

Neural Network

A neural network comprises neurons connected in an acyclic graph

The outputs of neurons can become inputs to other neurons

Neural networks typically contain multiple layers of neurons

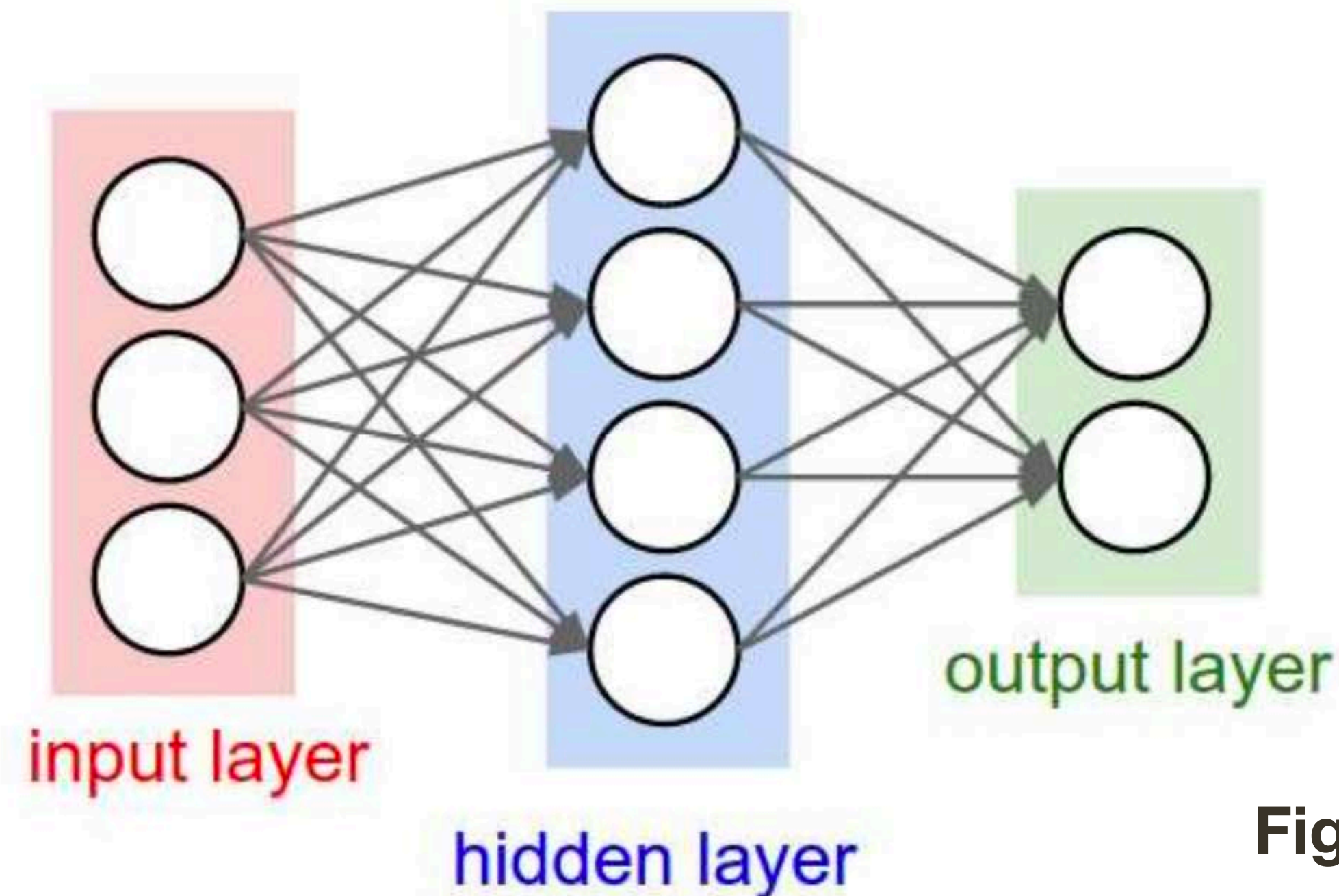
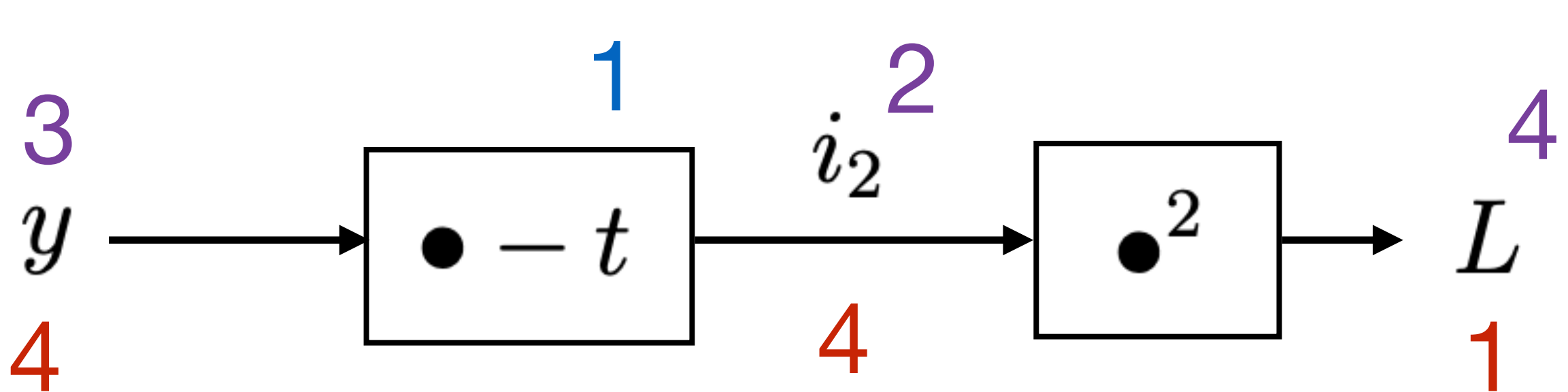
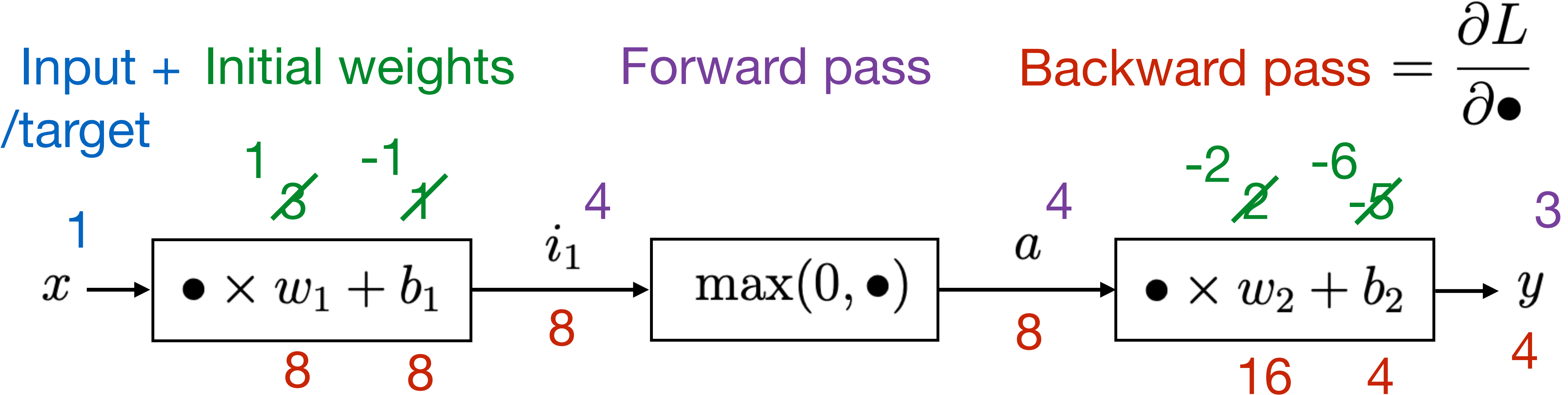


Figure credit: Fei-Fei and Karpathy

Example of a neural network with three inputs, a single hidden layer of four neurons, and an output layer of two neurons

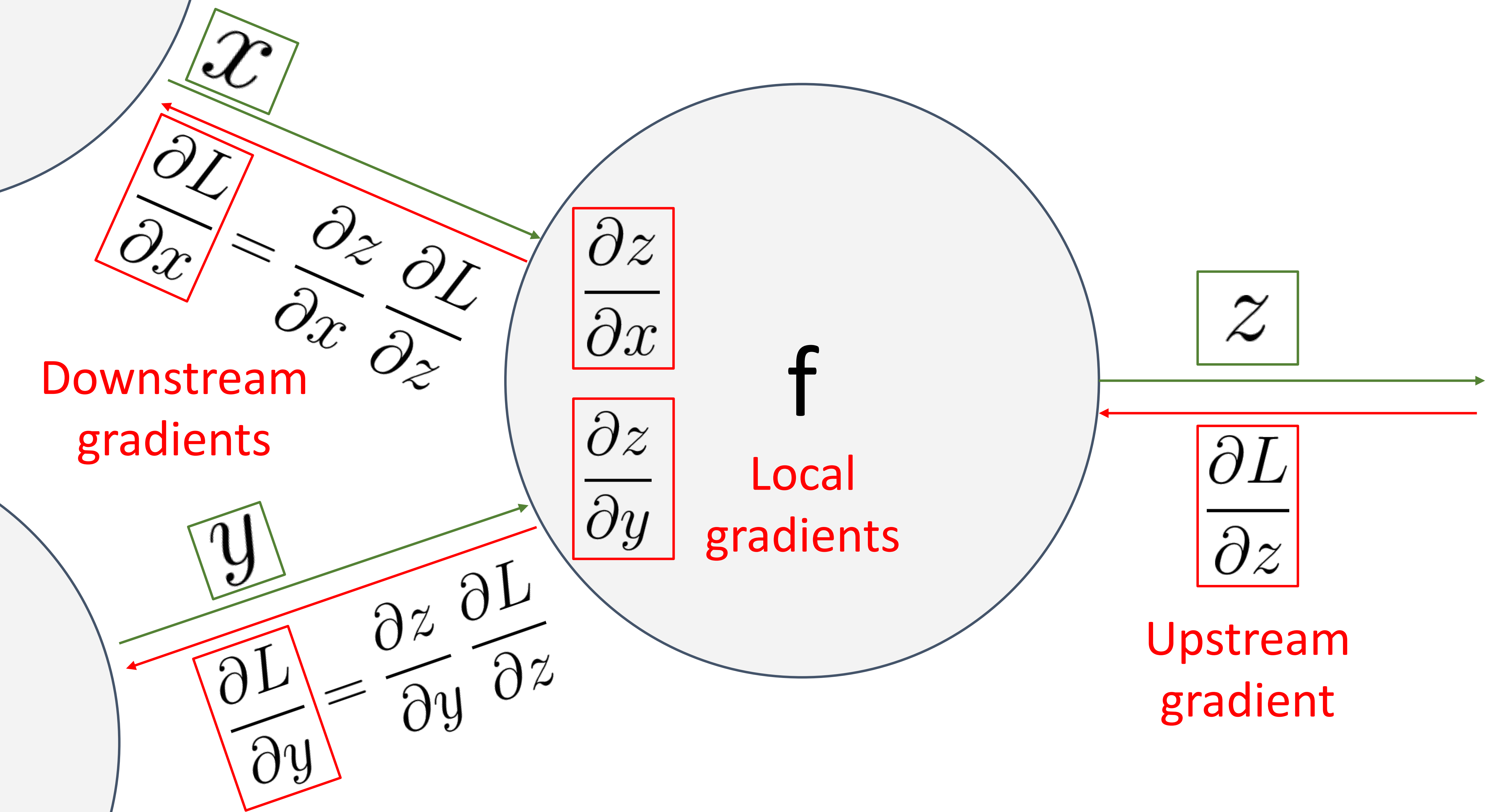
2-Layer **Neural** Network — 1 hidden, 1 input/output



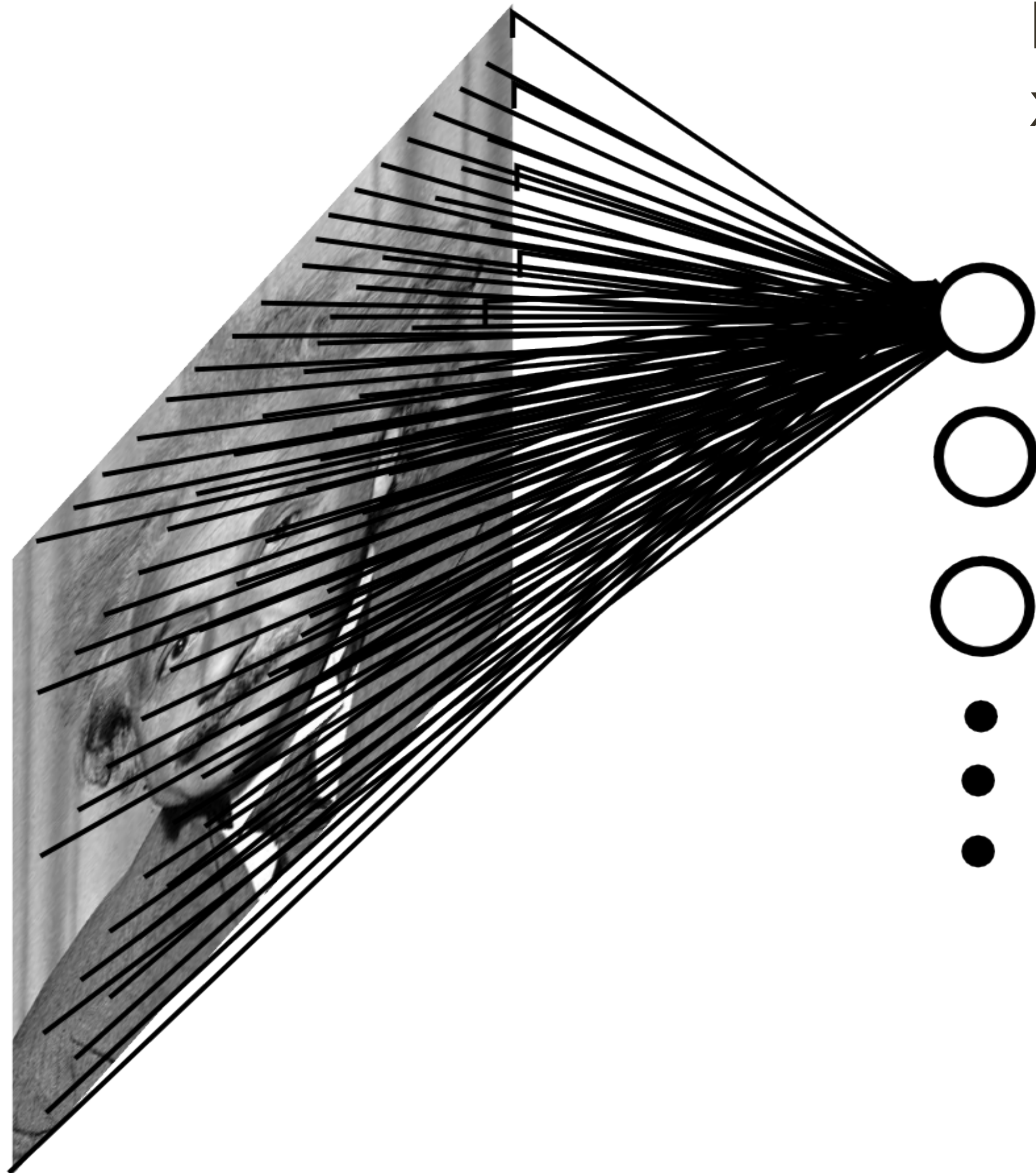
$$\begin{bmatrix} w_1 \\ b_1 \\ w_2 \\ b_2 \end{bmatrix} \rightarrow \begin{bmatrix} 3 \\ 1 \\ 2 \\ -5 \end{bmatrix} - \alpha \begin{bmatrix} 8 \\ 8 \\ 16 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \\ -2 \\ -6 \end{bmatrix}$$

+ update weights

Repeat: +Input/target, Forward, Backward, Update until convergence!



Fully Connected Layer



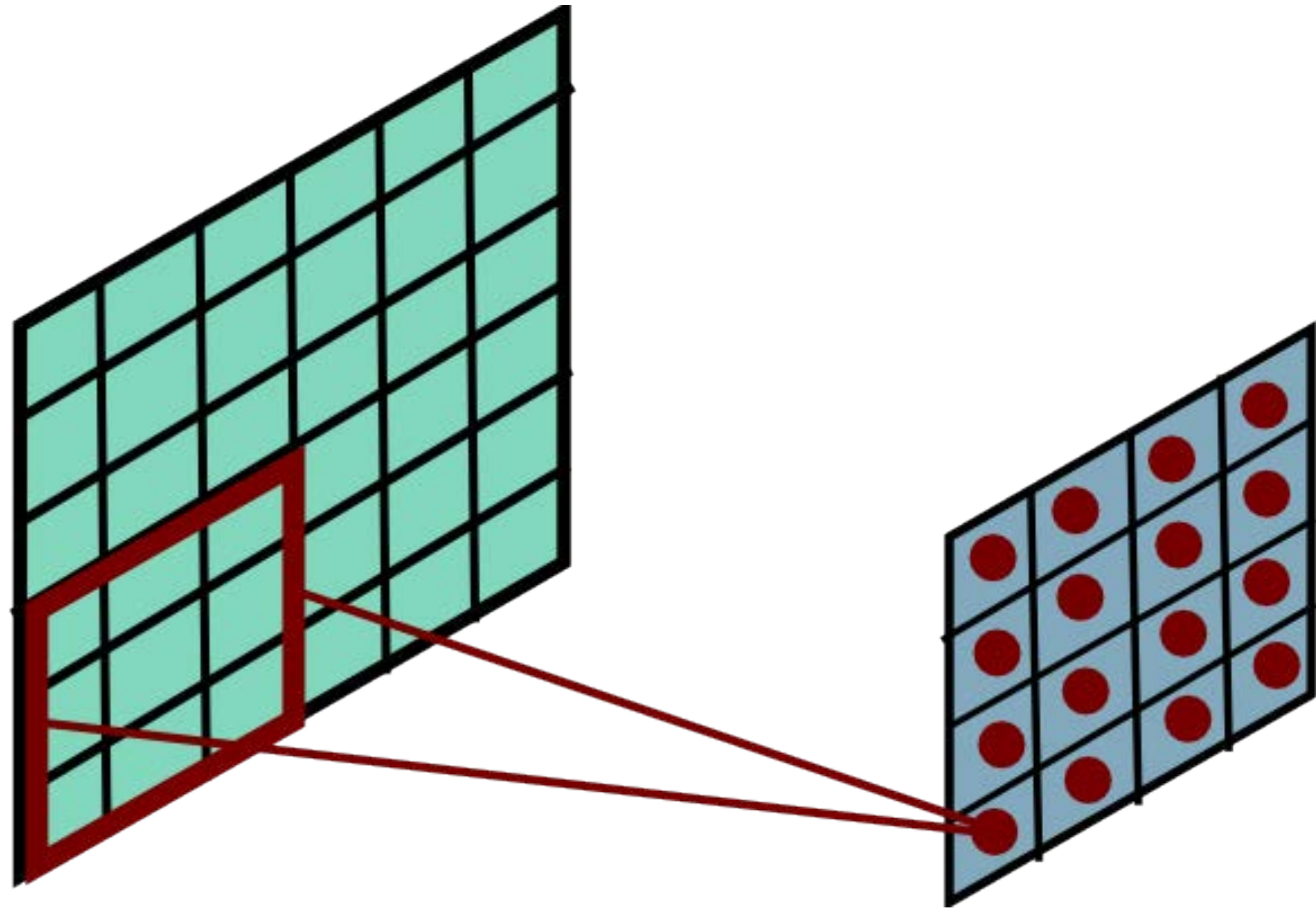
Example: 200 x 200 image (small)
x 40K hidden units (same size)

= **1.6 Billion** parameters (for one layer!)

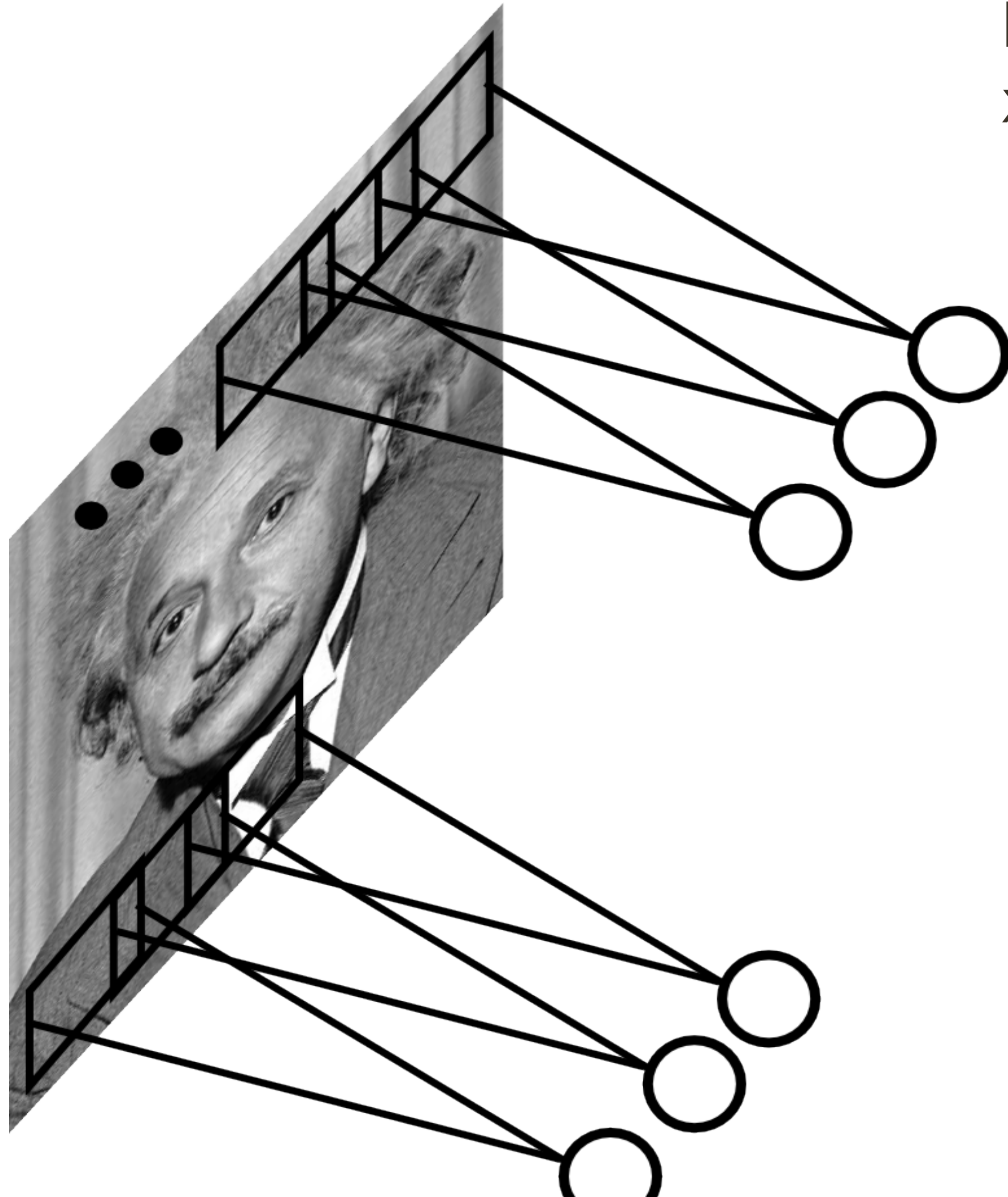
Spatial correlations are generally local

