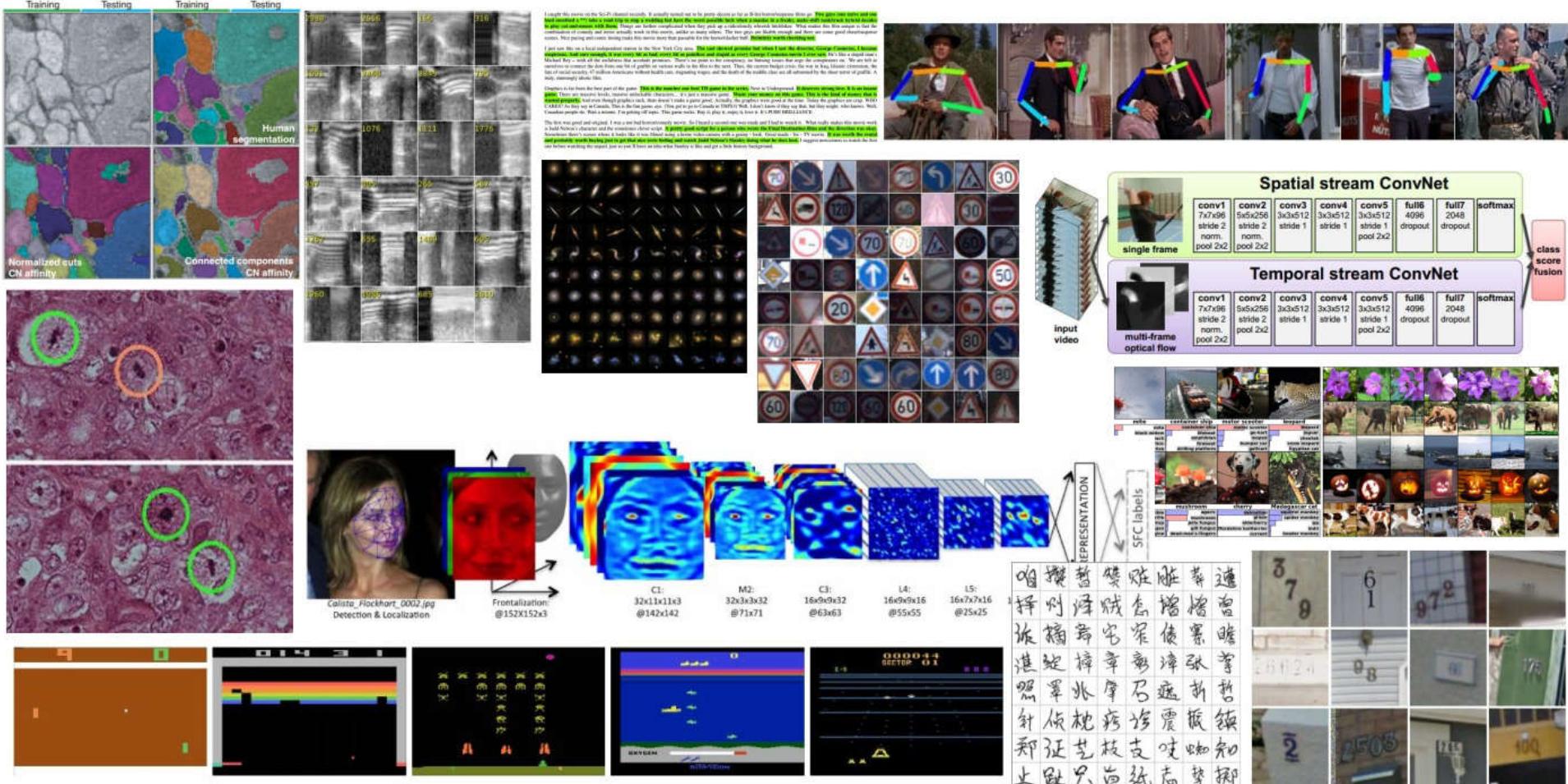


Lecture 09 – Understanding and Visualizing Convolutional Neural Networks

박사과정 김성빈 chengbinjin@inha.edu,
지도교수 김학일 교수 hikim@inha.ac.kr
인하대학교 컴퓨터비전 연구실



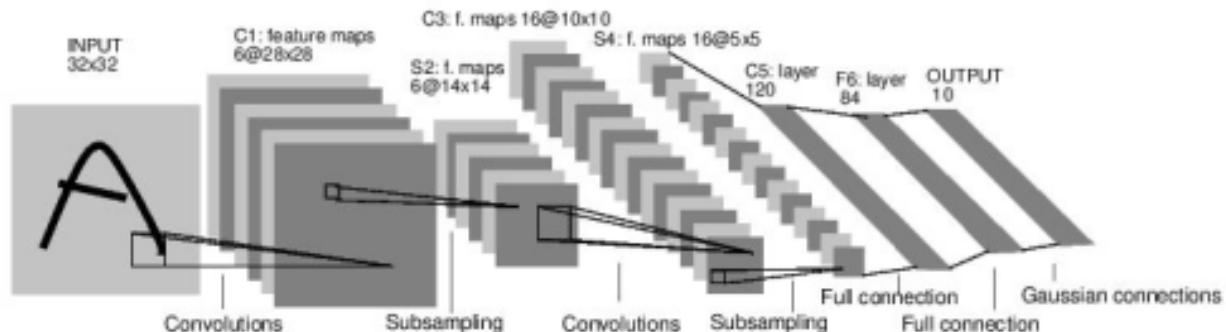
CNNs in a Wide Variety of Applications



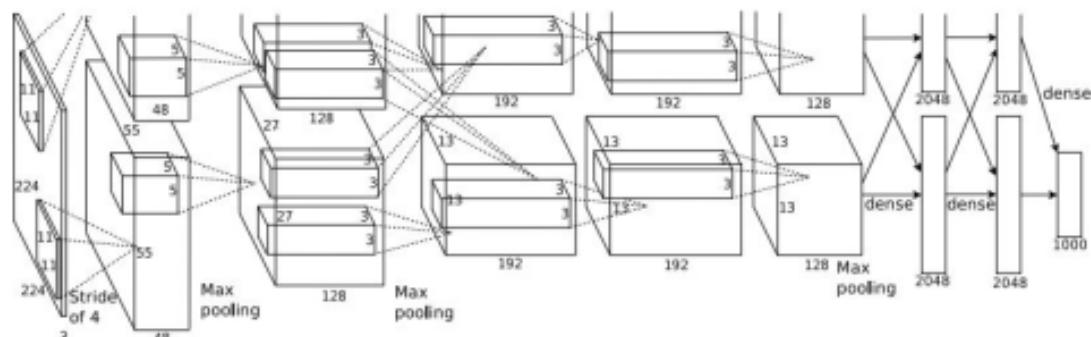
ConvNets Case Study

Case Studies

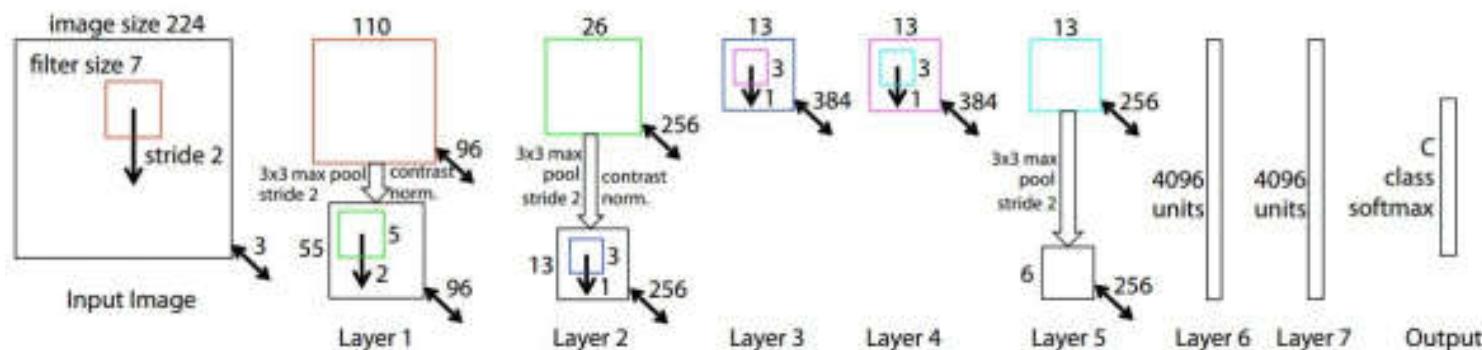
LeNet
 (1998)



AlexNet
 (2012)



ZFNet
 (2013)



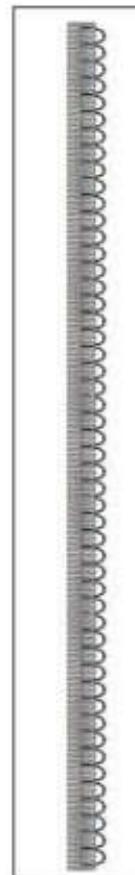
ConvNets Case Study

Case Studies

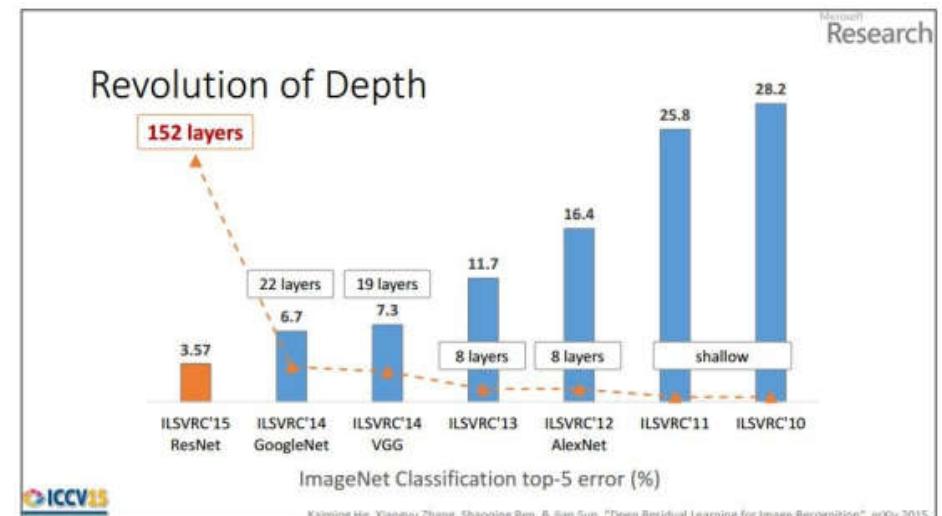
D	E
16 weight layers	19 weight layers
conv3-64	conv3-64
conv3-64	conv3-64
conv3-128	conv3-128
conv3-128	conv3-128
conv3-256	conv3-256
conv3-256	conv3-256
conv3-256	conv3-256
conv3-256	conv3-256
conv3-512	conv3-512
conv3-512	conv3-512
conv3-512	conv3-512
conv3-512	conv3-512
conv3-512	conv3-512
conv3-512	conv3-512
maxpool	
FC-4096	
FC-4096	
FC-1000	
soft-max	



VGG
(2014)

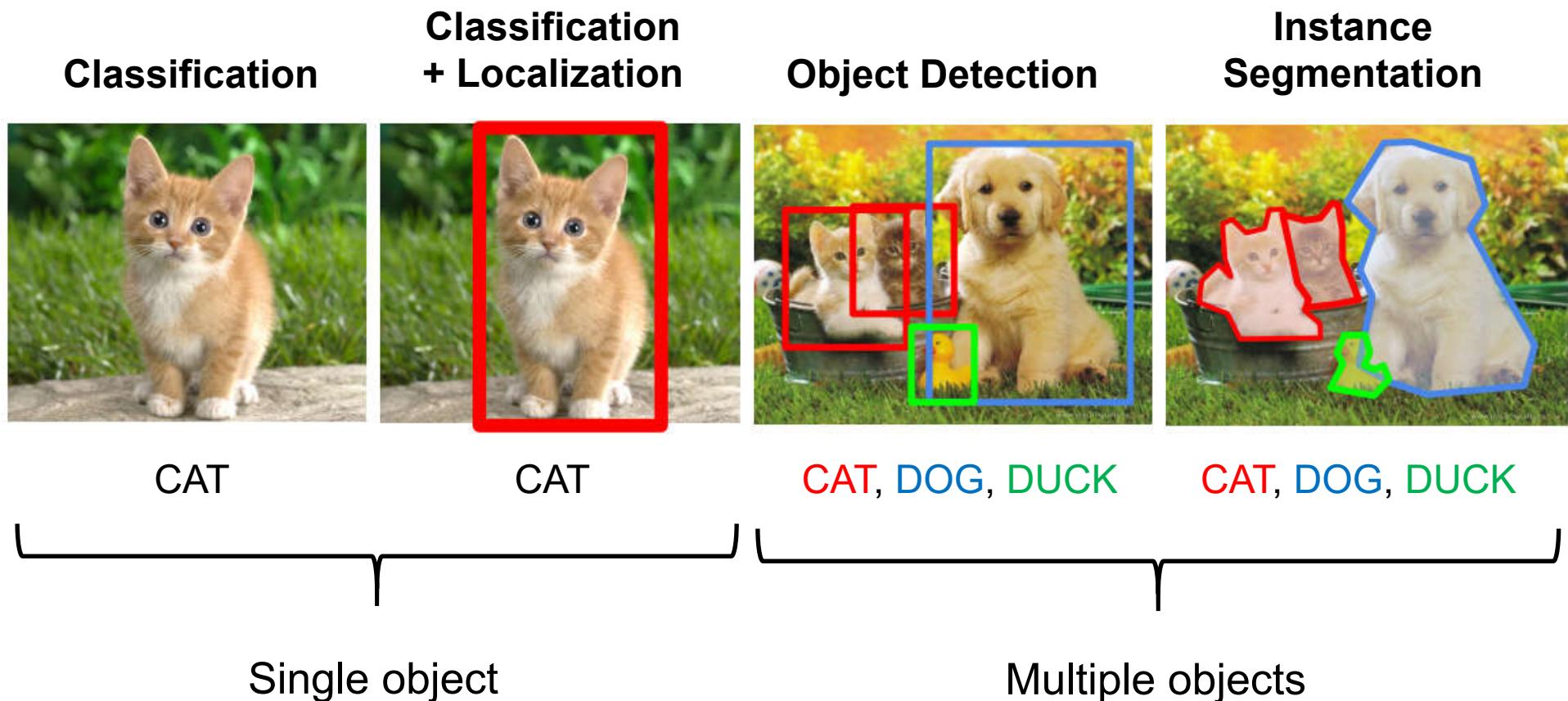


GoogLeNet
(2014)



ResNet
(2014)

Computer Vision Tasks



Understanding ConvNets

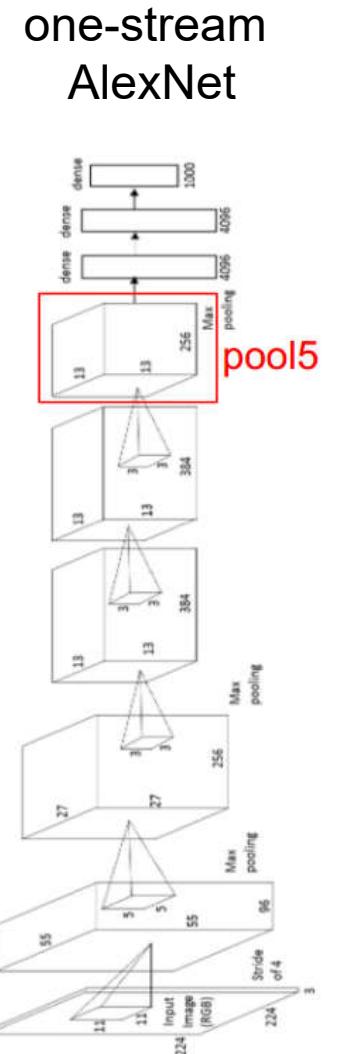
- I. Visualize patches that maximally activate neurons
- II. Visualize the weights
- III. Visualize the representation space (e.g. with t-SNE)
- IV. Occlusion experiments
- V. Deconv approaches (single backward pass)
- VI. Optimization over image approaches (optimization)

I. Visualize Patches that Maximally Activate Neurons

Visualize Patches that Maximally Activate Neurons



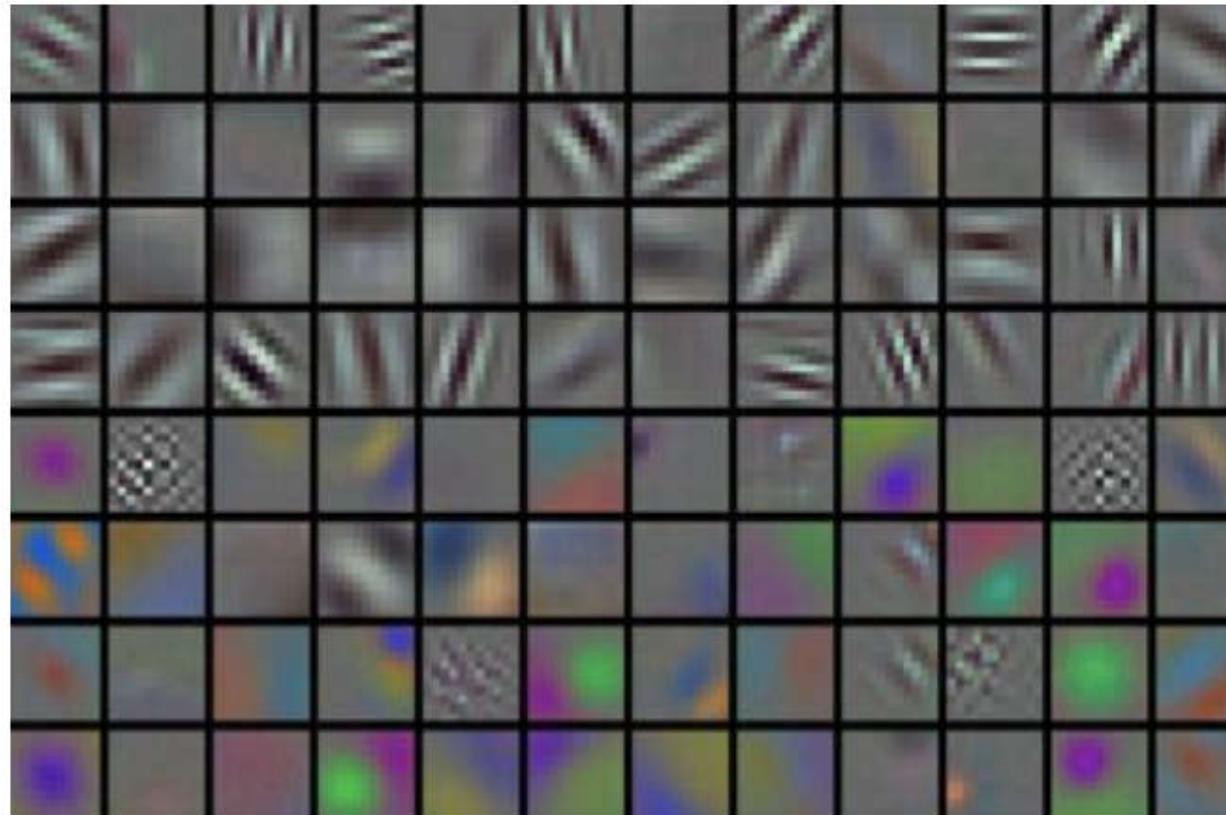
Figure 4: Top regions for six pool_5 units. Receptive fields and activation values are drawn in white. Some units are aligned to concepts, such as people (row 1) or text (4). Other units capture texture and material properties, such as dot arrays (2) and specular reflections (6).



