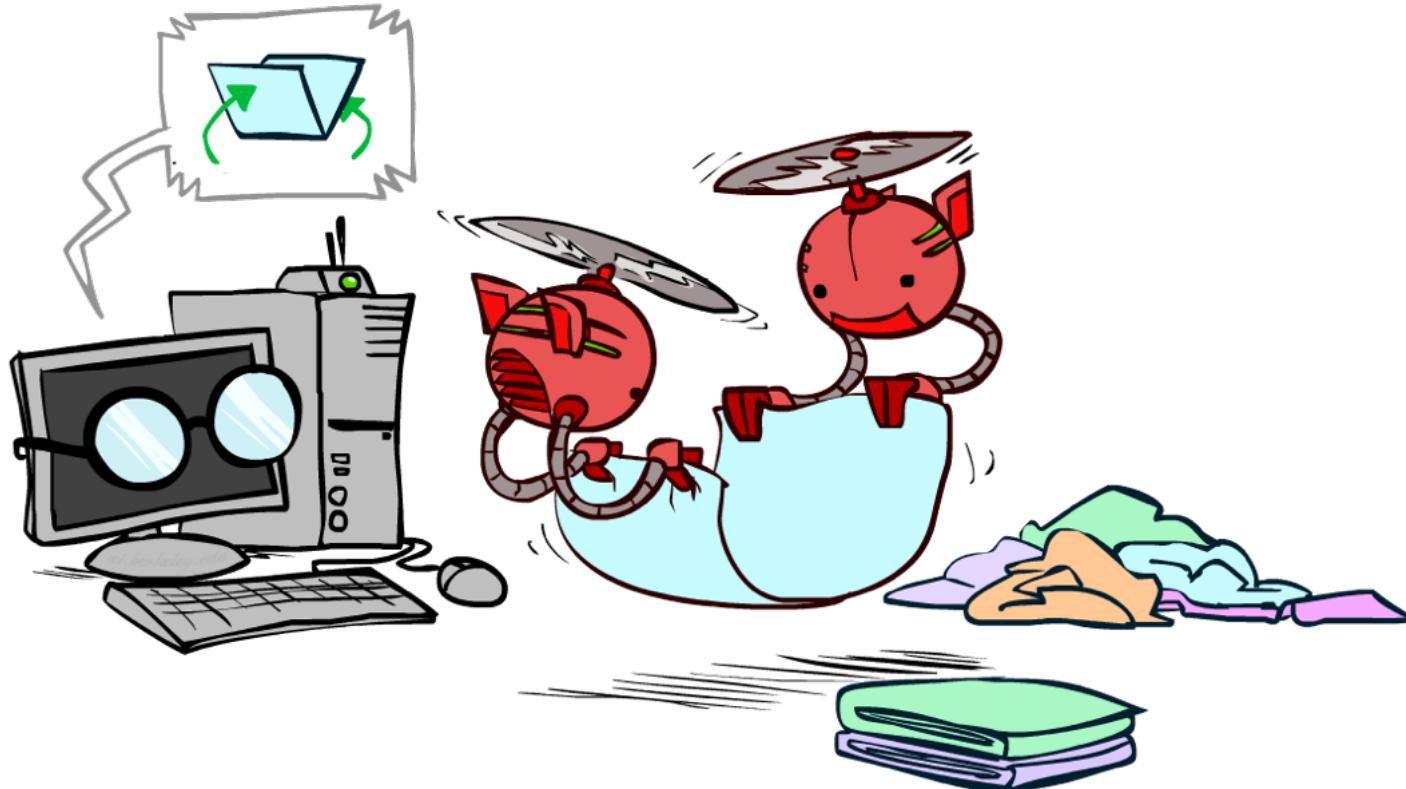


# CSE 5522: Artificial Intelligence II

## Advanced Applications: Computer Vision \*



Instructor: Wei Xu

Ohio State University

[These slides were adapted from CS231n Computer Vision at Stanford, and CS188 Intro to AI at UC Berkeley]

# Face Detection

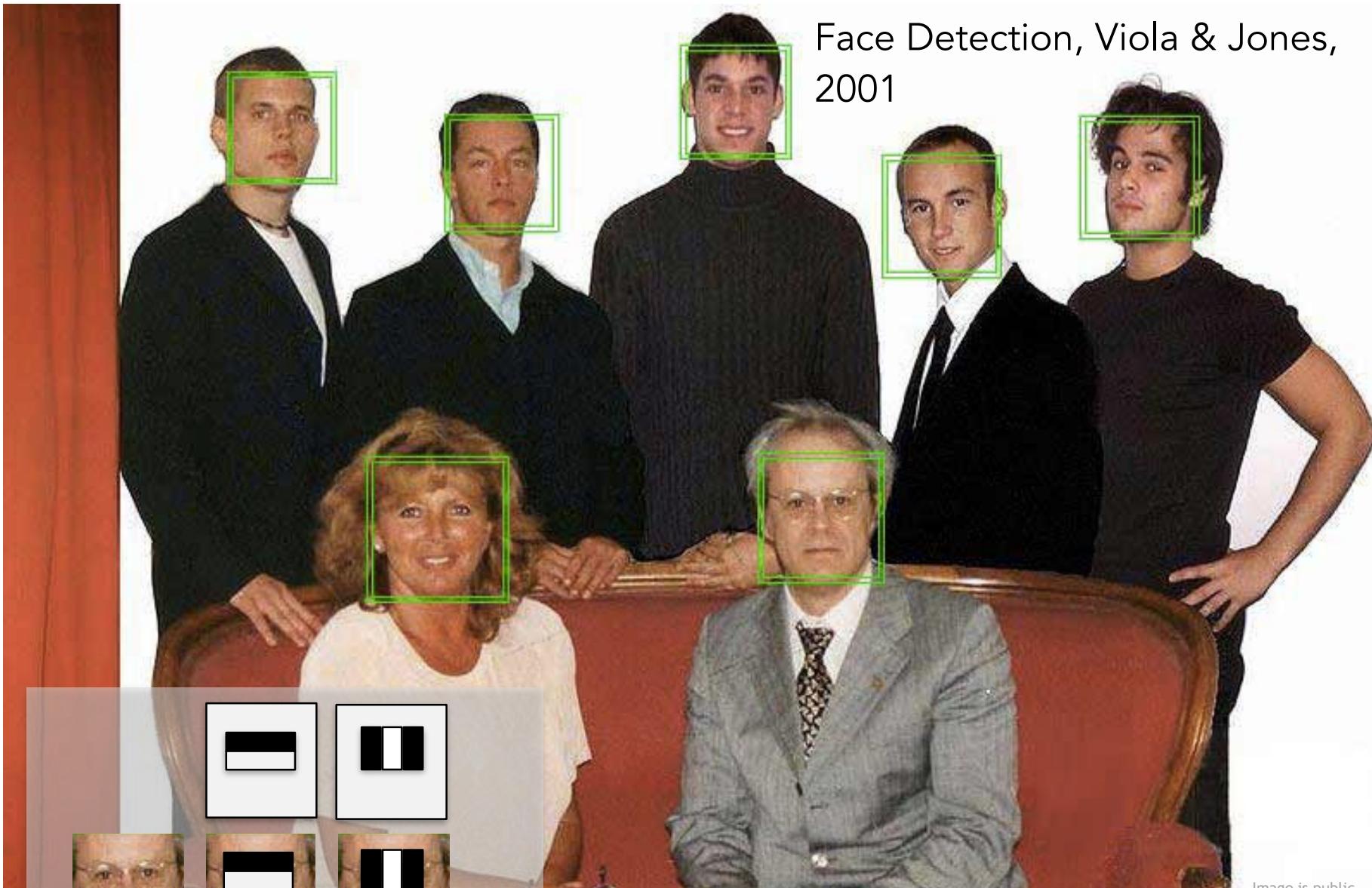


Image is public

# Object Detection & Image Segmentation

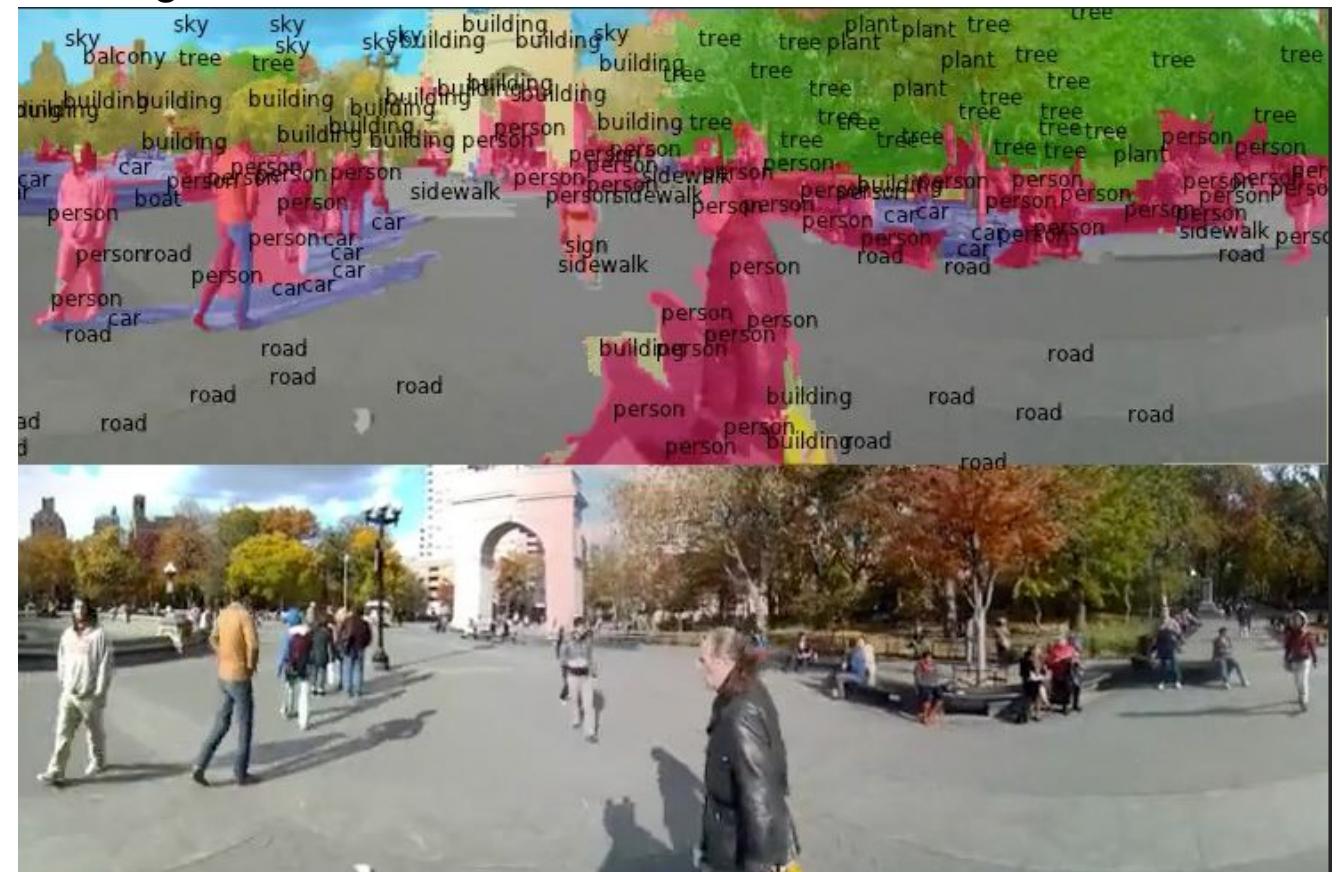
# Detection



Figures copyright Shaoqing Ren, Kaiming He, Ross Girshick, Jian Sun, 2015. Reproduced with permission.

## [Faster R-CNN: Ren, He, Girshick, Sun 2015]

## Segmentation

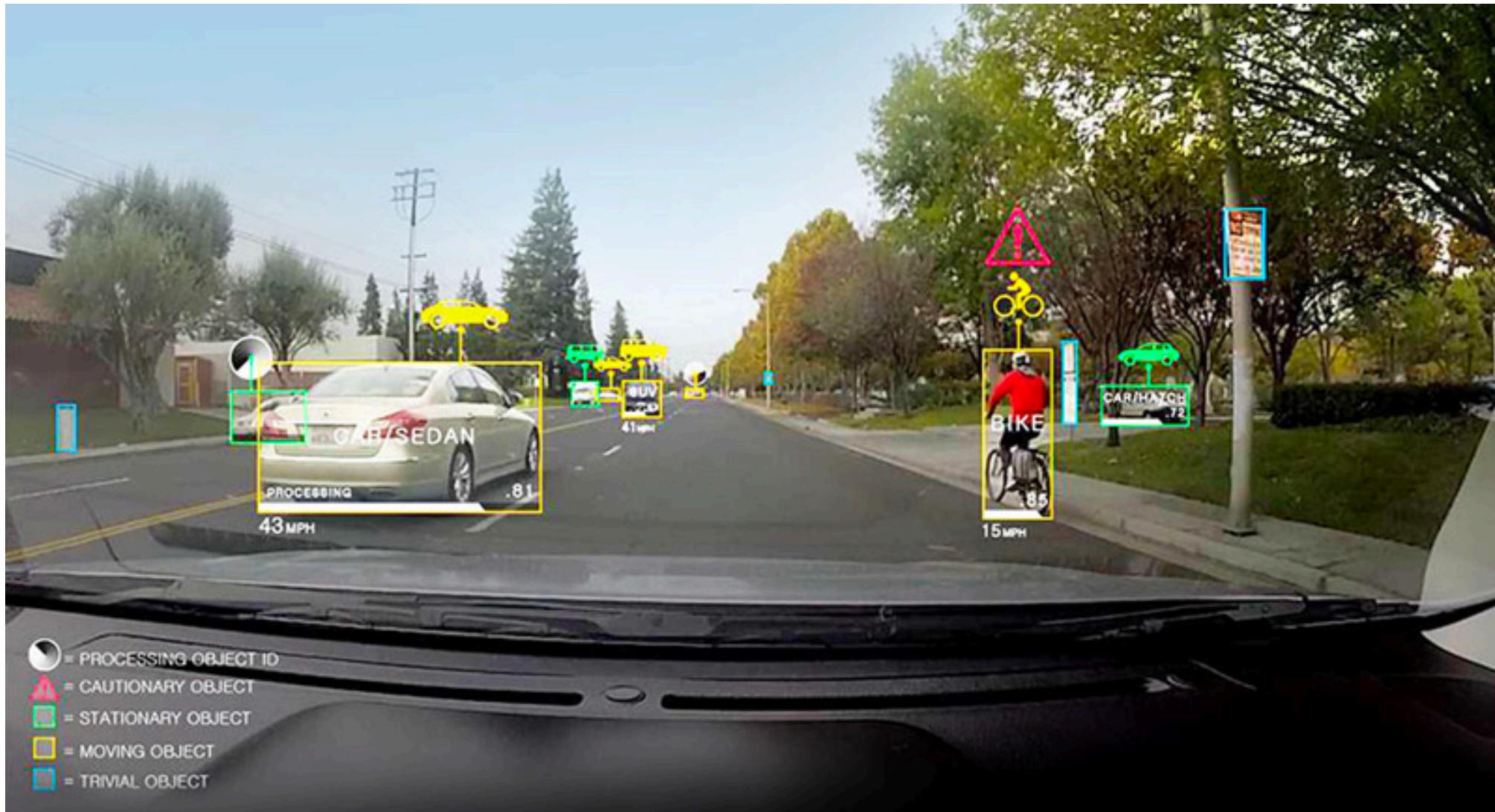


---

Figures copyright Clement Farabet, 2012.  
Reproduced with permission.

[Farabet et al., 2012]

# Self-driving Cars



NVIDIA DRIVE PX

# Whale and Road Recognition

[This image](#) by Christin Khan is in the public domain and originally came from the U.S. NOAA.



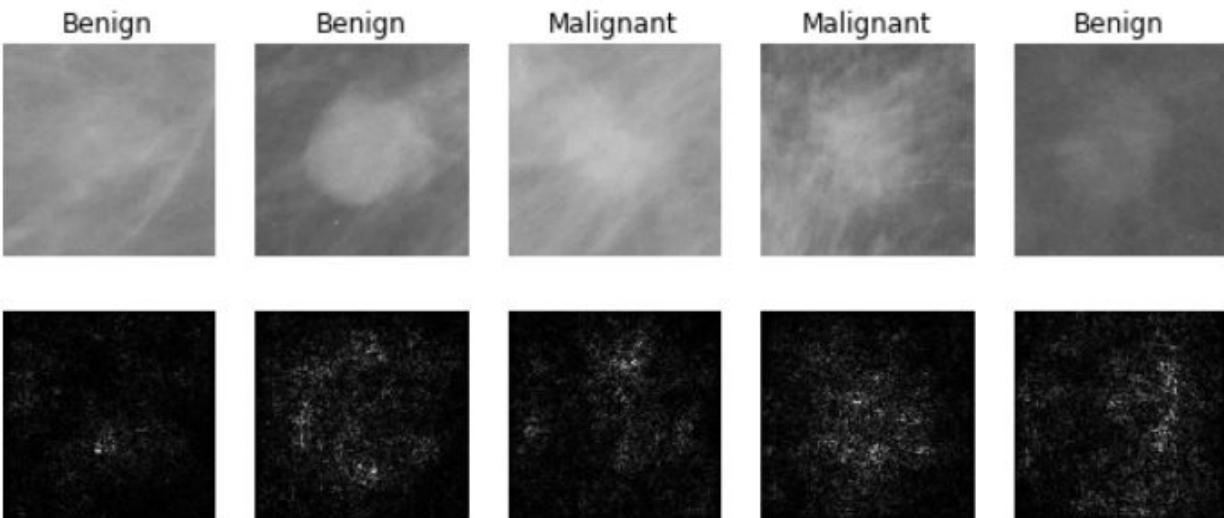
*Whale recognition, Kaggle Challenge*



*Mnih and Hinton, 2010*

Photo and figure by Lane McIntosh; not actual example from Mnih and Hinton, 2010 paper.

# And Many More ...



[Levy et al. 2016]

Figure copyright Levy et al. 2016.  
Reproduced with permission.



[Dieleman et al. 2014]

From left to right: [public domain by NASA](#), usage [permitted](#) by  
ESA/Hubble, [public domain by NASA](#), and [public domain](#).



[Sermanet et al. 2011]  
[Ciresan et al.]

Photos by Lane McIntosh.  
Copyright CS231n 2017.

# Galaxy Morphology Prediction

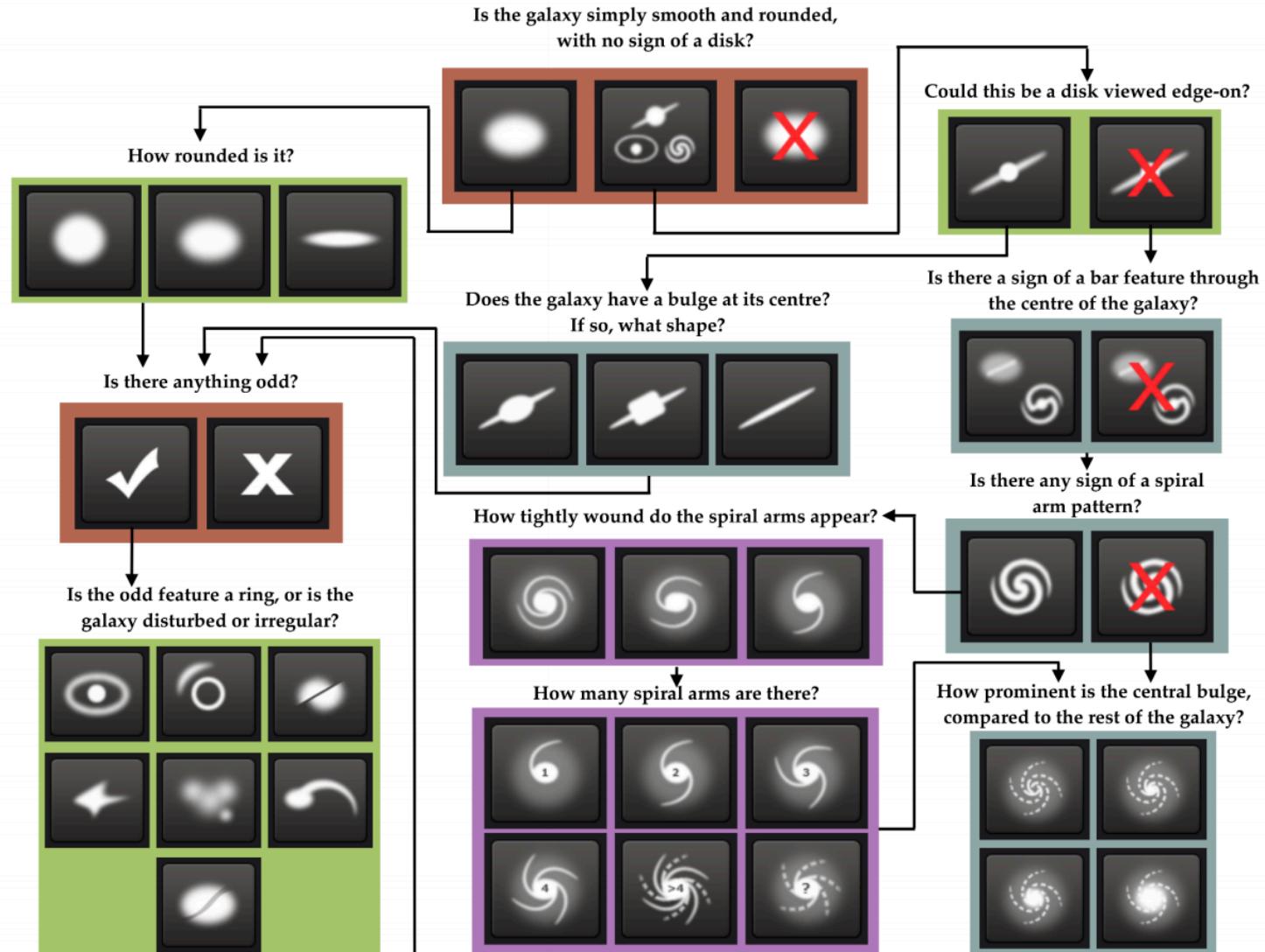


Figure 1. The Galaxy Zoo 2 decision tree. Reproduced from Figure 1 in Willett et al. (2013).

# Image and Video Captioning

No errors



*A white teddy bear sitting in the grass*



*A man riding a wave on top of a surfboard*

Minor errors



*A man in a baseball uniform throwing a ball*



*A cat sitting on a suitcase on the floor*

Somewhat related



*A woman is holding a cat in her hand*



*A woman standing on a beach holding a surfboard*

## Image Captioning

[Vinyals et al., 2015]  
[Karpathy and Fei-Fei, 2015]

All images are CC0 Public domain:  
<https://pixabay.com/en/luggage-antique-cat-1643010/>  
<https://pixabay.com/en/teddy-plush-bears-cute-teddy-bear-1623436/>  
<https://pixabay.com/en/surf-wave-summer-sport-litoral-1668716/>  
<https://pixabay.com/en/woman-female-model-portrait-adult-983967/>  
<https://pixabay.com/en/handstand-lake-meditation-496008/>  
<https://pixabay.com/en/baseball-player-shortstop-infield-1045263/>

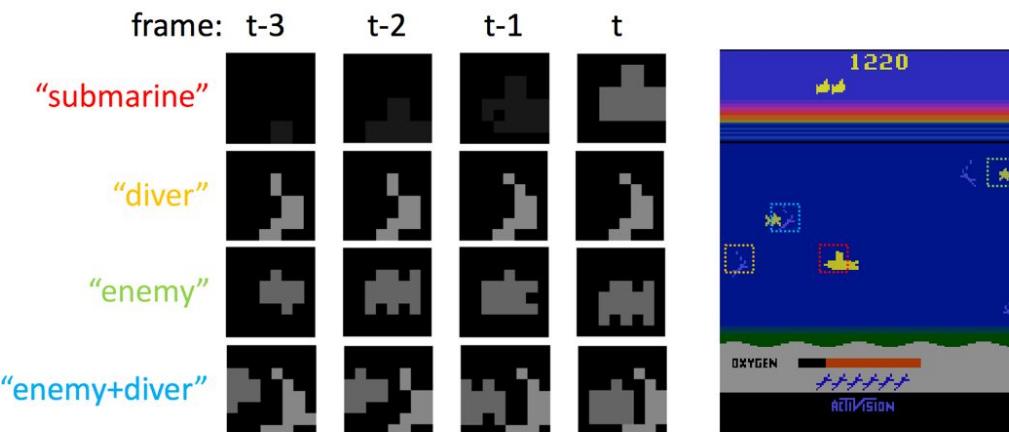
Captions generated by Justin Johnson using [Neuraltalk2](#)

# Pose Estimation & Atari Games

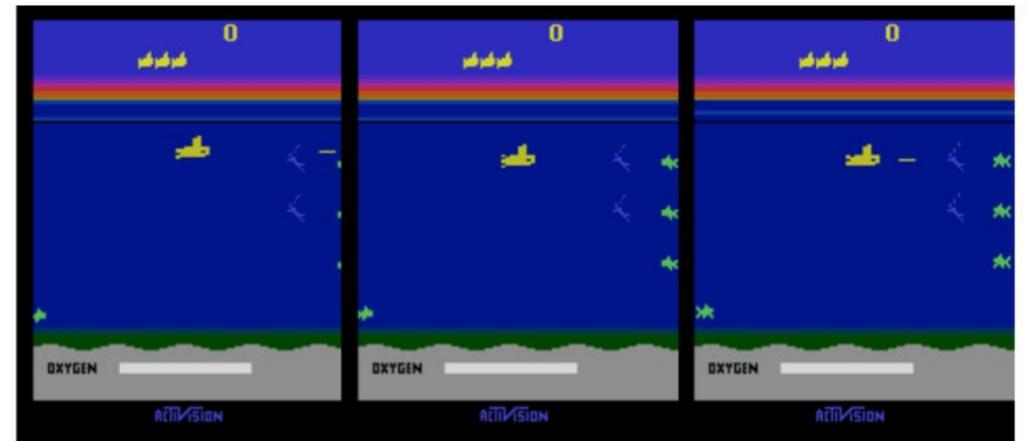


Images are examples of pose estimation, not actually from Toshev & Szegedy 2014. Copyright Lane McIntosh.

