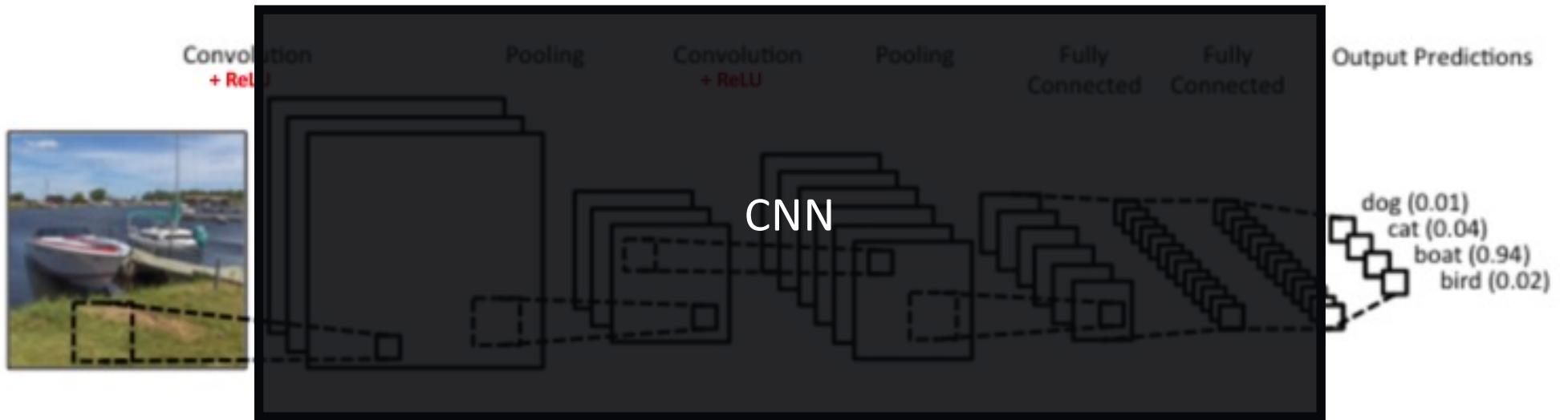


Week 08: Visualizing CNN

Dr. Hongping Cai
Department of Computer Science
University of Bath

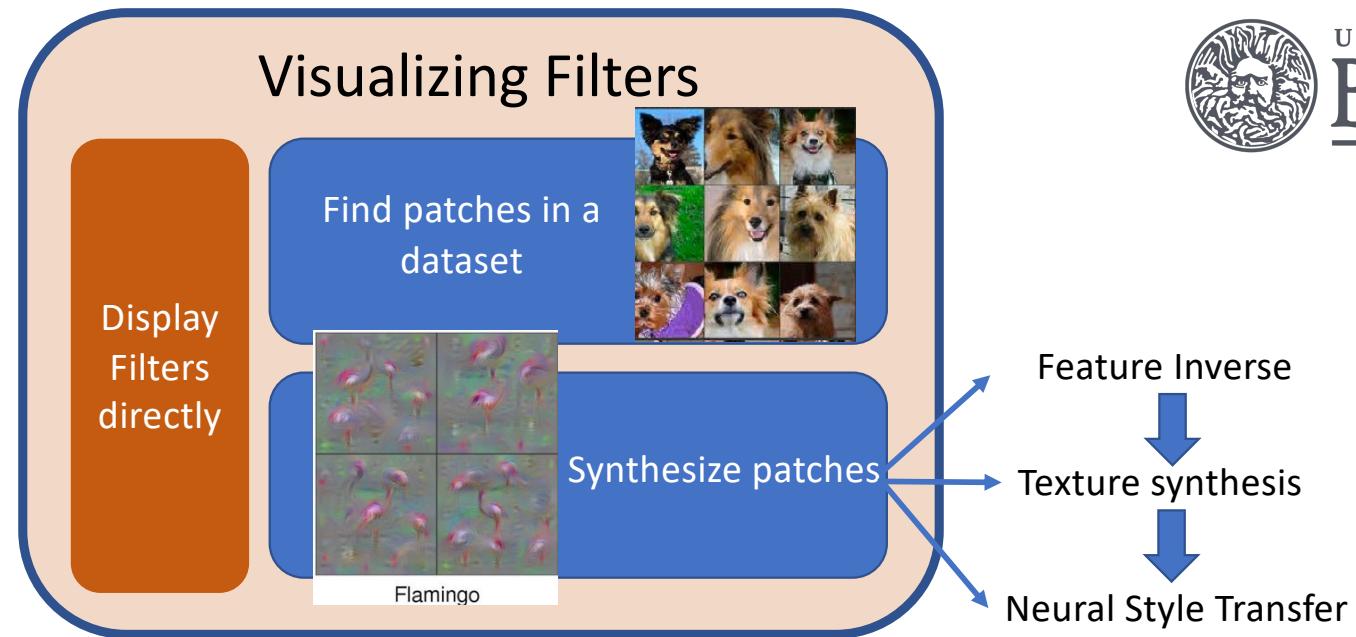
Topic 1: Visualizing Activation Maps



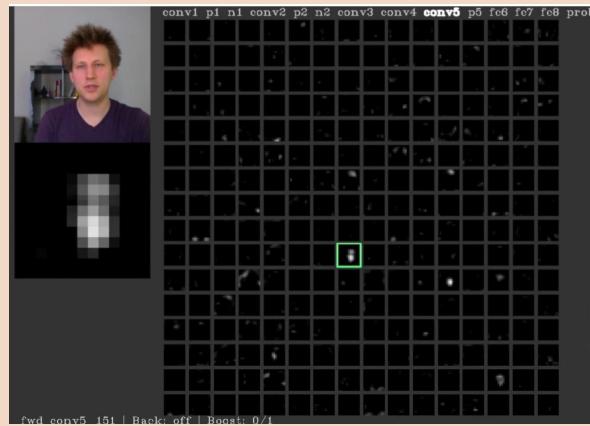
Q: What does CNN learn?

A: Visualizing CNN

Different ways



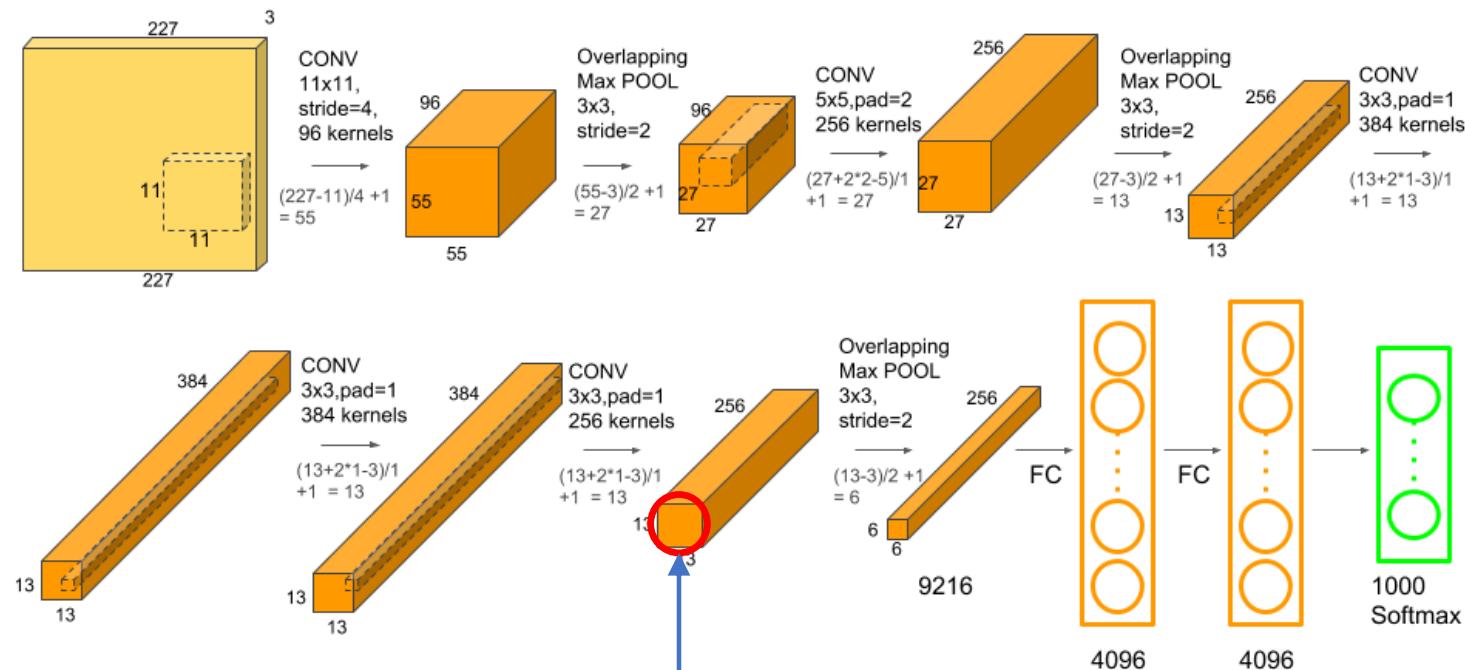
Visualizing Activation Maps



Visualizing Heatmaps

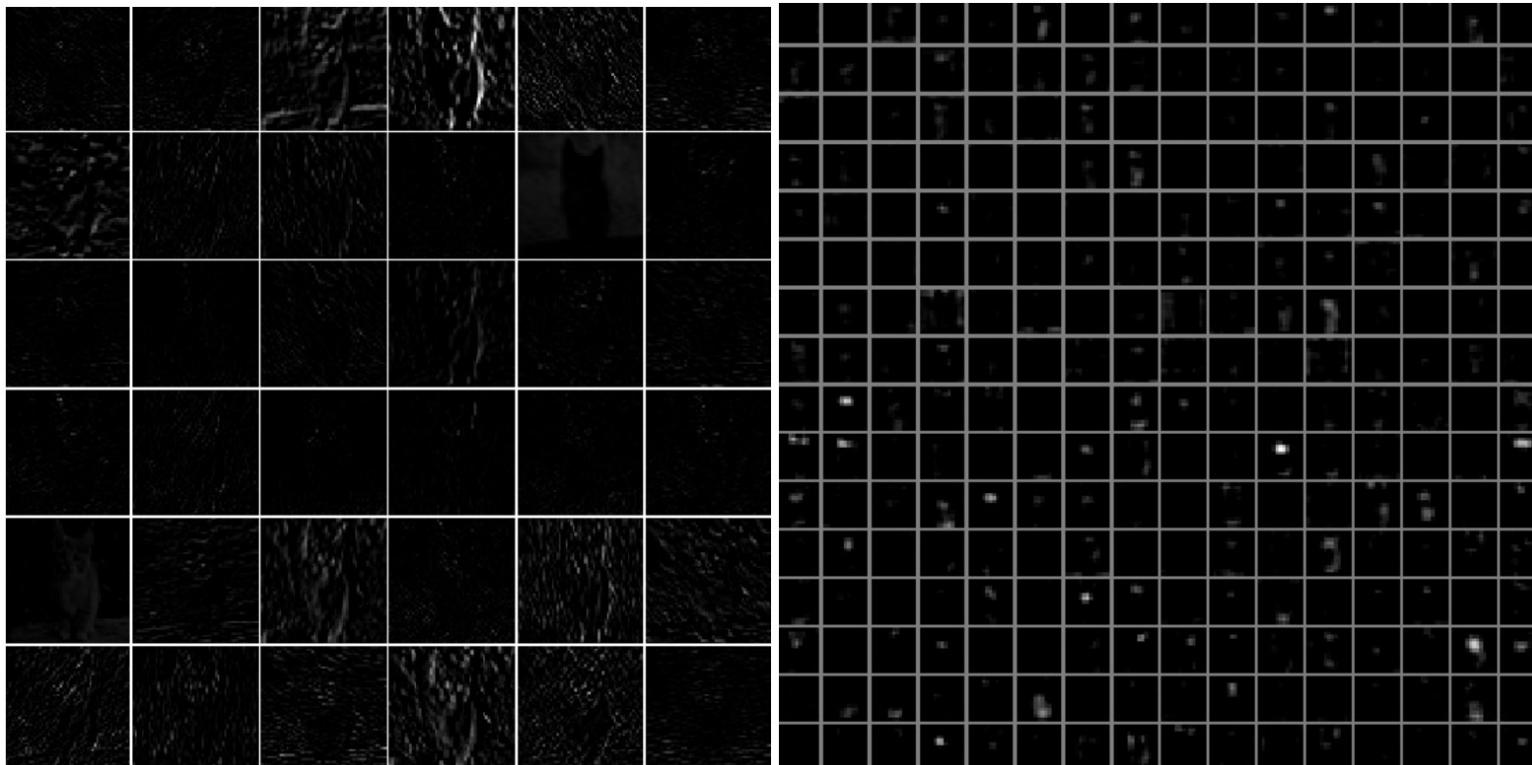


Visualizing the activation maps



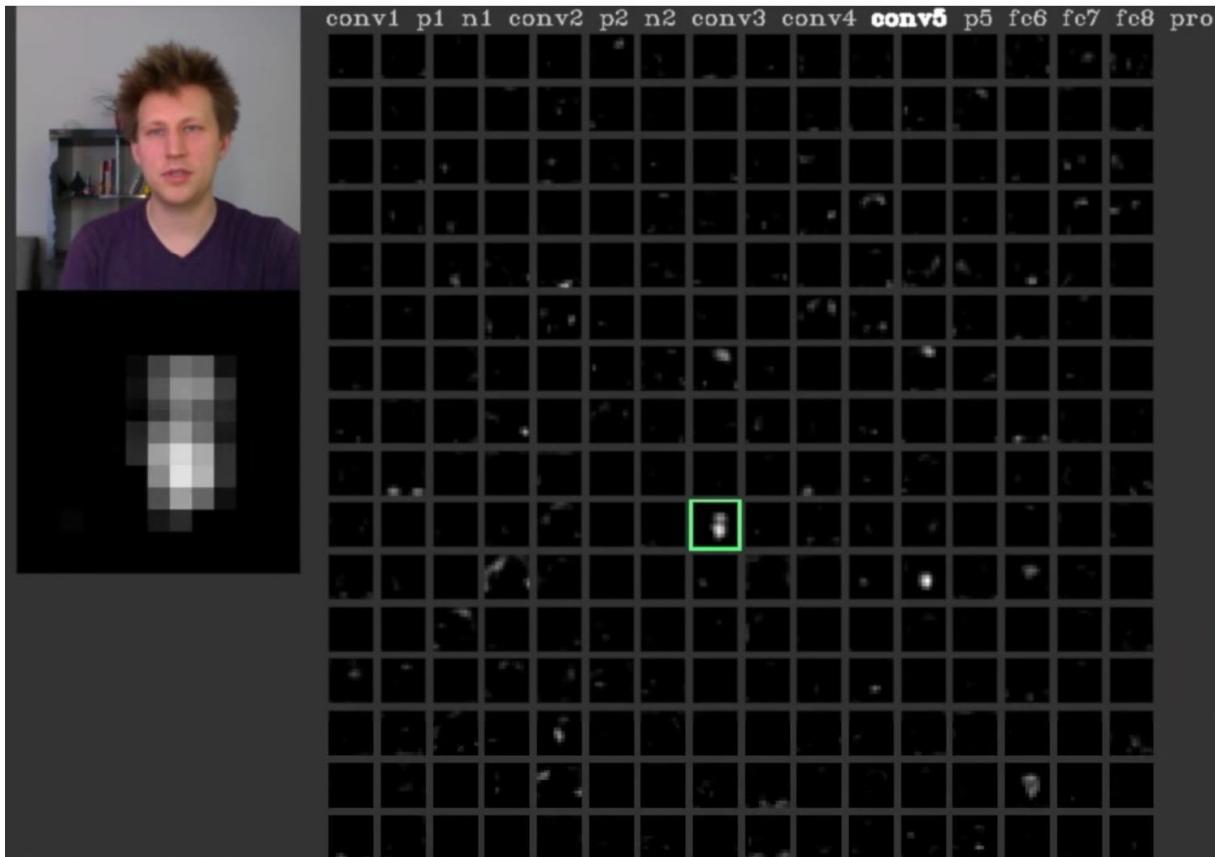
Visualizing the activation maps

Visualizing the activation maps



Activation maps on the first CONV layer (left), and the 5th CONV layer (right)
of a trained AlexNet of a cat image

Visualizing the activation maps



Q: What does it mean that many higher-layer feature maps are blank?