Girshick, et al., Rich feature hierarchies for accurate object detection and semantic segmentation, **CVPR2014**

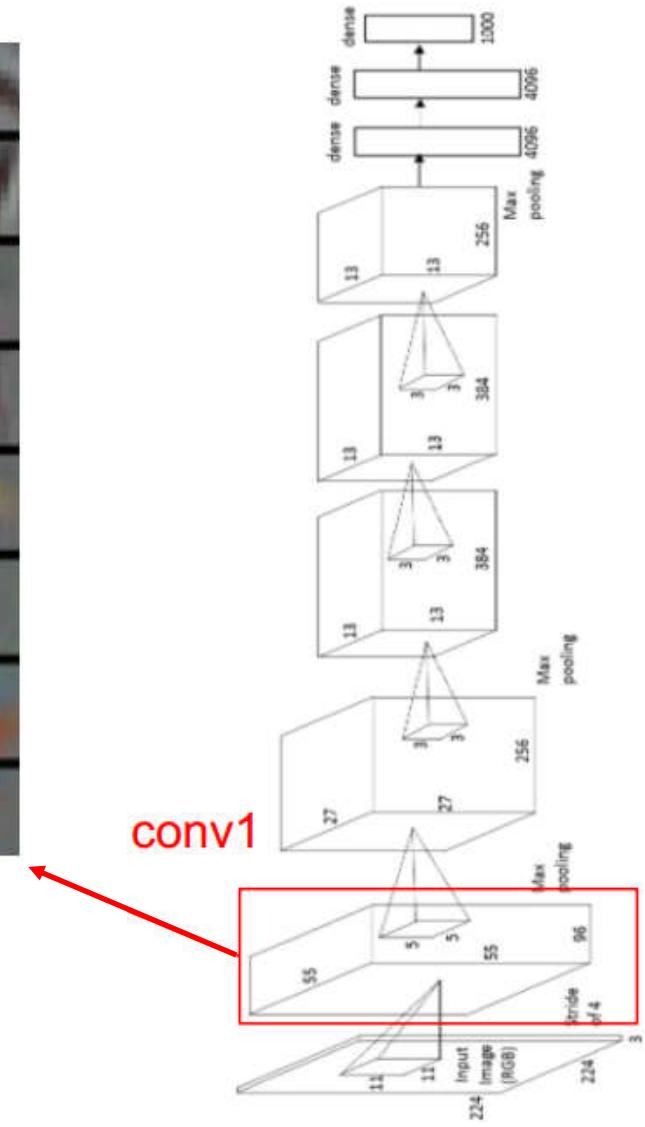
II. Visualize the Weights

Visualize the Filters/Kernels (Raw Weights)



Different rotated Gabor like features

- Only interpretable on the first layer ☹



Visualize the Filter/Kernels (Raw Weights)

You can still do it for higher layers, it's just not that interesting

(these are taken from ConvNetJS CIFAR-10 demo)

This doesn't make as much sense and it's not as nice to look at



Layer 1 weights



Layer 2 weights

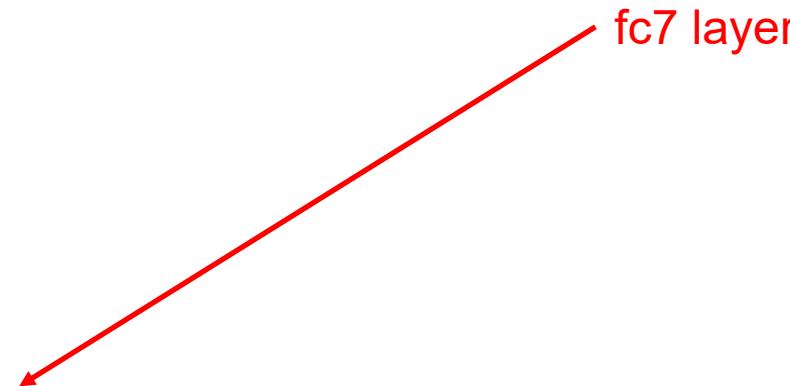


Layer 3 weights

III. Visualize the Representation Space

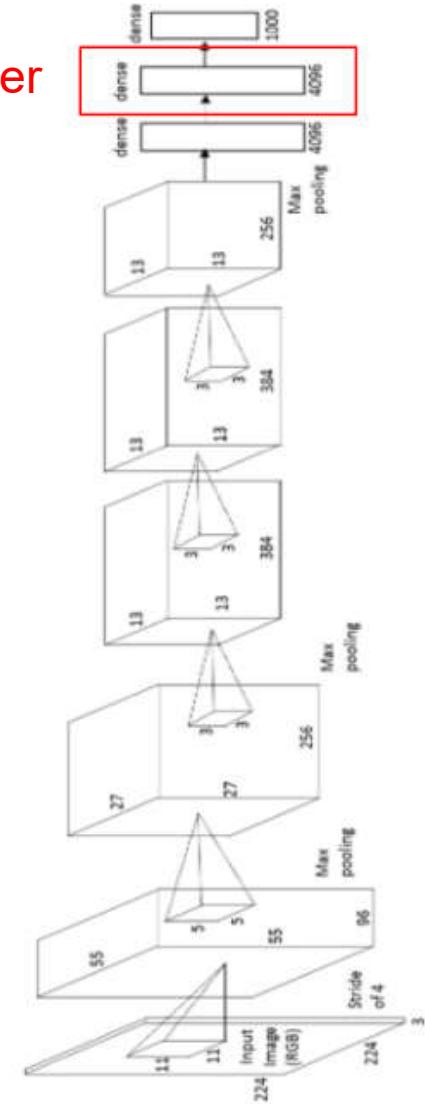
Visualizing the Representation

fc7 layer



4096-dimensional “code” for an image (layer immediately before the classifier)

Can collect the code for many images



Visualizing the Representation

t-SNE visualization

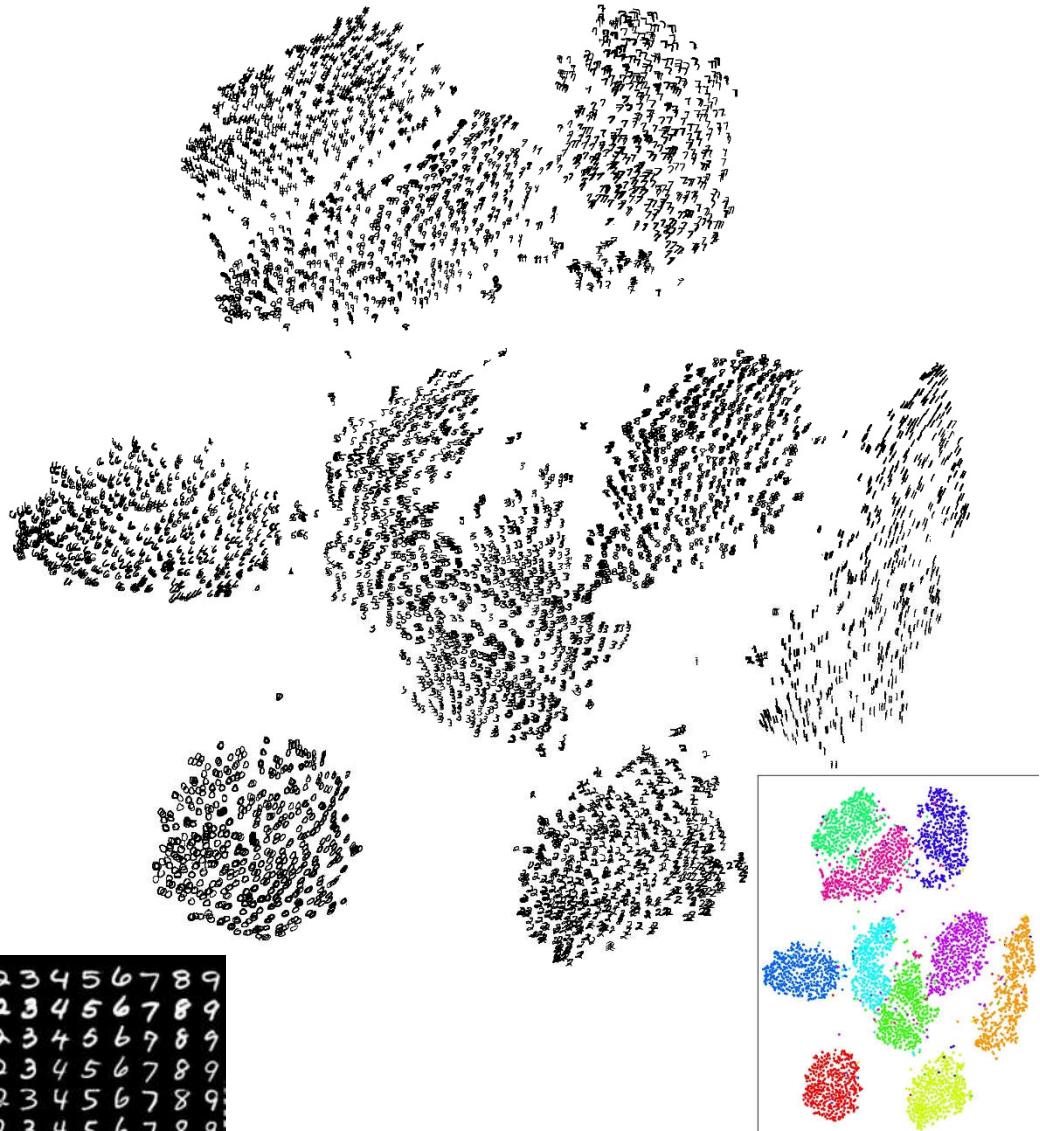
[van der Maaten & Hinton]

Embed high-dimensional points so
that locally, pairwise distances are
conserved

i.e. similar things end up in similar
places, dissimilar things end up
wherever

Right: Example embedding of
MNIST digits (0-9) in 2D

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

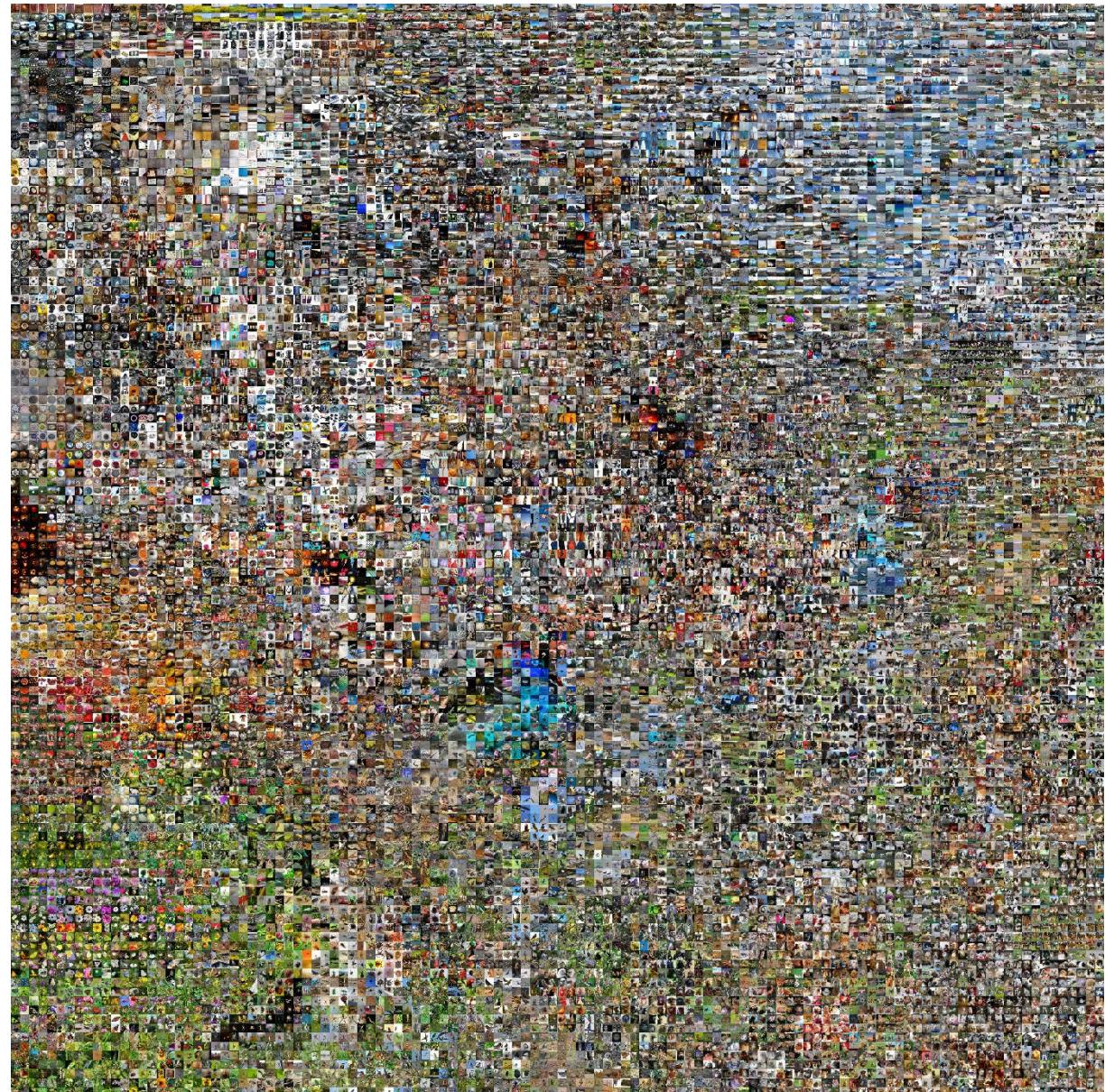


t-SNE Visualization

Two images are placed nearby if their CNN codes are close.

See more:

<http://cs.stanford.edu/people/karpathy/cnnembed/>



IV. Occlusion Experiments

Occlusion Experiments

[Zeiler & Fergus 2013]

(a) Input image



(d) Classifier, probability
of correct class

(as a function of
the position of the
square of zeros
in the original
image)

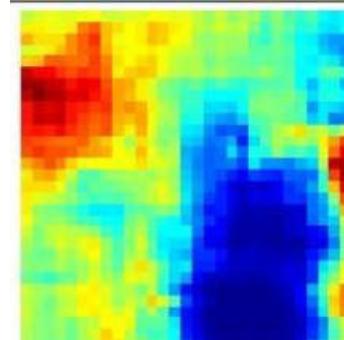
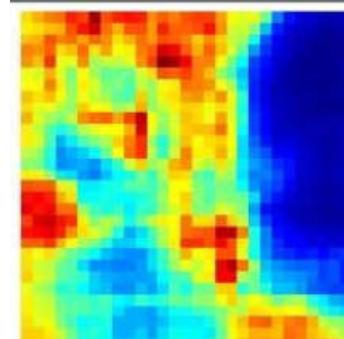
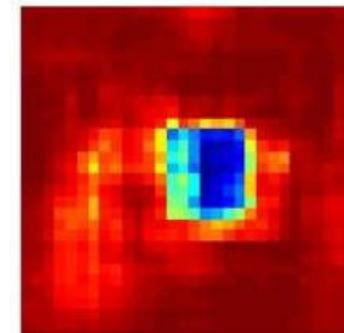
Occlusion Experiments

[Zeiler & Fergus 2013]

(a) Input image



(d) Classifier, probability of correct class

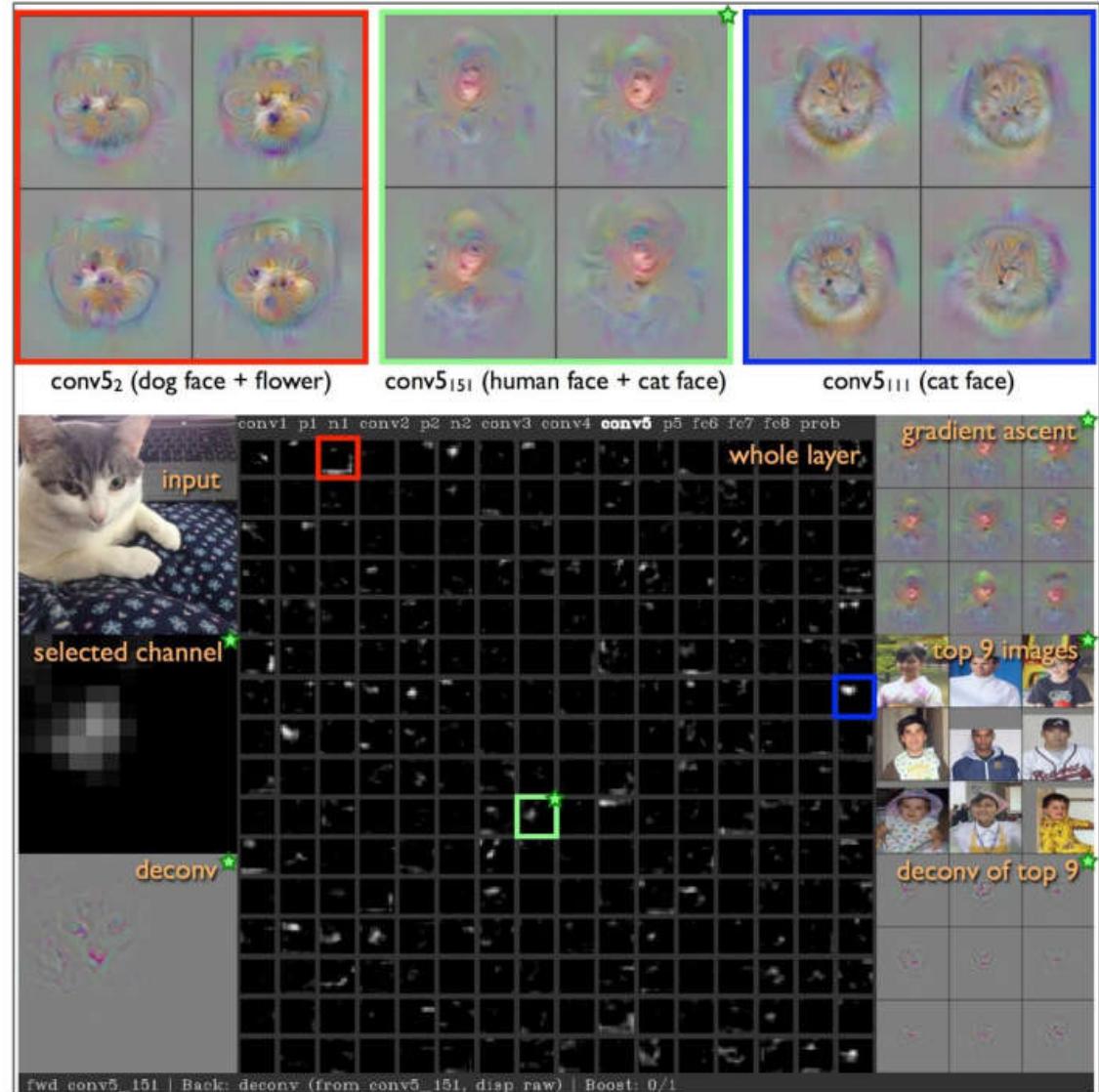


(as a function of the position of the **square of zeros** in the original image)

Visualizing Activations

<http://yosinski.com/deepvis>

- Deconv-based approach
- Optimization based approach



YouTube video: (4 mins)

<https://www.youtube.com/watch?v=AgkfIQ4IGaM>

Visualizing Activations



Deep Visualization Toolbox

yosinski.com/deepvis

#deepvis



Jason Yosinski



Jeff Clune



Anh Nguyen



Thomas Fuchs



Hod Lipson



Cornell University

UNIVERSITY
OF WYOMING

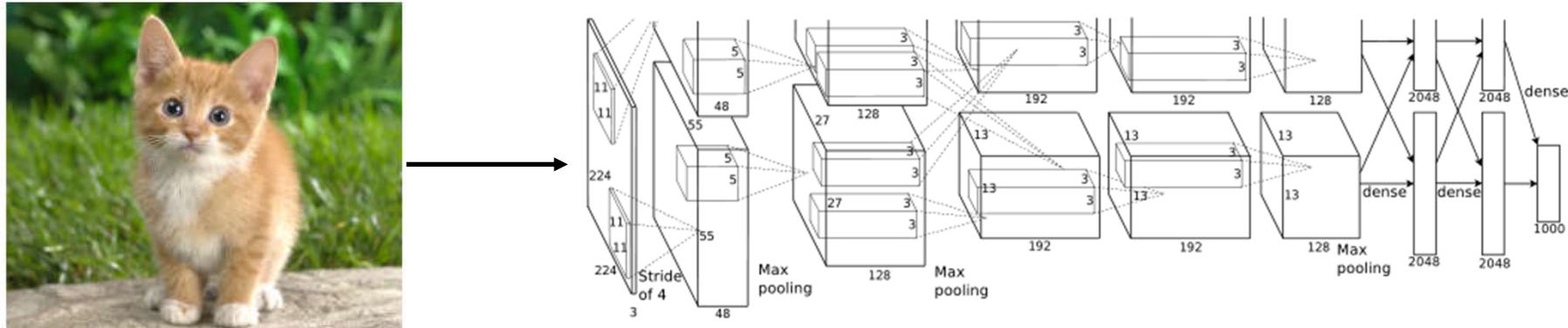


Jet Propulsion Laboratory
California Institute of Technology

V. Deconv-Based Visualization

Deconv Approach

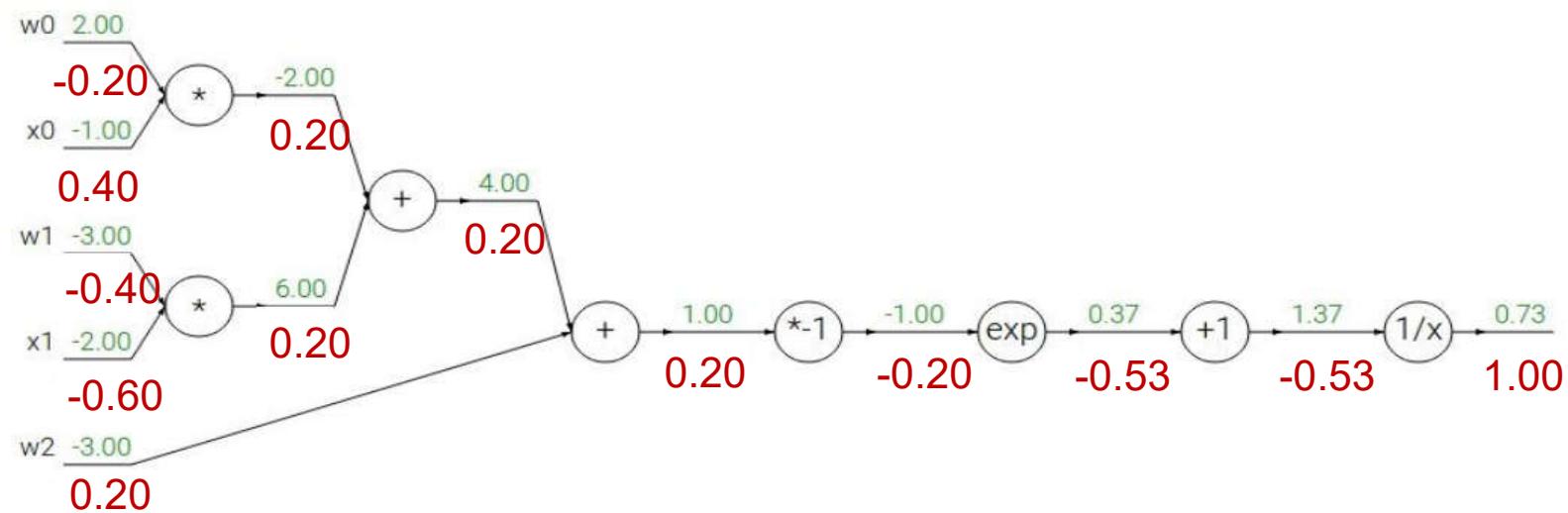
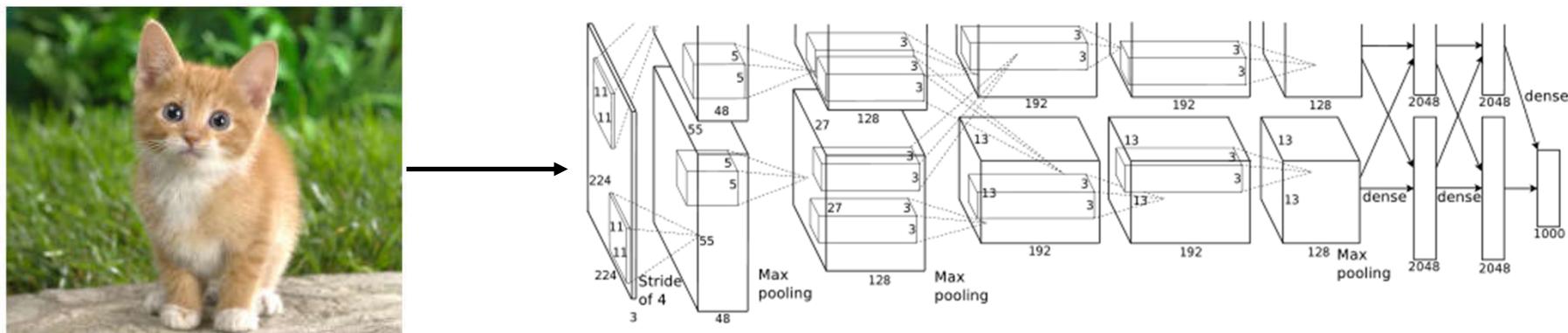
1. Feed image into net