Waste of resources + we don't have
enough data to train networks this large

Convolutional Layer



Convolutional Layer

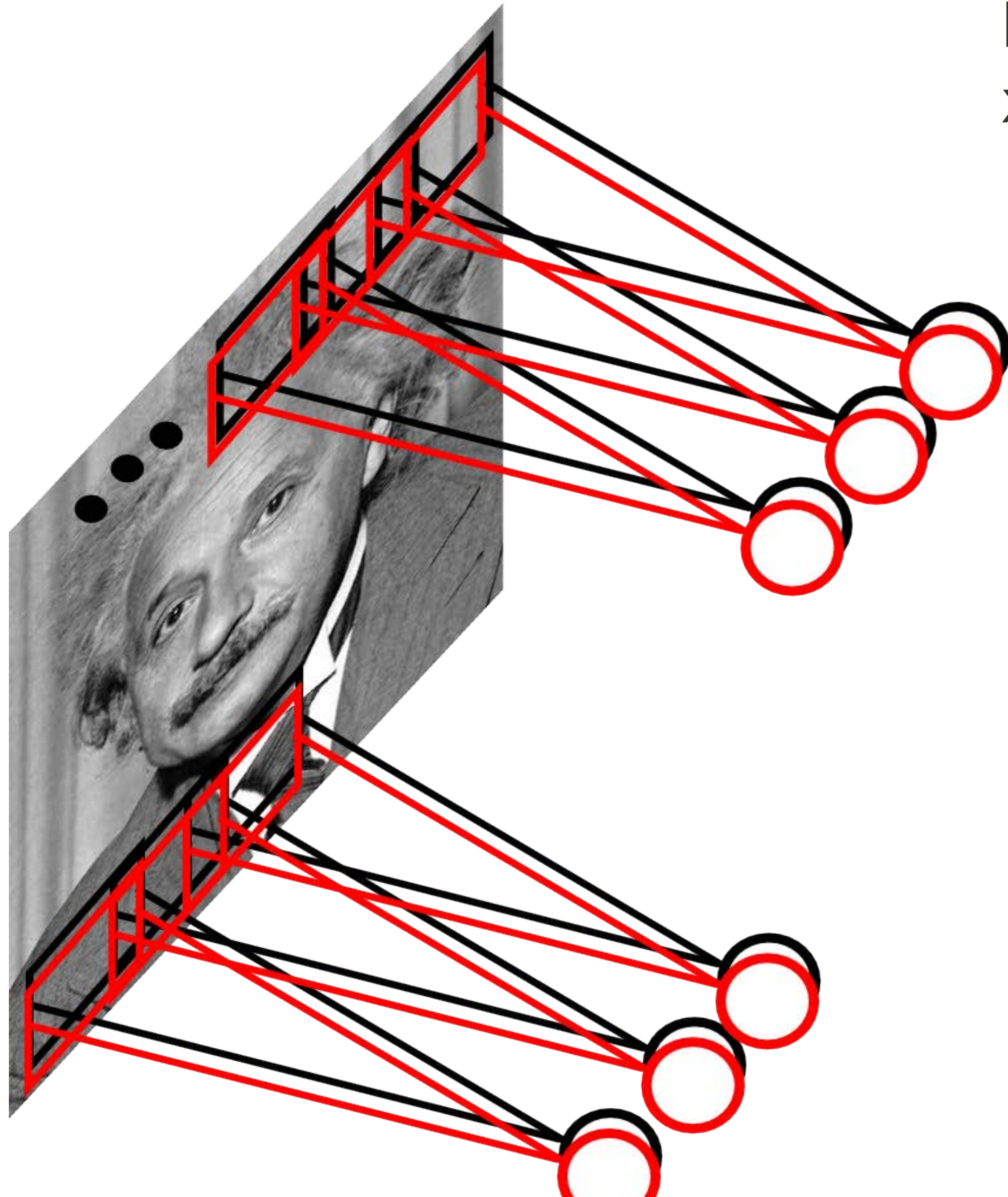


Example: 200 x 200 image (small)
x 40K hidden units (same size)

Filter size: 10 x 10
= 100 parameters

Share the same parameters across the locations (assuming input is stationary)

Convolutional Layer



Example: 200 x 200 image (small)
x 40K hidden units (same size)

Filter size: 10 x 10

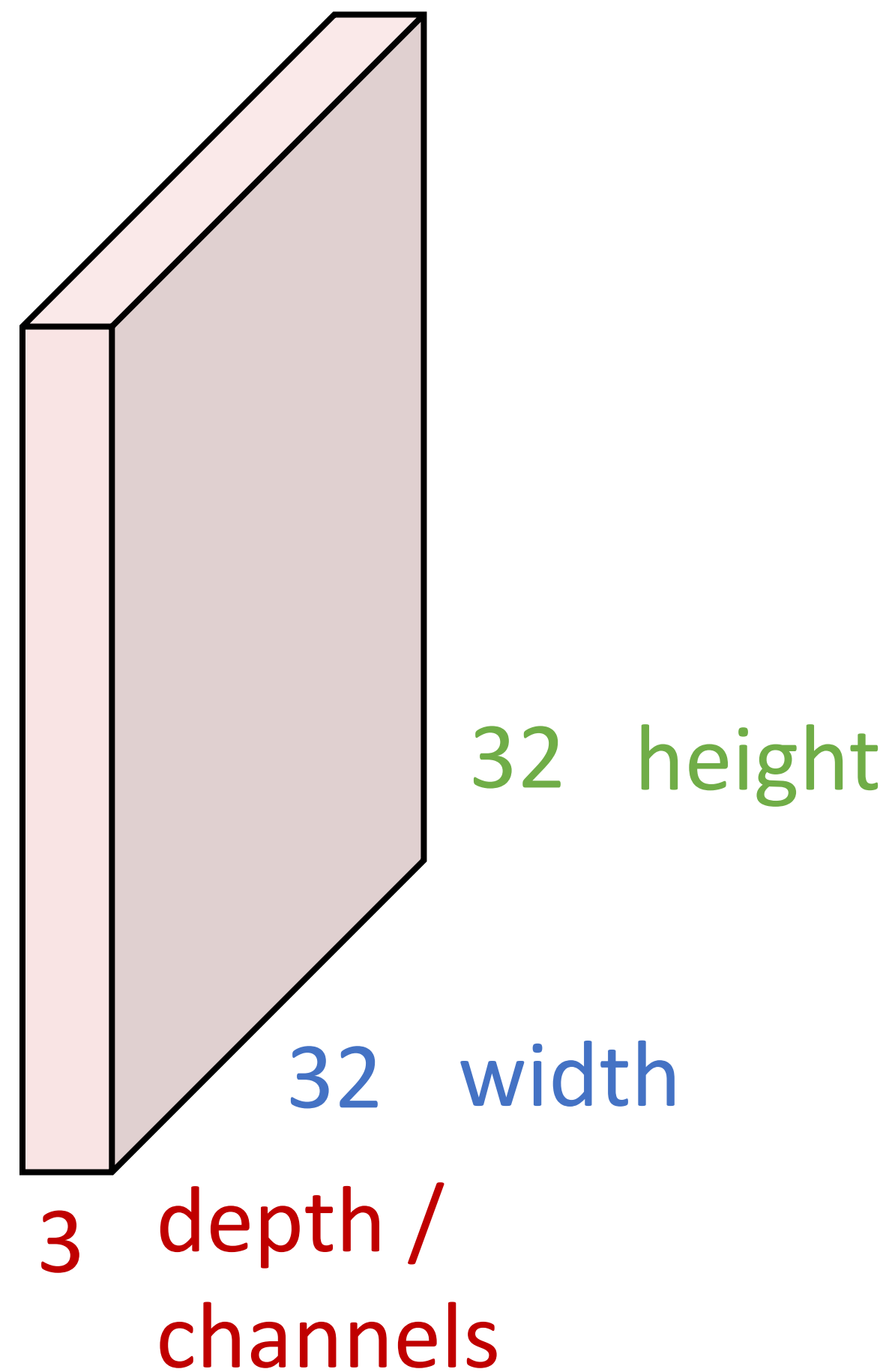
of filters: 20

= 2000 parameters

Learn **multiple filters**
→ **multiple output channels**

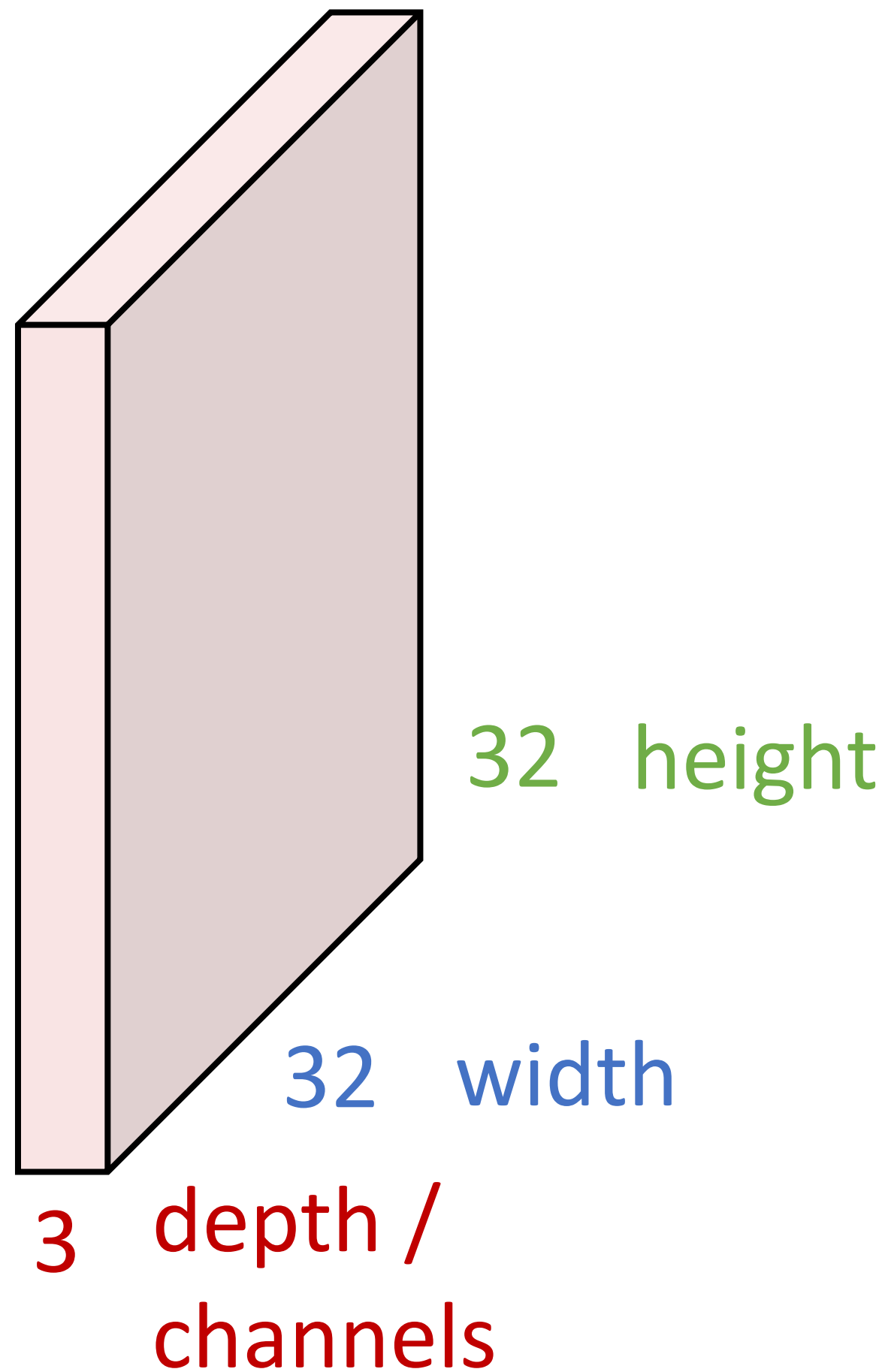
Convolution Layer

3x32x32 image: preserve spatial structure



Convolution Layer

3x32x32 image



Filters always extend the full depth of the input volume

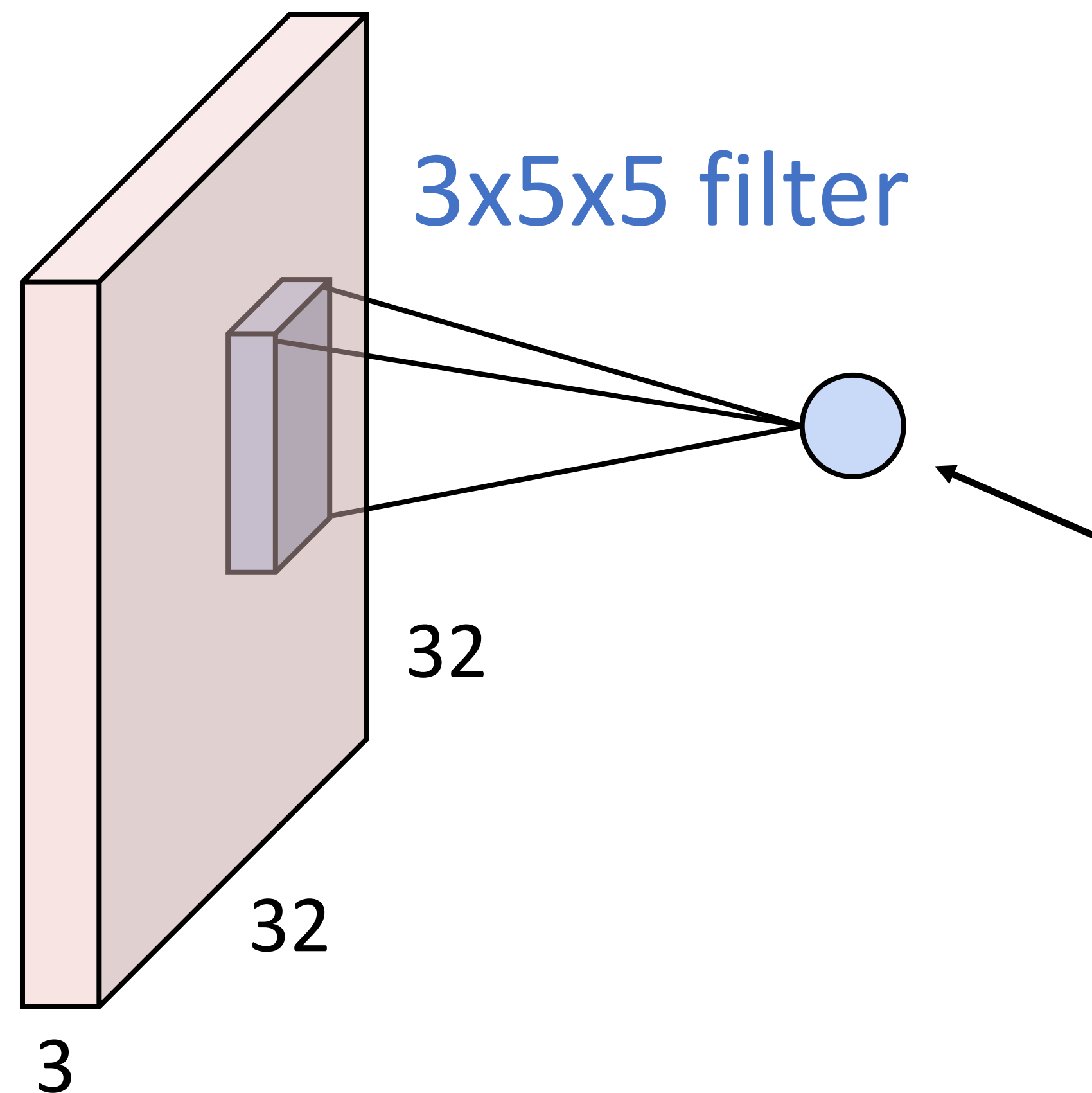
3x5x5 filter



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

Convolution Layer

3x32x32 image



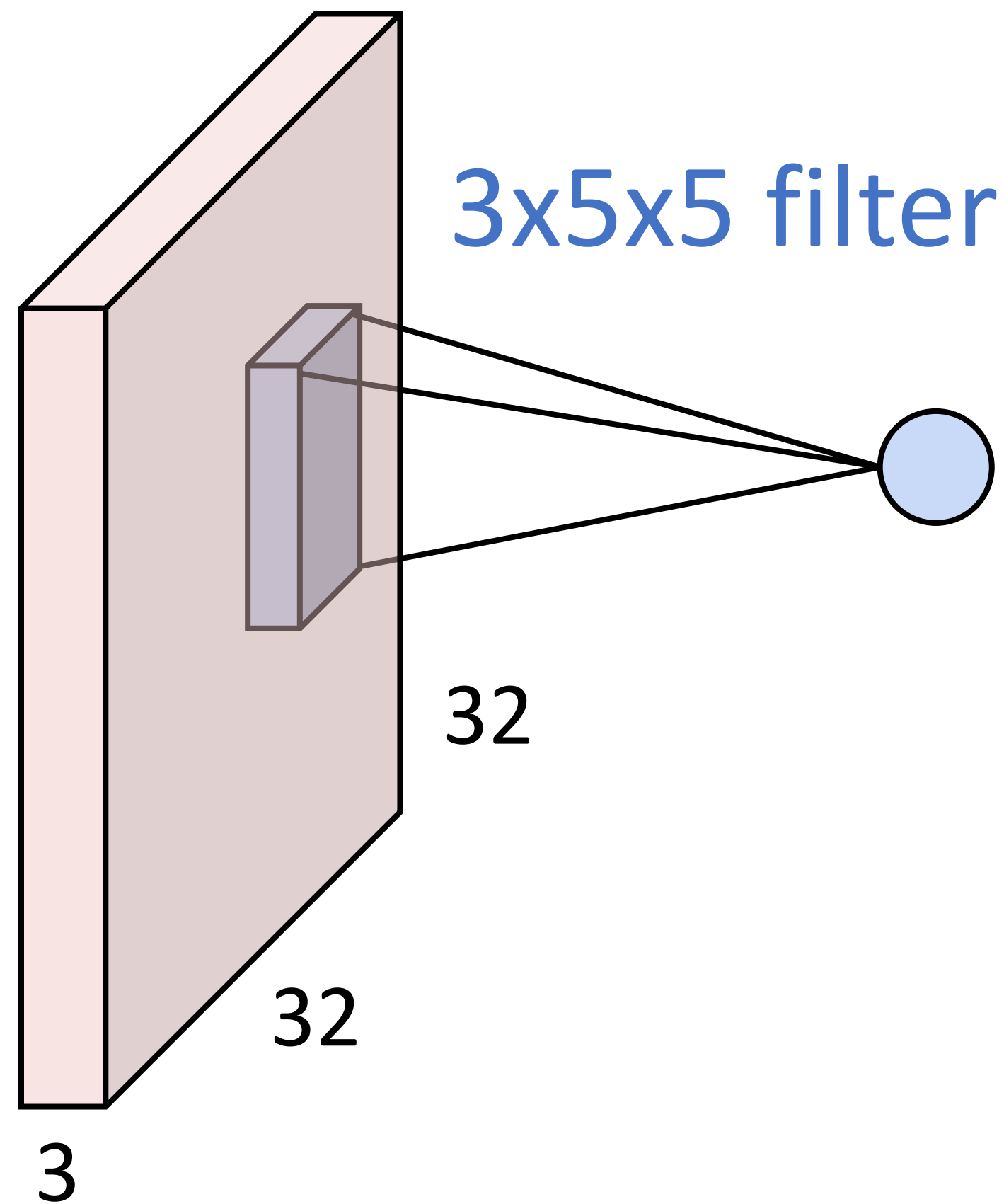
1 number:

the result of taking a dot product between the filter and a small 3x5x5 chunk of the image (i.e. $3*5*5 = 75$ -dimensional dot product + bias)

