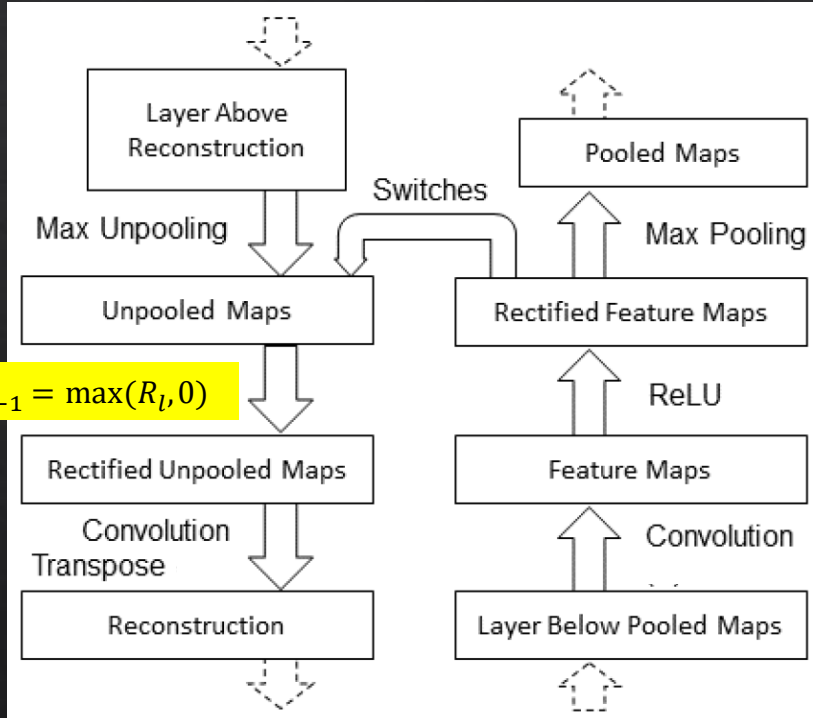
4. Image specific saliency



4. Image specific saliency



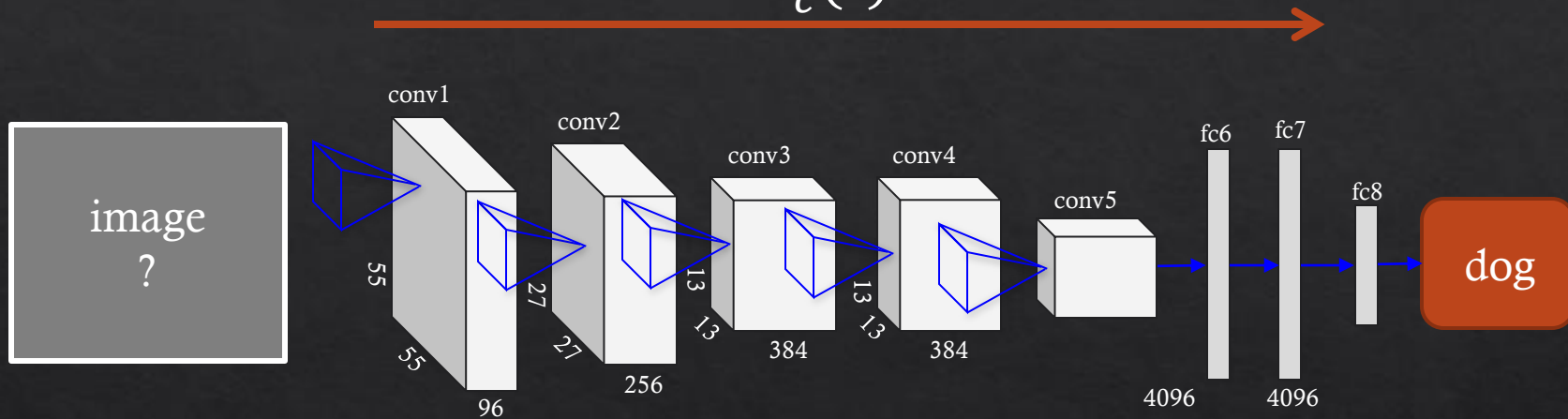
Deconv net vs Saliency net



$$M = \mathbf{1}(R_l > 0)$$