Question: how can we compute the gradient of any arbitrary neuron in the network w.r.t the image?

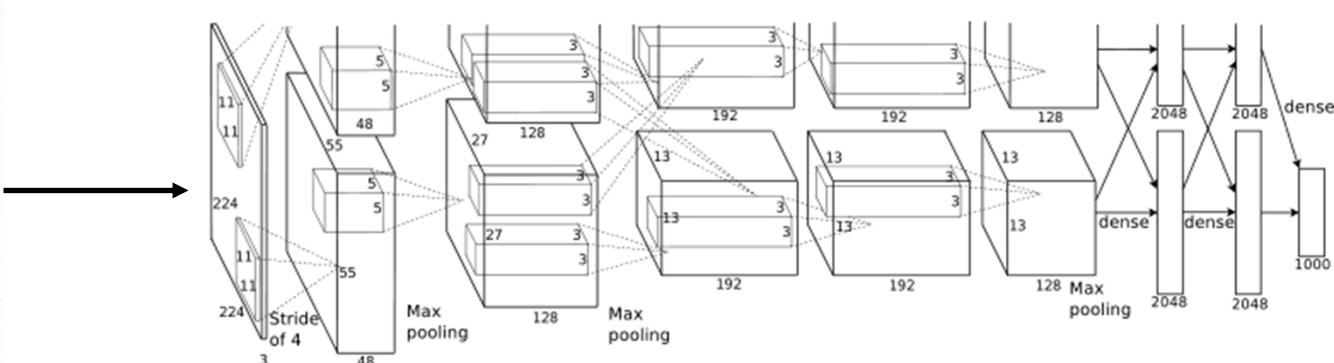
Deconv Approach

1. Feed image into net

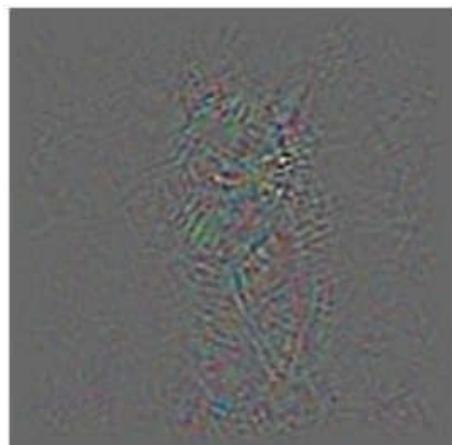


Deconv Approach

1. Feed image into net



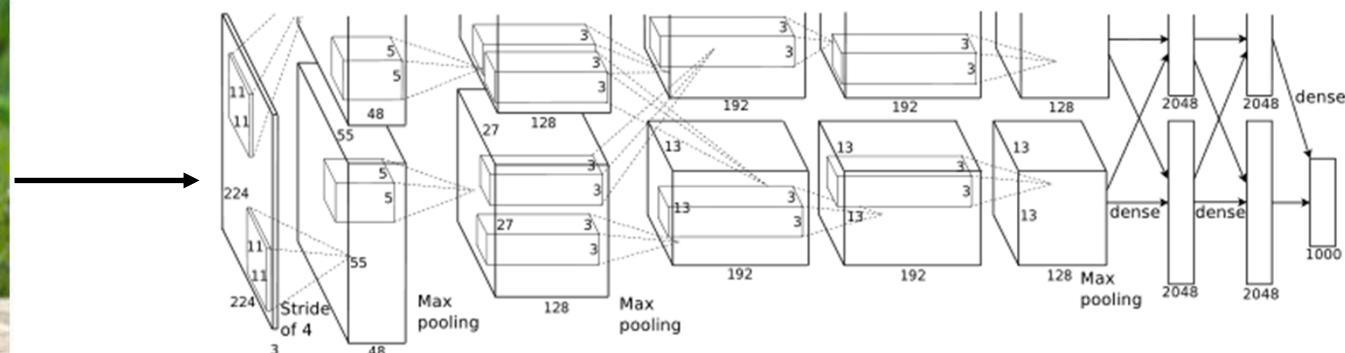
2. Pick a layer, set the gradient there to be all zeros except for one 1 for some neuron of interest



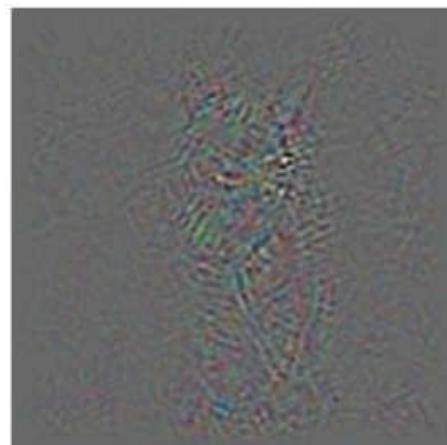
3. Backprop to image:

Deconv Approach

1. Feed image into net



2. Pick a layer, set the gradient there to be all zeros except for one 1 for some neuron of interest



**“Guided
backpropagation:”**
instead

3. Backprop to image:

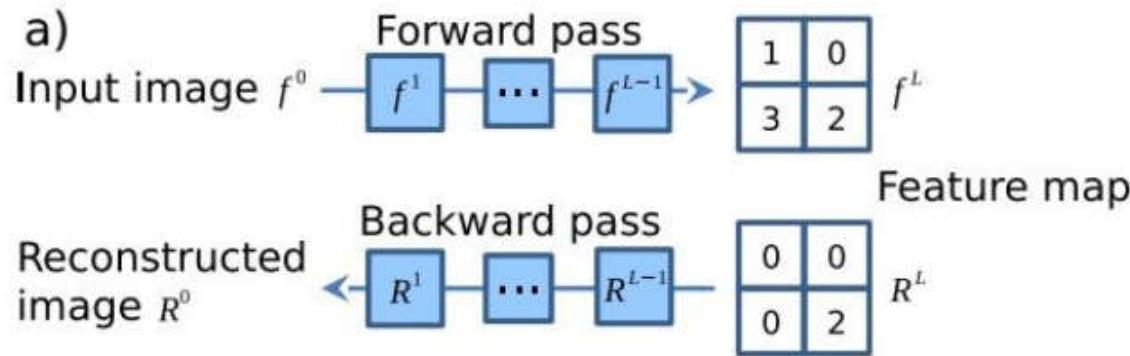


Deconv Approach

[Visualizing and Understanding Convolutional Networks, Zeiler and Fergus 2013]

[Deep Inside Convolutional Networks: Visualizing Image Classification Models and Saliency Maps, Simonyan et al., 2014]

[Striving for Simplicity: The all convolutional net, Springenberg, Dosovitskiy, et al., 2015]

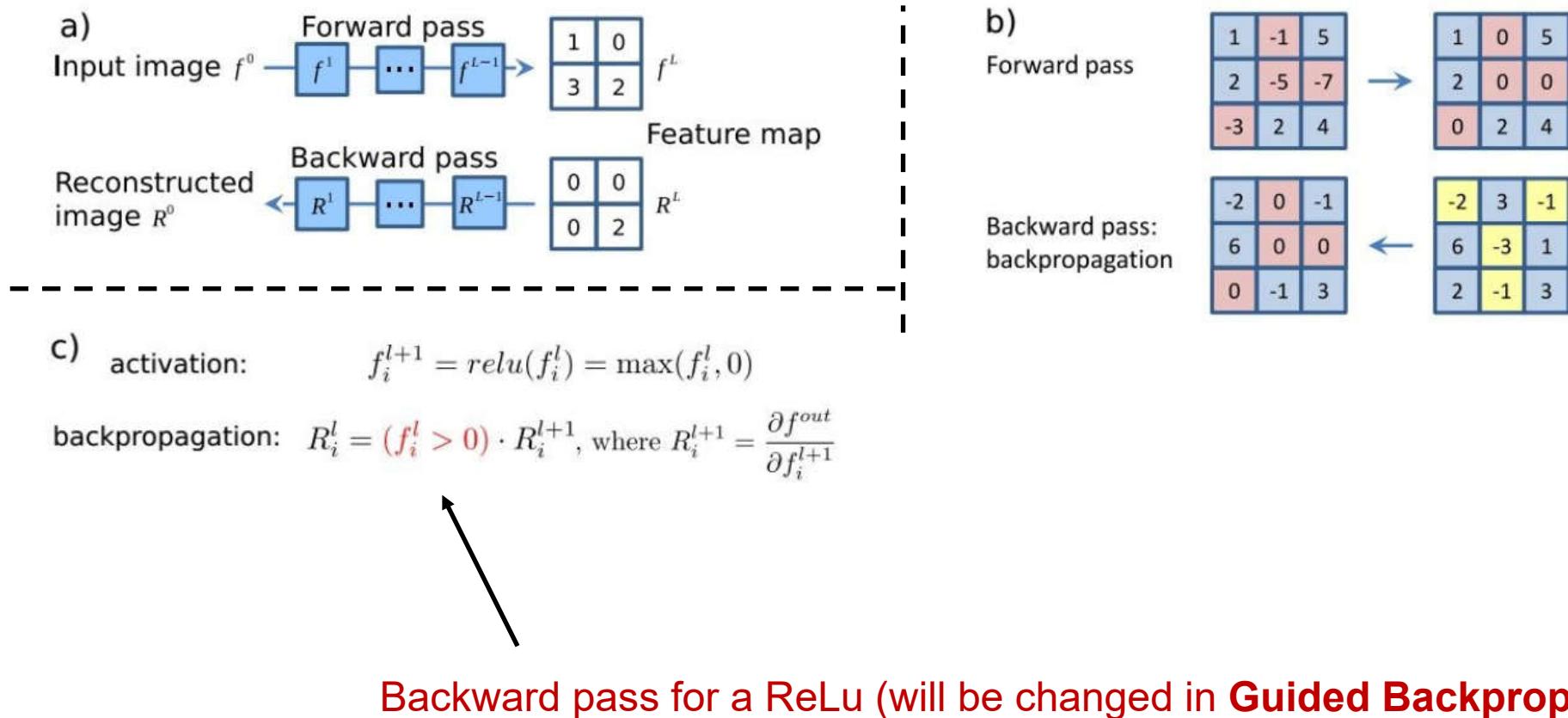


Deconv Approach

[Visualizing and Understanding Convolutional Networks, Zeiler and Fergus 2013]

[Deep Inside Convolutional Networks: Visualizing Image Classification Models and Saliency Maps, Simonyan et al., 2014]

[Striving for Simplicity: The all convolutional net, Springenberg, Dosovitskiy, et al., 2015]

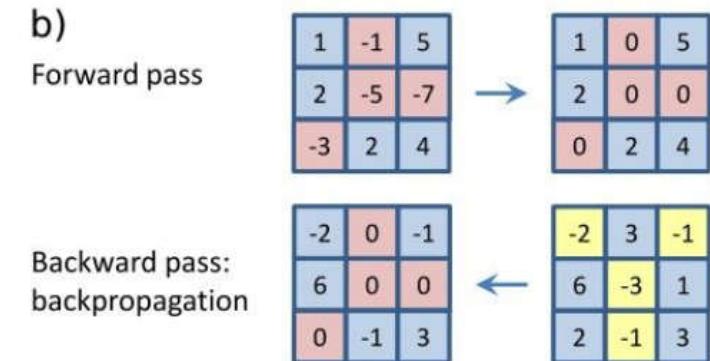
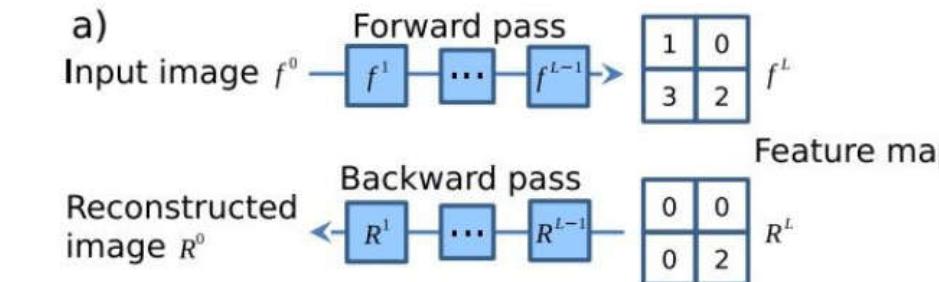


Deconv Approach

[Visualizing and Understanding Convolutional Networks, Zeiler and Fergus 2013]

[Deep Inside Convolutional Networks: Visualizing Image Classification Models and Saliency Maps, Simonyan et al., 2014]

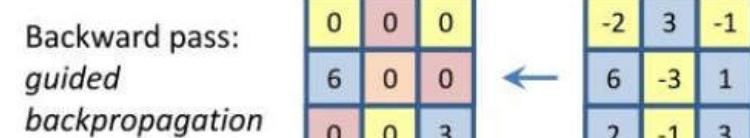
[Striving for Simplicity: The all convolutional net, Springenberg, Dosovitskiy, et al., 2015]



c) activation: $f_i^{l+1} = \text{relu}(f_i^l) = \max(f_i^l, 0)$

backpropagation: $R_i^l = (\textcolor{red}{f_i^l > 0}) \cdot R_i^{l+1}$, where $R_i^{l+1} = \frac{\partial f^{out}}{\partial f_i^{l+1}}$

guided backpropagation: $R_i^l = (\textcolor{red}{f_i^l > 0}) \cdot (\textcolor{yellow}{R_i^{l+1} > 0}) \cdot R_i^{l+1}$



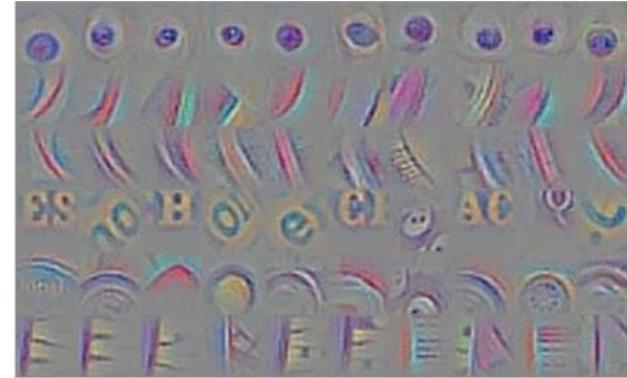
Deconv Approach

Visualization of patterns learned by the **layer conv6 (top)** and **layer conv9 (bottom)** of the network trained on ImageNet.

Each row corresponds to one filter.

The visualization using “guided backpropagation” is based on the top 10 image patches activating this filter taken from the ImageNet dataset.

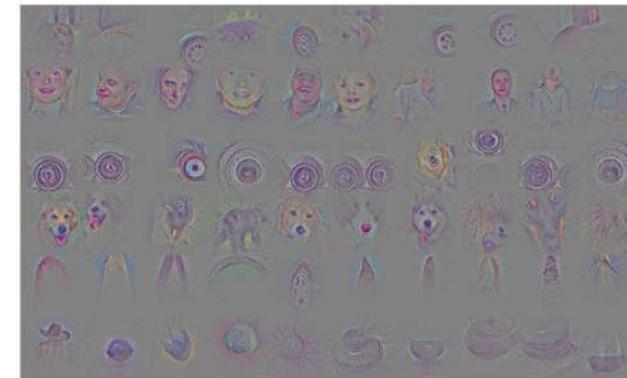
Guided backpropagation



Corresponding image crops



Guided backpropagation



Corresponding image crops

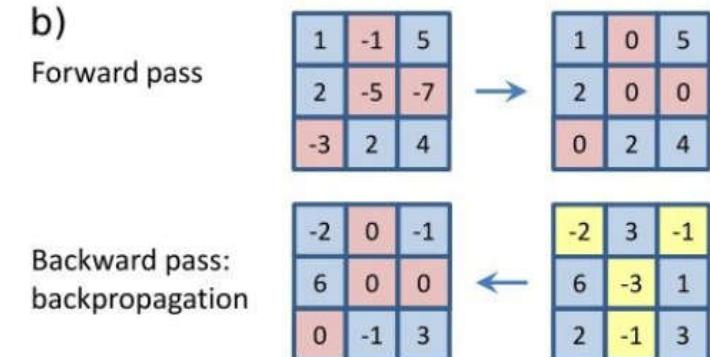
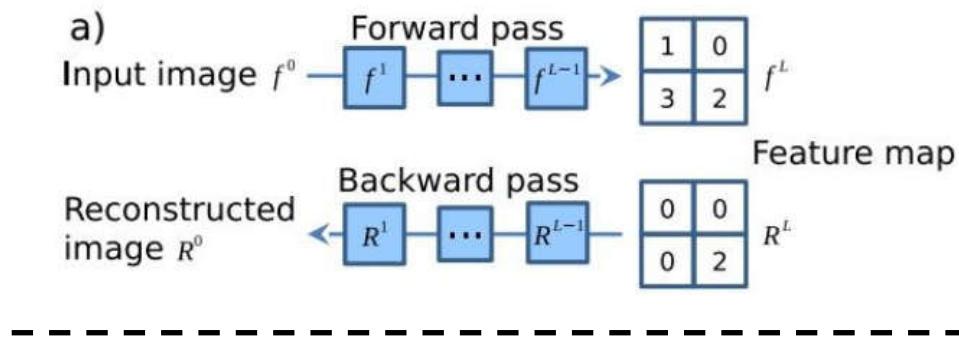


Deconv Approach

[Visualizing and Understanding Convolutional Networks, Zeiler and Fergus 2013]

[Deep Inside Convolutional Networks: Visualizing Image Classification Models and Saliency Maps, Simonyan et al., 2014]

[Striving for Simplicity: The all convolutional net, Springenberg, Dosovitskiy, et al., 2015]