$$w^T x + b$$

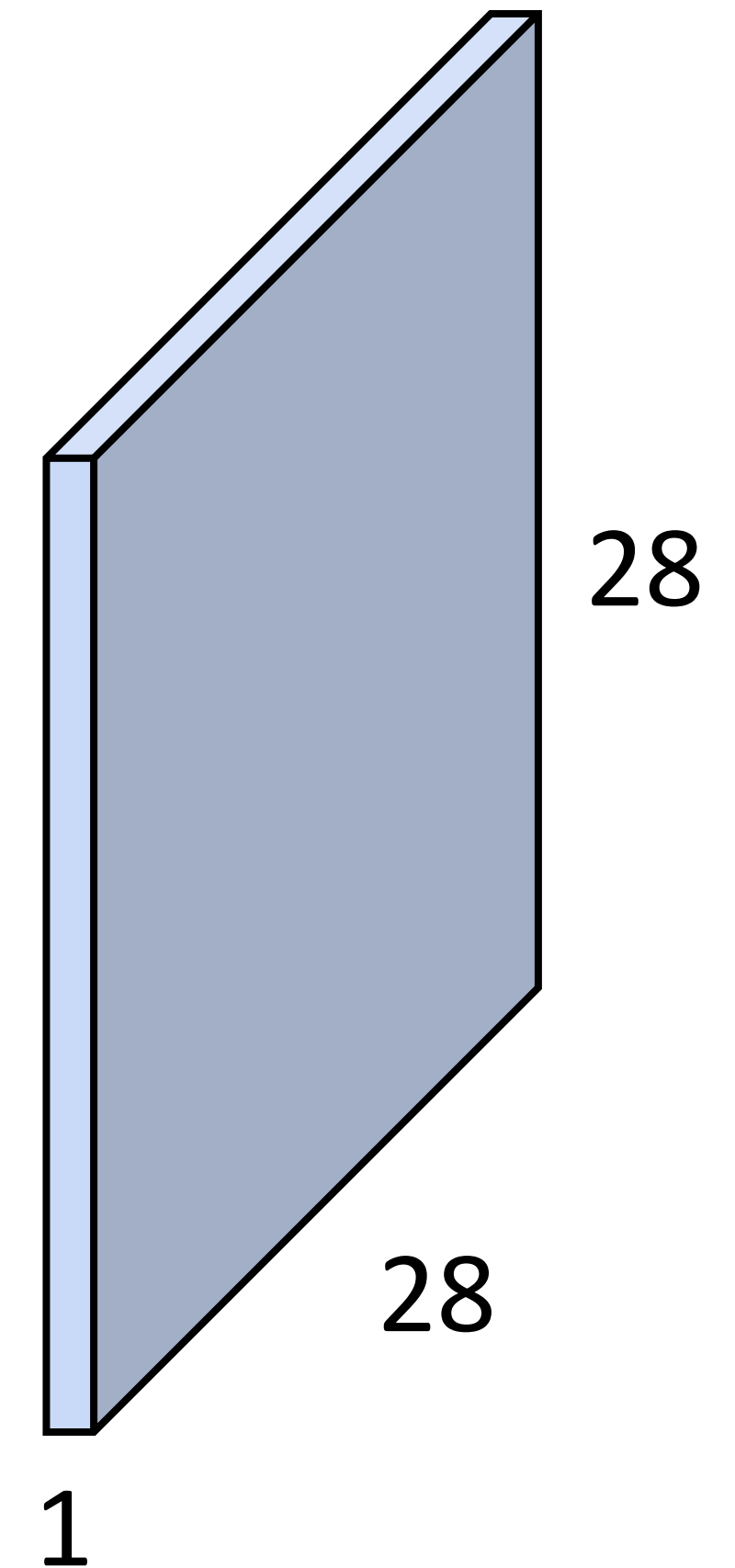
Convolution Layer

3x32x32 image



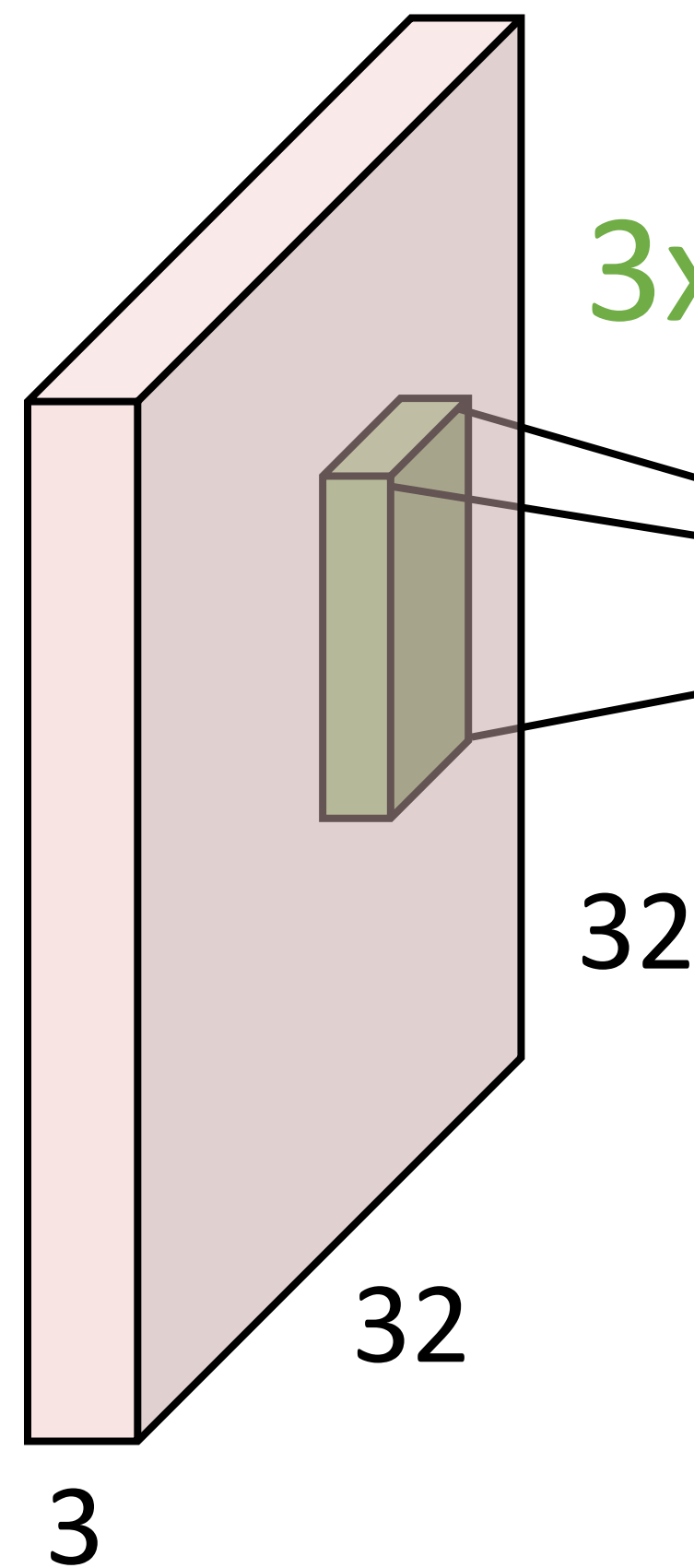
convolve (slide) over
all spatial locations

1x28x28
activation map

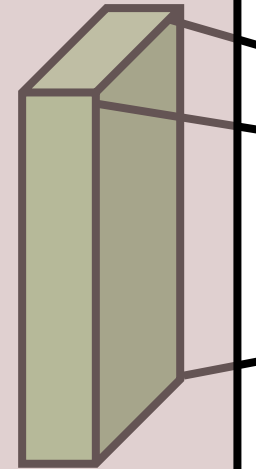


Convolution Layer

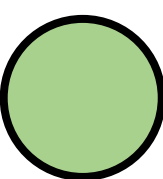
3x32x32 image



3x5x5 filter

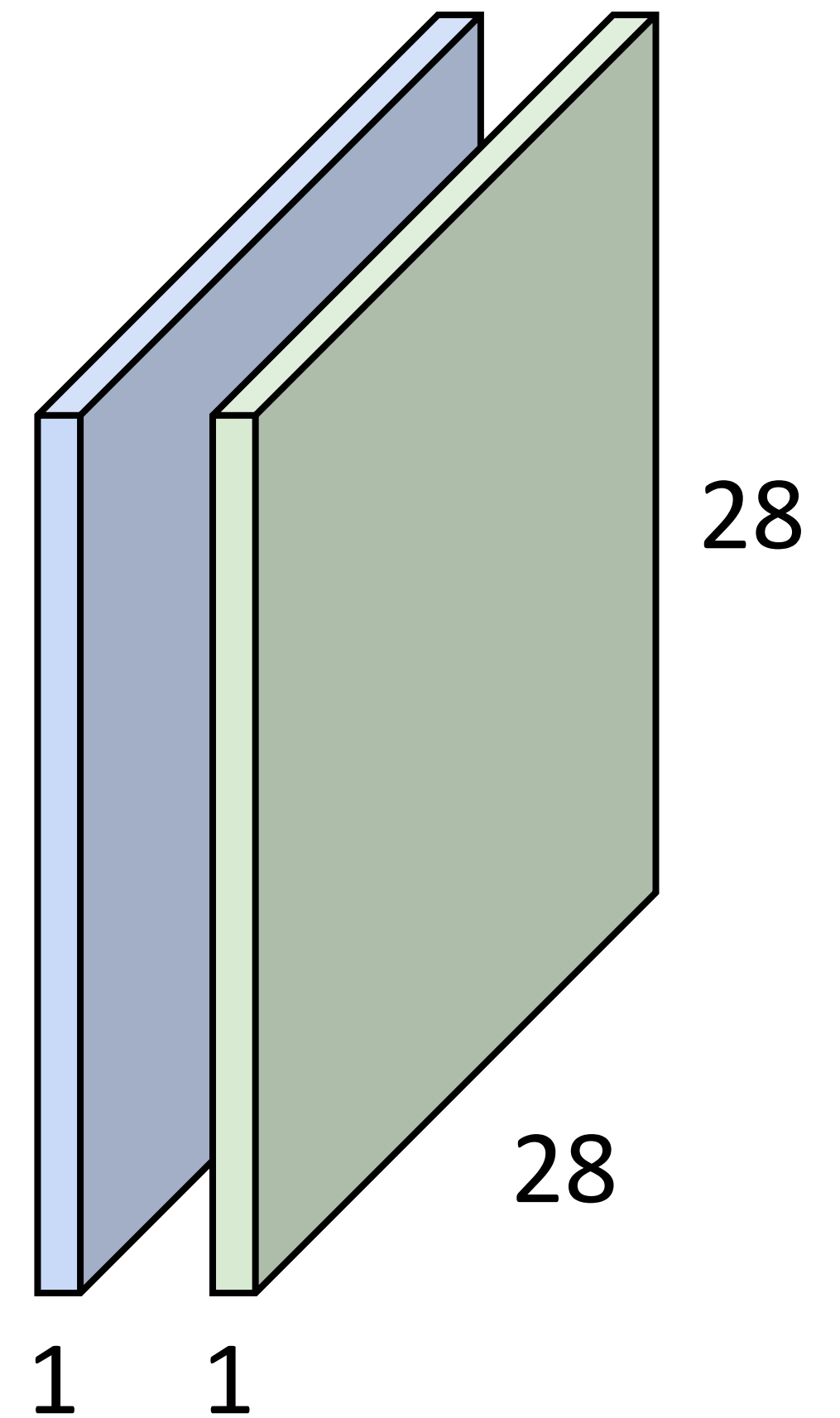


Consider repeating with a second (green) filter:



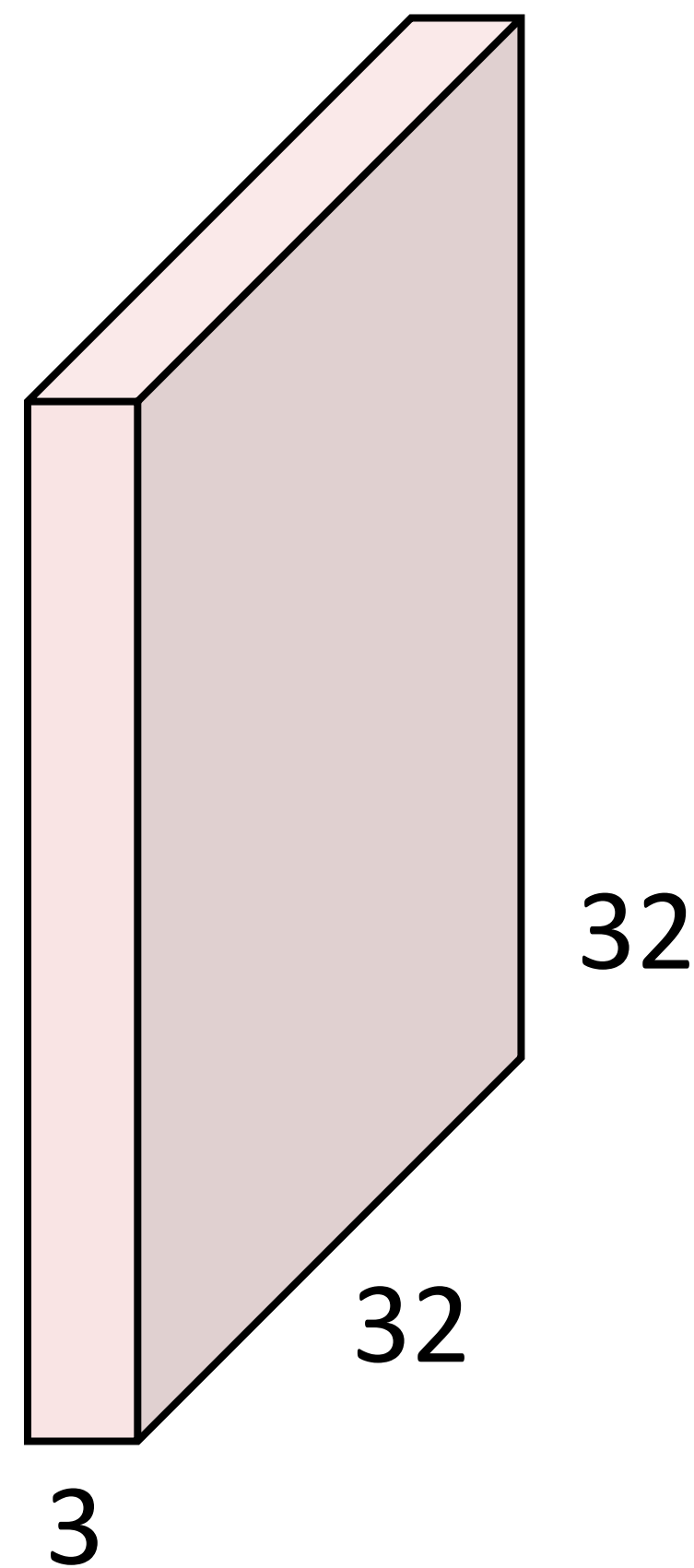
convolve (slide) over all spatial locations

two 1x28x28 activation map

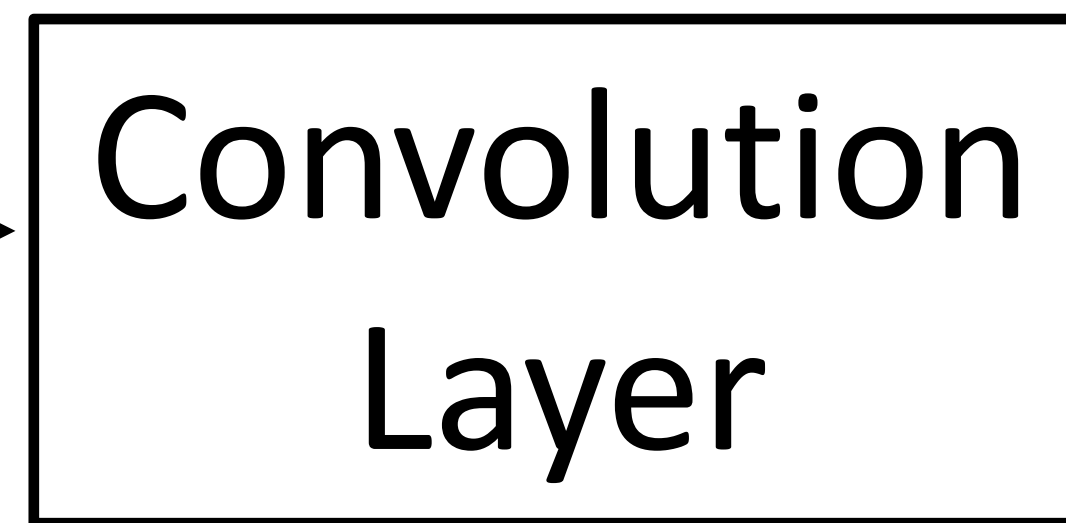


Convolution Layer

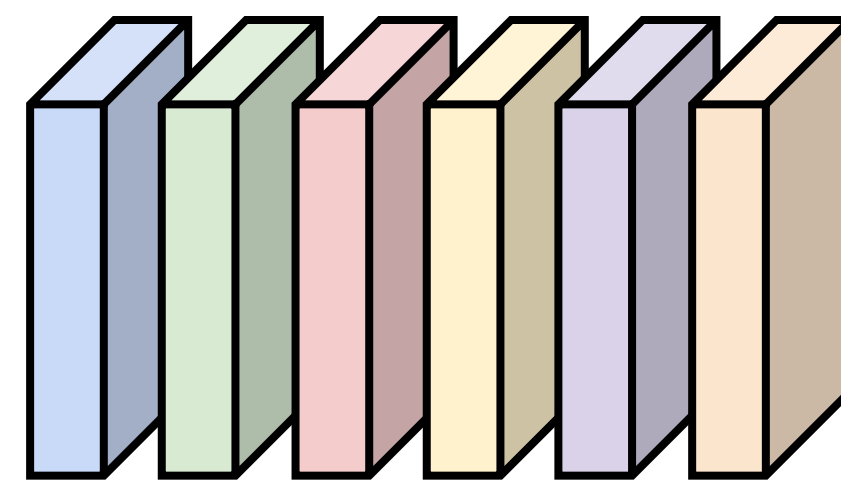
3x32x32 image



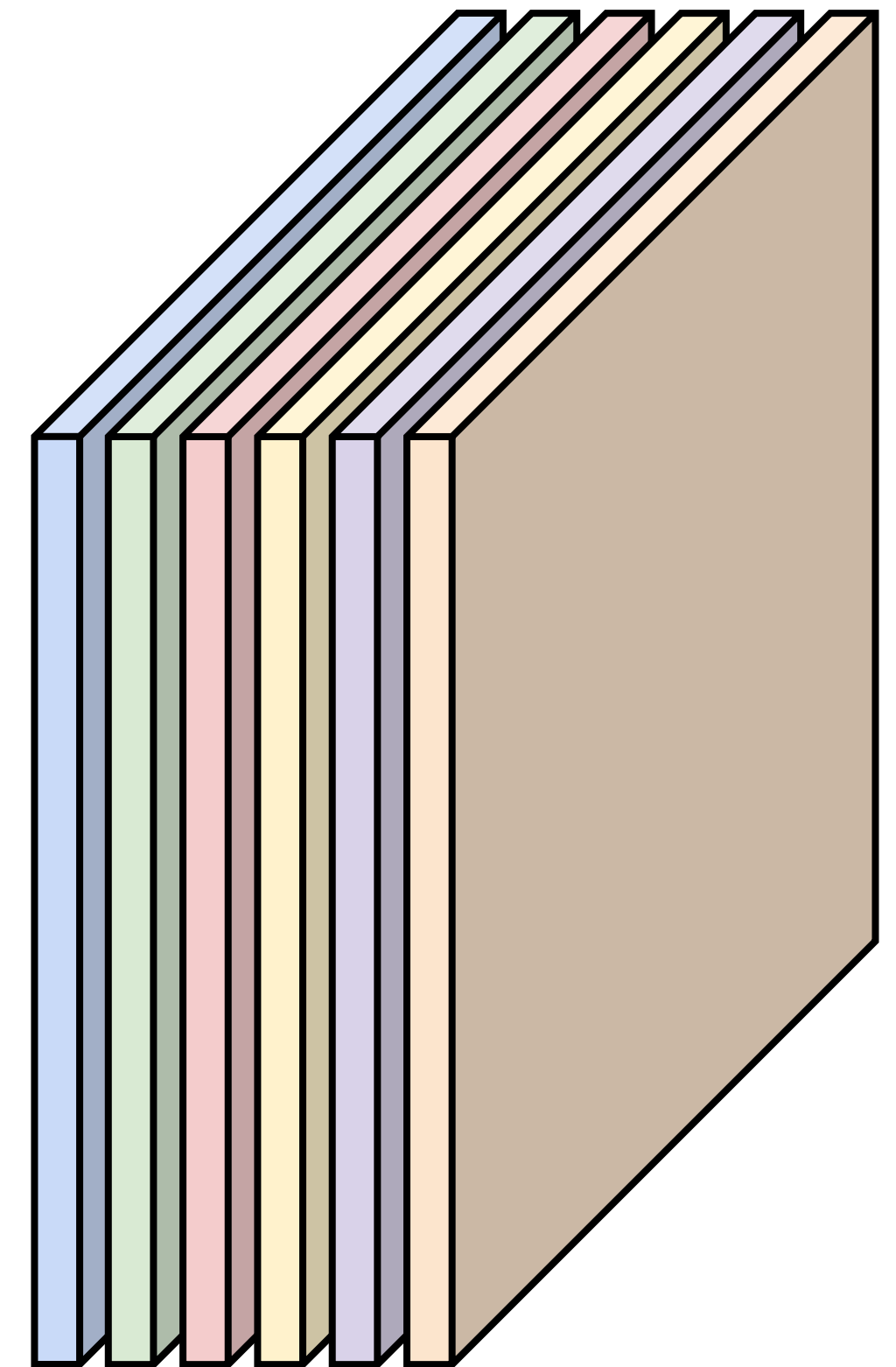
Consider 6 filters,
each 3x5x5



6x3x5x5
filters



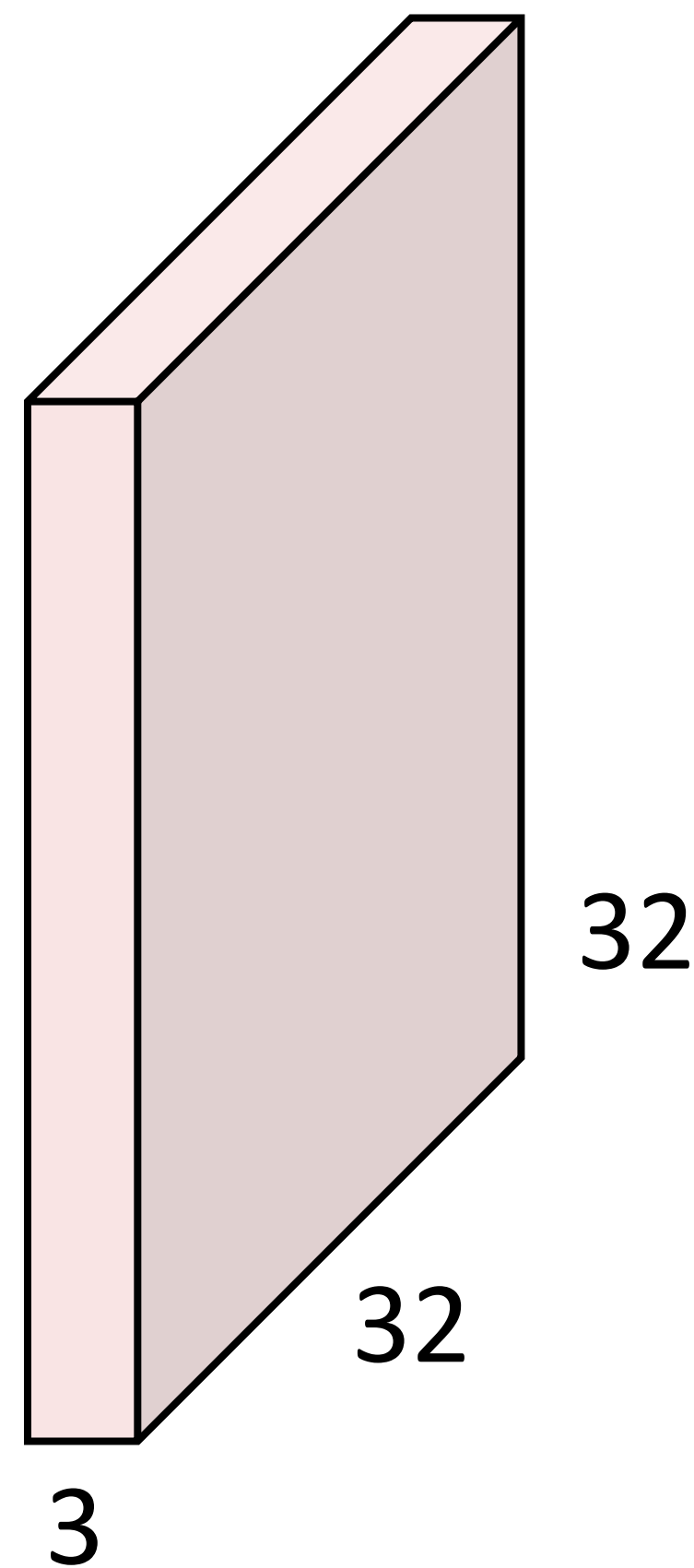
6 activation maps,
each 1x28x28



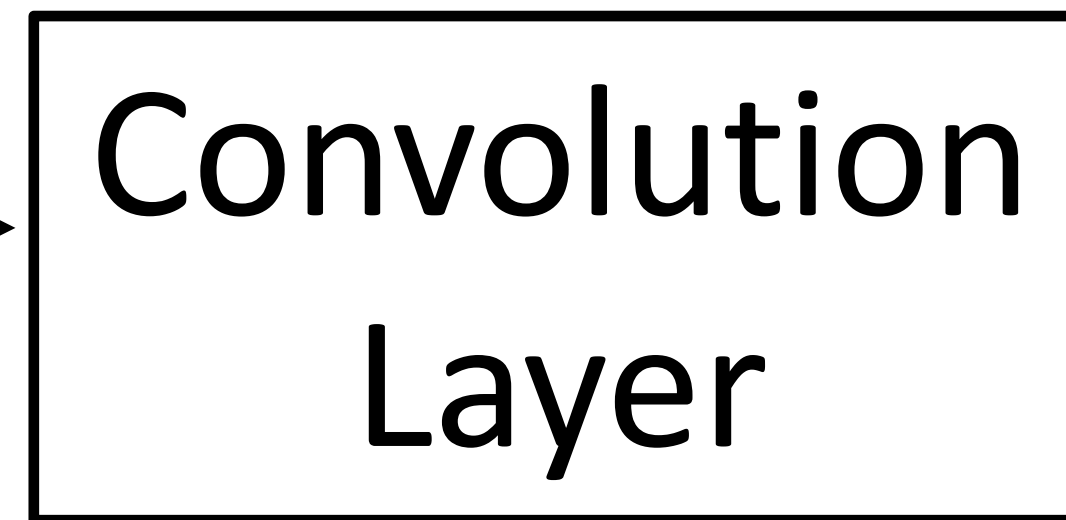
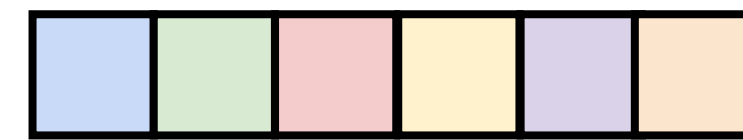
Stack activations to get a
6x28x28 output image!