[Toshev, Szegedy 2014]

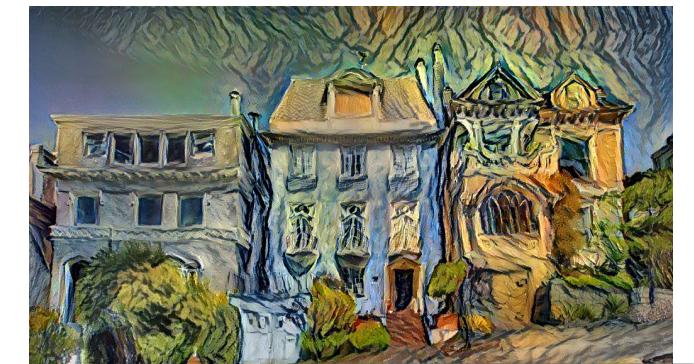
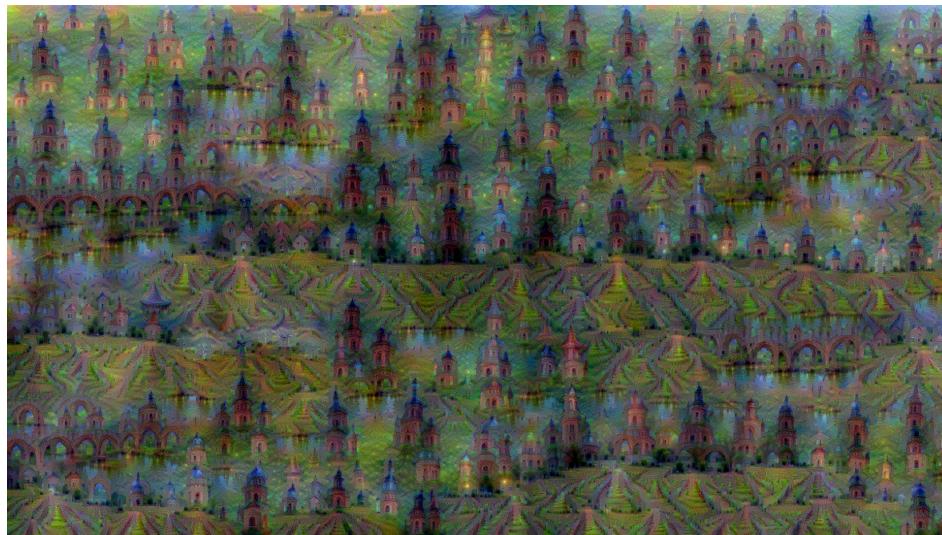
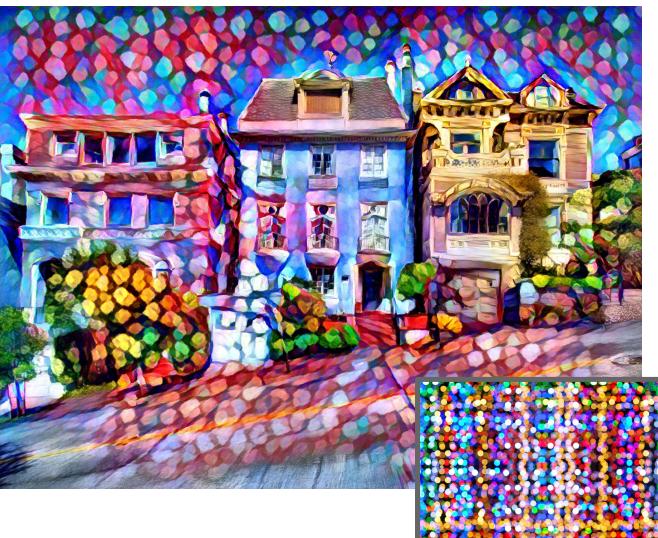
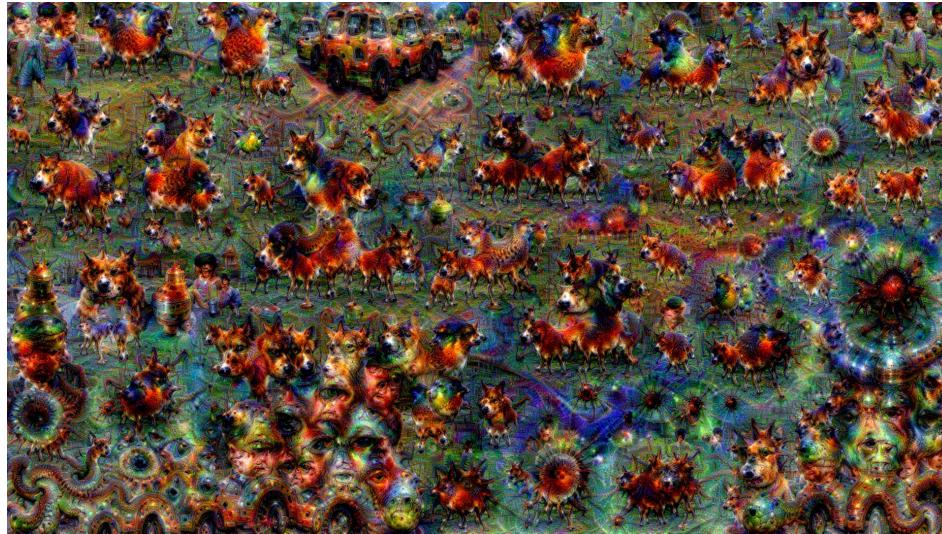


[Guo et al. 2014]



Figures copyright Xiaoxiao Guo, Satinder Singh, Honglak Lee, Richard Lewis, and Xiaoshi Wang, 2014. Reproduced with permission.

# DeepDream and Style Transfer



[Original image](#) is CC0 public domain

[Starry Night](#) and [Tree Roots](#) by Van Gogh are in the public domain

[Bokeh image](#) is in the public domain

Stylized images copyright Justin Johnson, 2017;  
reproduced with permission

Figures copyright Justin Johnson, 2015. Reproduced with permission. Generated using the Inceptionism approach from a [blog post](#) by Google Research.

Gatys et al, "Image Style Transfer using Convolutional Neural Networks", CVPR 2016

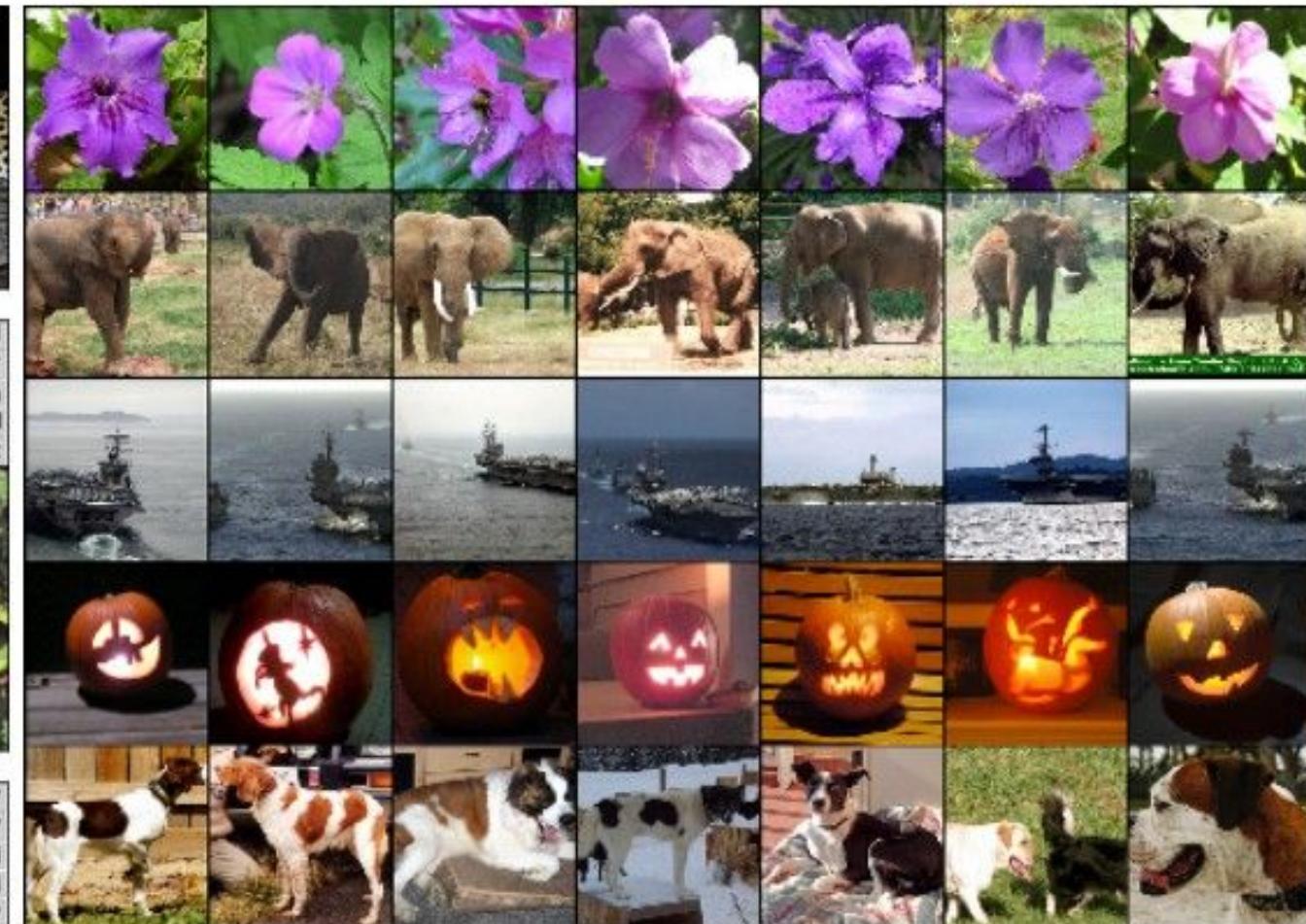
Gatys et al, "Controlling Perceptual Factors in Neural Style Transfer", CVPR 2017

# Image Classification & Retrieval

Classification



Retrieval



# PASCAL Visual Object Challenge

(20 object categories)

[Everingham et al. 2006-2012]

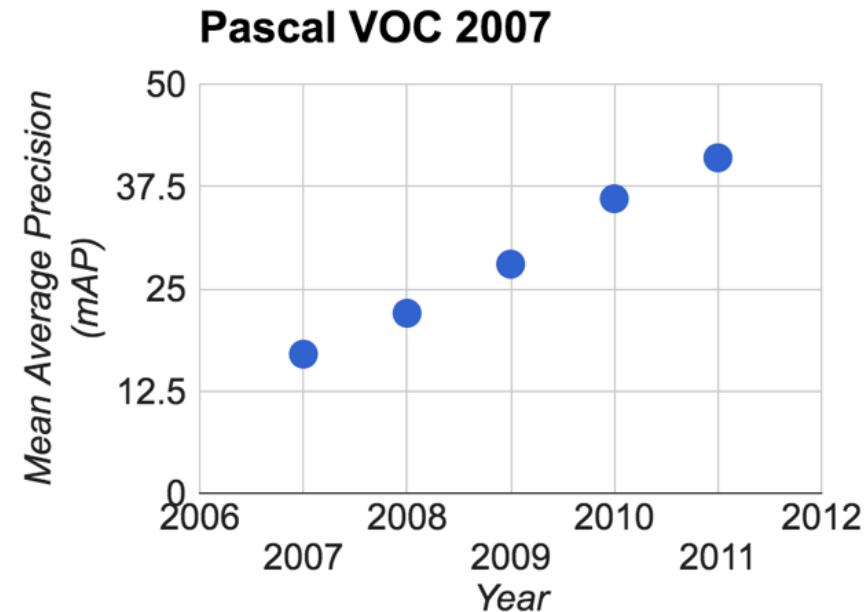
[Image is CC BY-SA 3.0](#)



[Image is CC0 1.0 public domain](#)



[This image is licensed under  
CC BY-SA 2.0; changes made](#)



# ImageNet



[www.image-net.org](http://www.image-net.org)

**22K** categories and **14M** images

- Animals
  - Bird
  - Fish
  - Mammal
  - Invertebrate
- Plants
  - Tree
  - Flower
  - Food
  - Materials
- Structures
  - Artifact
  - Tools
  - Appliances
  - Structures
- Person
- Scenes
  - Indoor
  - Geological Formations
- Sport Activities



# ImageNet Large Scale Visual Recognition Challenge

**Steel drum**

The Image Classification Challenge:

1,000 object classes

1,431,167 images



**Output:**

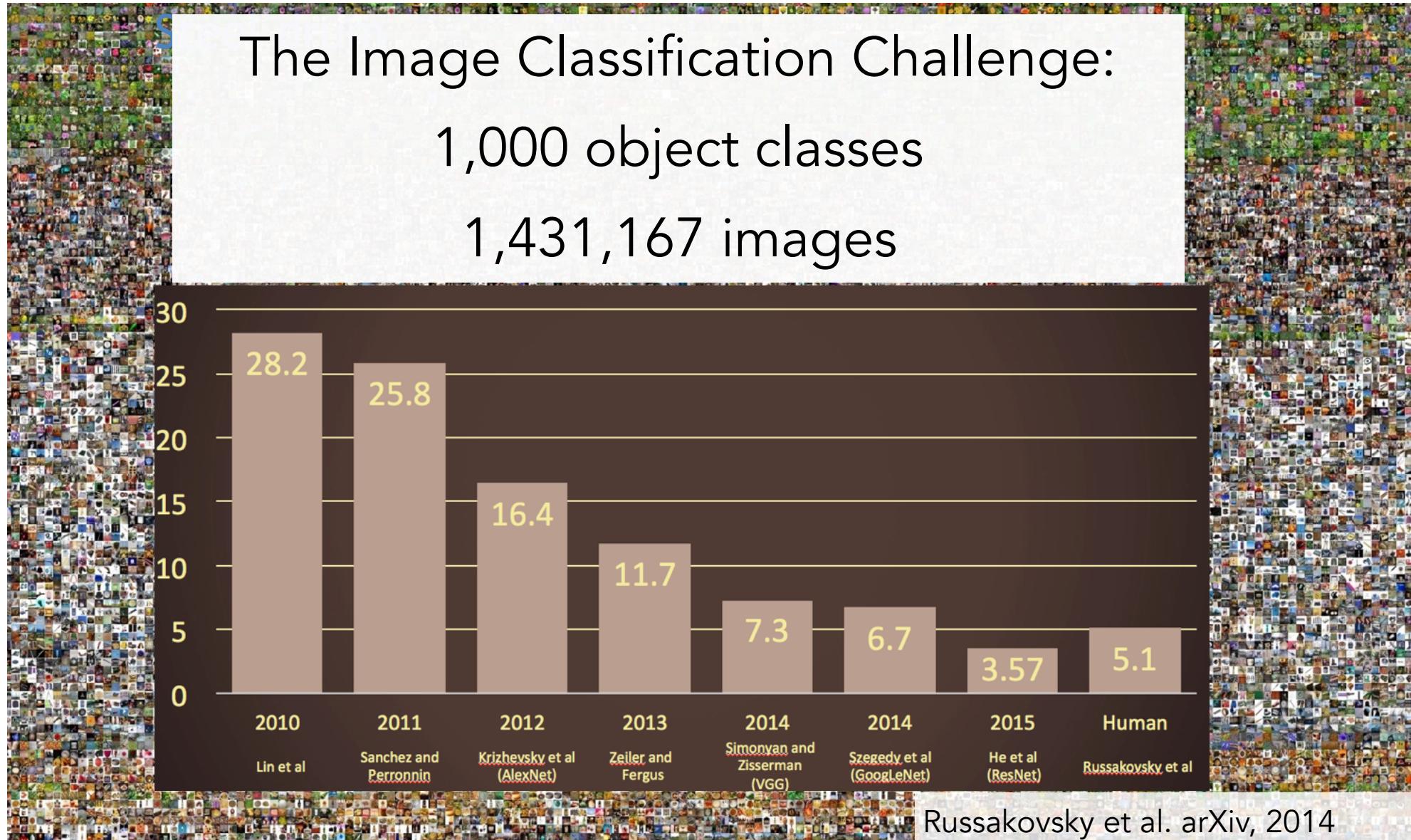
- Scale
- T-shirt
- Steel drum
- Drumstick
- Mud turtle

**Output:**

- Scale
- T-shirt
- Giant panda
- Drumstick
- Mud turtle

Russakovsky et al. arXiv, 2014

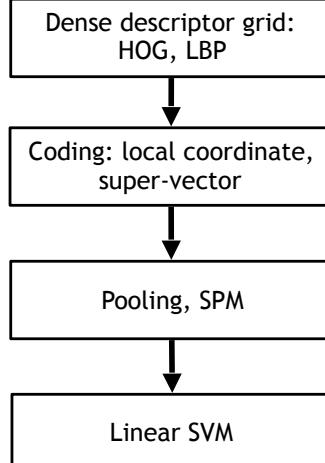
# ImageNet Large Scale Visual Recognition Challenge



# ImageNet Large Scale Visual Recognition Challenge

Year 2010

NEC-UIUC

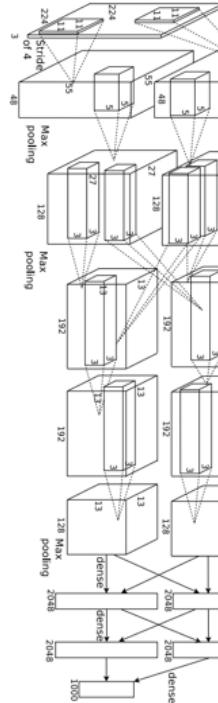


[Lin CVPR 2011]

[Lion image](#) by Swissfrog is licensed under [CC BY 3.0](#)

Year 2012

SuperVision

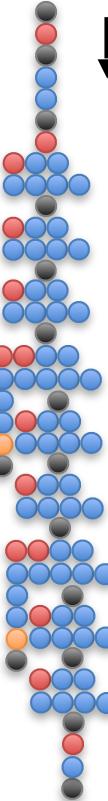


[Krizhevsky NIPS 2012]

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

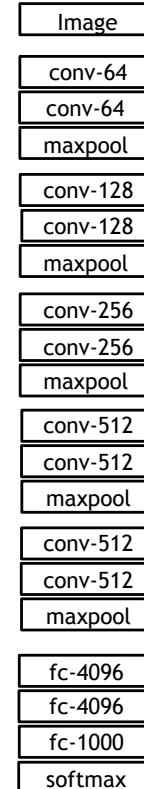
Year 2014

GoogLeNet



[Szegedy arxiv 2014]

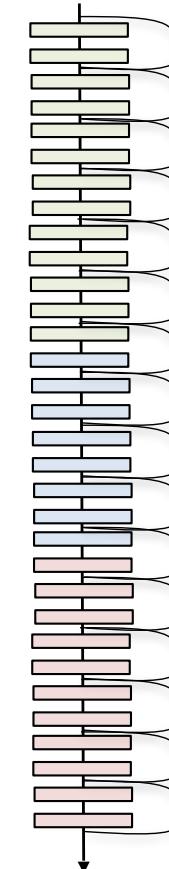
VGG



[Simonyan arxiv 2014]

Year 2015

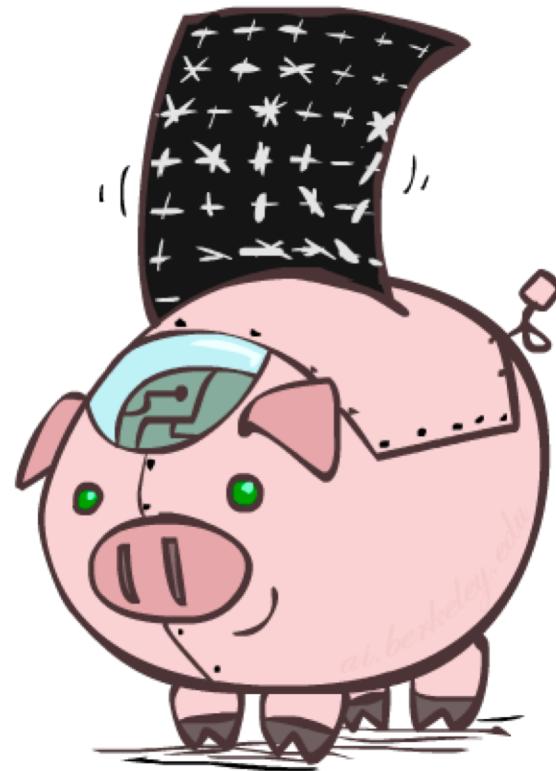
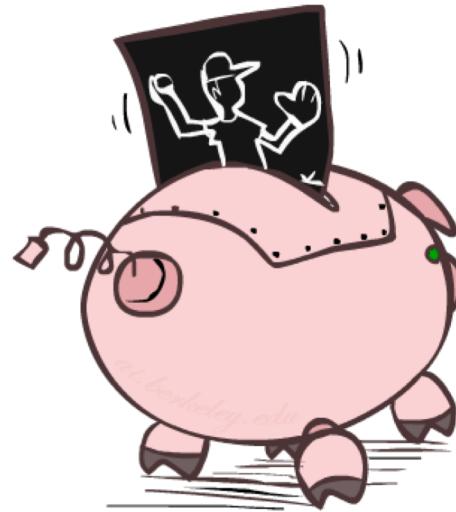
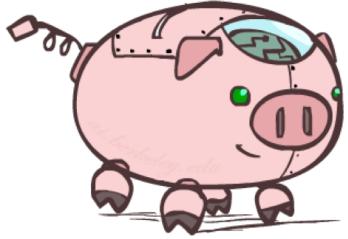
MSRA



[He ICCV 2015]

# Traditional Approach

---

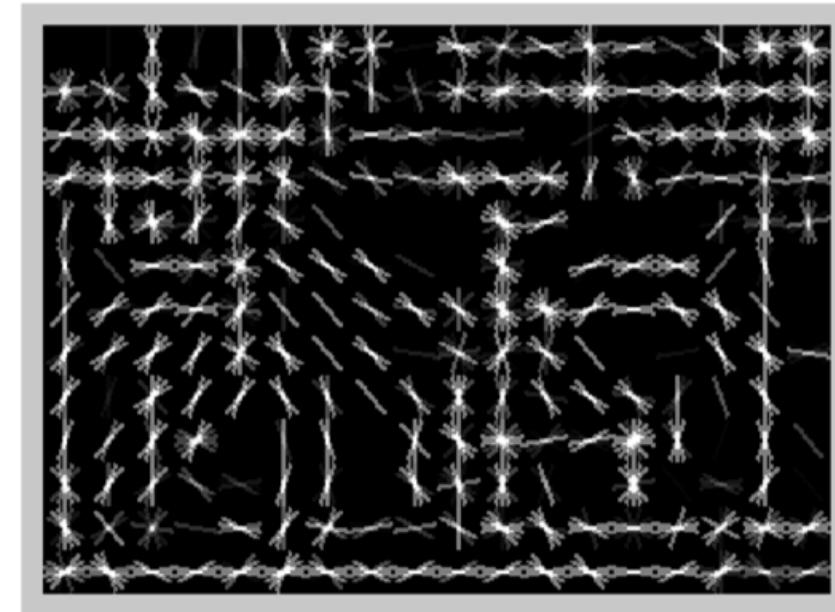


# Features and Generalization

---



Image

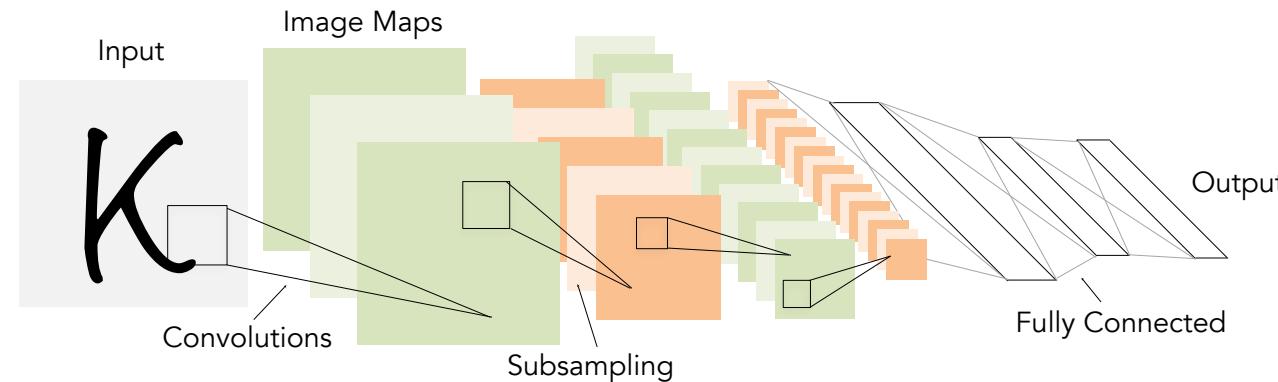


HoG

# Convolutional Neural Networks (CNN)

1998

LeCun et al.



# of transistors



$10^6$

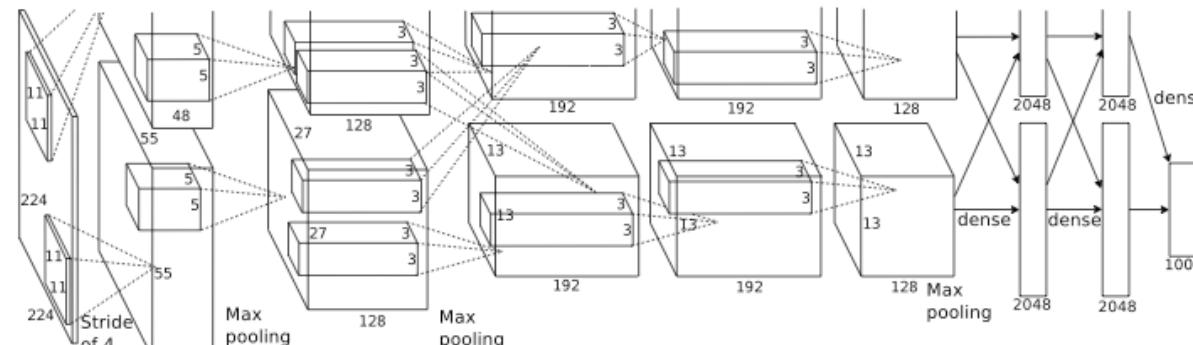
pentium<sup>®</sup> II

# of pixels used in training

$10^7$

2012

Krizhevsky et al.



# of transistors



$10^9$

GPUs



# of pixels used in training

$10^{14}$

# Image Classification

---



This image by [Nikita](#) is  
licensed under [CC-BY 2.0](#)

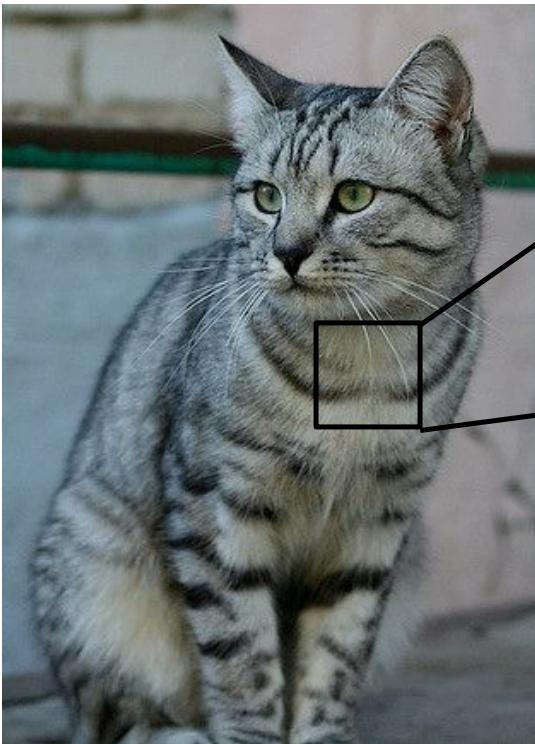
(assume given set of discrete labels)  
{dog, cat, truck, plane, ...}



cat

# Image Representation

## The Problem: Semantic Gap



This image by Nikita is  
licensed under CC-BY 2.0

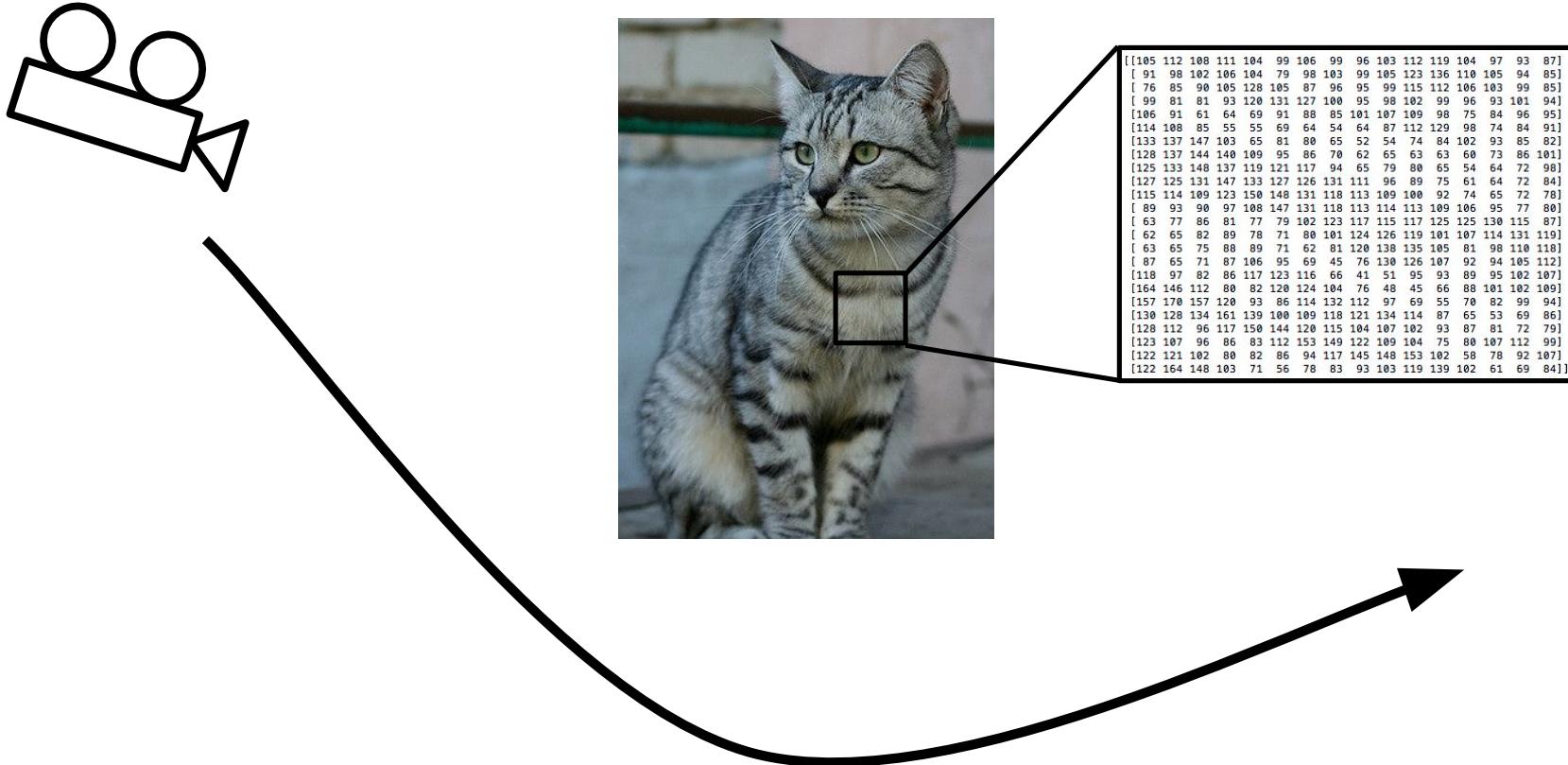
[105 112 108 111 104 99 106 99 96 103 112 119 104 97 93 87]
[ 91 98 102 106 104 79 98 103 99 105 123 136 110 105 94 85]
[ 76 85 90 105 128 105 87 96 95 99 115 112 106 103 99 85]
[ 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94]
[106 91 61 64 69 91 88 85 101 107 109 98 75 84 96 95]
[114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 91]
[133 137 147 103 65 81 80 65 52 54 74 84 102 93 85 82]
[128 137 144 140 109 95 86 70 62 65 63 63 60 73 86 101]
[125 133 148 137 119 121 117 94 65 79 80 65 54 64 72 98]
[127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]
[115 114 109 123 150 148 131 118 113 109 100 92 74 65 72 78]
[ 89 93 90 97 108 147 131 118 113 114 113 109 106 95 77 80]
[ 63 77 86 81 77 79 102 123 117 115 117 125 125 130 115 87]
[ 62 65 82 89 78 71 80 101 124 126 119 101 107 114 131 119]
[ 63 65 75 88 89 71 62 81 120 138 135 105 81 98 110 118]
[ 87 65 71 87 106 95 69 45 76 130 126 107 92 94 105 112]
[118 97 82 86 117 123 116 66 41 51 95 93 89 95 102 107]
[164 146 112 80 82 120 124 104 76 48 45 66 88 101 102 109]
[157 170 157 120 93 86 114 132 112 97 69 55 70 82 99 94]
[130 128 134 161 139 100 109 118 121 134 114 87 65 53 69 86]
[128 112 96 117 150 144 120 115 104 107 102 93 87 81 72 79]
[123 107 96 86 83 112 153 149 122 109 104 75 80 107 112 99]
[122 121 102 80 82 86 94 117 145 148 153 102 58 78 92 107]
[122 164 148 103 71 56 78 83 93 103 119 139 102 61 69 84]

What the computer sees

An image is just a big grid of numbers between [0, 255]:

e.g. 800 x 600 x 3  
(3 channels RGB)

# Challenges: Viewpoint variation



All pixels change when  
the camera moves!