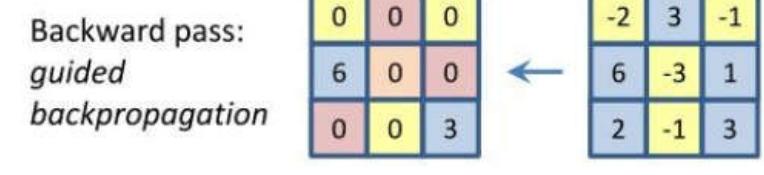
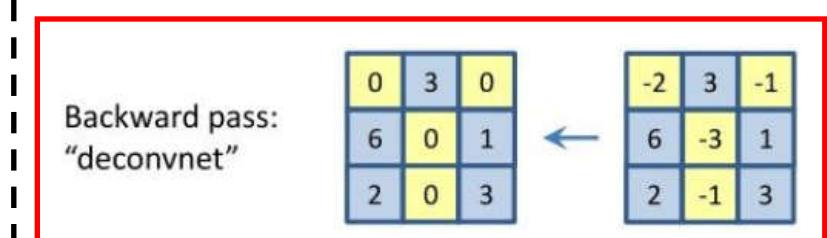


c) activation: $f_i^{l+1} = \text{relu}(f_i^l) = \max(f_i^l, 0)$

backpropagation: $R_i^l = (\textcolor{red}{f_i^l > 0}) \cdot R_i^{l+1}$, where $R_i^{l+1} = \frac{\partial f^{out}}{\partial f_i^{l+1}}$

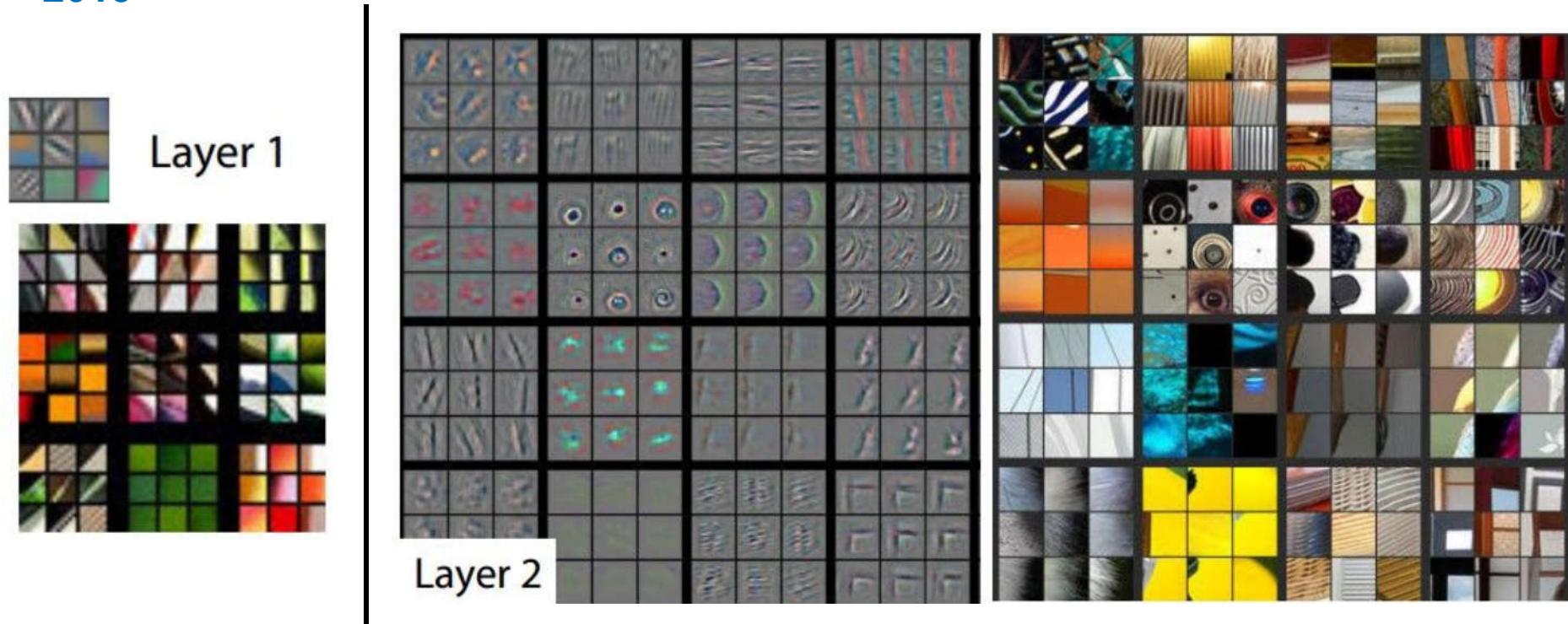
backward 'deconvnet': $\textcolor{yellow}{R_i^l = (R_i^{l+1} > 0)} \cdot R_i^{l+1}$

guided backpropagation: $R_i^l = (\textcolor{red}{f_i^l > 0}) \cdot (\textcolor{yellow}{R_i^{l+1} > 0}) \cdot R_i^{l+1}$



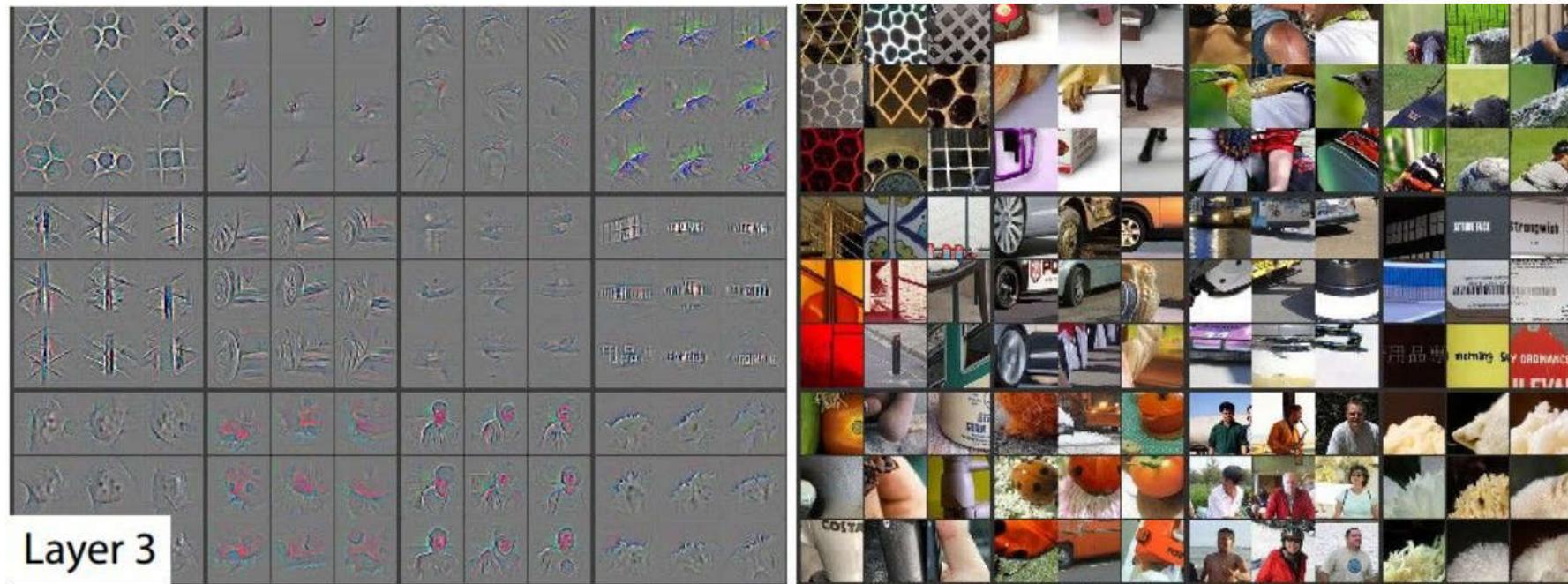
Visualizing Arbitrary Neurons Along the Way to the Top...

Visualizing and Understanding Convolutional Networks, Zeiler & Fergus, 2013

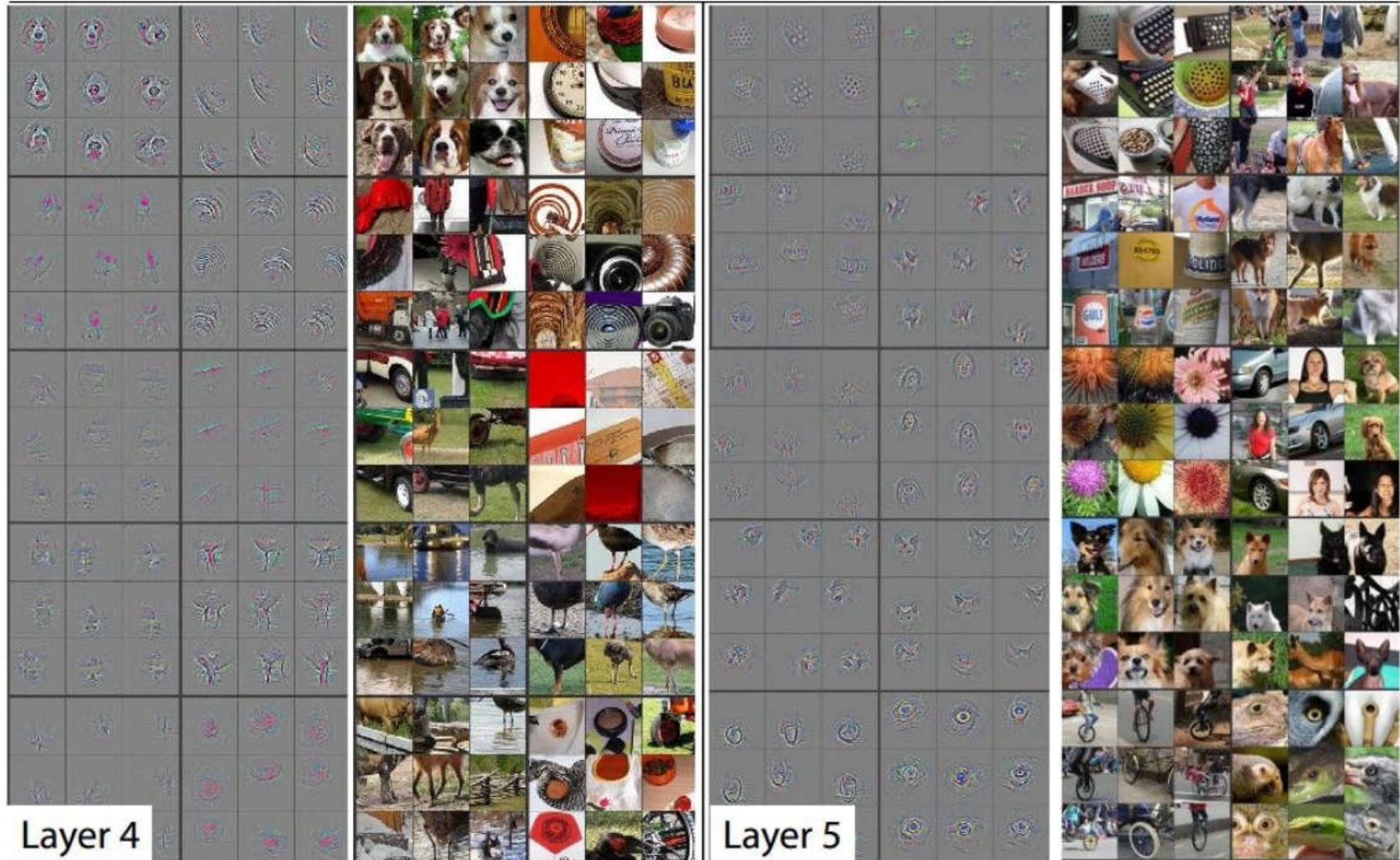


- It produces **much cleaner pictures**, because we only through the positive influences through the backward chain

Visualizing Arbitrary Neurons Along the Way to the Top...



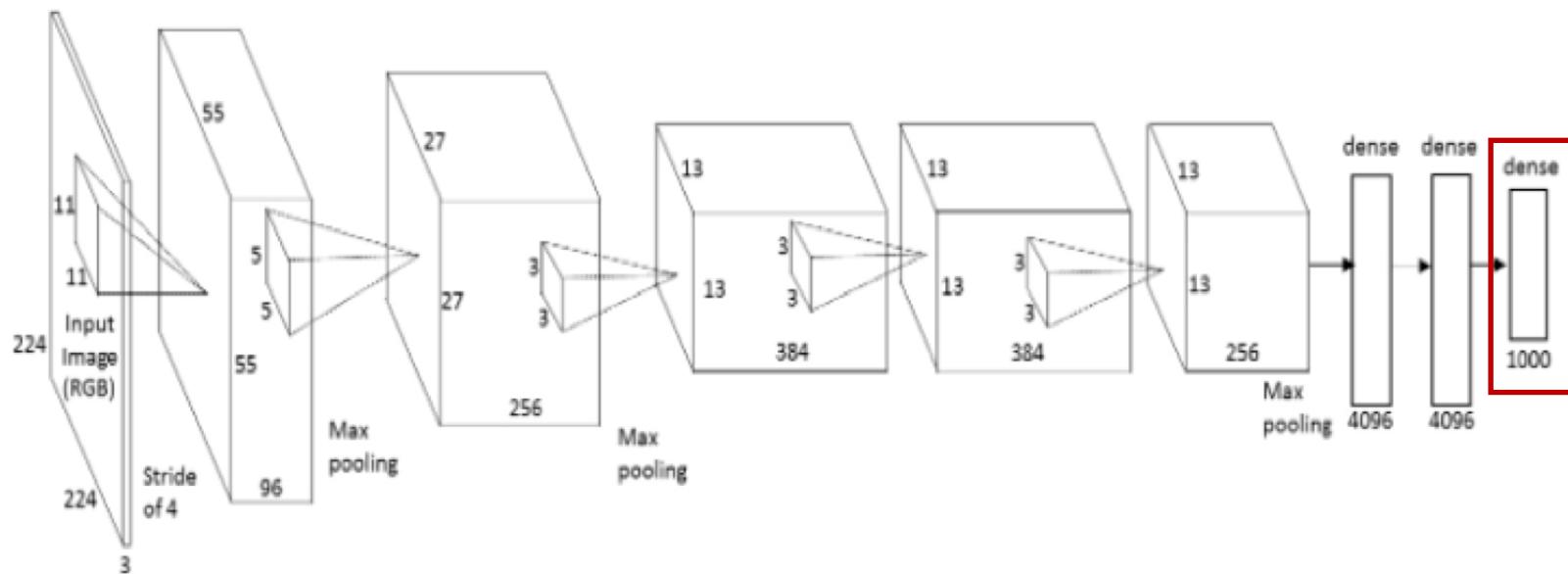
Visualizing Arbitrary Neurons Along the Way to the Top...



- Optimization based visualization which makes much more sense and are much more intuitive in terms of what's going on

VI. Optimization-Based Visualization

Optimization to Image

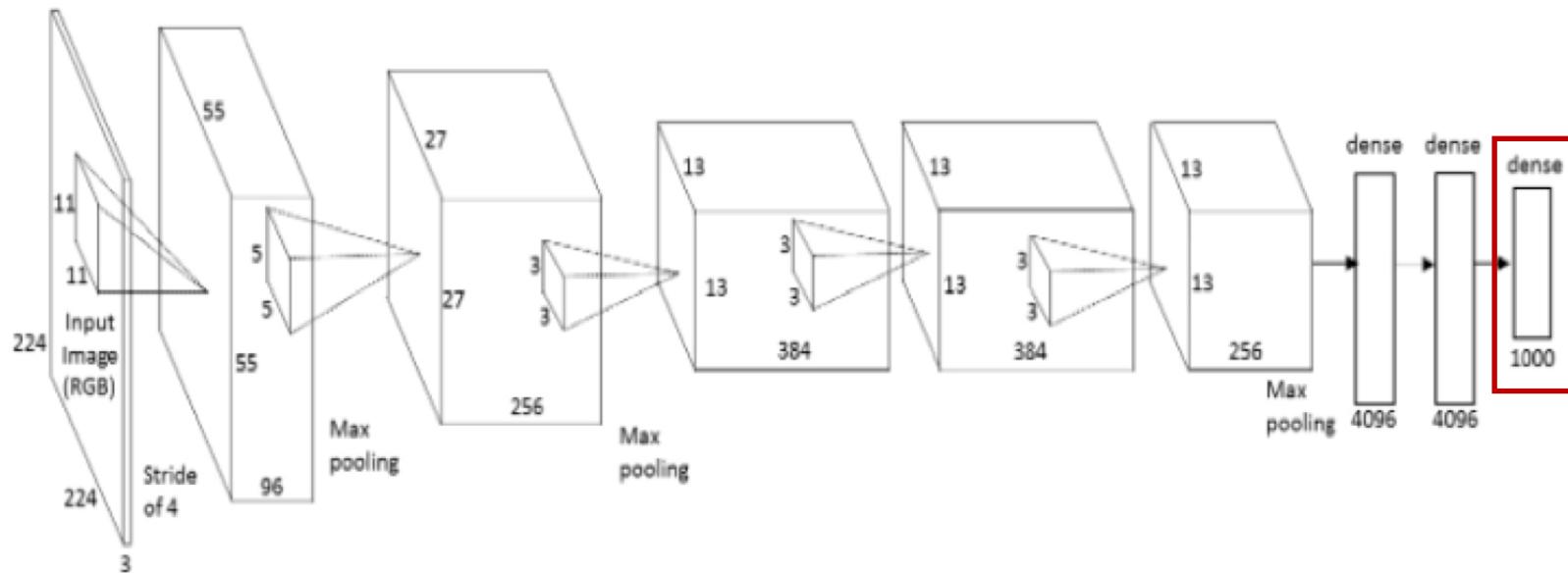


Question: can we find an image that maximizes some class score?

Optimization to Image

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

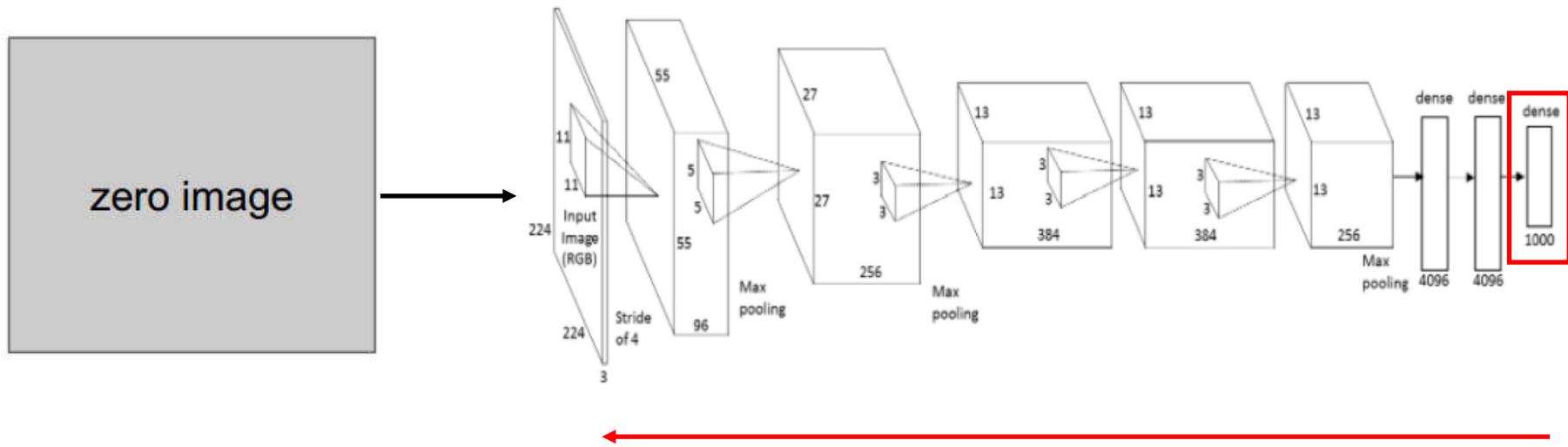
Score for class c (before Softmax)



Question: can we find an image that maximizes some class score?

Optimization to Image

1. Feed in zeros

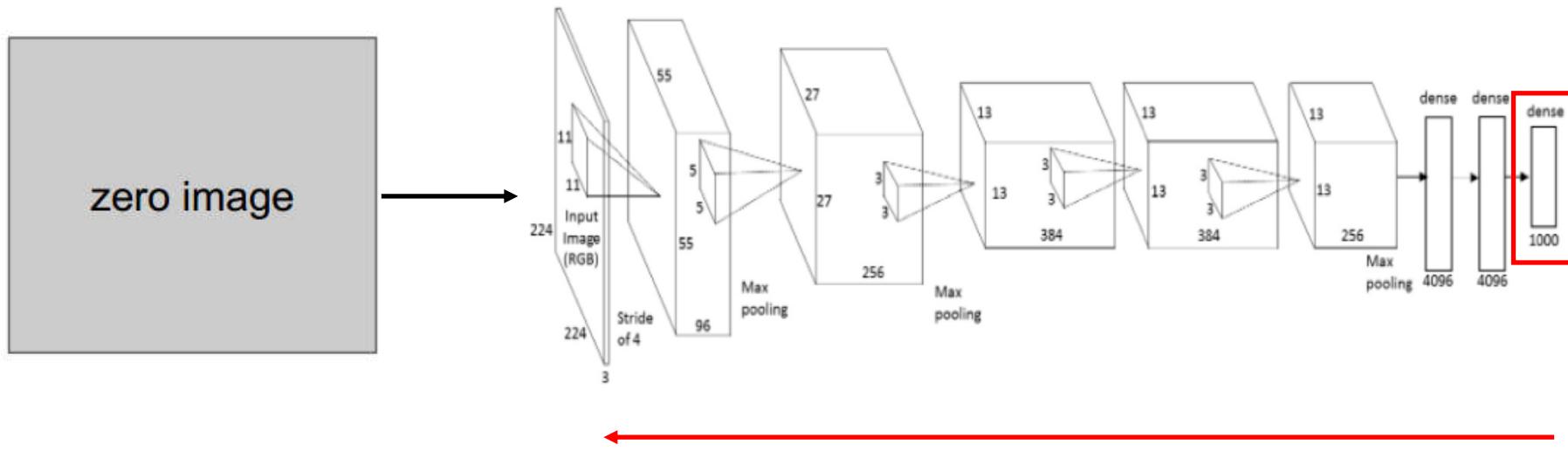


2. Set the gradient of the scores vector to be $[0, 0, \dots, 1, \dots, 0]$, then backprop to image

Question: Set zero image as initialization... is it can be a problem?