Convolution Layer

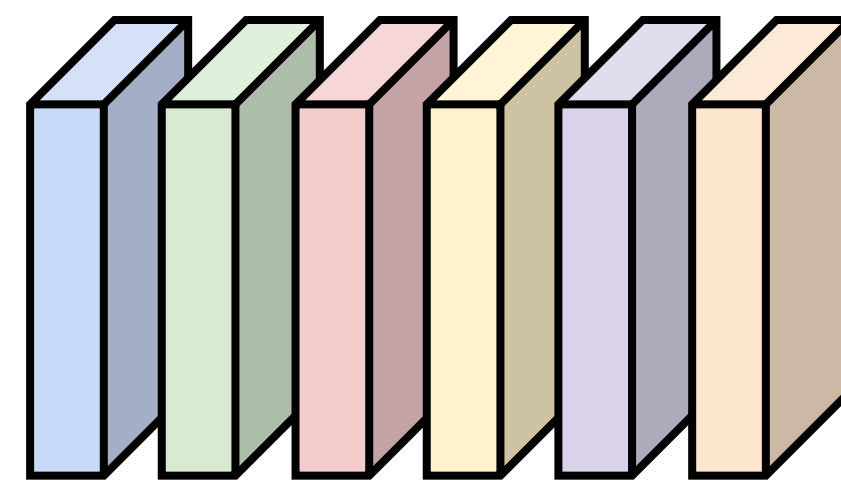
3x32x32 image



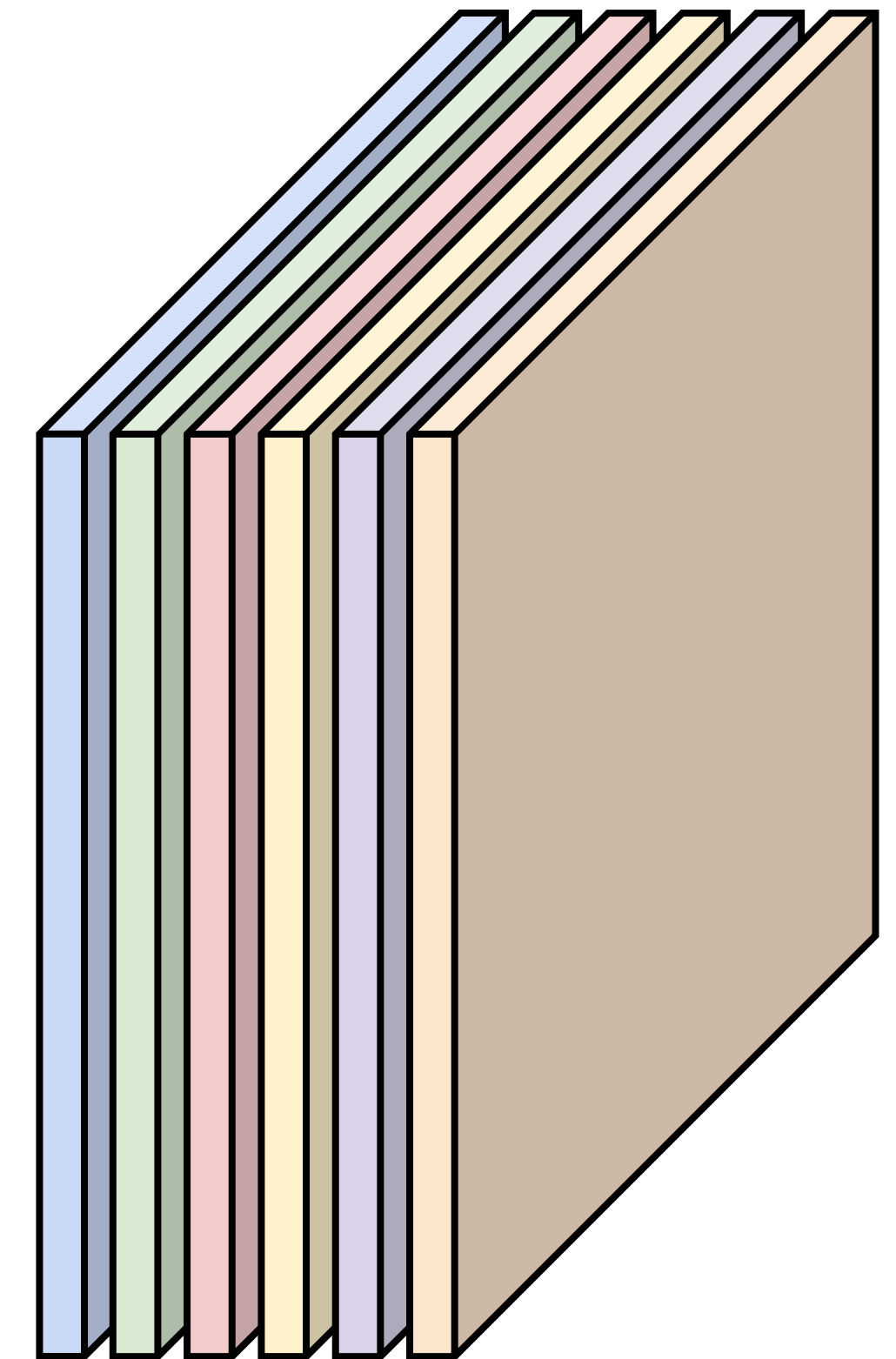
Also 6-dim bias vector:



6x3x5x5 filters



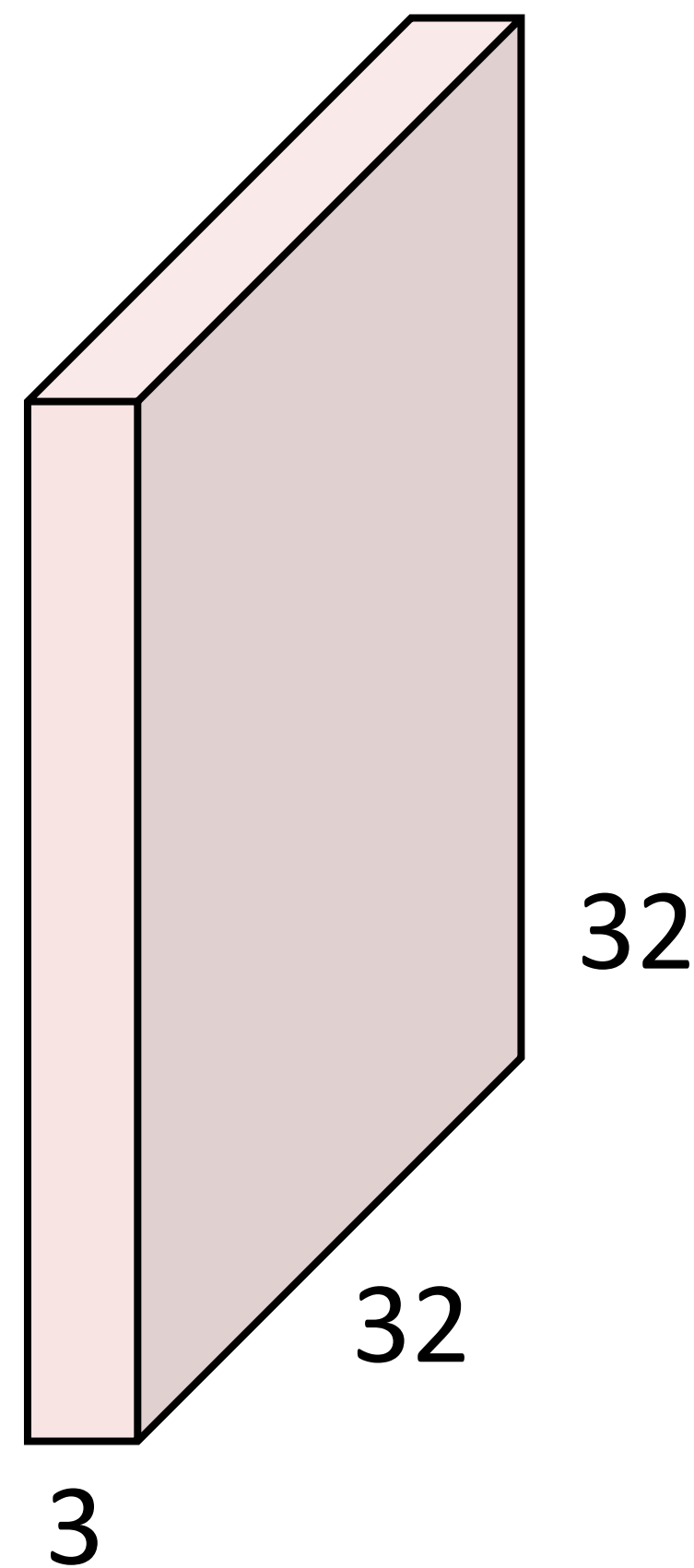
6 activation maps,
each 1x28x28



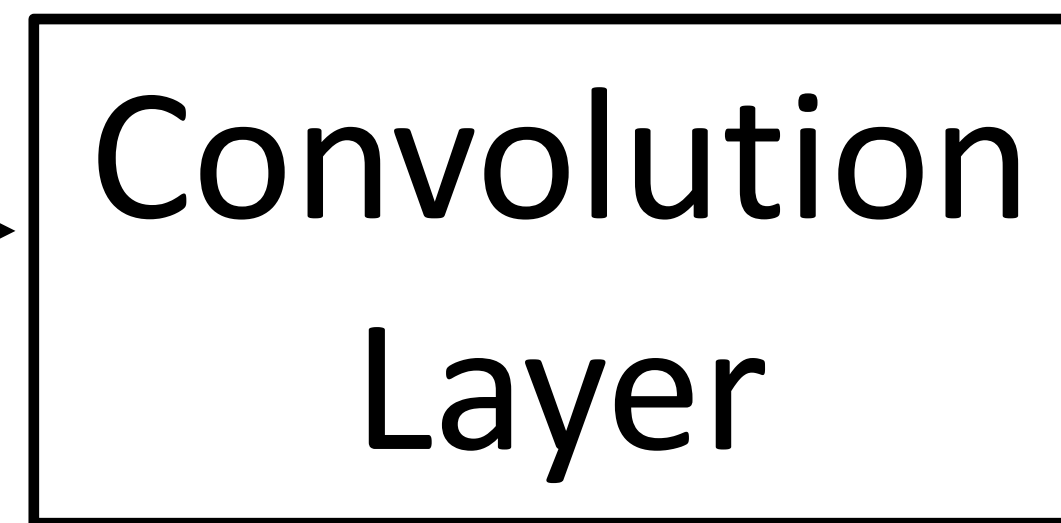
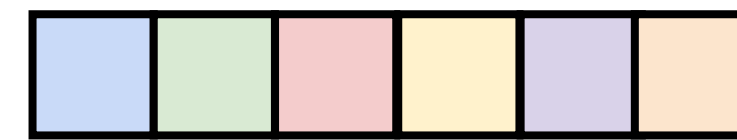
Stack activations to get a
6x28x28 output image!

Convolution Layer

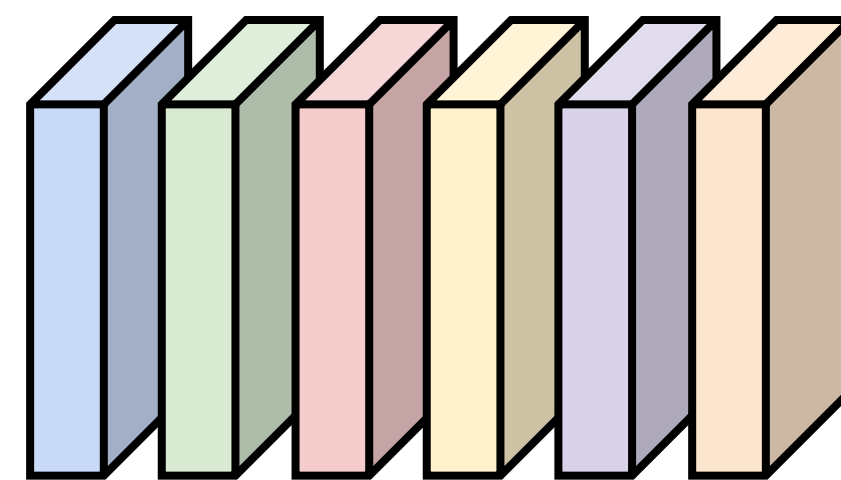
3x32x32 image



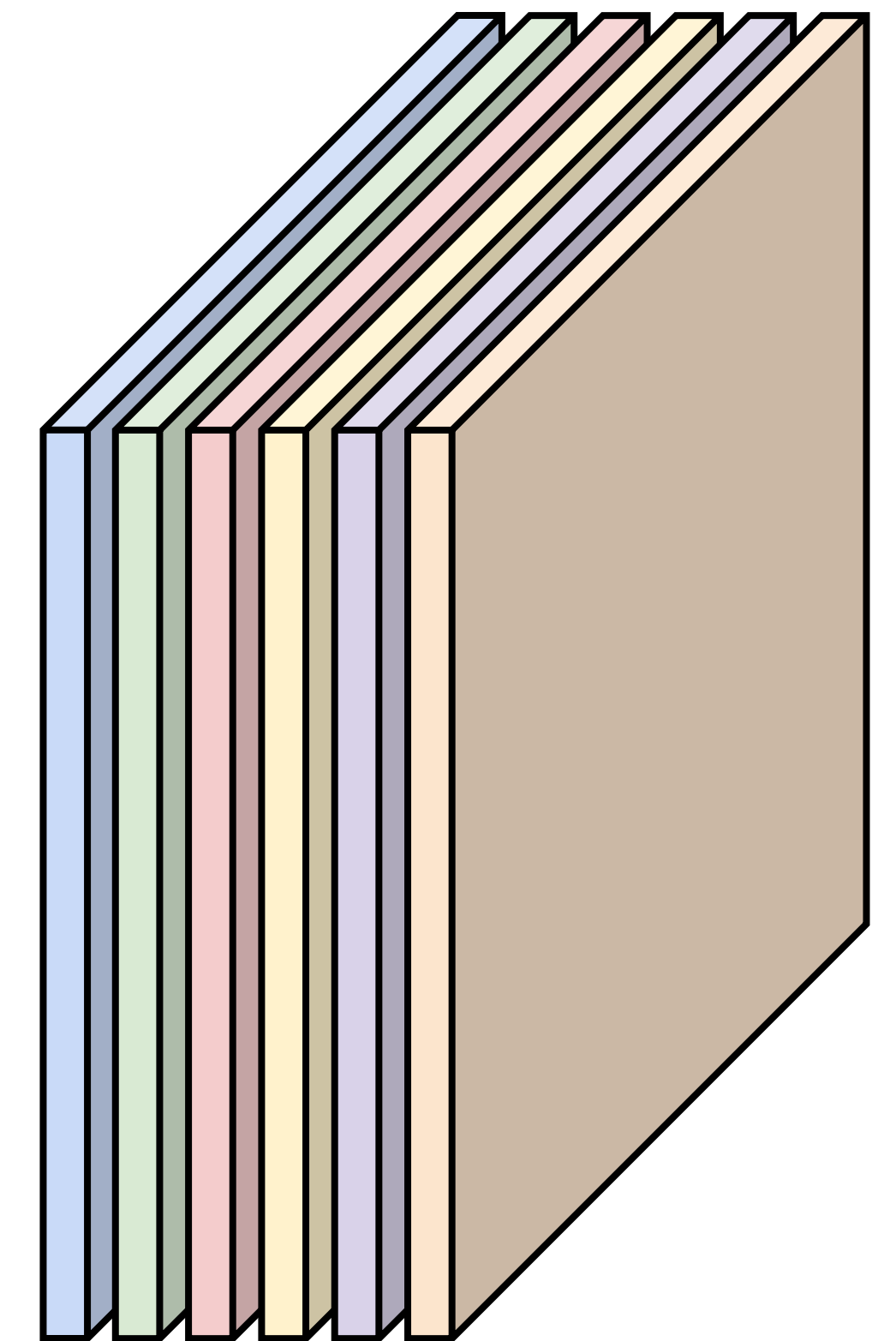
Also 6-dim bias vector:



6x3x5x5 filters



28x28 grid, at each point a 6-dim vector

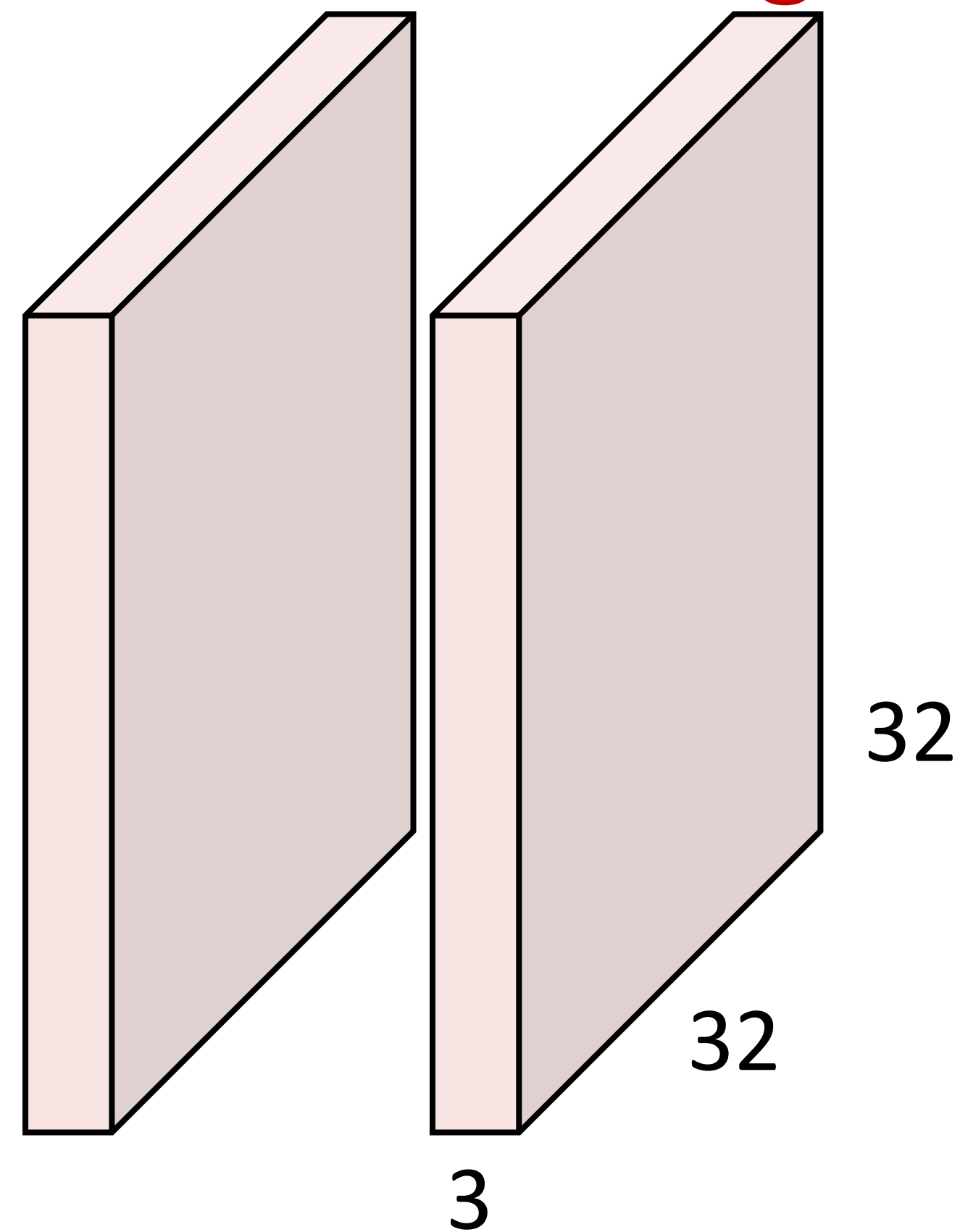


Stack activations to get a 6x28x28 output image!

