

Computer Vision

DATA.ML.300, 5 study credits

Esa Rahtu
Unit of Computing Sciences, Tampere University

Image classification

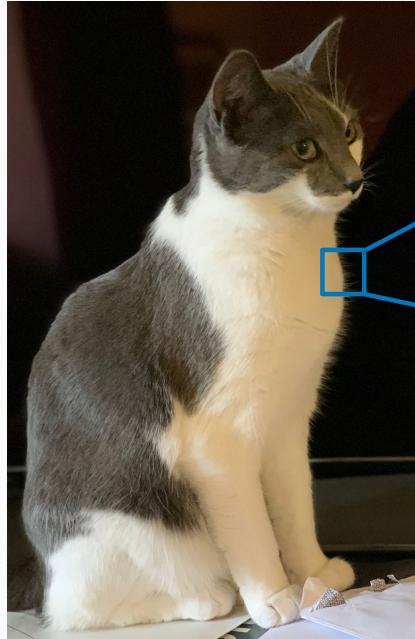
Image classification – Fundamental task in vision



Classify image using given set of discrete labels {dog, cat, truck, plane, ...}

cat

Challenges – Semantic gap



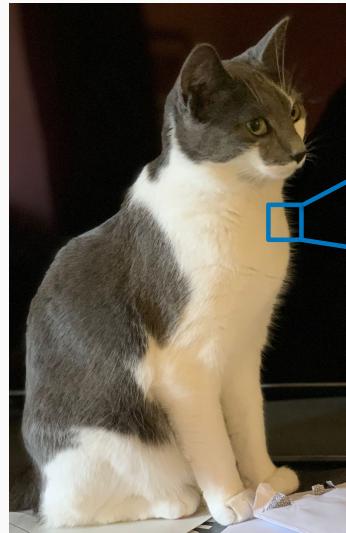
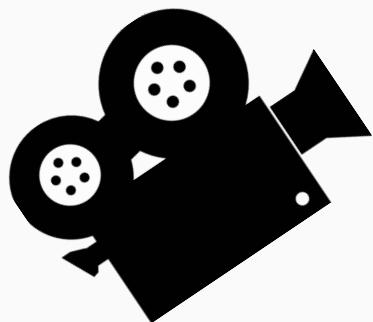
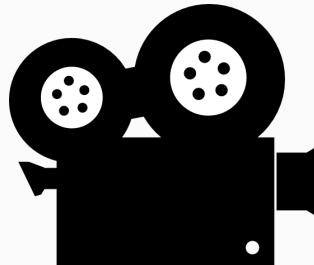
What you see

08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	91	08
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	48	04	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	53	88	30	03	49	13	36	65
52	70	95	23	04	60	11	42	69	24	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	63	89	41	92	36	54	22	40	40	28	66	33	13	80
24	47	32	60	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40

What the computer sees

For a computer the image is just an array of numbers (pixels) between [0, 255].
E.g. 1024 x 1024 x 3 (RGB)

Challenges – Viewpoint



49	99	40	17	81	18	57	60	87	17	40	98	43	69	48	04	56	62	00	
49	08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	91
70	49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	48	04	56	62
31	81	49	31	73	55	79	14	29	93	71	40	67	53	88	30	03	49	13	36
47	52	70	95	23	04	60	11	42	69	24	68	56	01	32	56	71	37	02	36
98	22	31	16	71	51	67	63	89	41	92	36	54	22	40	40	28	66	33	13
26	24	47	32	60	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12
55	32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64
36	67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94
17	24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63
39	21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33
56	78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56
80	16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85
	86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17

All pixel values change if camera moves!

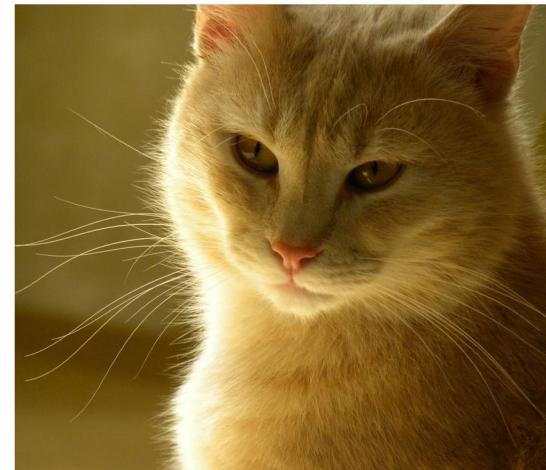
Challenges – Illumination



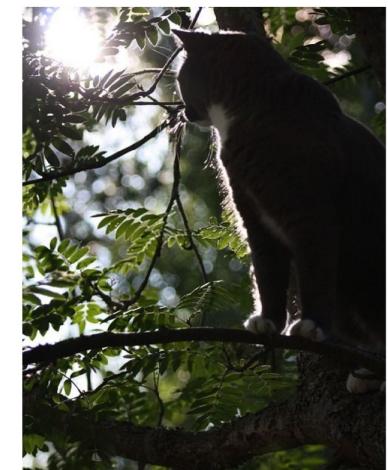
[This image is CC0 1.0 public domain](#)



[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 to build image classifier?

- Extremely hard task!
- No obvious way to "hard-code" the algorithm for recognizing image categories (e.g. cat). (c.f. other problems like sorting a list of numbers)

```
def classify_image(image):  
    # Some magic here?  
    return class_label
```

How to build image classifier?

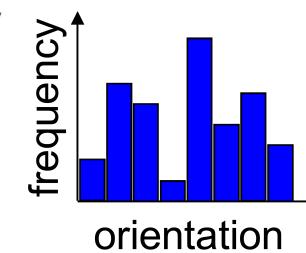
- Extremely hard task!
- No obvious way to "hard-code" the algorithm for recognizing image categories (e.g. cat). (c.f. other problems like sorting a list of numbers)
- Many attempts have been made over the years:



Edge detection
→



Count frequency
of orientations
→



Classify
→ Flower

Image classification with machine learning

- Instead of hard coding, learn from examples!
- Data driven approach:

Image classification with machine learning

- Instead of hard coding, learn from examples!
- Data driven approach:
 1. Collect a dataset of images and labels

Example training set

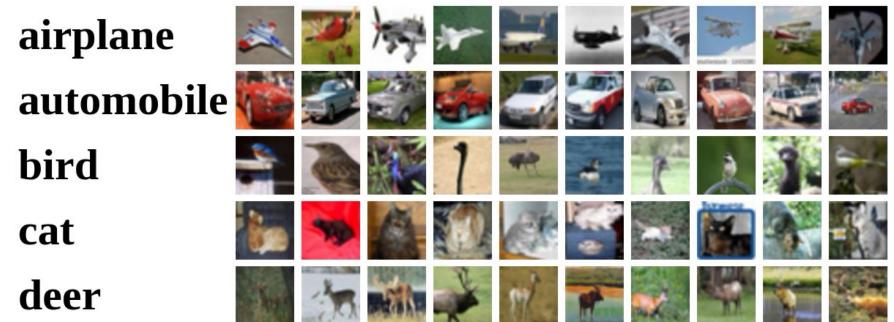
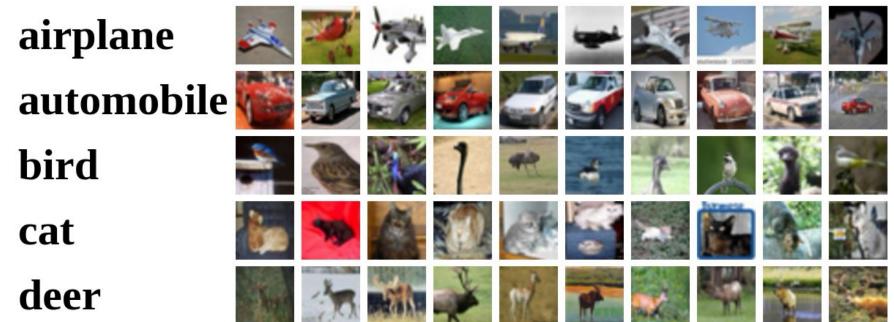


Image classification with machine learning

- Instead of hard coding, learn from examples!
- Data driven approach:
 1. Collect a dataset of images and labels
 2. Use machine learning to train a classifier

Example training set



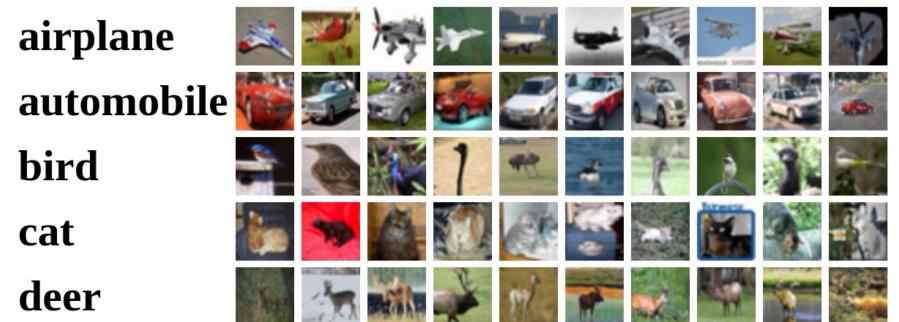
```
def train(images, labels):  
    # Machine learning!  
    return model
```

Image classification with machine learning

- Instead of hard coding, learn from examples!
- 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

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

Example training set



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

Nearest neighbour classifier

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

Nearest neighbour classifier

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



Memorize all data and labels

Nearest neighbour classifier

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



Memorize all data and labels

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



Predict the label of the most
similar training image

Example: CIFAR 10

10 classes
50 000 training images
10 000 testing images



How to determine the most similar image?

- Determine distance metric between images

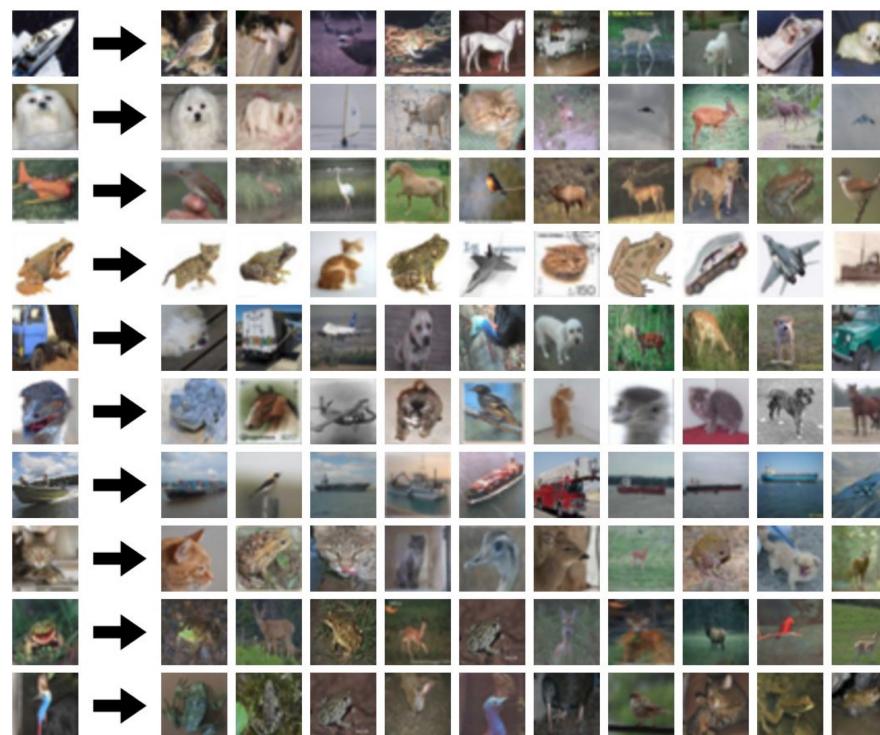
- L1 distance:
$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

Test image				Training image				Pixel-wise absolute difference			
56	32	10	18	10	20	24	17	46	12	14	1
90	23	128	133	8	10	89	100	82	13	39	33
24	26	178	200	12	16	178	170	12	10	0	30
2	0	255	220	4	32	233	112	2	32	22	108

- = add → 456

Figure from Fei-Fei Li

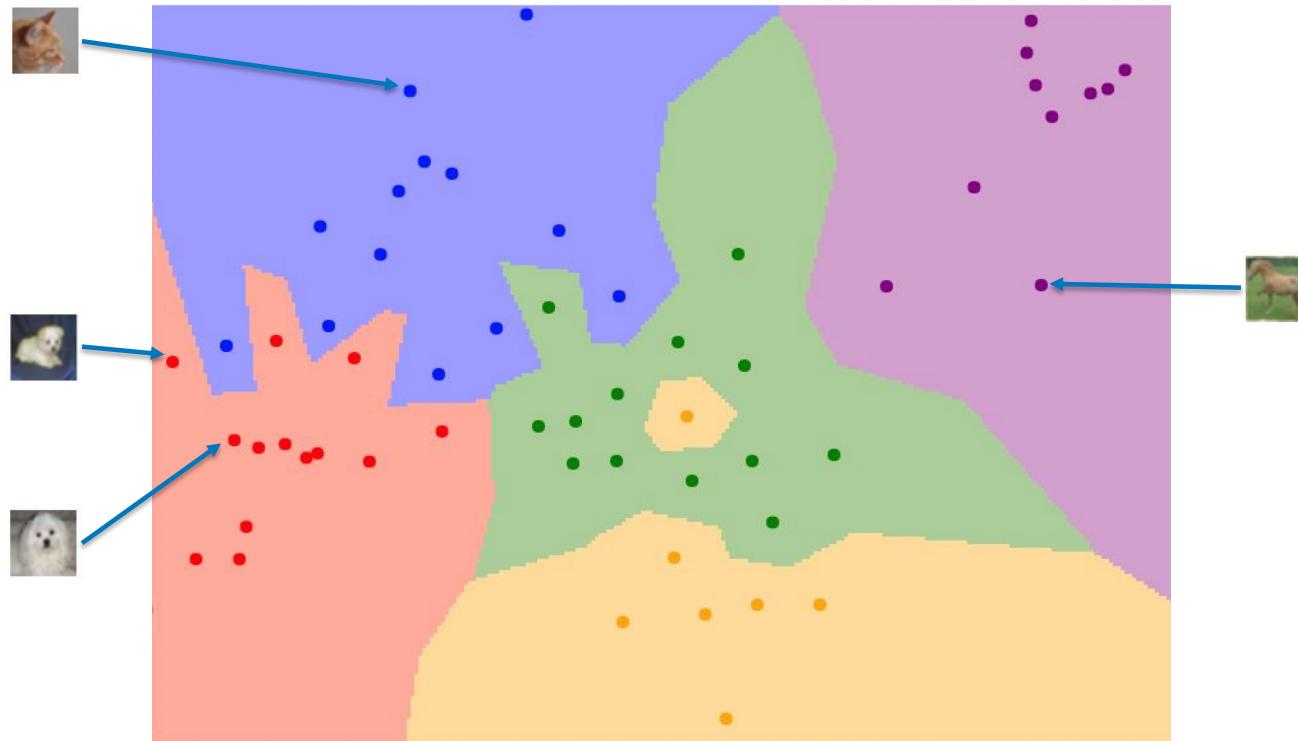
How does this look like?



How does this look like?

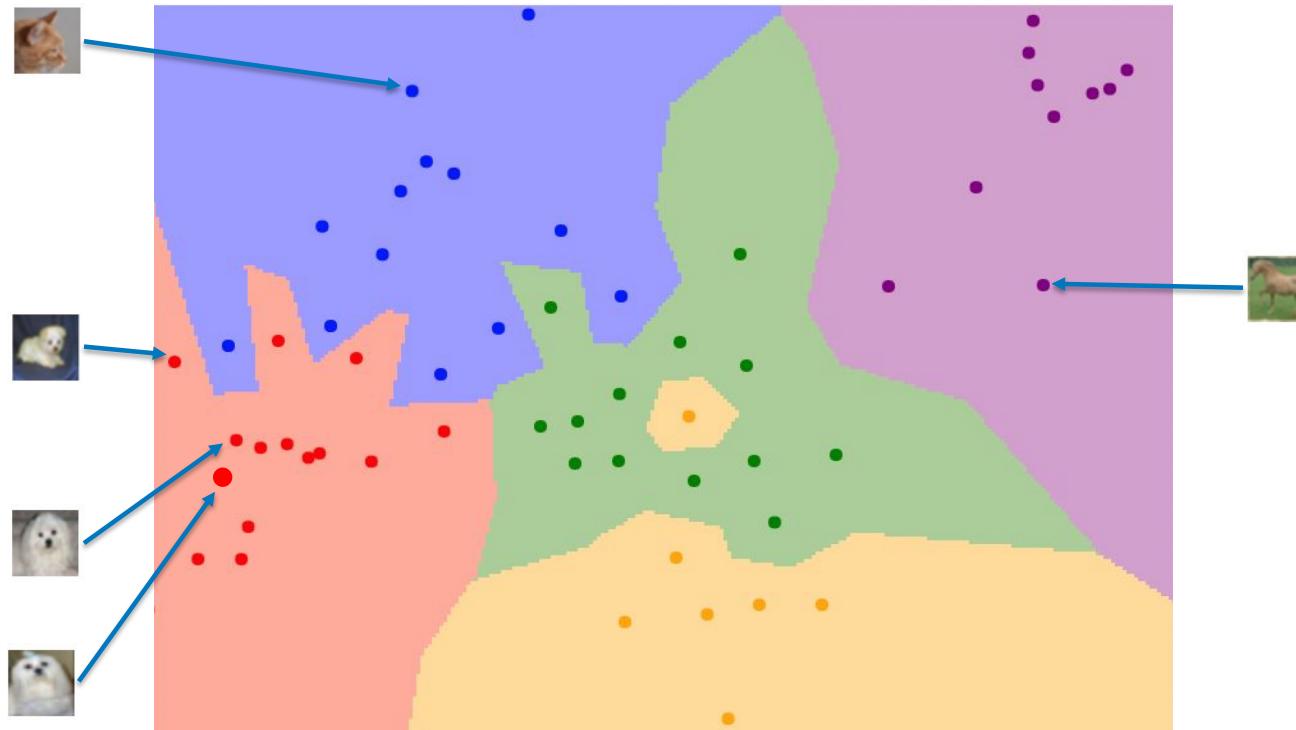


How does this look like?



Credits: Fei-Fei Li

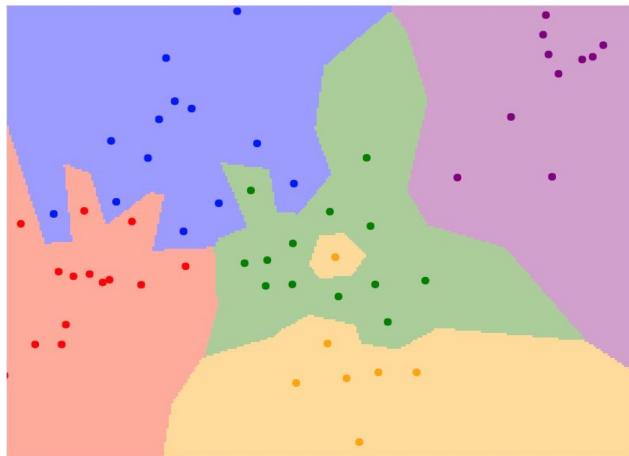
How does this look like?



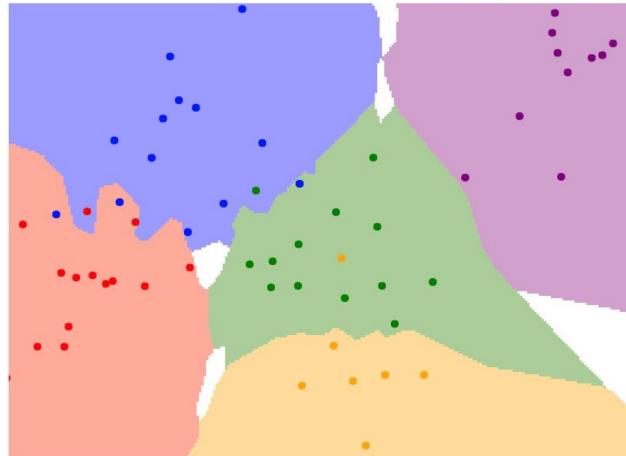
Credits: Fei-Fei Li

K-Nearest neighbours

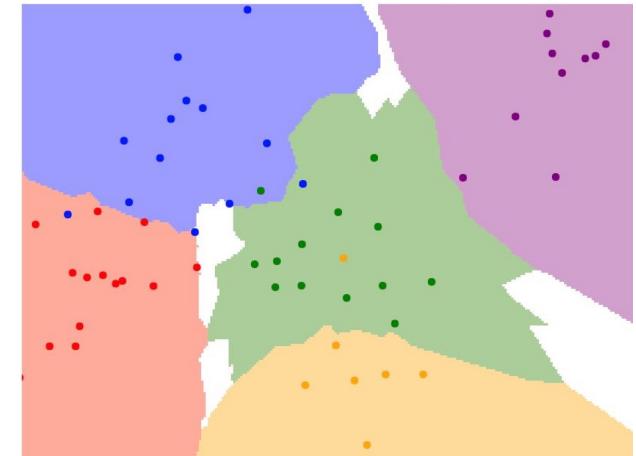
- Instead of copying label from the nearest sample, take **majority vote** from K closest



$K = 1$



$K = 3$



$K = 5$

How to select hyper parameters?