Optimization to Image

1. Feed in zeros



2. Set the gradient of the scores vector to be $[0, 0, \dots, 1, \dots, 0]$, then backprop to image
3. Do a small “image update”
4. Forward the image through the network.
5. Go back to 2.

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

Score for class c (before Softmax)

Optimization to Image

Deep Inside Convolutional networks: Visualizing Image Classification Models and Saliency Maps, Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, 2014

1. Find images that maximize some class score:



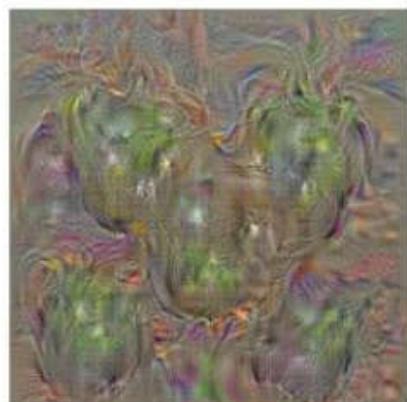
dumbbell



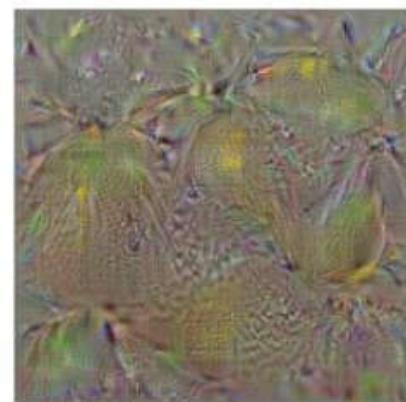
cup



dalmatian



bell pepper



lemon

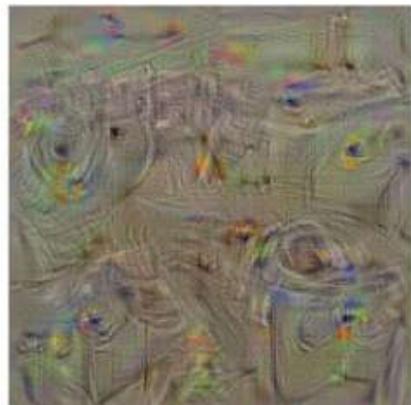


husky

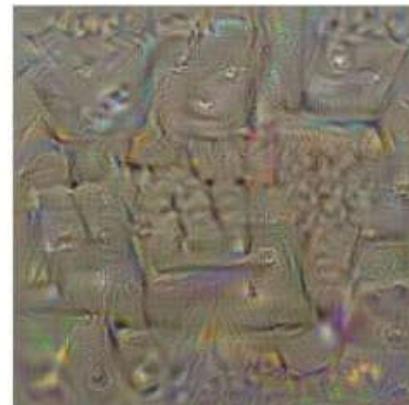
Optimization to Image

Deep Inside Convolutional networks: Visualizing Image Classification Models and Saliency Maps, Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, 2014

1. Find images that maximize some class score:



washing machine



computer keyboard



kit fox



goose



ostrich



limousine

Optimization to Image

Deep Inside Convolutional networks: Visualizing Image Classification Models and Saliency Maps, Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, 2014

2. Visualize the Data gradient:



M = ?

Note that the gradient on data has three channels. Here they visualize M, s.t.:

Hit map

$$M_{ij} = \max_c |w_h(i,j,c)|$$

(at each pixel take abs val, and max over channels)

Optimization to Image

Deep Inside Convolutional networks: Visualizing Image Classification Models and Saliency Maps, Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, 2014

2. Visualize the Data gradient:

Note that the gradient on data has three channels. Here they visualize M, s.t.:



Hit map

$$M_{ij} = \max_c |w_h(i,j,c)|$$

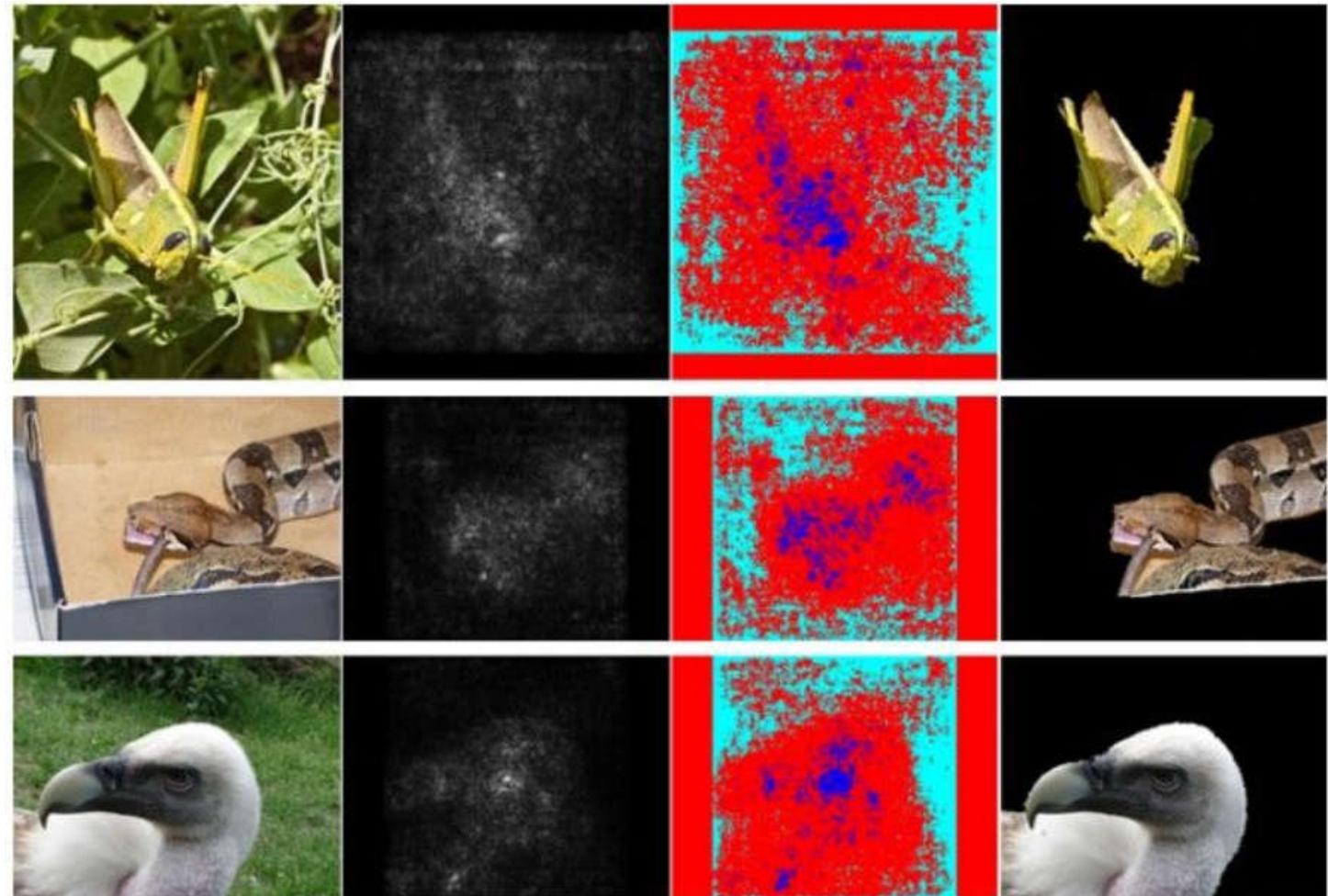
(at each pixel take abs val, and max over channels)



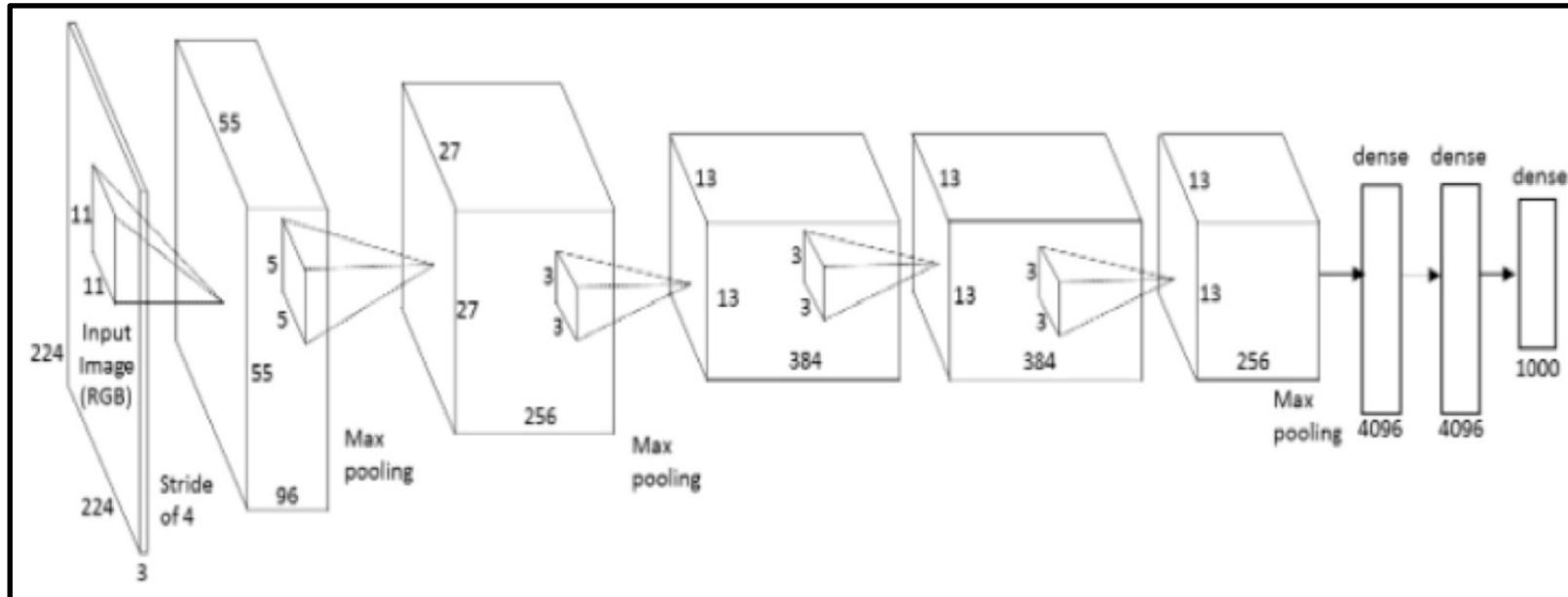
Optimization to Image

Deep Inside Convolutional networks: Visualizing Image Classification Models and Saliency Maps, Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, 2014

- Use **grabcut** for segmentation



We Can In Fact Do This for Arbitrary Neurons Along the ConvNet



Repeat:

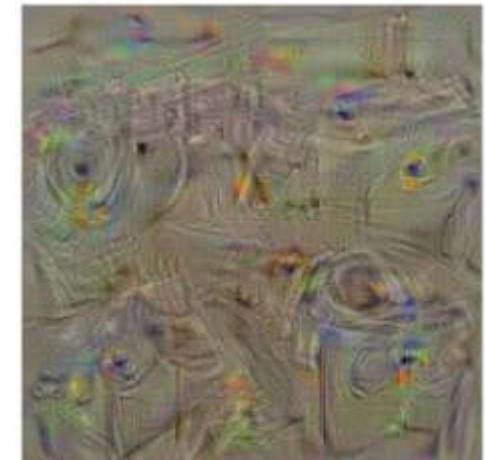
1. Forward an image
2. Set activations in layer of interest to all zero, except for a 1.0 for a neuron of interest
3. Backprop to image
4. Do an “image update”

We Can In Fact Do This for Arbitrary Neurons Along the ConvNet

[Understanding Neural Networks Through Deep Visualization, Yosinski et al., 2015]

Proposed a different form of regularizing the image

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$



More explicit scheme:

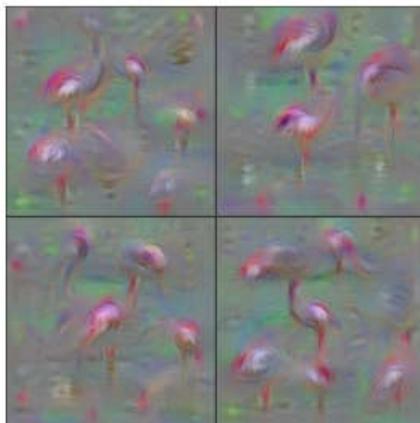
Repeat:

- Update the image \mathbf{x} with gradient from some unit of interest
- Blur \mathbf{x} a bit
- Take any pixel with small norm to zero (to encourage sparsity)

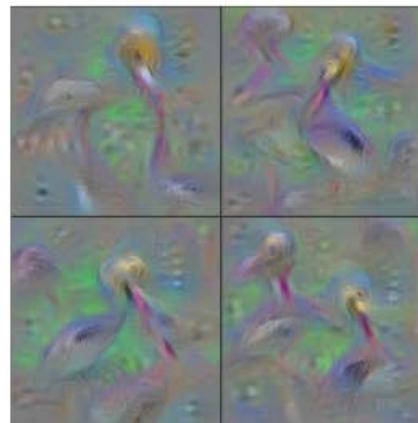
We Can In Fact Do This for Arbitrary Neurons Along the ConvNet

[Understanding Neural Networks Through Deep Visualization, Yosinski et al., 2015]