A: The pattern encoded by these filters were not found in the input image.

For higher layers, the feature maps become less interpretable.

References

- Online Lecture: “Visualizing what ConvNets learn”.
<https://cs231n.github.io/understanding-cnn/>
- Video lecture in University of Stanford. “Visualizing and Understanding”. 2017.
<https://www.youtube.com/watch?v=6wcs6szJWMY>
- Blog by Jason Yosinski. “Understanding Neural Networks Through Deep Visualization”. 2015. <https://yosinski.com/deepvis>

Topic 2: Visualizing Filters

Visualizing Filters

Visualizing the filters directly

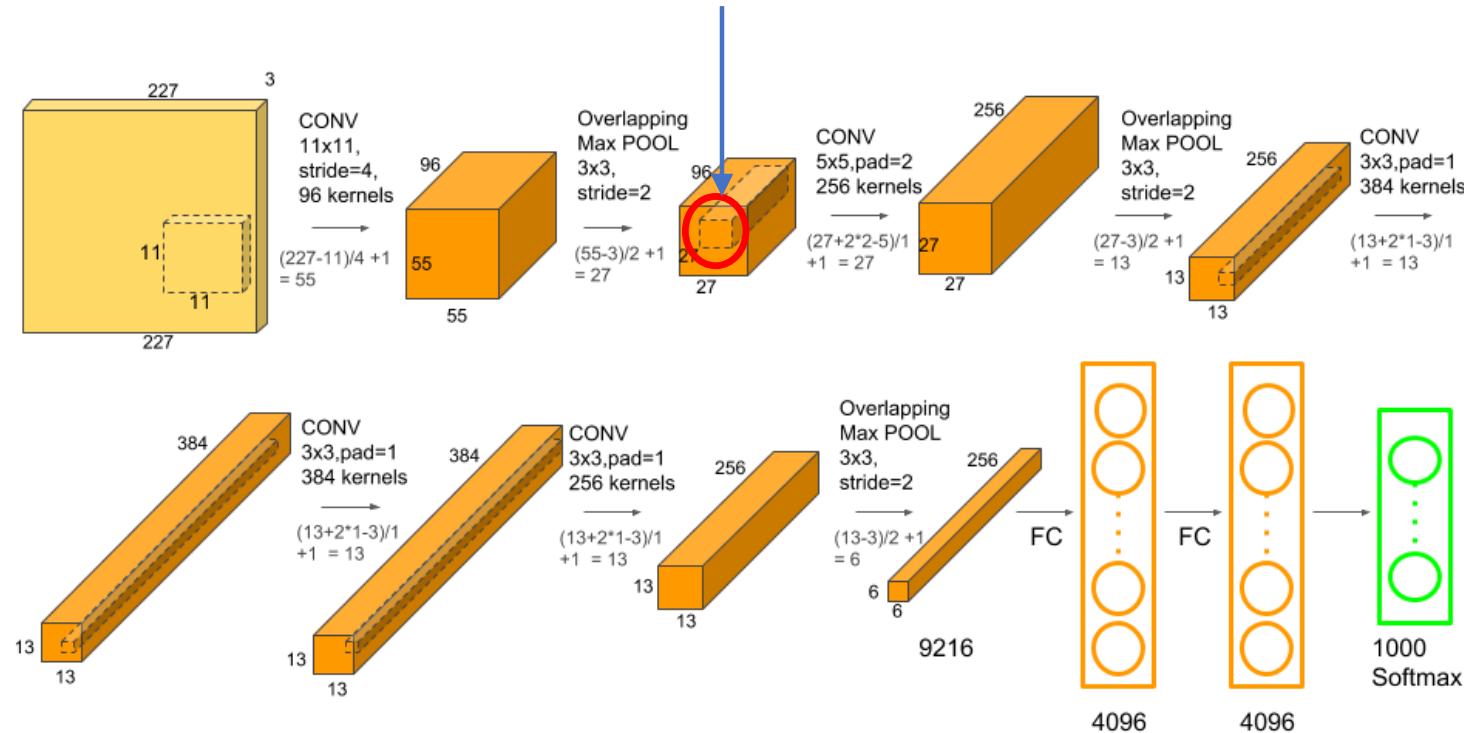
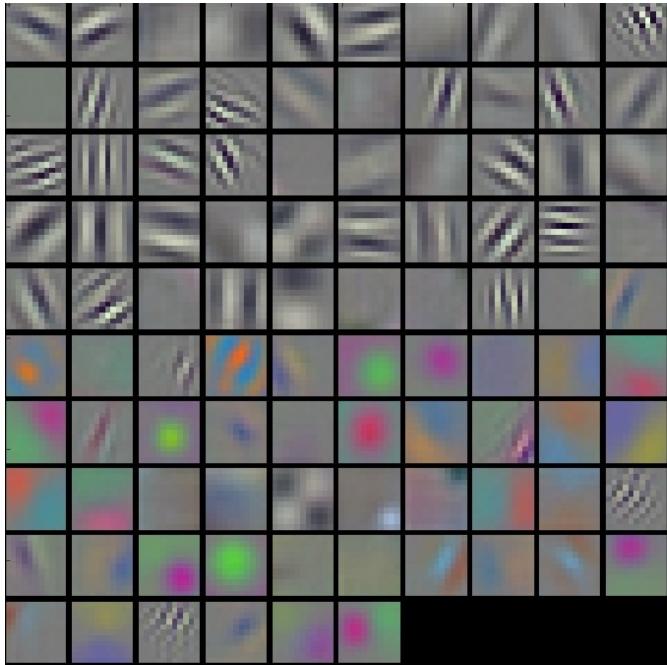


Image from: <https://learnopencv.com/understanding-alexnet/>

Visualizing Filters

1st Conv Layer of AlexNet



2nd Conv Layer of AlexNet

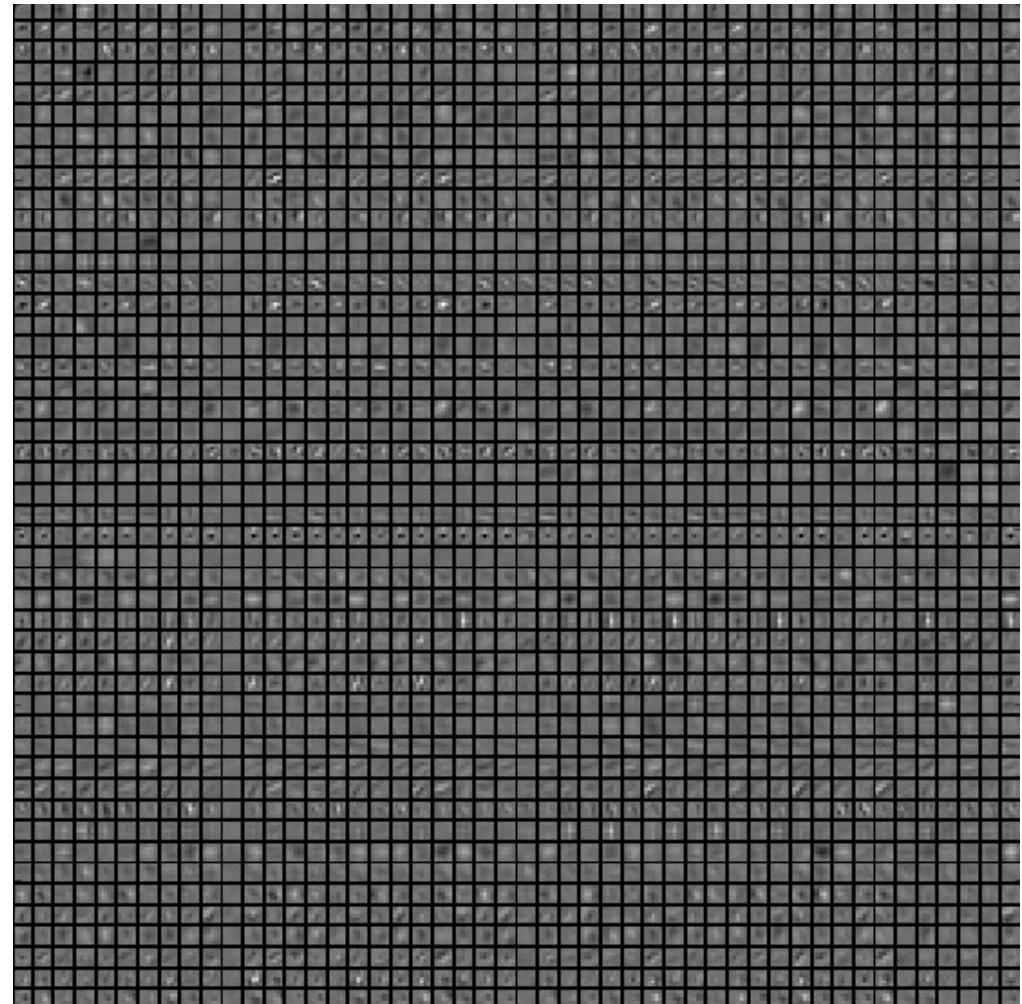
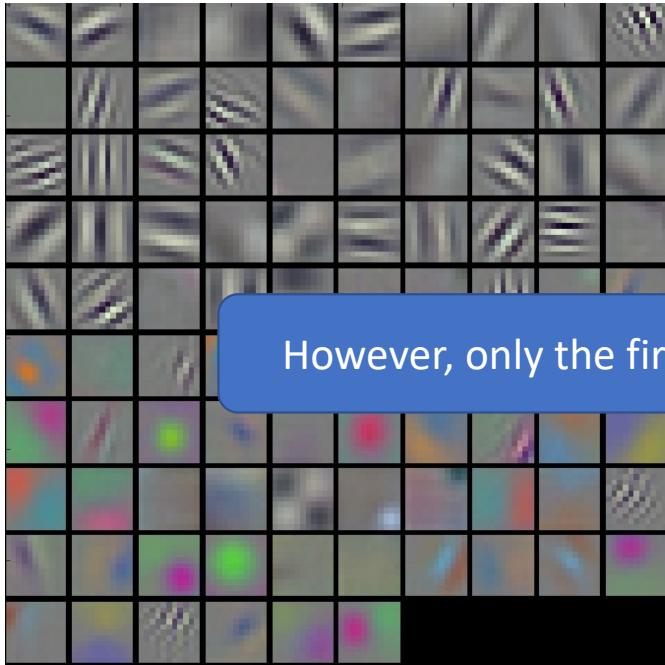


Image from: <https://cs231n.github.io/understanding-cnn/>

- Helpful to check if the network is trained properly (no noisy pattern)

Visualizing Filters

1st Conv Layer of AlexNet



However, only the first layer is interpretable by visualizing filters directly.

- Helpful to check if the network is trained properly (no noisy pattern)

2nd Conv Layer of AlexNet

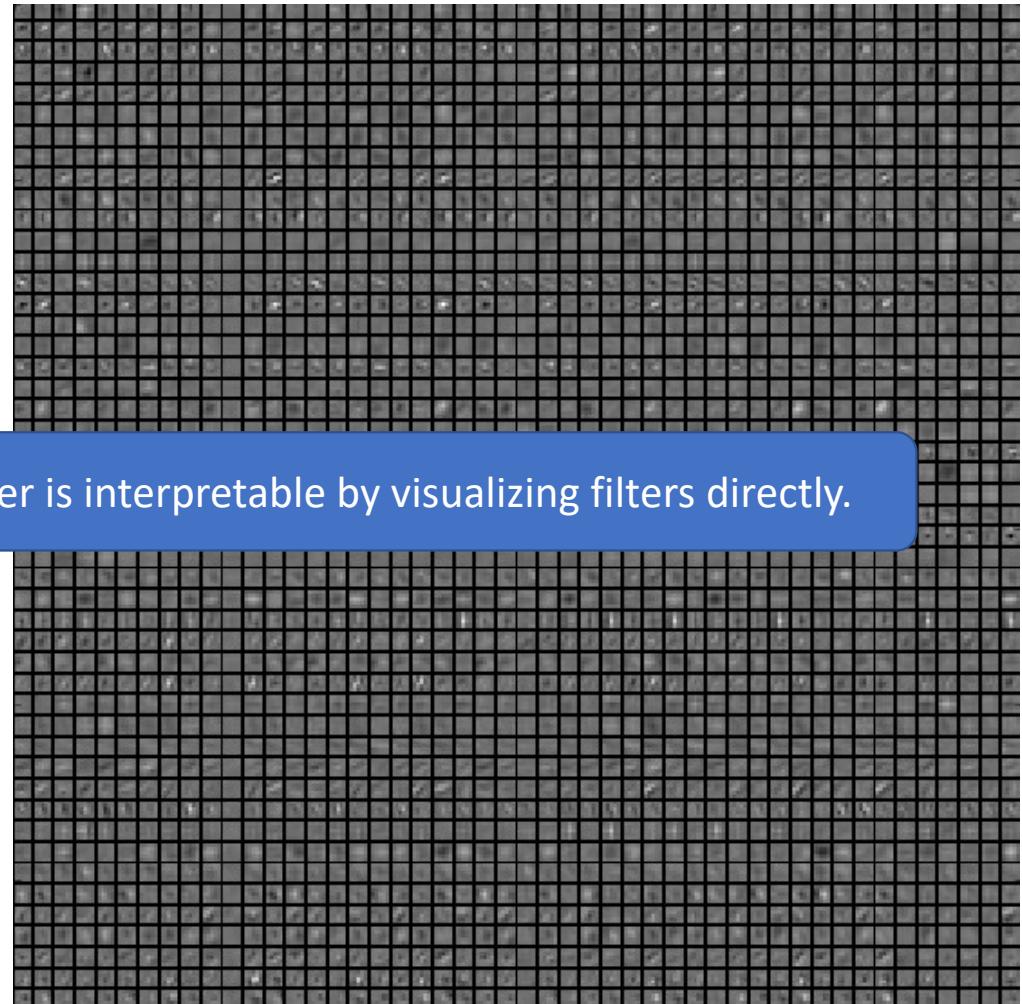
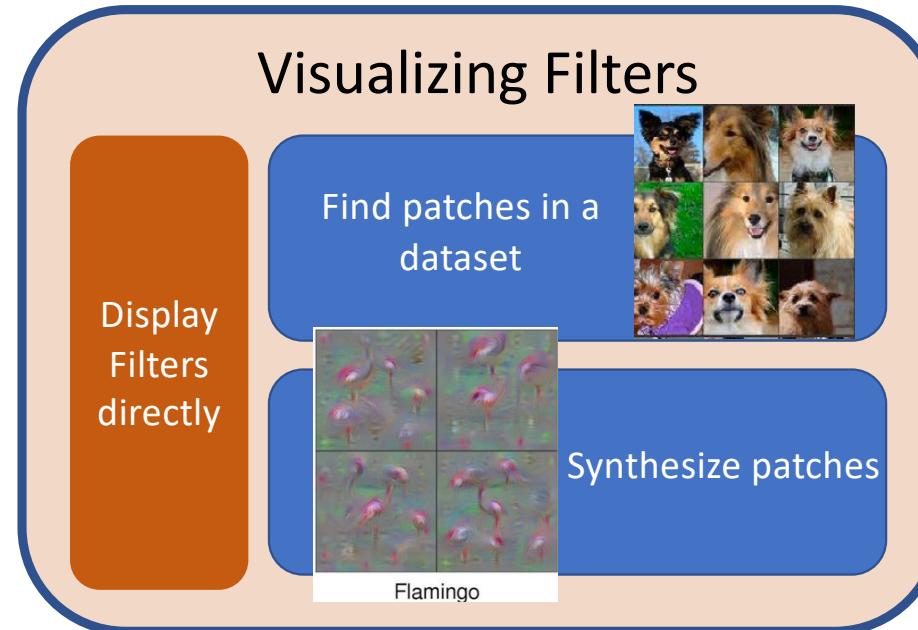


Image from: <https://cs231n.github.io/understanding-cnn/>

Visualizing Filters implicitly



- Instead of visualizing filters directly, visualize what sort of input maximizes the activation of one filter.
- We get an idea of what pattern each filter has learnt to extract.