- How to choose an optimal value for K?
- How to choose optimal distance metric?
- These are **hyperparameters**: choices about the algorithm that we set rather than learn.
- Very problem dependent and often need to optimized by trial and error.

How to select hyper parameters?

- Common approach for parameter setting:
 1. Split data into training, validation, and test sets
 2. Choose hyperparameters on validation set
 3. Evaluate generalization using test set



Image classification with neural networks

Recap: Supervised learning

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

Recap: Supervised learning

Functions \mathcal{F}

$$\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$$

Training data

$$\{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}$$

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

find $\hat{f} \in \mathcal{F}$
s.t. $y_i \approx \hat{f}(x_i)$



Learning machine

Recap: Supervised learning

Functions \mathcal{F}

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

Training data

$$\{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}$$

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

find $\hat{f} \in \mathcal{F}$
s.t. $y_i \approx \hat{f}(x_i)$



Learning machine

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

$$y = \hat{f}(x)$$

New data

$$x$$

Recap: Supervised learning

Functions \mathcal{F}

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

Training data

$$\{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}$$

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

find $\hat{f} \in \mathcal{F}$
s.t. $y_i \approx \hat{f}(x_i)$

Find f that minimises
given **loss** for data
 $\frac{1}{N} \sum_{i=1}^N [y_i \neq f(x_i)]$
loss function