# Challenges: Illumination

---



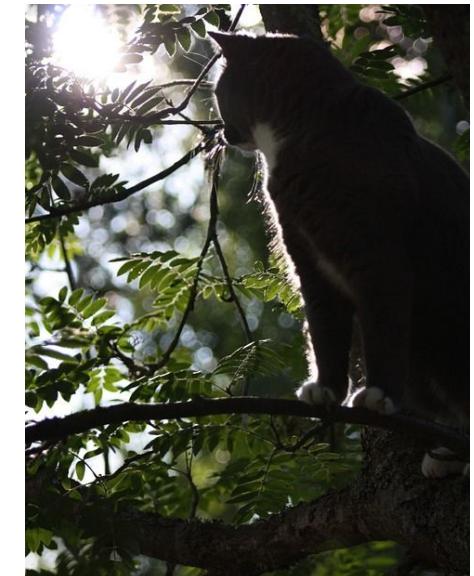
[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)

# Challenges: Deformation

---



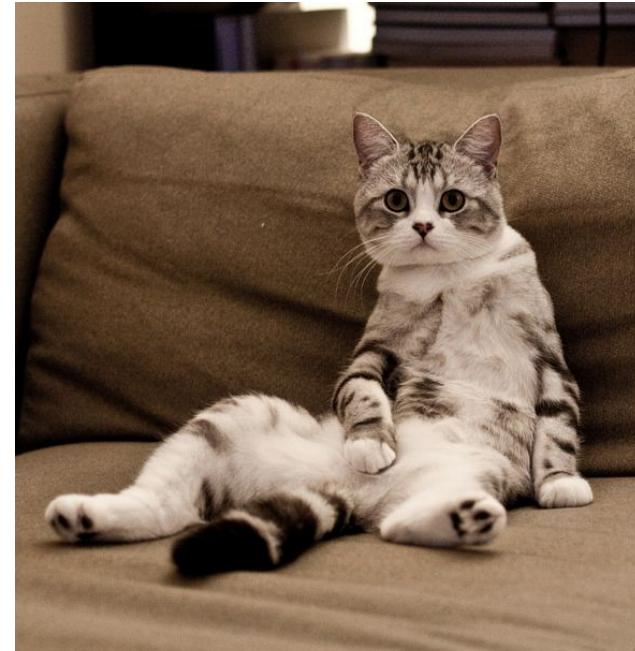
[This image by Umberto Salvagnin](#)  
is licensed under [CC-BY 2.0](#)



[This image by Umberto Salvagnin](#)  
is licensed under [CC-BY 2.0](#)



[This image by sare bear](#) is  
licensed under [CC-BY 2.0](#)



[This image by Tom Thai](#) is  
licensed under [CC-BY 2.0](#)

# Challenges: Occlusion

---



This image is [CC0 1.0 public domain](#)



This image is [CC0 1.0 public domain](#)



This image by [jonsson](#) is licensed  
under [CC-BY 2.0](#)

# Challenges: Background clutter

---



This image is [CC0 1.0 public domain](#)



This image is [CC0 1.0 public domain](#)

# Challenges: Intraclass variation

---



This image is CC0 1.0 public domain



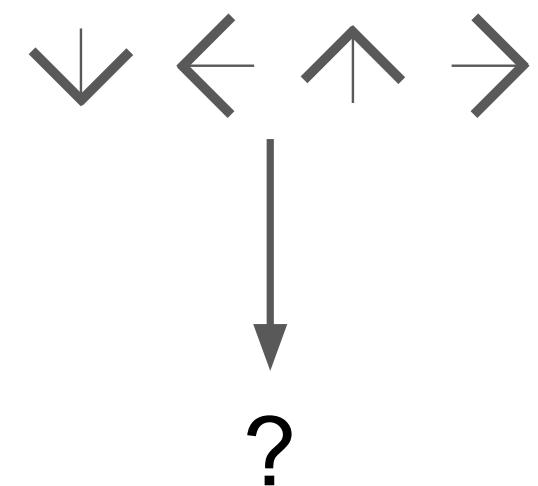
# How do we do this?



Find edges



Find corners



# Machine Learning

## Data-Driven Approach

1. Collect a dataset of images and labels
2. Use Machine Learning to train a classifier
3. Evaluate the classifier on new images

Example training set

```
def train(images, labels):
    # Machine learning!
    return model
```

```
def predict(model, test_images):
    # Use model to predict labels
    return test_labels
```

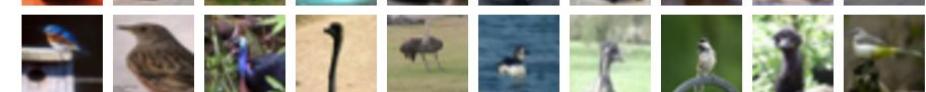
**airplane**



**automobile**



**bird**



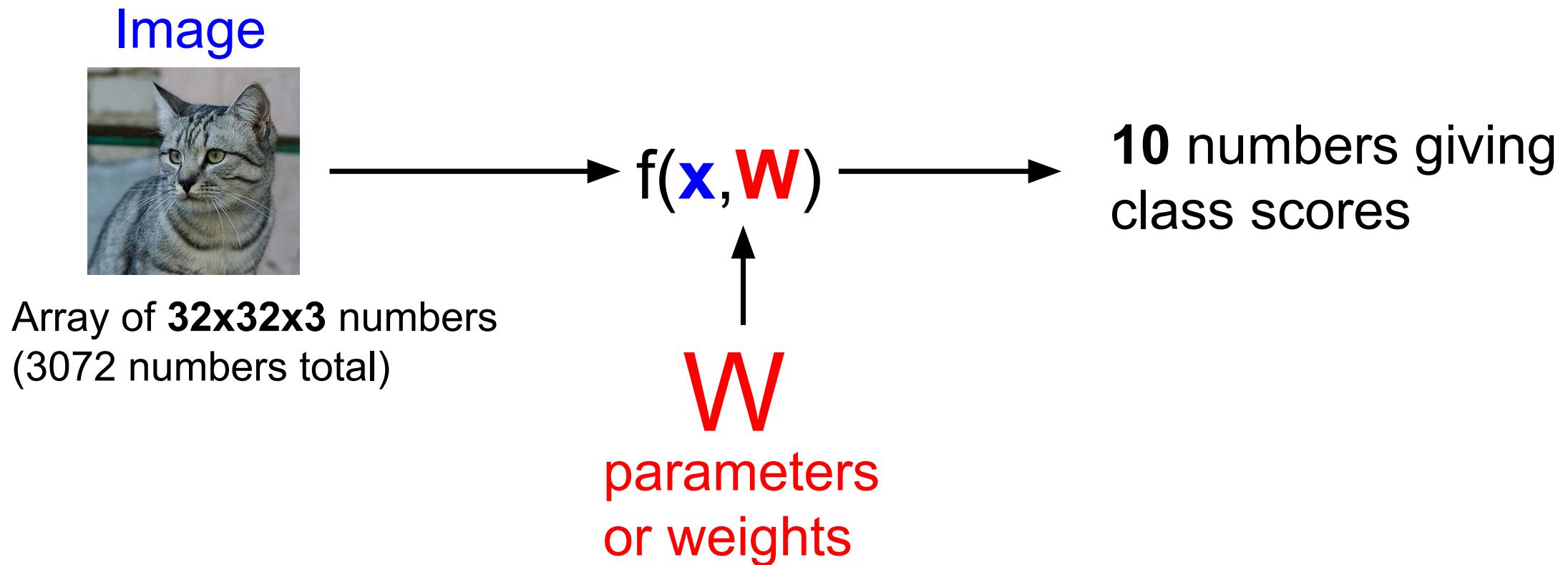
**cat**



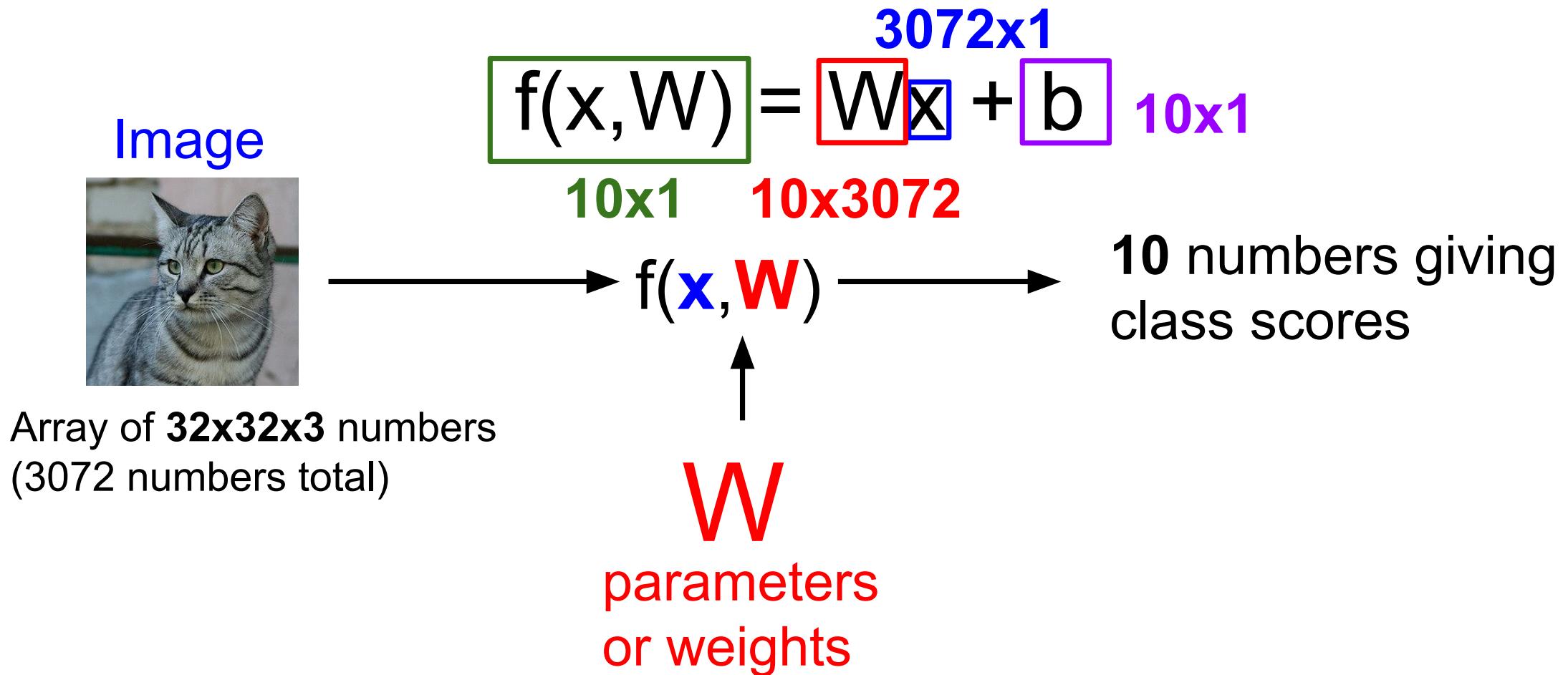
**deer**



# Linear Classifier

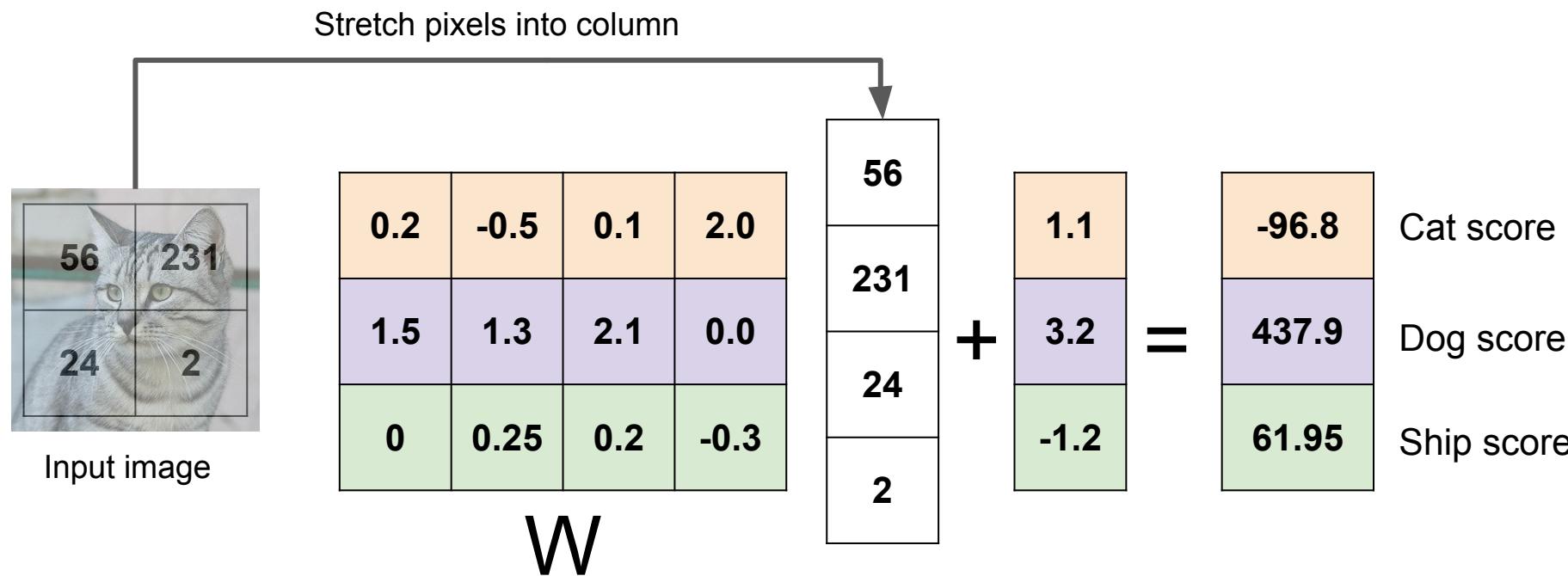


# Linear Classifier



# Linear Classifier

Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**)

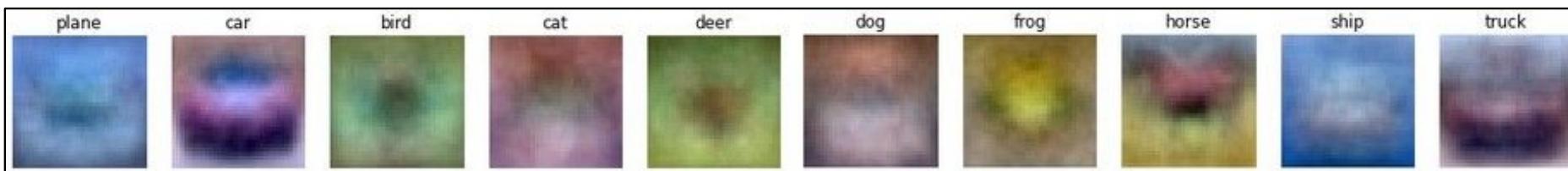


# Linear Classifier

---

Linear score function:

$$f = Wx$$



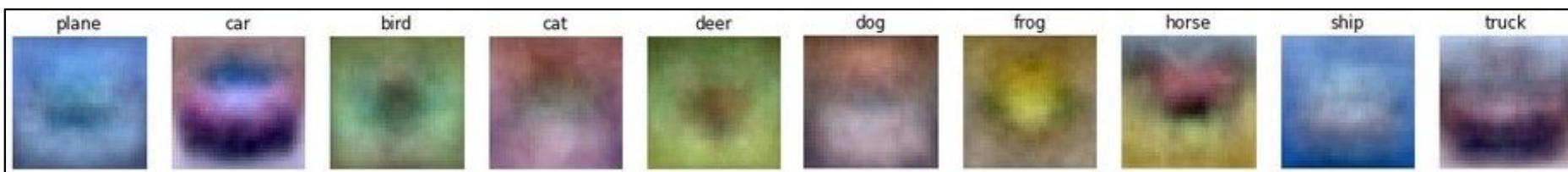
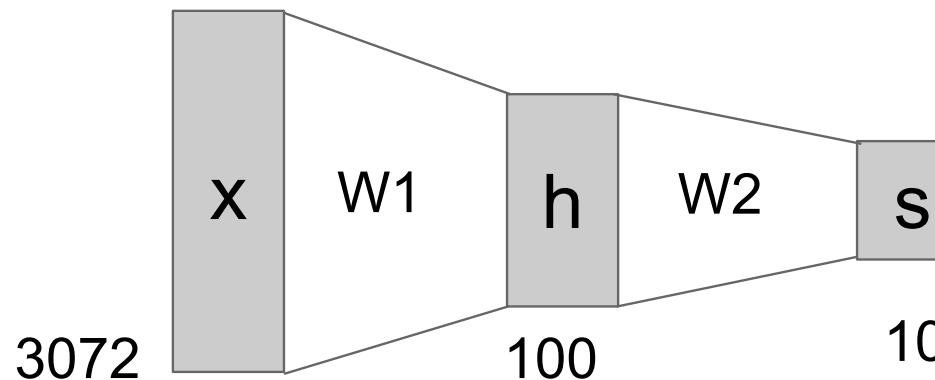
# Neural Networks

Linear score function:

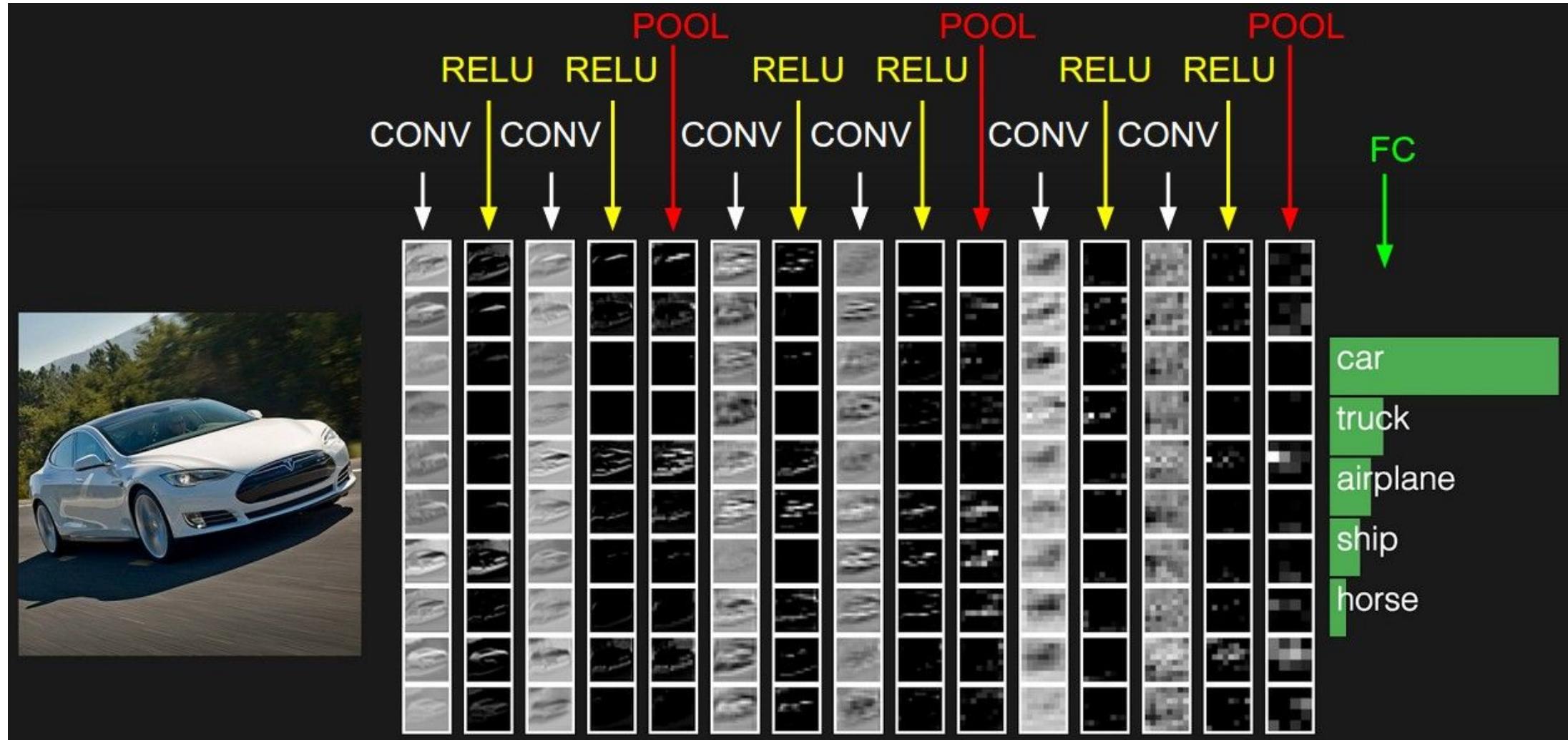
$$f = Wx$$

2-layer Neural Network

$$f = W_2 \max(0, W_1 x)$$

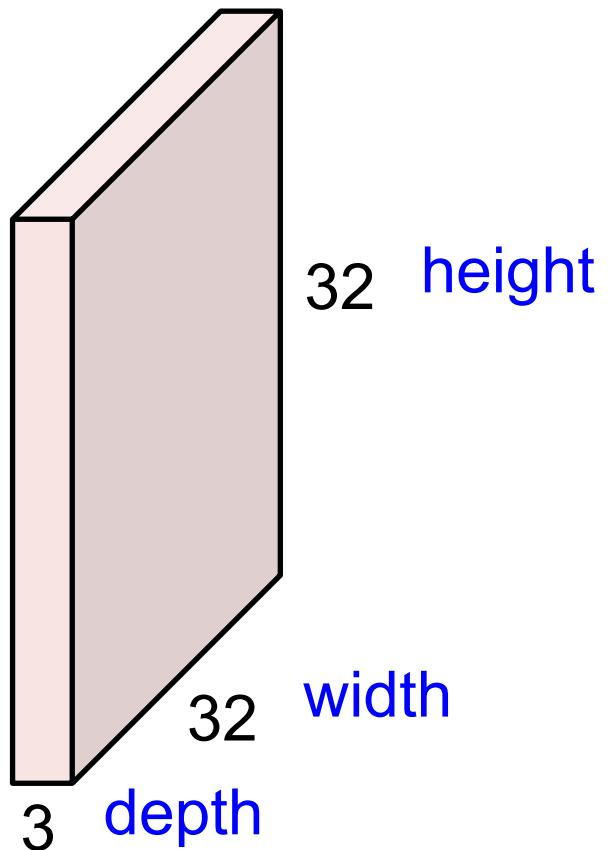


# Convolutional Neural Networks (CNN)



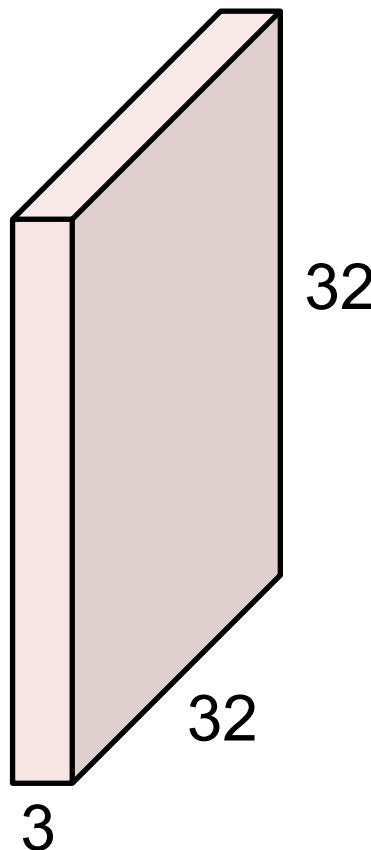
# Convolution Layer

32x32x3 image

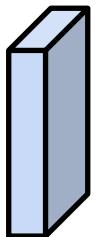


# Convolution Layer

32x32x3 image

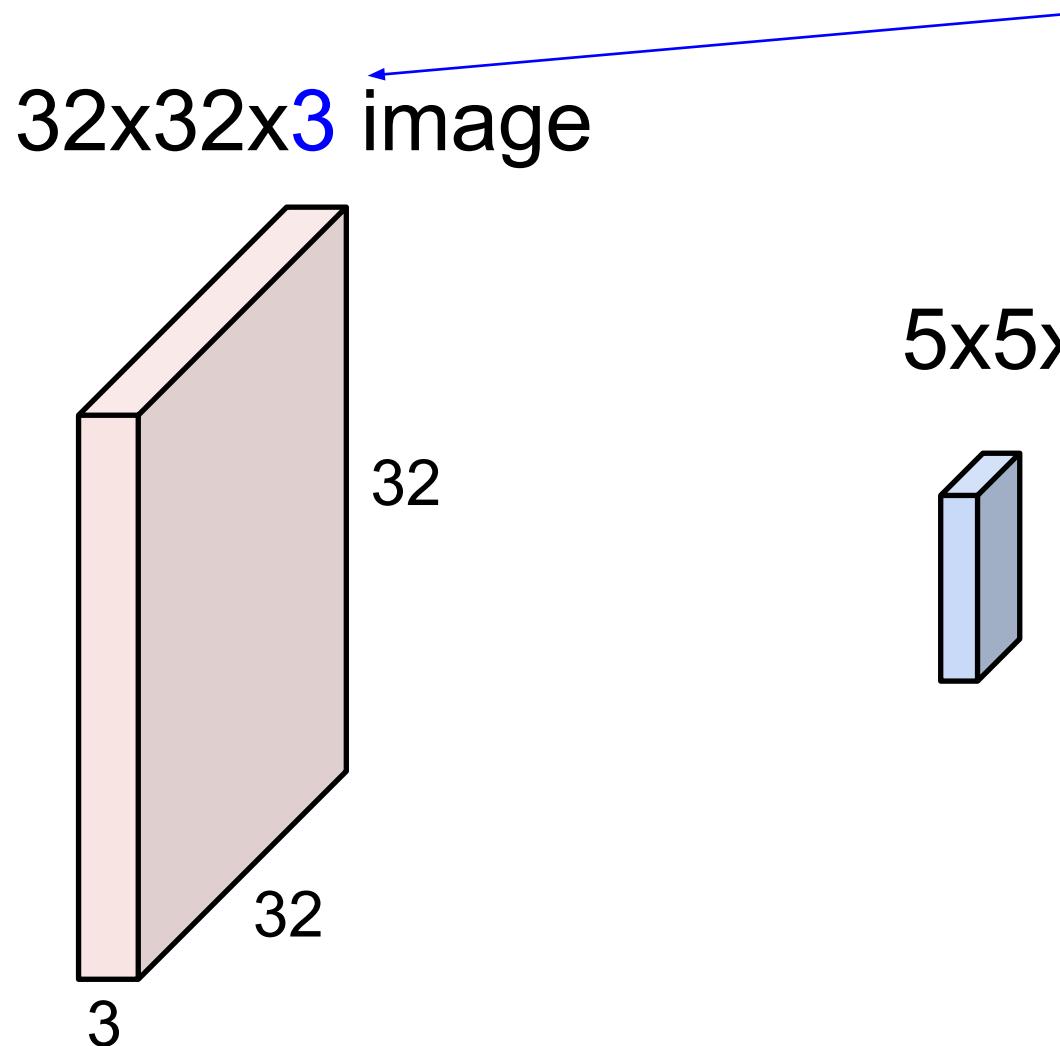


5x5x3 filter



**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolution Layer

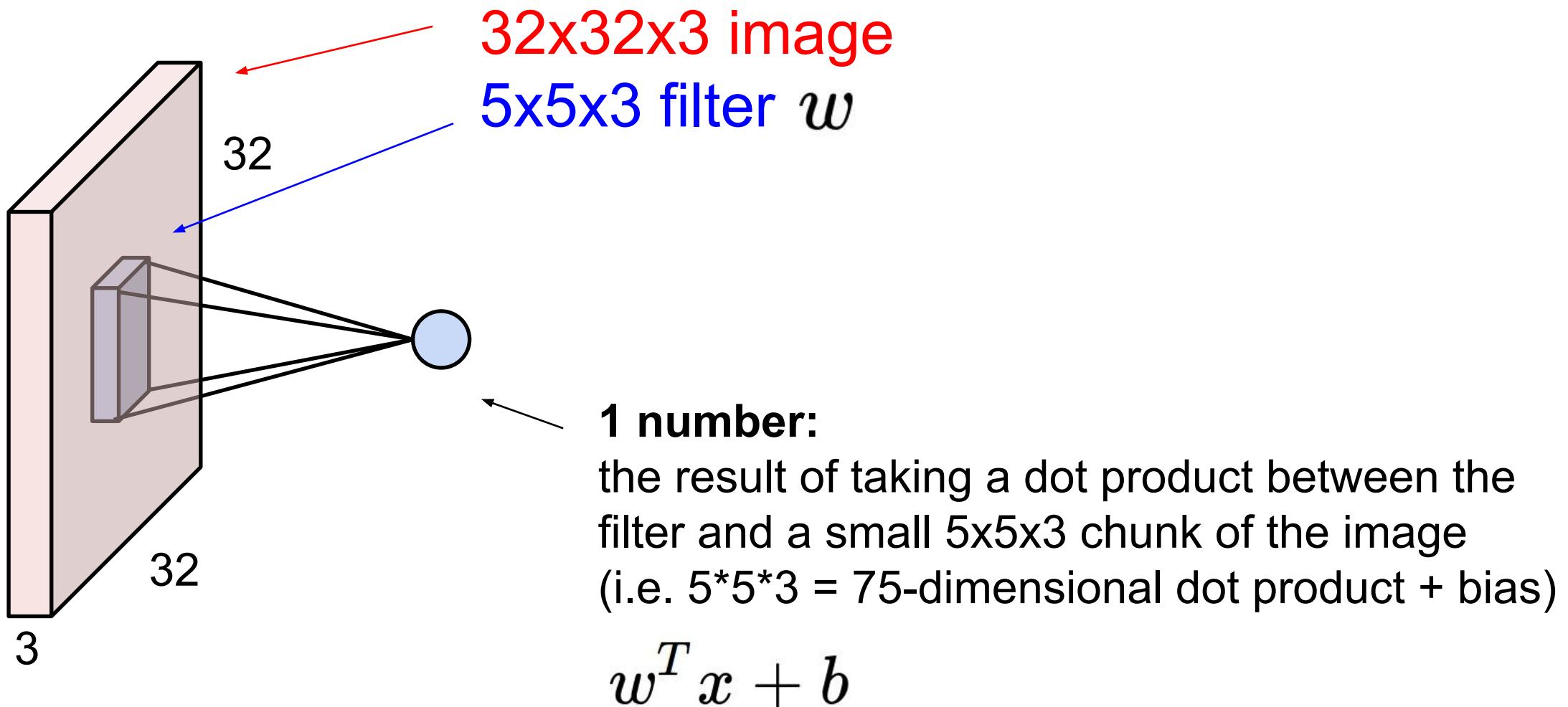


Filters always extend the full depth of the input volume

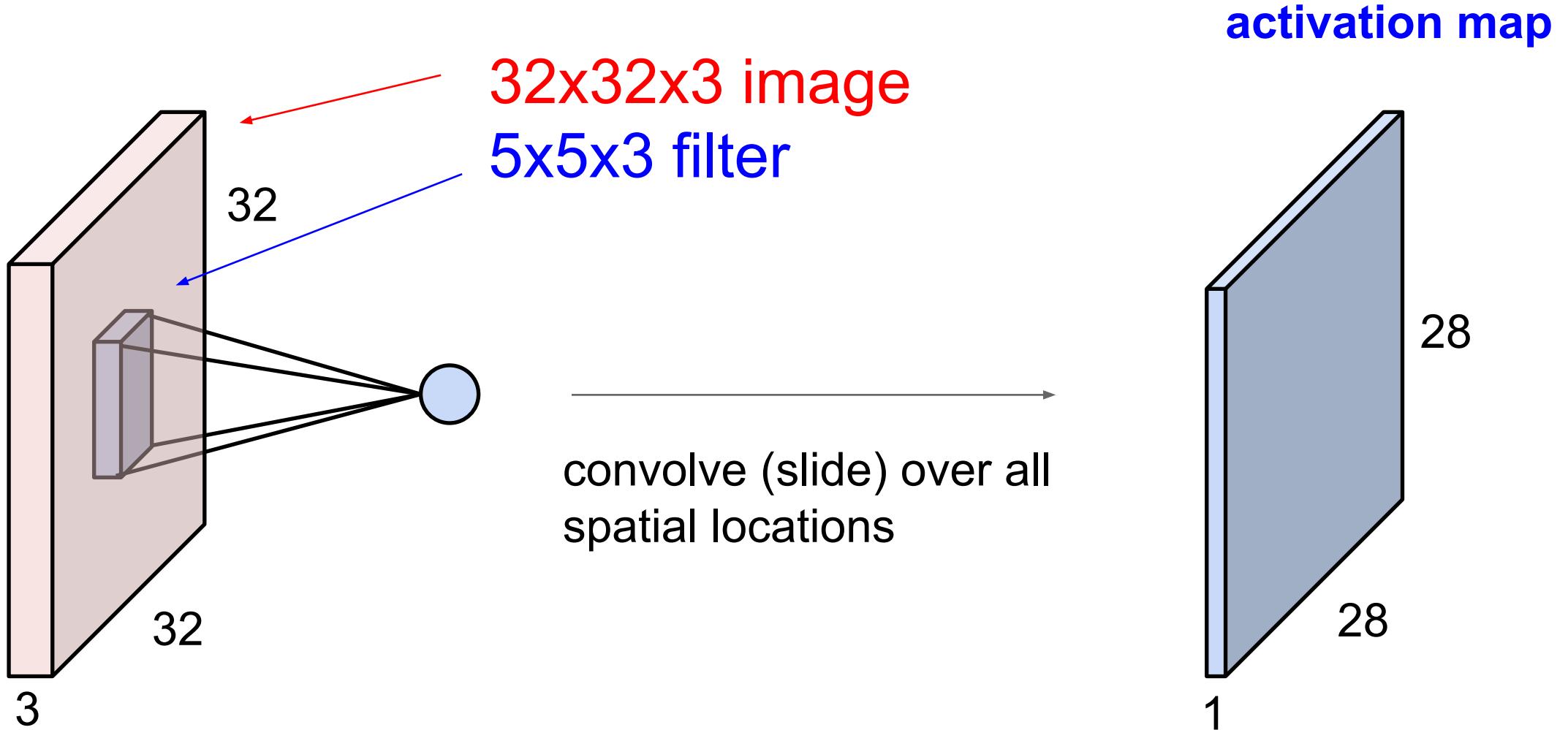
5x5x3 filter

**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolution Layer

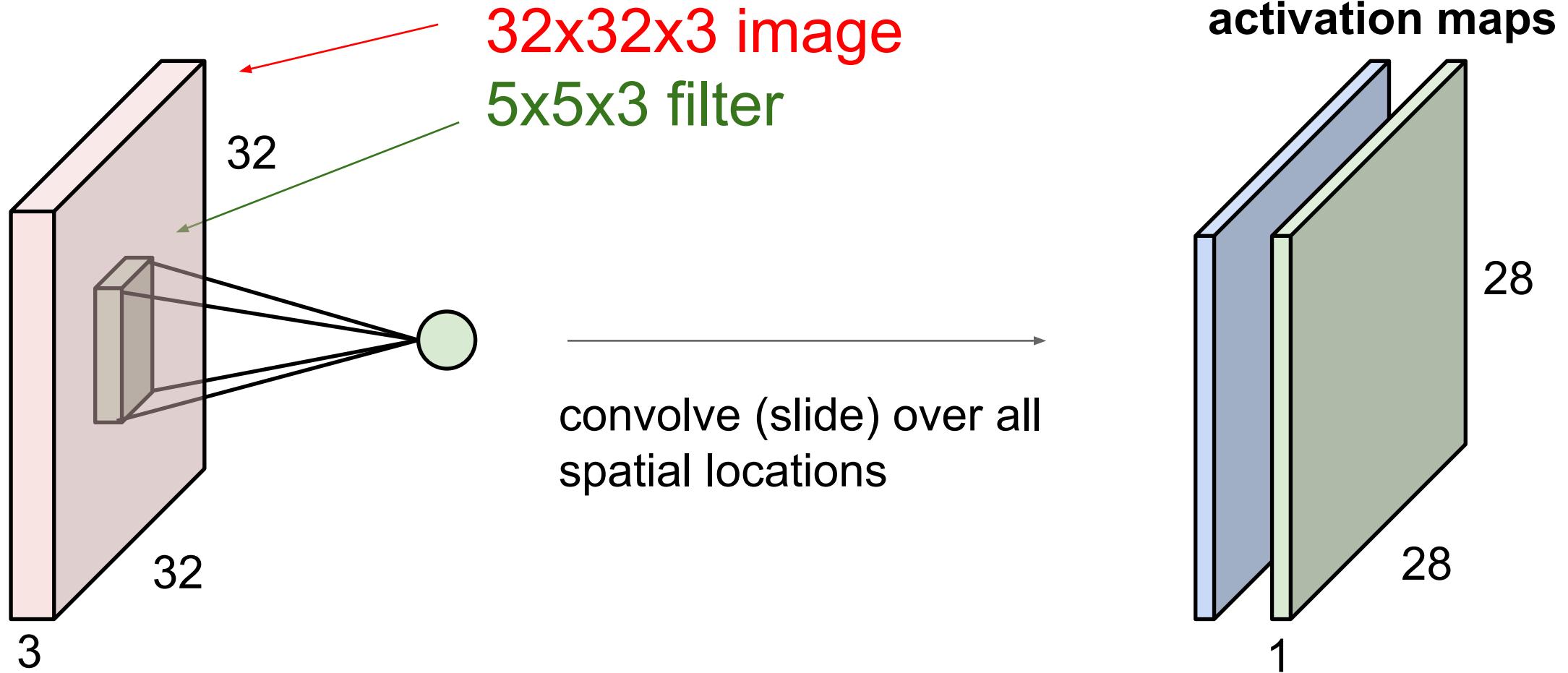


# Convolution Layer

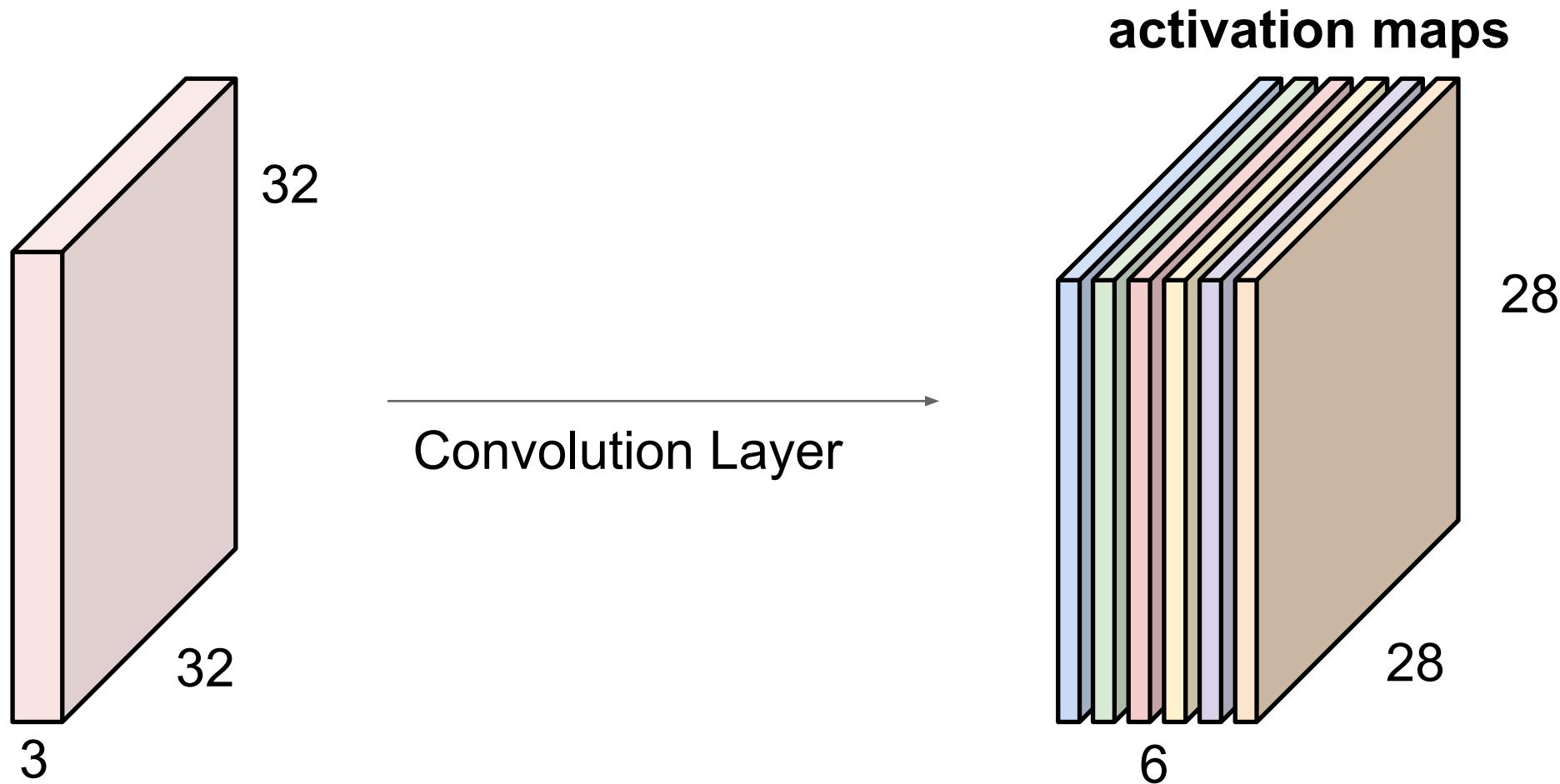


# Convolution Layer

consider a second, green filter

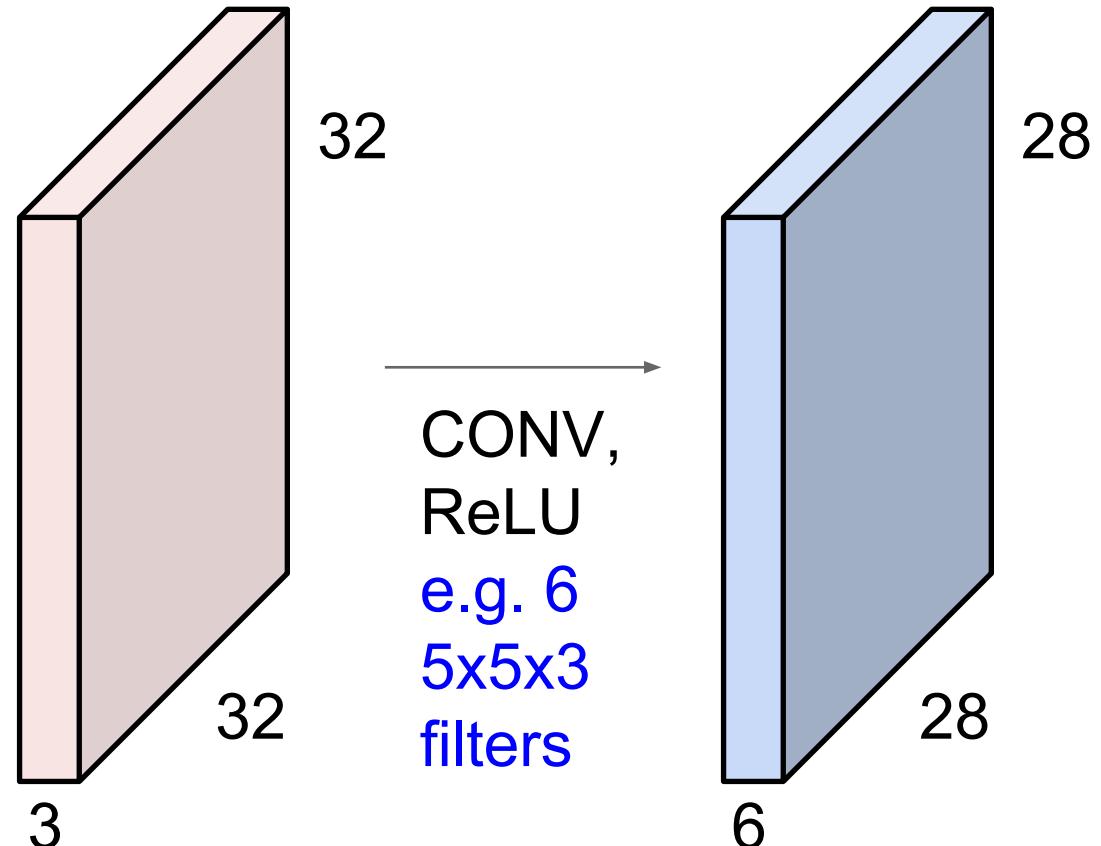


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

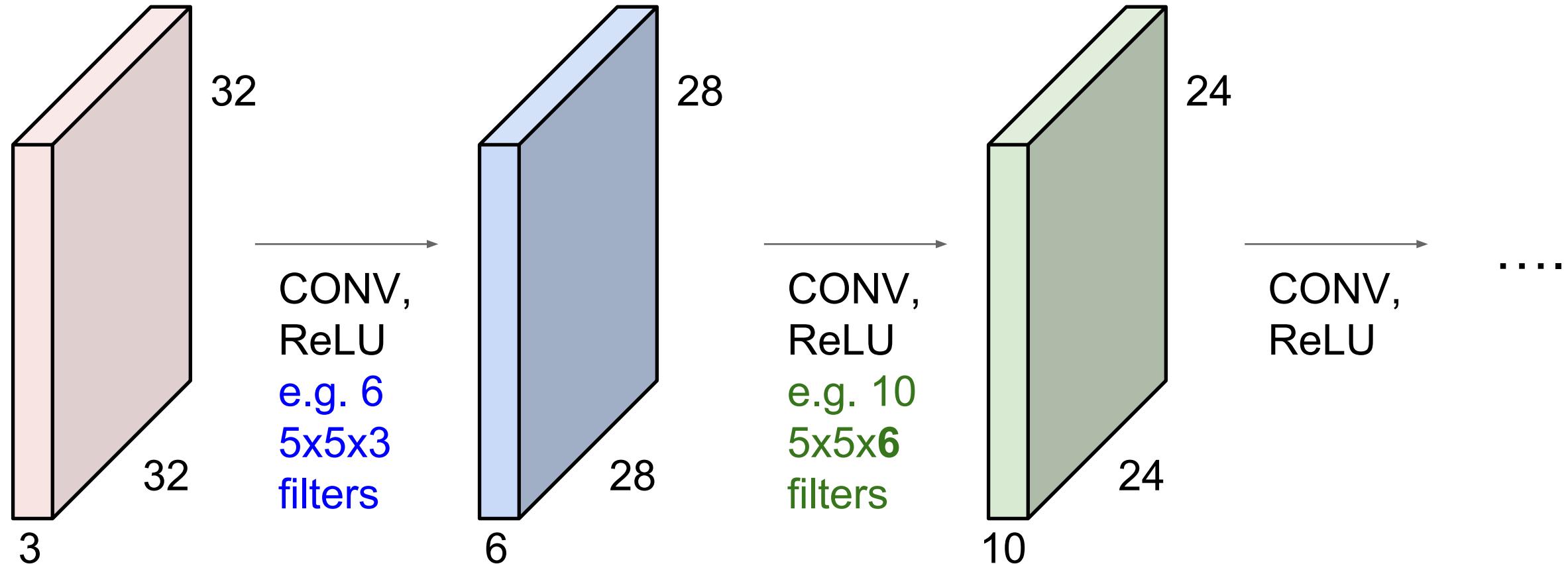


We stack these up to get a “new image” of size  $28 \times 28 \times 6$ !

**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions

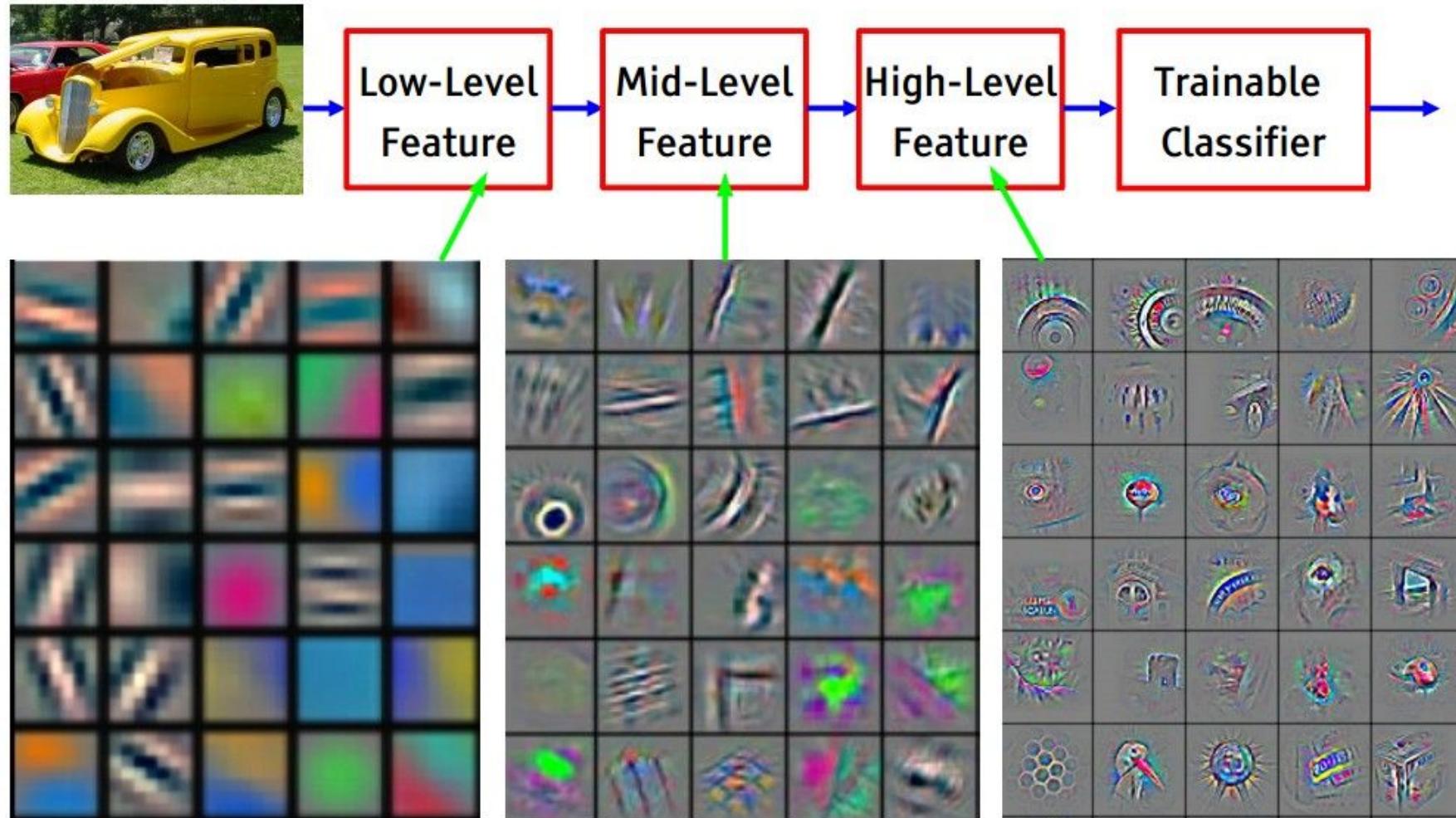


**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



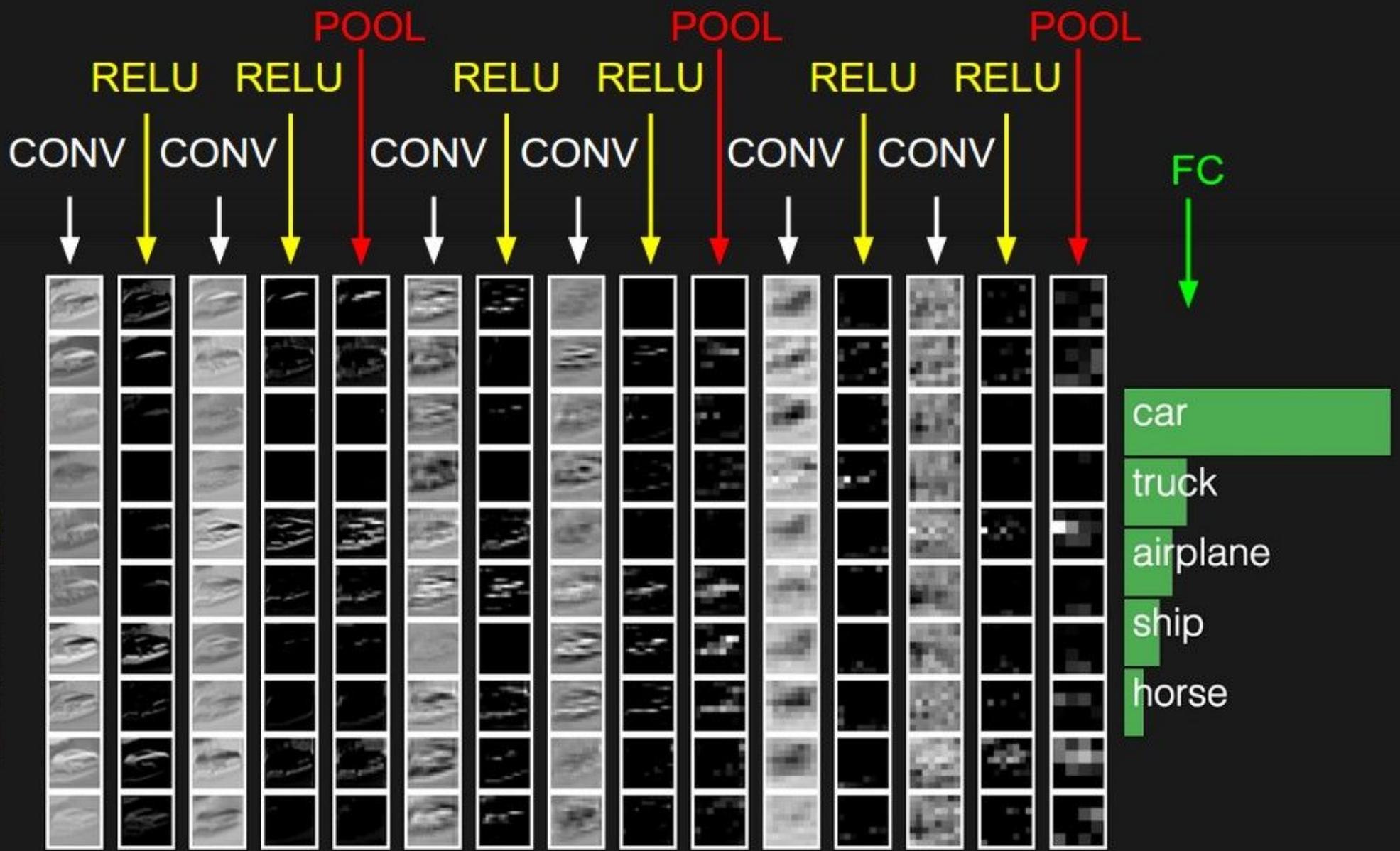
# Preview

[From recent Yann LeCun slides]