Maximally Activating Patches

- Find patches produce maximal activation for a certain filter

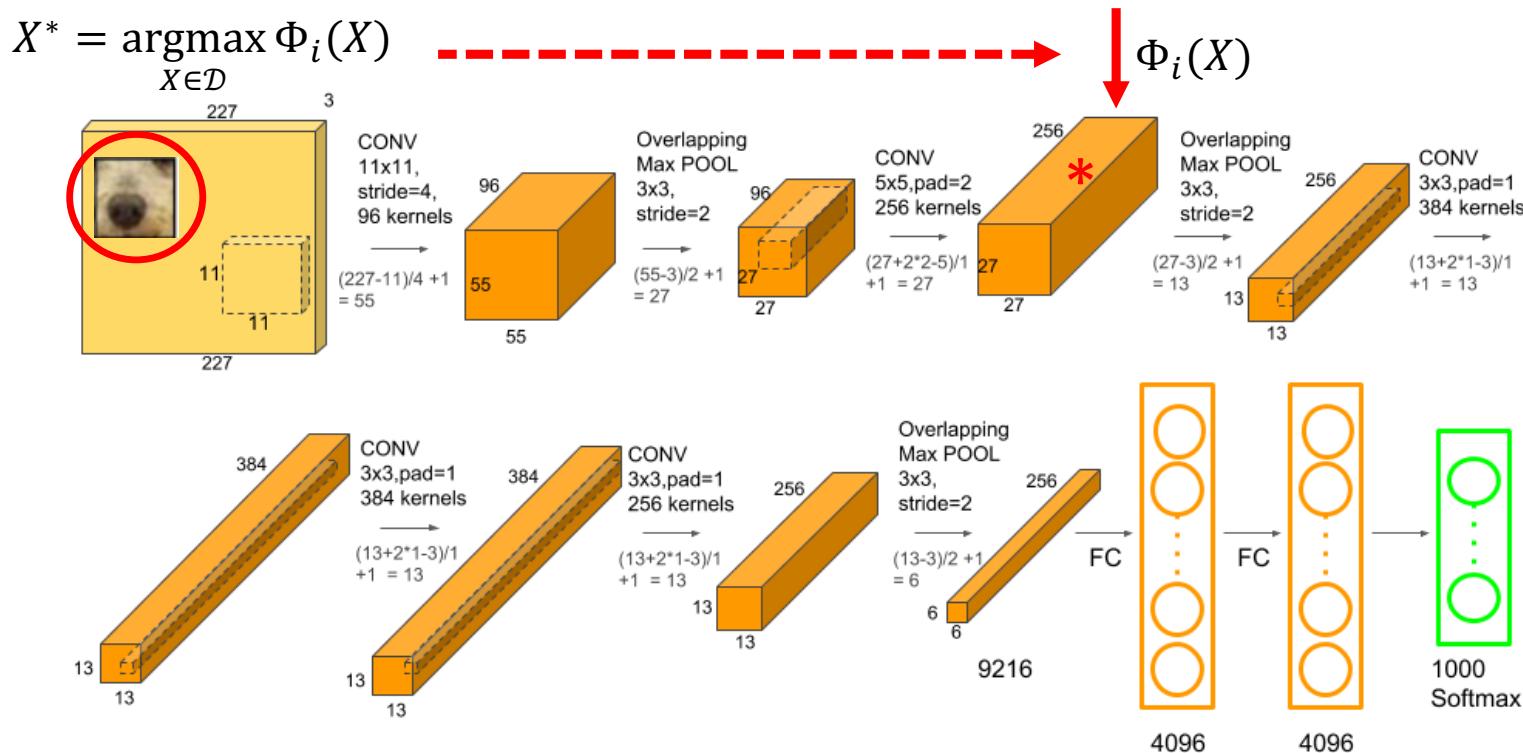
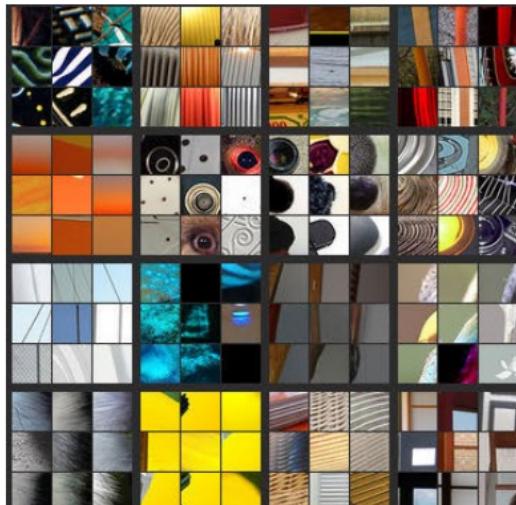


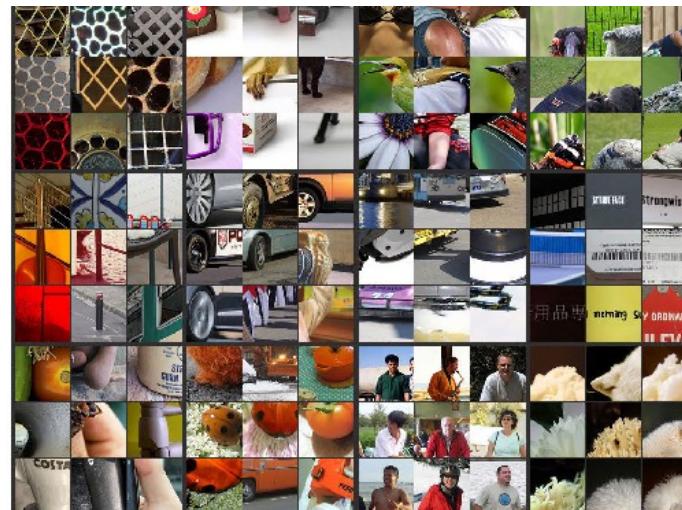
Image from: <https://learnopencv.com/understanding-alexnet/>

Maximally Activating Patches

Top 9 image patches (from ImageNet) maximally activating each filter



Layer 2



Layer 3

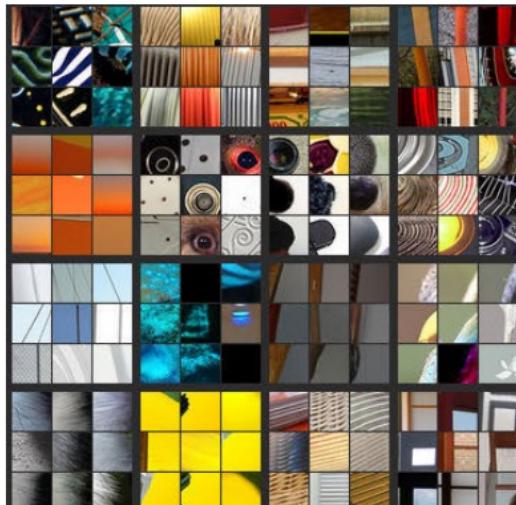


Layer 5

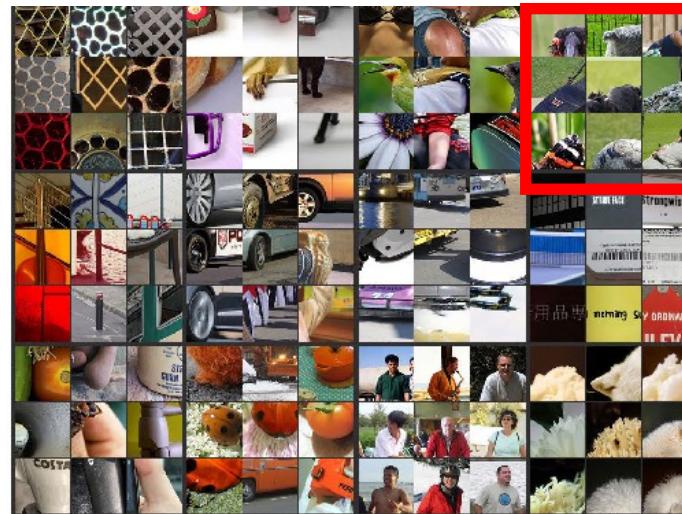
Zeiler, Fergus. "Visualizing and Understanding Convolutional Networks". ECCV14

Maximally Activating Patches

Top 9 image patches (from ImageNet) maximally activating each filter



Layer 2



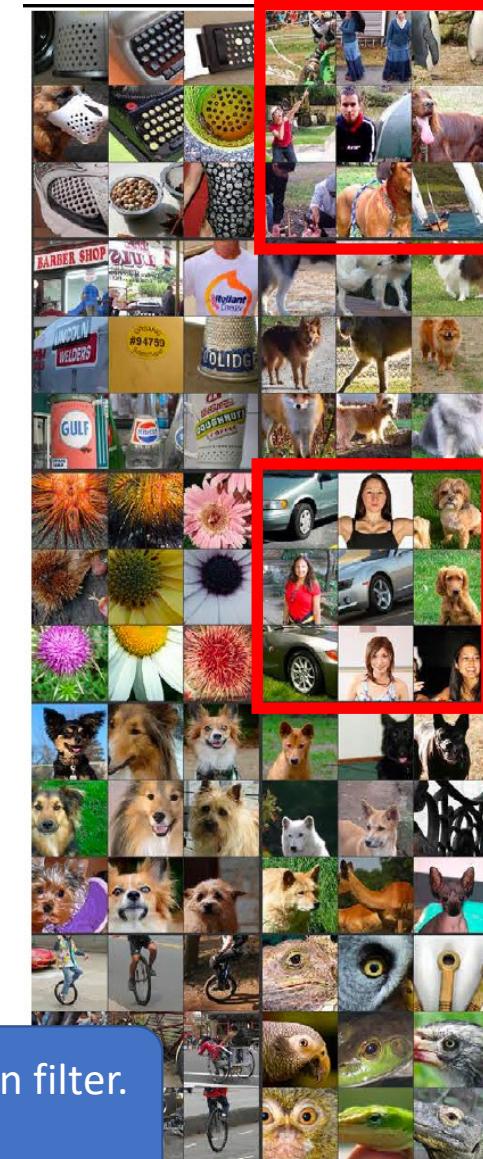
Layer 3

Q1: What are these 9 patches in common?

Q2: What is the discriminative part for these 9 patches?

We need also visualize which pixels in the patch “excite” a given filter.

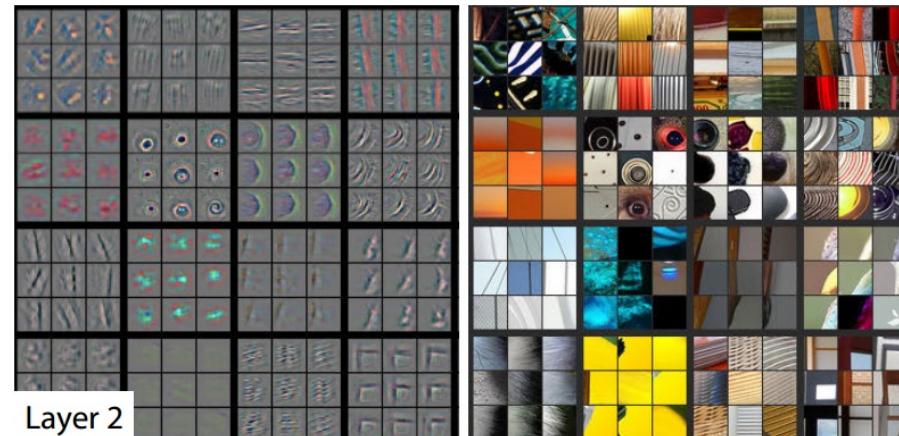
Solution: **Deconvnet; Guided-backpropagation**



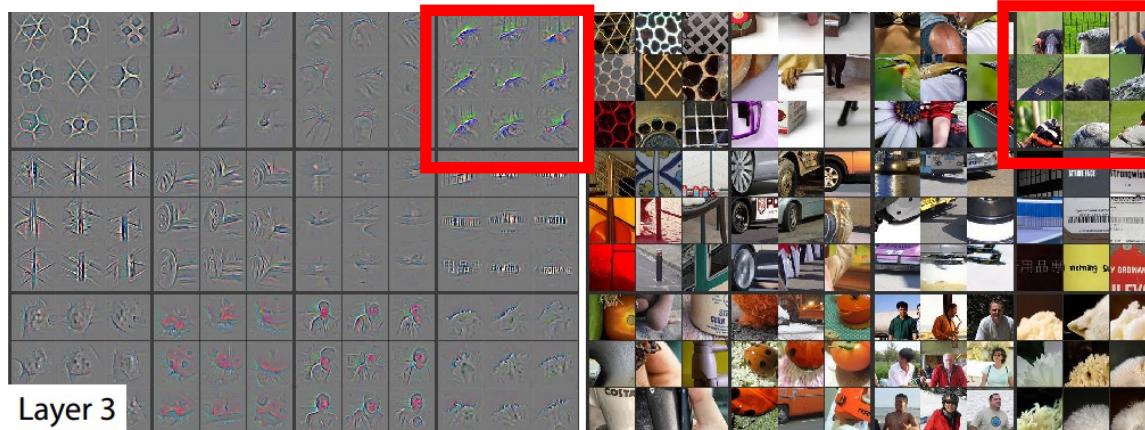
Layer 5

Decovnet [Zeiler-ECCV14]

Deconvolutional Network (Decovnet) can be thought of as a ConvNet that uses the same components (conv, pooling) but in reverse, so instead of mapping pixels to features, deconvnet projects the feature activations back to the input pixel space.