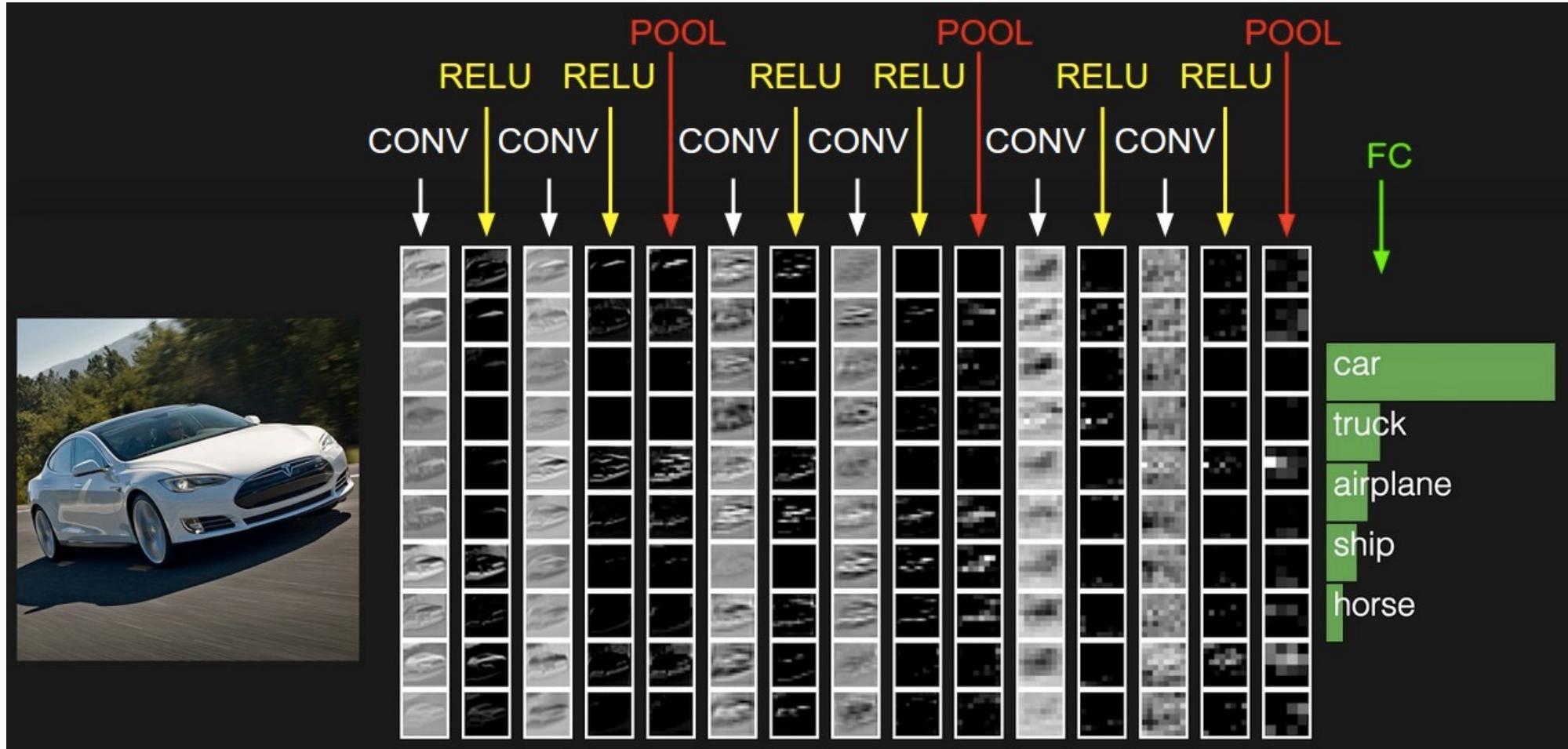
```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```

$$y = \hat{f}(x)$$

New data

$$x$$

Deep convolutional networks for image classification



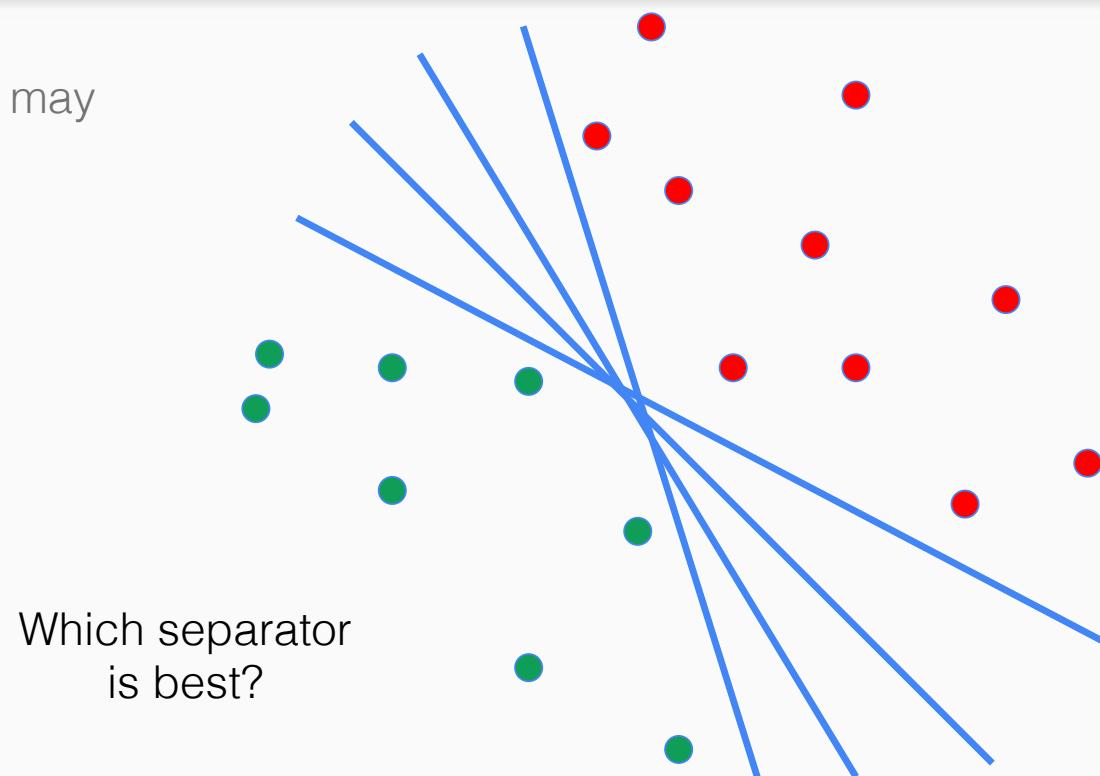
Deep learning

- Learn a feature hierarchy all the way from pixels to classifier
- Each layer extracts features from the output of previous layer
- Train all layers jointly



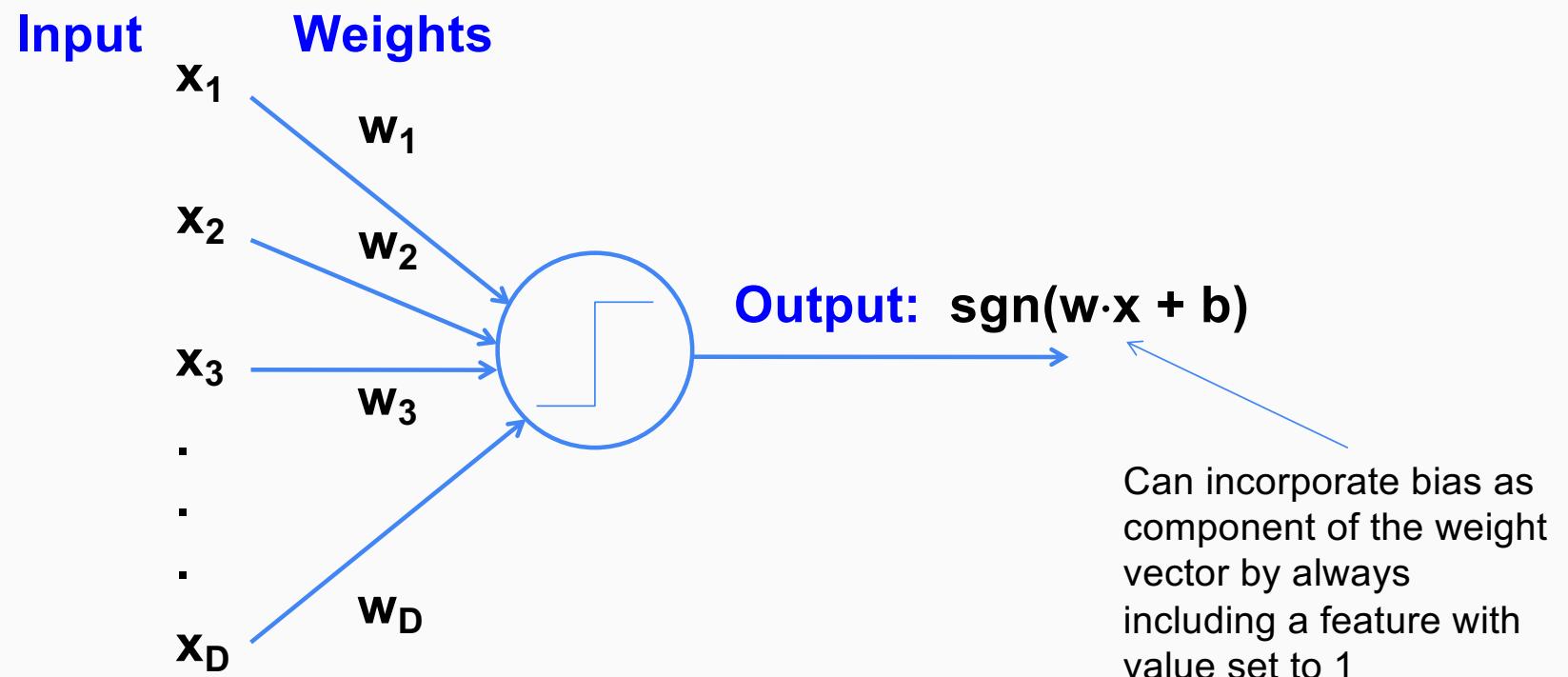
Background: Linear classifiers

- When the data is linearly separable, there may be more than one separator (hyperplane)

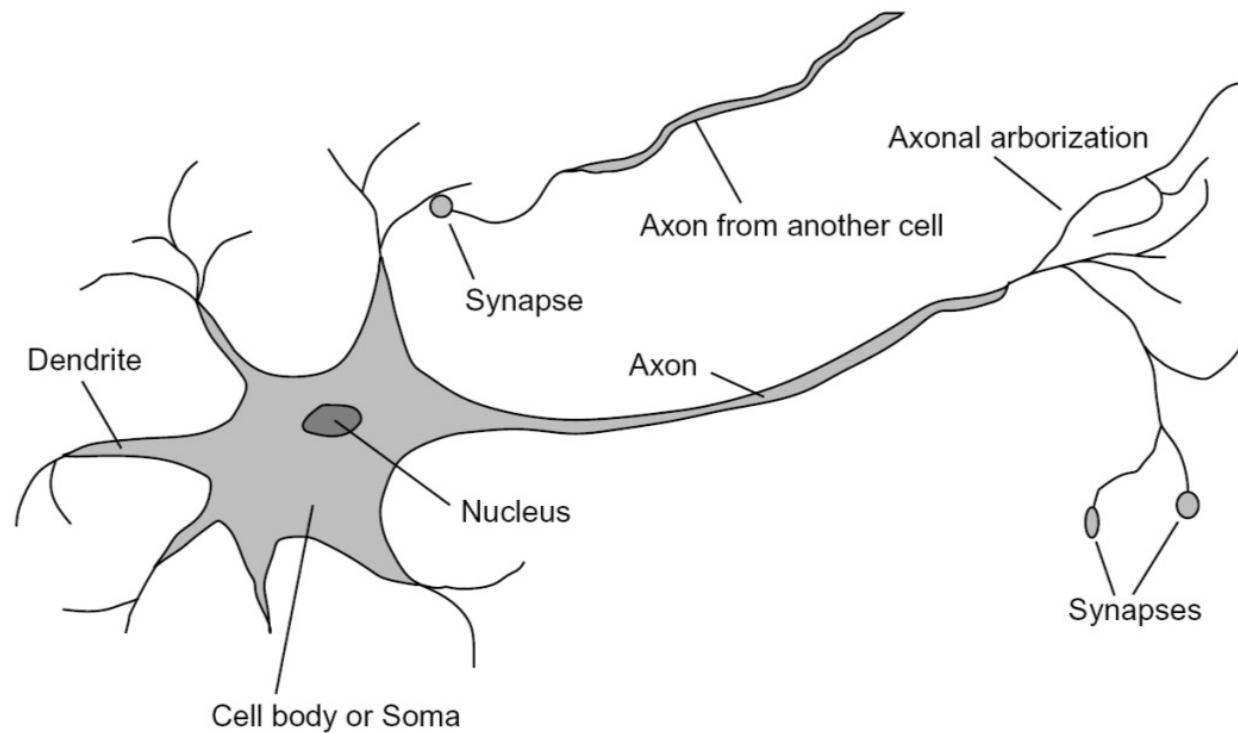


Perceptron

Perceptron



Loose inspiration: Human neurons



History

NEW NAVY DEVICE LEARNS BY DOING

Psychologist Shows Embryo of Computer Designed to Read and Grow Wiser

WASHINGTON, July 7 (UPI)—The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.

The embryo—the Weather Bureau's \$2,000,000 "704" computer—learned to differentiate between right and left after fifty attempts in the Navy's demonstration for newsmen.

The service said it would use this principle to build the first of its Perceptron thinking machines that will be able to read and write. It is expected to be finished in about a year at a cost of \$100,000.

Dr. Frank Rosenblatt, designer of the Perceptron, conducted the demonstration. He said the machine would be the first device to think as the human brain. As do human be-

ings, Perceptron will make mistakes at first, but will grow wiser as it gains experience, he said.

Dr. Rosenblatt, a research psychologist at the Cornell Aeronautical Laboratory, Buffalo, said Perceptrons might be fired to the planets as mechanical space explorers.

Without Human Controls

The Navy said the perceptron would be the first non-living mechanism "capable of receiving, recognizing and identifying its surroundings without any human training or control."

The "brain" is designed to remember images and information it has perceived itself. Ordinary computers remember only what is fed into them on punch cards or magnetic tape.

Later Perceptrons will be able to recognize people and call out their names and instantly translate speech in one language to speech or writing in another language, it was predicted.

Mr. Rosenblatt said in principle it would be possible to build brains that could reproduce themselves on an assembly line and which would be conscious of their existence.

1958 New York Times...

In today's demonstration, the "704" was fed two cards, one with squares marked on the left side and the other with squares on the right side.

Learns by Doing

In the first fifty trials, the machine made no distinction between them. It then started registering a "Q" for the left squares and "O" for the right squares.