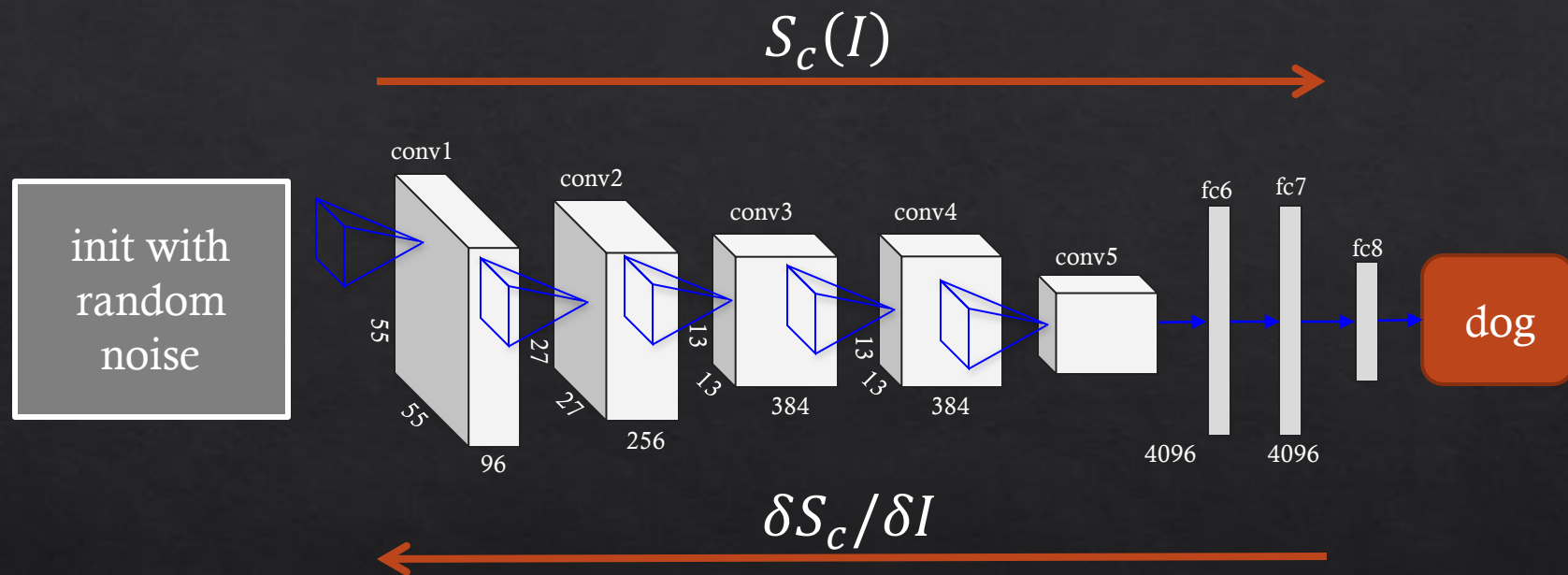
4. Generic class saliency maps

$$S_c(I)$$



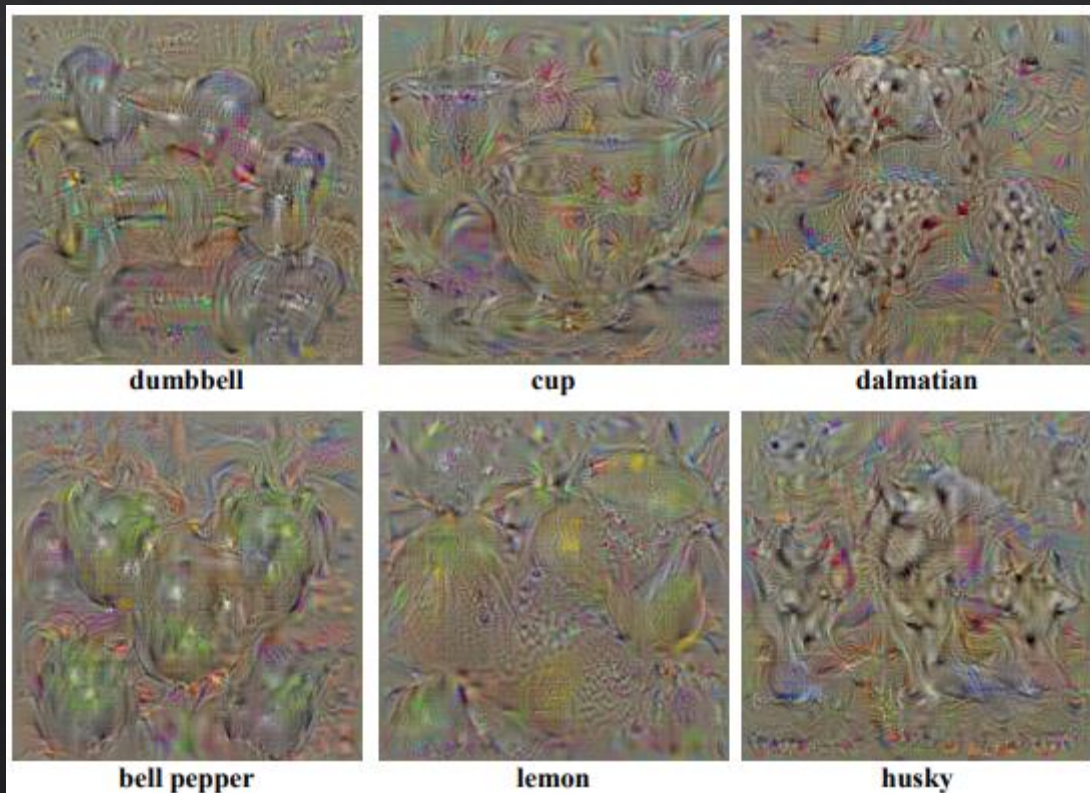
Can we generate an image that outputs high score for dog?