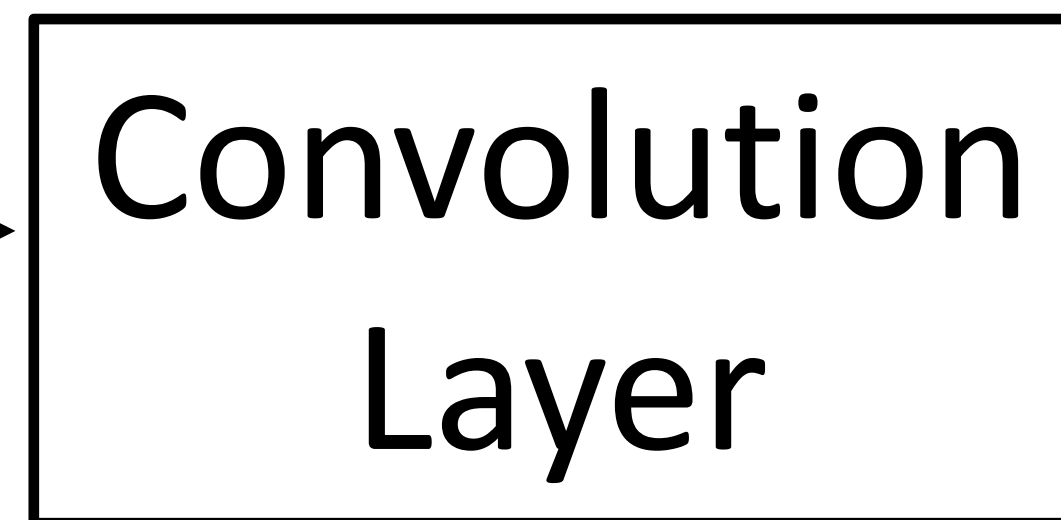
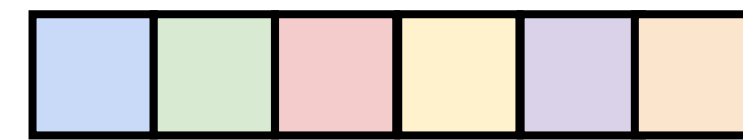
Convolution Layer

$2 \times 3 \times 32 \times 32$

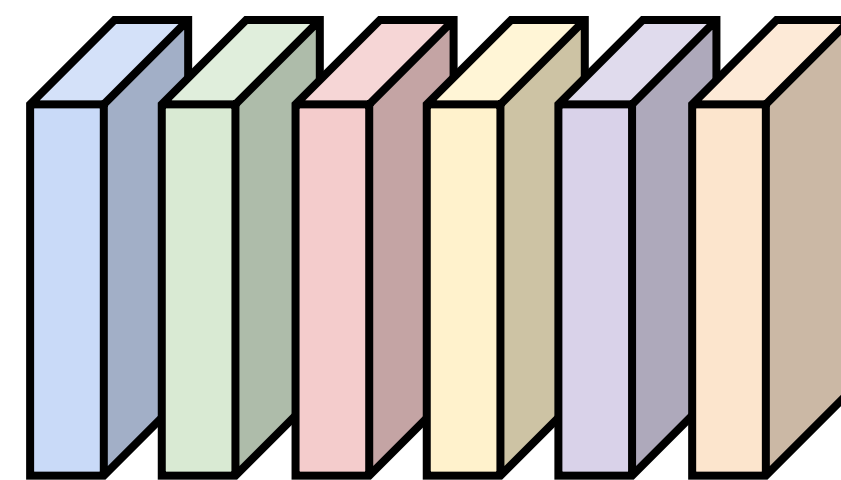
Batch of images



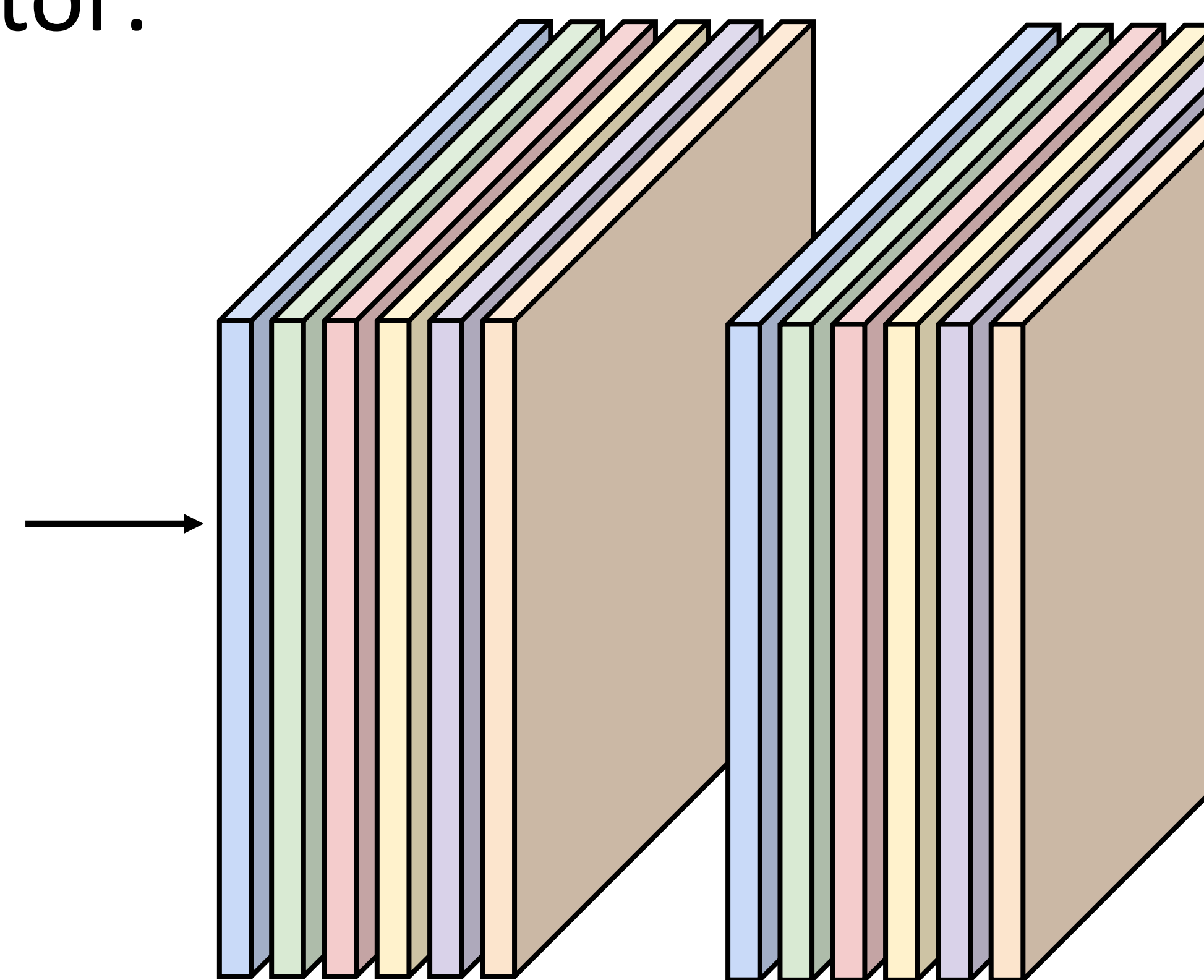
Also 6-dim bias vector:



$6 \times 3 \times 5 \times 5$
filters

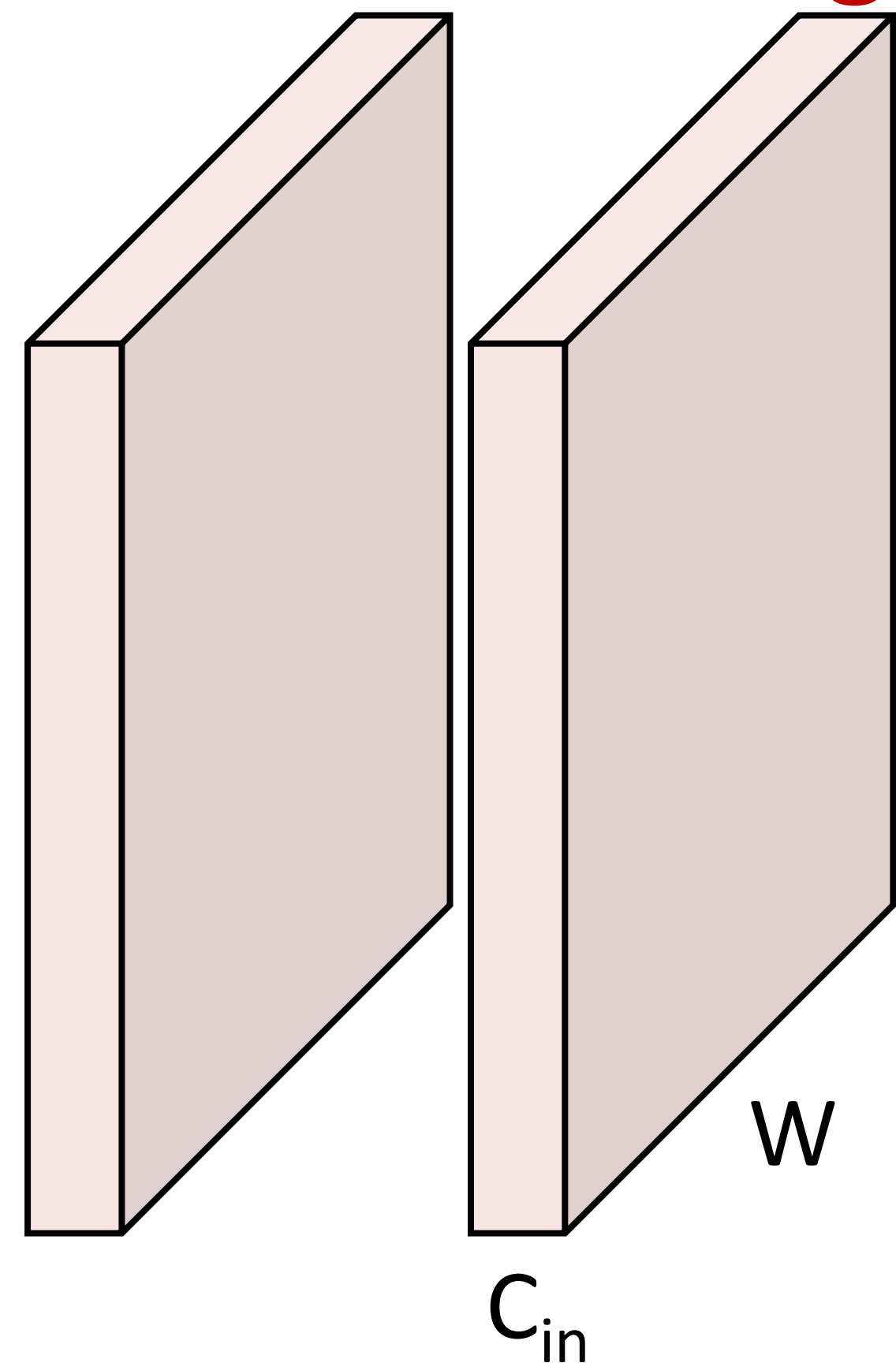


$2 \times 6 \times 28 \times 28$
Batch of outputs

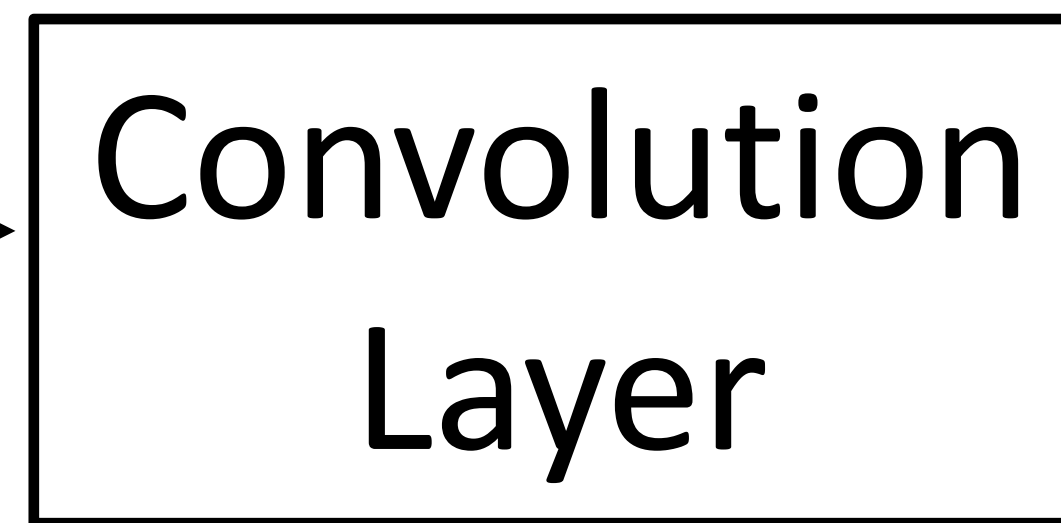
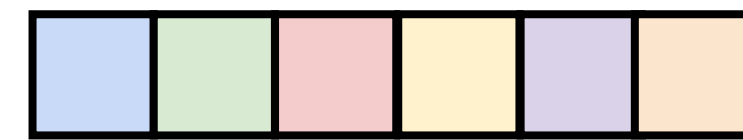


Convolution Layer

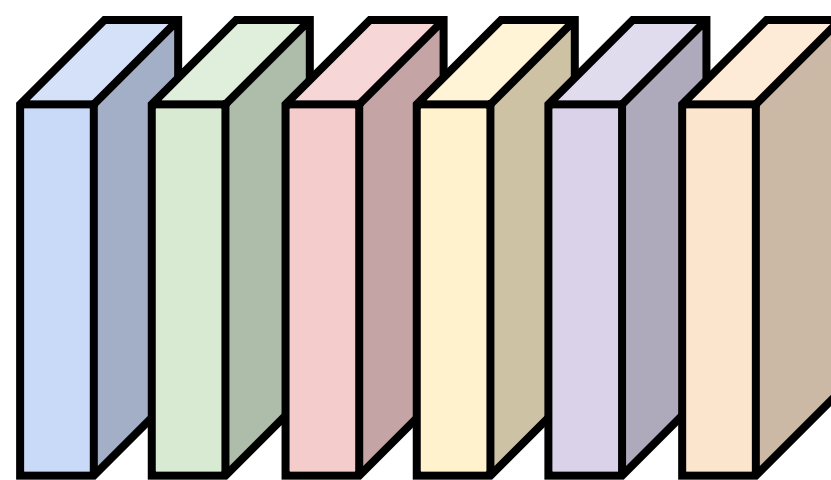
$N \times C_{in} \times H \times W$
Batch of images



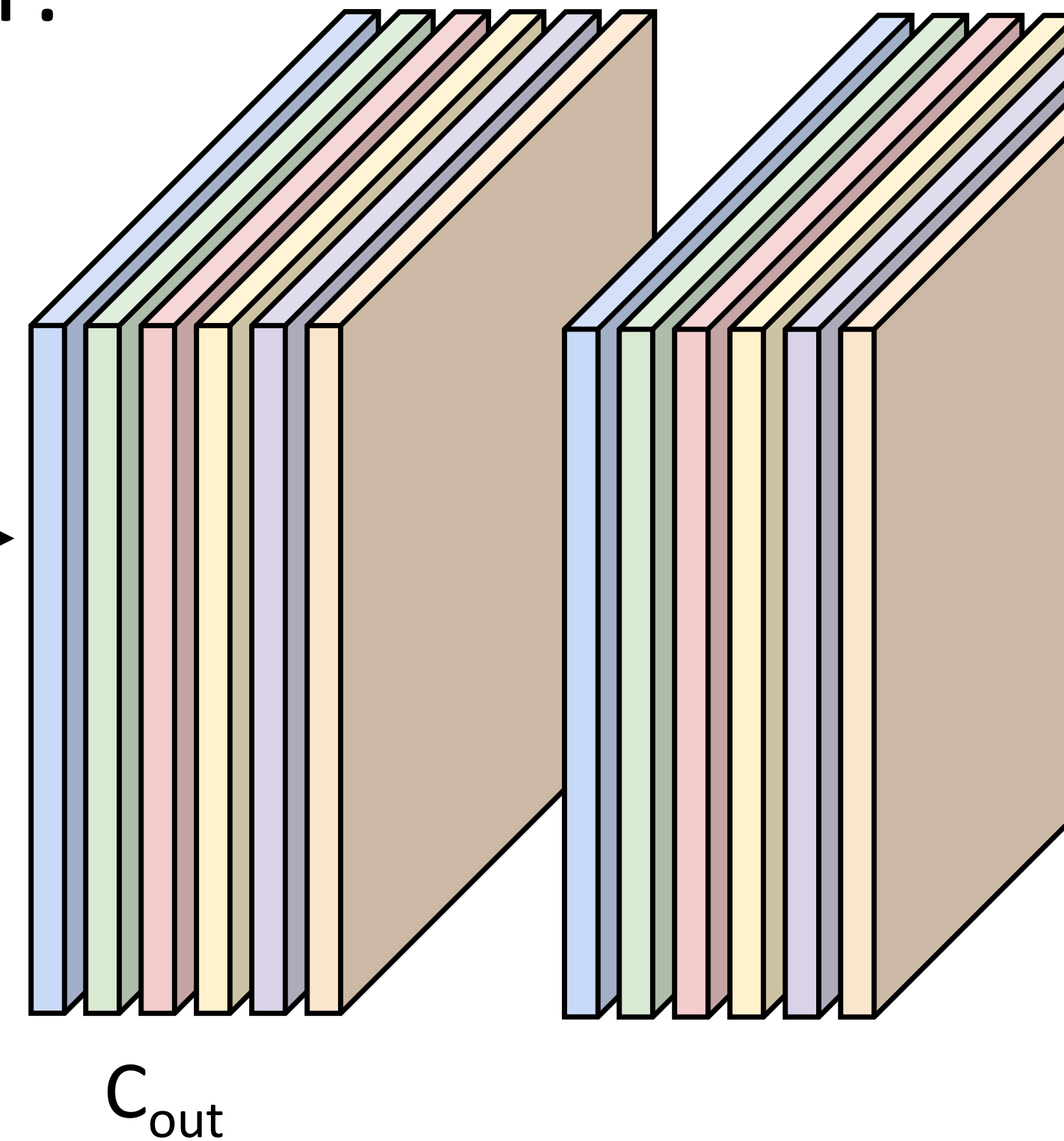
Also C_{out} -dim bias vector:



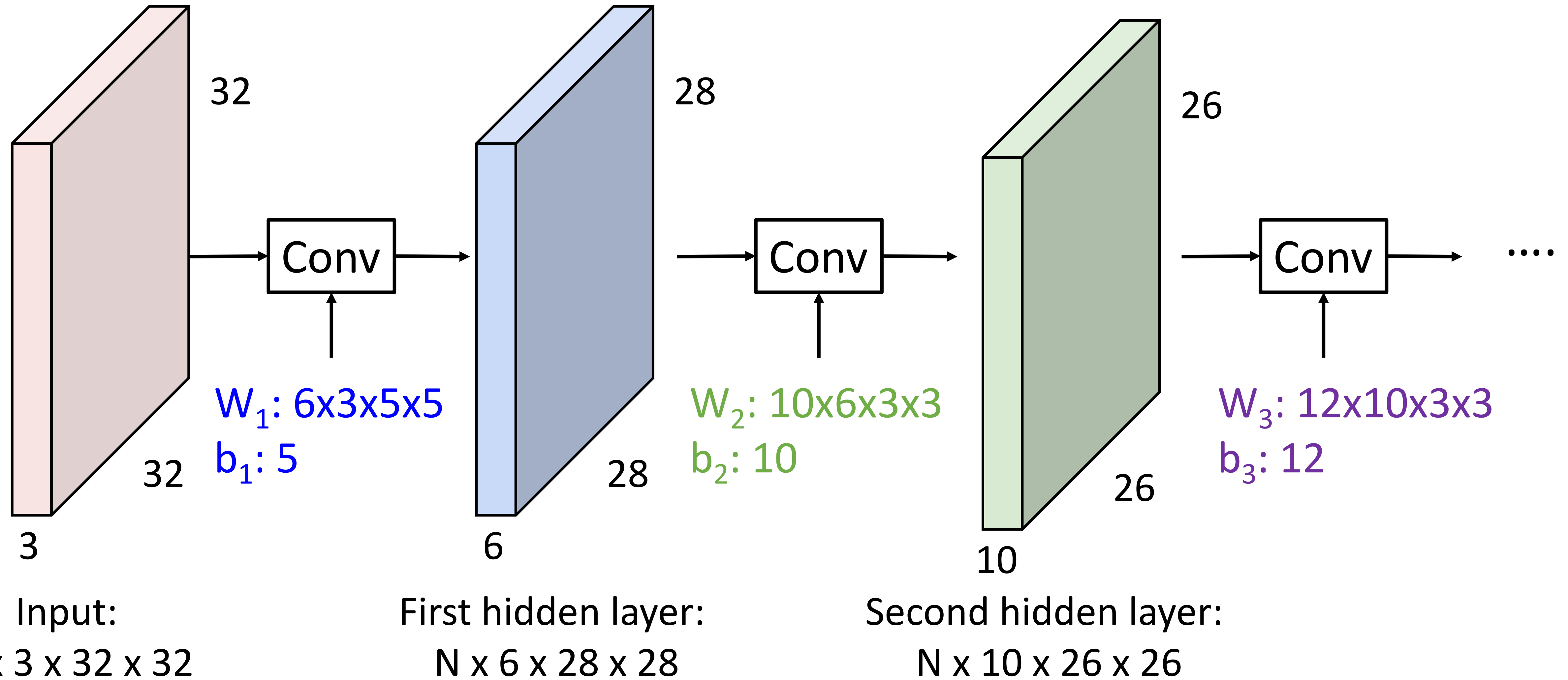
$C_{out} \times C_{in} \times K_w \times K_h$
filters



$N \times C_{out} \times H' \times W'$
Batch of outputs



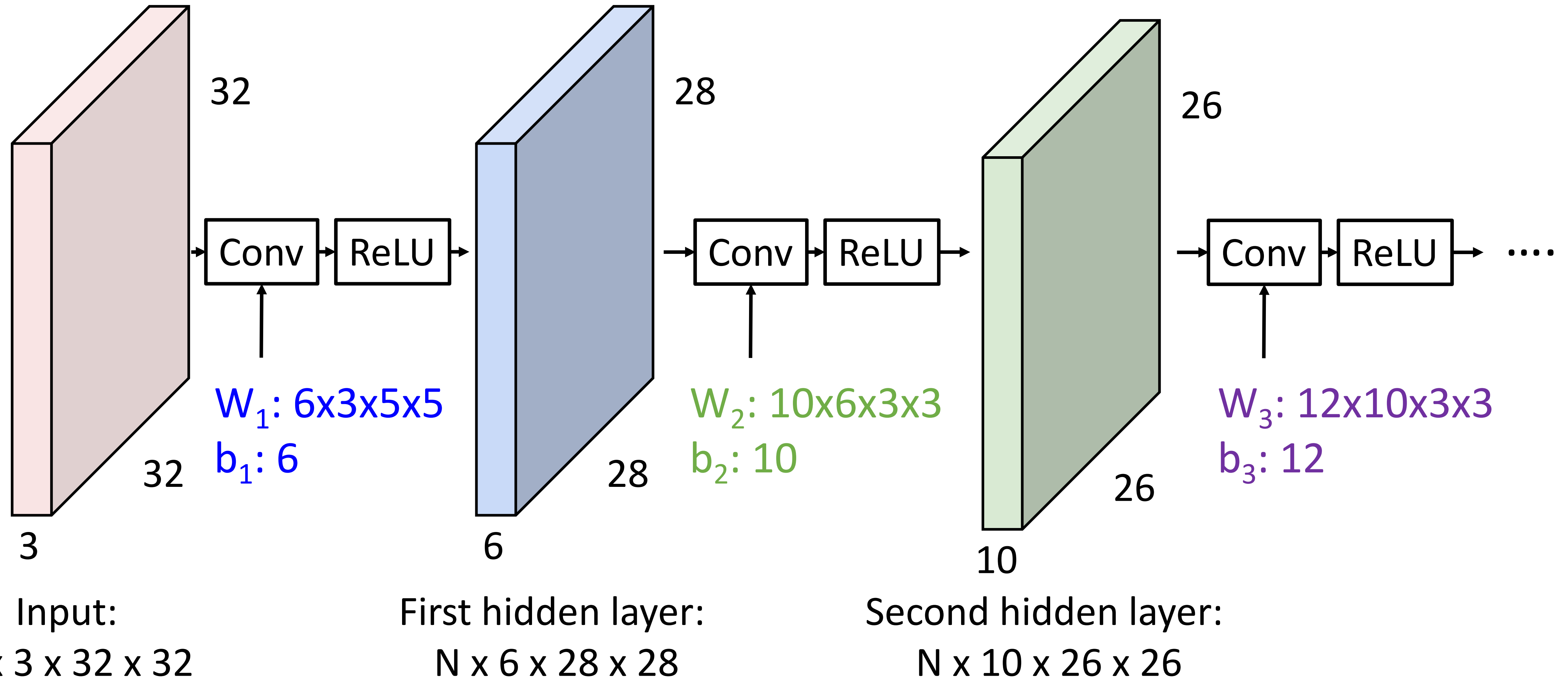
Stacking Convolutions



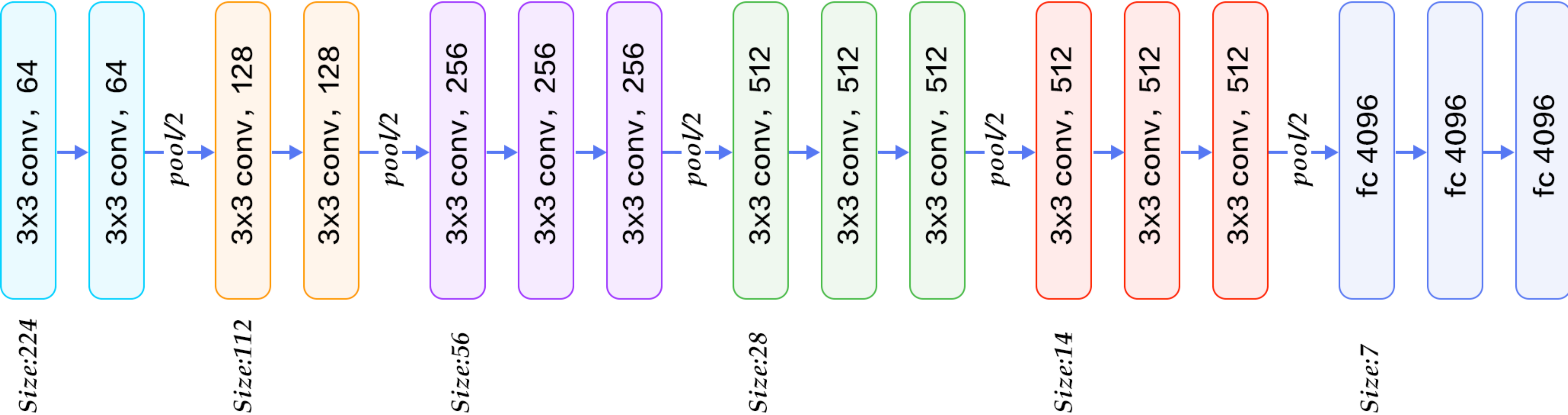
Stacking Convolutions

Q: What happens if we stack two convolution layers? (Recall $y=W_2W_1x$ is a linear classifier)

A: We get another convolution!