Dr. Rosenblatt said he could explain why the machine learned only in highly technical terms. But he said the computer had undergone a "self-induced change in the wiring diagram."

The first Perceptron will have about 1,000 electronic "association cells" receiving electrical impulses from an eye-like scanning device with 400 photo-cells. The human brain has 10,000,000,000 responsive cells, including 100,000,000 connections with the eyes.

Perceptron update rule

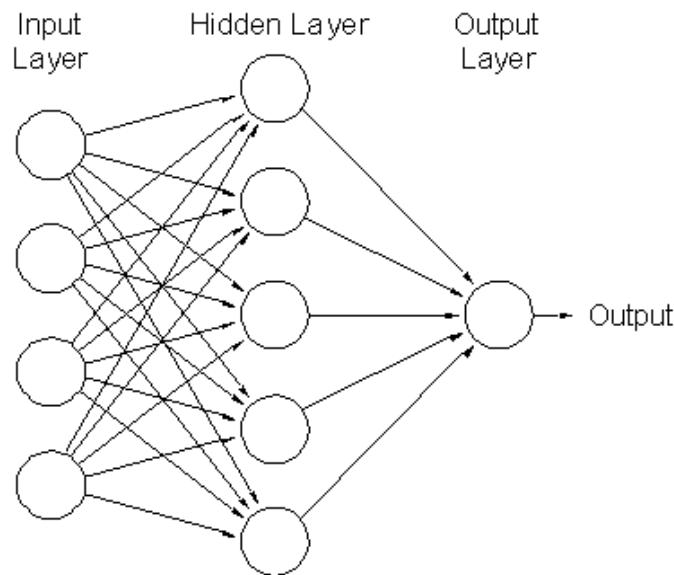
- Initialize weights randomly
- Cycle through training examples in multiple passes (epochs)
- For each training instance x with label y :
 - Classify with current weights: $y' = \text{sgn}(w \cdot x)$
 - Update weights: $w \leftarrow w + \alpha(y - y')x$
- α is a learning rate that should decay as $1/t$ (t is the epoch)
- What happens if y' is correct?
- Otherwise, consider what happens to individual weights
 $w_i \leftarrow w_i + \alpha(y - y')x_i$
 - If $y = 1$ and $y' = -1$, w_i will be increased if x_i is positive or decreased if x_i is negative $\rightarrow w \cdot x$ will get bigger
 - If $y = -1$ and $y' = 1$, w_i will be decreased if x_i is positive or increased if x_i is negative $\rightarrow w \cdot x$ will get smaller

Convergence of perceptron update rule

- **Linearly separable data:** converges to a perfect solution
- **Non-separable data:** converges to a minimum-error solution assuming learning rate decays as $O(1/t)$ and examples are presented in random sequence

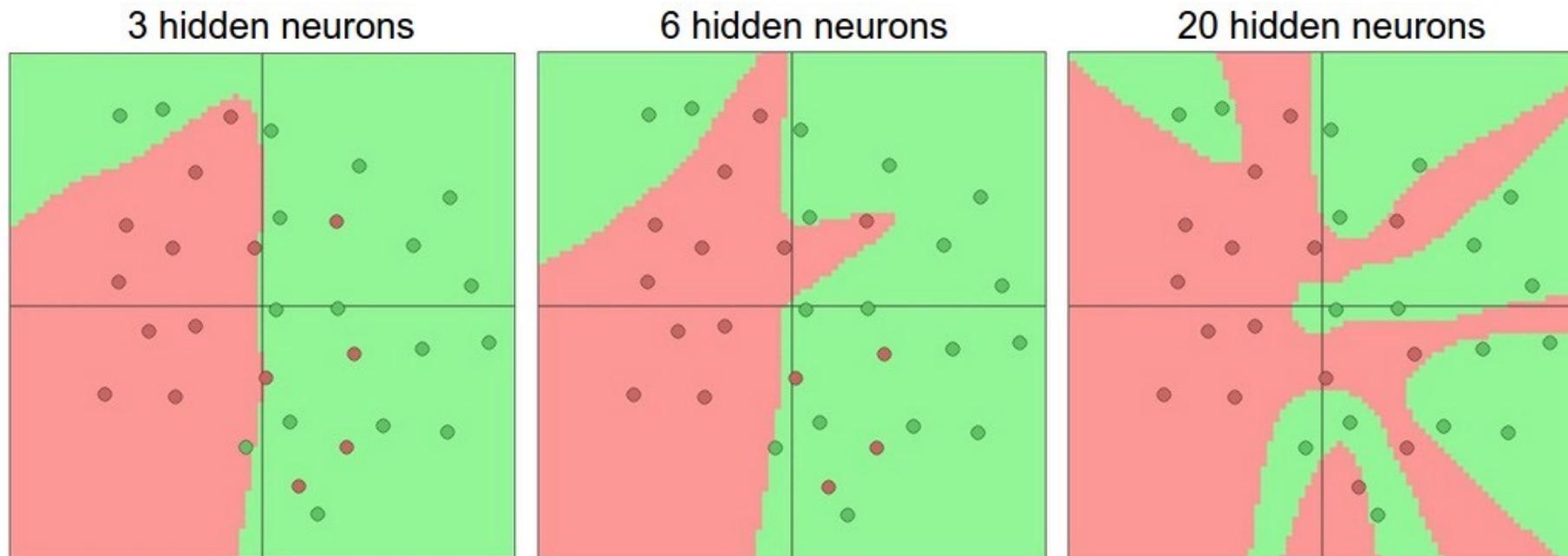
Multi-Layer Neural Networks

- Network with a hidden layer:



- Can represent nonlinear functions (provided each perceptron has a nonlinearity)

Multi-Layer Neural Networks



Source: <http://cs231n.github.io/neural-networks-1/>

Multi-Layer Neural Networks

- Beyond a single hidden layer:

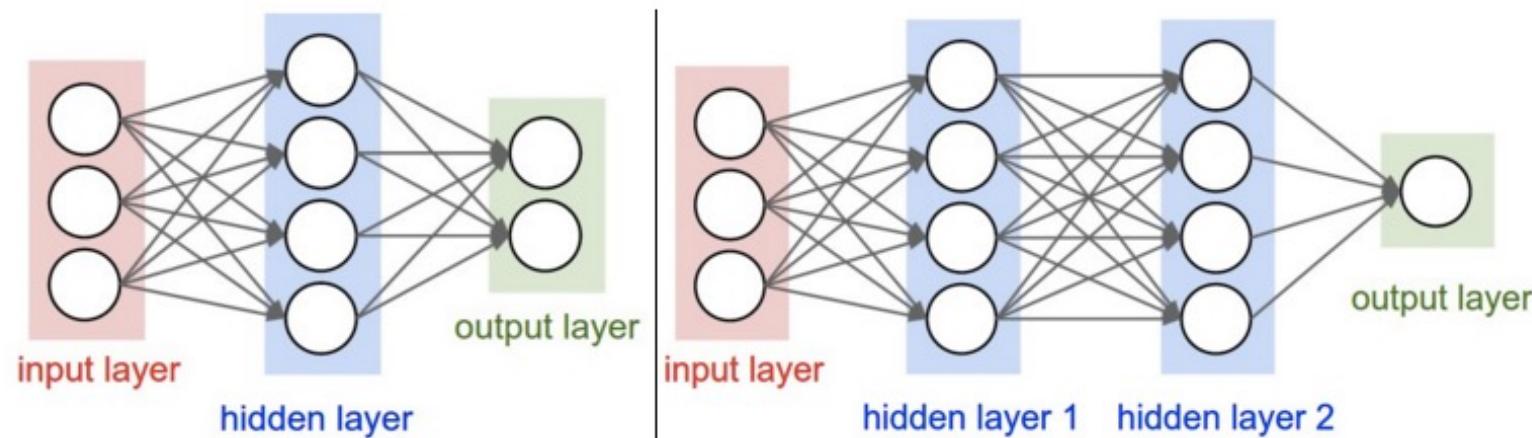
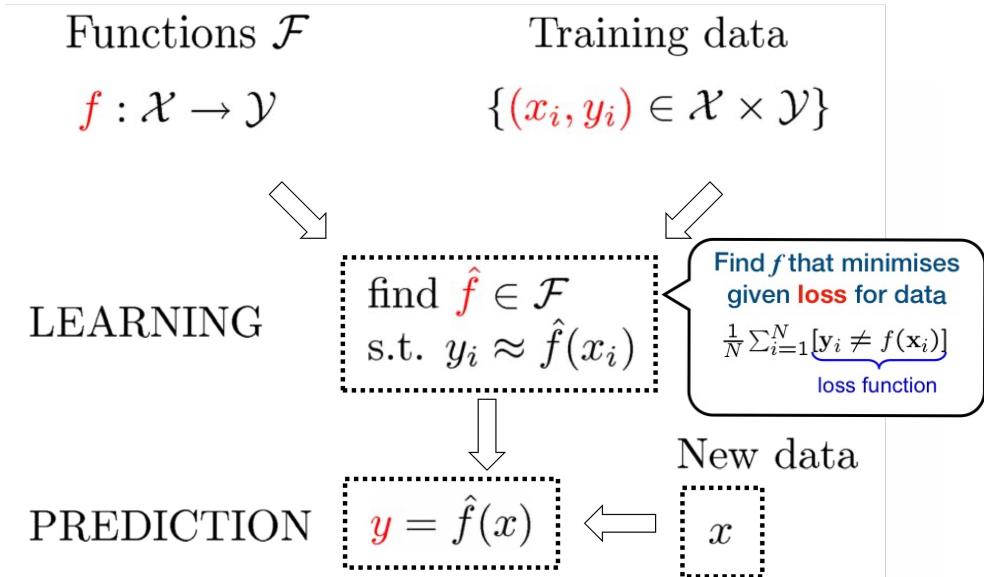


Figure source: <http://cs231n.github.io/neural-networks-1/>

Training of multi-layer networks

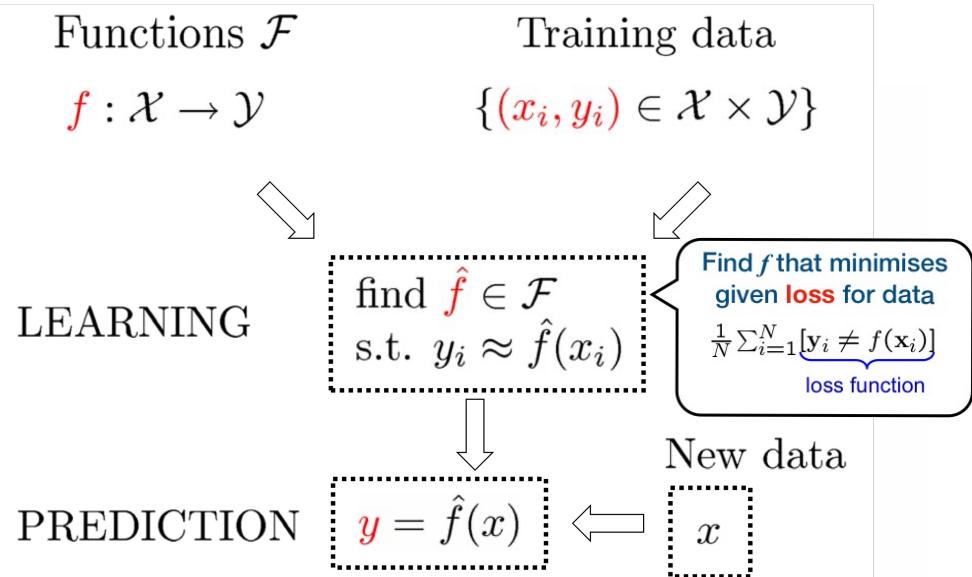
- Find network weights to minimize the error between true and estimated labels of training examples:



Training of multi-layer networks

- Find network weights to minimize the error between true and estimated labels of training examples:

$$E(\mathbf{w}) = \sum_{j=1}^N (y_j - f_{\mathbf{w}}(\mathbf{x}_j))^2$$



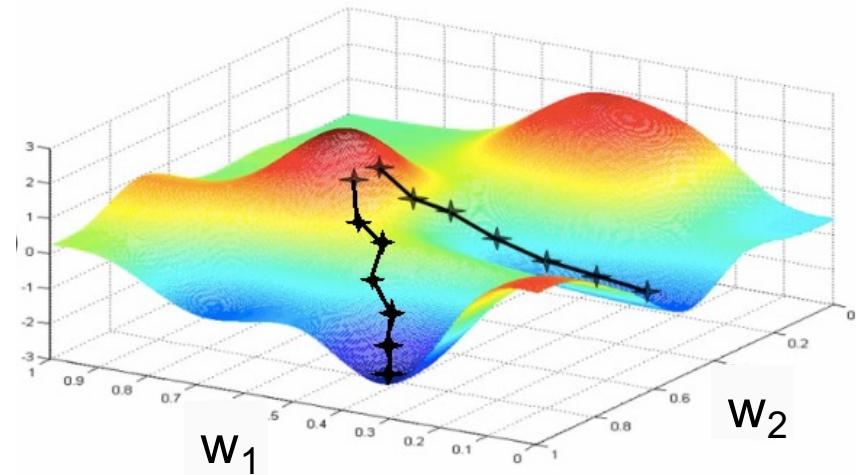
Training of multi-layer networks

- Find network weights to minimize the error between true and estimated labels of training examples:

$$E(\mathbf{w}) = \sum_{j=1}^N (y_j - f_{\mathbf{w}}(\mathbf{x}_j))^2$$

- Update weights by gradient descent:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial E}{\partial \mathbf{w}}$$



Training of multi-layer networks

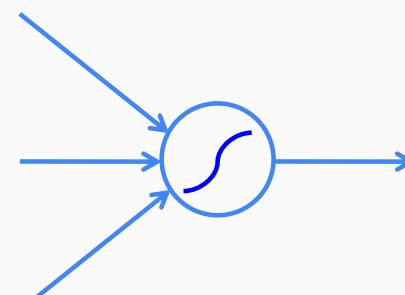