4. Generic class saliency maps



- $\operatorname{argmax}_I S_c(I) - \lambda \|I\|_2^2$
- Maximize “dogness” by modifying pixel values

4. Generic class saliency maps



4. Image and generic class saliency (Deep dream)



"Admiral Dog!"



"The Pig-Snail"

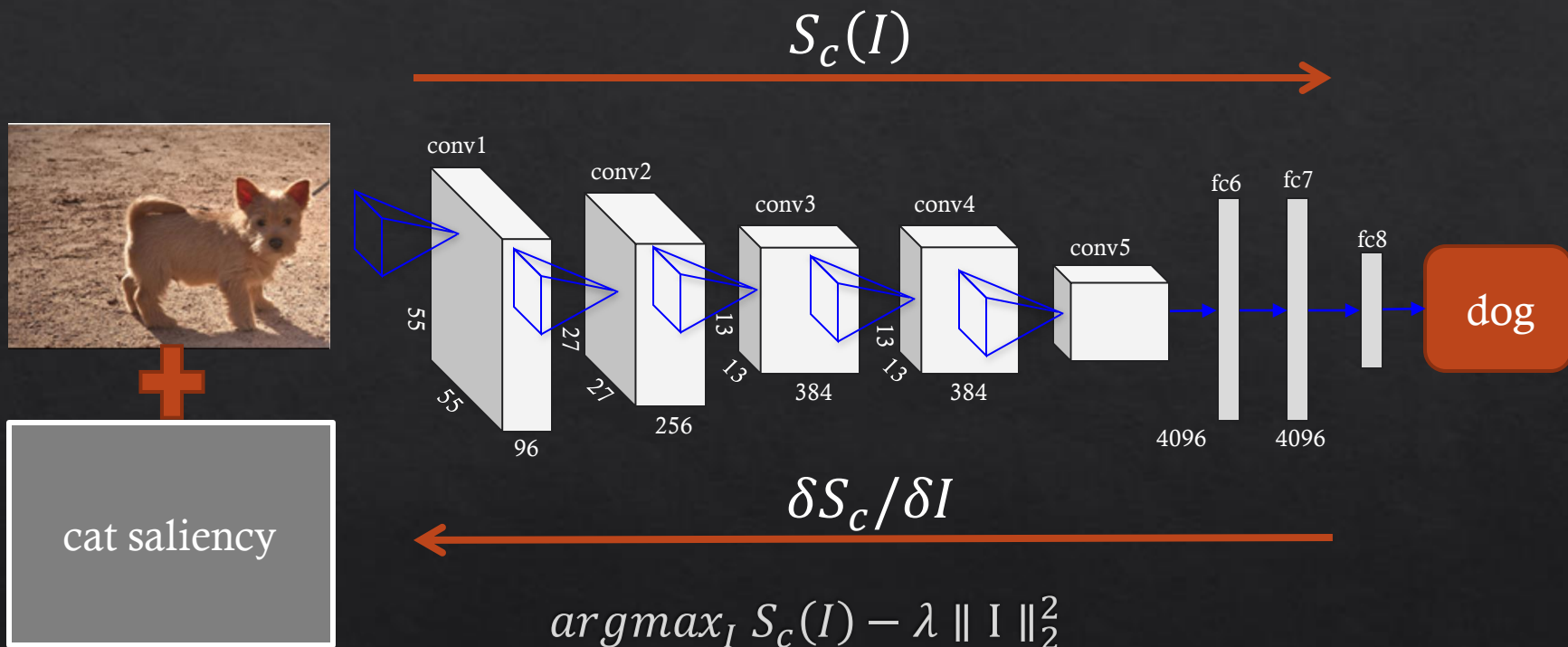


"The Camel-Bird"