Convolutional Neural Networks



VGG-16 Network

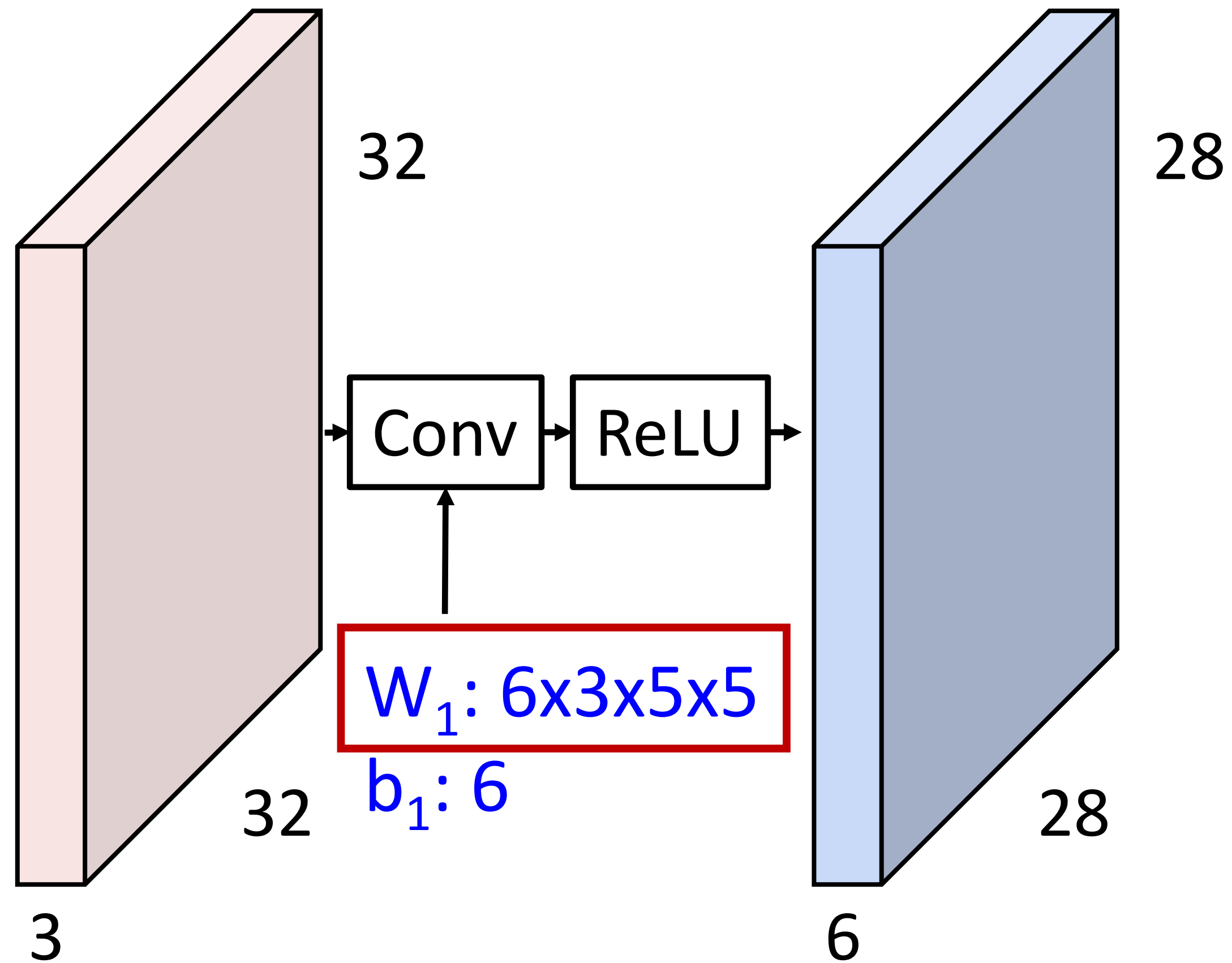
Backward Pass for Some Common Layers

Convolutional layer



20.3

What do convolutional filters learn?



Input:

$N \times 3 \times 32 \times 32$

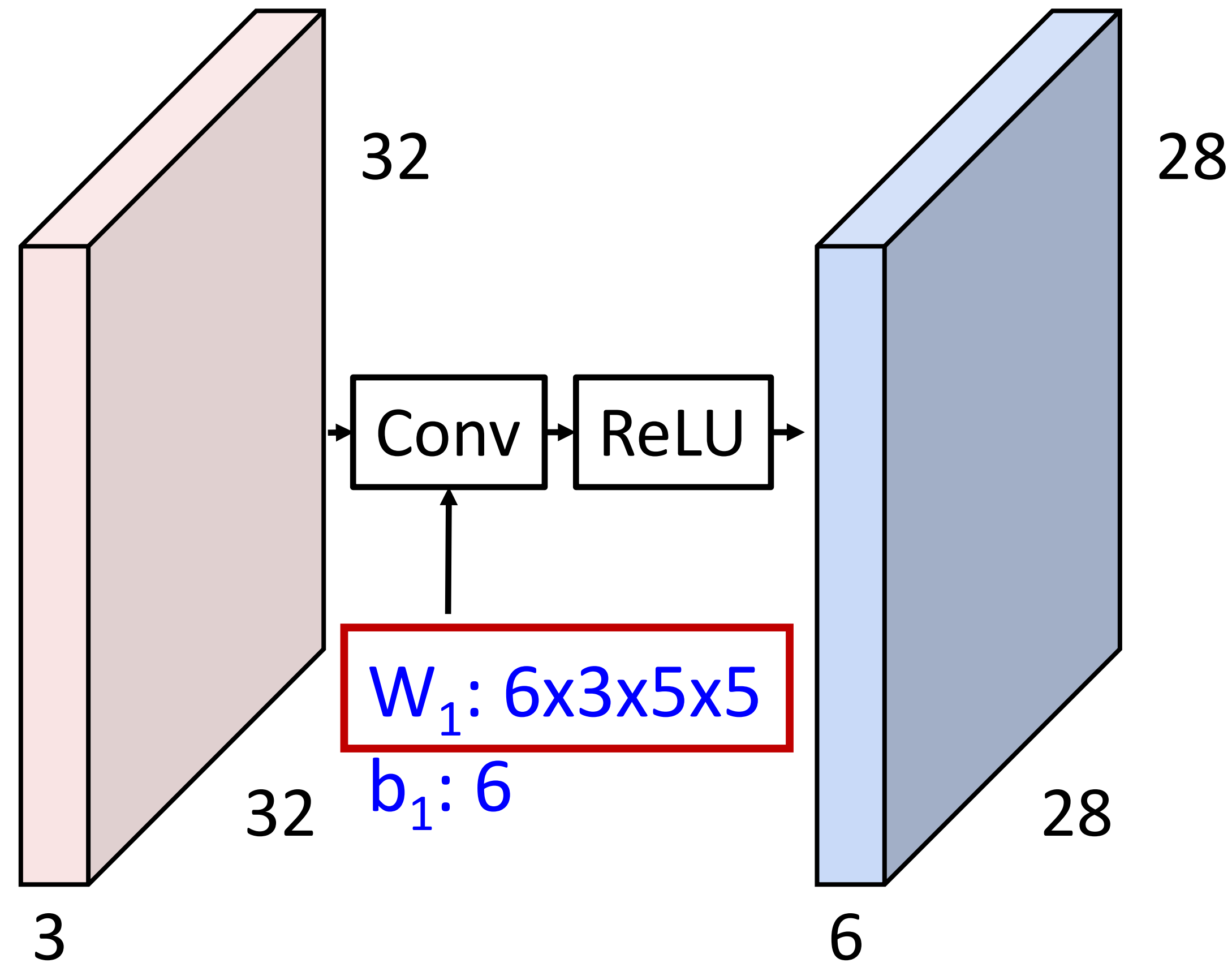
First hidden layer:

$N \times 6 \times 28 \times 28$

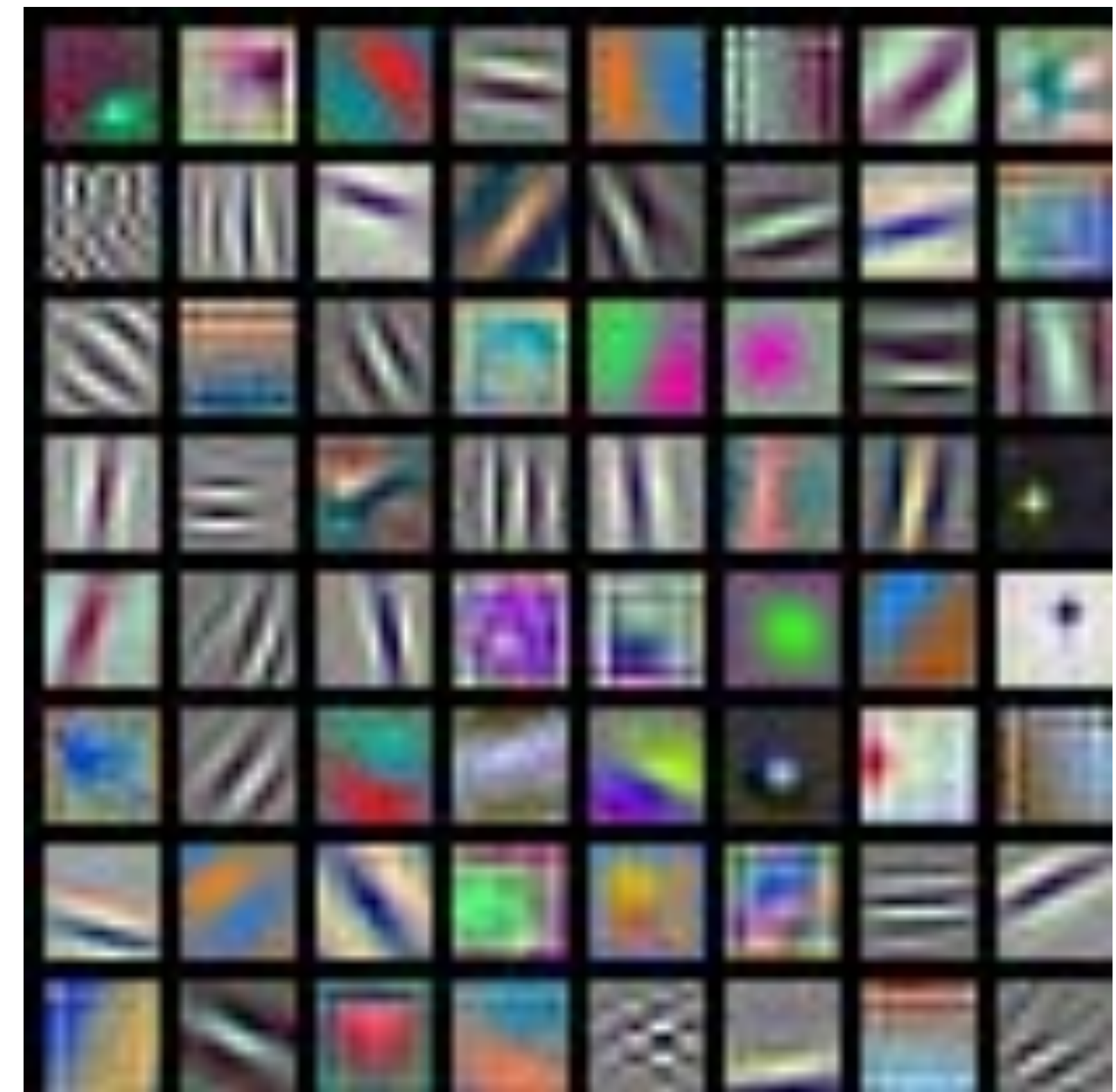
Linear classifier: One template per class



What do convolutional filters learn?

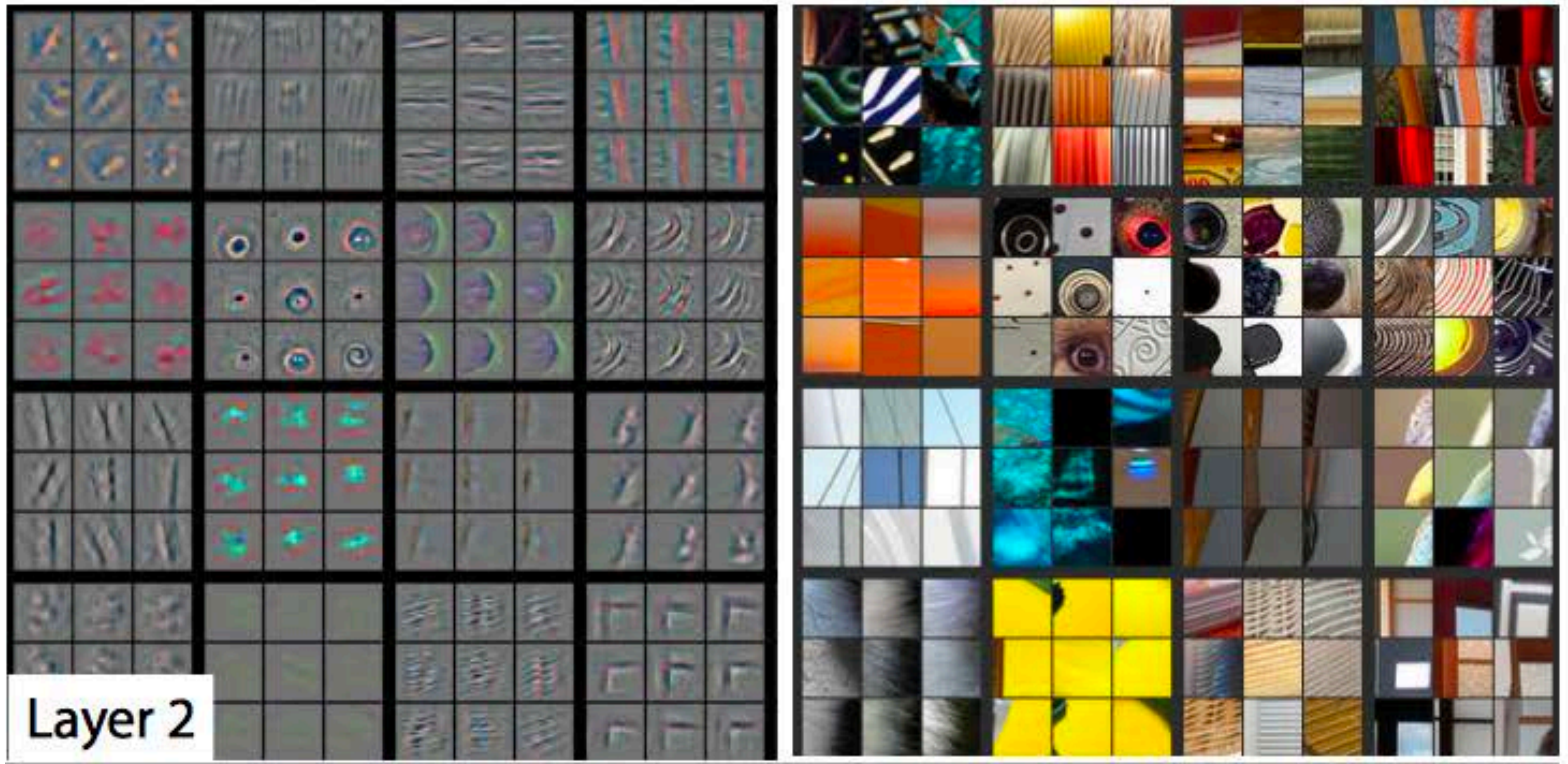


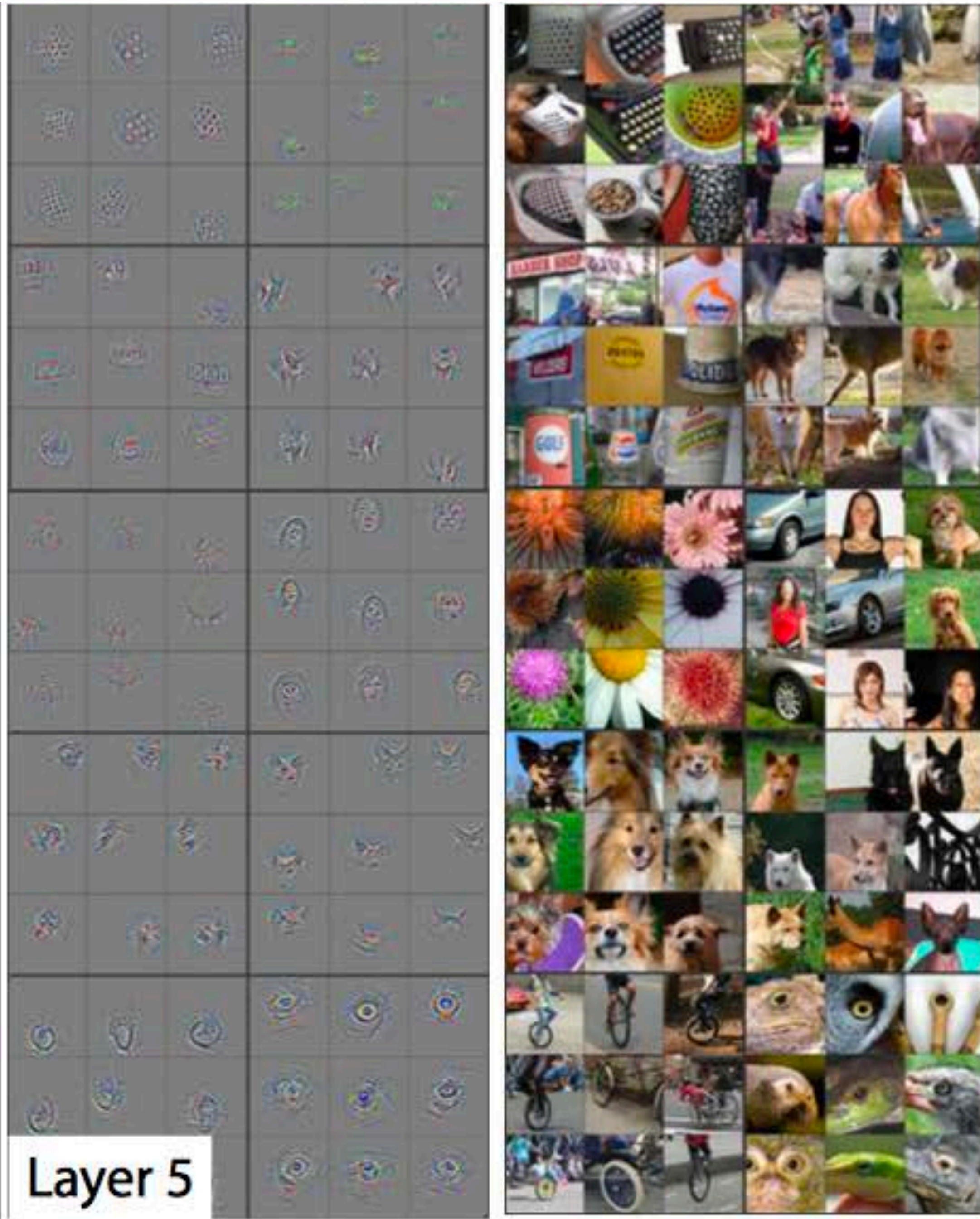
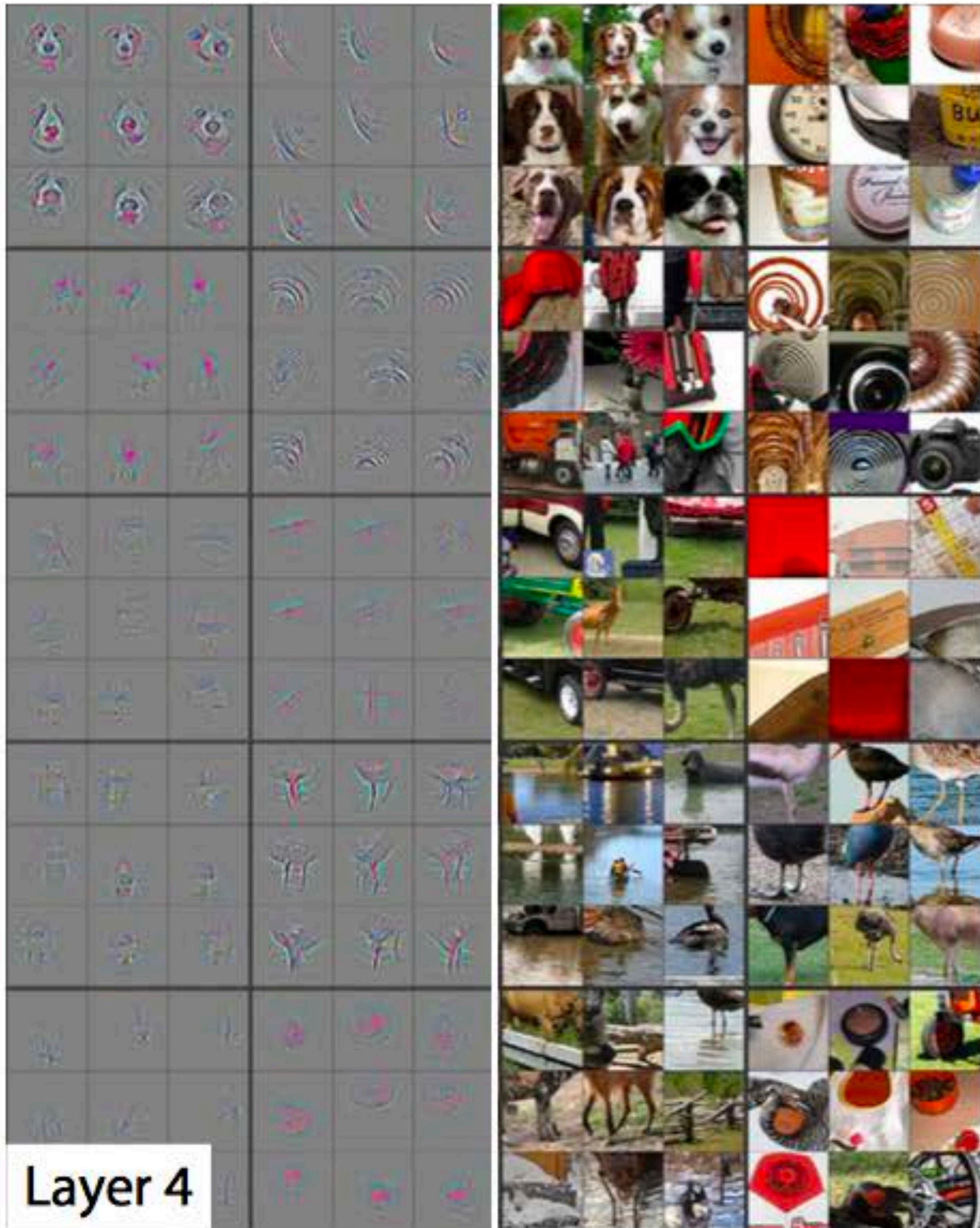
First-layer conv filters: local image templates
(Often learns oriented edges, opposing colors)



AlexNet: 64 filters, each $3 \times 11 \times 11$

What **filters** do networks learn?