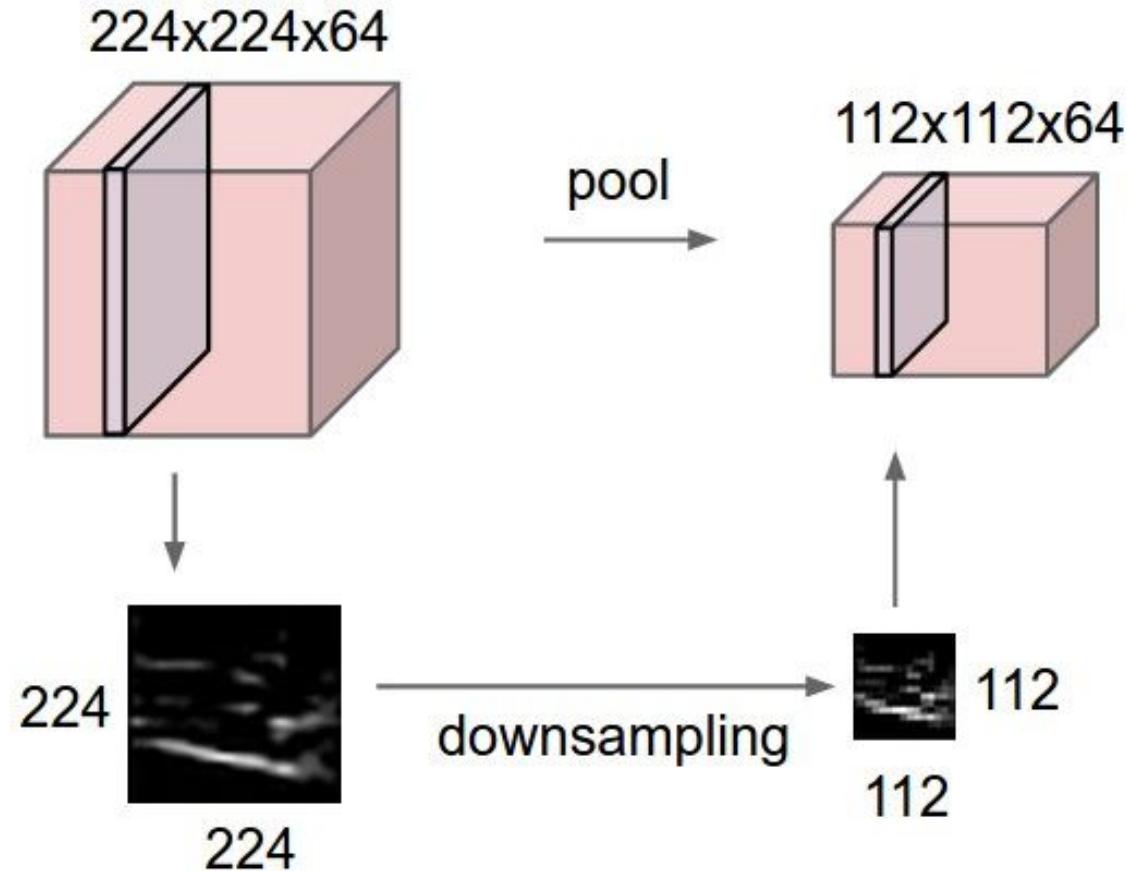
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

two more layers to go: POOL/FC



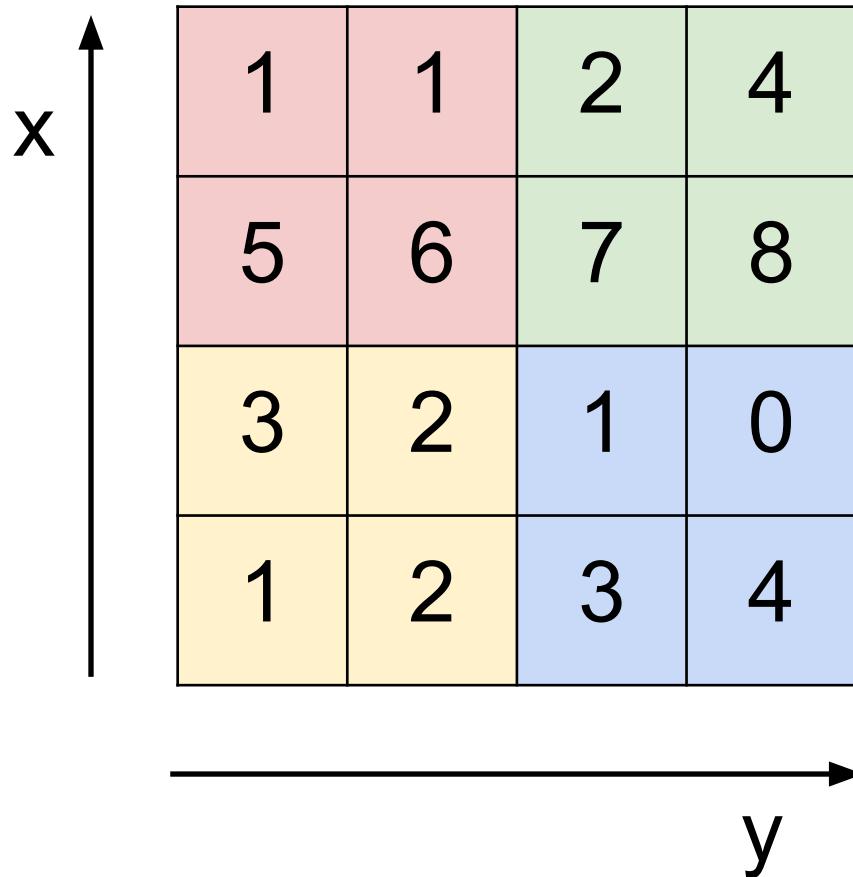
# Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:

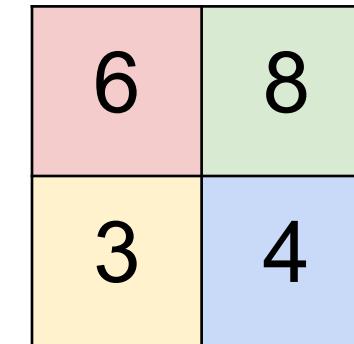


# MAX POOLING

Single depth slice

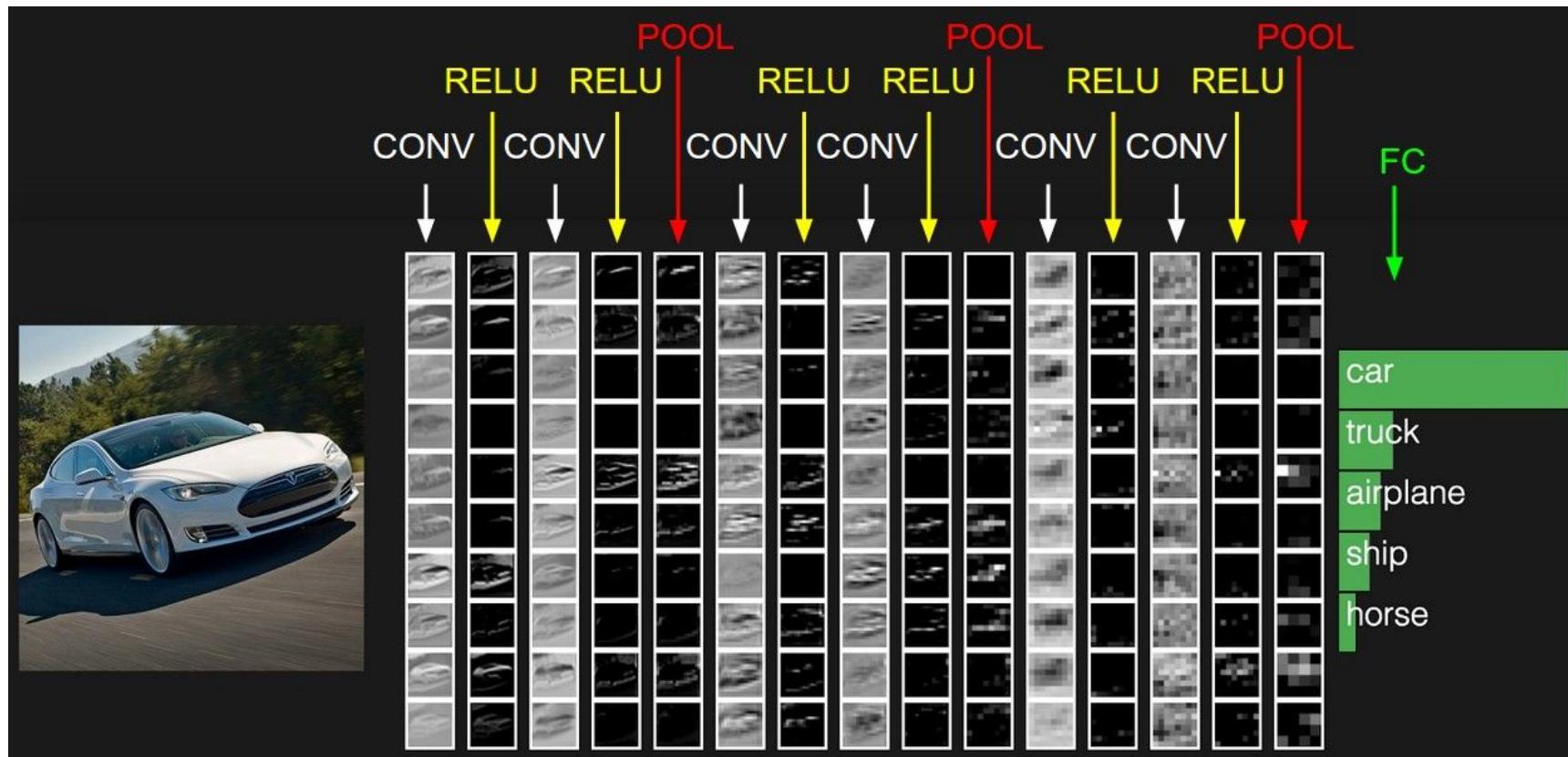


max pool with 2x2 filters  
and stride 2



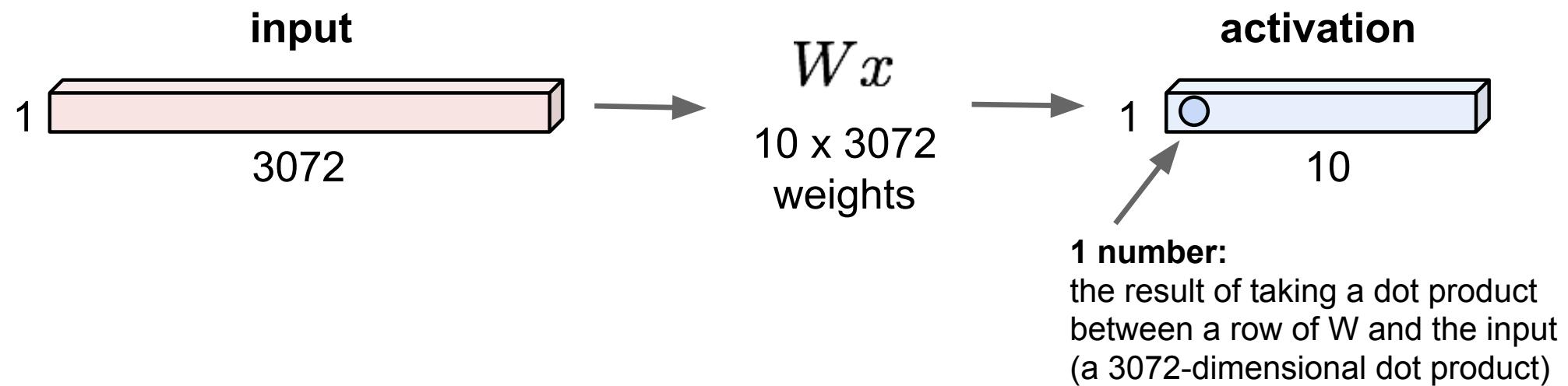
# Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



# Fully Connected Layer

32x32x3 image -> stretch to  $3072 \times 1$



# Summary of ConvNets

---

- ConvNets stack CONV,POOL,FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like

**$[(\text{CONV-RELU})^*N - \text{POOL?}]^*M - (\text{FC-RELU})^*K, \text{SOFTMAX}$**

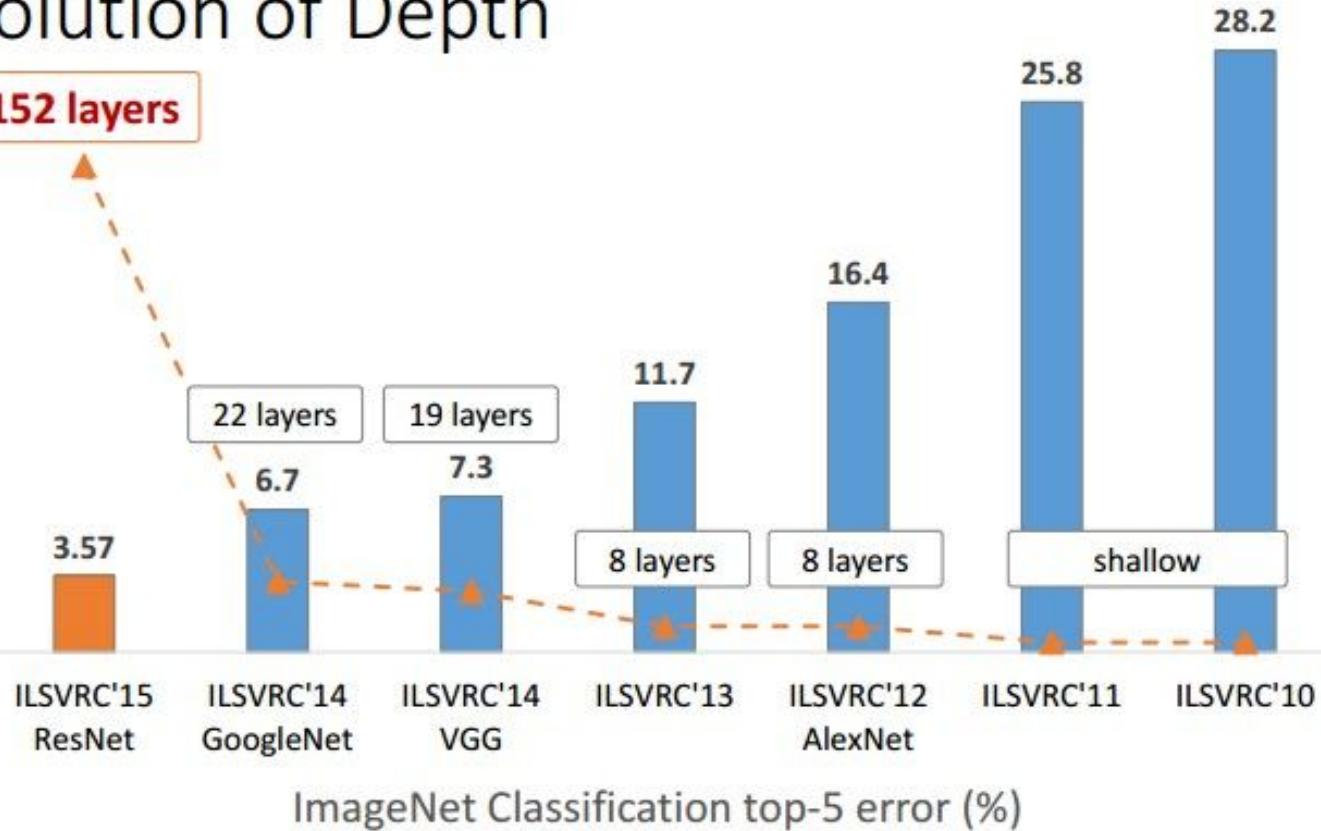
where N is usually up to ~5, M is large,  $0 \leq K \leq 2$ .

- but recent advances such as ResNet/GoogLeNet challenge this paradigm

# Deep Learning

Microsoft  
Research

## Revolution of Depth



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

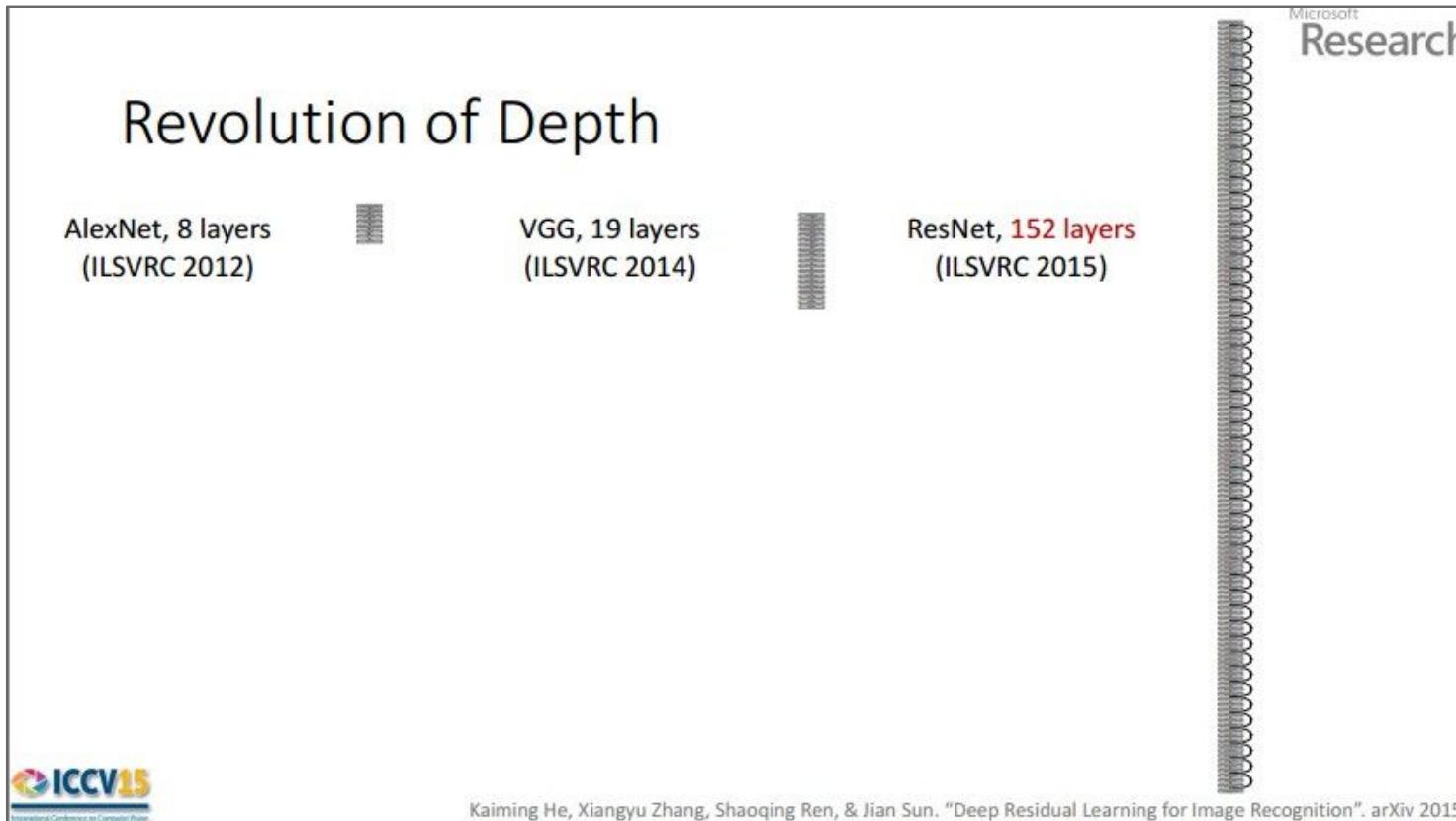
(slide from Kaiming He's recent presentation)

# Deep Learning

## Case Study: ResNet

[He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)

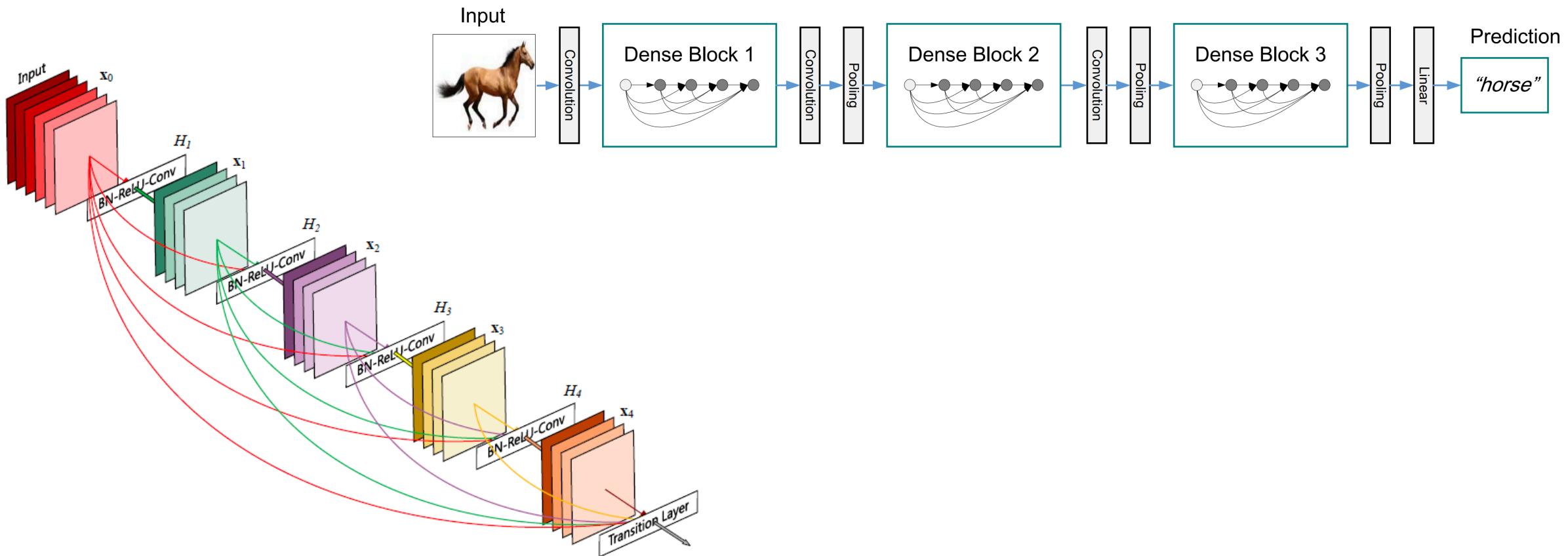


2-3 weeks of training  
on 8 GPU machine

at runtime: faster  
than a VGGNet!  
(even though it has  
8x more layers)

(slide from Kaiming He's recent presentation)

# Densely Connected Convolutional Networks



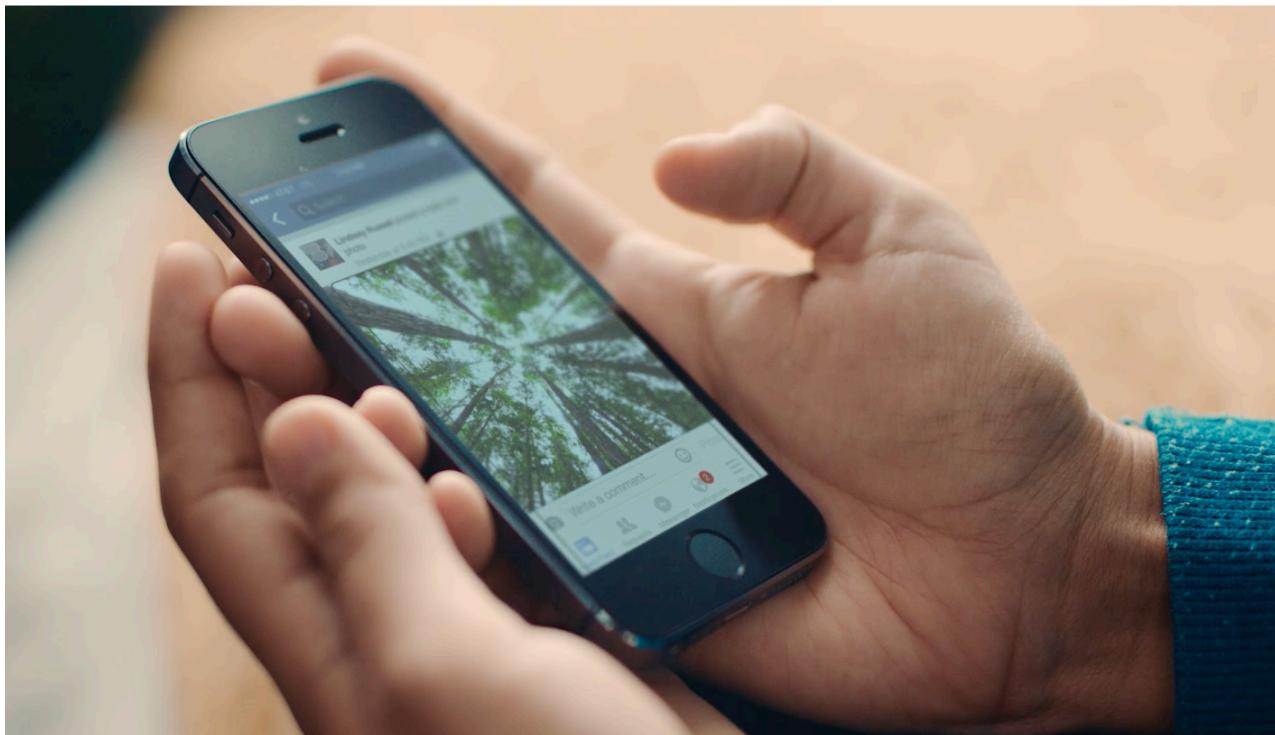
Figures copyright Gao Huang, Zhuang Liu, Laurens van Der Maaten, Kilian Q. Weinberger, 2017.

[Huang et al. 2017 CVPR]

# Facebook's Blind Users

April 4, 2016

## Using Artificial Intelligence to Help Blind People 'See' Facebook



By [Shaomei Wu](#), Software Engineer and [Hermes Pique](#), Software Engineer on iOS and [Jeffrey Wieland](#), Head of Accessibility

<https://www.youtube.com/watch?v=5fZLch5DjZc>

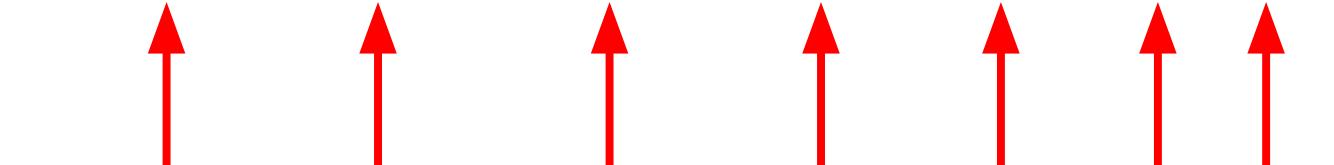
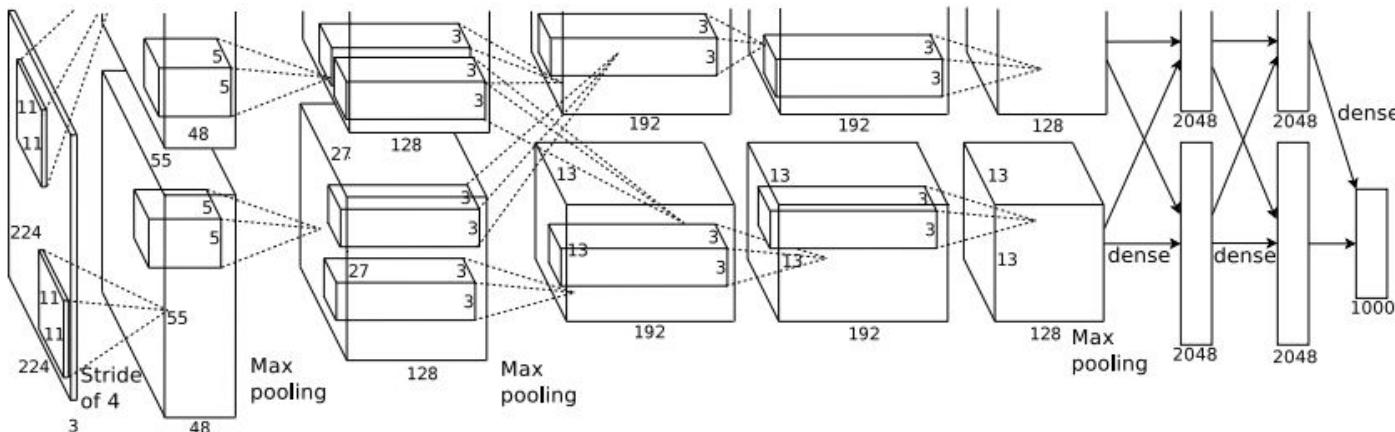
<https://newsroom.fb.com/news/2016/04/using-artificial-intelligence-to-help-blind-people-see-facebook/>

# What's going on inside ConvNets?

This image is CC0 public domain



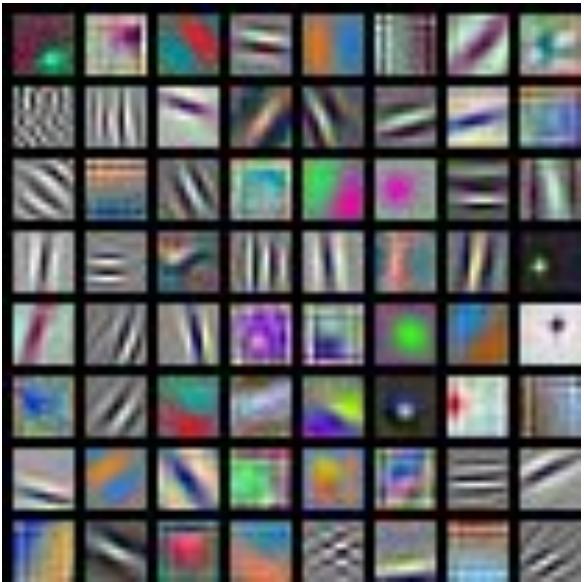
Input Image:  
 $3 \times 224 \times 224$



What are the intermediate features looking for?