- Find network weights to minimize the error between true and estimated labels of training examples:

$$E(\mathbf{w}) = \sum_{j=1}^N (y_j - f_{\mathbf{w}}(\mathbf{x}_j))^2$$

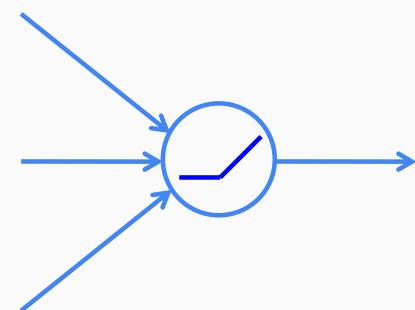
- Update weights by gradient descent:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial E}{\partial \mathbf{w}}$$

- This requires perceptrons with a differentiable nonlinearity



Sigmoid: $g(t) = \frac{1}{1 + e^{-t}}$



Rectified linear unit (ReLU):
 $g(t) = \max(0, t)$

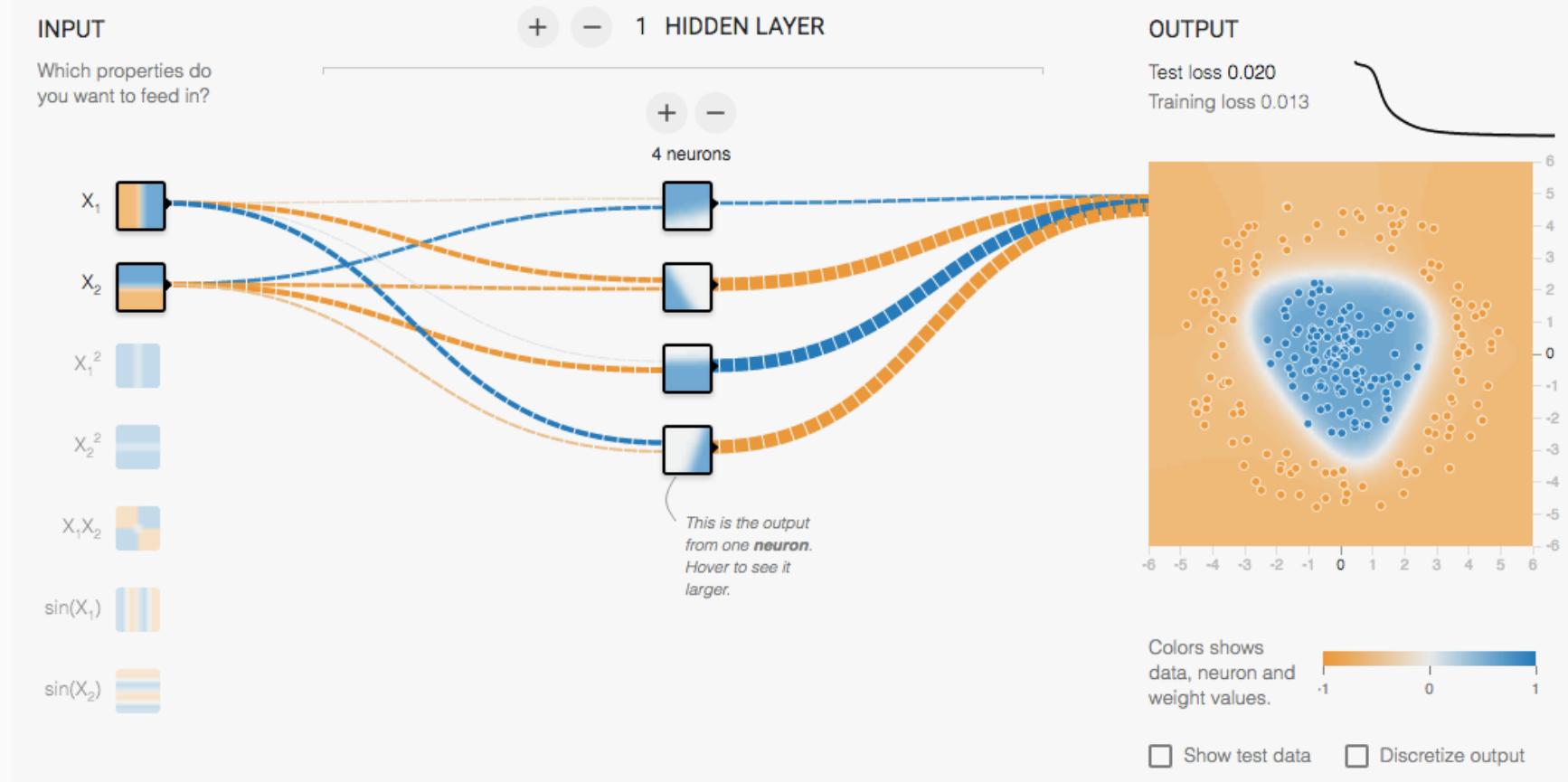
Training of multi-layer networks

- Find network weights to minimize the error between true and estimated labels of training examples:

$$E(\mathbf{w}) = \sum_{j=1}^N (y_j - f_{\mathbf{w}}(\mathbf{x}_j))^2$$

- Update weights by gradient descent: $\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial E}{\partial \mathbf{w}}$
- **Back-propagation:**
gradients are computed from output to input layers and combined using chain rule
- **Stochastic gradient descent:**
compute the weight update w.r.t. one training example (or a small batch of examples) at a time, cycle through training examples in random order in multiple epochs

Multi-Layer Network Demo



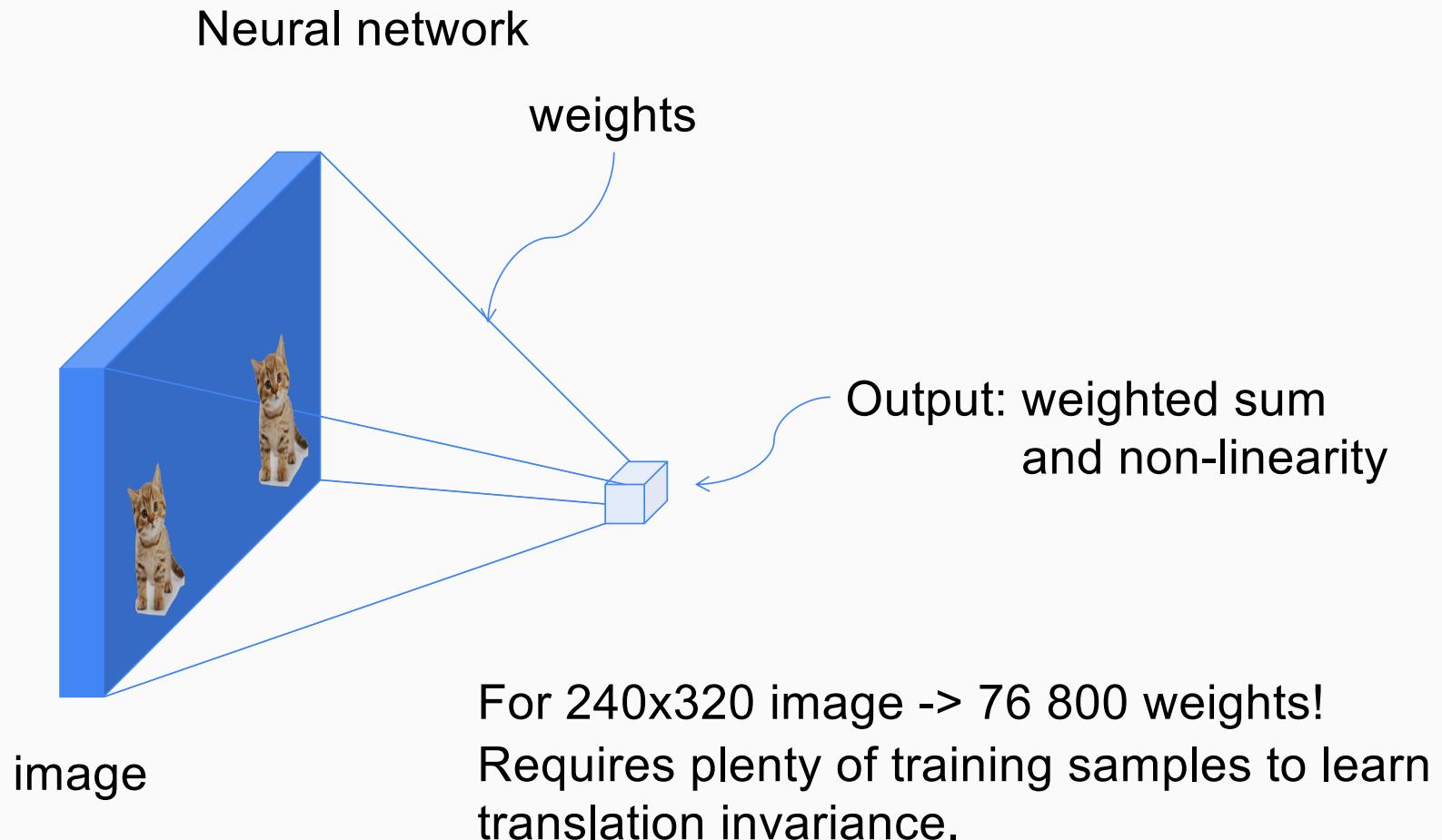
<http://playground.tensorflow.org/>

Neural networks: Pros and cons

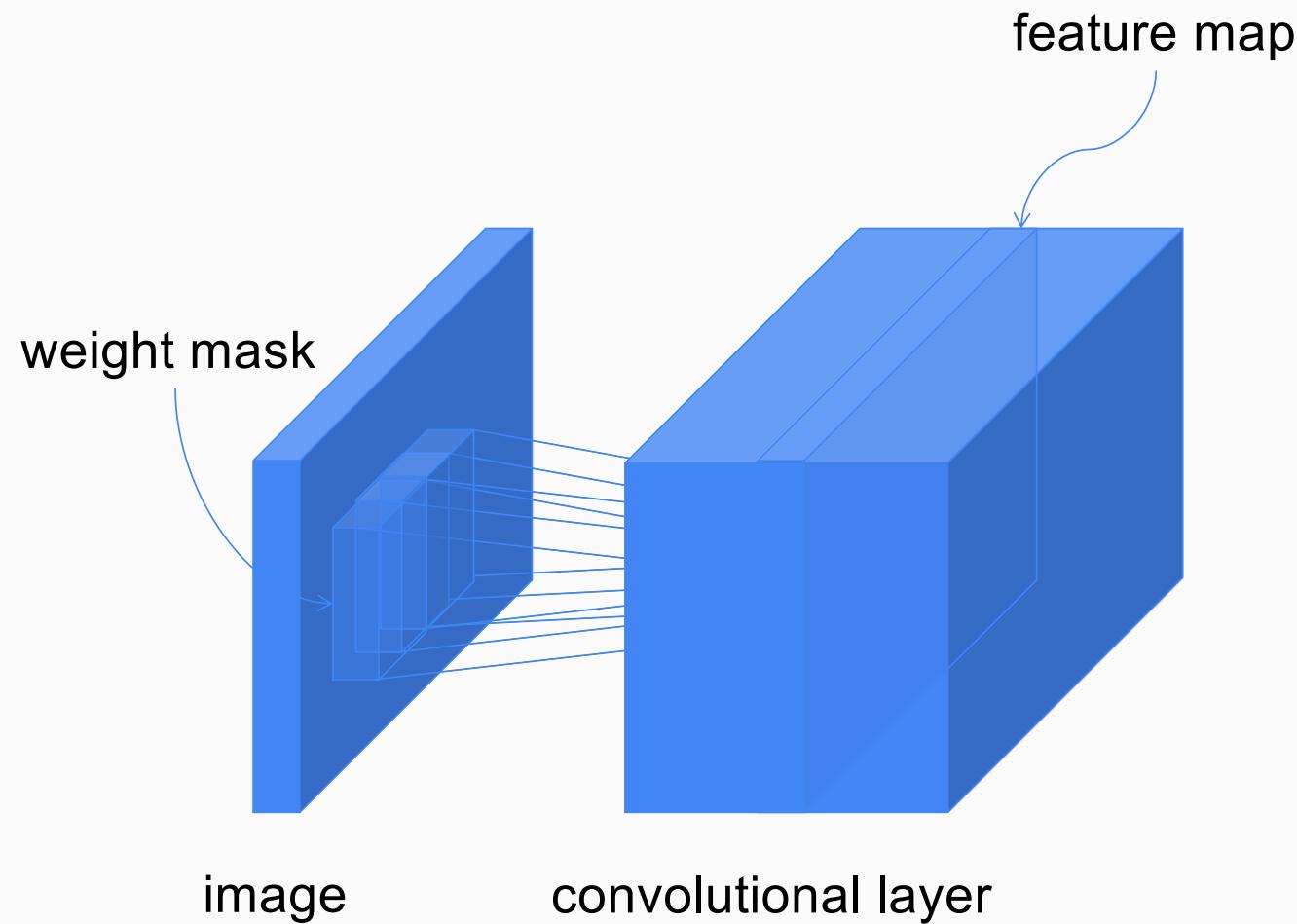
- Pros
 - Flexible and general function approximation framework
 - Can build extremely powerful models by adding more layers
- Cons
 - Hard to analyze theoretically (e.g., training is prone to local optima)
 - Huge amount of training data and computing power may be required to get good performance
 - The space of implementation choices is huge (network architectures, parameters)

Convolutional neural network

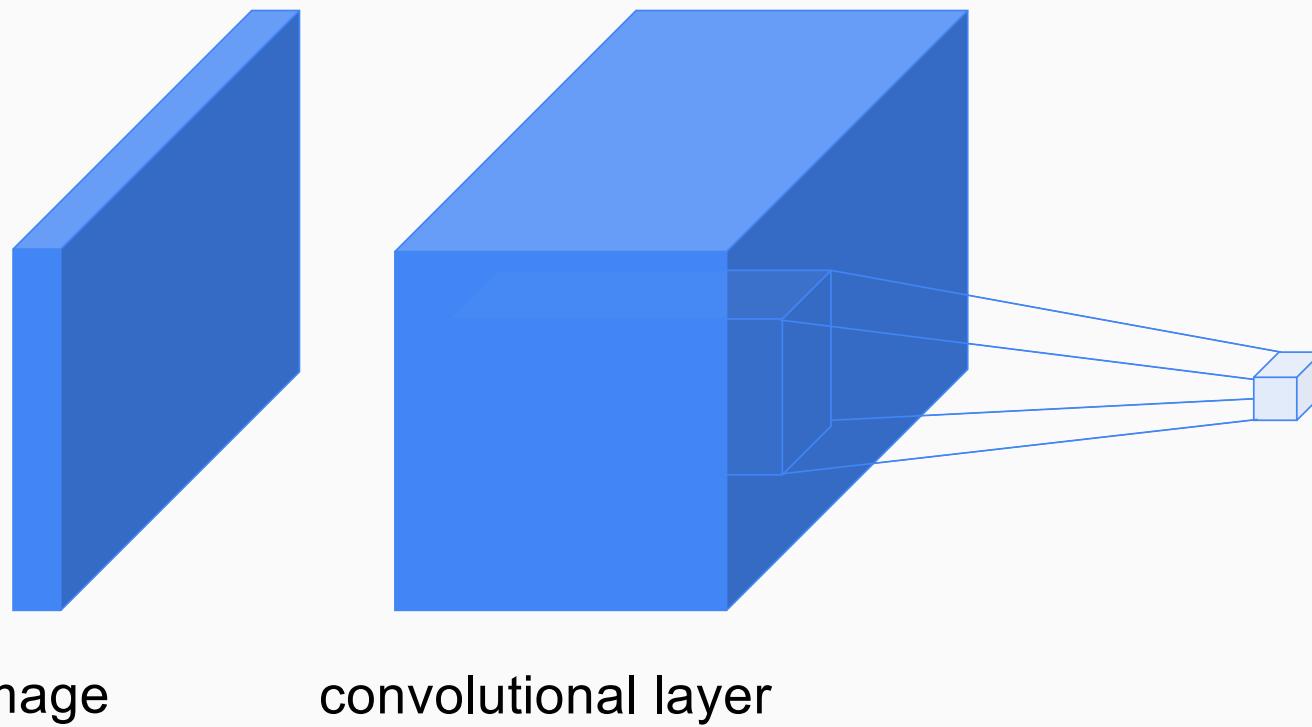
Neural networks for images



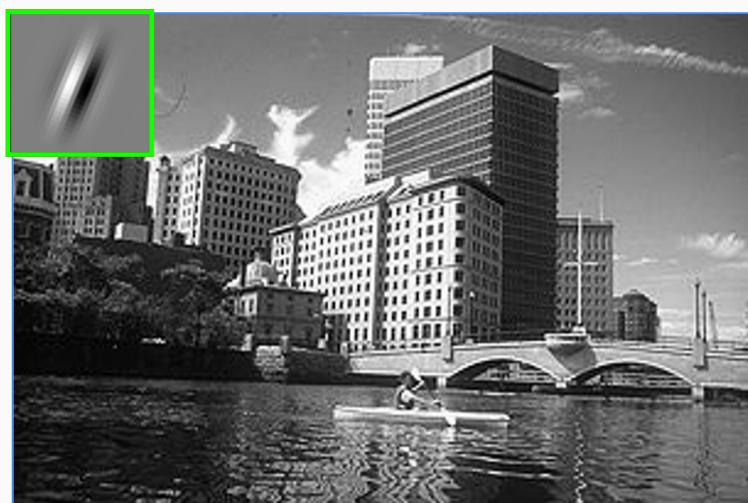
Neural networks for images



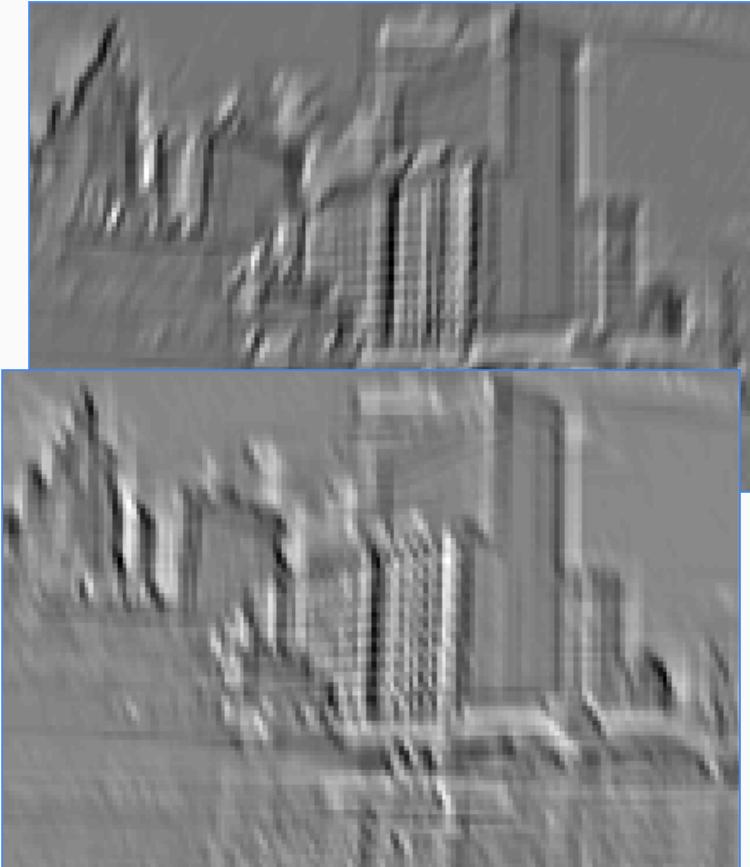
Neural networks for images



Convolution as feature extraction



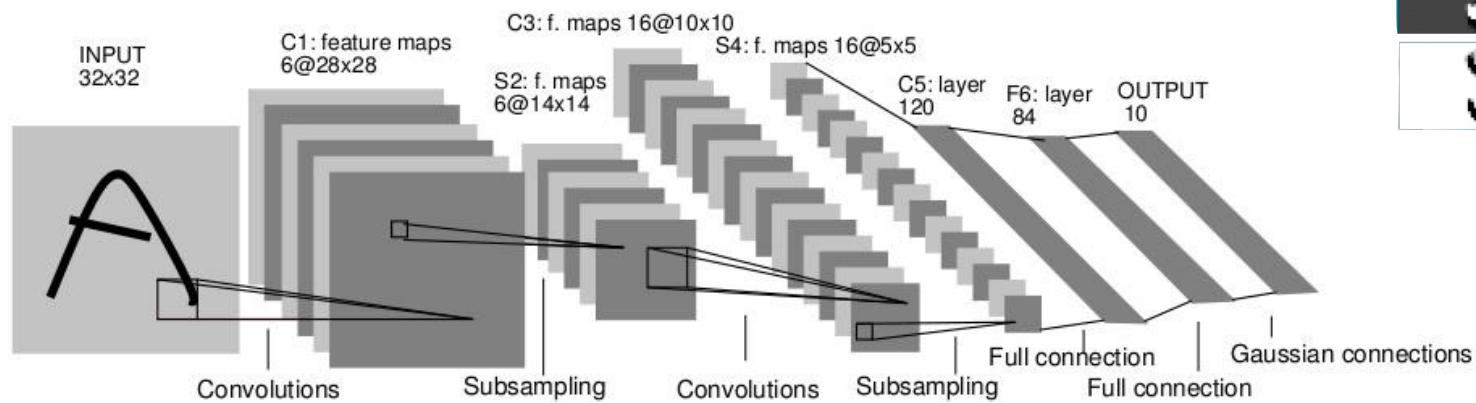
Input



Feature Map

Convolutional Neural Networks

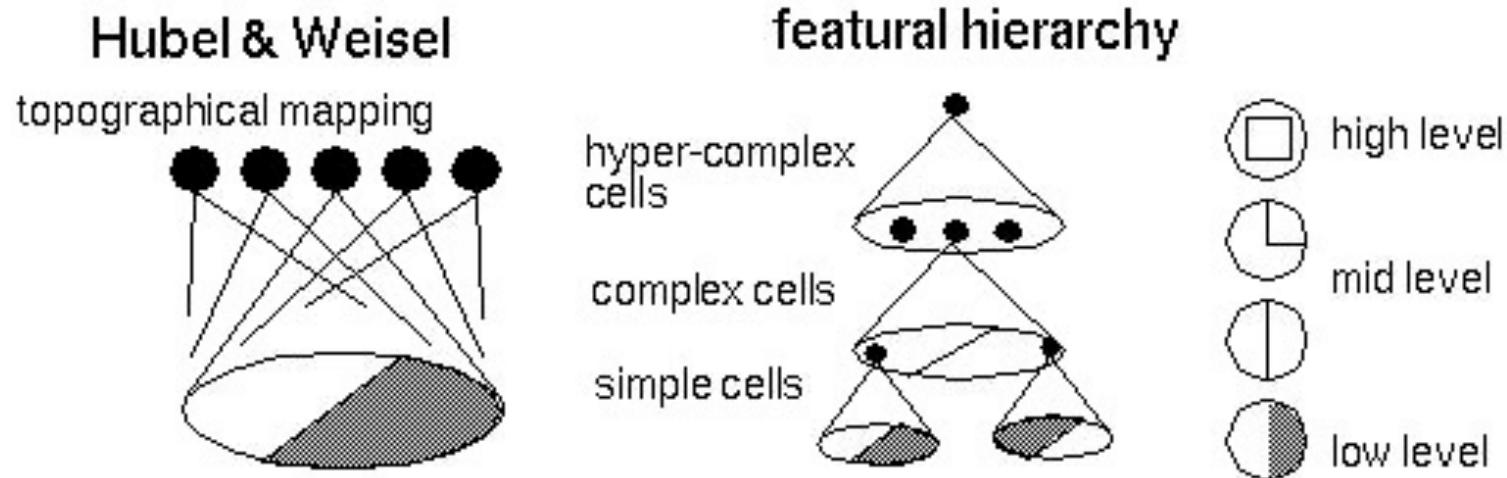
- Neural network with specialized connectivity structure
- Stack multiple stages of feature extractors
- Higher stages compute more global, more invariant features
- Classification layer at the end



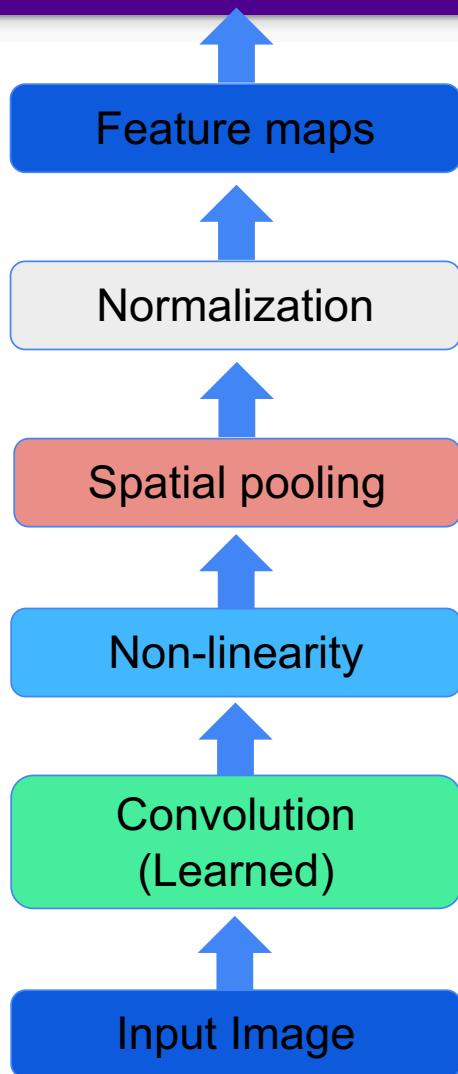
Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, [Gradient-based learning applied to document recognition](#), Proc. IEEE 86(11): 2278–2324, 1998.

Biological inspiration

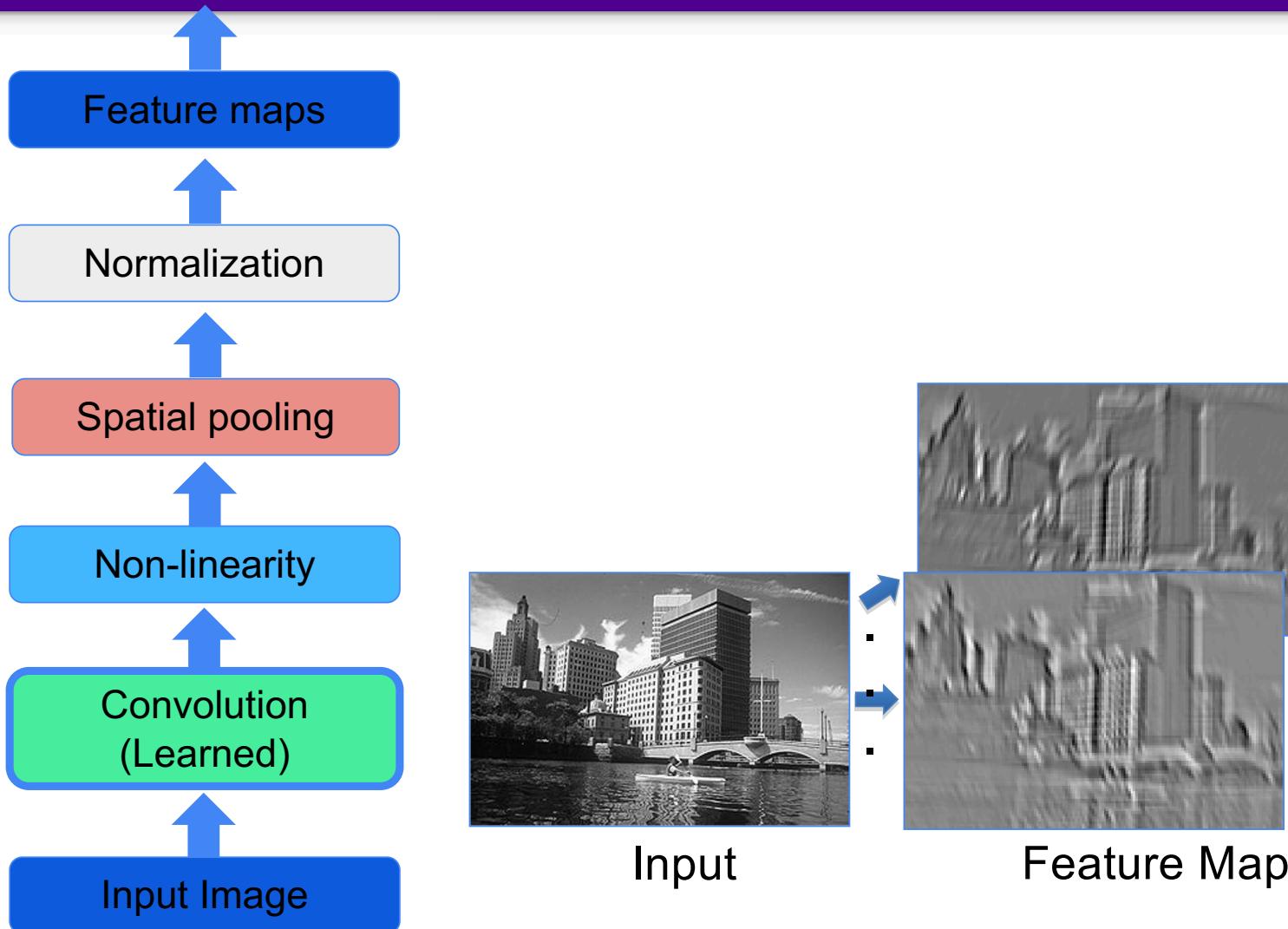
- D. Hubel and T. Wiesel (1959, 1962, Nobel Prize 1981)
 - Visual cortex consists of a hierarchy of simple, complex, and hyper-complex cells



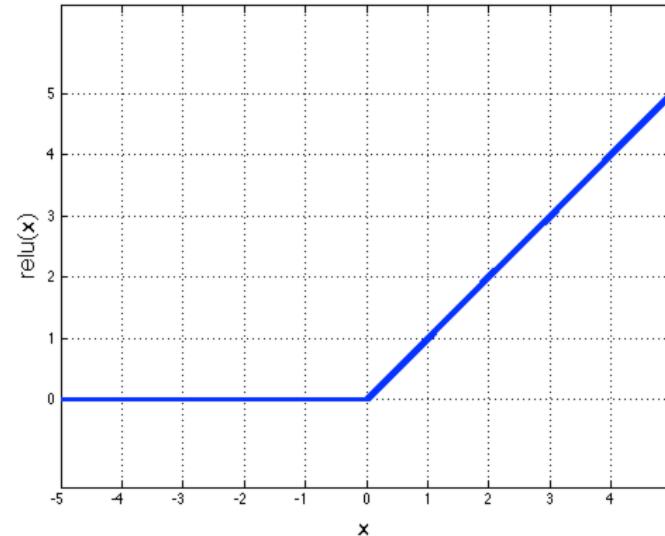
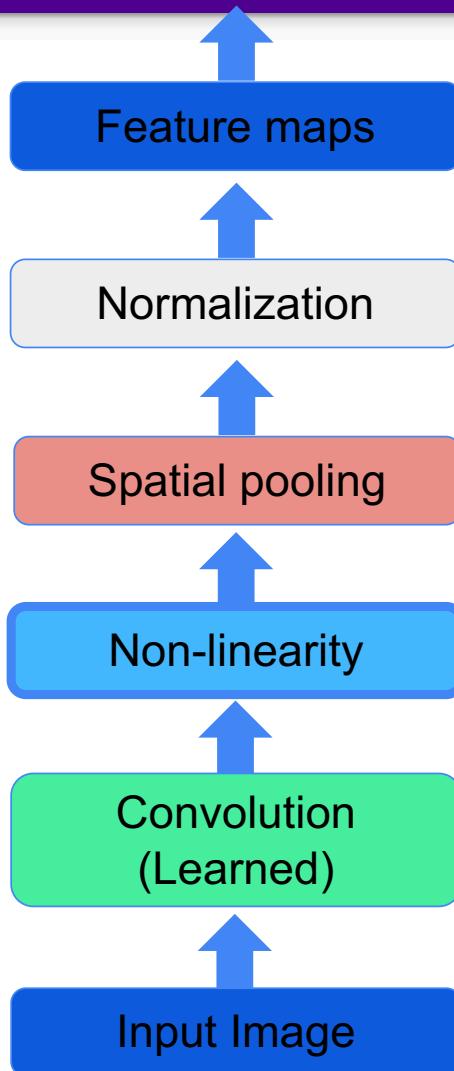