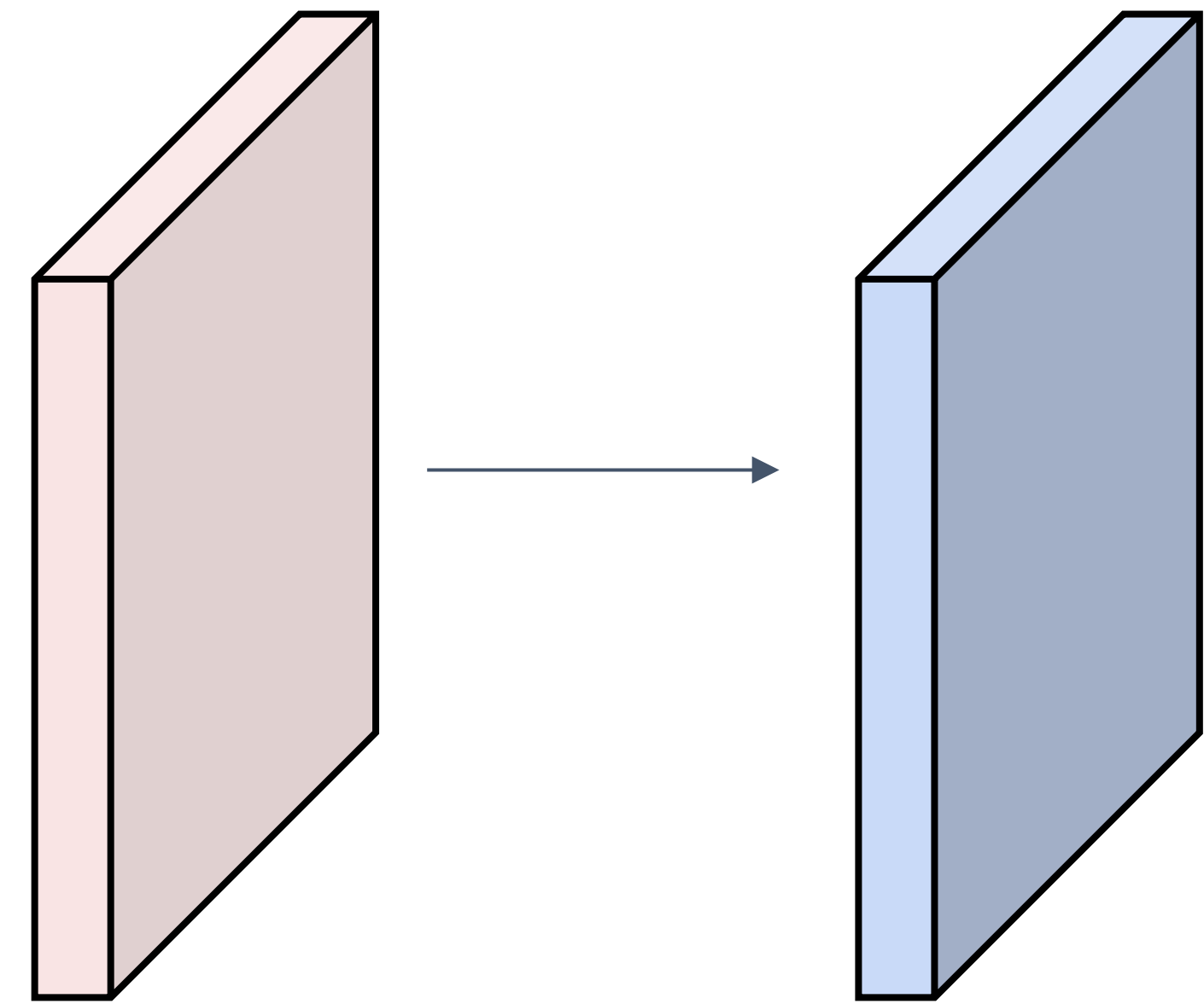




Convolution Example

Input volume: $3 \times 32 \times 32$
10 5×5 filters with stride 1, pad 2

Output volume size: ?



Convolution Example

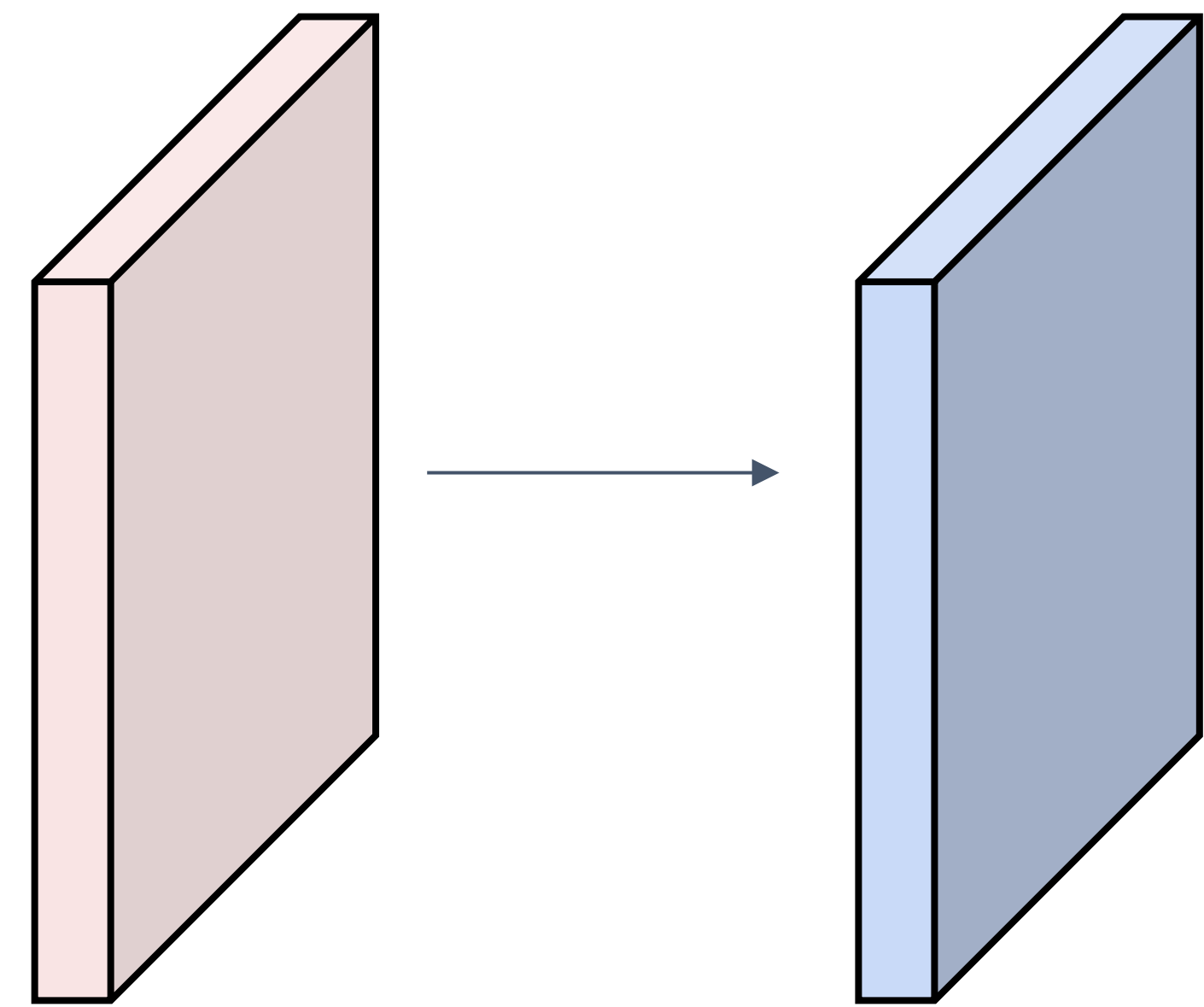
Input volume: 3 x **32** x **32**

10 **5x5** filters with stride **1**, pad **2**

Output volume size:

$(\mathbf{32} + 2 * \mathbf{2} - \mathbf{5}) / \mathbf{1} + 1 = 32$ spatially, so

10 x 32 x 32



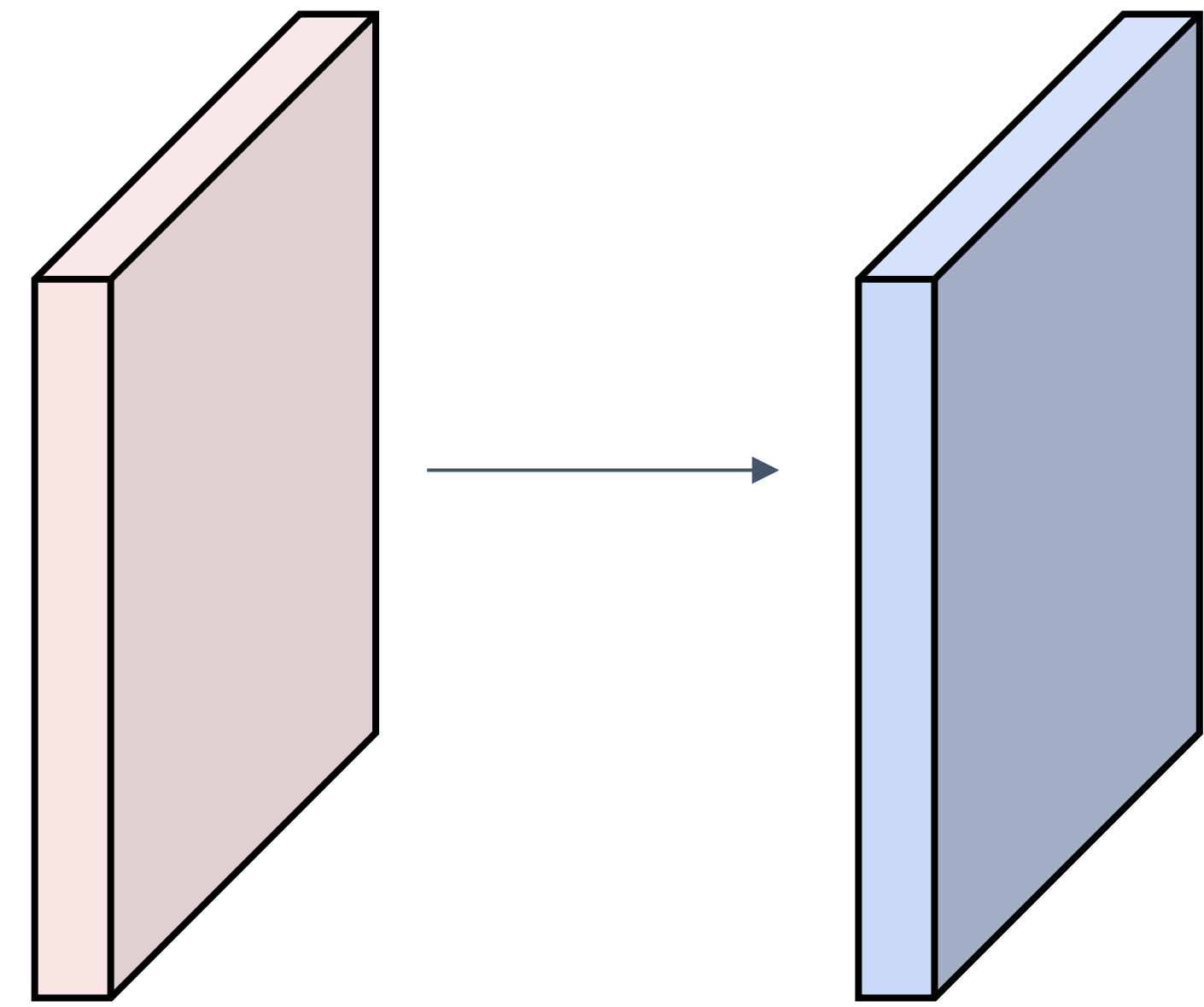
Convolution Example

Input volume: $3 \times 32 \times 32$

10 5×5 filters with stride 1, pad 2

Output volume size: $10 \times 32 \times 32$

Number of learnable parameters: ?



Convolution Example

Input volume: **3** x 32 x 32

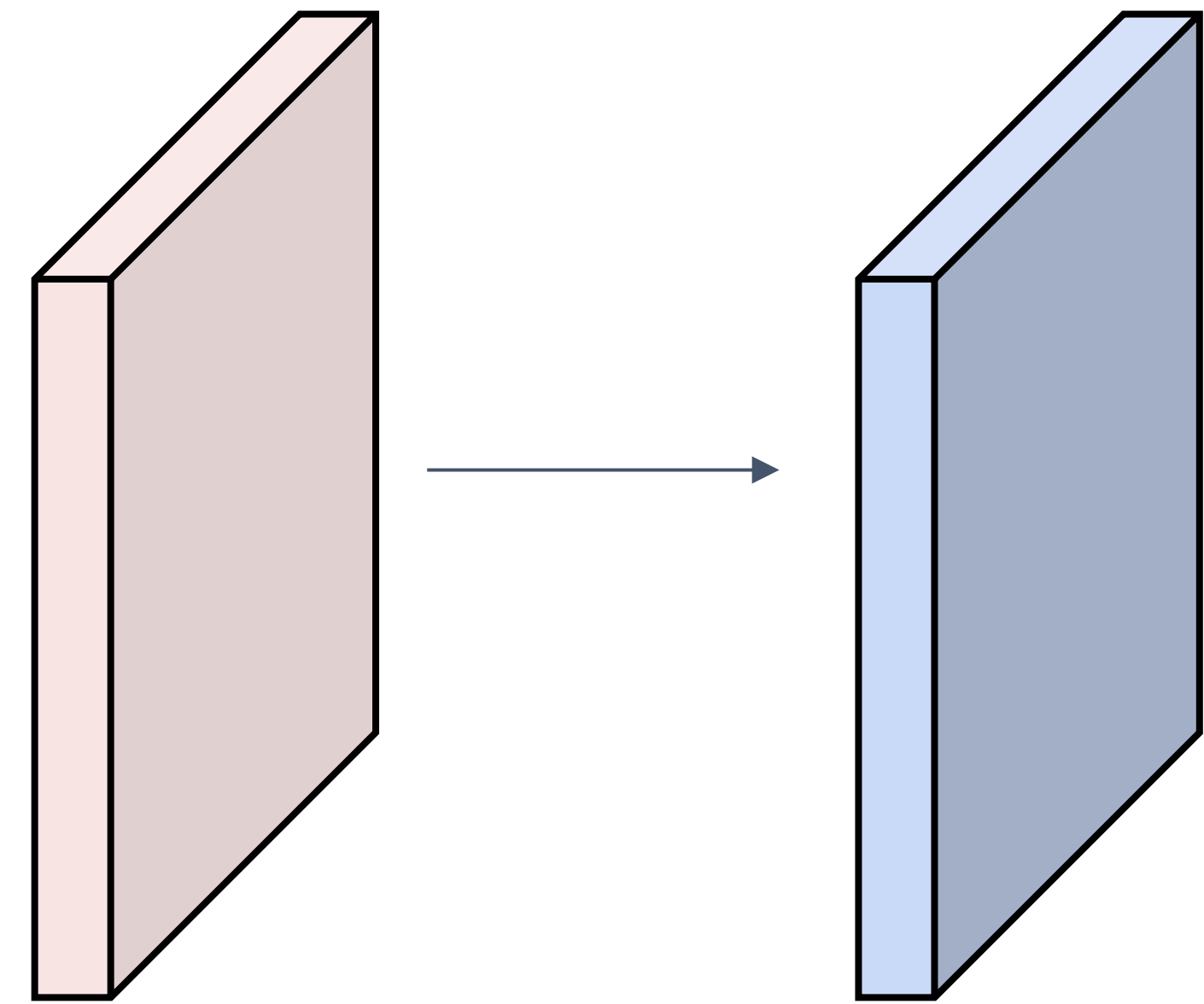
10 **5x5** filters with stride 1, pad 2

Output volume size: 10 x 32 x 32

Number of learnable parameters: **760**

Parameters per filter: **3*****5*****5** + 1 (for bias) = **76**

10 filters, so total is **10** * **76** = **760**



Convolution Example

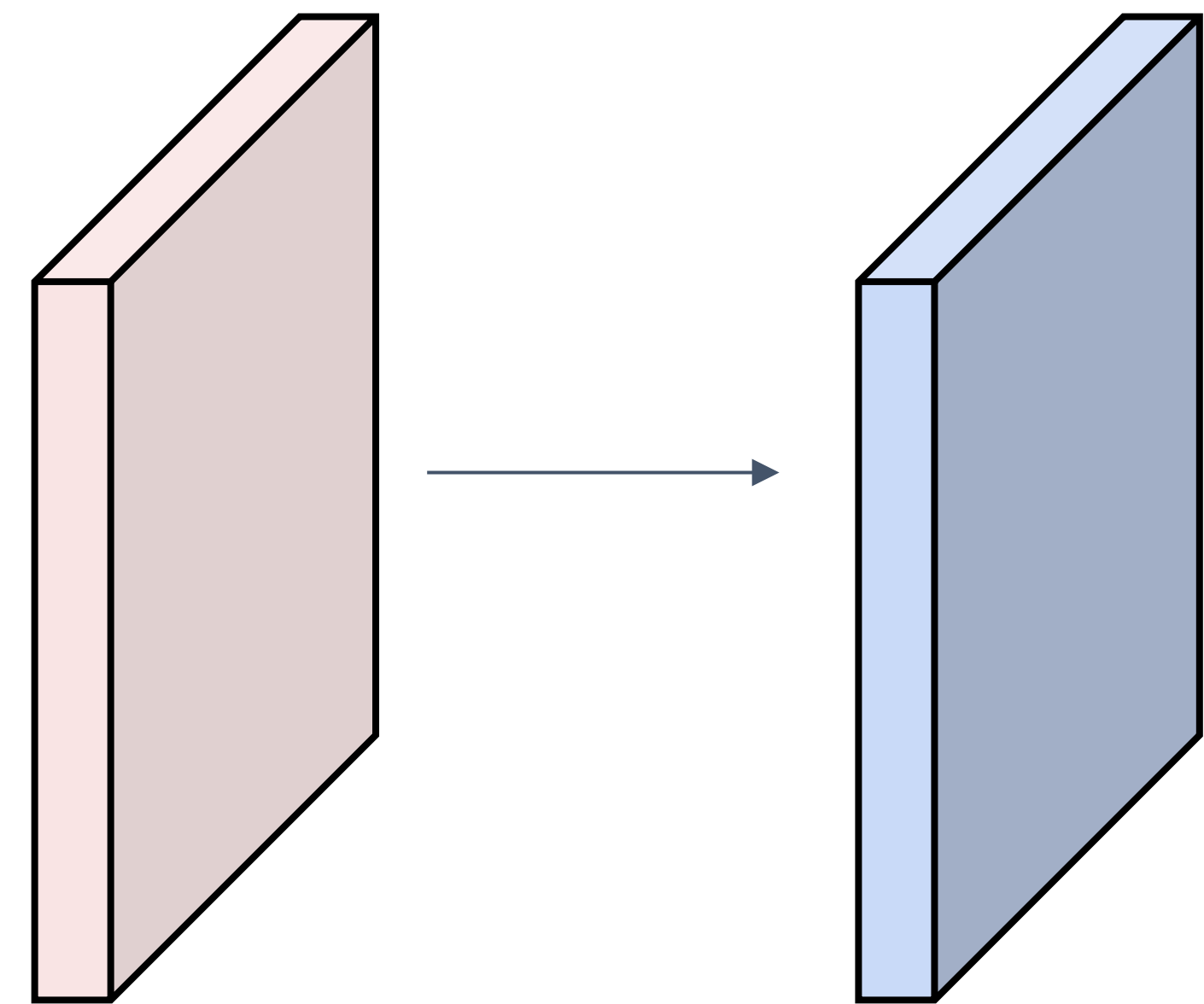
Input volume: $3 \times 32 \times 32$

10 5×5 filters with stride 1, pad 2

Output volume size: $10 \times 32 \times 32$

Number of learnable parameters: 760

Number of multiply-add operations: ?