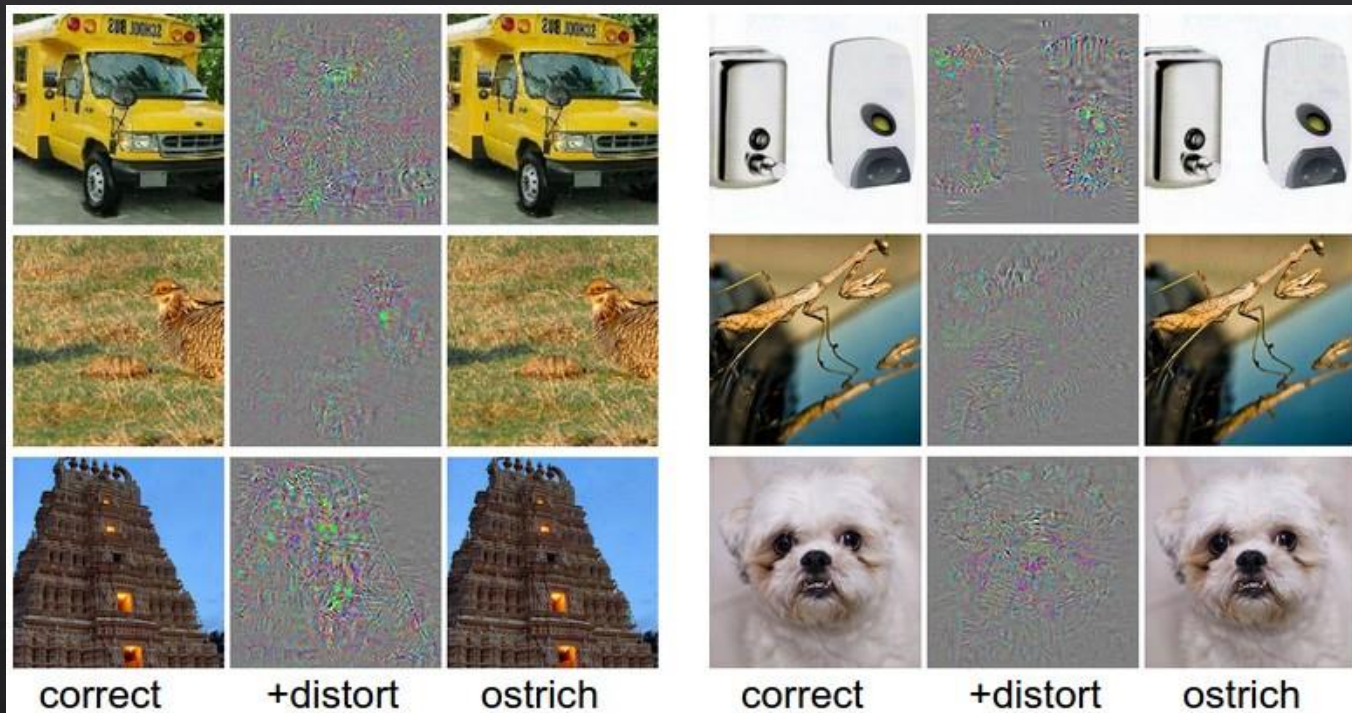
"The Dog-Fish"

4. Image and generic class saliency maps



What happens if we add saliency for another class?

Adversarial Examples



Problem common to any discriminative method!

Summary

Visualize CNNs

- ◇ Filters
- ◇ Highest activations
- ◇ Deconv network
- ◇ Saliency network
- ◇ Generating adversarial samples

Reading material

Recommended

- [Zeiler & Fergus, Visualizing and Understanding Convolutional Networks, ECCV'14](#)

Extra

- [Simonyan, Vedaldi, Zisserman, Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, ICLR'14](#)
- [Szegedy et al. Intriguing properties of neural networks, ICLR'14](#)
- [Nice summary of adversarial techniques by Karpathy](#)
- Try to generate adversarial examples or interesting pictures!