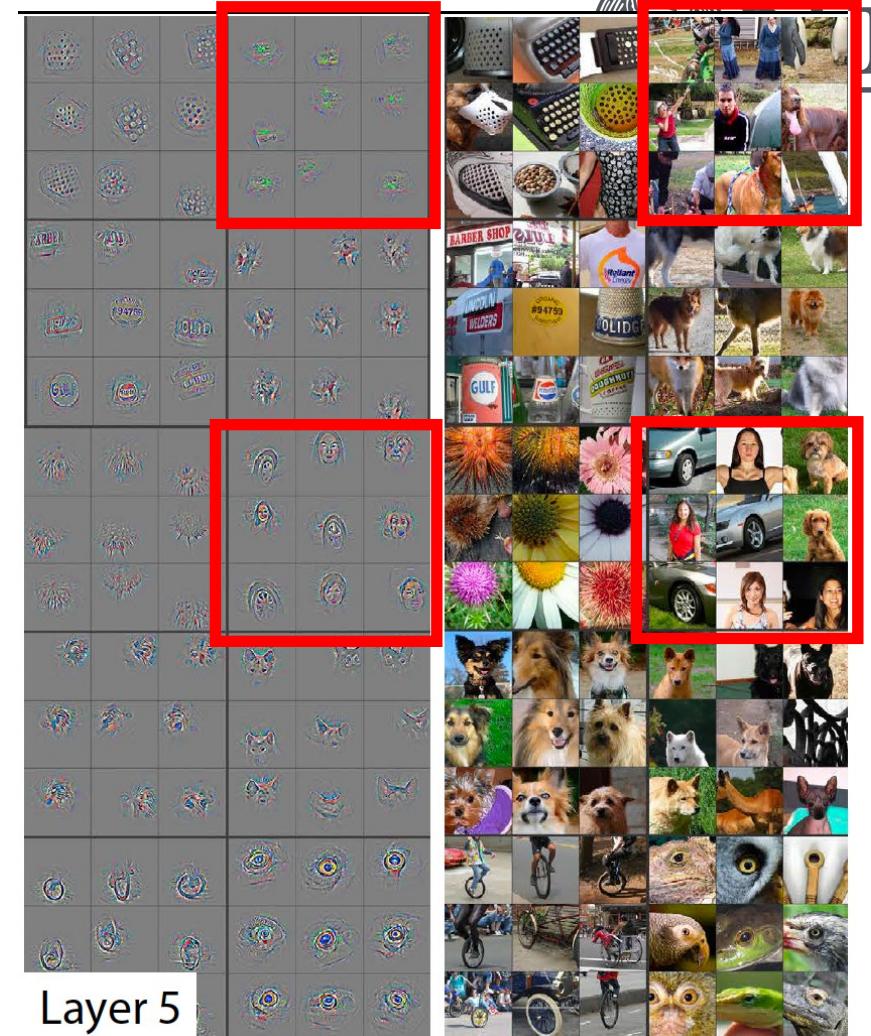
Layer 2



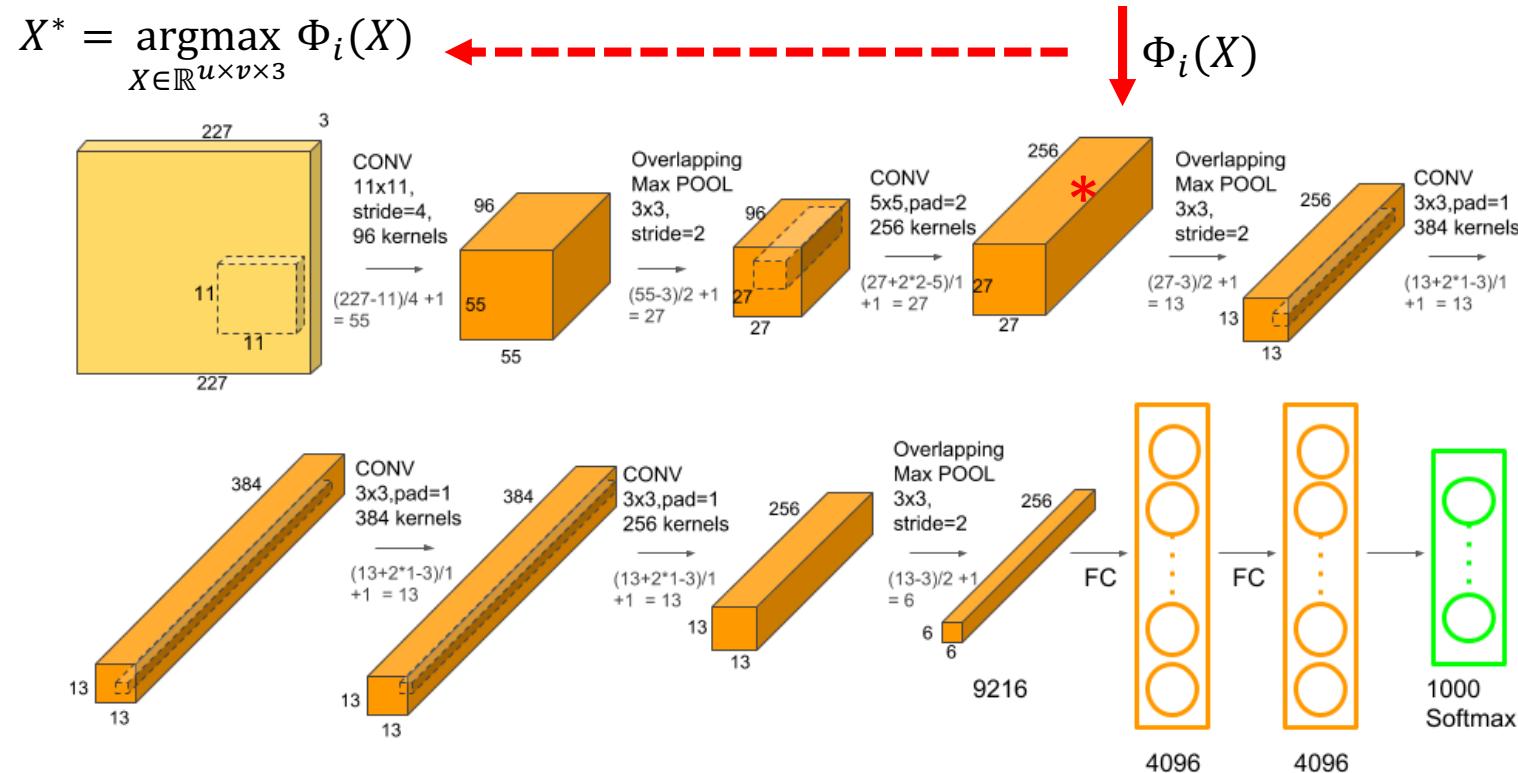
Layer 3

Decovnet [Zeiler-ECCV14]

- Strong grouping within each feature map
- Greater invariance at higher layers
- Exaggeration of discriminative parts of the patch



Synthesize an image to maximize the activation

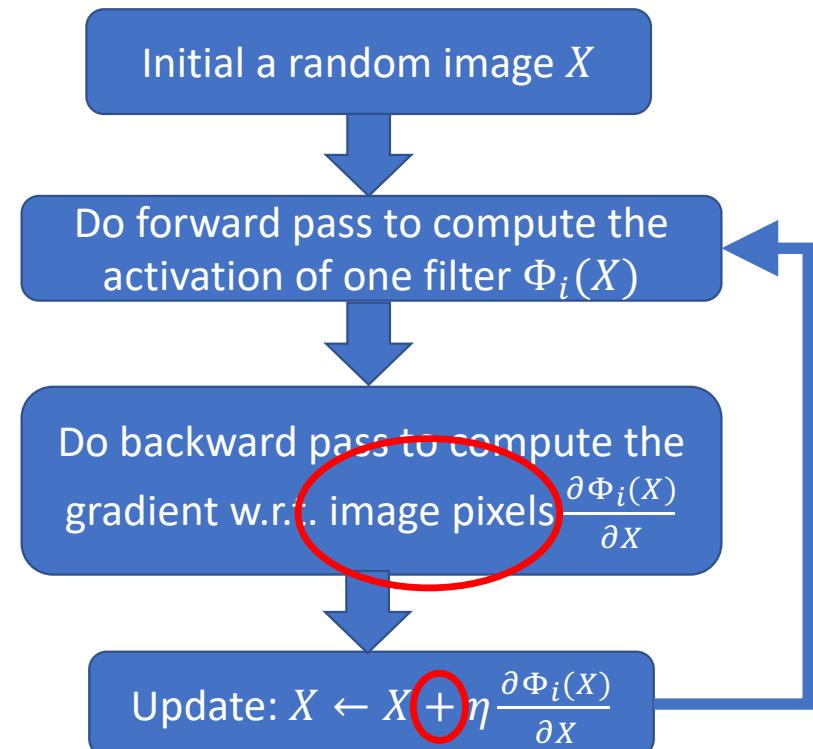


Synthesize an image to maximize the activation

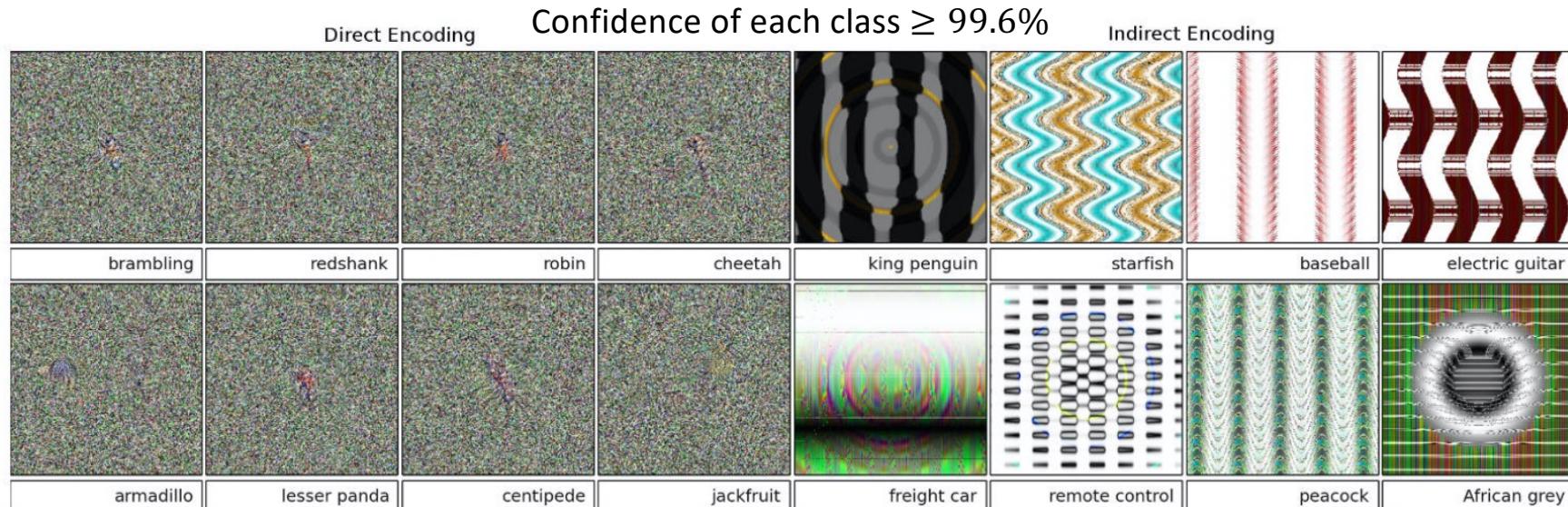
$$X^* = \underset{X \in \mathbb{R}^{u \times v \times 3}}{\operatorname{argmax}} \Phi_i(X)$$

Using **gradient ascent** to generate an image to maximize the activation of one filter.

Q: What is the difference from gradient descent of training neural networks?



Synthesize an image to maximize the activation



But, these unrecognizable images do not help much for understanding filters.

That's because it is lack of **natural image prior**.

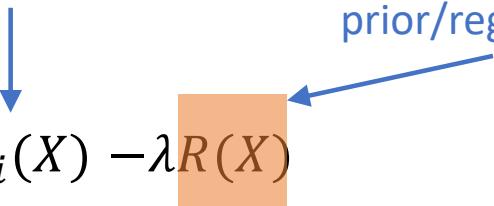
Visualizing filters with regularizations

- To produce more recognizable images, we need

$$X^* = \underset{X \in \mathbb{R}^{u \times v \times 3}}{\operatorname{argmax}} \Phi_i(X) - \lambda R(X)$$

The activation value

Natural image prior/regularizer



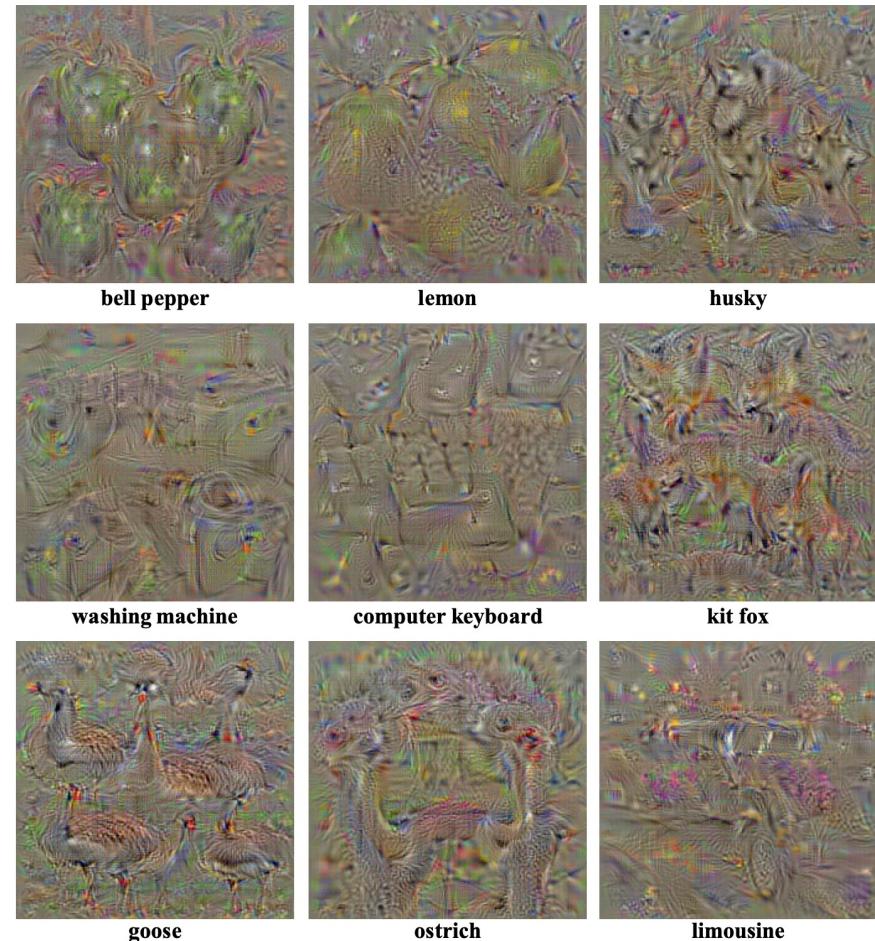
Visualizing filters with regularizations

- L_2 -regularizer [Simonyan-2014]

The score of the class c
(before softmax)

$$I^* = \underset{I \in \mathbb{R}^{H \times W \times 3}}{\operatorname{argmax}} S_c(I) - \lambda \|I\|_2^2$$

L_2 -regularizer



Visualizing filters with regularizations

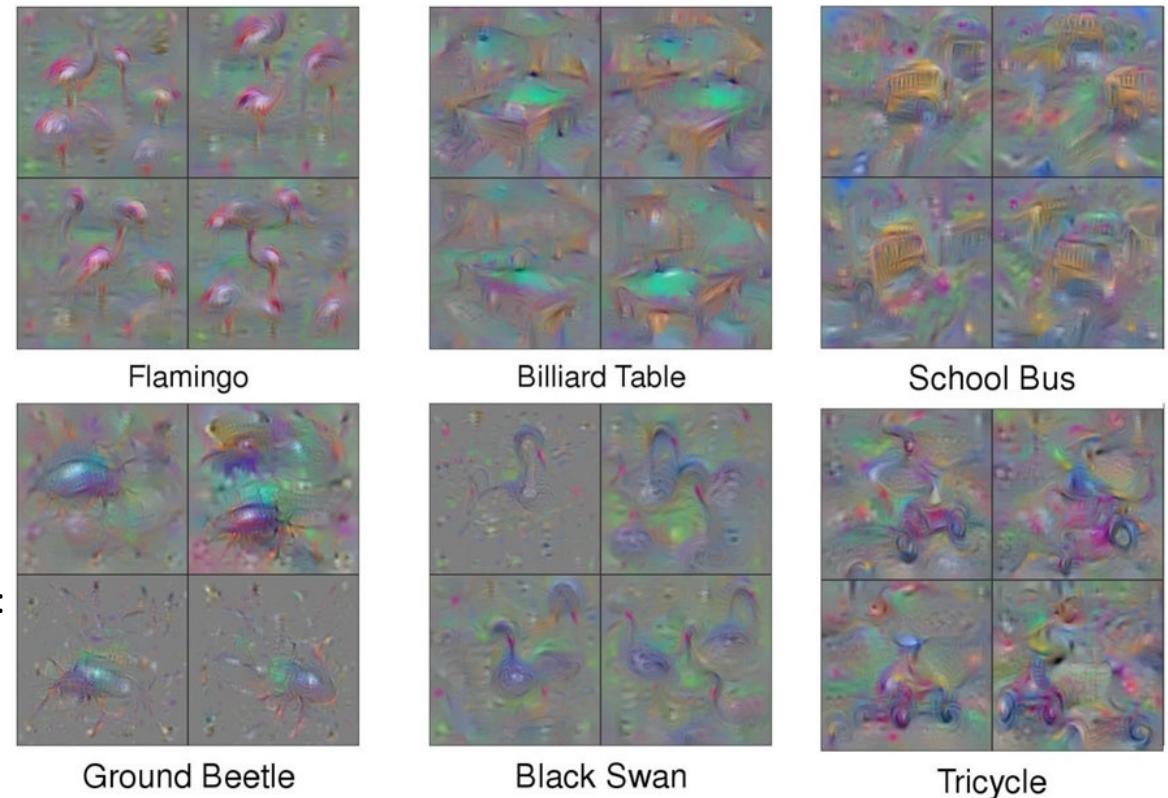
- Better regulariser [Yosinski-2015]

$$X^* = \underset{X \in \mathbb{R}^{u \times v \times 3}}{\operatorname{argmax}} \Phi_i(X) - \lambda \|X\|_2^2$$