Convolutional Neural Networks



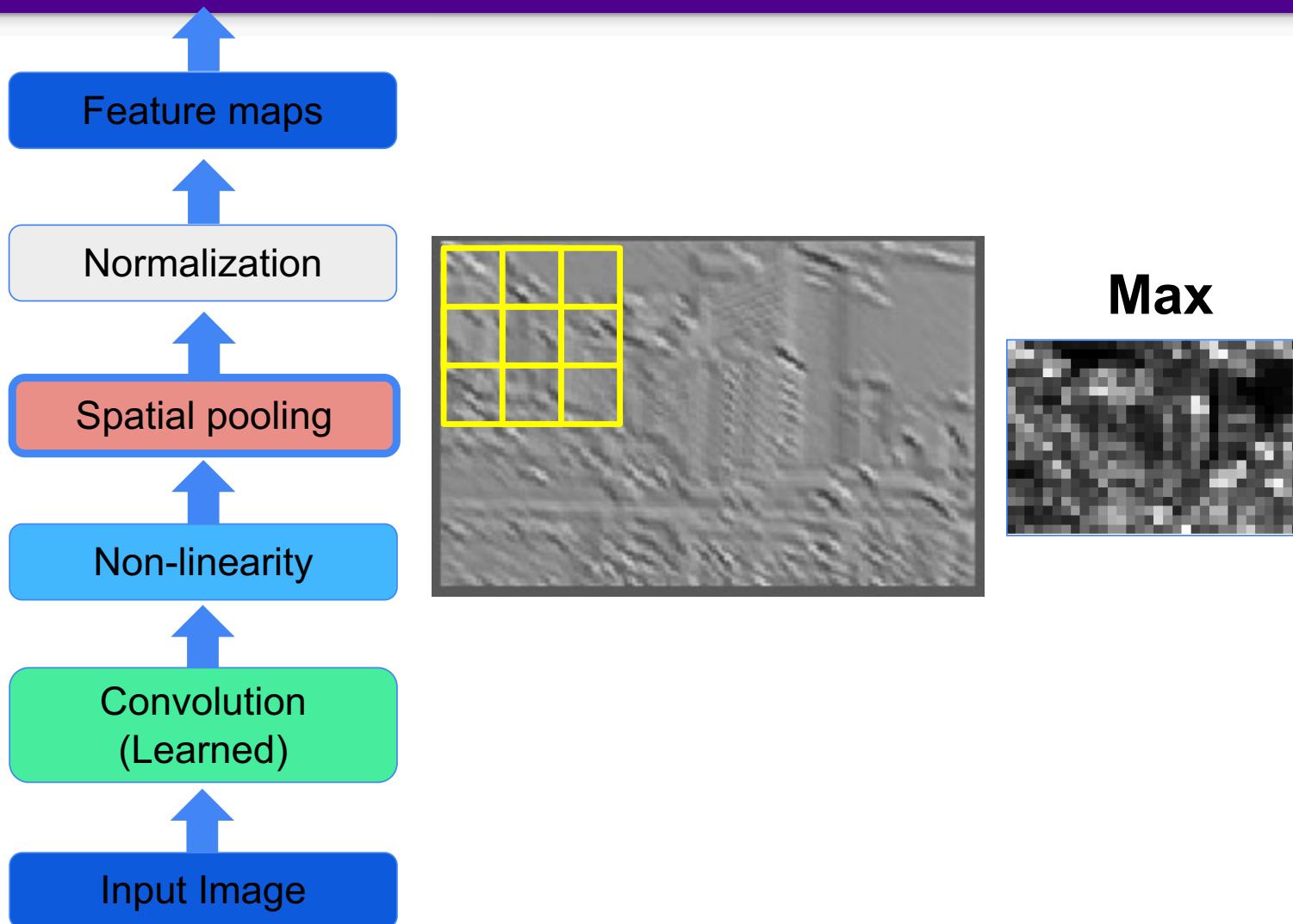
Convolutional Neural Networks



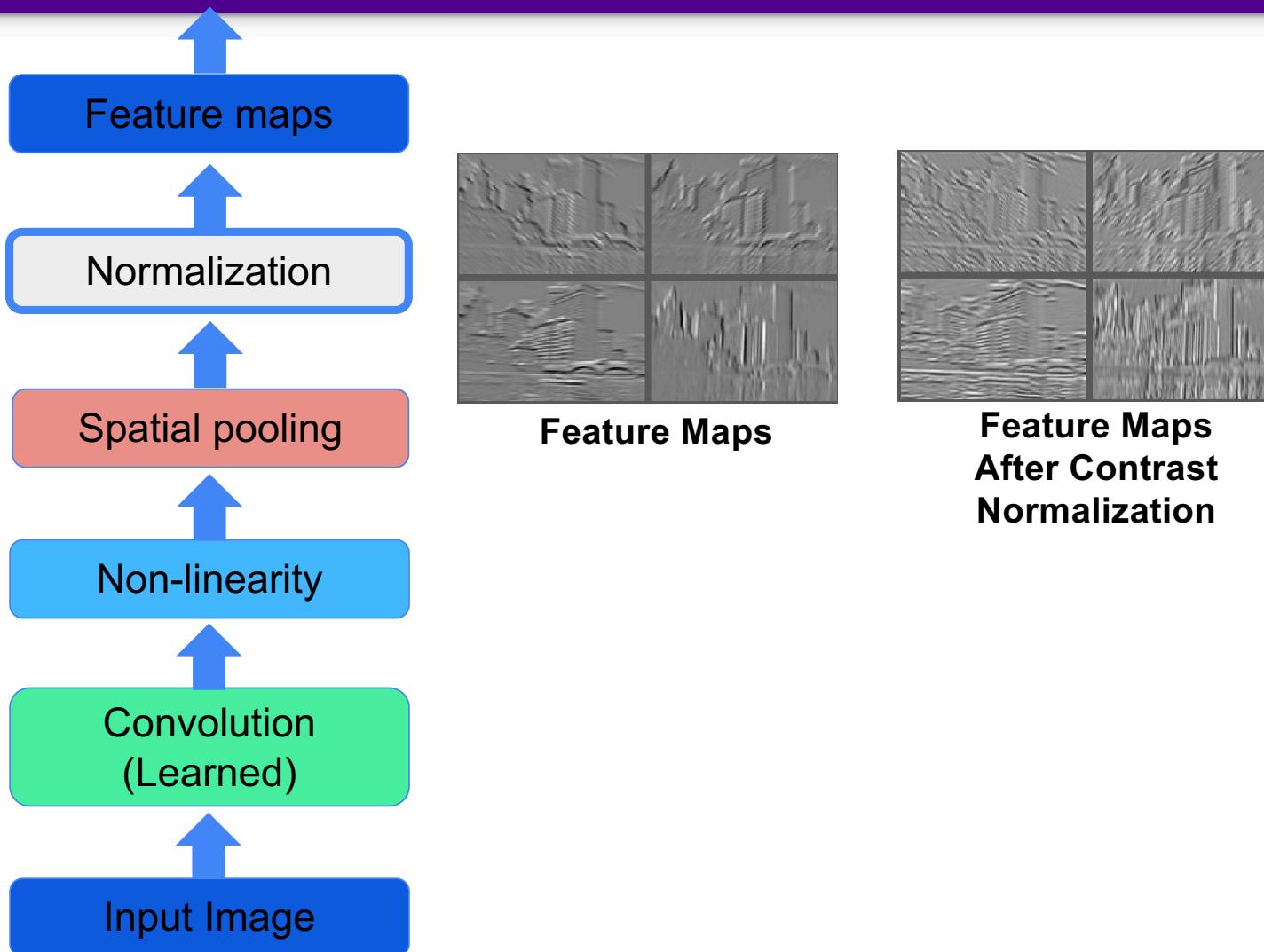
Convolutional Neural Networks



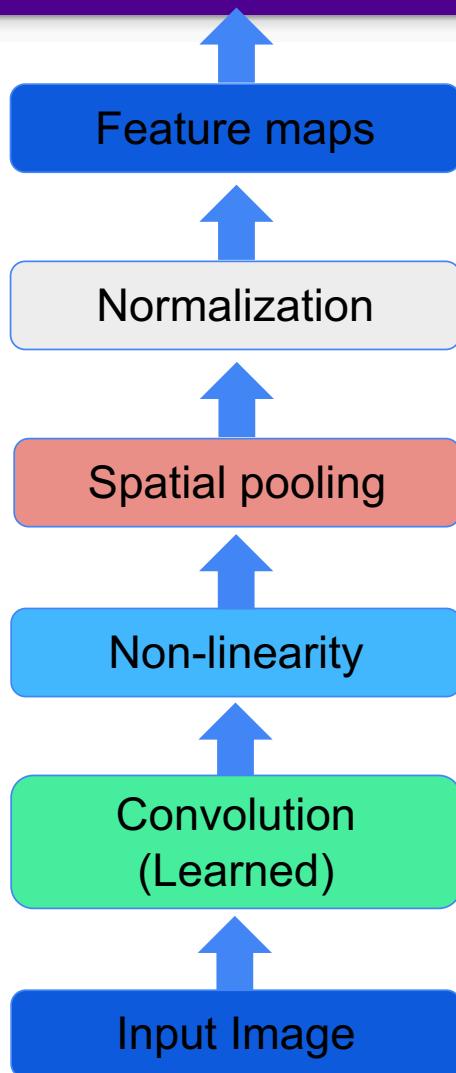
Convolutional Neural Networks



Convolutional Neural Networks

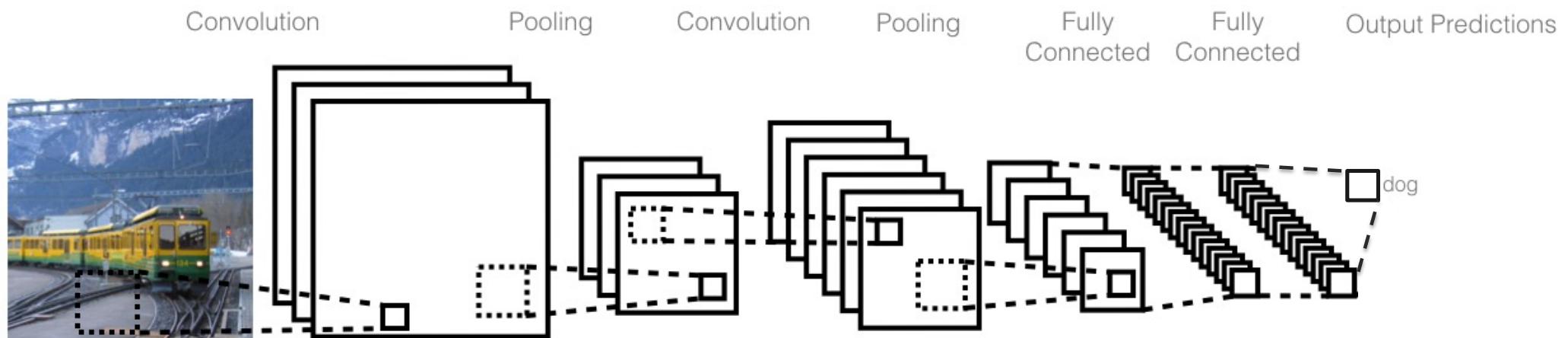


Convolutional Neural Networks

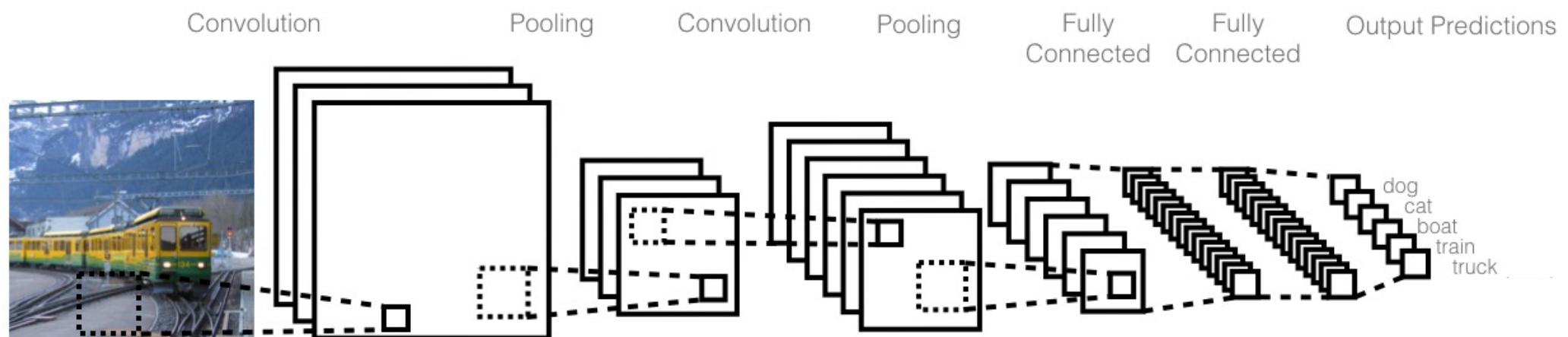


Convolutional filters are trained in a supervised manner by back-propagating classification error

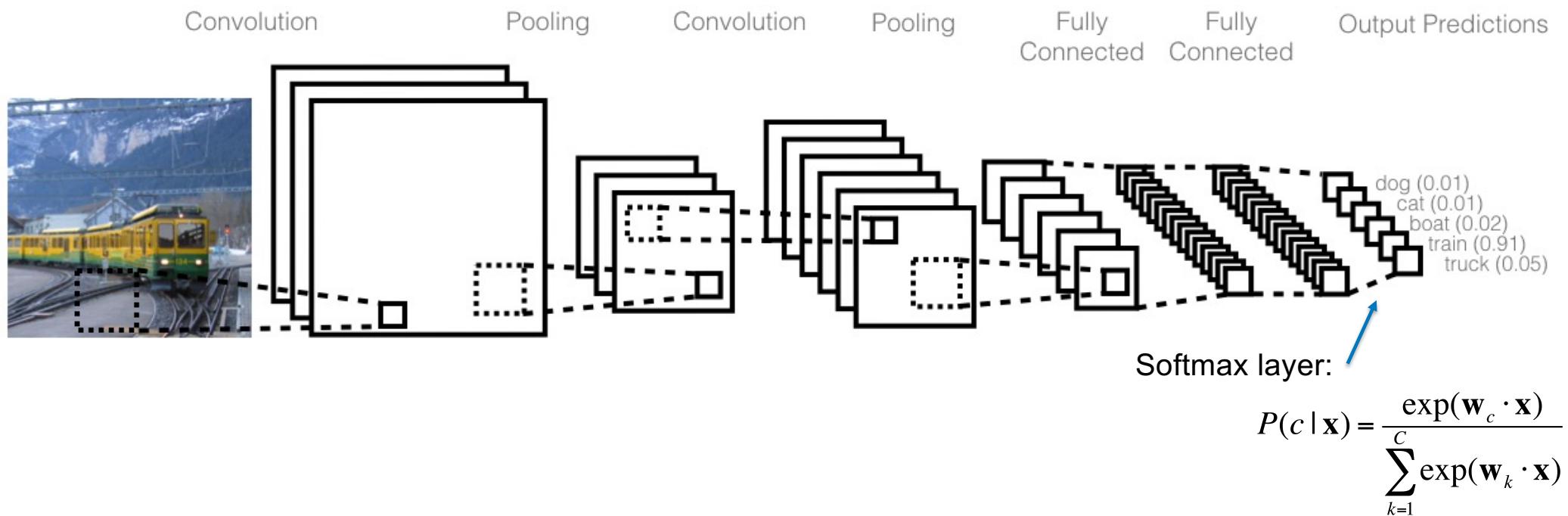
Simplified architecture



Simplified architecture

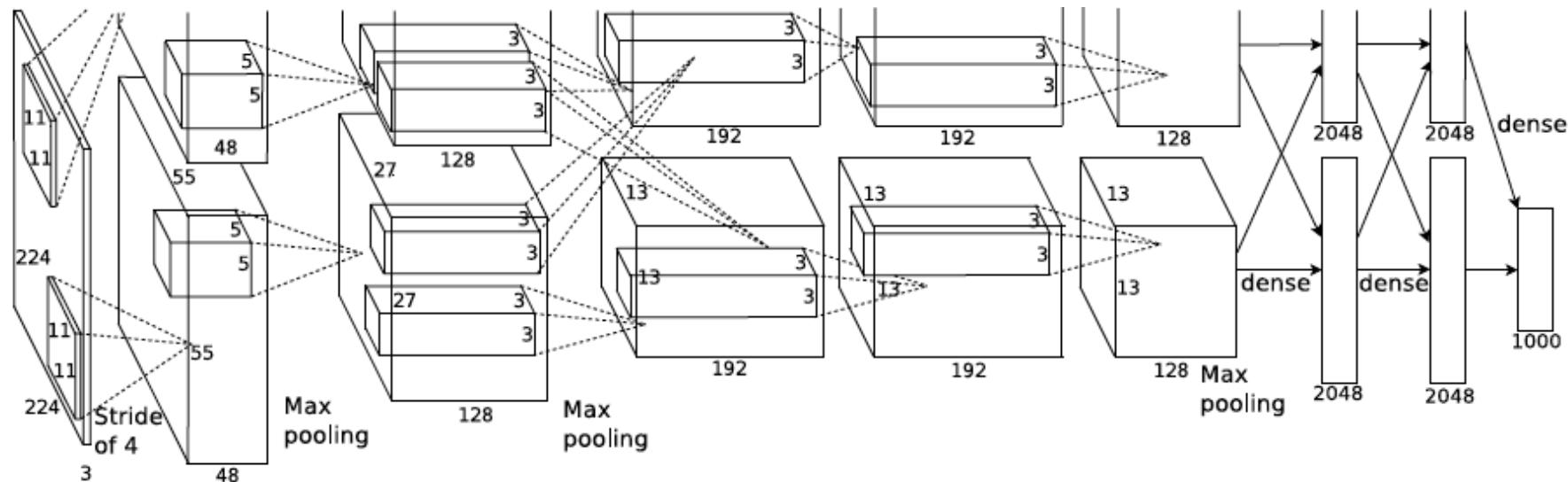


Simplified architecture



AlexNet

- Similar framework to LeCun'98, but:
 - Larger model (7 hidden layers, 650 000 units, 60 000 000 parameters)
 - More data (10^6 vs. 10^3 images)
 - GPU implementation (50x speedup over CPU, trained on two GPUs for a week)



A. Krizhevsky, I. Sutskever, and G. Hinton, [ImageNet Classification with Deep Convolutional Neural Networks](#), NIPS 2012

Refinement of AlexNet architecture

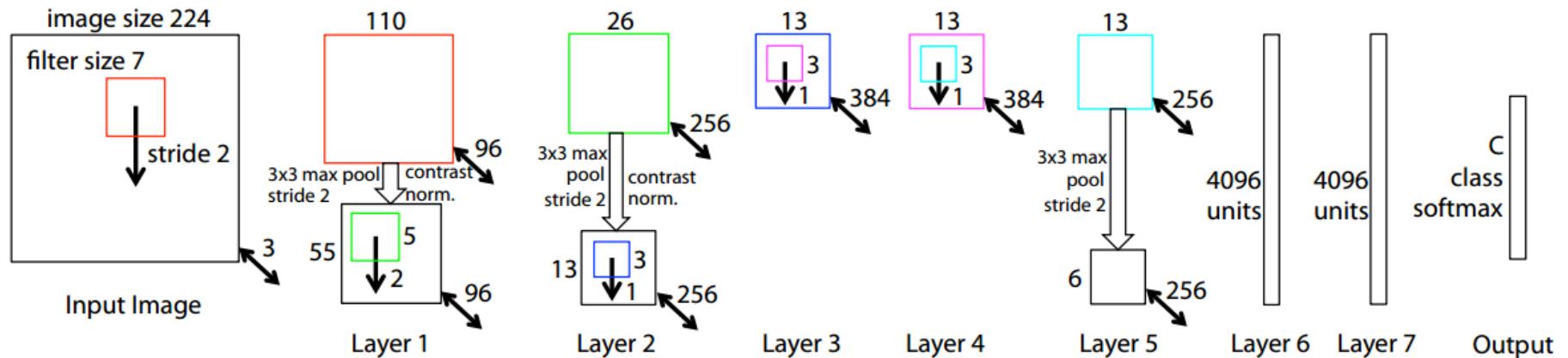
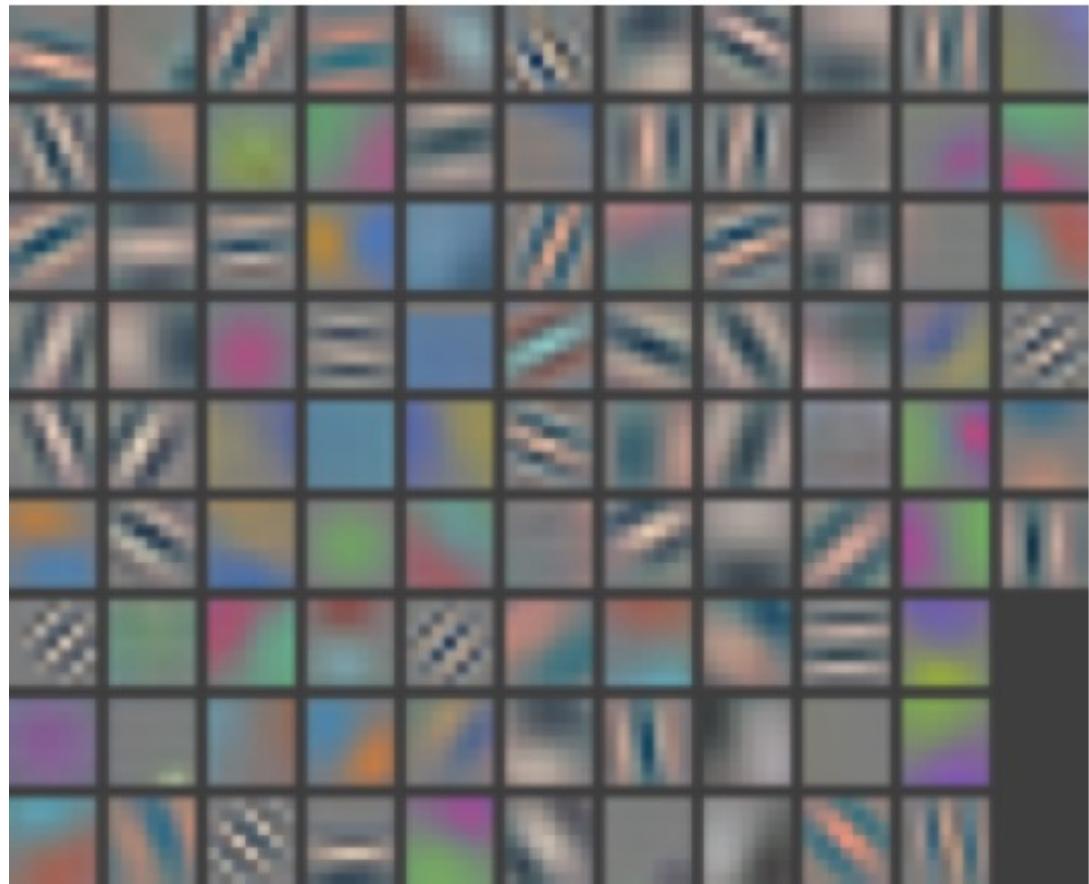
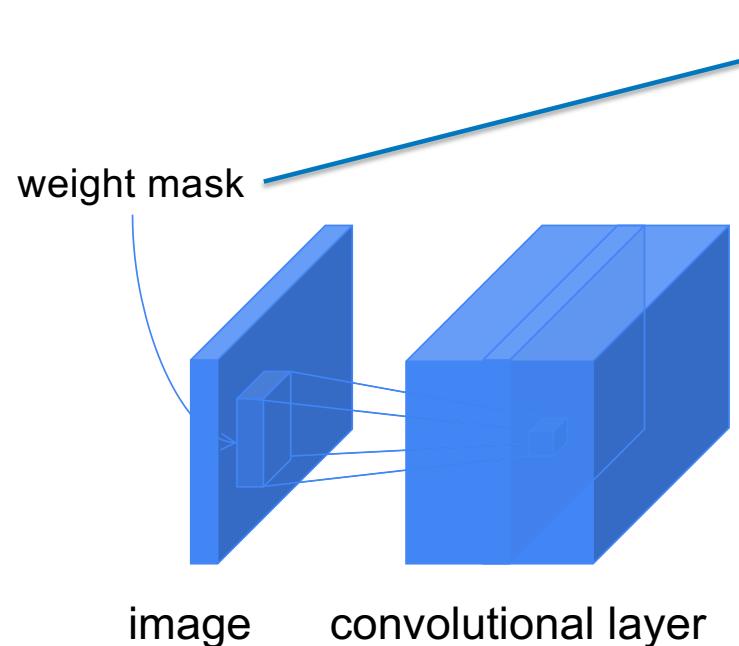


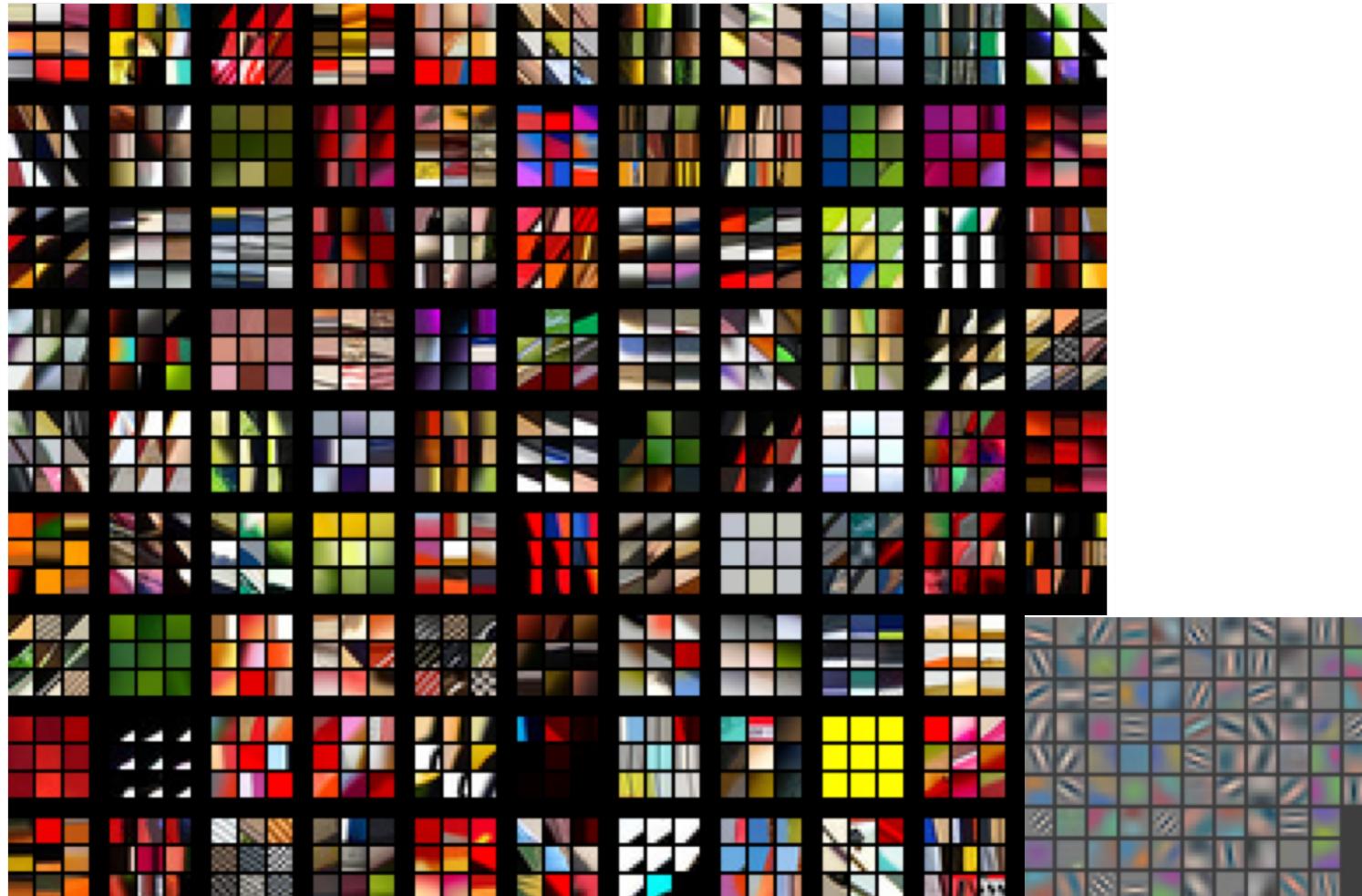
Figure 3. Architecture of our 8 layer convnet model. A 224 by 224 crop of an image (with 3 color planes) is presented as the input. This is convolved with 96 different 1st layer filters (red), each of size 7 by 7, using a stride of 2 in both x and y. The resulting feature maps are then: (i) passed through a rectified linear function (not shown), (ii) pooled (max within 3x3 regions, using stride 2) and (iii) contrast normalized across feature maps to give 96 different 55 by 55 element feature maps. Similar operations are repeated in layers 2,3,4,5. The last two layers are fully connected, taking features from the top convolutional layer as input in vector form ($6 \cdot 6 \cdot 256 = 9216$ dimensions). The final layer is a C -way softmax function, C being the number of classes. All filters and feature maps are square in shape.

Layer 1 Filters

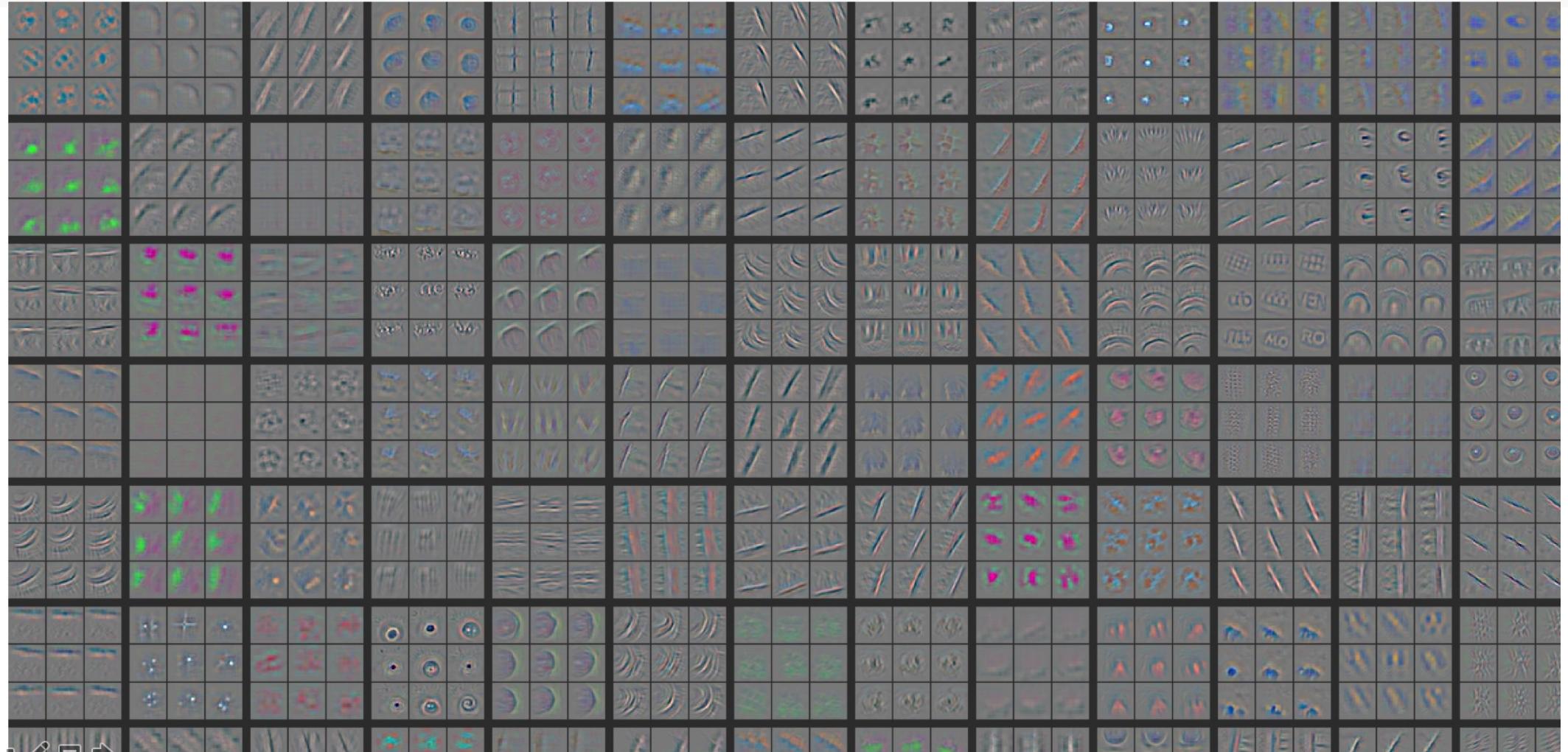


M. Zeiler and R. Fergus, [Visualizing and Understanding Convolutional Networks](#), ECCV 2014 (Best Paper Award winner)

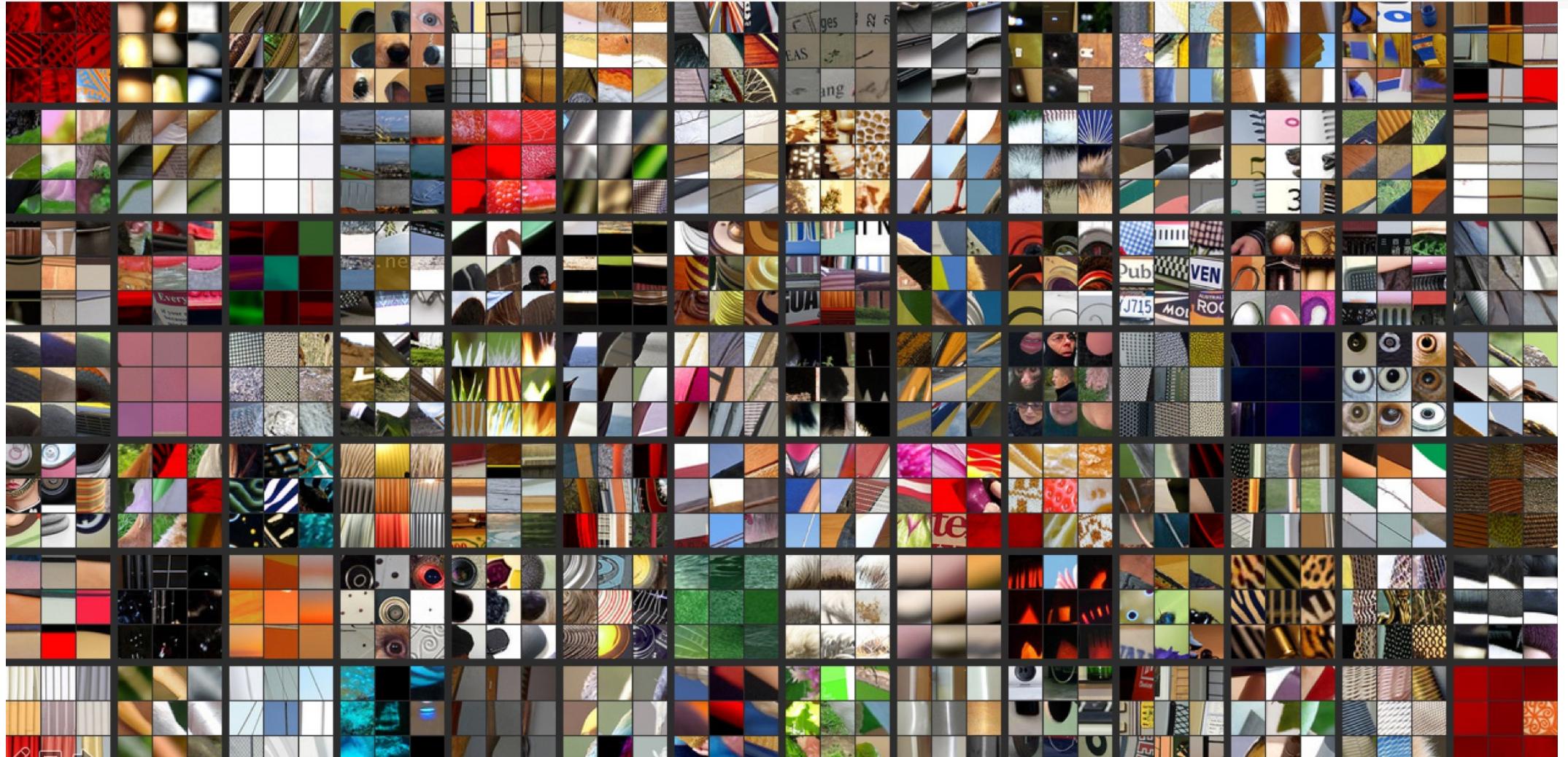
Layer 1: Top-9 Patches (patches that give maximal activation)



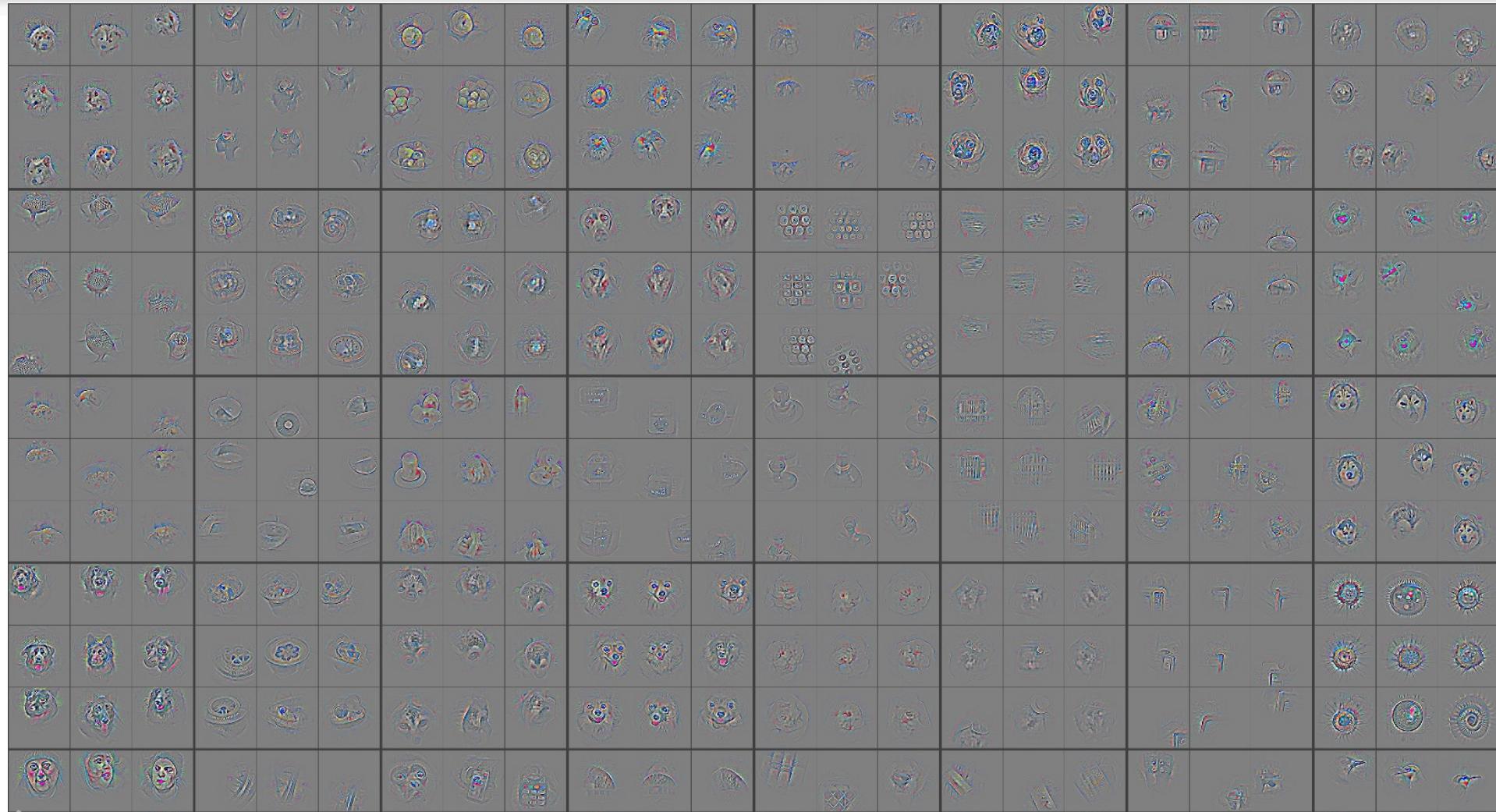
Layer 2: Filters



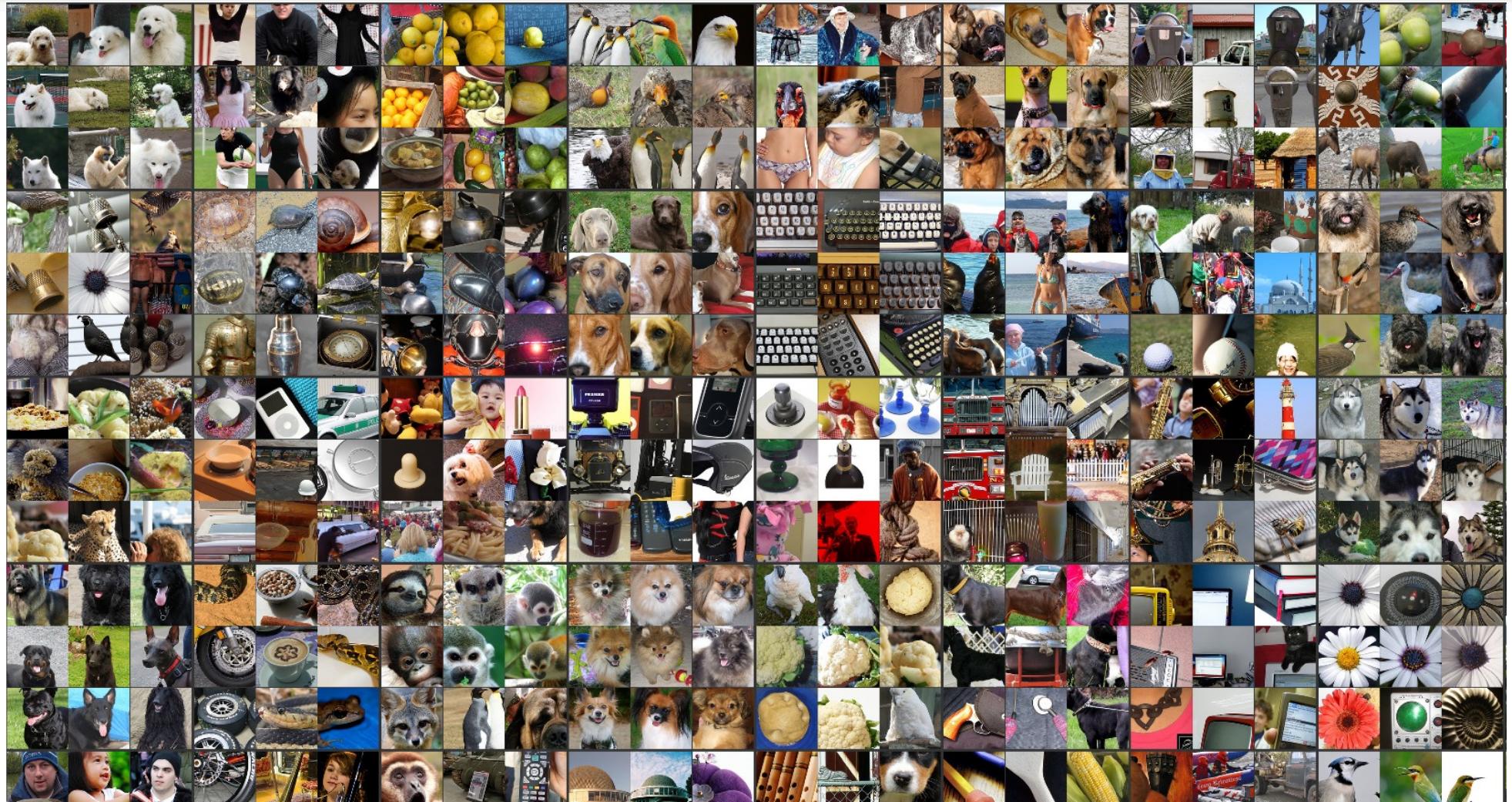
Layer 2: Top-9 Patches



Layer 5: Filters

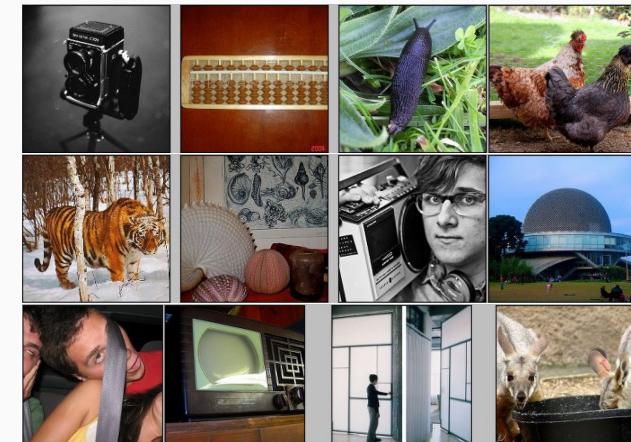


Layer 5: Top-9 Patches



ImageNet Challenge

- Approx. 14 million labelled images, 20k classes
- Image gathered from Internet
- Human labels via Amazon Mturk
- Challenge:
1.2 million training images and 1000 classes
- Metrics: top-5 and top-1 errors



www.image-net.org/challenges/LSVRC/