Class Scores:  
1000 numbers

# First Layer: Visualize Filters



# AlexNet: 64 x 3 x 11 x 11



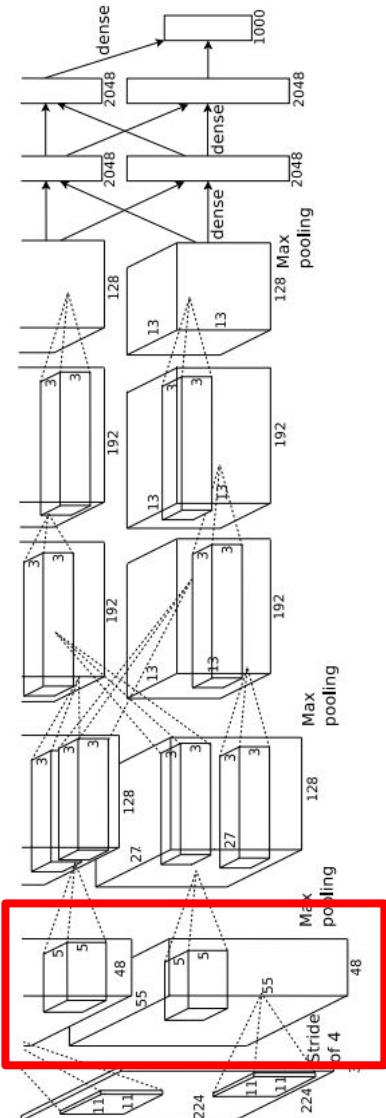
# ResNet-18: 64 x 3 x 7 x 7



# ResNet-101: 64 x 3 x 7 x 7



# DenseNet-121: 64 x 3 x 7 x 7



# Visualize the filters/kernels (raw weights)

We can visualize filters at higher layers, but not that interesting

(these are taken from ConvNetJS CIFAR-10 demo)



layer 1 weights

$16 \times 3 \times 7 \times 7$



layer 2 weights

$20 \times 16 \times 7 \times 7$

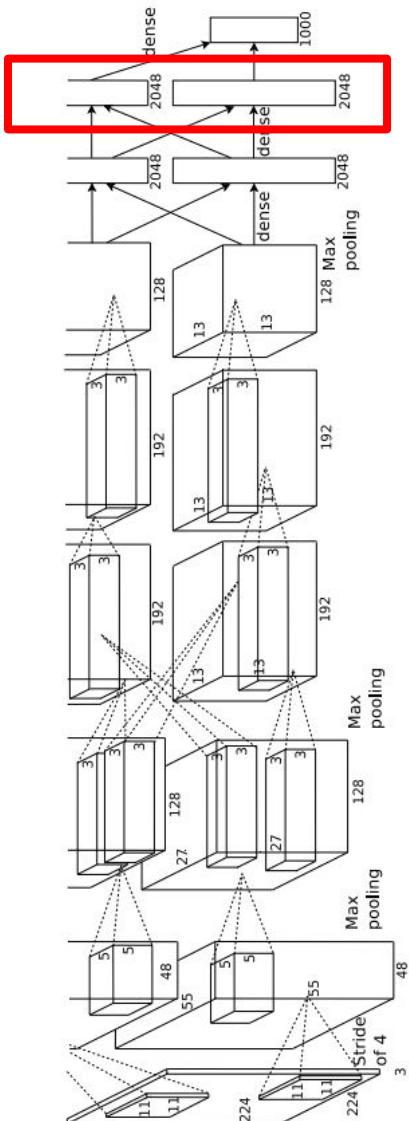


layer 3 weights

$20 \times 20 \times 7 \times 7$

# Last Layer

FC7 layer



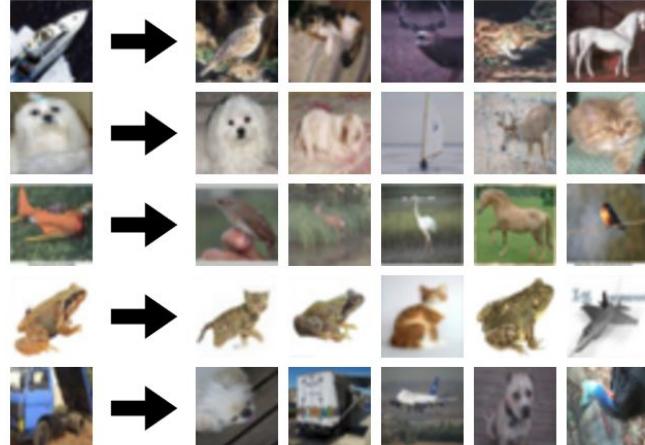
4096-dimensional feature vector for an image  
(layer immediately before the classifier)

Run the network on many images, collect the  
feature vectors

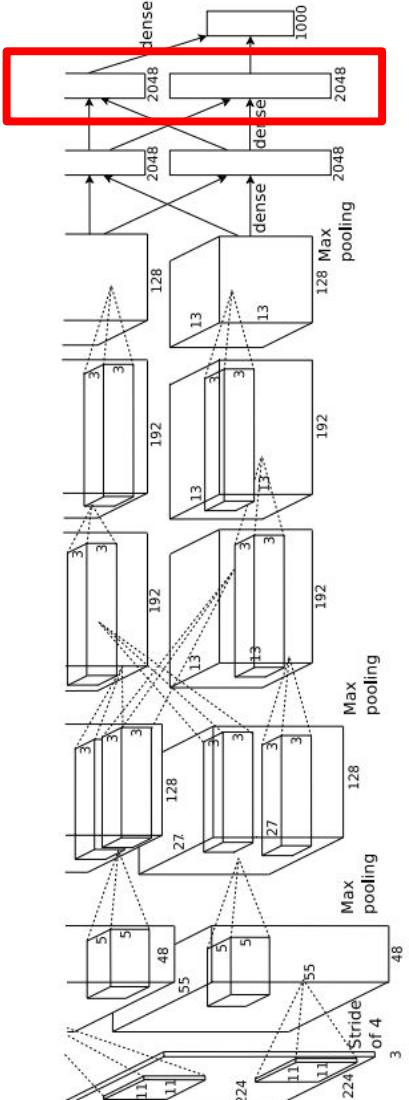
# Last Layer: Nearest Neighbors

4096-dim vector

Recall: Nearest neighbors  
in pixel space



Test image L2 Nearest neighbors in feature space



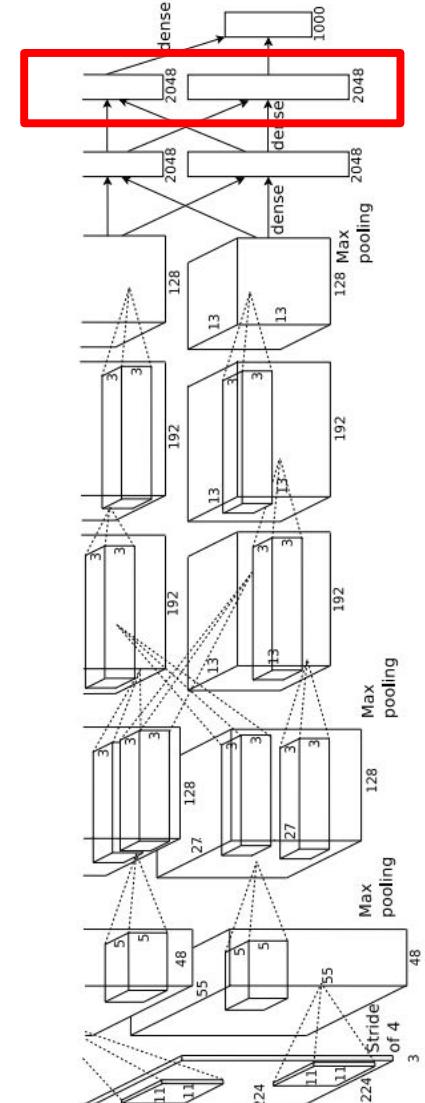
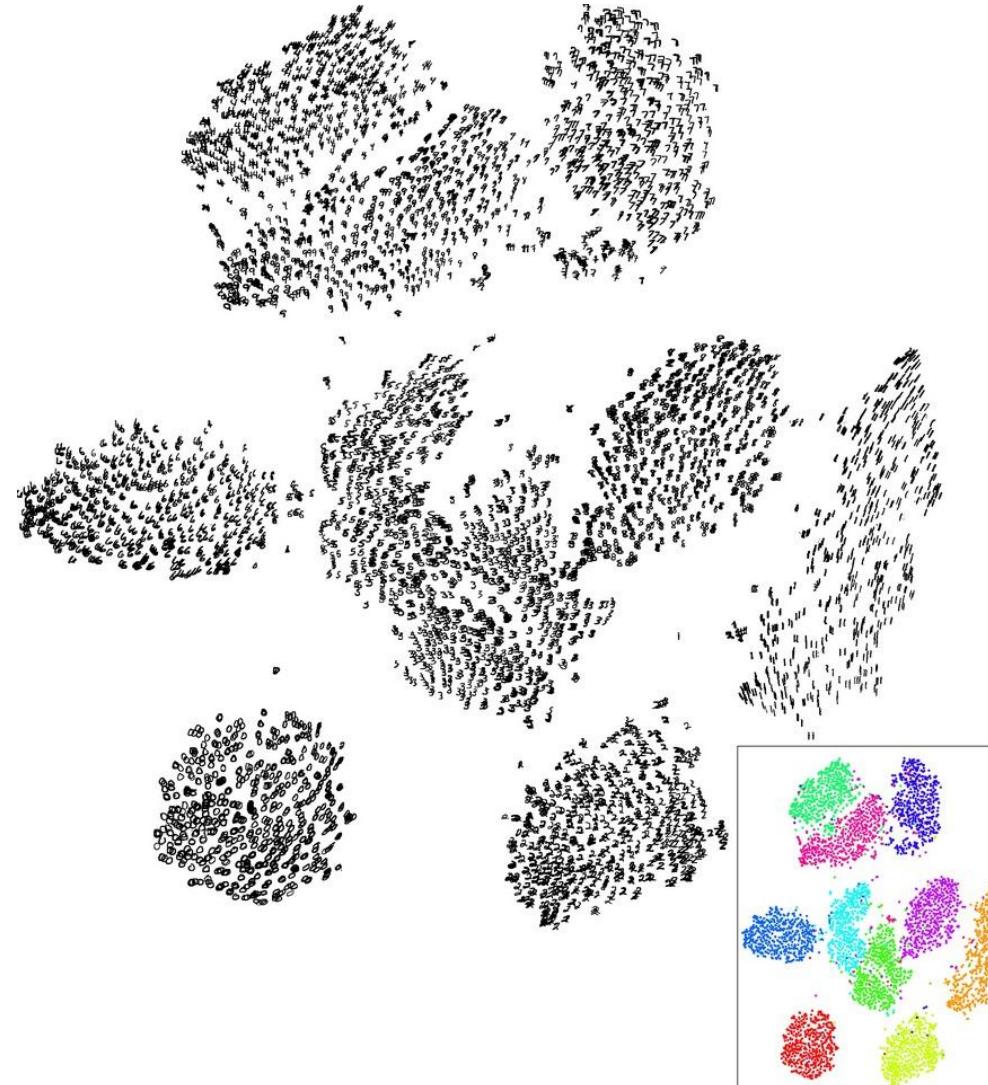
Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.  
Figures reproduced with permission.

# Last Layer: Dimensionality Reduction

Visualize the “space” of FC7 feature vectors by reducing dimensionality of vectors from 4096 to 2 dimensions

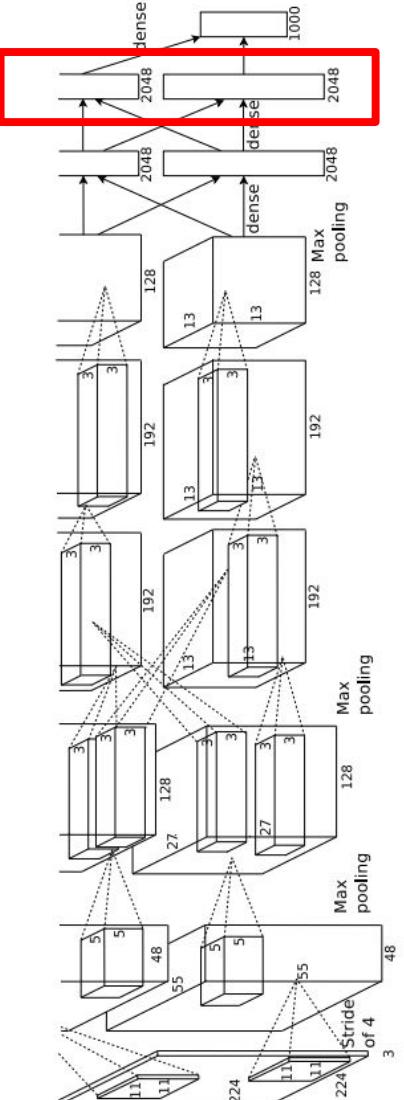
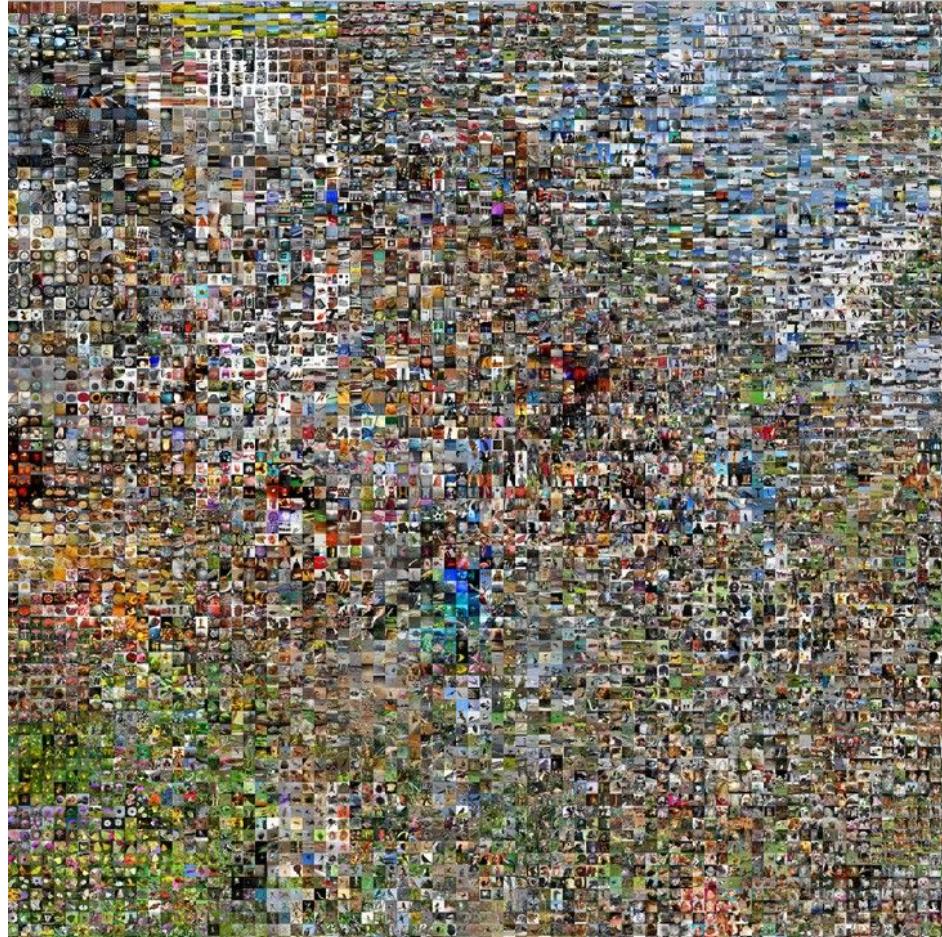
Simple algorithm: Principle Component Analysis (PCA)

More complex: t-SNE



Van der Maaten and Hinton, “Visualizing Data using t-SNE”, JMLR 2008  
Figure copyright Laurens van der Maaten and Geoff Hinton, 2008. Reproduced with permission.

# Last Layer: Dimensionality Reduction



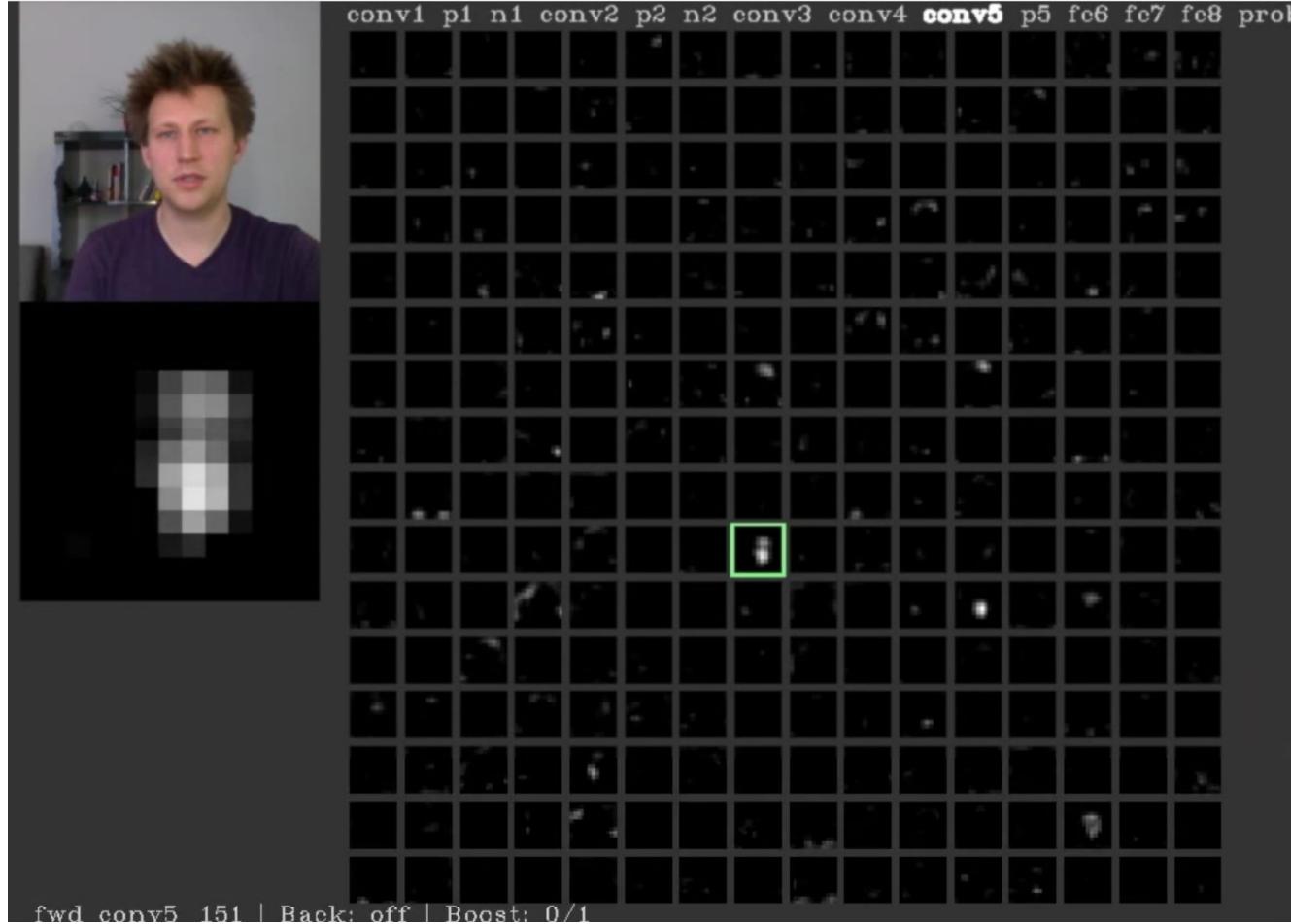
Van der Maaten and Hinton, "Visualizing Data using t-SNE", JMLR 2008  
Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.  
Figure reproduced with permission.

See high-resolution versions at  
<http://cs.stanford.edu/people/karpathy/cnnembed/>

<https://cs.stanford.edu/people/karpathy/cnnembed/>  
<https://projector.tensorflow.org/>

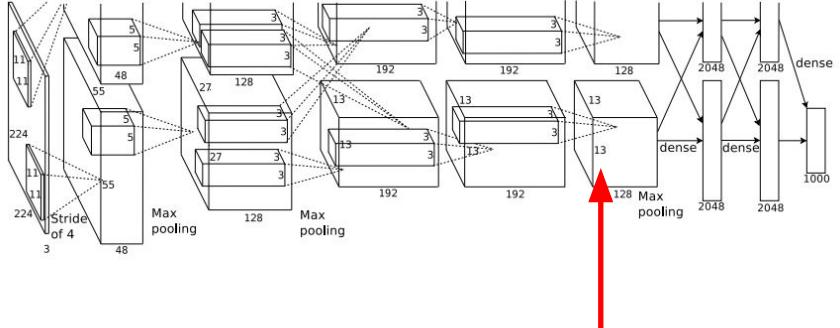
# Visualizing Activations

conv5 feature map is  
128x13x13; visualize  
as 128 13x13  
grayscale images



Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.  
Figure copyright Jason Yosinski, 2014. Reproduced with permission.

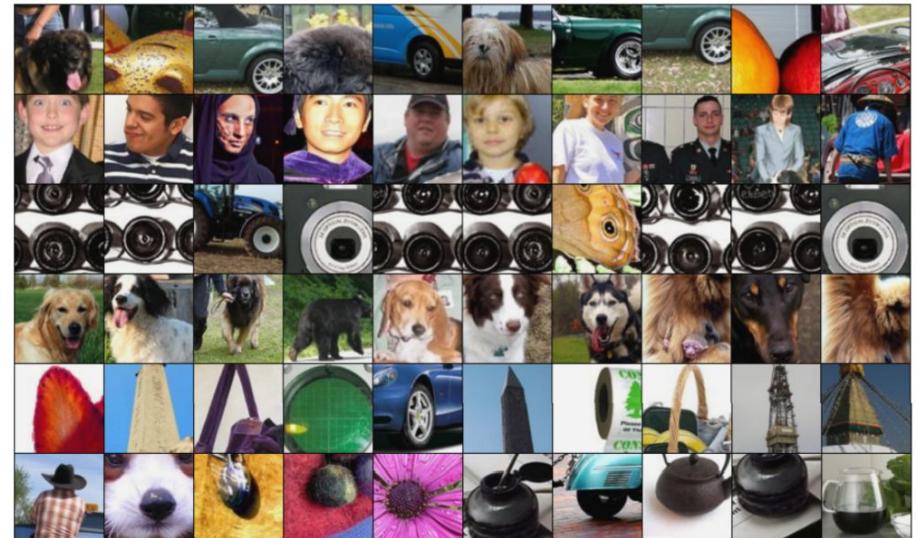
# Maximally Activating Patches



Pick a layer and a channel; e.g. conv5 is  $128 \times 13 \times 13$ , pick channel 17/128

Run many images through the network,  
record values of chosen channel

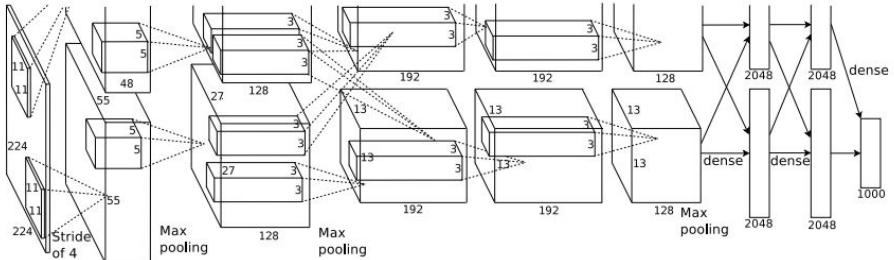
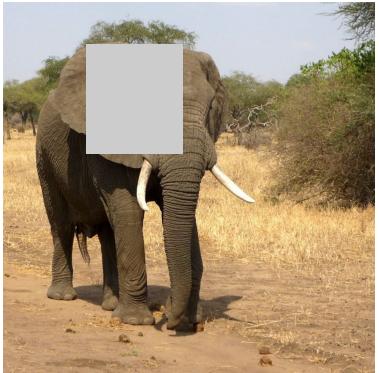
Visualize image patches that correspond  
to maximal activations



Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015  
Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015;  
reproduced with permission.

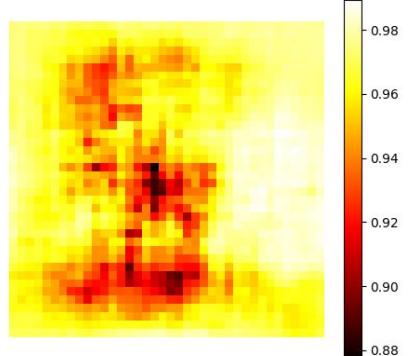
# Occlusion Experiments

Mask part of the image before feeding to CNN, draw heatmap of probability at each mask location

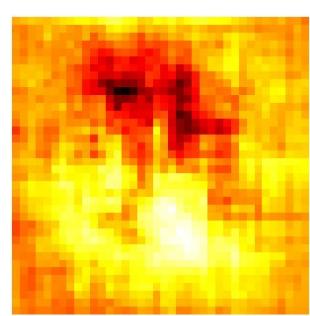


Boat image is [CC0 public domain](#)  
Elephant image is [CC0 public domain](#)  
Go-Karts image is [CC0 public domain](#)

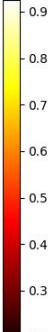
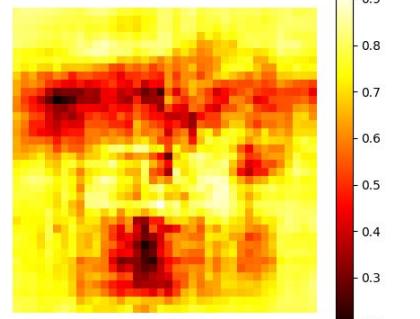
schooner



African elephant, *Loxodonta africana*



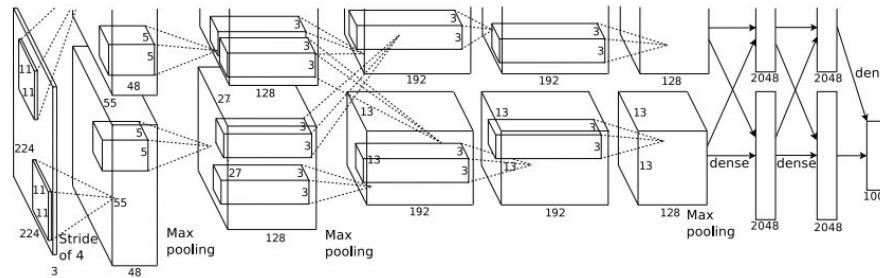
go-kart



Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

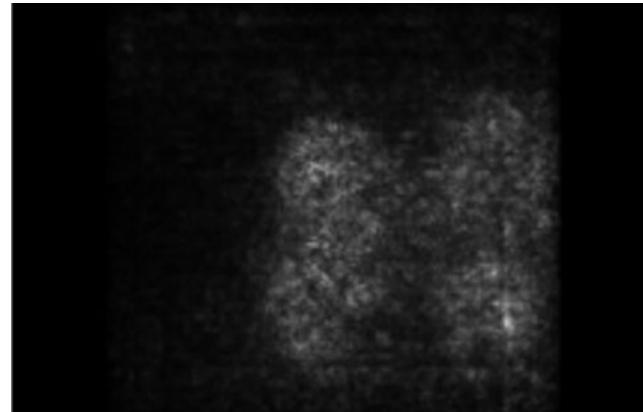
# Saliency Maps

How to tell which pixels matter for classification?