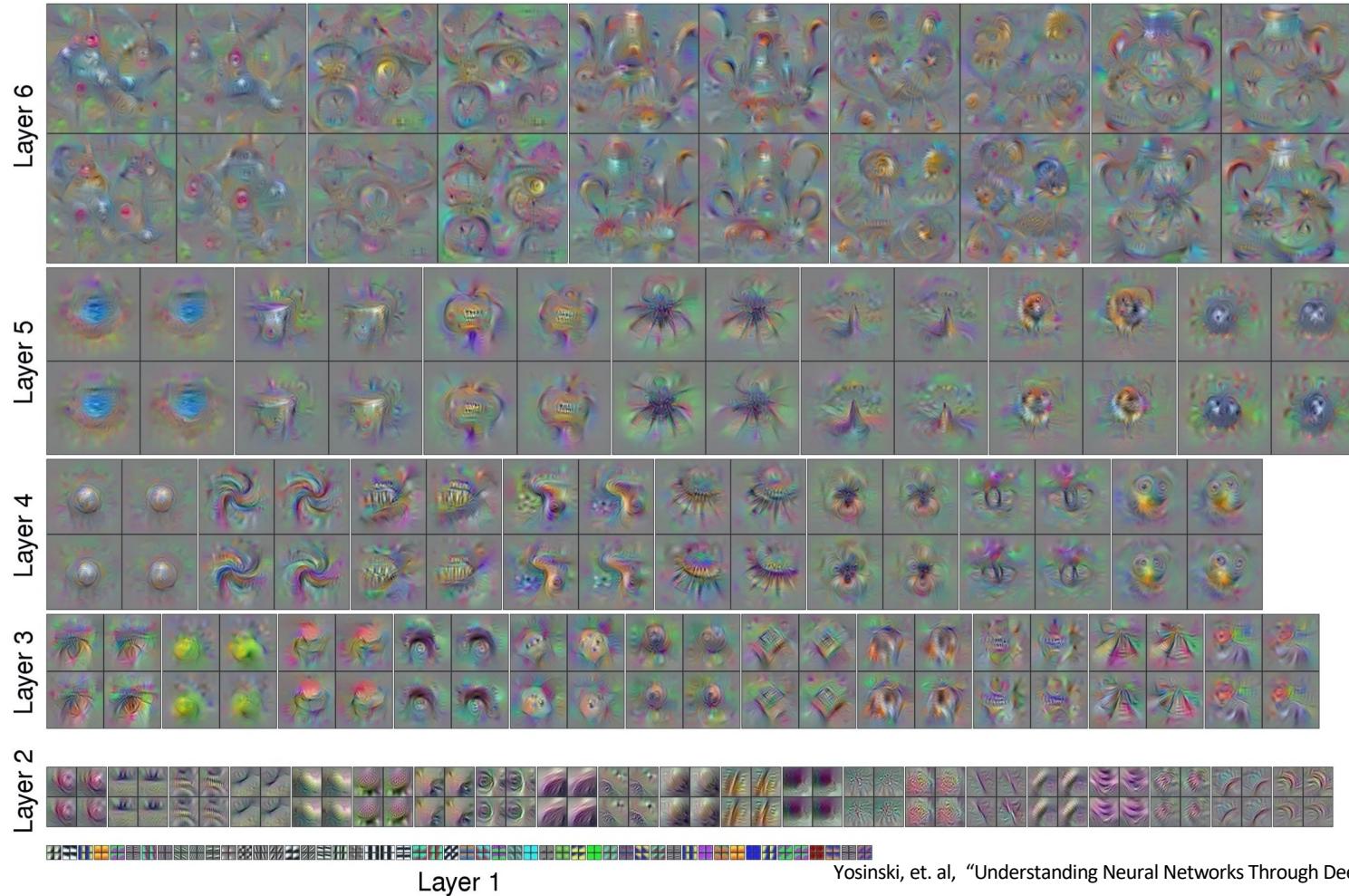
The activation value
↓
 L_2 -regularizer ↑

Not only penalize L2 norm of the image, also periodically:

- Gaussian blur image
- Clip pixels with small values to 0
- Clip pixels with small gradients to 0

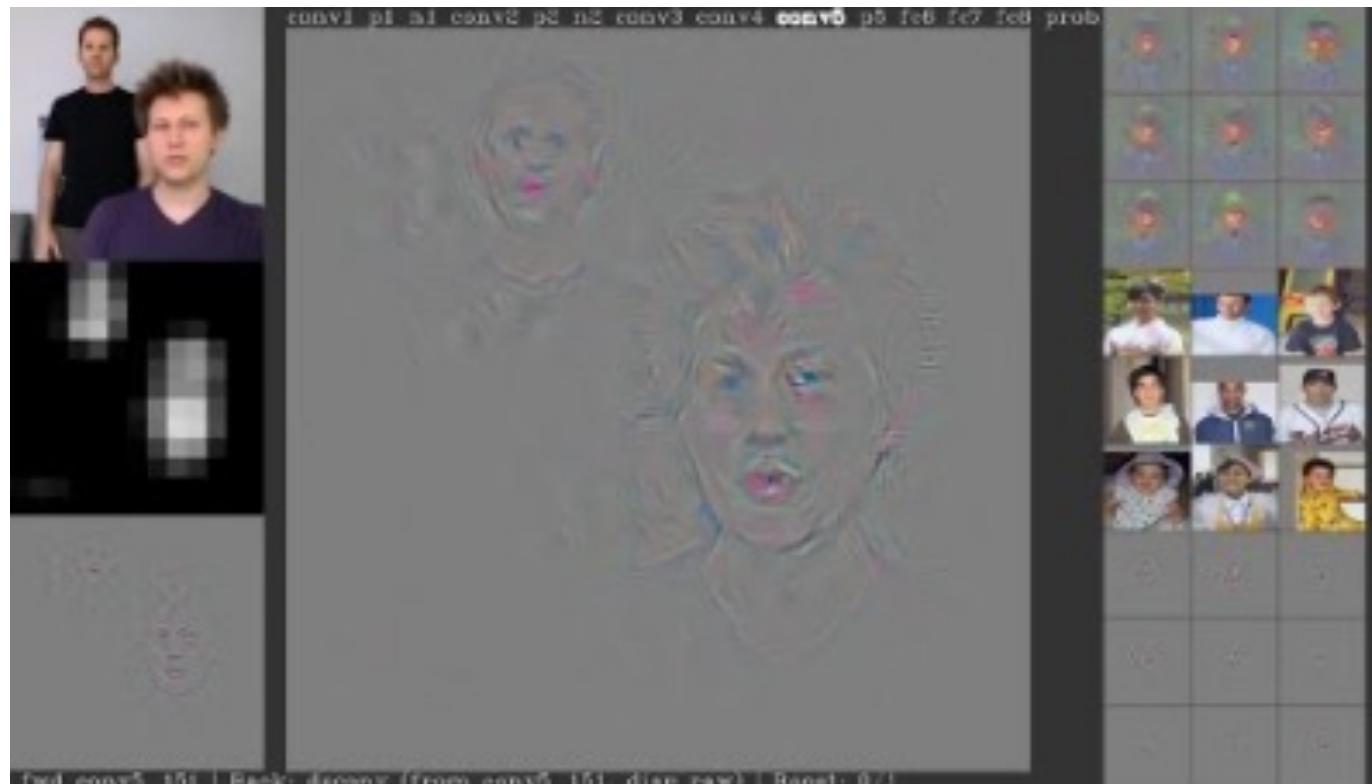


Visualizing filters with regularizations



Yosinski, et. al, "Understanding Neural Networks Through Deep Visualization ", ICLR Workshop 2015.

Deep Visualization Toolbox



Video from: <https://yosinski.com/deepvis>

Summary

Find patches in the dataset



$$X^* = \operatorname{argmax}_{X \in \mathcal{D}} \Phi_i(X)$$

Synthesize patches

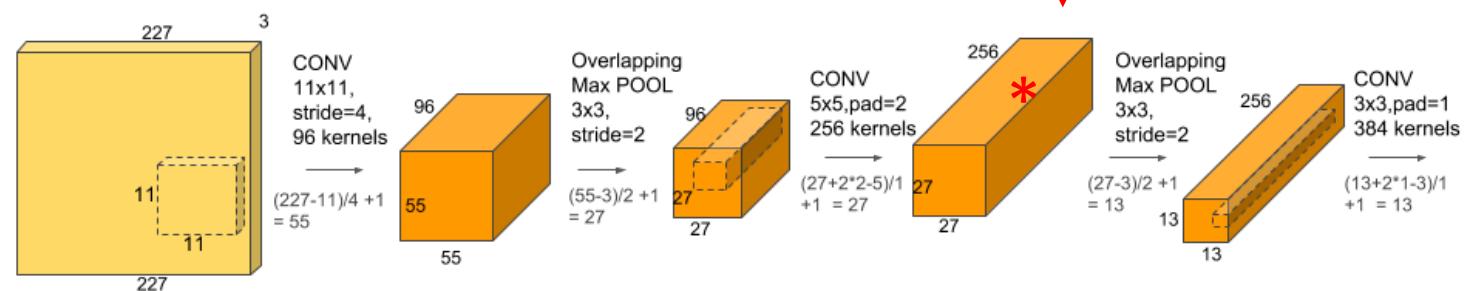


$$X^* = \operatorname{argmax}_{X \in \mathbb{R}^{u \times v \times 3}} \Phi_i(X)$$

Method: Gradient accent

Trick: Adding regularizer

$$X^* = \operatorname{argmax}_X \Phi_i(X)$$



Reference

- Online Lecture: “Visualizing what ConvNets learn”.
<https://cs231n.github.io/understanding-cnn/>
- Video lecture in University of Stanford. “Visualizing and Understanding”. 2017.
<https://www.youtube.com/watch?v=6wcs6szJWMY>
- Blog by Jason Yosinski. “Understanding Neural Networks Through Deep Visualization”. 2015. <https://yosinski.com/deepvis>

Topic 3: Neural Style Transfer

The activation value

$$X^* = \operatorname{argmax}_{X \in \mathbb{R}^{u \times v \times 3}} \Phi_i(X) - \lambda R(X)$$

Natural image regularizer/prior

Visualizing Filters
Visualize what sort of input maximizes the activation of each filter

$$X^* = \operatorname{argmin}_{X \in \mathbb{R}^{H \times W \times 3}} \ell(\Phi(X), \Phi(X_0)) + \lambda R(X)$$

Loss (matches the feature)
 $\Phi(X): \mathbb{R}^{H \times W \times C} \rightarrow \mathbb{R}^d$ is an image representation of a certain layer
 (a feature vector)

Natural image regularizer/prior

Feature Inversion
Reconstruct an image whose feature representation best matches the original image

Feature Inversion

- Reconstruct an image whose feature representation best matches the original image's.

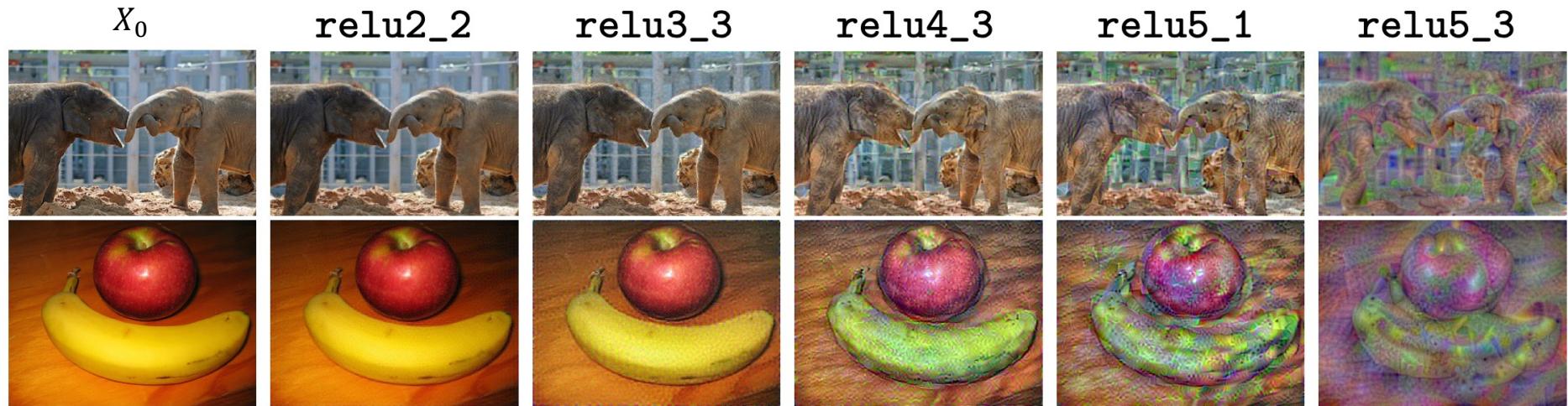
$$X^* = \underset{X \in \mathbb{R}^{H \times W \times 3}}{\operatorname{argmin}} \ell(\Phi(X), \Phi(X_0)) + \lambda R(X)$$

↗
↗

$\ell(\Phi(X), \Phi(X_0)) = \|\Phi(X) - \Phi(X_0)\|^2$
Feature distance

$R_{V^\beta} = \sum_{i,j} \left((x_{i,j+1} - x_{ij})^2 + (x_{i+1,j} - x_{ij})^2 \right)^{\frac{\beta}{2}}$
Total Variation regularizer
(encourages spatial smoothness)

Feature Inversion



As we reconstruct from higher layers, image content and overall spatial structure are preserved, but color, texture, and exact shape are not.

Neural Texture Synthesis

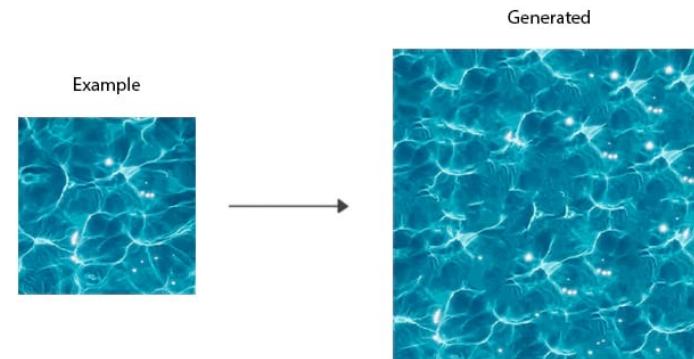
$$X^* = \underset{X \in \mathbb{R}^{H \times W \times 3}}{\operatorname{argmin}} \ell(\Phi(X), \Phi(X_0))$$



$$X^* = \underset{X \in \mathbb{R}^{H \times W \times 3}}{\operatorname{argmin}} \ell(G(X), G(X_0))$$



Texture representation using CNN



Neural Texture Synthesis

Gram Matrix

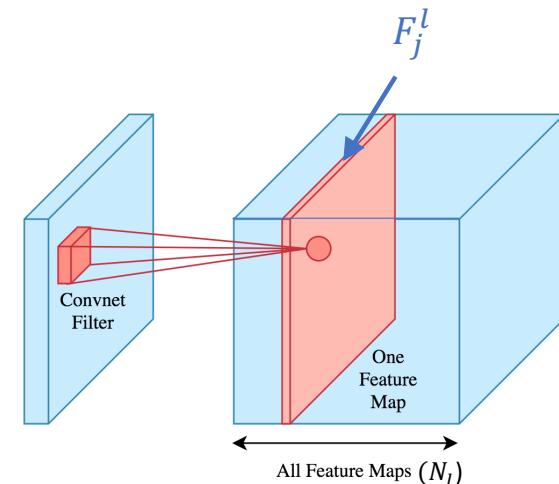
F_j^l : j -th feature map at layer l

F_i^l : i -th feature map at layer l

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l$$

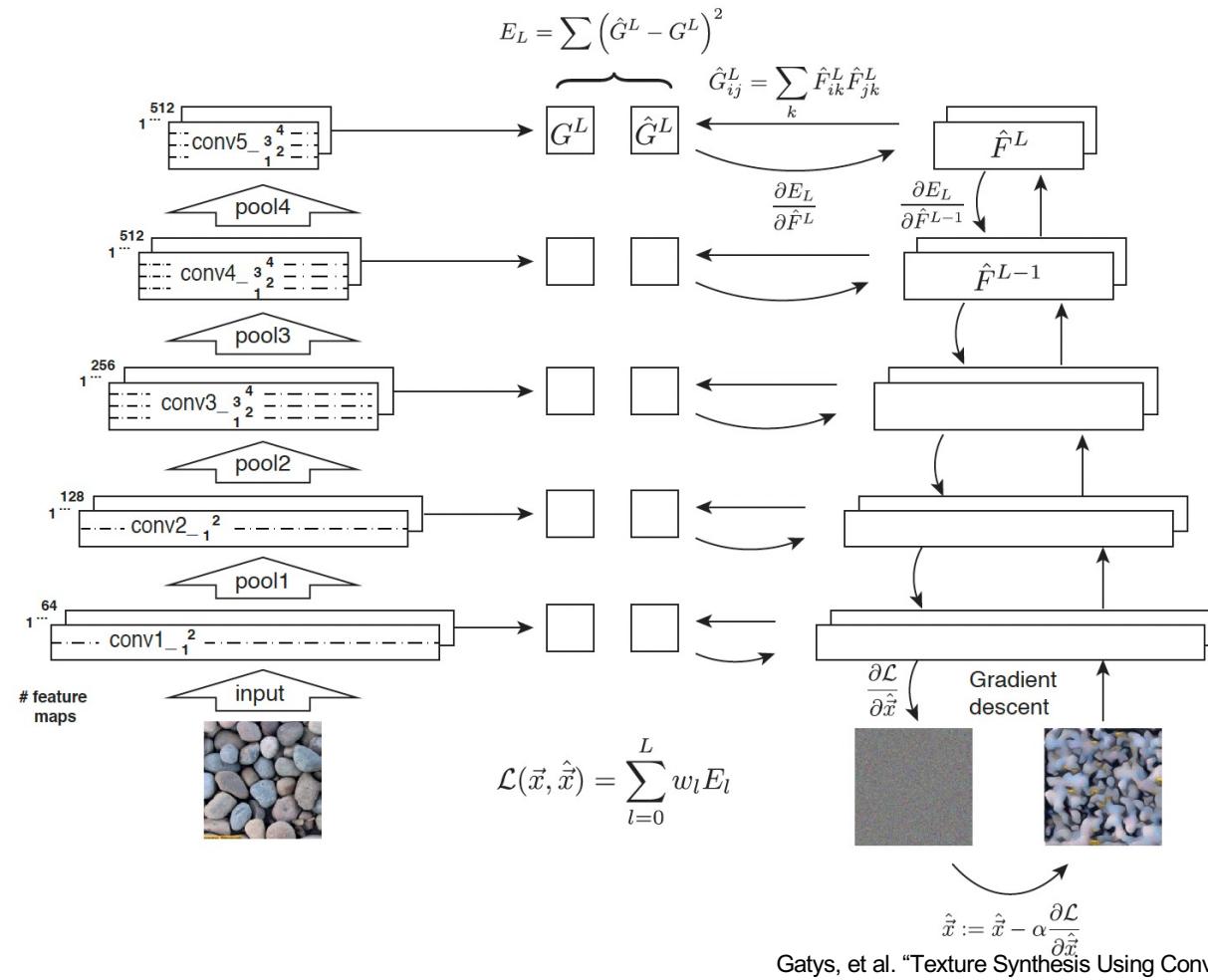
the inner product between two feature maps.