Typical results



ImageNet Challenge 2012-2013

Team	Year	Place	Error (top-5)	External data
SuperVision – Toronto (7 layers)	2012	-	16.4%	no
SuperVision	2012	1st	15.3%	ImageNet 22k
Clarifai – NYU (7 layers)	2013	-	11.7%	no
Clarifai	2013	1st	11.2%	ImageNet 22k

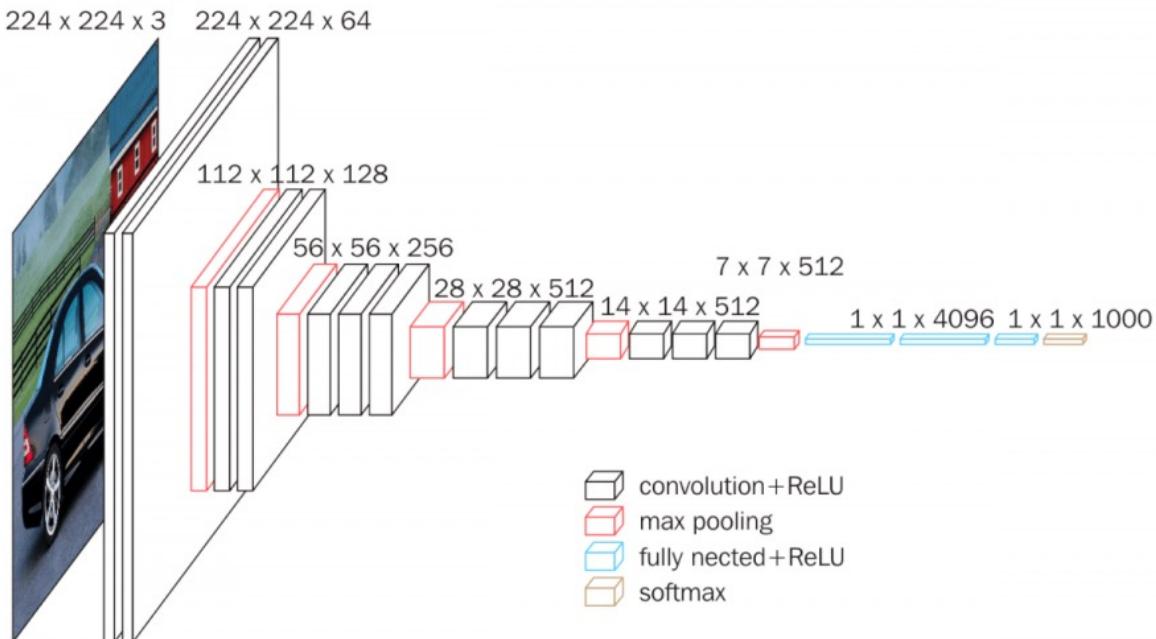
<http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/>

ImageNet Challenge 2014

Team	Year	Place	Error (top-5)	External data
SuperVision – Toronto (7 layers)	2012	-	16.4%	no
SuperVision	2012	1st	15.3%	ImageNet 22k
Clarifai – NYU (7 layers)	2013	-	11.7%	no
Clarifai	2013	1st	11.2%	ImageNet 22k
VGG – Oxford (16 layers)	2014	2nd	7.32%	no
GoogLeNet (19 layers)	2014	1st	6.67%	no
Human expert*			5.1%	

<http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/>

VGG-16 (Simonyan & Zisserman 2014)



138 Million parameters
Top-5 error in ImageNet 7%

What about after 2014?

AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)



Deep Residual Nets

Revolution of Depth

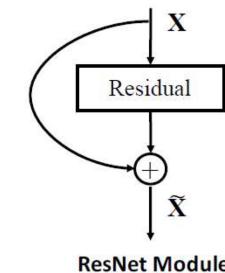
AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)



ResNet, 152 layers
(ILSVRC 2015)

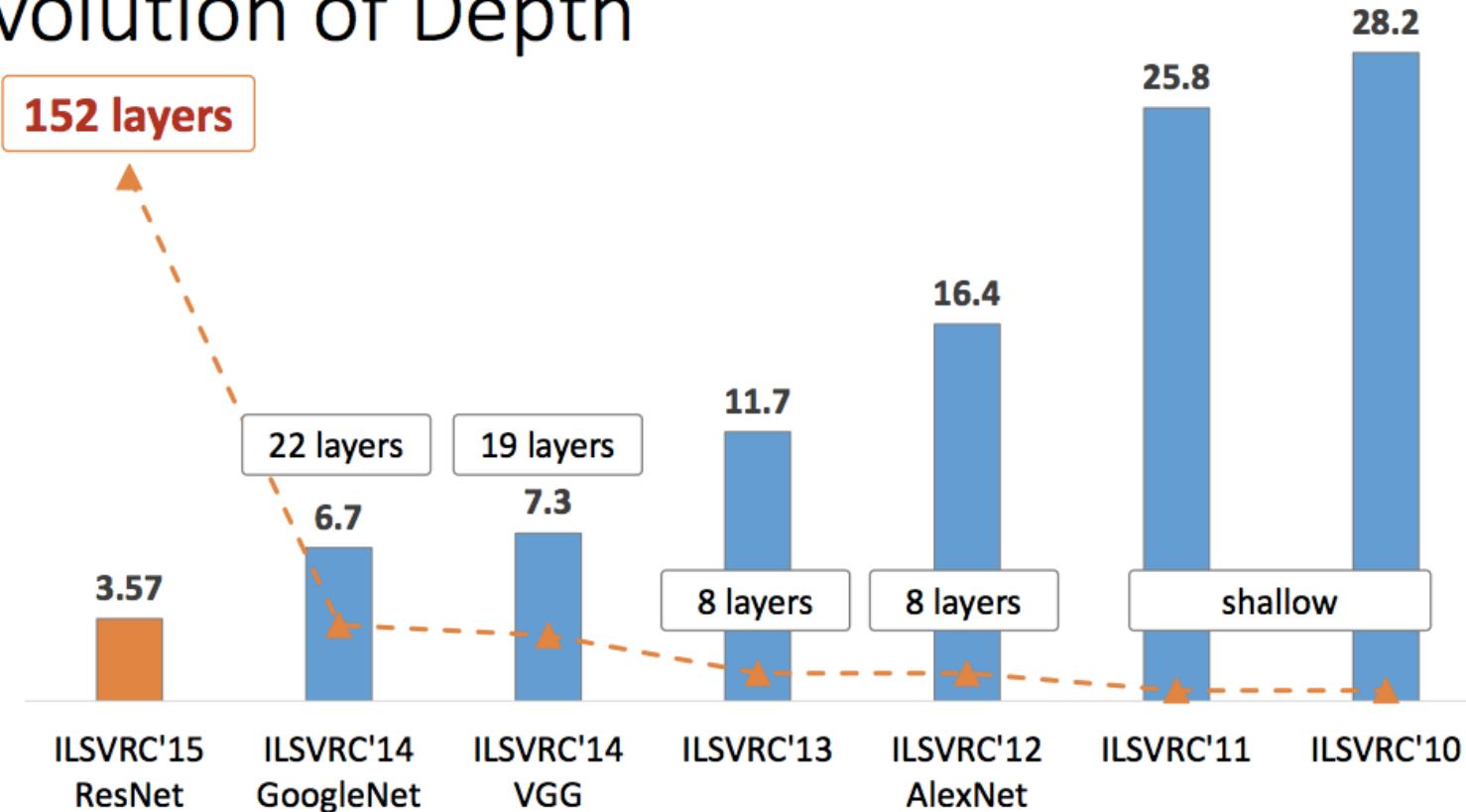


Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, [Deep Residual Learning for Image Recognition](#), arXiv 2015



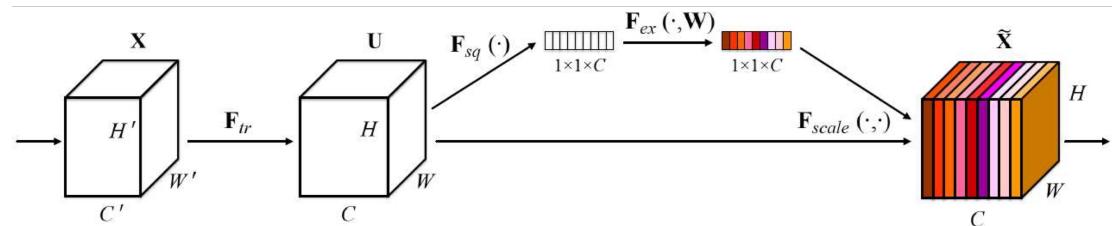
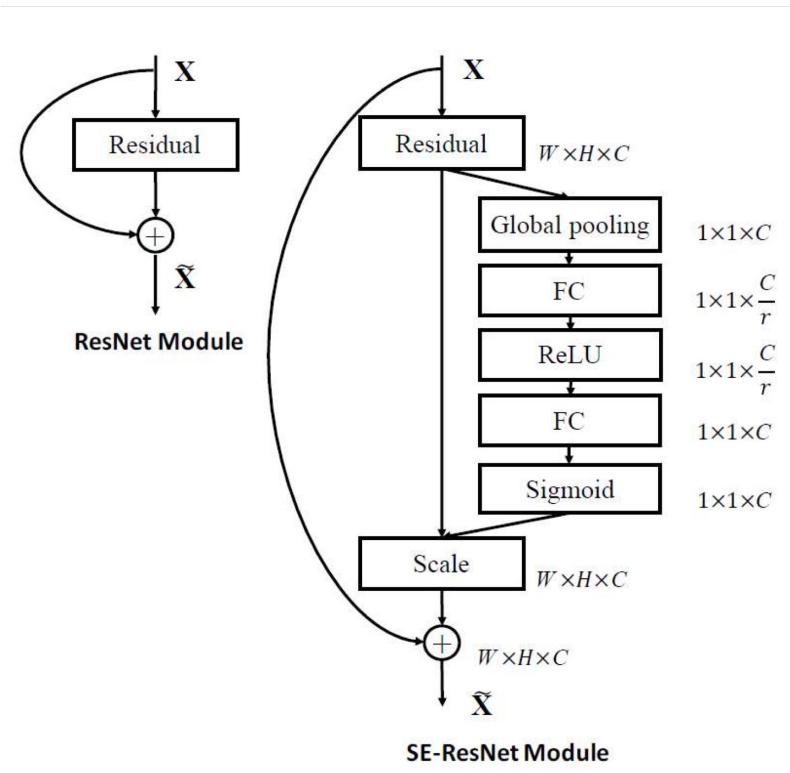
Deep Residual Nets

Revolution of Depth

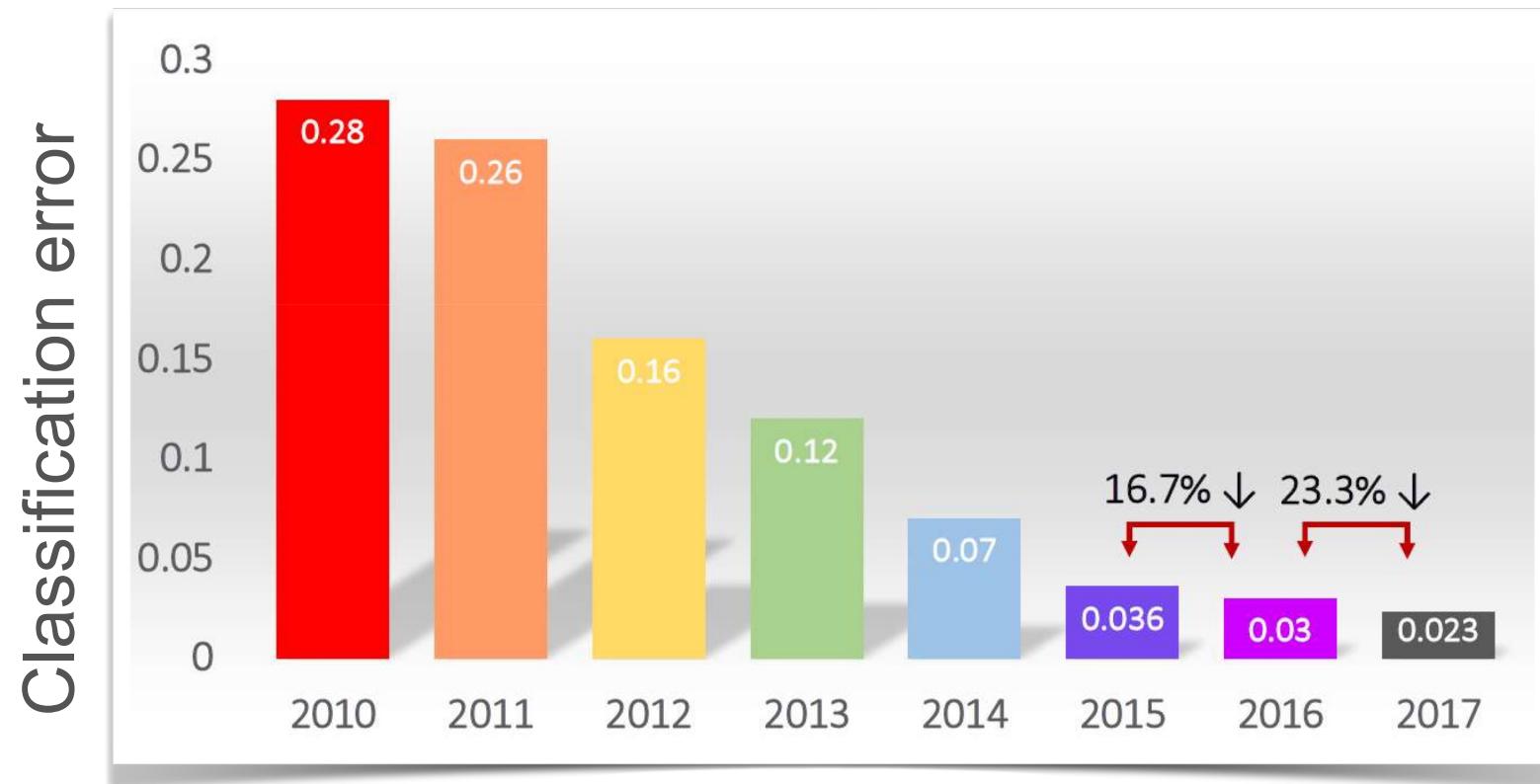


Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, [Deep Residual Learning for Image Recognition](#), arXiv 2015

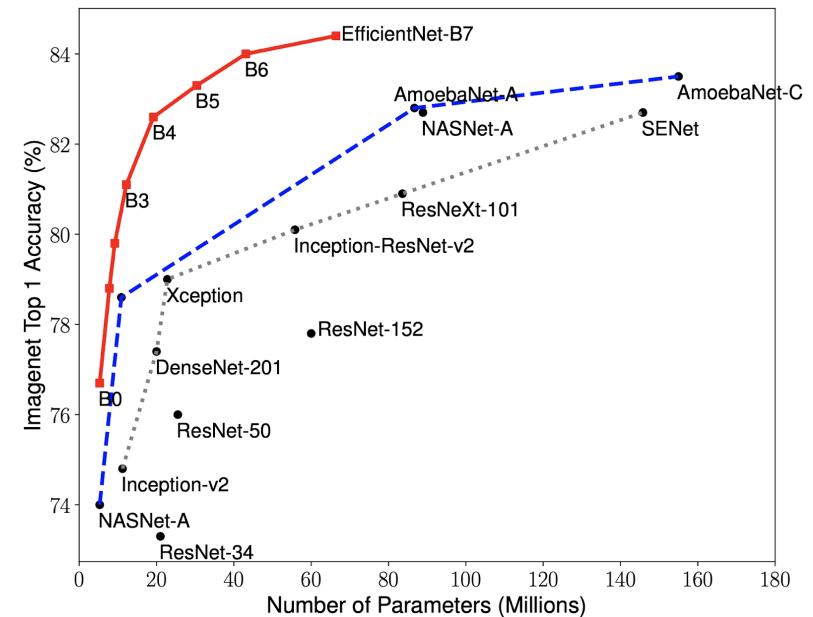
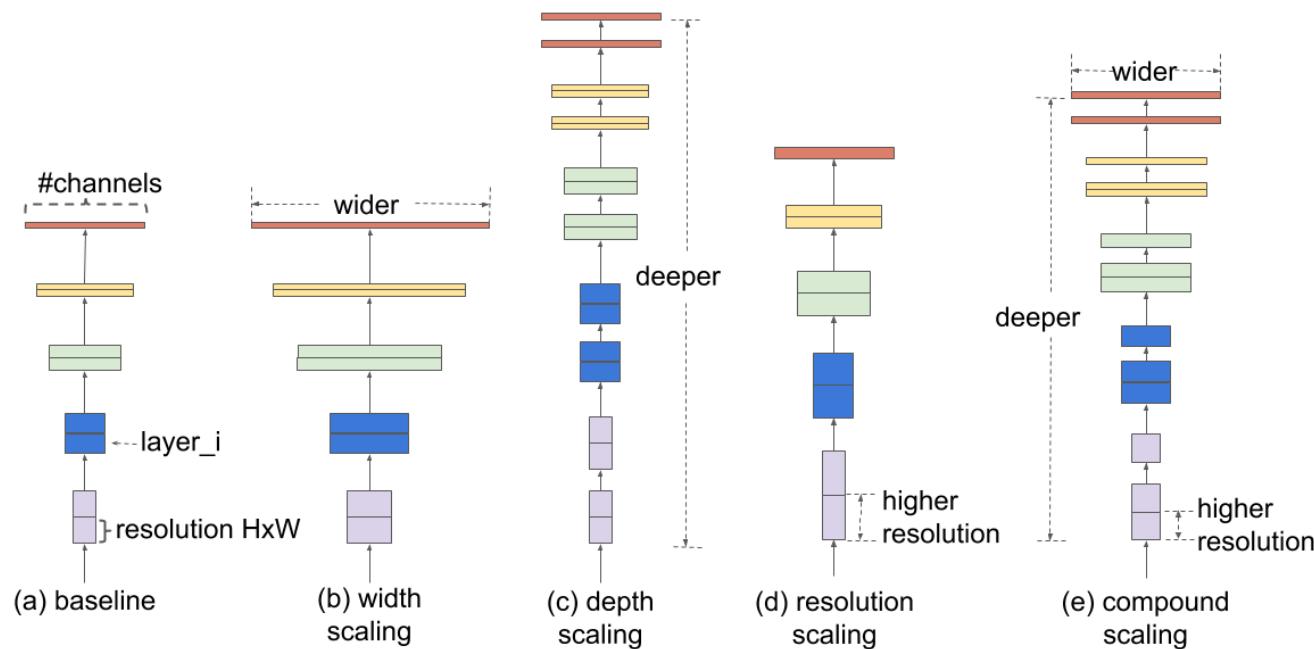
Squeeze & Excitation



ImageNet classification results (CLS)

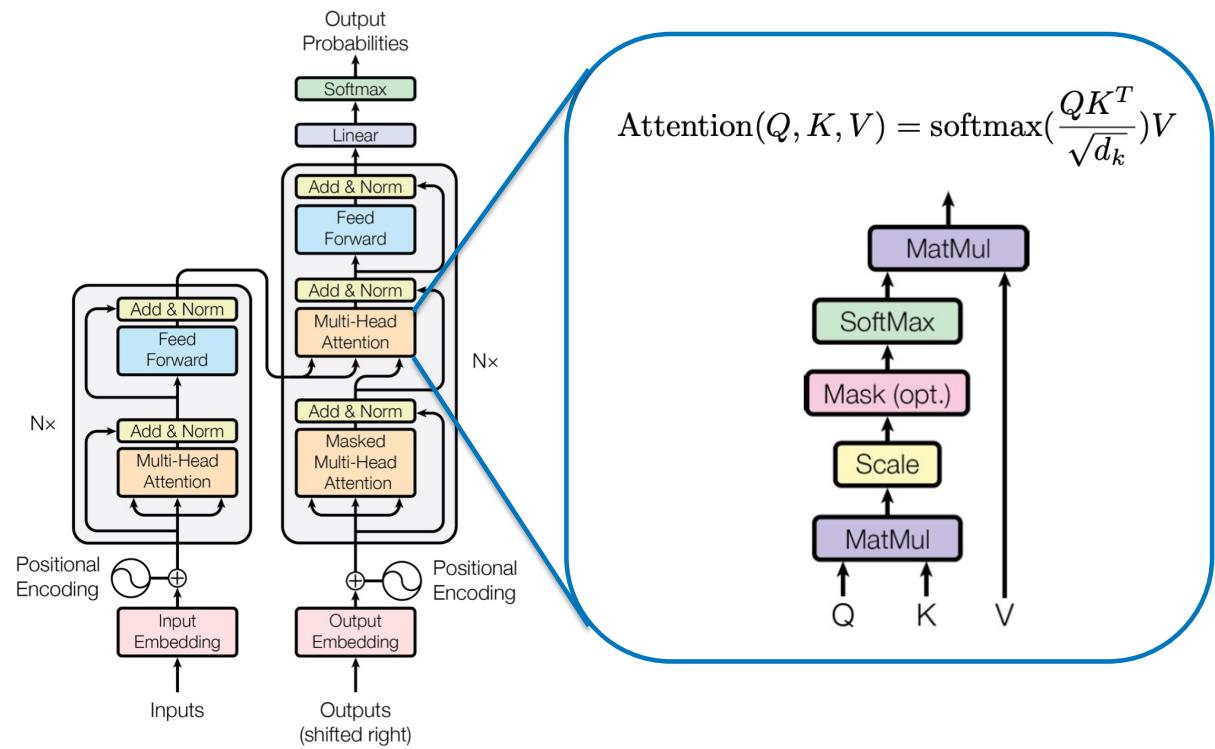


EfficientNet



Beyond Convolution Networks - Transformer

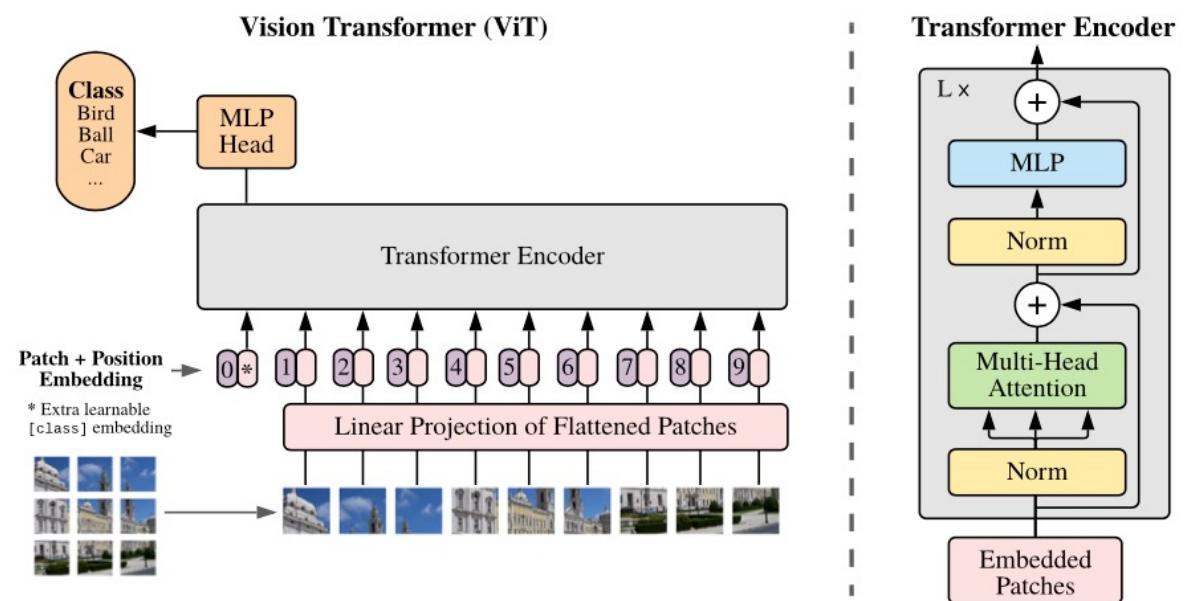
- Alternative to CNNs
- Transformers were originally developed for language processing (translation)
- Currently dominating approach in language modelling.
E.g. GPT-3 (175 billion parameters, trained with 45 TB of text data)
- Substantial and rapid development



Vaswani et al. Attention is all you need, NeurIPS 2017
Brown et al. Language models few-shot learners, NeurIPS 2020

Transformer

- Transformer models regularly used in vision tasks, particularly those with large amount of training data.
- Rapid development over 2021 and beyond.



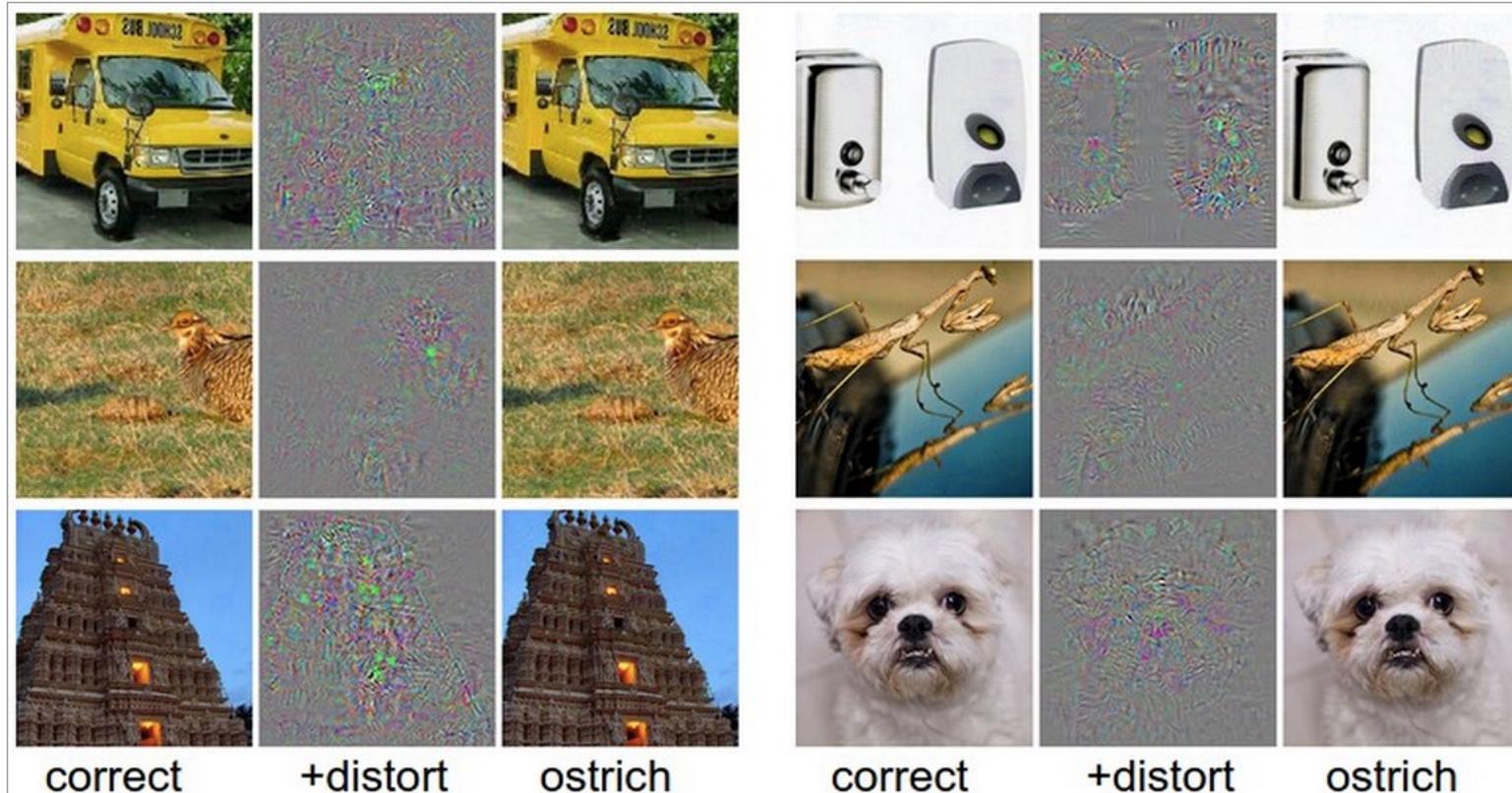
Deep learning packages

- PyTorch
- TensorFlow
- Keras
- Caffe
- Theano
- Matconvnet
- ...

https://en.wikipedia.org/wiki/Comparison_of_deep-learning_software

Breaking CNNs

Breaking CNNs

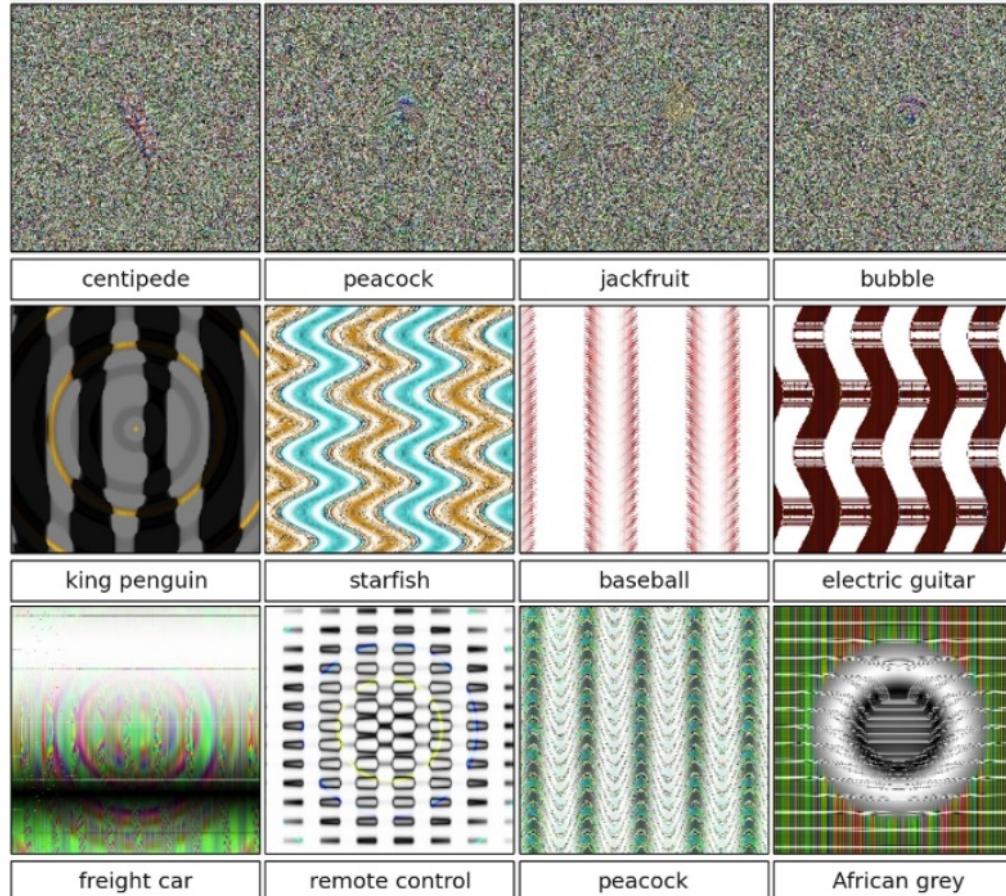


Take a correctly classified image (left image in both columns), and add a tiny distortion (middle) to fool the ConvNet with the resulting image (right).