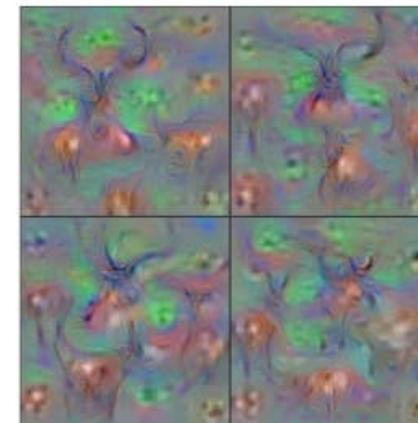
<http://yosinski.com/deepvis>



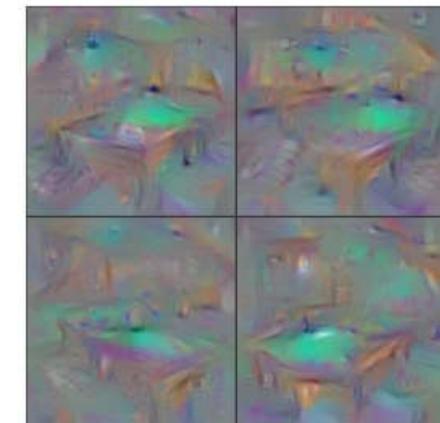
Flamingo



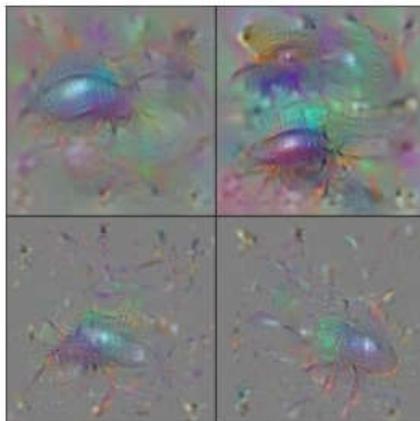
Pelican



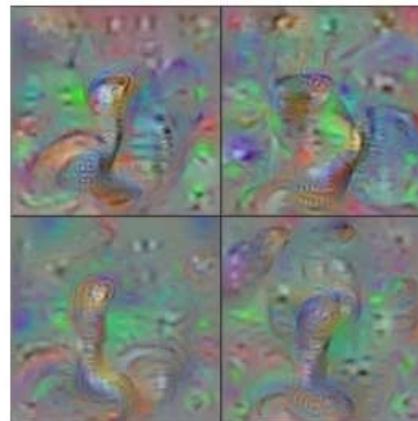
Hartebeest



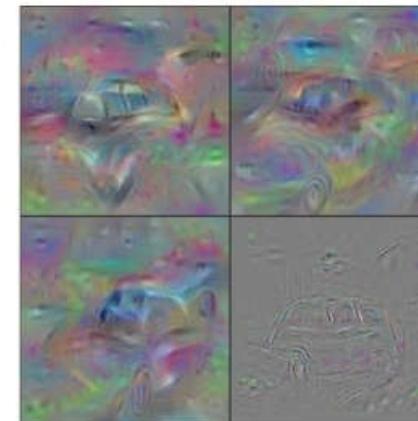
Billiard Table



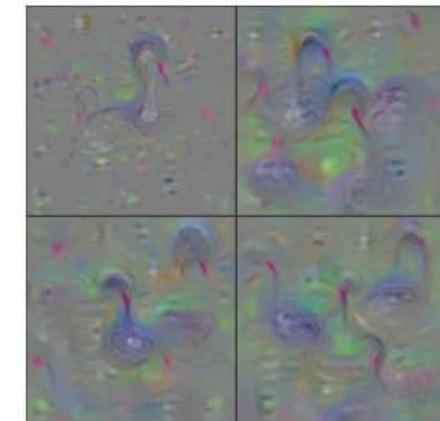
Ground Beetle



Indian Cobra



Station Wagon

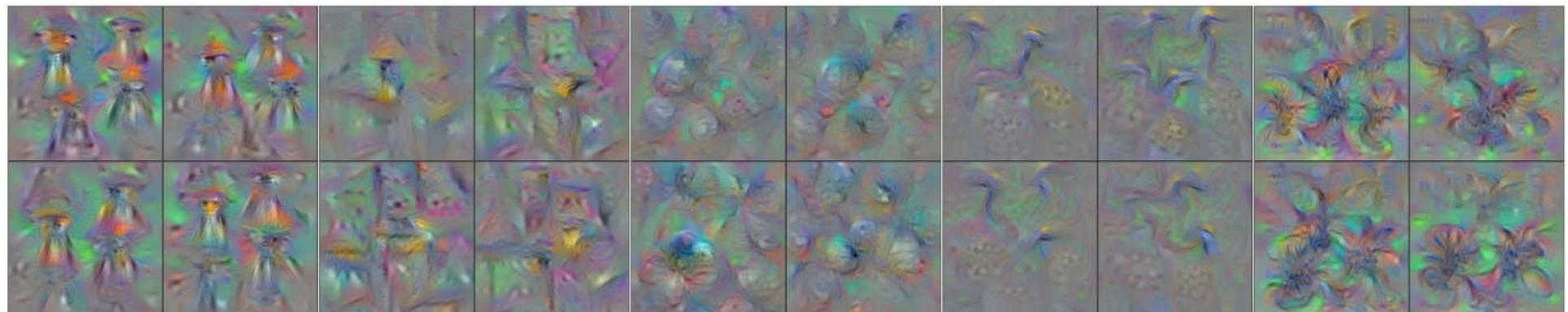


Black Swan

We Can In Fact Do This for Arbitrary Neurons Along the ConvNet

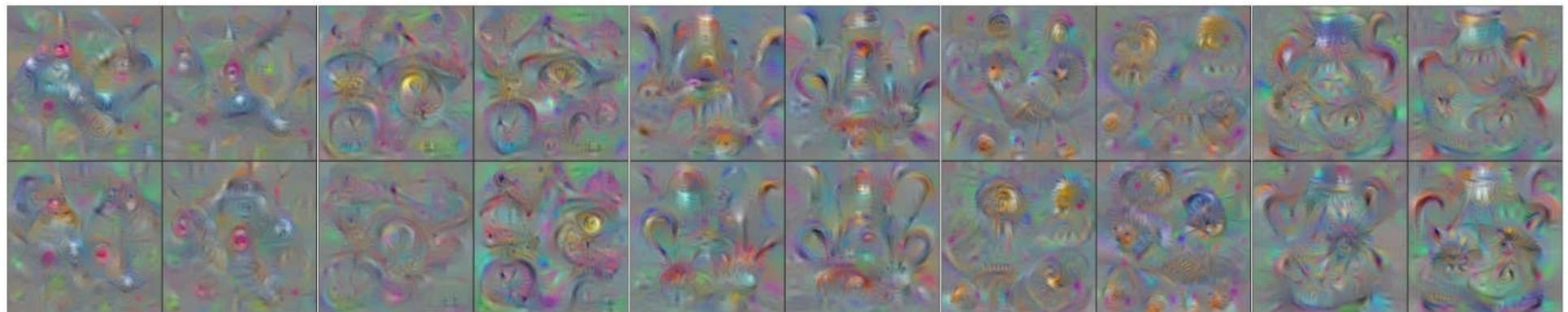


Layer 7

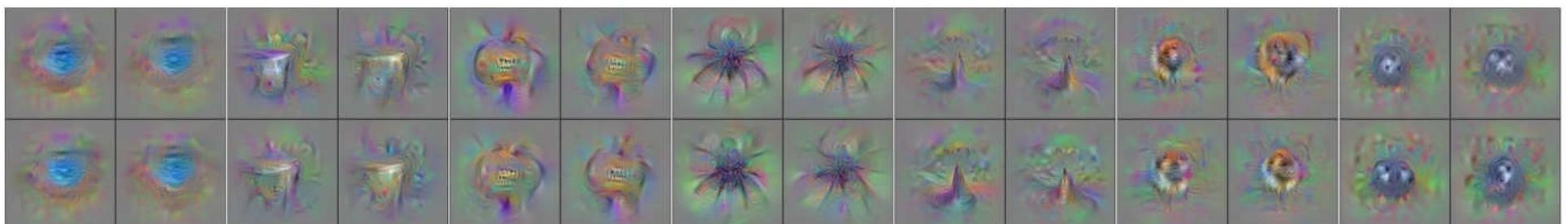


We Can In Fact Do This for Arbitrary Neurons Along the ConvNet

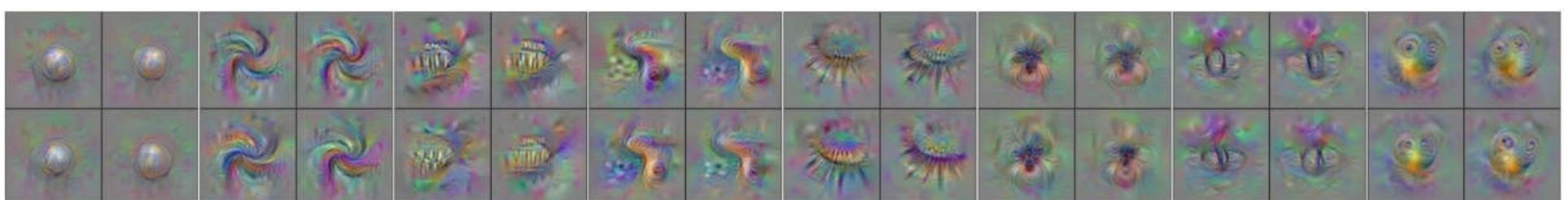
Layer 6



Layer 5

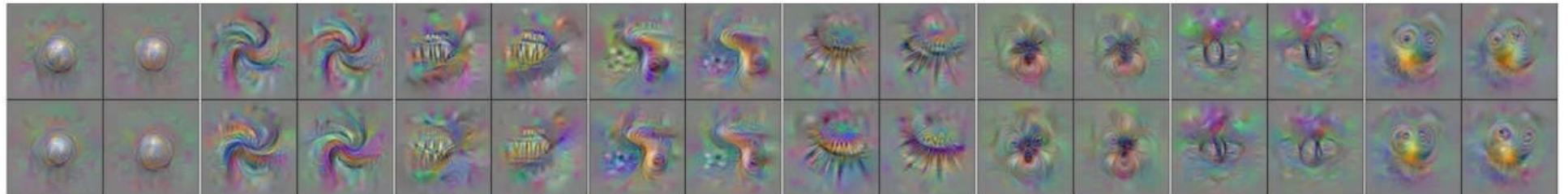


Layer 4

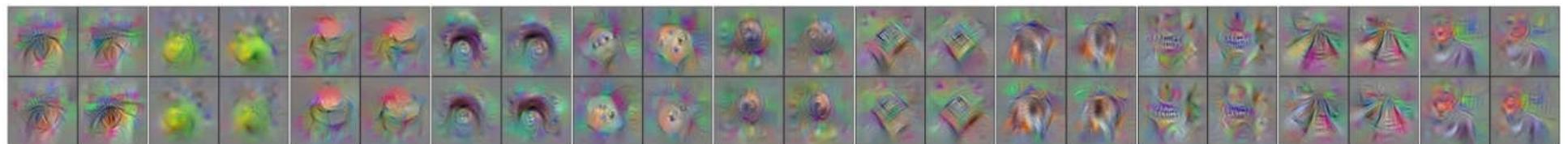


We Can In Fact Do This for Arbitrary Neurons Along the ConvNet

Layer 4



Layer 3



Layer 2



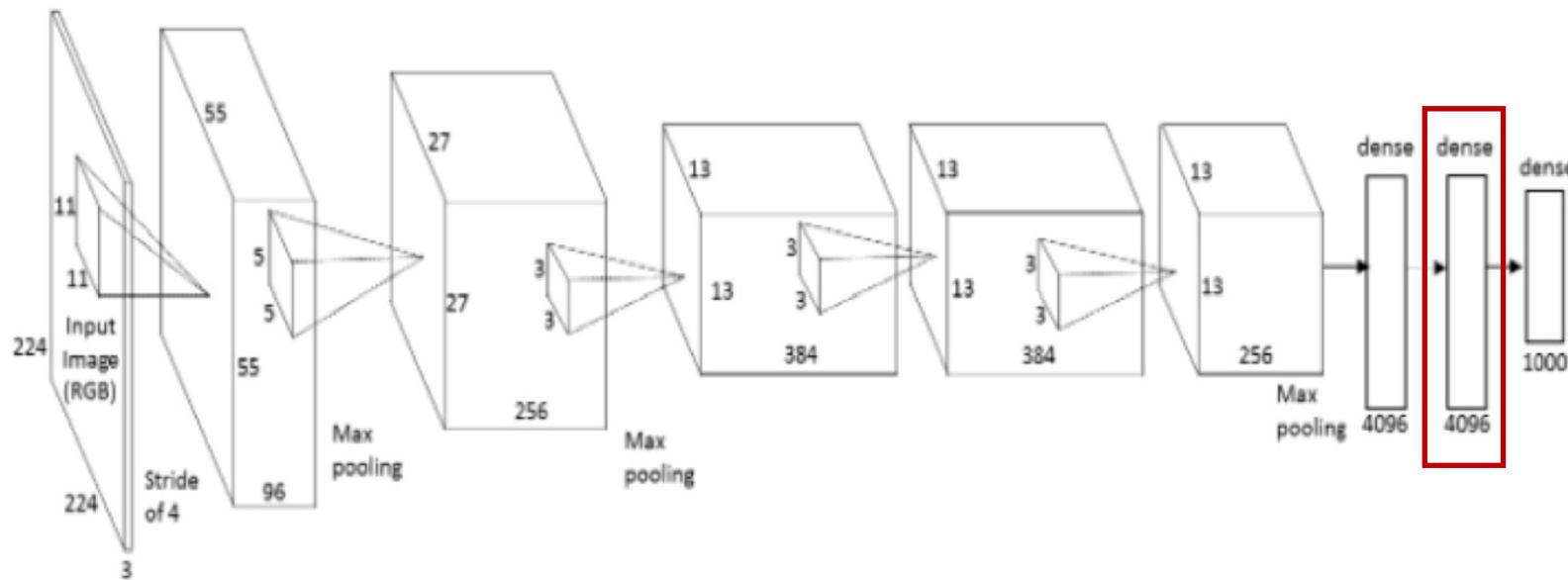
Layer 1



Reconstruction from CNN Code

Reconstruction from CNN Code

Question: Given a CNN code, is it possible to reconstruct the original image?



- Some companies might try to only store the features not the raw images

Reconstruction from CNN Code

Find an image such that:

- Its code is similar to a given code
- It “looks natural” (image prior regularization)

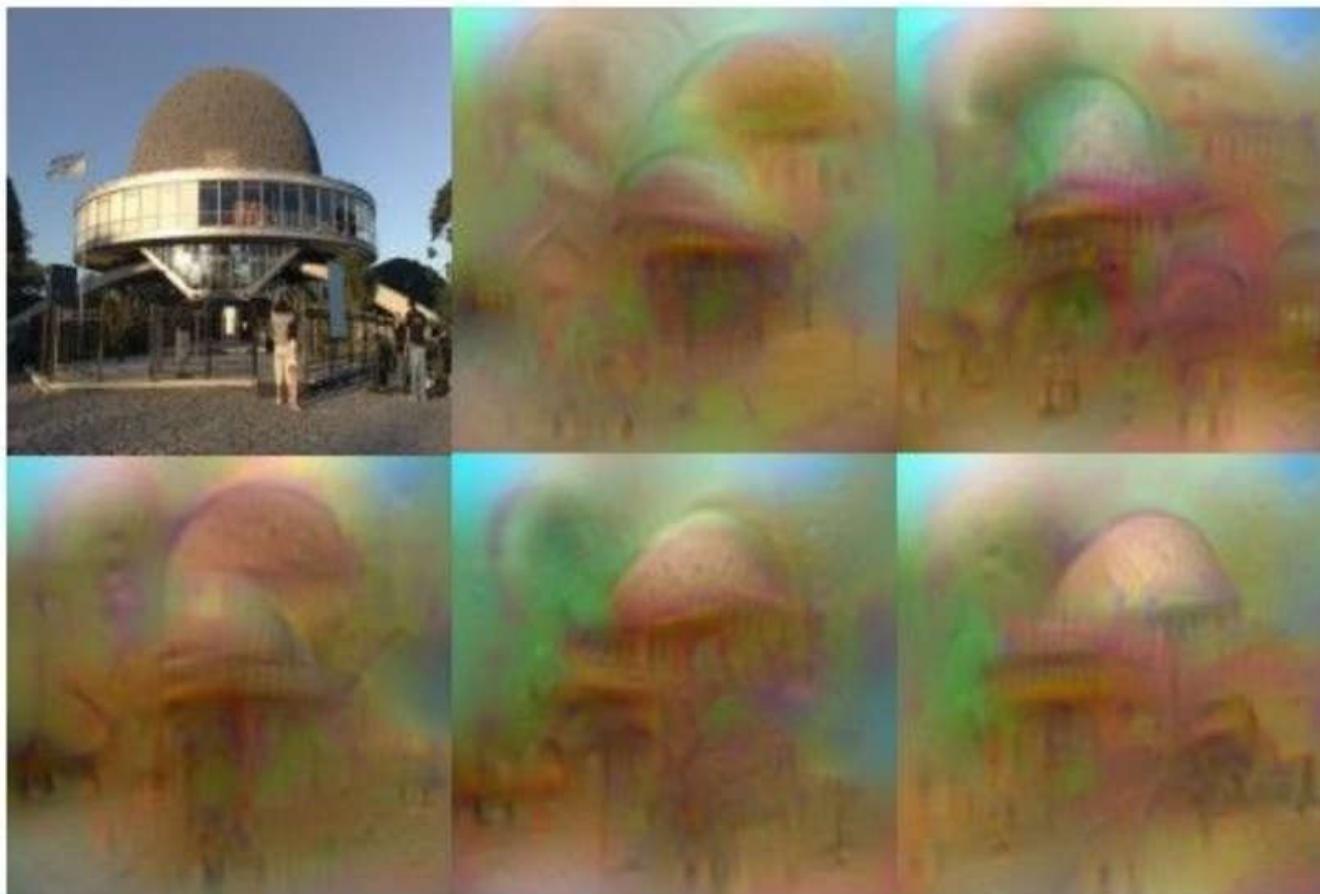
$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbf{R}^{H \times W \times C}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda R(\mathbf{x})$$

$$\ell(\Phi(\mathbf{x}), \Phi_0) = \|\Phi(\mathbf{x}) - \Phi_0\|^2$$

Reconstruction from CNN Code

[Understanding Deep Image Representations by Inverting Them, Mahendran and Vedaldi, 2014]

Original image



Reconstructions
from the 1000 log
probabilities for
ImageNet
(ILSVRC classes)

Reconstruction from CNN Code

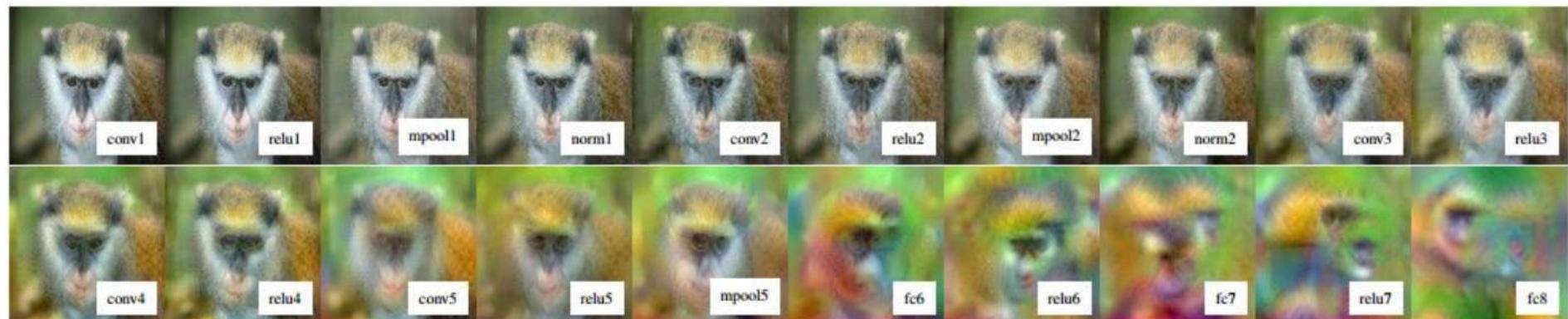
Reconstructions from the representation after last pooling layer
(immediately before the first Fully Connected layer)



Reconstruction from CNN Code



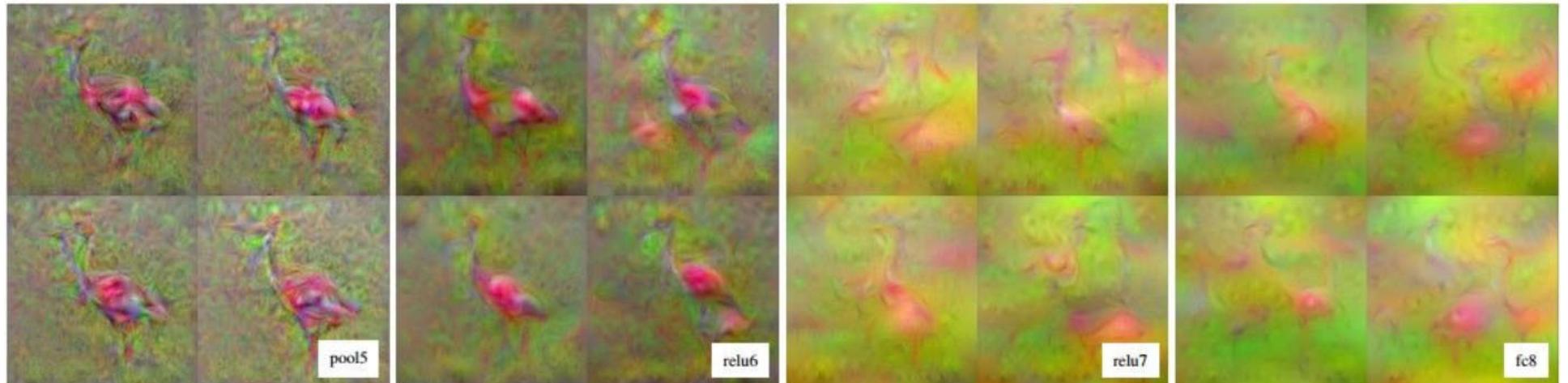
Reconstructions from intermediate layers



Reconstruction from CNN Code

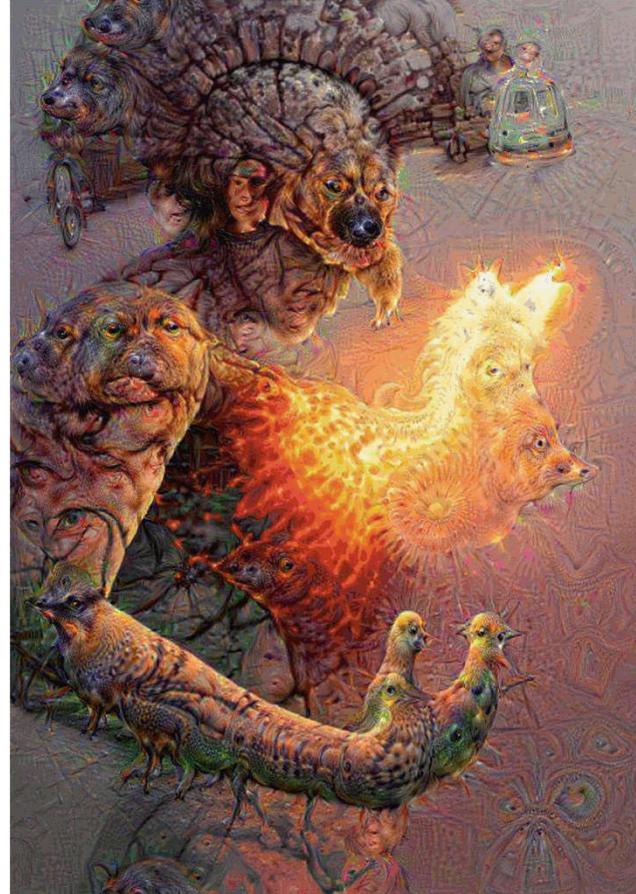
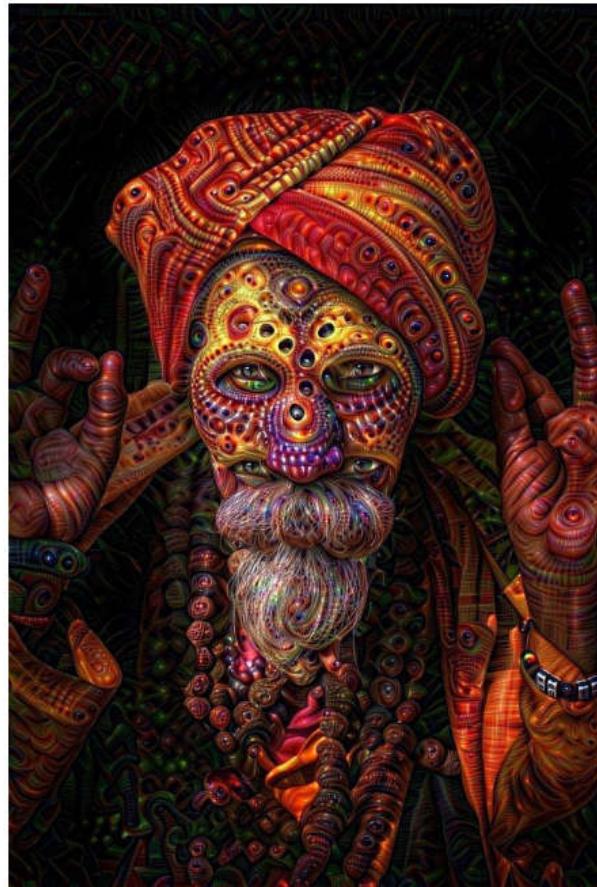


Multiple reconstructions. Images in quadrants all “look” the same to the CNN (same code)



DeepDream

DeepDream



DeepDream <https://github.com/google/deepdream>

Slide from CS231n 58

DeepDream

```

def objective_L2(dst):
    dst.diff[:] = dst.data  DeepDream: set  $\delta x = x$  :)

def make_step(net, step_size=1.5, end='inception_4c/output',
             jitter=32, clip=True, objective=objective_L2):
    '''Basic gradient ascent step.'''

    src = net.blobs['data'] # input image is stored in Net's 'data' blob
    dst = net.blobs[end]

    ox, oy = np.random.randint(-jitter, jitter+1, 2)
    src.data[0] = np.roll(np.roll(src.data[0], ox, -1), oy, -2) # apply jitter shift

    net.forward(end=end)
    objective(dst) # specify the optimization objective
    net.backward(start=end)
    g = src.diff[0]
    # apply normalized ascent step to the input image
    src.data[:] += step_size/np.abs(g).mean() * g

    src.data[0] = np.roll(np.roll(src.data[0], -ox, -1), -oy, -2) # unshift image

if clip:
    bias = net.transformer.mean['data']
    src.data[:] = np.clip(src.data, -bias, 255-bias)

```

Jitter regularizer

“image update”

DeepDream



Inception_4c/output

DeepDream modifies the image in a way that “boosts” all activations, at any layer