Gram matrix: $G^l \in \mathbb{R}^{N_l \times N_l}$

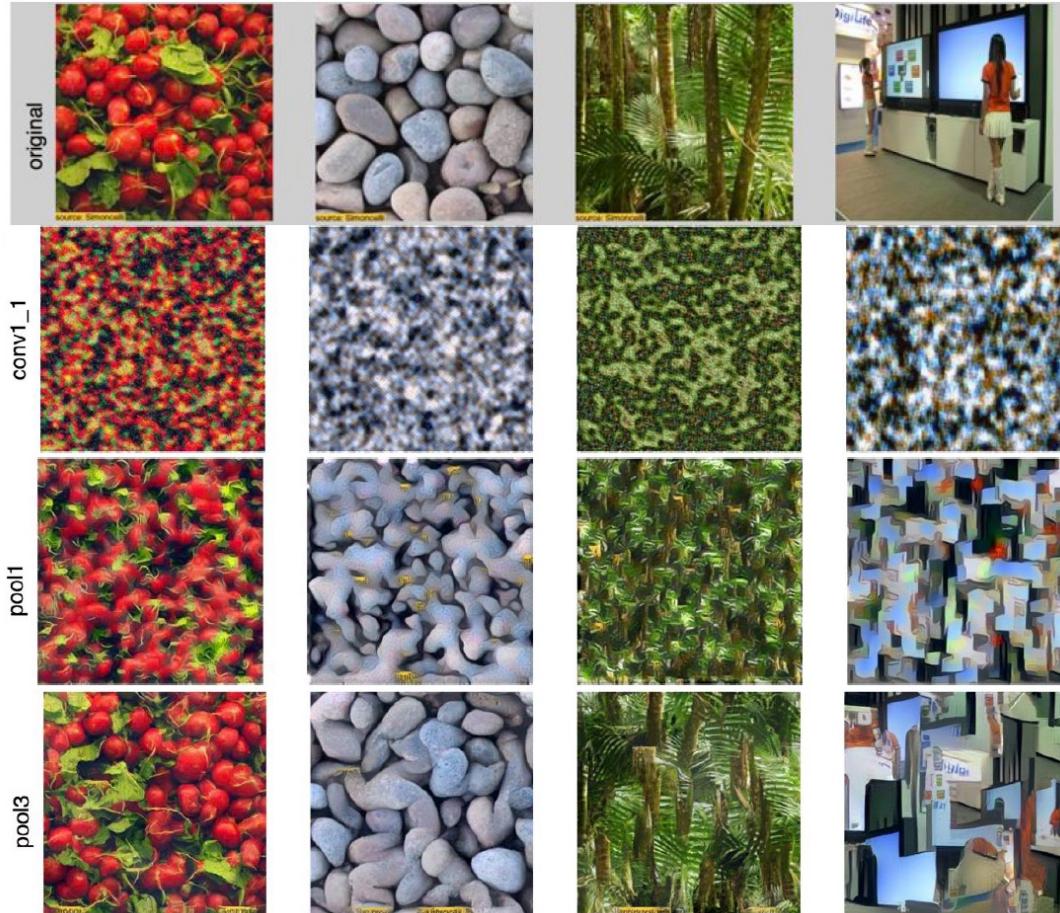


Texture features could be represented by a set of
Gram Matrixes of conv layers: $\{G^1, G^2, \dots, G^L\}$

Neural Texture Synthesis

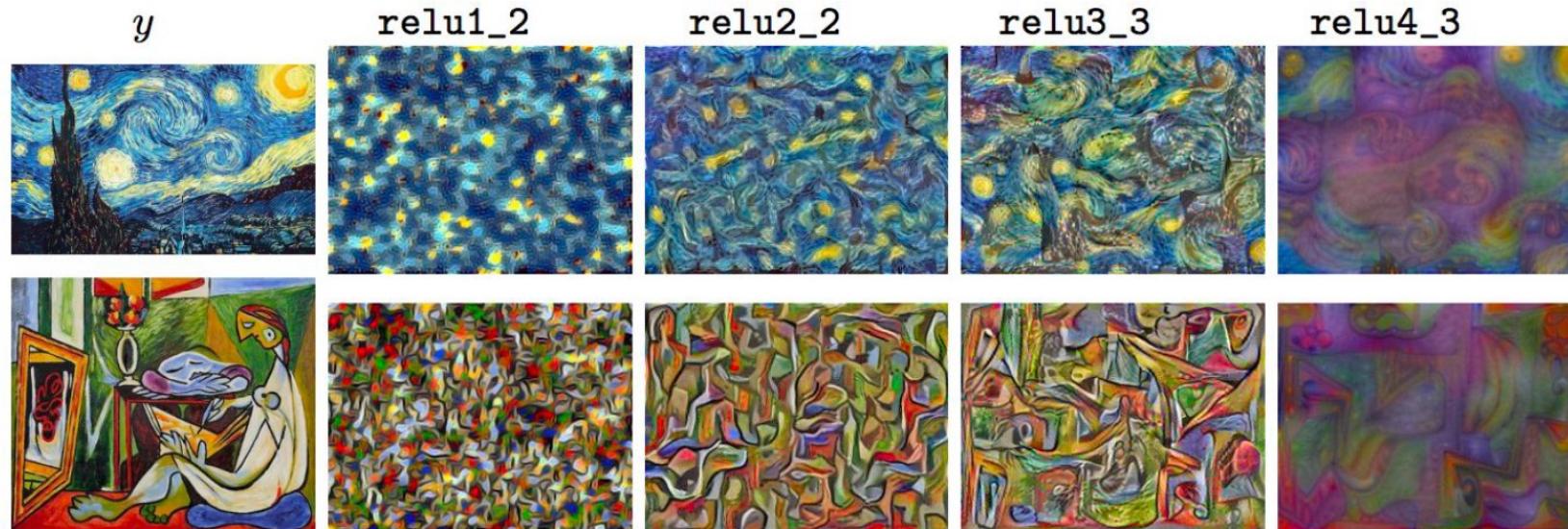


Neural Texture Synthesis



Gatys, et al. "Texture Synthesis Using Convolutional Neural Networks". NIPS2015.

When texture synthesis applies to art ...

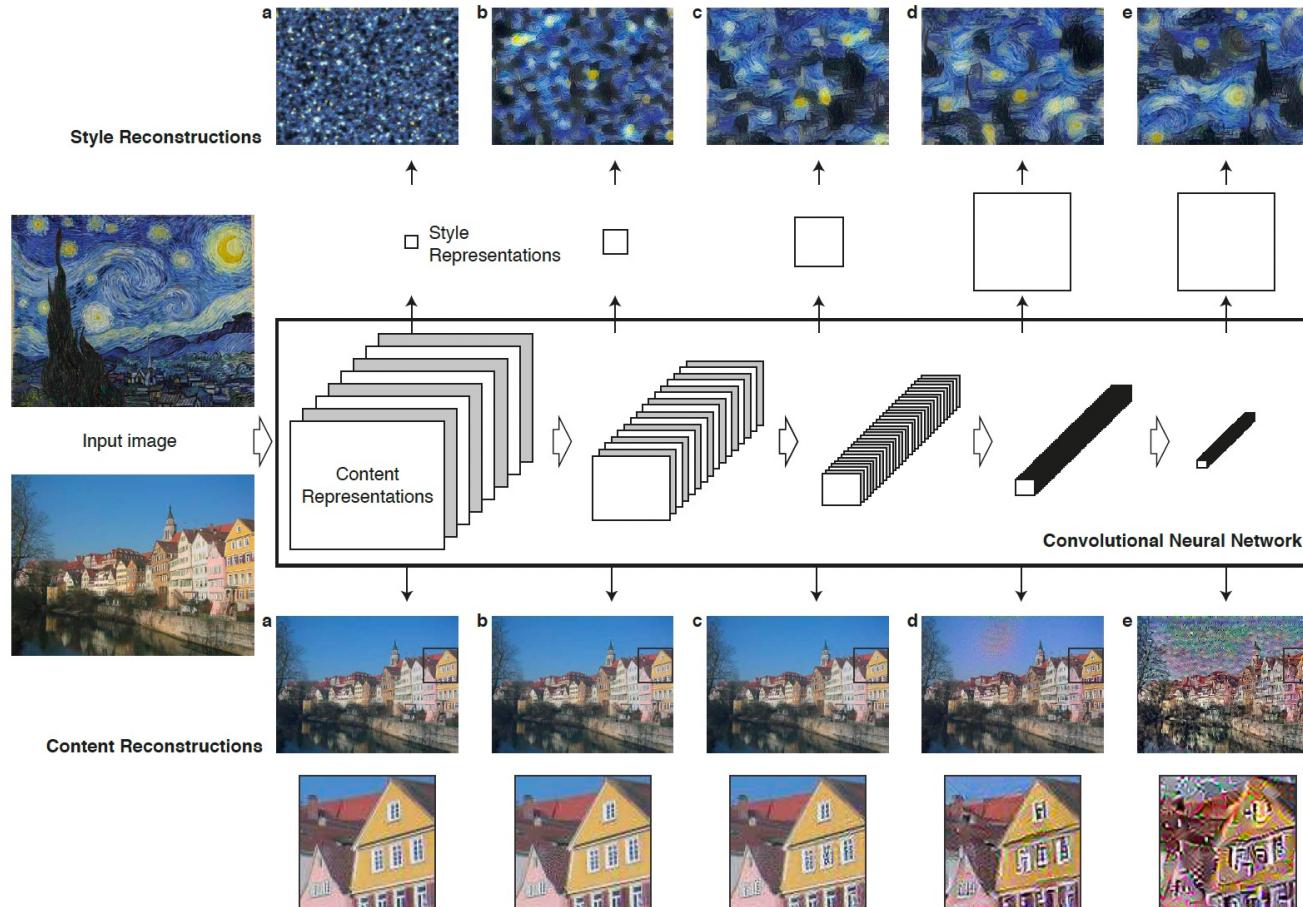


Q: Can we transfer a particular art style to a new image?



Gatys, et al. "Image style transfer using convolutional neural networks", CVPR 2016.
Gatys, et al. "Texture Synthesis Using Convolutional Neural Networks". NIPS2015.

Neural Style Transfer



$$X^* = \underset{X \in \mathbb{R}^{H \times W \times 3}}{\operatorname{argmin}} \ell(G(X), G(X_0))$$

Texture synthesis

Reconstruct an image whose texture representation best matches the original image

$$X^* = \underset{X \in \mathbb{R}^{H \times W \times 3}}{\operatorname{argmin}} \ell(\Phi(X), \Phi(X_0))$$

Feature Inversion

Reconstruct an image whose feature representation best matches the original image

Neural Style Transfer

Content Image



Style Image


 $\mathcal{L}_{\text{style}}$
 A_{ij}^l
 G_{ij}^l

Style Transfer!

=


 F_{ij}^l
 $\mathcal{L}_{\text{content}}$
 P_{ij}^l

Match the texture (Gram-matrix):

$$\sum_{l=1}^L \sum_{ij} (G_{ij}^l - A_{ij}^l)^2$$

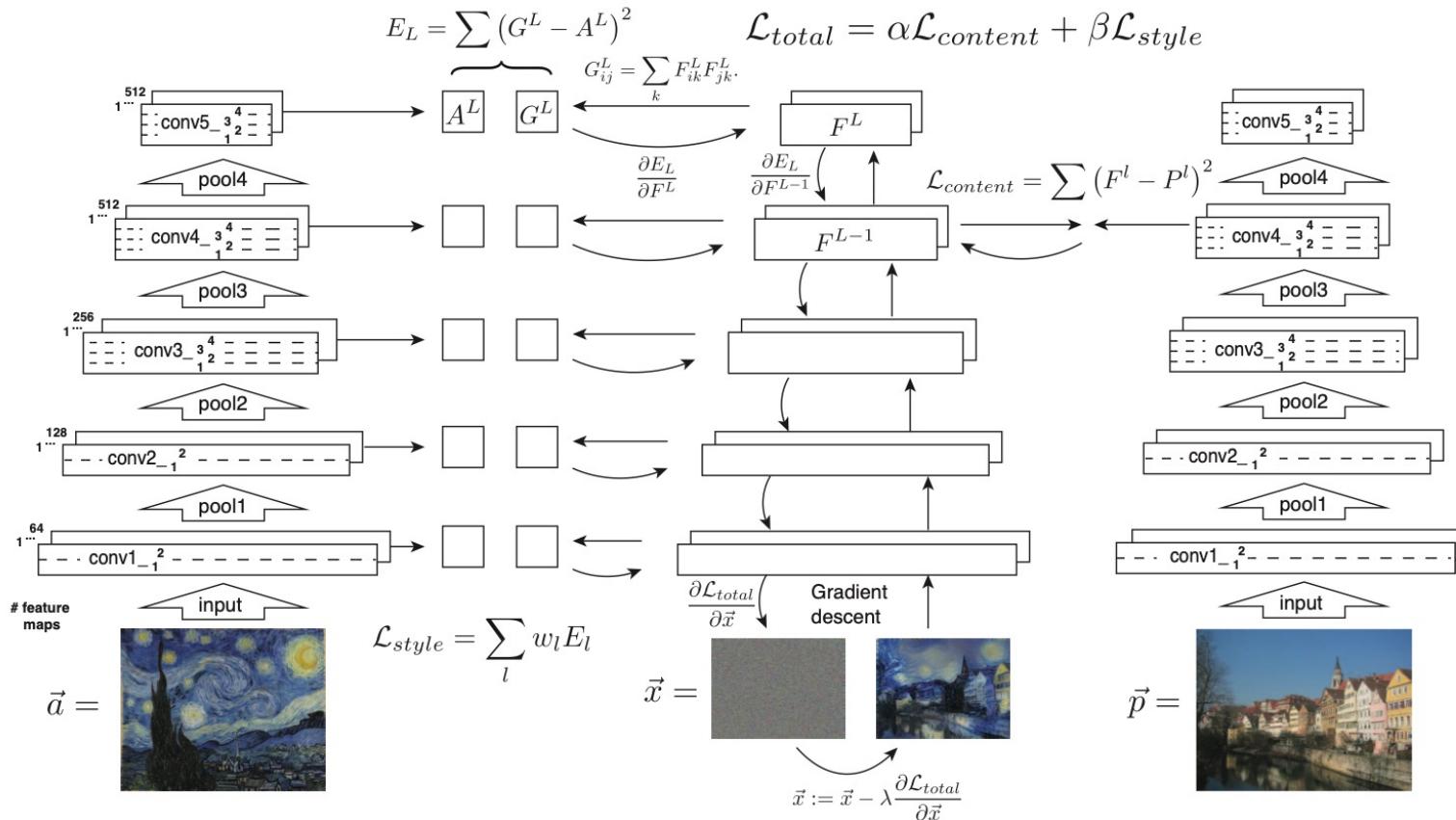
$$\mathcal{L}_{\text{total}} = \alpha \mathcal{L}_{\text{style}} + \beta \mathcal{L}_{\text{content}}$$

Match the feature:

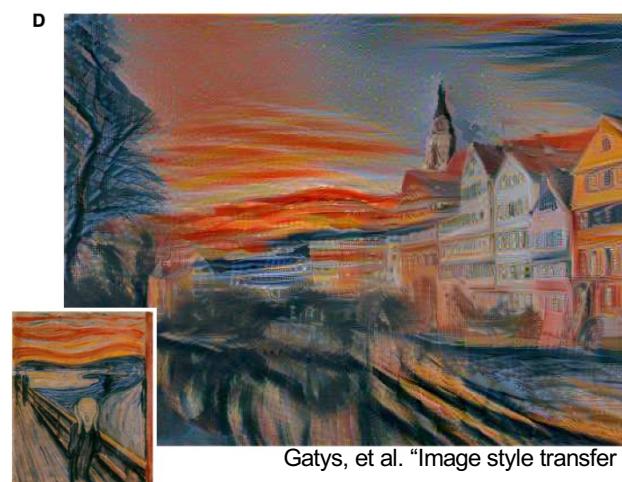
$$\sum_{ij} (F_{ij}^l - P_{ij}^l)^2$$

Gatys, et al. "Image style transfer using convolutional neural networks", CVPR 2016.
Gatys, et al. "Texture Synthesis Using Convolutional Neural Networks". NIPS2015.