Dog

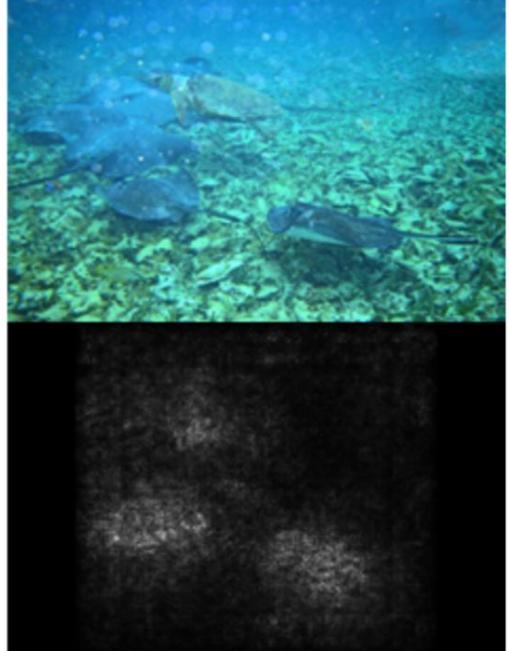
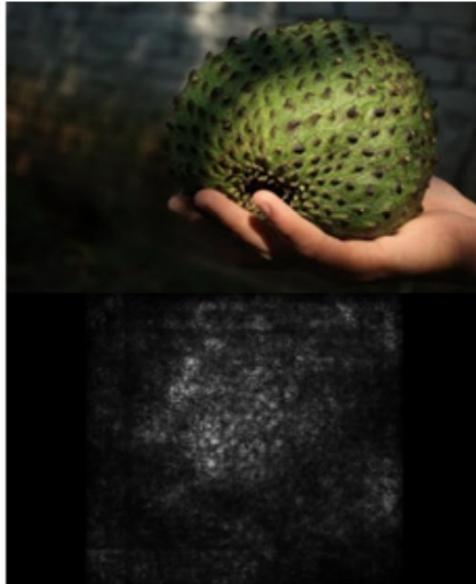
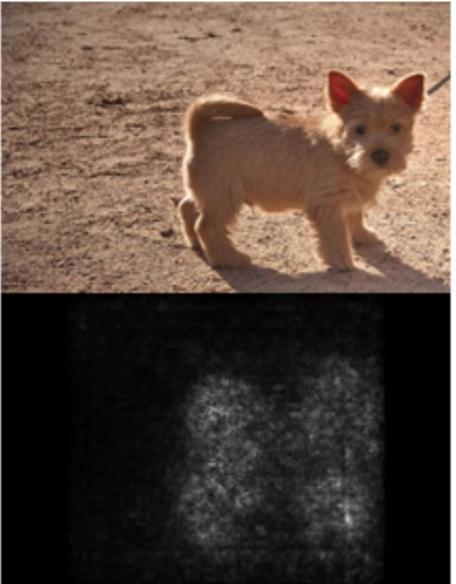
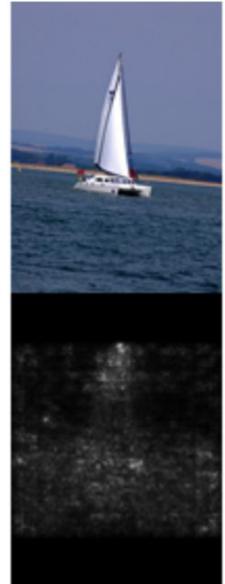
Compute gradient of (unnormalized) class score with respect to image pixels, take absolute value and max over RGB channels



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

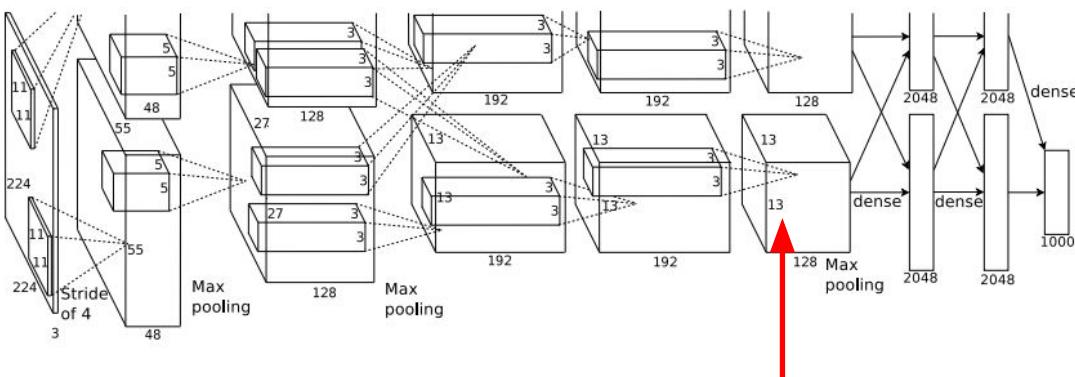
# Saliency Maps



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

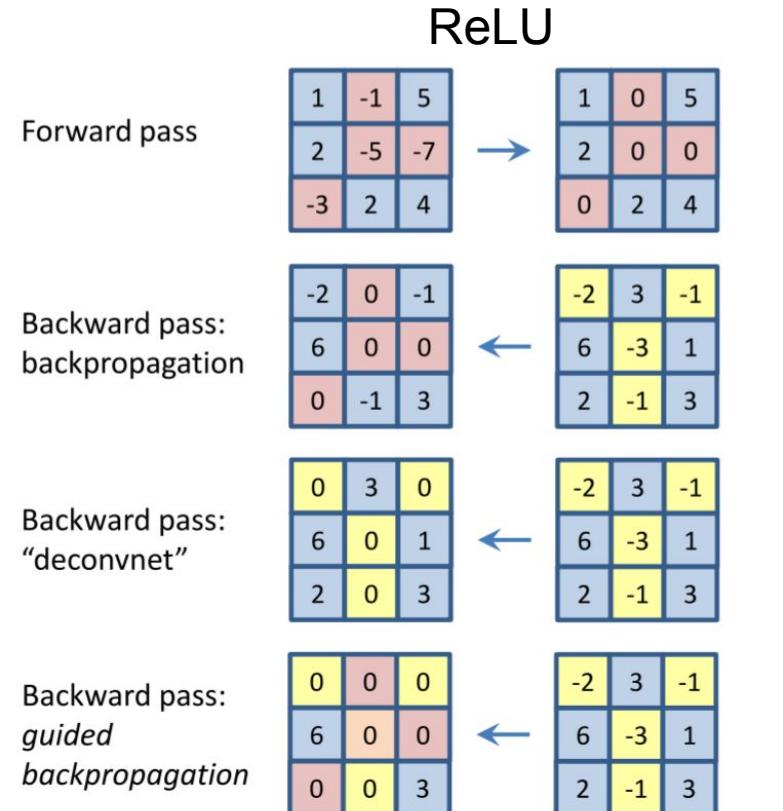
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

# Intermediate features via (guided) backprop



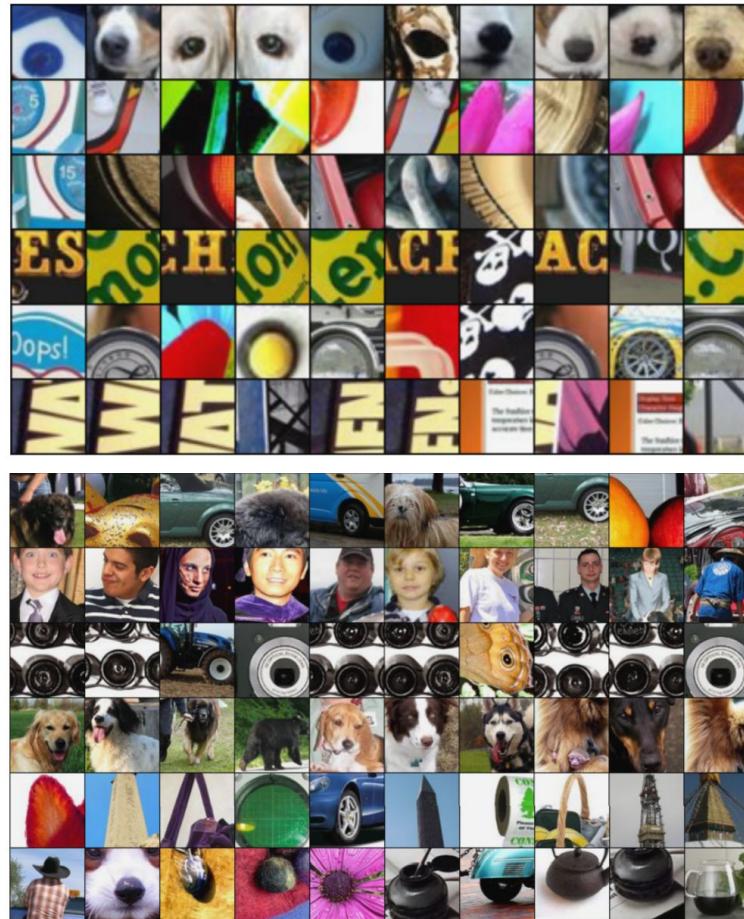
Pick a single intermediate neuron, e.g. one value in  $128 \times 13 \times 13$  conv5 feature map

Compute gradient of neuron value with respect to image pixels



Images come out nicer if you only backprop positive gradients through each ReLU (guided backprop)

# Intermediate features via (guided) backprop



Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014  
Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015

Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

# Visualizing CNN features: Gradient Ascent

**(Guided) backprop:**

Find the part of an image that a neuron responds to

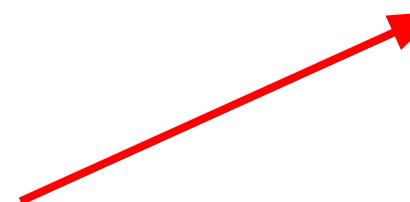
**Gradient ascent:**

Generate a synthetic image that maximally activates a neuron

$$I^* = \arg \max_I f(I) + R(I)$$

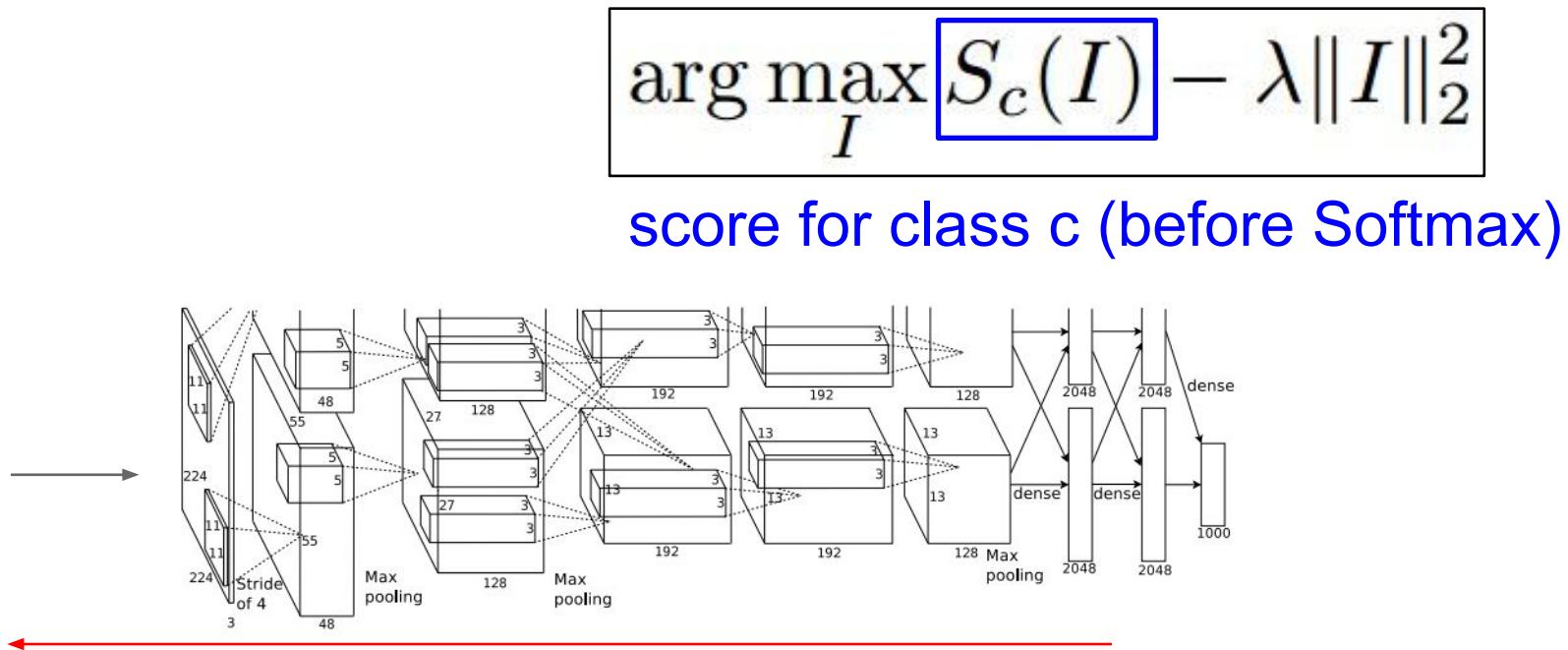
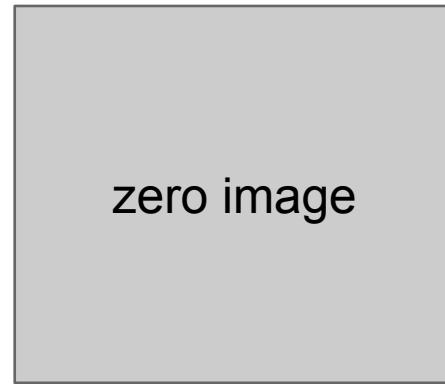
Neuron value

Natural image regularizer



# Visualizing CNN features: Gradient Ascent

1. Initialize image to zeros



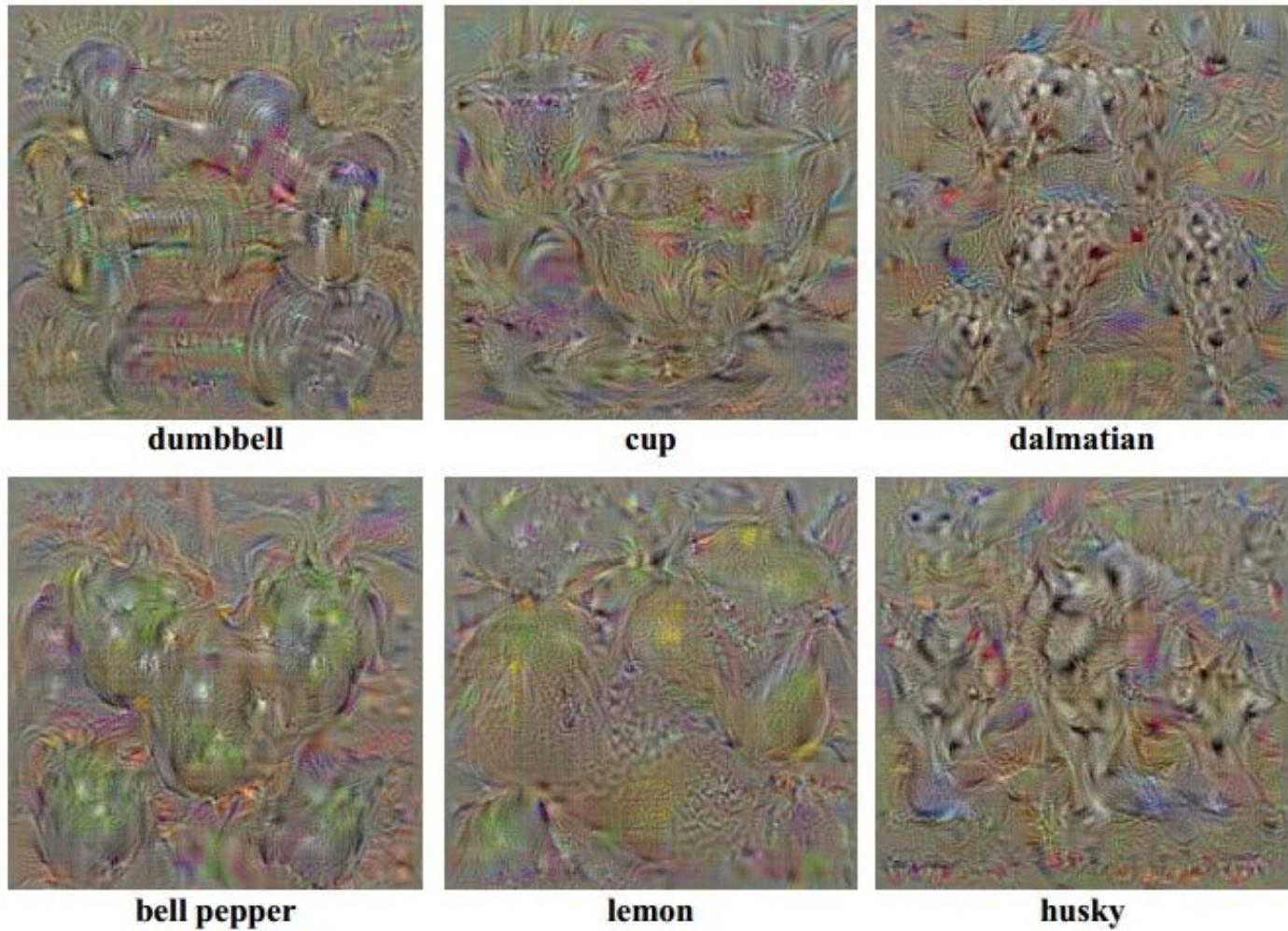
Repeat:

2. Forward image to compute current scores
3. Backprop to get gradient of neuron value with respect to image pixels
4. Make a small update to the image

# Visualizing CNN features: Gradient Ascent

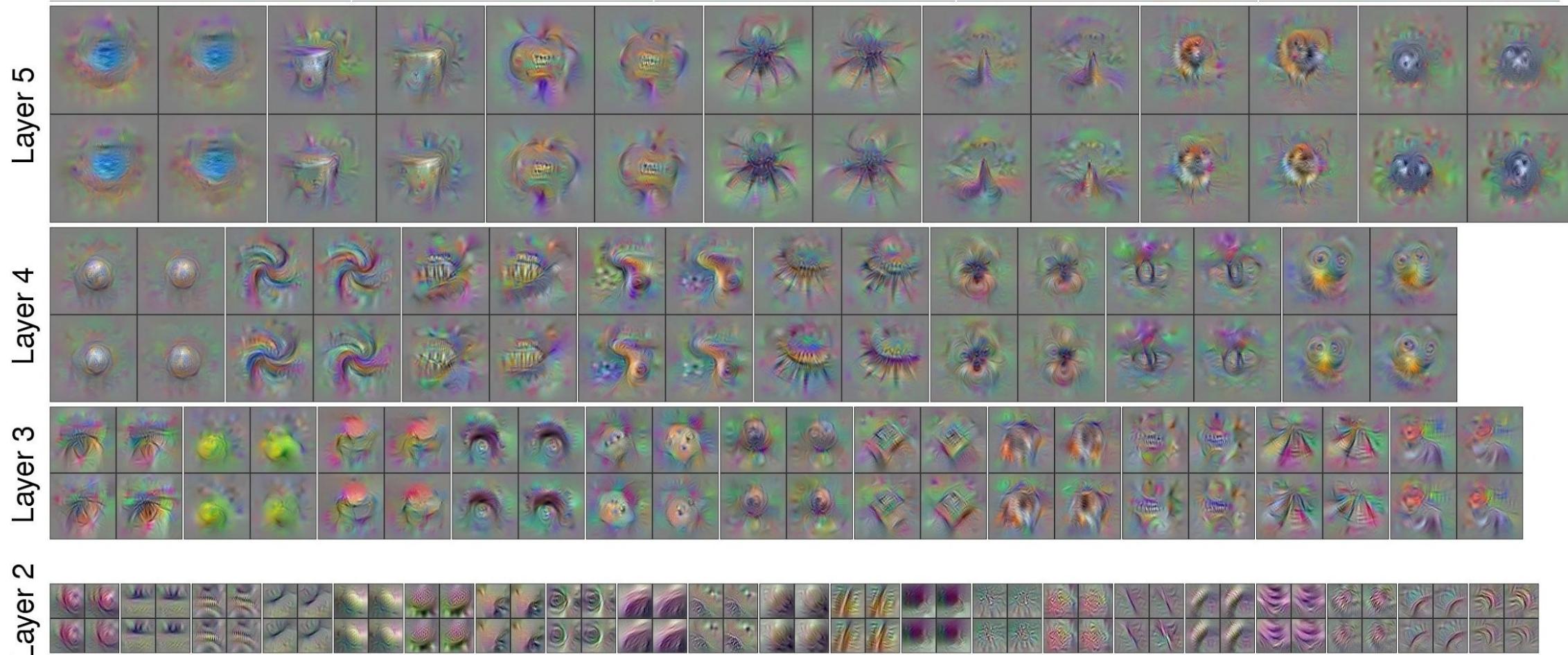
$$\arg \max_I S_c(I) - \boxed{\lambda \|I\|_2^2}$$

Simple regularizer: Penalize L2 norm of generated image



# Visualizing CNN features: Gradient Ascent

Use the same approach to visualize intermediate features



# Visualizing CNN features: Gradient Ascent

Adding “multi-faceted” visualization gives even nicer results:  
(Plus more careful regularization, center-bias)

Reconstructions of multiple feature types (facets) recognized by the same “grocery store” neuron



Corresponding example training set images recognized by the same neuron as in the “grocery store” class



# Visualizing CNN features: Gradient Ascent



Nguyen et al, "Multifaceted Feature Visualization: Uncovering the Different Types of Features Learned By Each Neuron in Deep Neural Networks", ICML Visualization for Deep Learning Workshop 2016.

Figures copyright Anh Nguyen, Jason Yosinski, and Jeff Clune, 2016; reproduced with permission.

# Visualizing CNN features: Gradient Ascent

Optimize in FC6 latent space instead of pixel space:

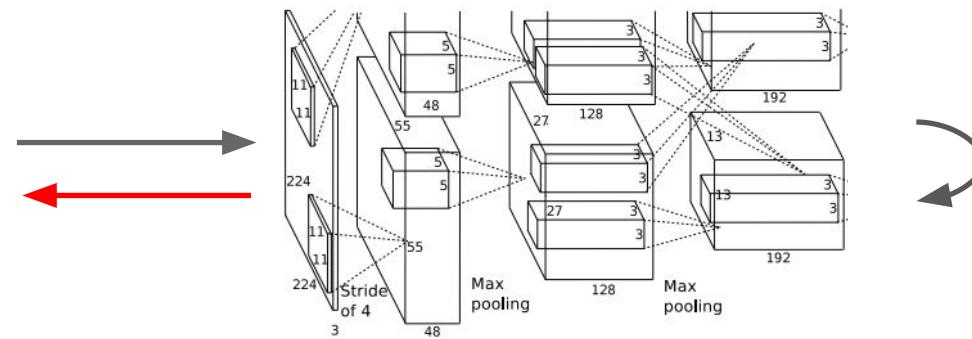


Nguyen et al, "Synthesizing the preferred inputs for neurons in neural networks via deep generator networks," NIPS 2016

Figure copyright Nguyen et al, 2016; reproduced with permission.

# DeepDream: Amplify existing features

Rather than synthesizing an image to maximize a specific neuron, instead try to **amplify** the neuron activations at some layer in the network

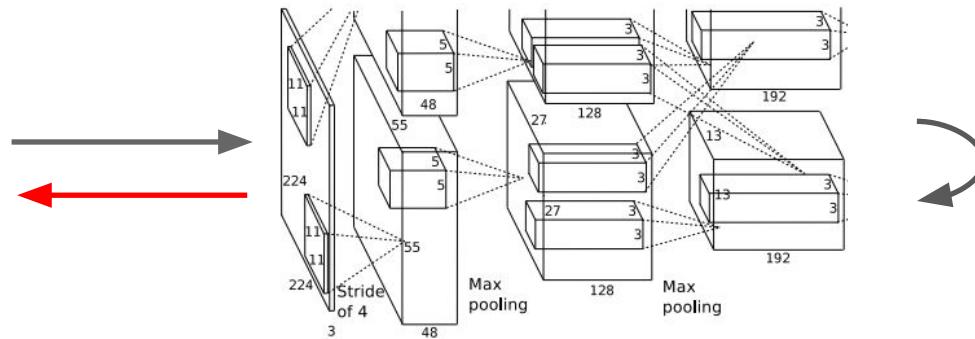


Choose an image and a layer in a CNN; repeat:

1. Forward: compute activations at chosen layer
2. Set gradient of chosen layer *equal to its activation*
3. Backward: Compute gradient on image
4. Update image

# DeepDream: Amplify existing features

Rather than synthesizing an image to maximize a specific neuron, instead try to **amplify** the neuron activations at some layer in the network



Choose an image and a layer in a CNN; repeat:

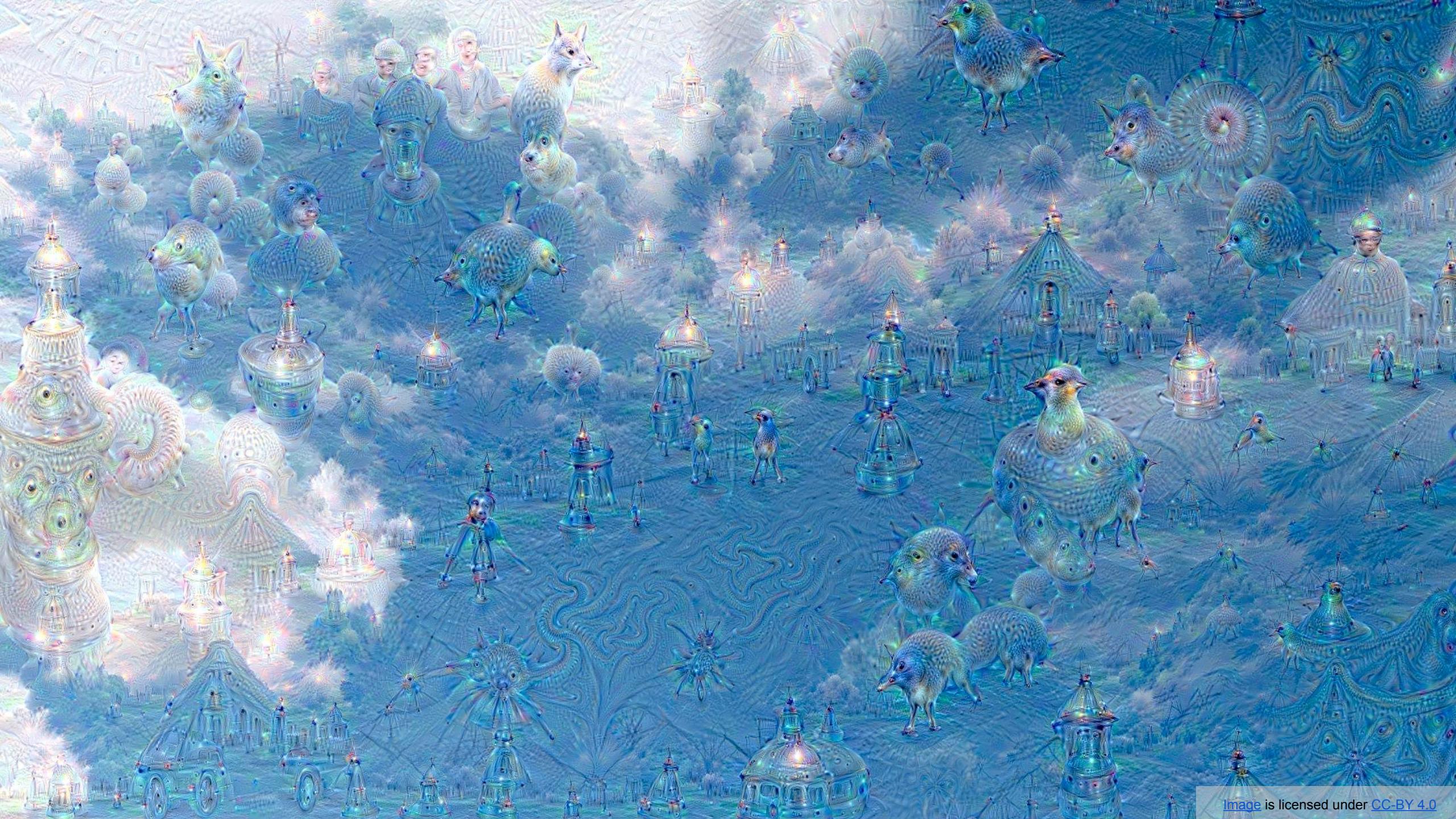
1. Forward: compute activations at chosen layer
2. Set gradient of chosen layer *equal to its activation*
3. Backward: Compute gradient on image
4. Update image

Equivalent to:

$$I^* = \arg \max_I \sum_i f_i(I)^2$$

Mordvintsev, Olah, and Tyka, "Inceptionism: Going Deeper into Neural Networks", [Google Research Blog](#). Images are licensed under CC-BY







"Admiral Dog!"