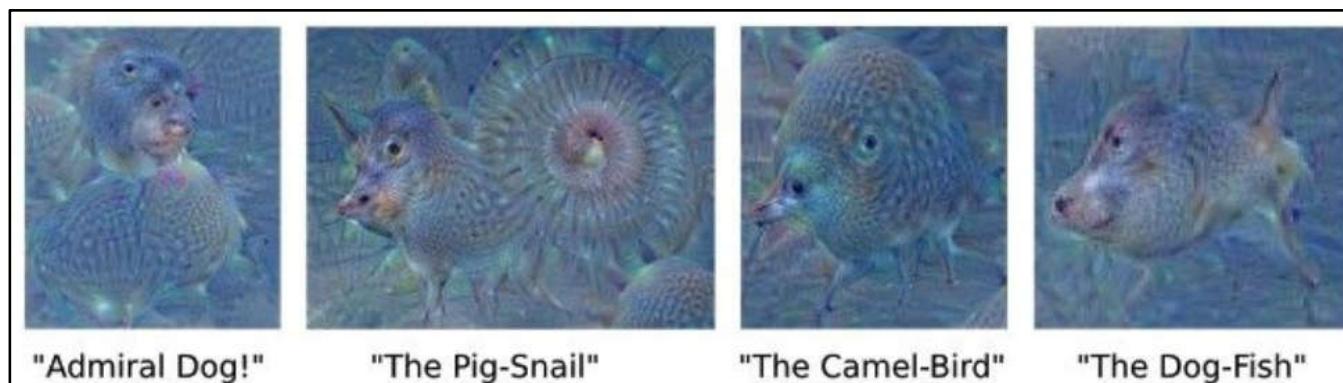
This creates a feedback loop: e.g. any slightly detected dog face will be made more and more dog like over time

DeepDream

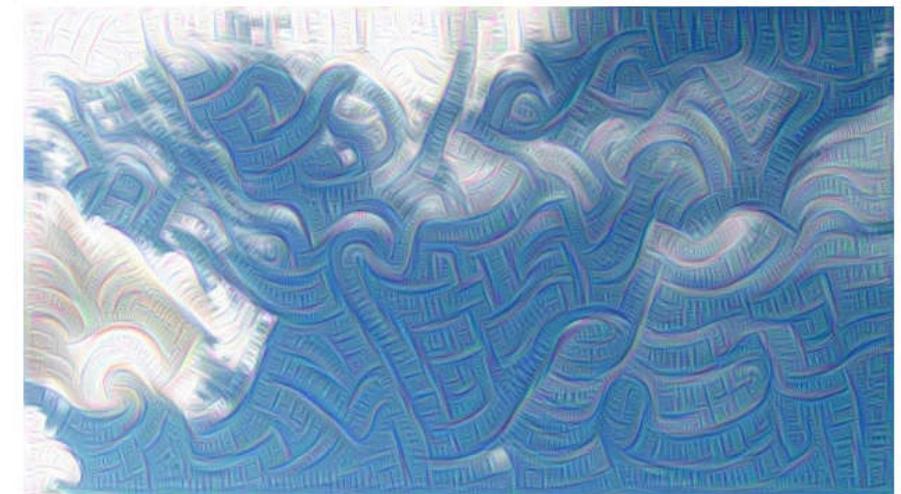
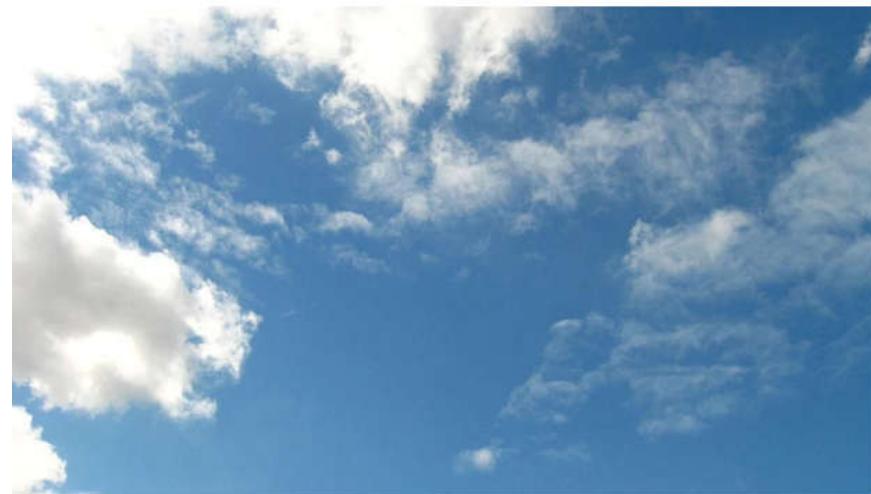


Inception_4c/output

DeepDream modifies the image in a way that “boosts” all activations, at any layer



DeepDream



Inception_3b/5x5_output

DeepDream modifies the image in a way that “boosts” all activations, at any layer

Bonus Videos

Deep Dream Grocery Trip: <https://www.youtube.com/watch?v=DgPaCWJL7XI>



Bonus Videos

Deep Dreaming Fear & Loathing in Las Vegas: the Great San Francisco Acid Wave:

<https://www.youtube.com/watch?v=oyxSerkkP4o>



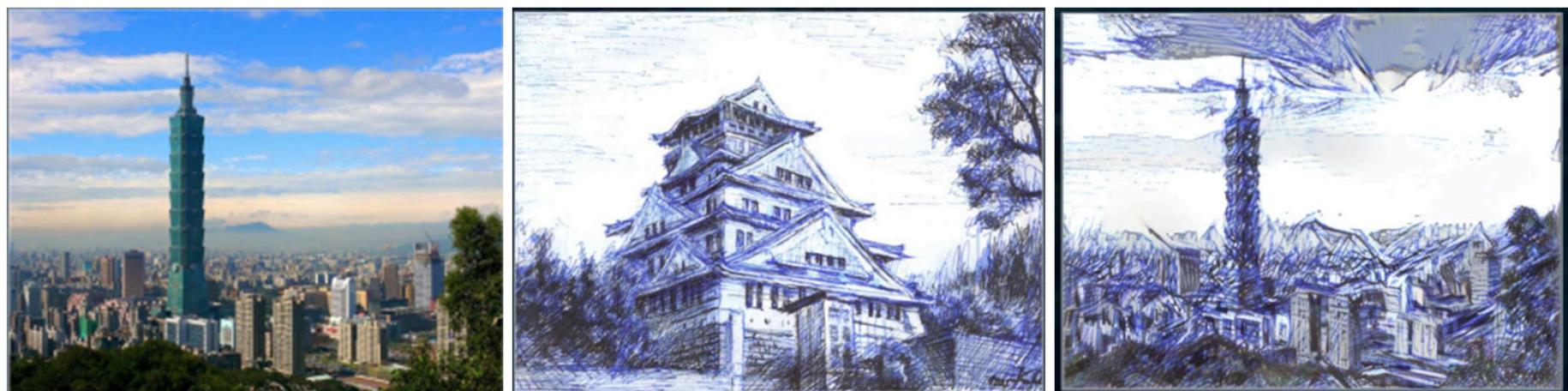
Neural Artistic Style

Neural Artistic Style

[L. A. Gatys, et al., Image Style Transfer
Using Convolutional Neural Networks,
CVPR2016]

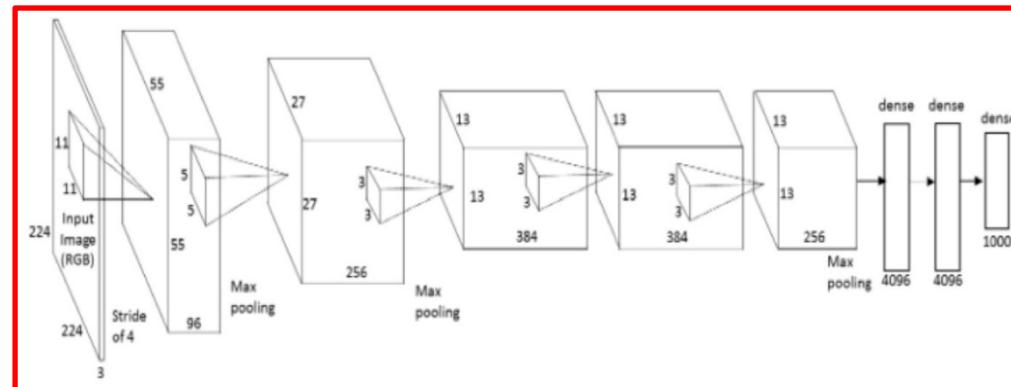


Neural Artistic Style



Neural Artistic Style

Step 1: Extract **content targets** (ConvNet activations of all layers for the given content image)



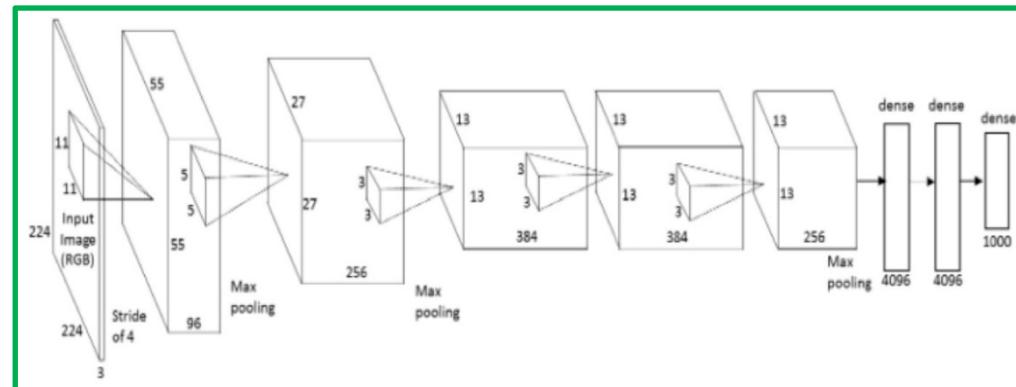
Content activations

e.g.

At **CONV5_1** layer we would have a **[14x14x512]** array of target activations

Neural Artistic Style

Step 2: Extract **style targets** (Gram matrices of ConvNet activations of all layers for the given style image)



Style gram matrices

e.g.

$$G = V^T V$$

At CONV1 layer (with **[224x224x64]** activations) would give a **[64x64]** Gram matrix of all pairwise activation covariances (summed across spatial locations)

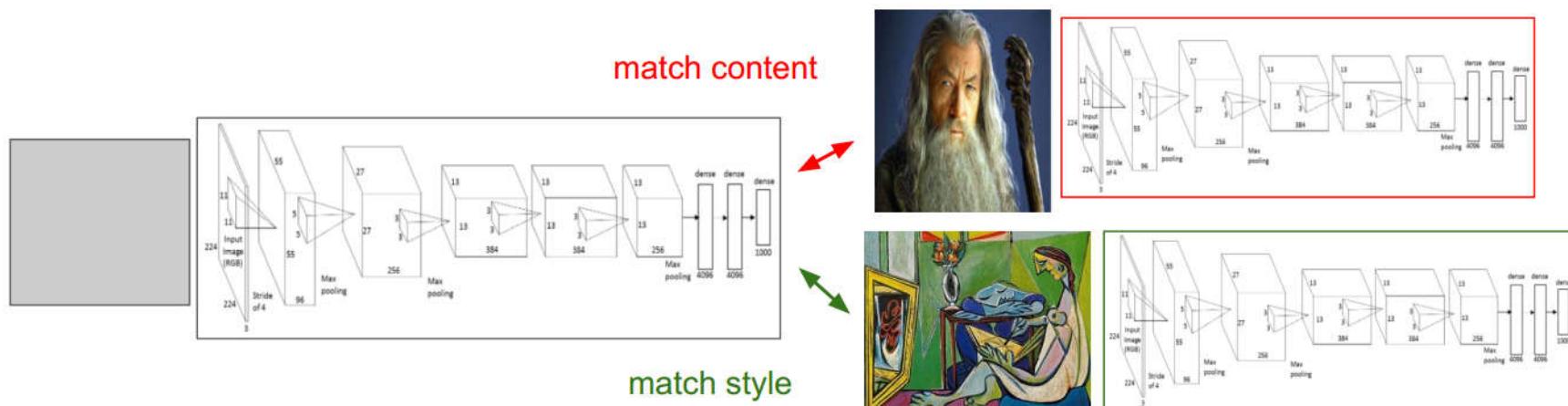
Neural Artistic Style

Step 3: optimize over image to have:

- The **content** of the content image (activations match content)
- The **style** of the style image (Gram matrices of activations match style)

$$L_{total} (\vec{p}, \vec{a}, \vec{x}) = \alpha L_{content} (\vec{p}, \vec{x}) + \beta L_{style} (\vec{a}, \vec{x})$$

(+ Total Variation regularization)



Fooling Neural Networks

Fooling Neural Networks

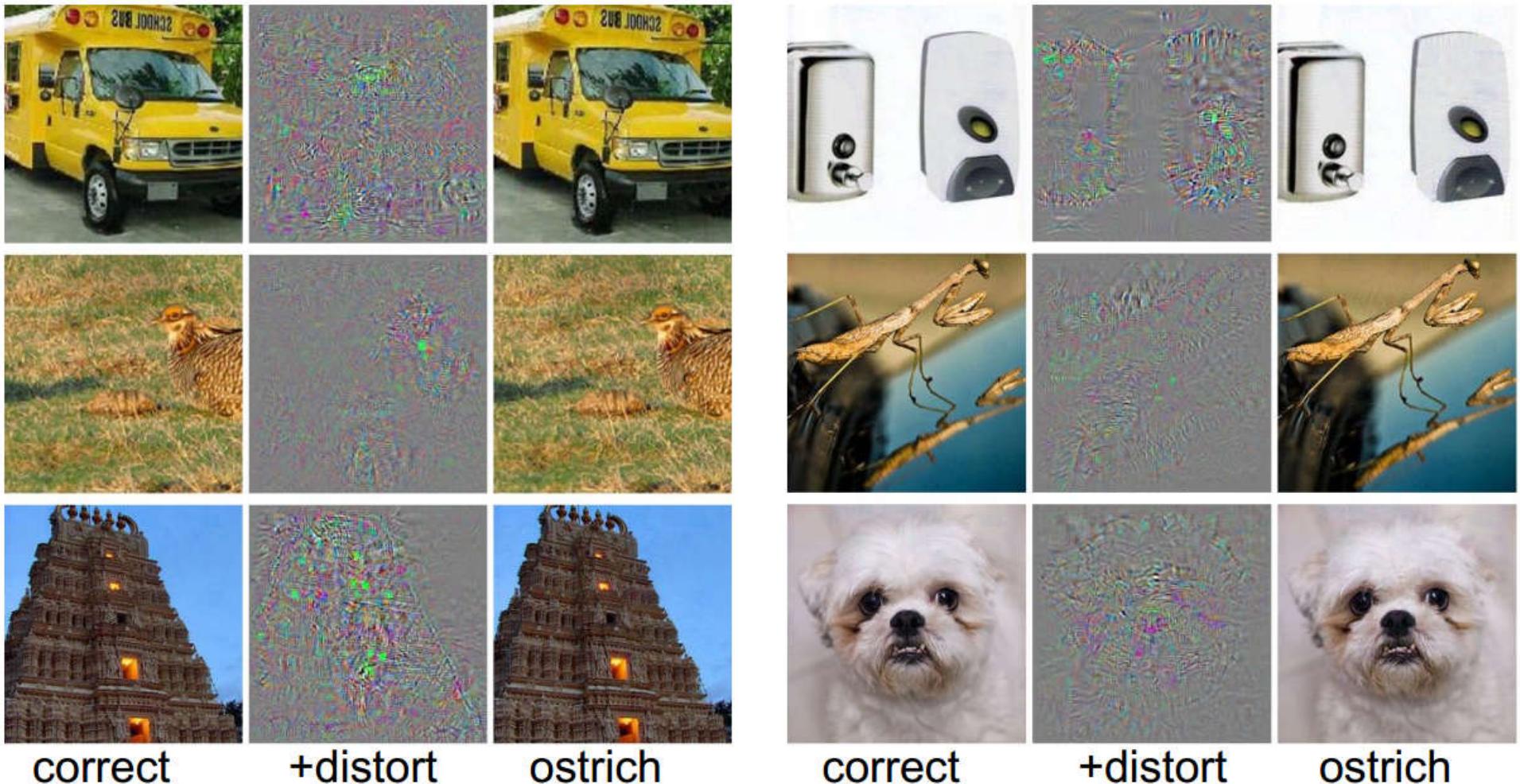
We can pose an optimization over the input image to maximize any class score. That seems useful.

Question: Can we use this to “fool” ConvNets?

spoiler alert: yeah

Fooling Neural Networks

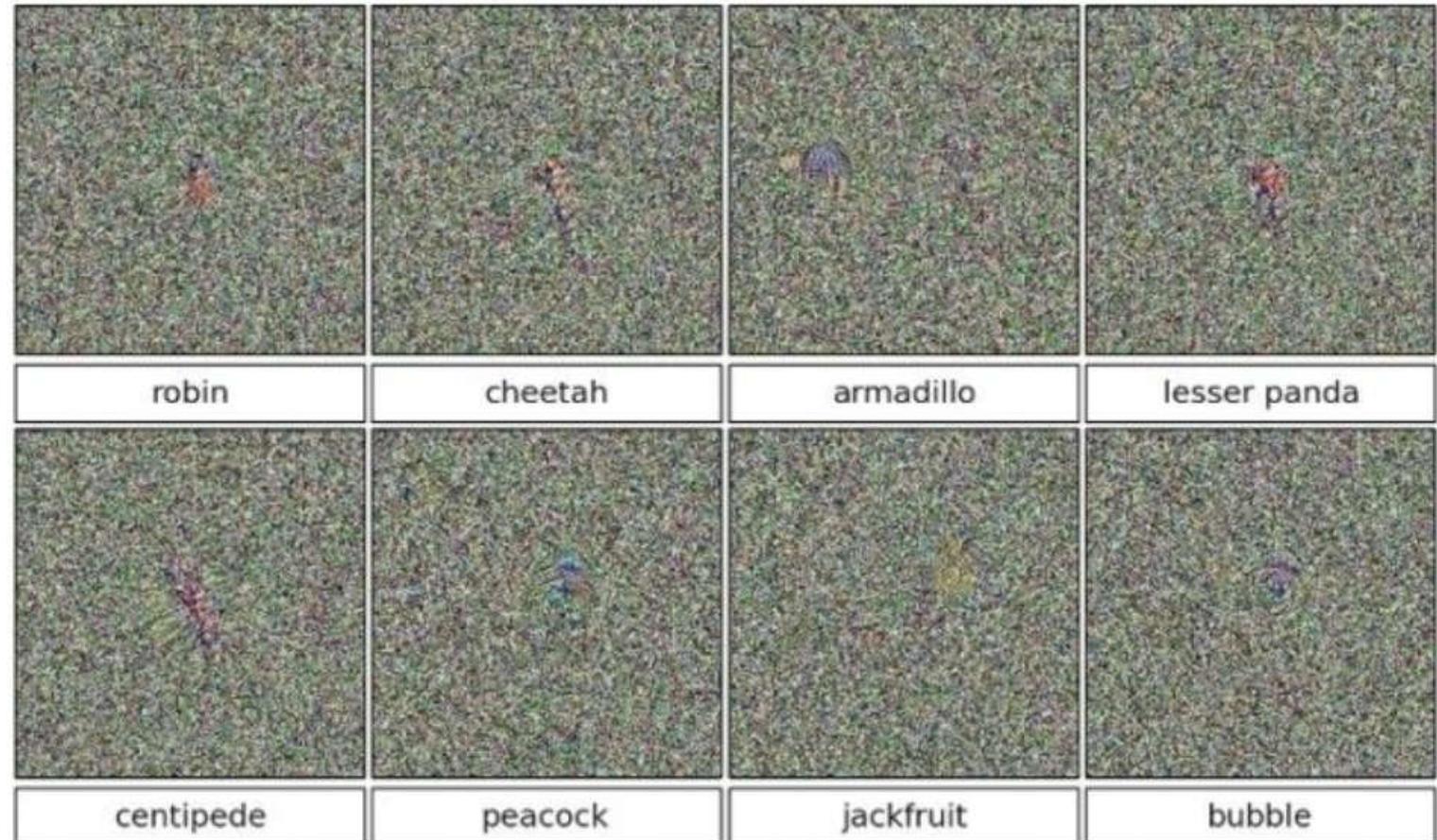
C. Szegedy, et al., Intriguing properties of neural networks, arXiv2013



Fooling Neural Networks

[Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images, Nguyen, Yosinski, Clune, 2014]