Neural Style Transfer



Neural Style Transfer



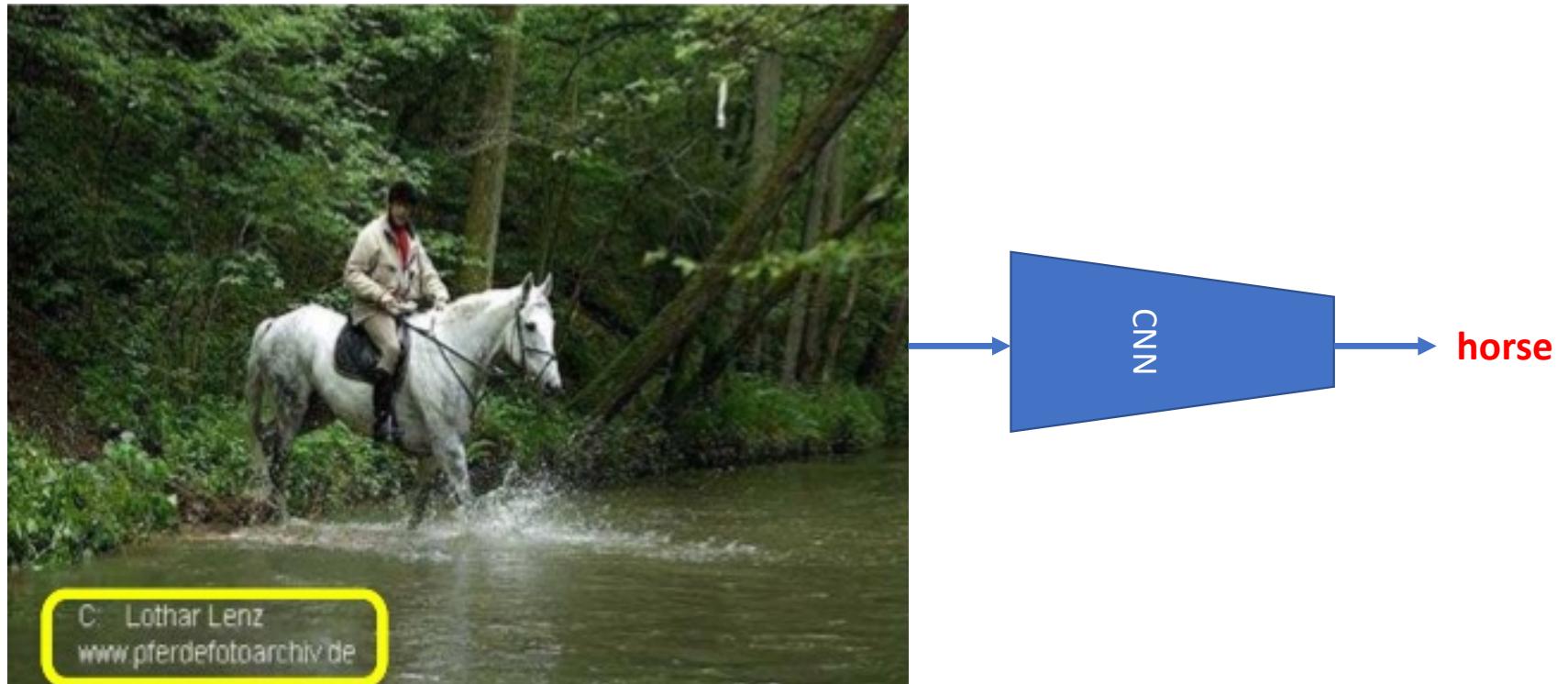
Gatys, et al. "Image style transfer using convolutional neural networks", CVPR 2016.

References

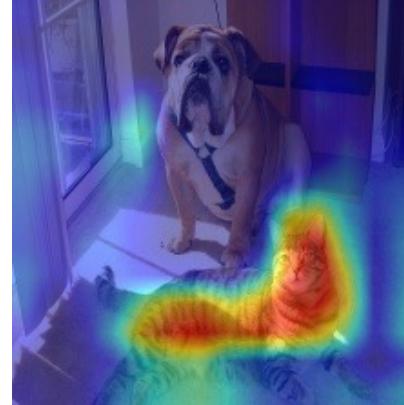
- Video lecture in University of Stanford. “Visualizing and Understanding”. 2017.
<https://www.youtube.com/watch?v=6wcs6szJWMY>

Topic 4: Visualizing Class-Discriminative Heatmap

Q: When a CNN model outputs a class for an image, which parts of the image contribute high to the decision?



Visualizing heat maps



CAM for the cat class

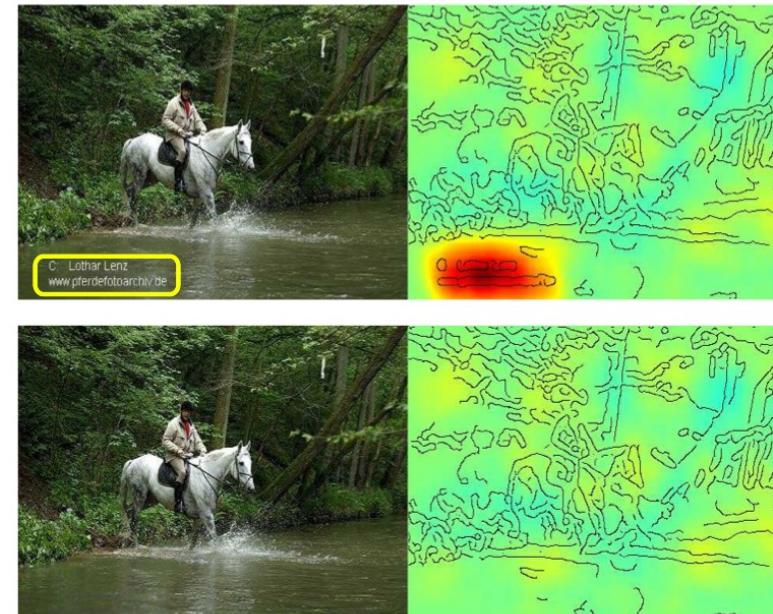
- **Class Activation Mapping (CAM)** is to produce heat maps to highlight class specific region of images, to indicate how important each location is with respect to the given class.

Image from: <https://glassboxmedicine.com/2019/06/11/cnn-heat-maps-class-activation-mapping-cam/>

Why does this matter?

- Help investigate the traits or features of the input images which had major contribution in producing a certain class
- Identify the weak-points of our models
- Identify dataset bias

Horse-picture from Pascal VOC data set



Source tag present

↓
Classified as horse

No source tag present

↓
Not classified as horse

Visualizing heat maps

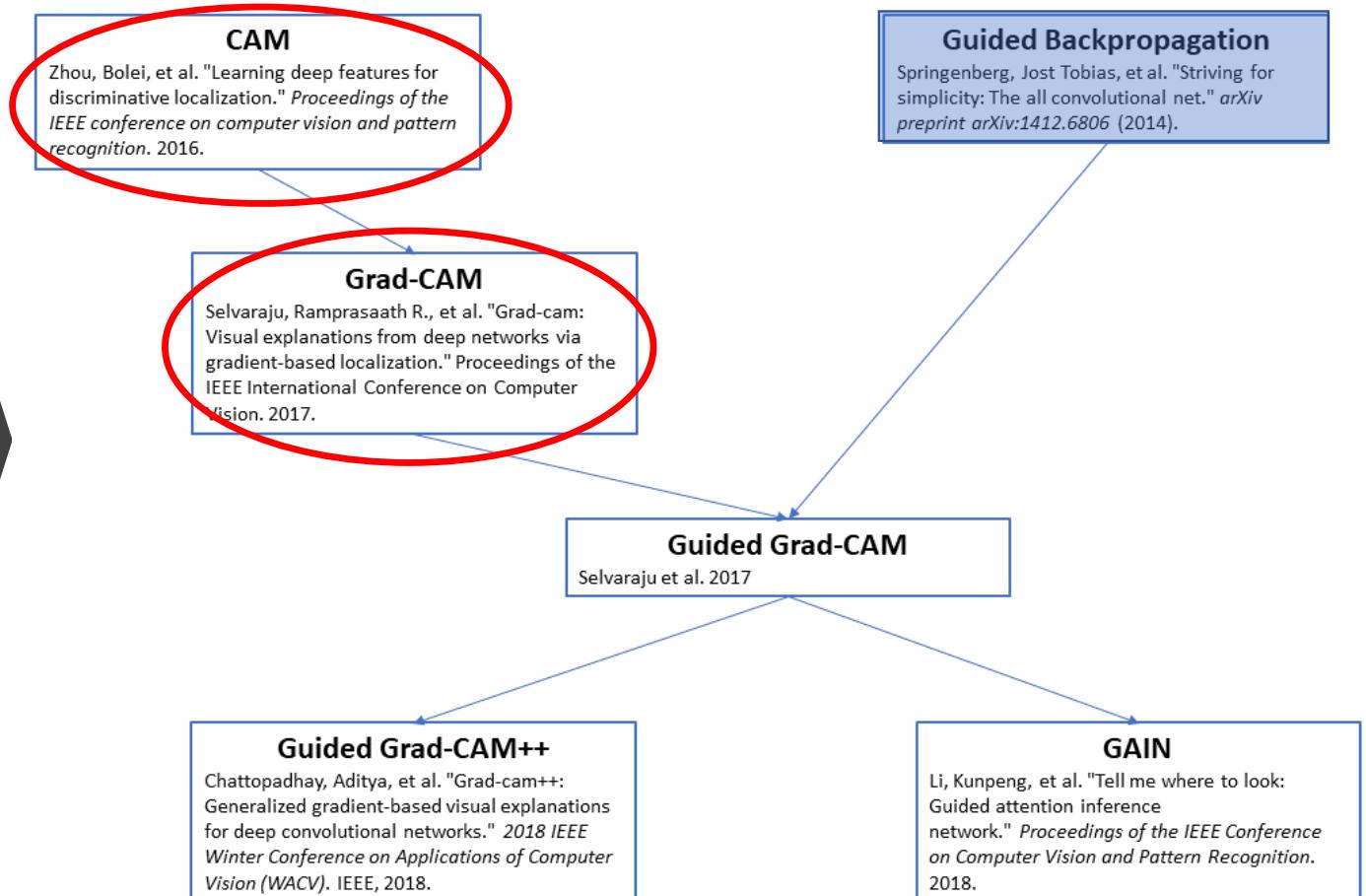
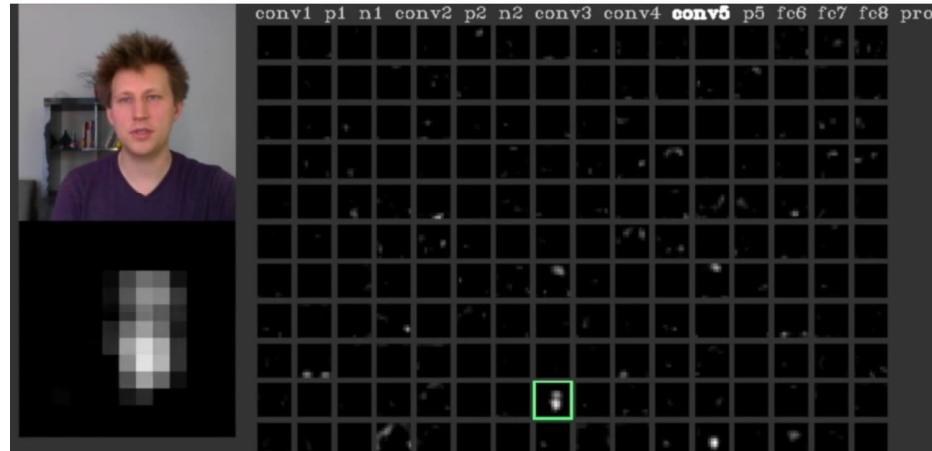
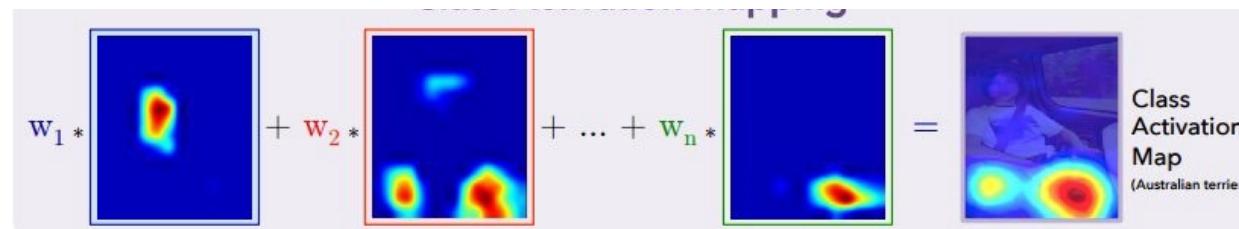


Image from: <https://glassboxmedicine.com/2019/06/11/cnn-heat-maps-class-activation-mapping-cam/>



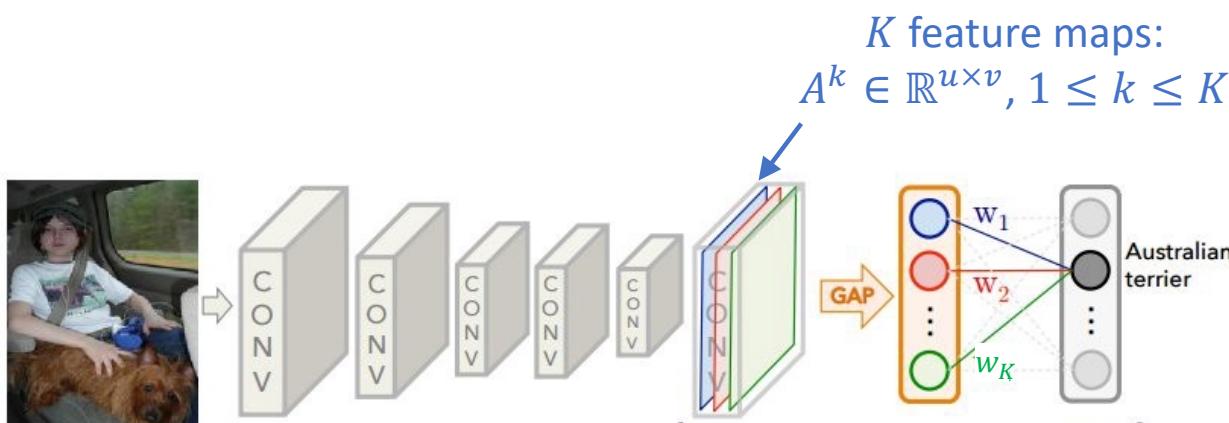
Main principle of CAM and Grad-CAM: Assign a significance value (weight) to each activation map, then sum up the weighted activation maps.

fwd conv5_151 | Back off | Boost: 0/1



Q: How to calculate the weights?