"The Pig-Snail"



"The Camel-Bird"



"The Dog-Fish"

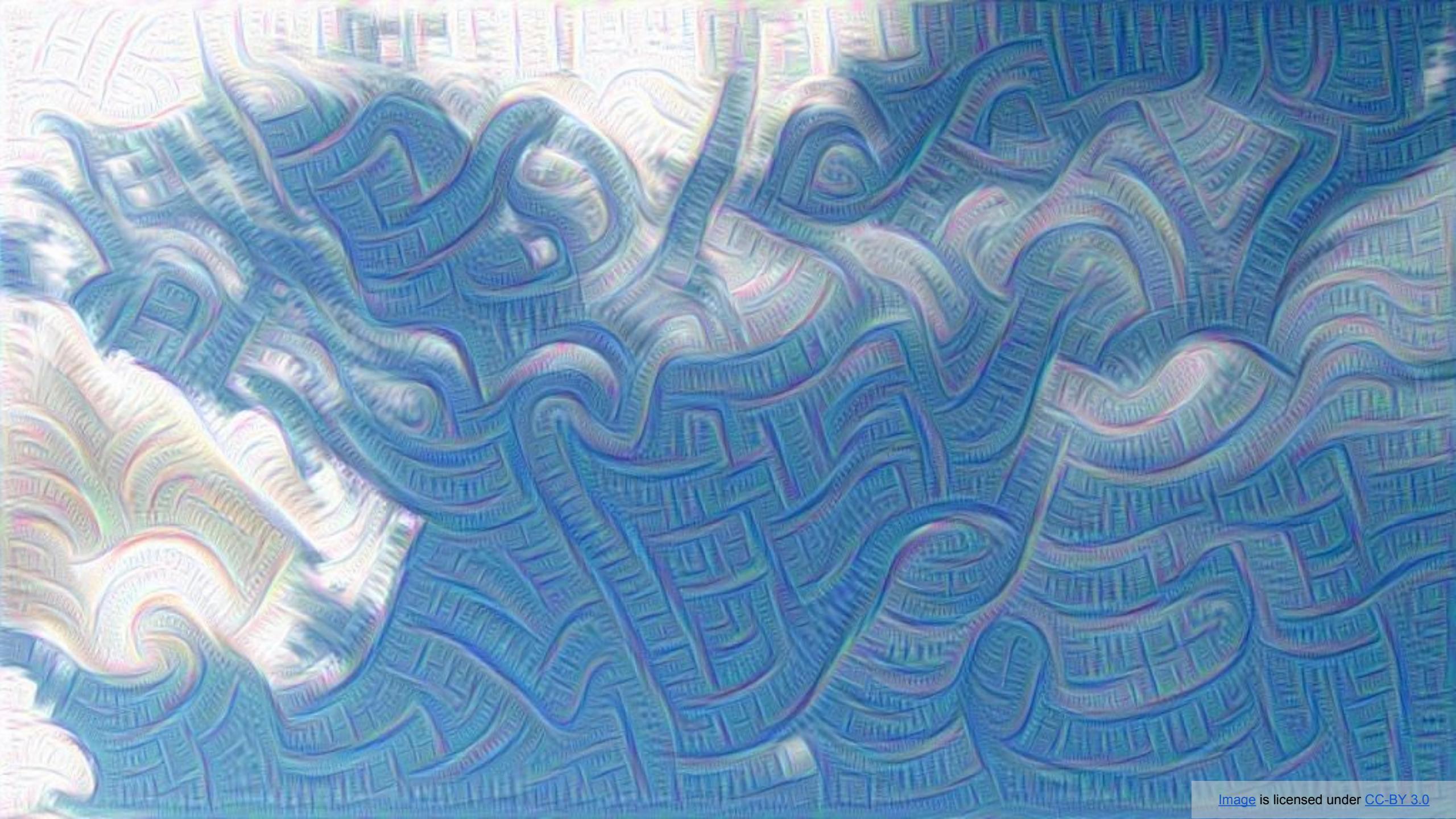




Image is licensed under CC-BY 3.0

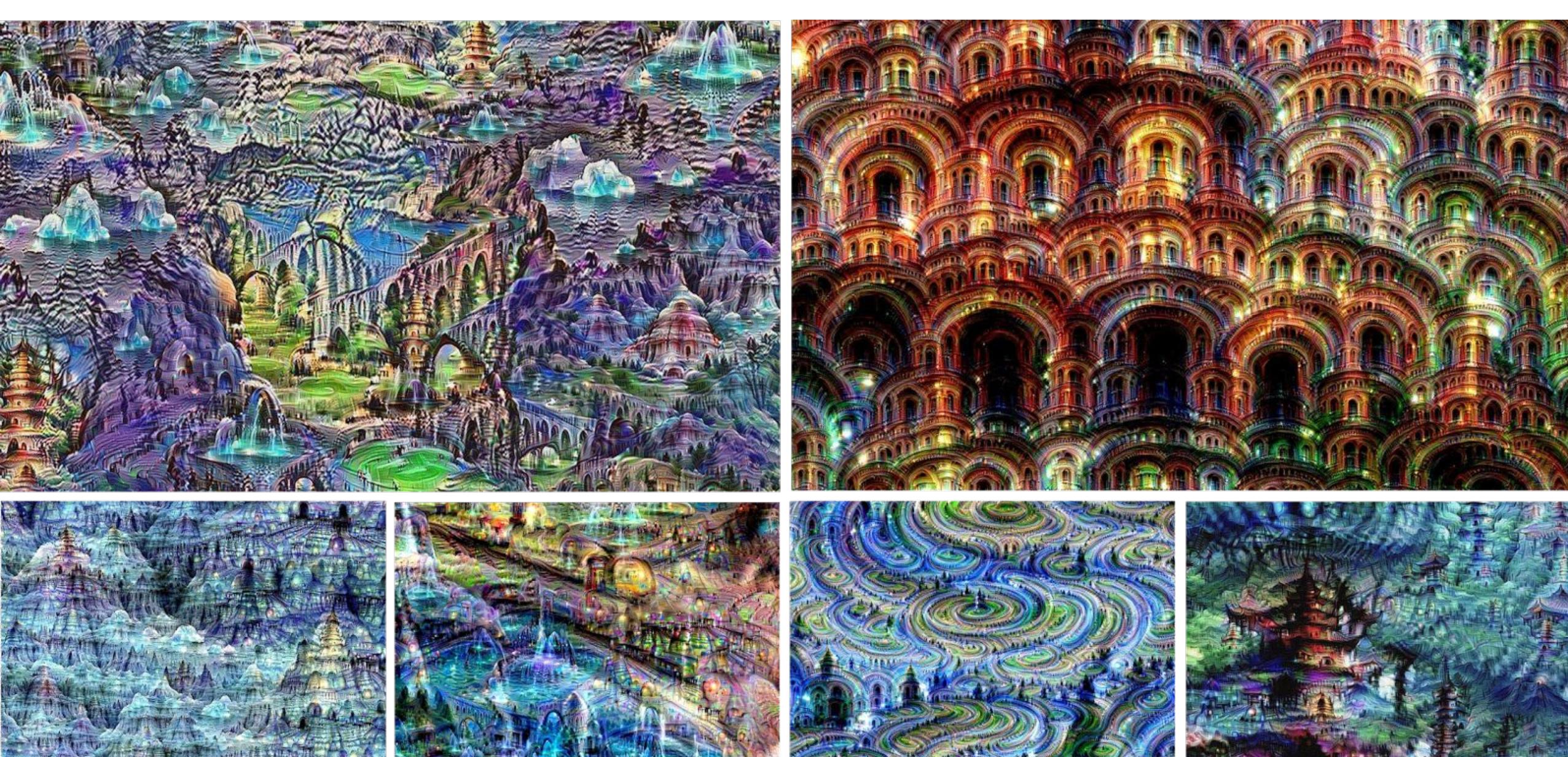


Image is licensed under [CC-BY 4.0](#)

# Neural Style Transfer

Content Image



+

Style Image



=

Style Transfer!



[This image](#) is licensed under CC-BY 3.0

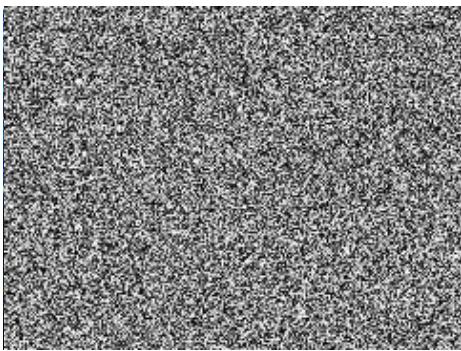
[Starry Night](#) by Van Gogh is in the public domain

[This image](#) copyright Justin Johnson, 2015. Reproduced with permission.

Style  
image



Output  
image  
(Start with  
noise)



Content  
image



Style Target

$$y_s$$

$$\ell_{style}^{\phi, relu1\_2}$$

$$\ell_{style}^{\phi, relu2\_2}$$

$$\ell_{style}^{\phi, relu3\_3}$$

$$\ell_{style}^{\phi, relu4\_3}$$

$$\hat{y}$$

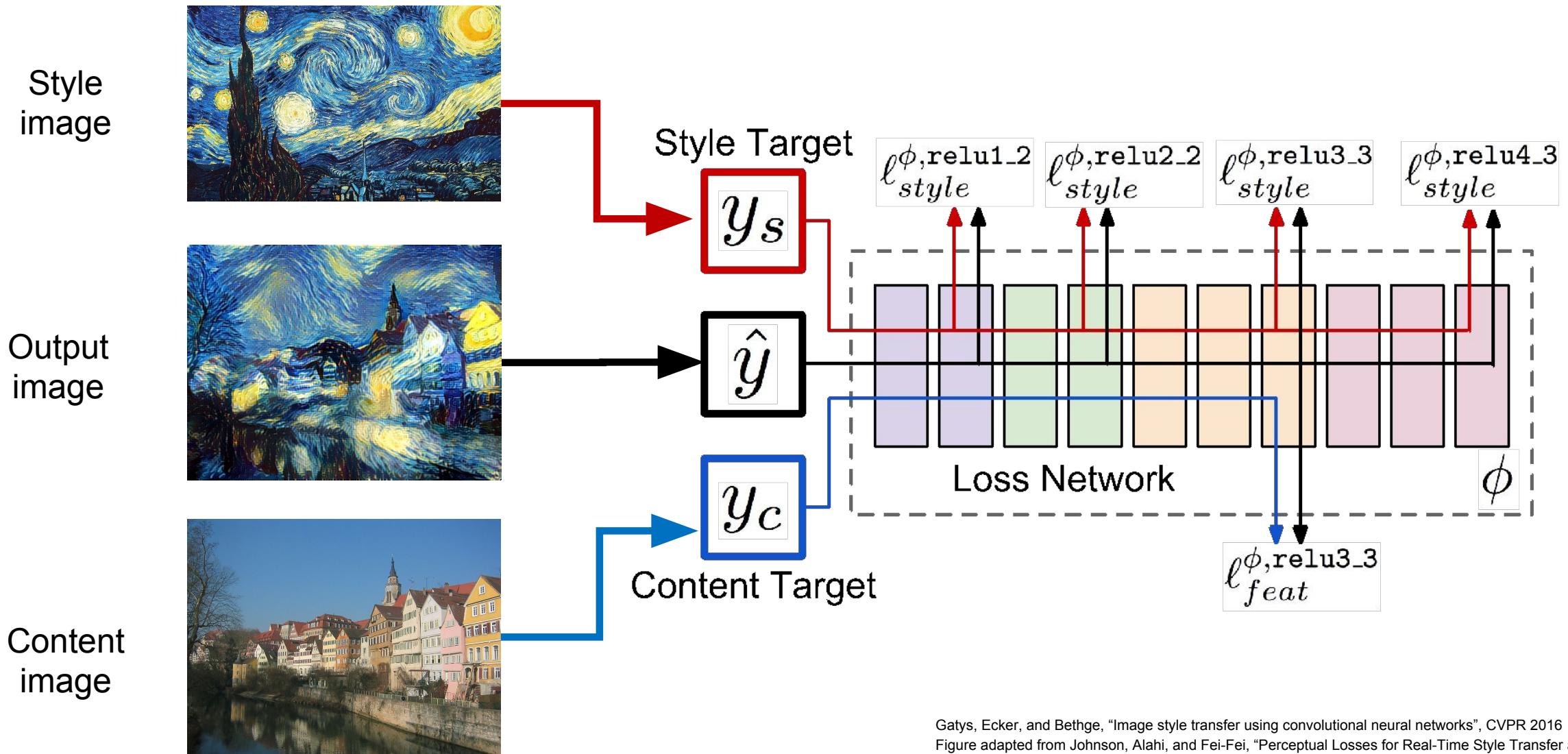
$$y_c$$

Loss Network

$$\phi$$

$$\ell_{feat}^{\phi, relu3\_3}$$

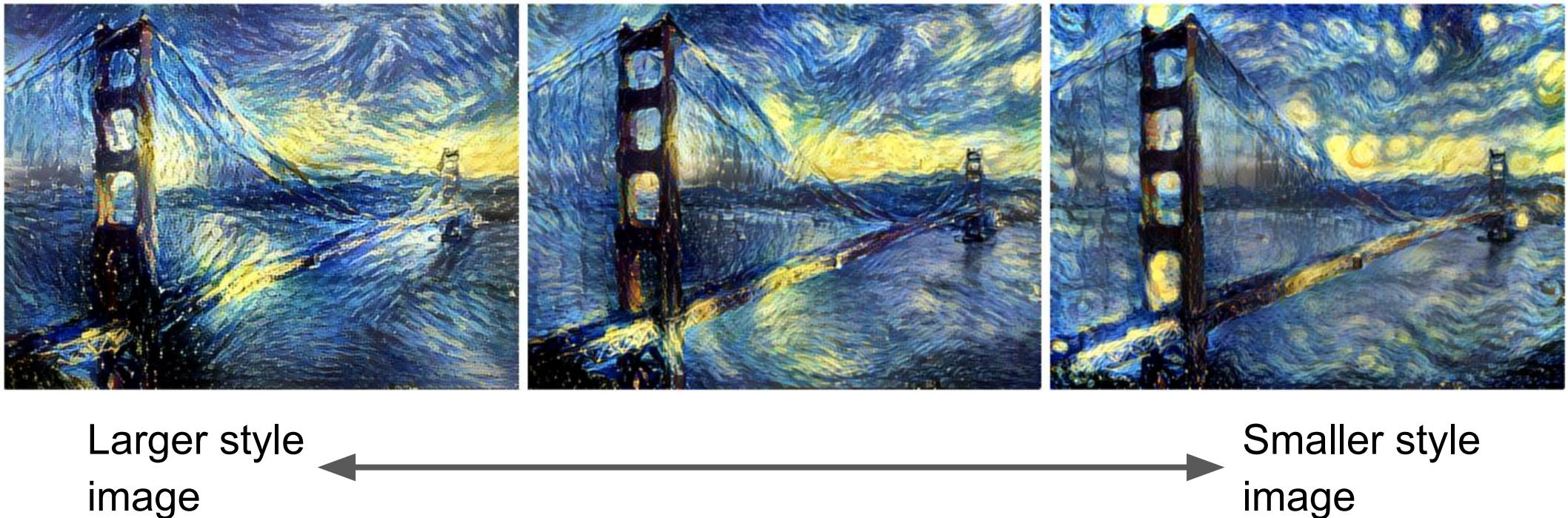
Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016  
Figure adapted from Johnson, Alahi, and Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and  
Super-Resolution", ECCV 2016. Copyright Springer, 2016. Reproduced for educational purposes.



Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016  
 Figure adapted from Johnson, Alahi, and Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution", ECCV 2016. Copyright Springer, 2016. Reproduced for educational purposes.

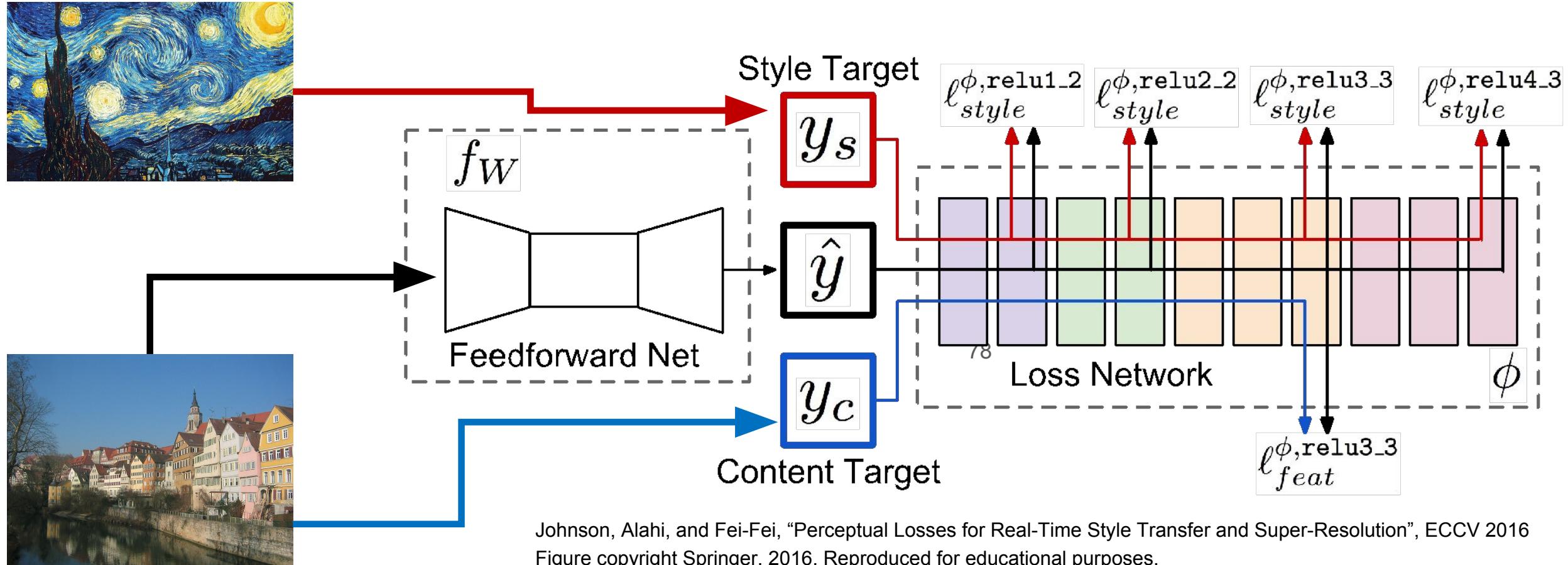
# Neural Style Transfer

Resizing style image before running style transfer algorithm can transfer different types of features



# Fast Style Transfer

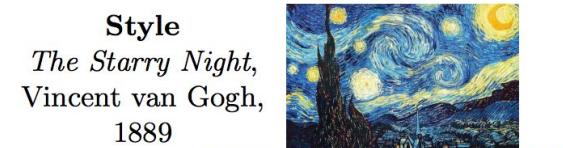
- (1) Train a feedforward network for each style
- (2) Use pretrained CNN to compute same losses as before
- (3) After training, stylize images using a single forward pass



Johnson, Alahi, and Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution", ECCV 2016  
Figure copyright Springer, 2016. Reproduced for educational purposes.

# Fast Style Transfer

**Style**  
*The Starry Night*,  
Vincent van Gogh,  
1889



Slow

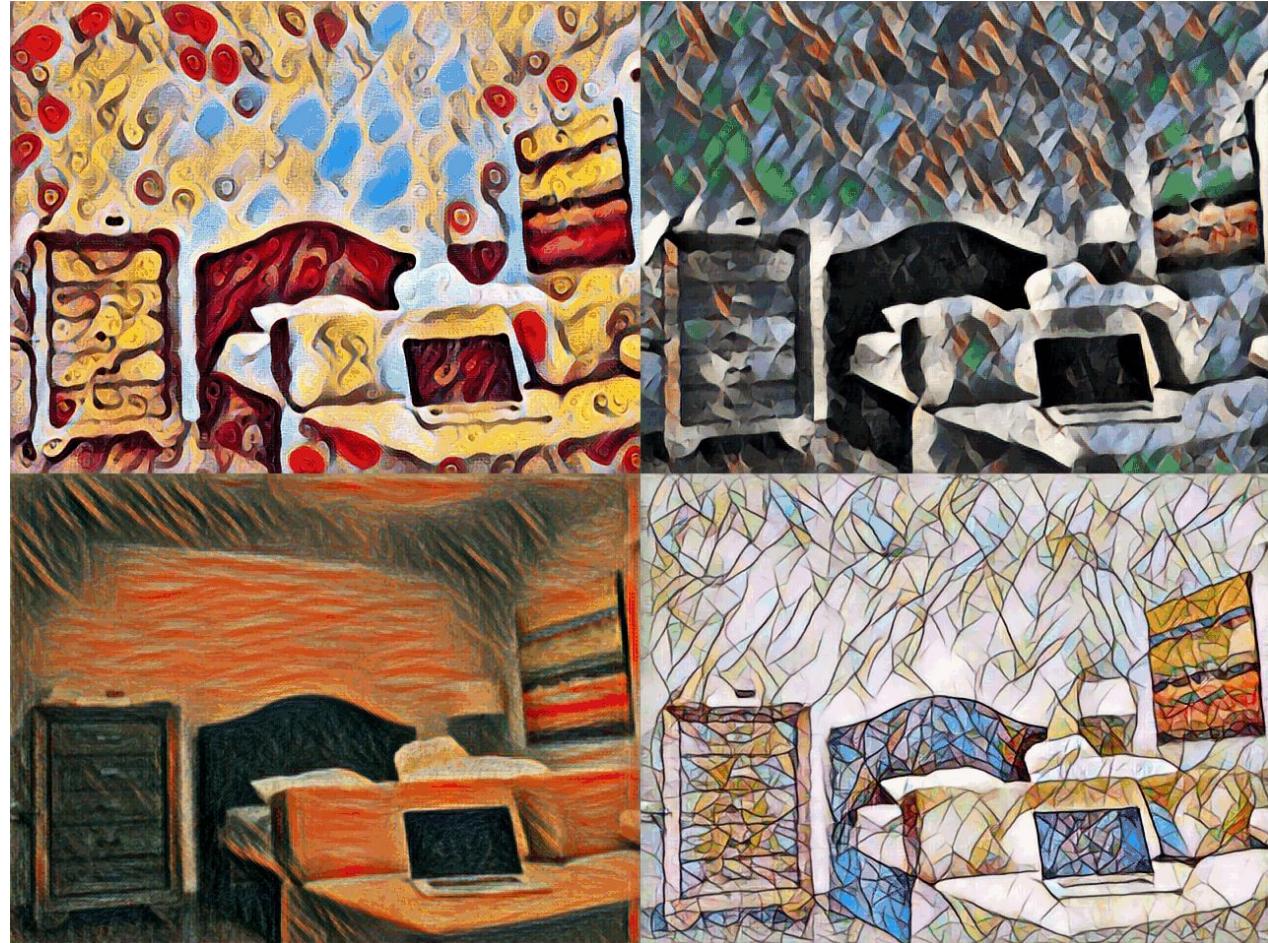
Fast

**Style**  
*The Muse*,  
Pablo Picasso,  
1935



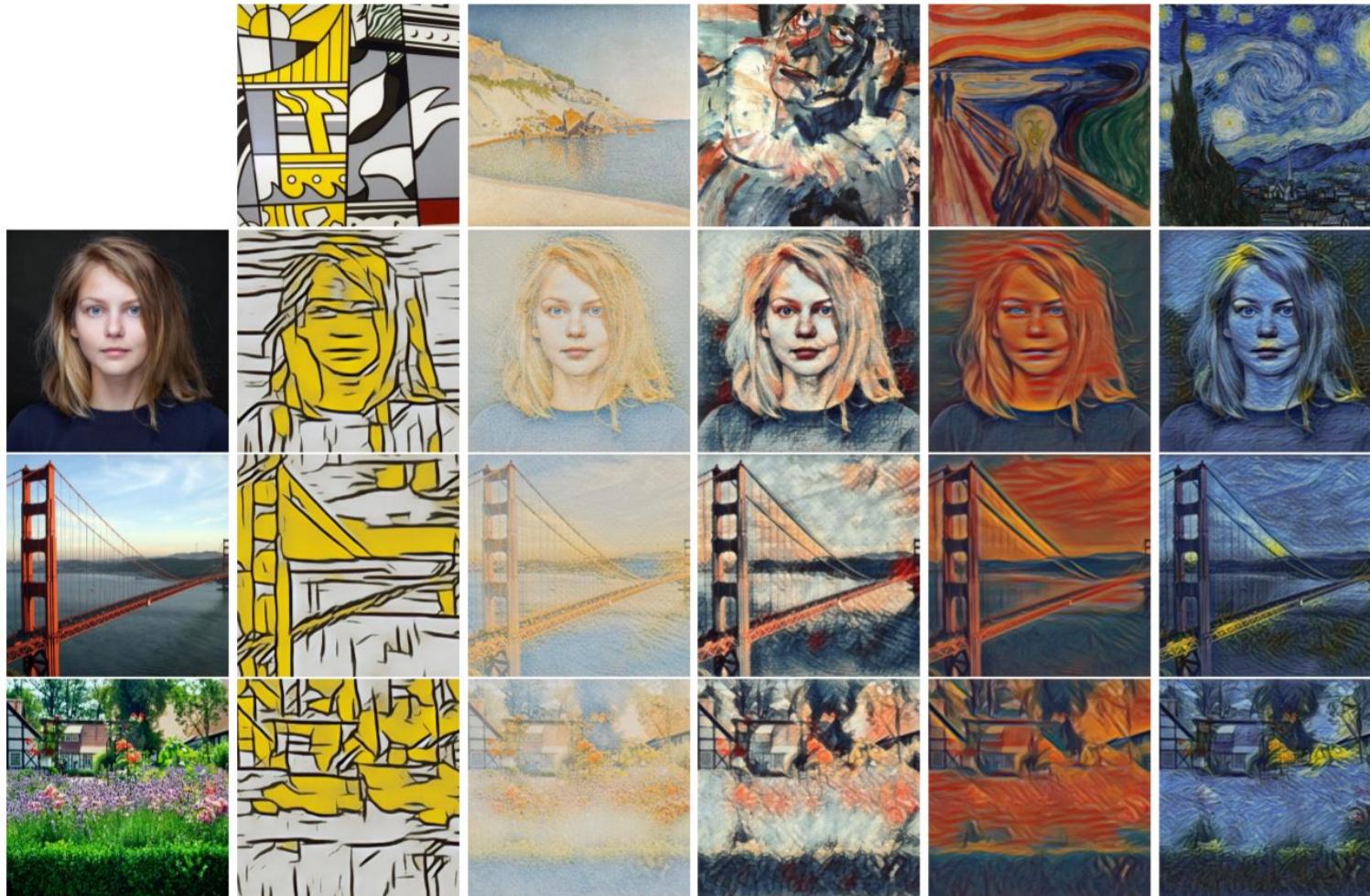
Slow

Fast



<https://github.com/jcjohnson/fast-neural-style>

# One Network, Many Styles



Dumoulin, Shlens, and Kudlur, "A Learned Representation for Artistic Style", ICLR 2017.

Figure copyright Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur, 2016; reproduced with permission.

# Some Real-world Commercial Applications

Facebook Accessibility: <https://www.youtube.com/watch?v=F2eD8Co0z34>

Blue River Technology: <https://www.youtube.com/watch?v=Z-jNFHcOPnc>

Tesla: <https://www.youtube.com/watch?v=zywlvINSIaI&t=1452s>