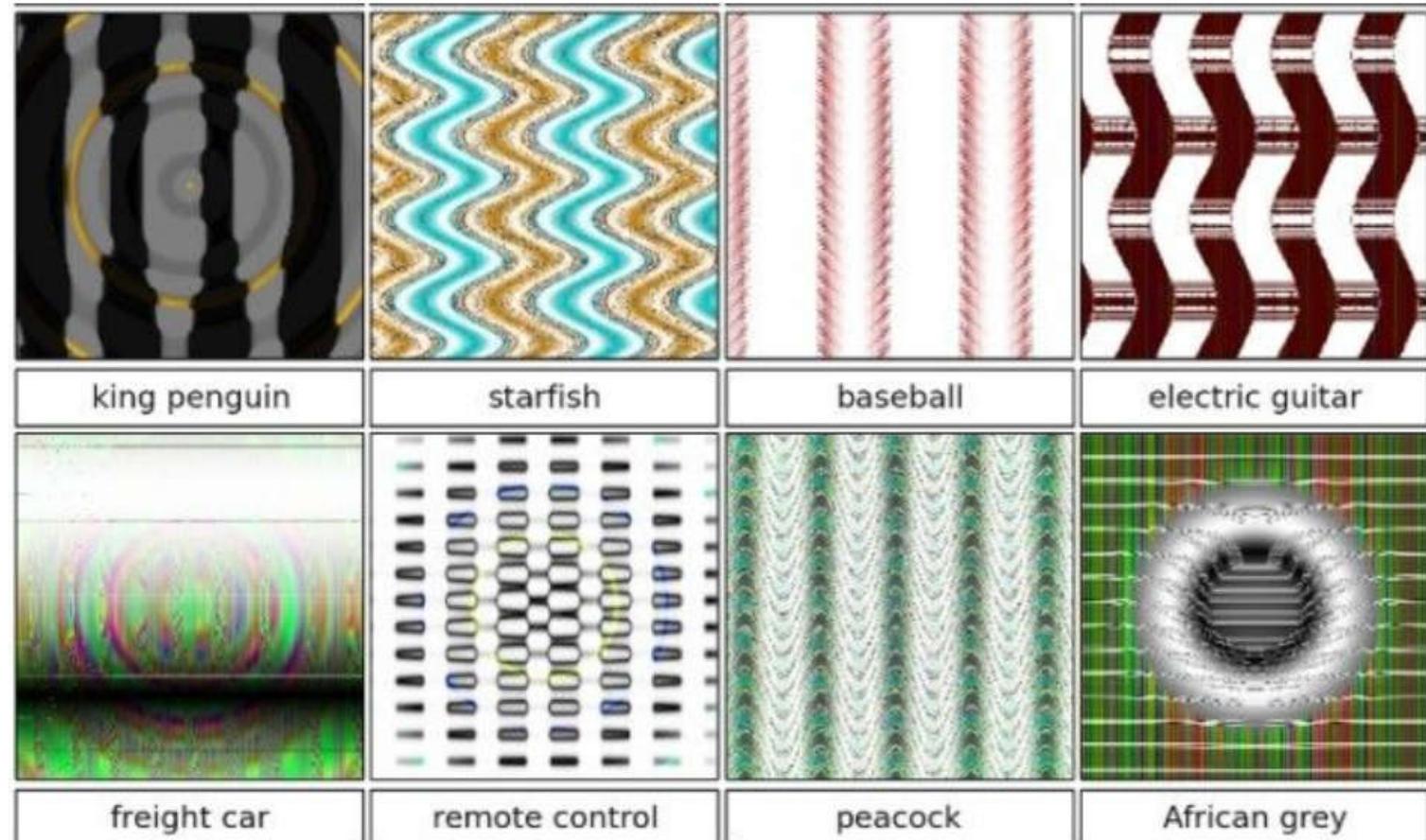
>99.6%
confidences



Fooling Neural Networks

[Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images, Nguyen, Yosinski, Clune, 2014]

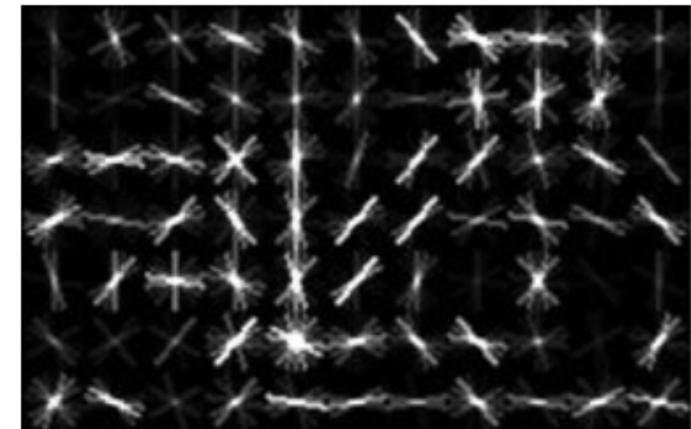
>99.6%
confidences



Fooling Neural Networks

These kinds of results were around even before ConvNets...

[Exploring the Representation Capabilities of the HOG Descriptor, Tatu et al., 2011]

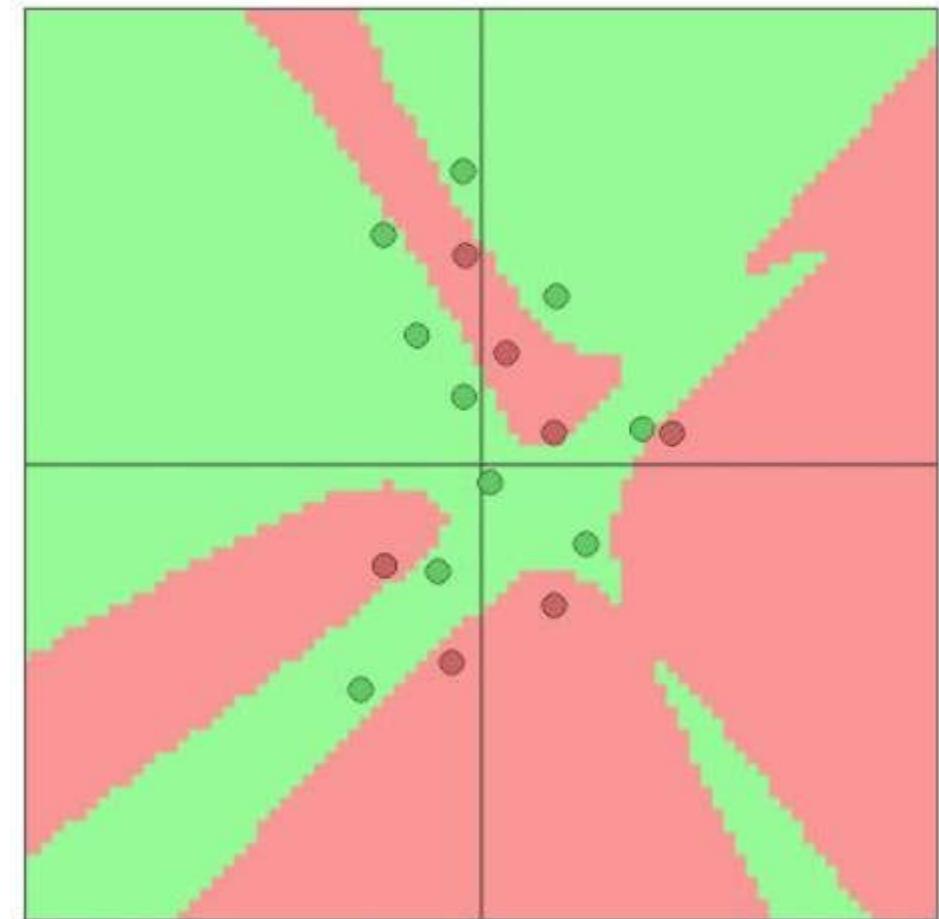


Identical HOG representation

Fooling Neural Networks

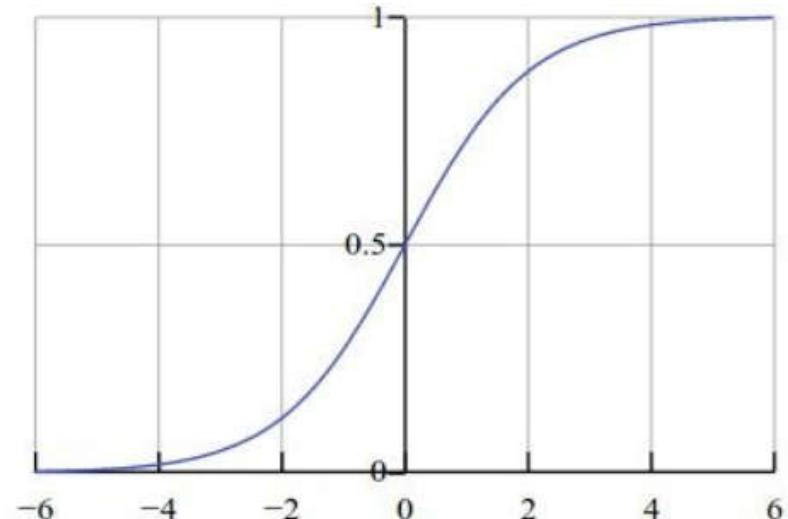
Explaining and harnessing adversarial examples, Goodfellow, Shlens & Szegedy, 2014

“Primary cause of neural networks’
vulnerability to adversarial
perturbation is their linear nature”



Lets Fool a Binary Linear Classifier (Logistic Regression)

$$P(y=1|x; w, b) = \frac{1}{1+e^{-(w^T x + b)}} = \sigma(w^T x + b)$$



Since the probabilities of class 1 and 0 sum to one, the probability for class 0 is $P(y=0 | x; w, b) = 1 - P(y=1 | x; w, b)$. Hence, an example is classified as a positive example ($y=1$) if $\sigma(w^T x + b) > 0.5$, or equivalently if the score $w^T x + b > 0$.

Lets Fool a Binary Linear Classifier (Logistic Regression)

x	2	-1	3	-2	2	2	1	-4	5	1	← input example
w	-1	-1	1	-1	1	-1	1	1	-1	1	← weights

Class 1 score = dot product:

$$= -2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3$$

$$\Rightarrow \text{Probability of class 1 is } 1/(1+e^{-(-3)}) = 0.0474$$

i.e. the classifier is **95%** certain that this is class 0 example.

$$P(y=1|x; w, b) = \frac{1}{1+e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

Lets Fool a Binary Linear Classifier (Logistic Regression)

x	2	-1	3	-2	2	2	1	-4	5	1	← input example
w	-1	-1	1	-1	1	-1	1	1	-1	1	← weights
	?	?	?	?	?	?	?	?	?	?	

Adversarial x

Class 1 score = dot product:

$$= -2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3$$

⇒ Probability of class 1 is $1/(1+e^{-(-3)}) = 0.0474$

i.e. the classifier is **95%** certain that this is class 0 example.

$$P(y=1|x; w, b) = \frac{1}{1+e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

Lets Fool a Binary Linear Classifier (Logistic Regression)

x	2	-1	3	-2	2	2	1	-4	5	1	← input example
w	-1	-1	1	-1	1	-1	1	1	-1	1	← weights
	1.5	-1.5	3.5	-2.5	2.5	1.5	1.5	-3.5	4.5	1.5	

Adversarial x

$$P(y=1|x; w, b) = \frac{1}{1+e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

Class 1 score before:

$$-2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3$$

⇒ Probability of class 1 is $1/(1+e^{-(-3)}) = 0.0474$

$$-1.5 + 1.5 + 3.5 + 2.5 + 2.5 - 1.5 + 1.5 - 3.5 - 4.5 + 1.5 = 2$$

⇒ Probability of class 1 is now $1/(1+e^{-(-2)}) = 0.88$

i.e. we improved the class 1 probability from 5% to 88%

Lets Fool a Binary Linear Classifier (Logistic Regression)

x	2	-1	3	-2	2	2	1	-4	5	1	← input example
w	-1	-1	1	-1	1	-1	1	1	-1	1	← weights
	1.5	-1.5	3.5	-2.5	2.5	1.5	1.5	-3.5	4.5	1.5	

Adversarial x

This was only with 10 input dimensions. A 224x224 input image has 150,528.

Class 1 score before:

$$-2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3$$

(It's significantly easier with more numbers, need smaller nudge for each)

$$\Rightarrow \text{Probability of class 1 is } 1/(1+e^{(-(-3))}) = 0.0474$$

$$-1.5 + 1.5 + 3.5 + 2.5 + 2.5 - 1.5 + 1.5 - 3.5 - 4.5 + 1.5 = 2$$

$$\Rightarrow \text{Probability of class 1 is now } 1/(1+e^{(-(2)}) = 0.88$$

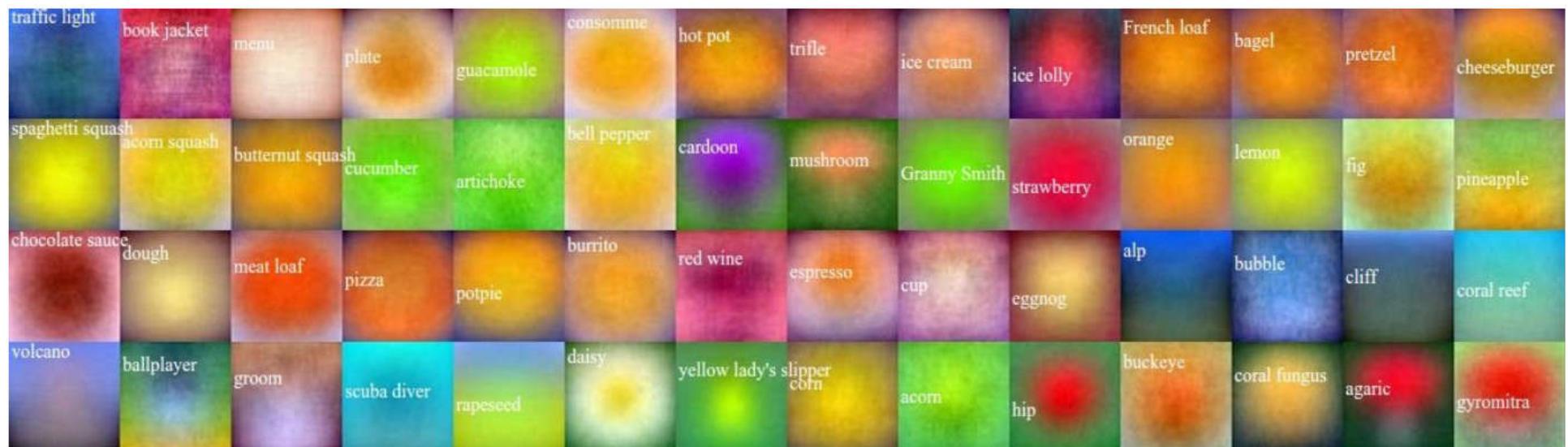
i.e. we improved the class 1 probability from 5% to 88%

Blog Post: Breaking Linear Classifiers on ImageNet

Recall CIFAR-10 linear classifiers:



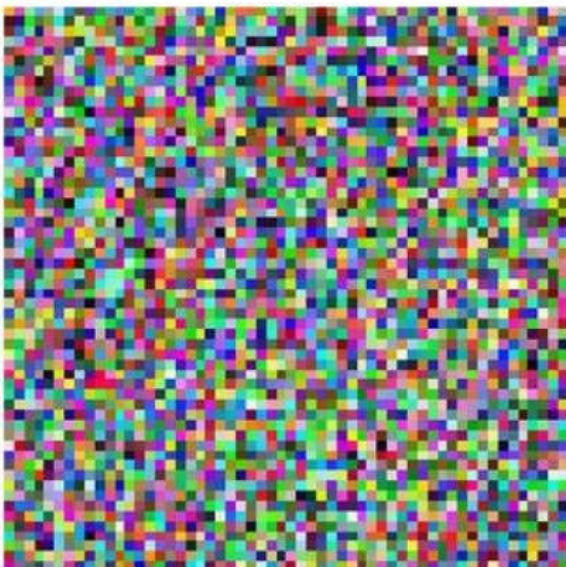
ImageNet classifiers:



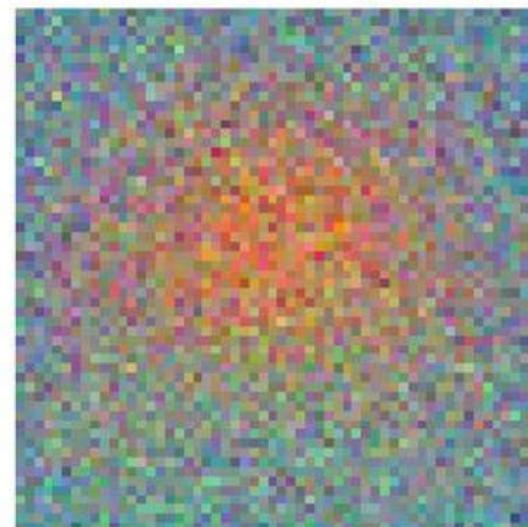
Blog Post: Breaking Linear Classifiers on ImageNet

Mix in a tiny bit of
Goldfish classifier weights

0.9% bobsled

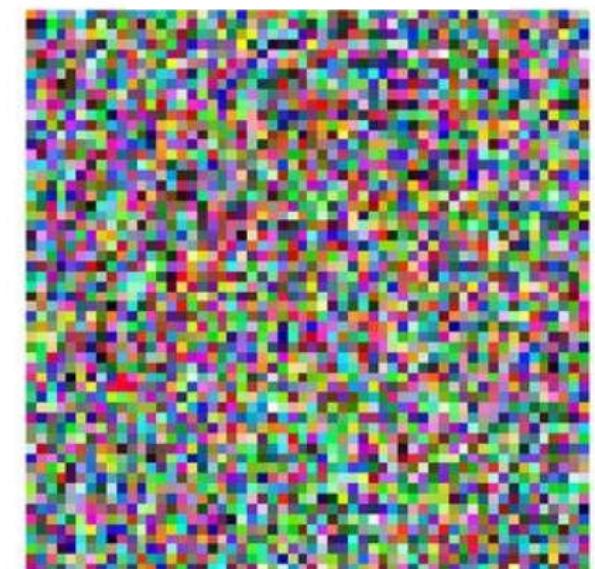


+



=

100.0% goldfish



100% Goldfish

Blog Post: Breaking Linear Classifiers on ImageNet

1.0% kit fox



8.0% goldfish



Blog Post: Breaking Linear Classifiers on ImageNet

1.0% kit fox



3.9% school bus



8.3% goldfish



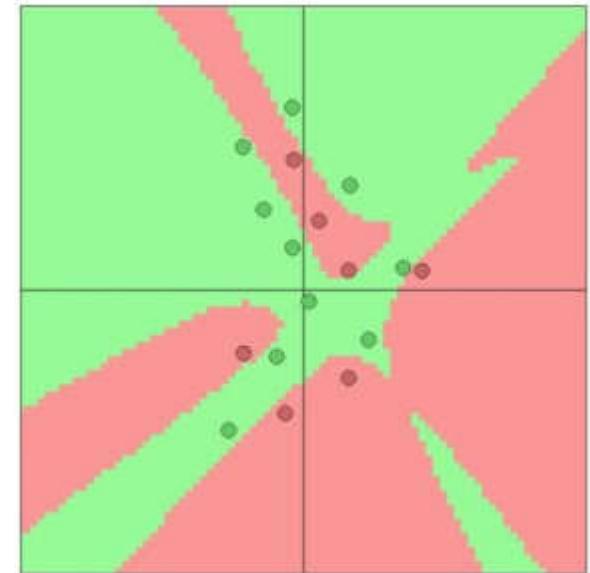
12.5% daisy



Fooling Deep Neural Networks

Explaining and Harnessing Adversarial Examples [Goodfellow, Shlens & Szegedy, 2014]

“Primary cause of neural networks’ vulnerability to adversarial perturbation is their linear nature” (and very high-dimensional, sparsely-populated input spaces)



In particular, this is not a problem with Deep Learning, and has little to do with ConvNets specifically. Same issue would come up with Neural Nets in any other modalities.

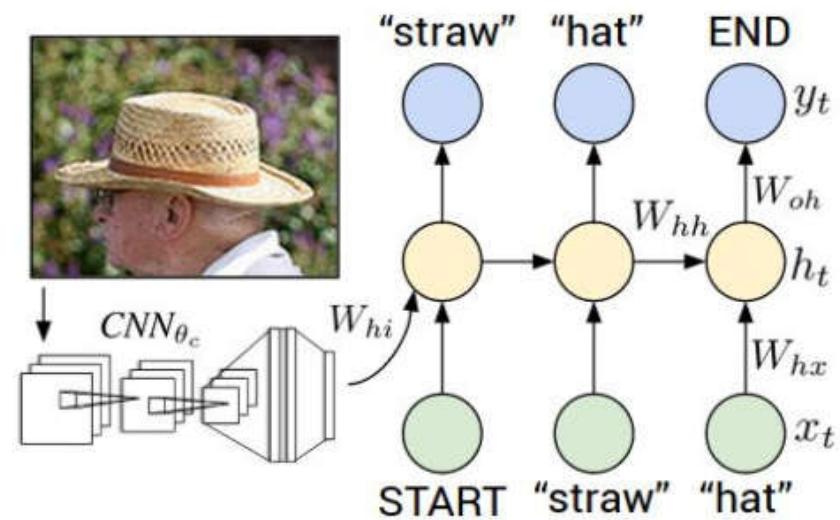
Summary

Backpropping to the image is powerful. It can be used for:

- **Understanding** (e.g. visualize optimal stimuli for arbitrary neurons)
- **Segmenting** objects in the image (kind of)
- **Inverting** codes and introducing privacy concerns
- **Fun** (NeuralStyle/DeepDream)
- **Confusion and chaos** (Adversarial examples)

Next Lecture

Image Captioning Recurrent Neural Networks RNN Language Models



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."

Thank you for your attention!