CAM [Zhou-CVPR16]



GAP (Global Average Pooling):

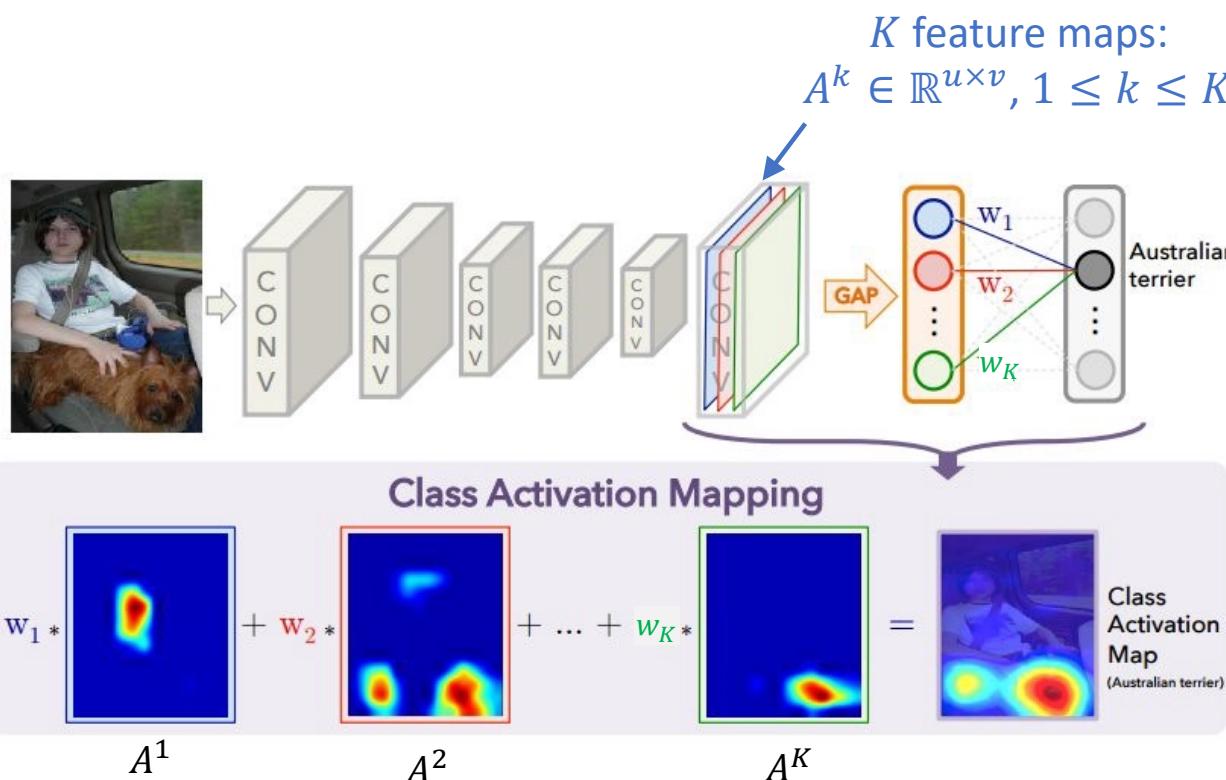
$$\alpha_k = \frac{1}{uv} \sum_{i=1}^u \sum_{j=1}^v A_{i,j}^k$$

Class score for c -th class:

$$\hat{y}^c = \sum_{k=1}^K w_k^c \alpha_k$$

w_k^c is learnt by
backpropagation

CAM [Zhou-CVPR16]



CAM [Zhou-CVPR16]

Mushroom



Penguin



Teapot



Caltech256

Cleaning the floor



Cooking



Fixing a car



Stanford Action40

Zhou et. Al. "Learning deep features for discriminative localization" CVPR 2016.

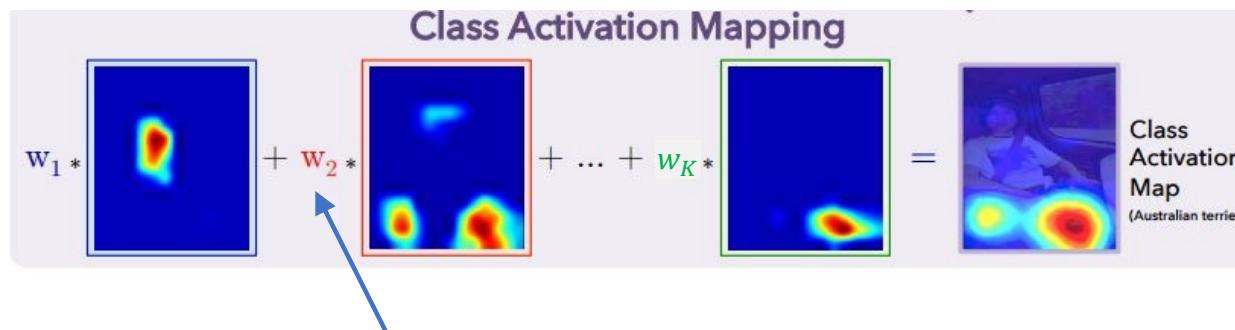
Limitations of CAM

- Need to have a GAP layer in the model architecture
 - Only applicable to this particular kind of CNN architecture (conv layers → GAP → FC layer (with Softmax)), leading to inferior accuracies.
- Can only visualize the final layer heatmap

Solution: A generalization of CAM: **Grad-CAM**
Compatible with any kind of CNN architecture

Grad-CAM [Selvaraju-ICCV17]

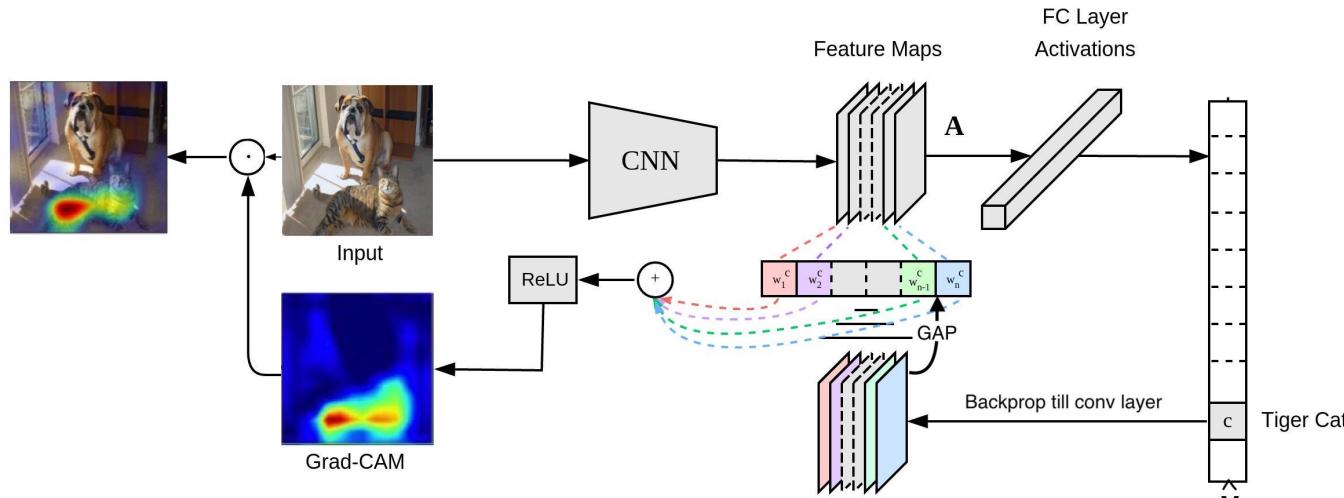
- **Intuition:** Gradients indicates its importance for making a decision by a model.



Q: How to calculate w_k ?

A: Using average of gradients w.r.t. a certain feature map.

Grad-CAM [Selvaraju-ICCV17]



K gradient maps:
 $\frac{\partial \hat{y}}{\partial A^k} \in \mathbb{R}^{u \times v}, 1 \leq k \leq K$

Gradient via backpropagation

$$\frac{\partial \hat{y}^c}{\partial A_{i,j}^k}$$

GAP (Global Average Pooling):

$$w_k^c = \frac{1}{uv} \sum_{i=1}^u \sum_{j=1}^v \frac{\partial \hat{y}^c}{\partial A_{i,j}^k}$$

Grad-CAM for c -th class:

$$L_{\text{Grad-CAM}}^c = \text{ReLU} \left(\sum_{k=1}^K w_k^c A^k \right)$$

Beyond classification

- Grad-CAM can be also used for other CNN-based models.

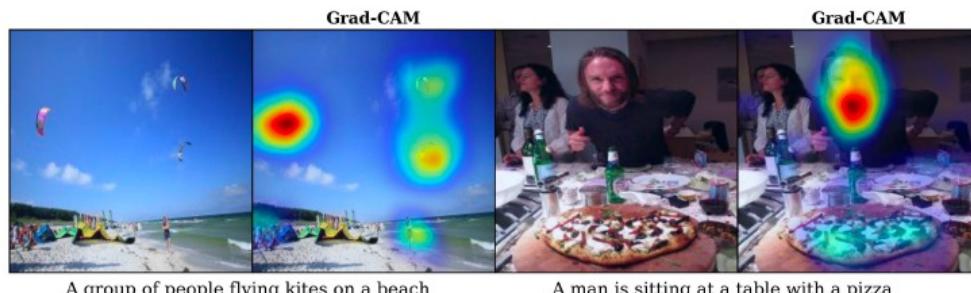
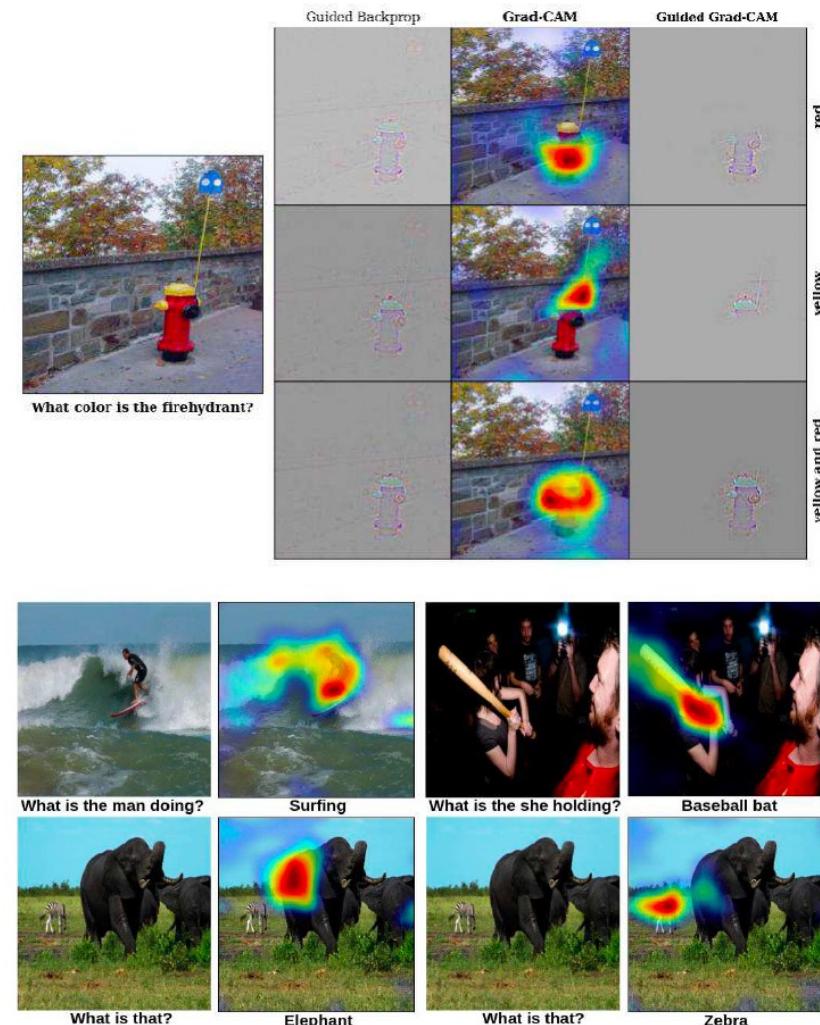


Image captioning

More on Grad-CAM Demo: <http://gradcam.cloudcv.org>

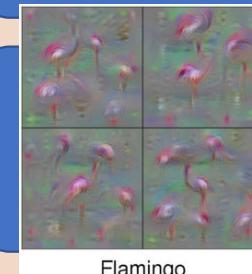
VQA (Visual Question Answering)



Visualizing Filters

Display
Filters
directly

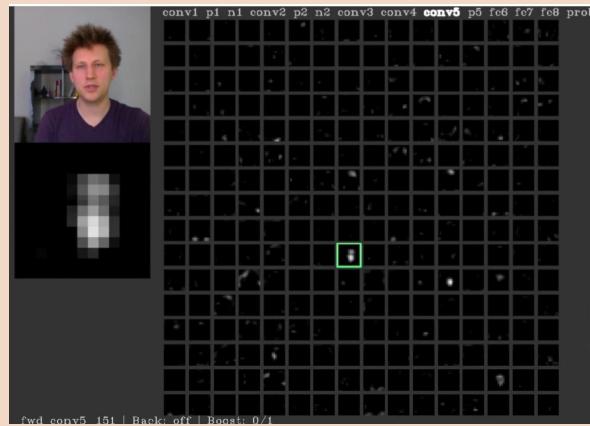
Find patches in a
dataset



Synthesize patches

Feature Inverse
Texture synthesis
Neural Style Transfer

Visualizing Activation Maps



Visualizing Heatmaps

CAM, Grad-CAM, ...



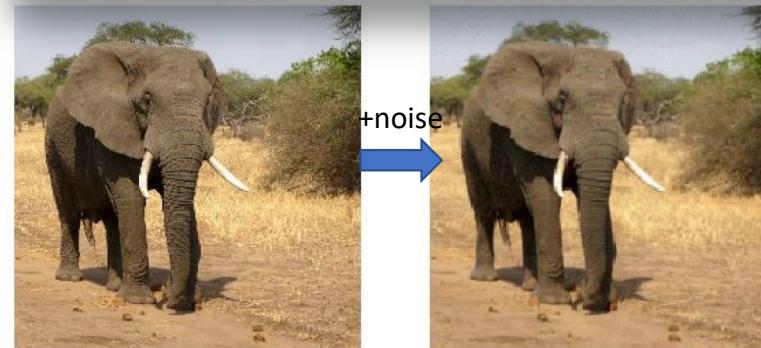
Benefits of visualizing CNN

- Helps to understand why a model they build works or does not work.
- Helps to debug models and improve training and testing strategies.
- Helps to do architectural changes or build new models
- Helps to identify dataset bias

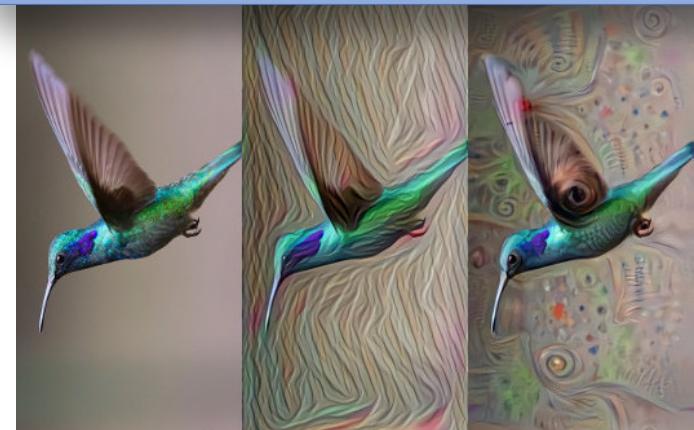
Style transfer



Fooling images



Deep Dream



<https://deeplearning.net/deepdreamgenerator.com>

More than just visualization ...

References

- Blog by Rachel Draelos. “CNN Heat maps: class activation mapping (GAM)”. <https://glassboxmedicine.com/2019/06/11/cnn-heat-maps-class-activation-mapping-cam/>
- Blog by Anil Matcha. “Class activation maps: visualizing neural network decision-making”, 2019. <https://heartbeat.comet.ml/class-activation-maps-visualizing-neural-network-decision-making-92efa5af9a33>
- Blog by Shubham Panchal. “Grab-CAM: A camera for your model’s decision”, 2021. <https://towardsdatascience.com/grad-cam-camera-for-your-models-decision-1ef69aae8fe7>