<http://arxiv.org/abs/1312.6199>

<http://karpathy.github.io/2015/03/30/breaking-convnets/>

Breaking CNNs



<http://arxiv.org/abs/1412.1897>

<http://karpathy.github.io/2015/03/30/breaking-convnets/>

What is going on?

- Recall gradient descent training: modify the weights to reduce classifier error

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial E}{\partial \mathbf{w}}$$

- Adversarial examples: modify the image to increase classifier error

$$\mathbf{x} \leftarrow \mathbf{x} + \alpha \frac{\partial E}{\partial \mathbf{x}}$$

<http://arxiv.org/abs/1412.6572>

<http://karpathy.github.io/2015/03/30/breaking-convnets/>

What is going on?



\mathbf{x}

$+ .007 \times$

“nematode”
8.2% confidence



$$\frac{\partial E}{\partial \mathbf{x}}$$

“gibbon”
99.3 % confidence



$$\mathbf{x} \leftarrow \mathbf{x} + \alpha \frac{\partial E}{\partial \mathbf{x}}$$

<http://arxiv.org/abs/1412.6572>

<http://karpathy.github.io/2015/03/30/breaking-convnets/>

Fooling a linear classifier

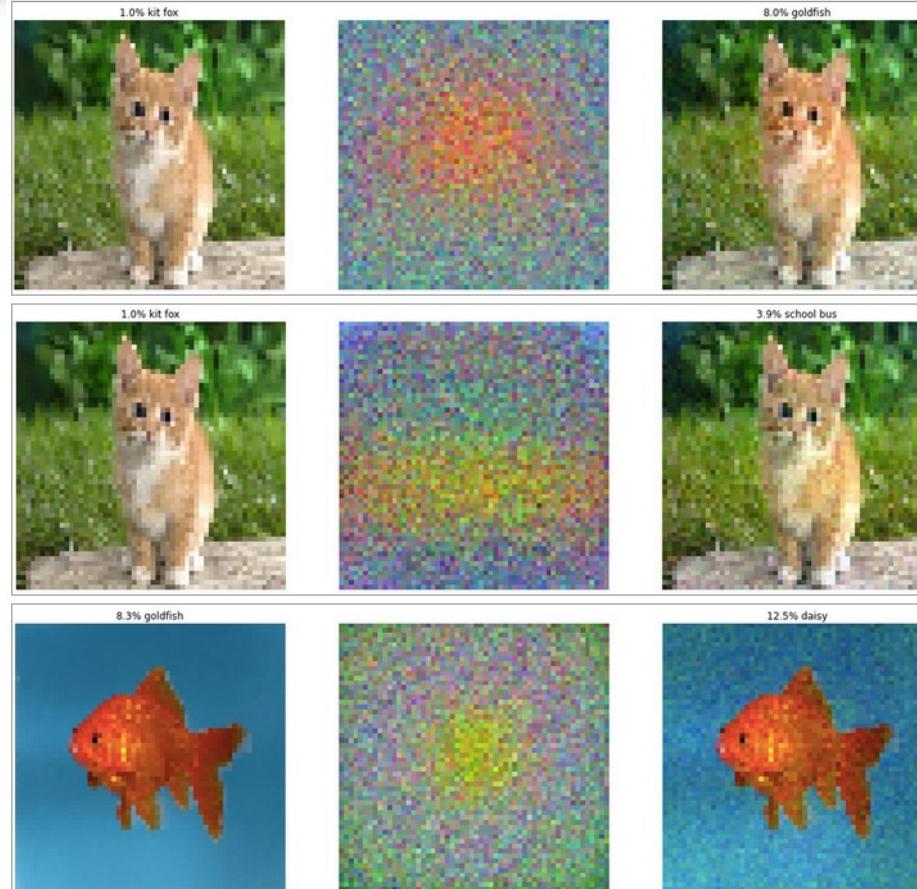
- Perceptron weight update: add a small multiple of the example to the weight vector:

$$w \leftarrow w + \alpha x$$

- To fool a linear classifier, add a small multiple of the weight vector to the training example:

$$x \leftarrow x + \alpha w$$

Fooling a linear classifier



Fooled linear classifier: The starting image (left) is classified as a kit fox. That's incorrect, but what can you expect from a linear classifier? However, if we add a small amount "goldfish" weights to the image (top row, middle), suddenly the classifier is convinced that it's looking at one with high confidence. We can distort it with the school bus template instead if we wanted to.

<http://karpathy.github.io/2015/03/30/breaking-convnets/>

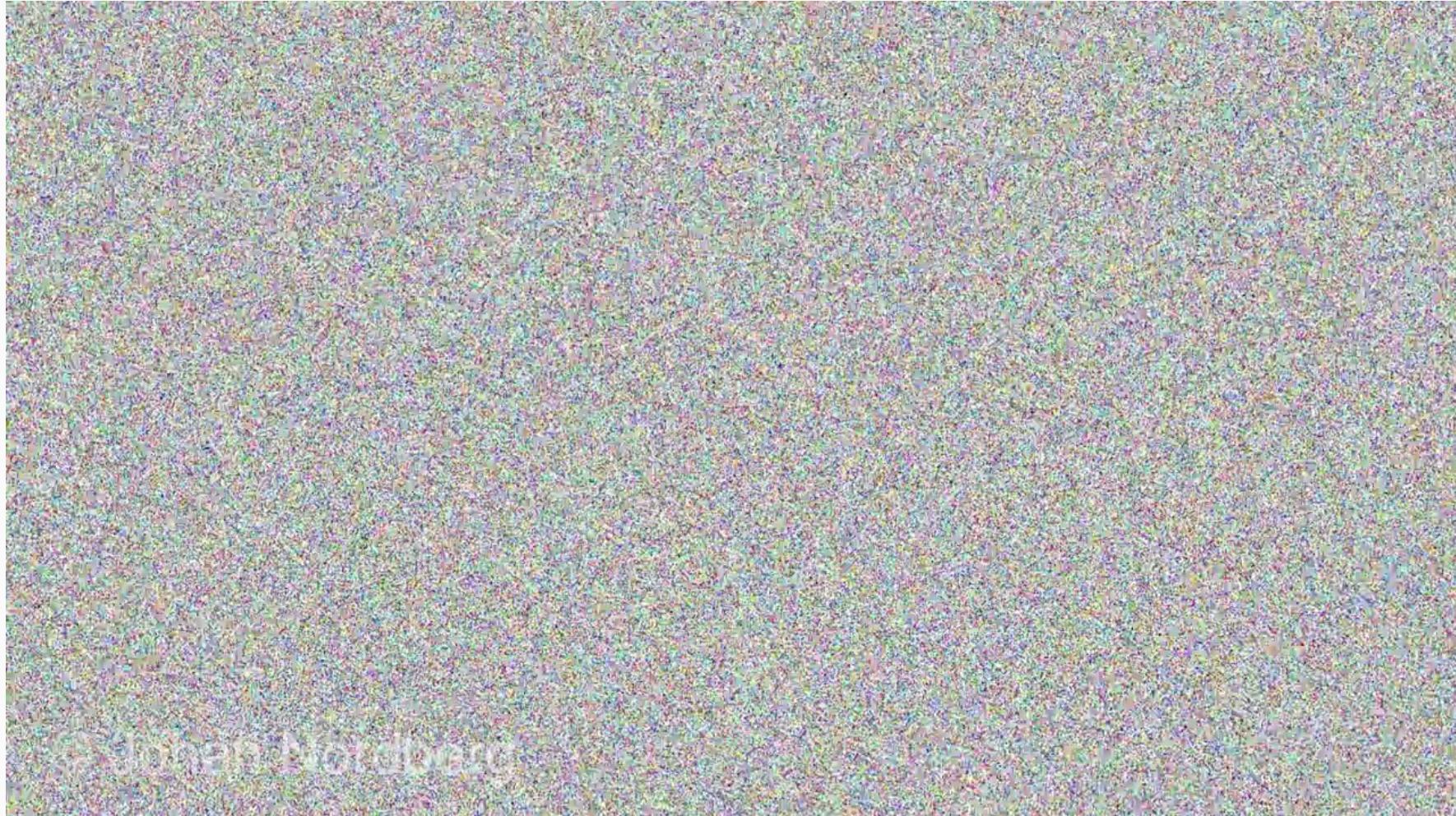
Google DeepDream

- Modify the image to maximize activations of units in a given layer



<https://github.com/google/deepdream/blob/master/dream.ipynb>

Google DeepDream



<https://github.com/google/deepdream/blob/master/dream.ipynb>

Summary

- Currently most applications of deep neural networks in computer vision are based on supervised learning
- Need lots of annotated training data (and GPUs)
- Training via back-propagation takes time and is prone to local minima (many tips and tricks)
- Manual design of feature detectors and descriptors is replaced by manual design of network architectures (-> better automation needed)
- Lately there has been progress in weakly supervised and unsupervised learning (e.g. generative adversarial networks)

Further reading

- LeCun, Bengio, Hinton: Deep learning. Nature 2015.

<http://pages.cs.wisc.edu/~dyer/cs540/handouts/deep-learning-nature2015.pdf>

- Goodfellow, Bengio, Courville:

<http://www.deeplearningbook.org/>