Convolution Example

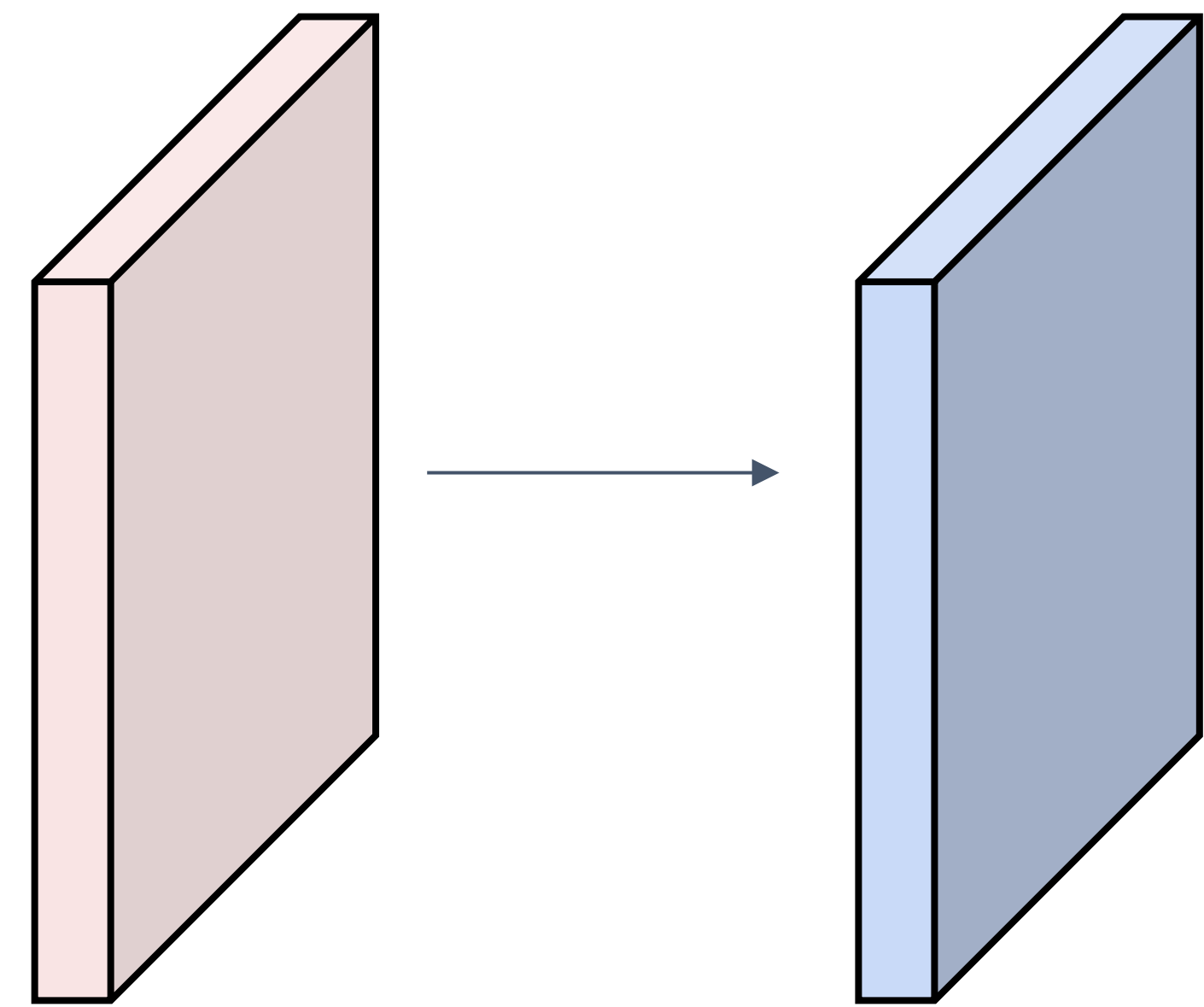
Input volume: **3** x 32 x 32
10 **5x5** filters with stride 1, pad 2

Output volume size: **10** x 32 x 32

Number of learnable parameters: 760

Number of multiply-add operations: **768,000**

10*32*32 = 10,240 outputs; each output is the inner product of two **3x5x5** tensors (75 elems); total = $75 * 10240 = 768K$



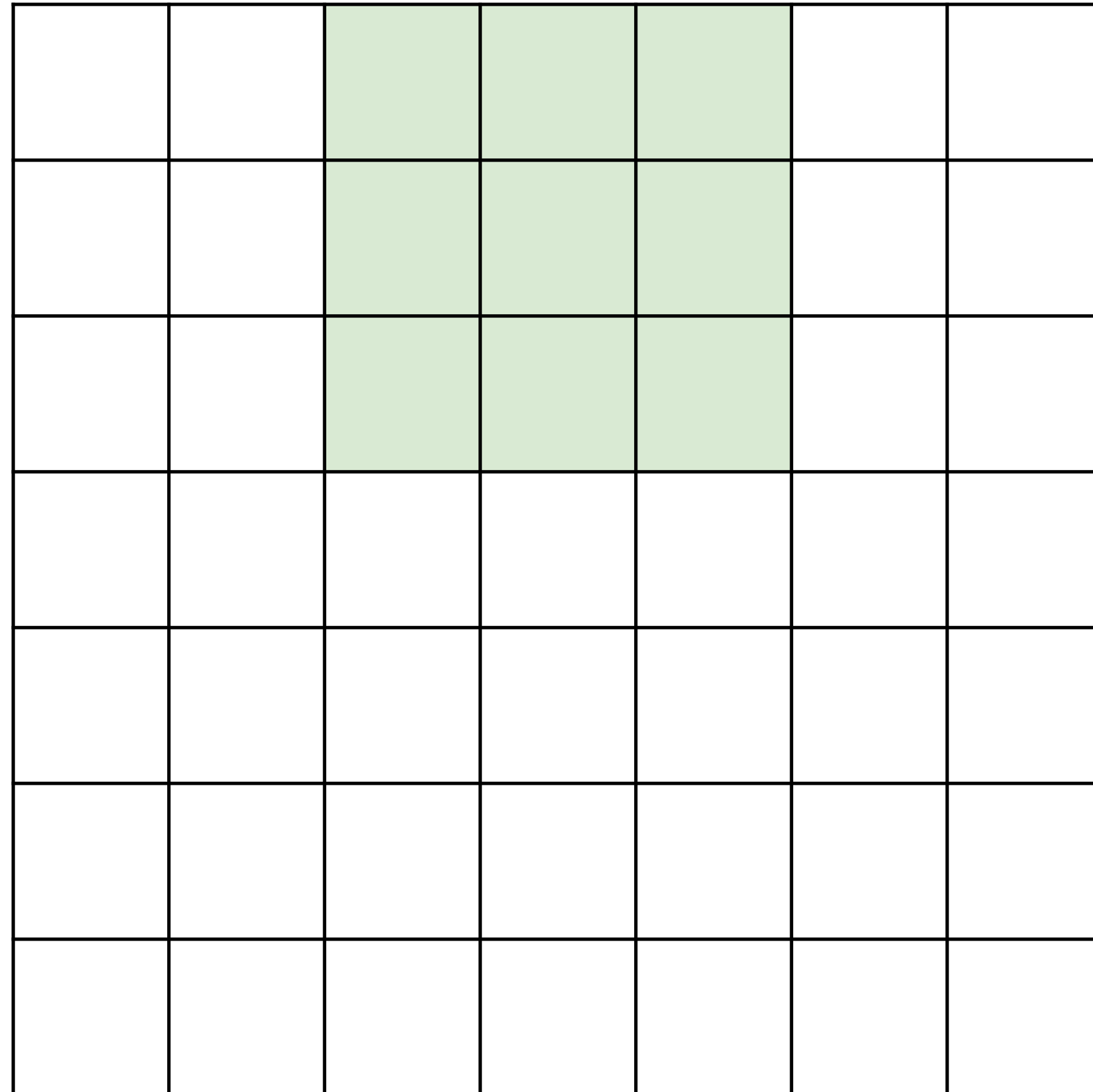
Strided Convolution

Input: 7×7

Filter: 3×3

Stride: 2

Strided Convolution

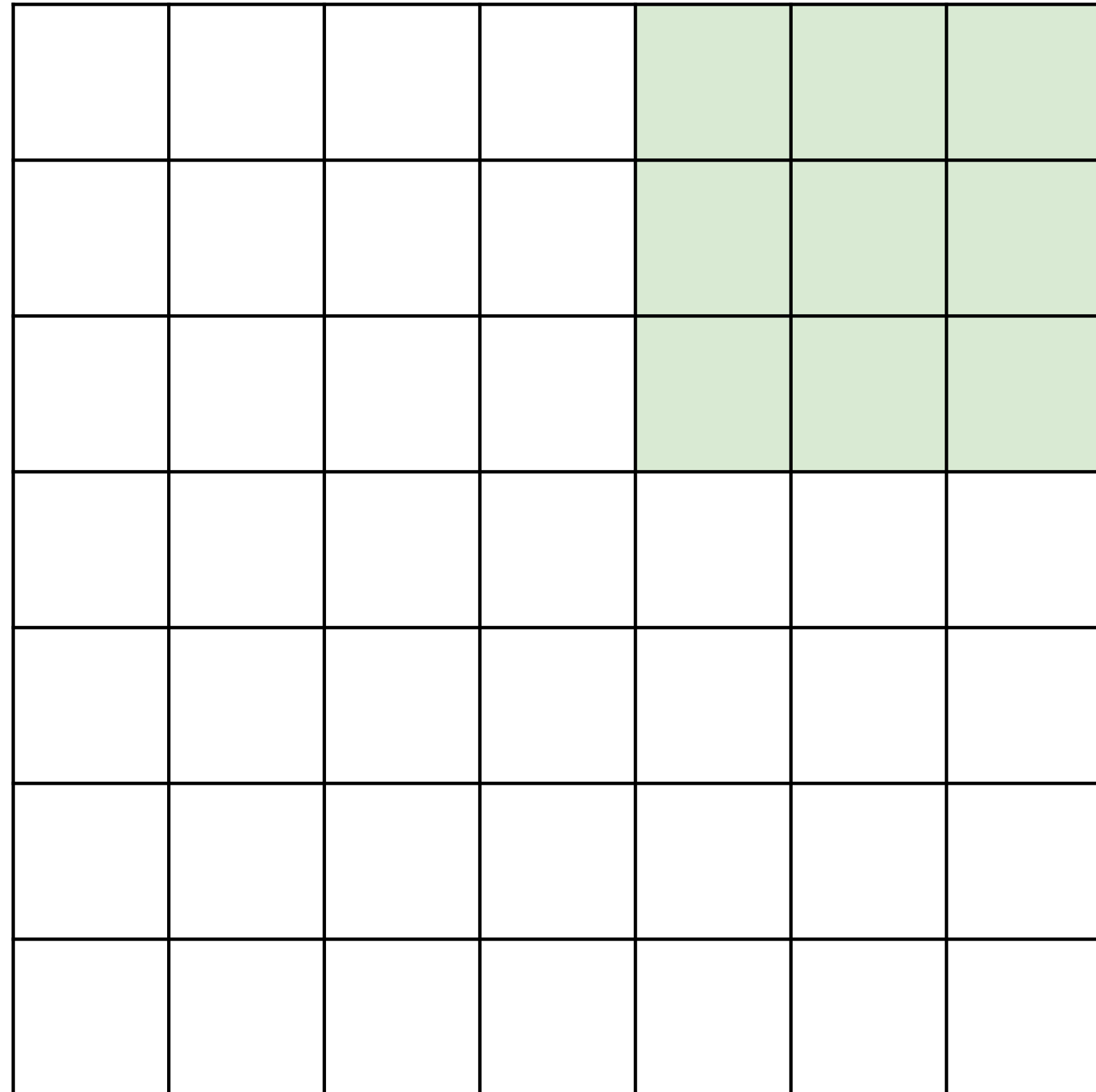


Input: 7×7

Filter: 3×3

Stride: 2

Strided Convolution



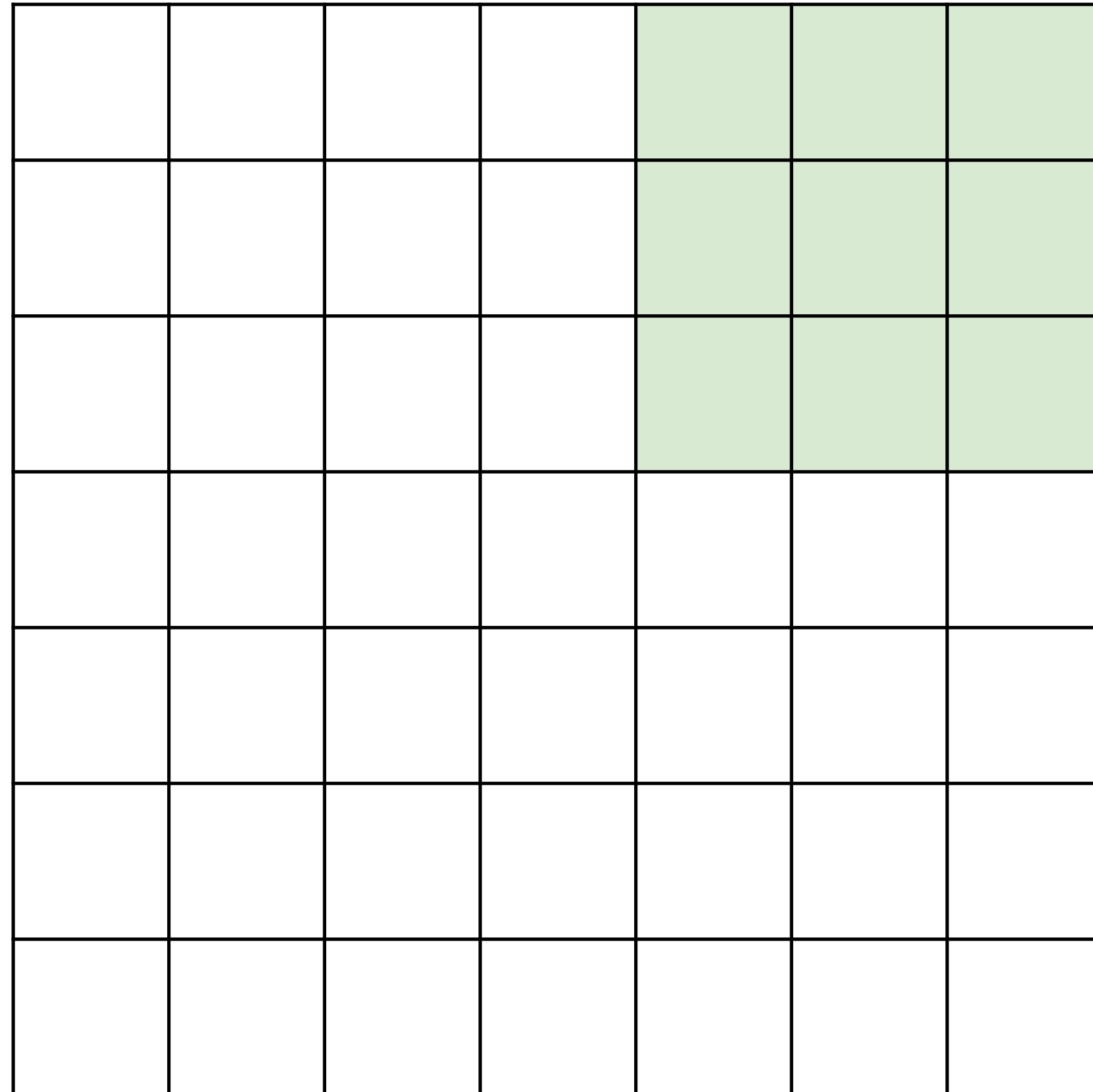
Input: 7x7

Filter: 3x3

Stride: 2

Output: 3x3

Strided Convolution



Input: 7x7

Filter: 3x3

Stride: 2

Output: 3x3

In general:

Input: W

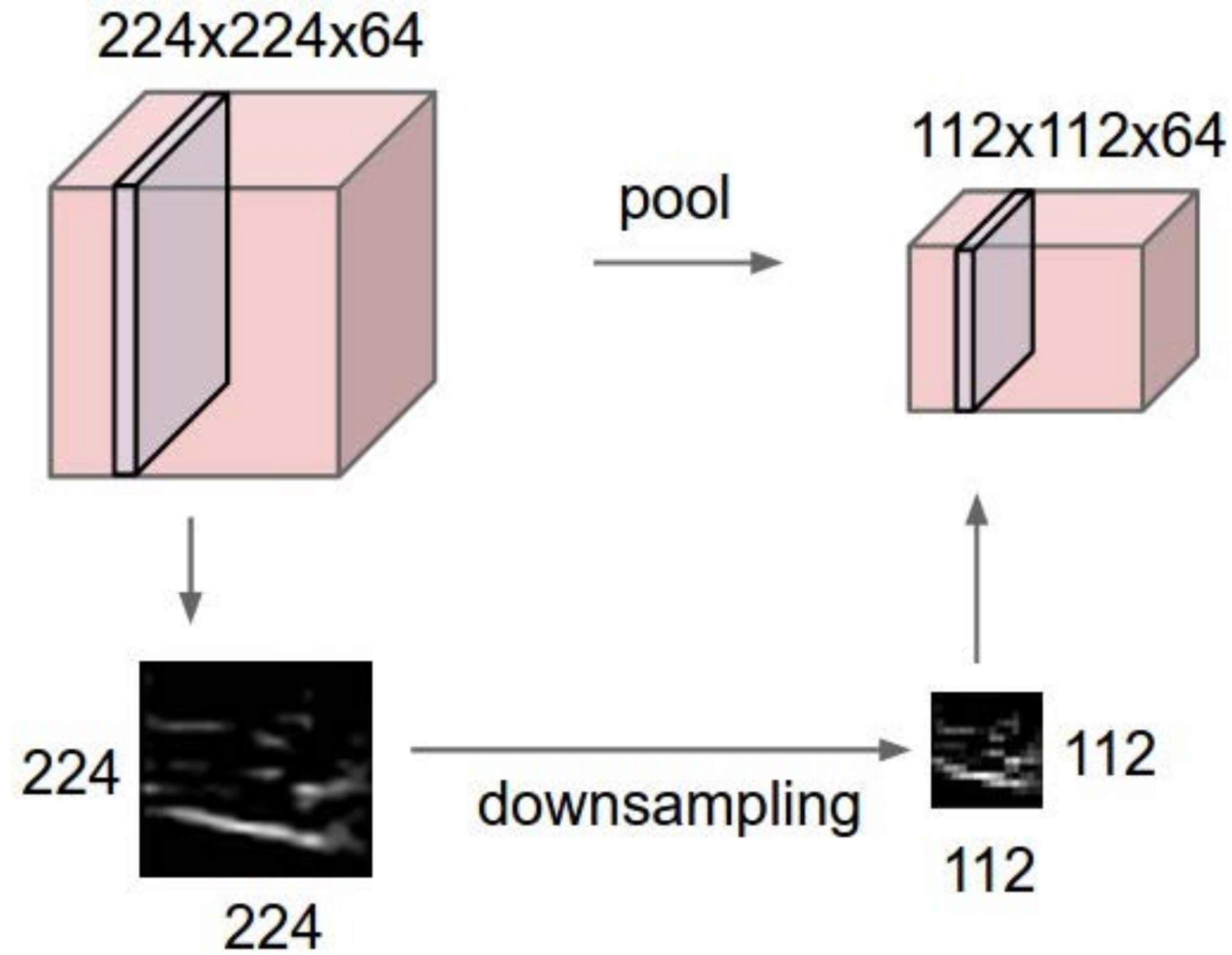
Filter: K

Padding: P

Stride: S

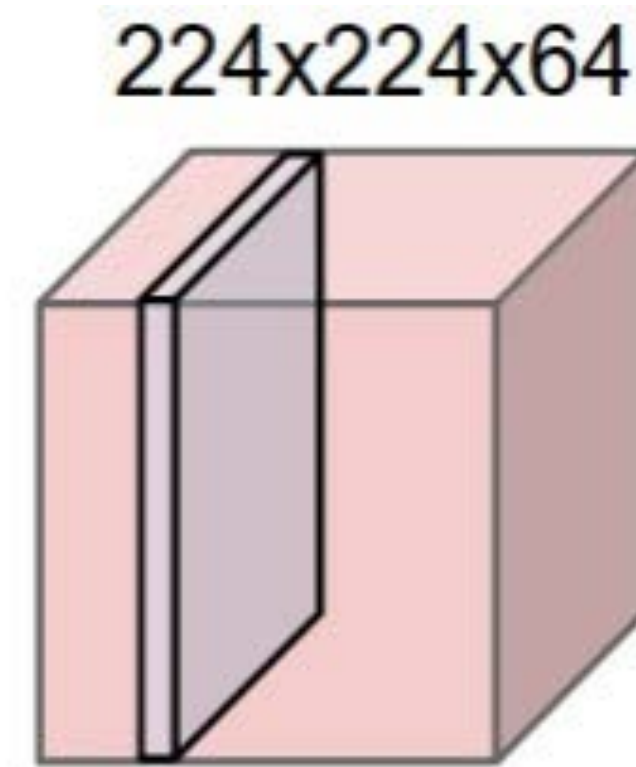
Output: $(W - K + 2P) / S + 1$

Pooling Layers: Another way to downsample

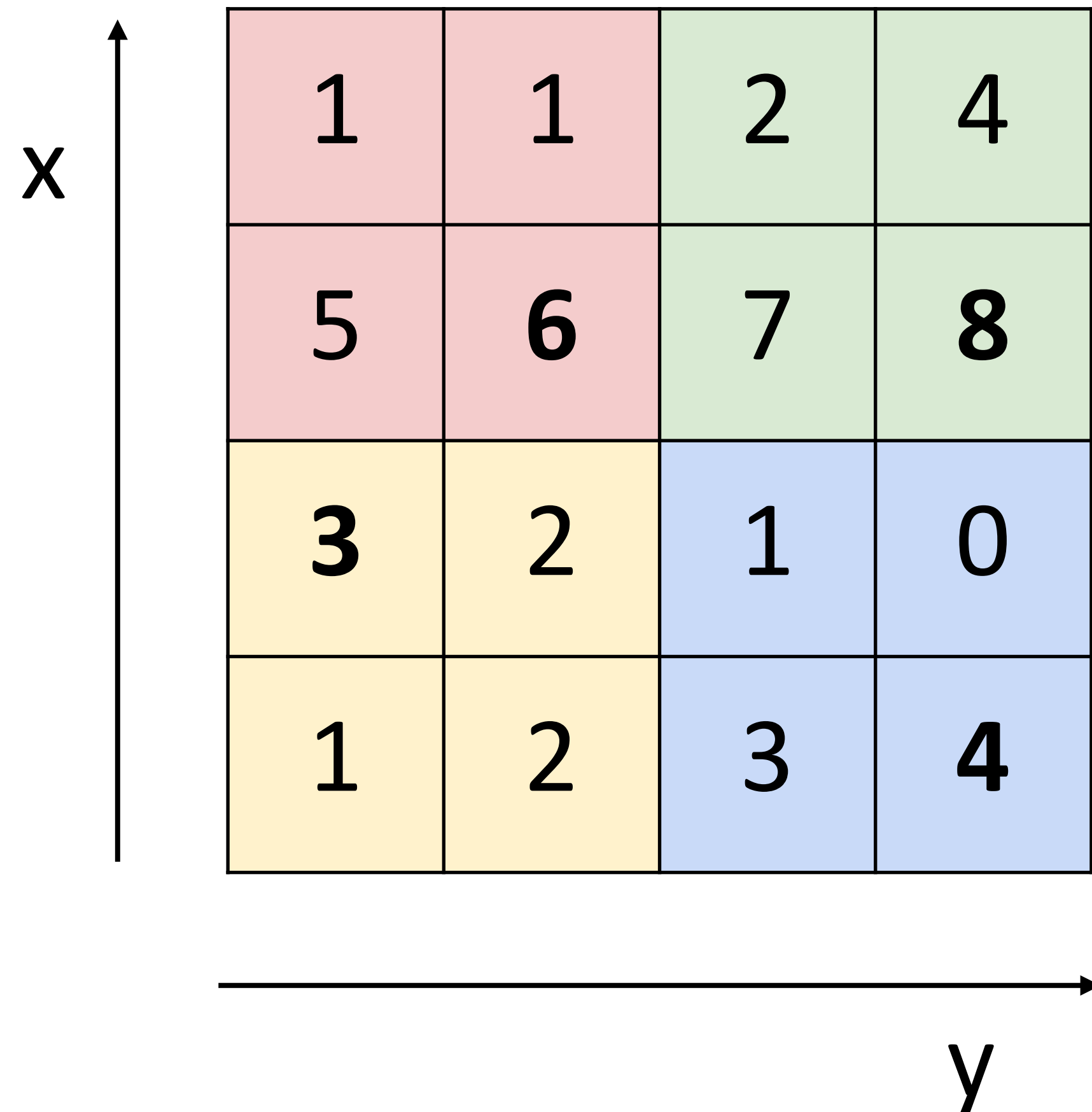


Hyperparameters:
Kernel Size
Stride
Pooling function

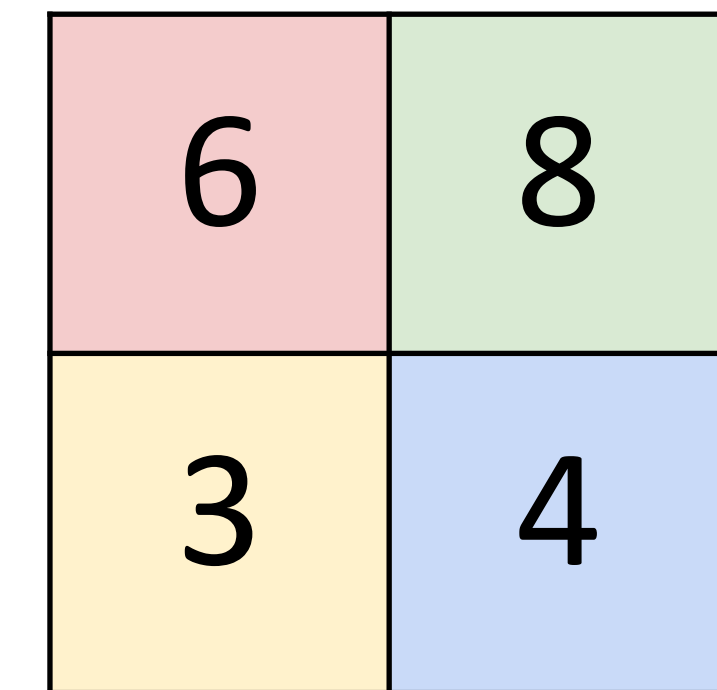
Max Pooling



Single depth slice



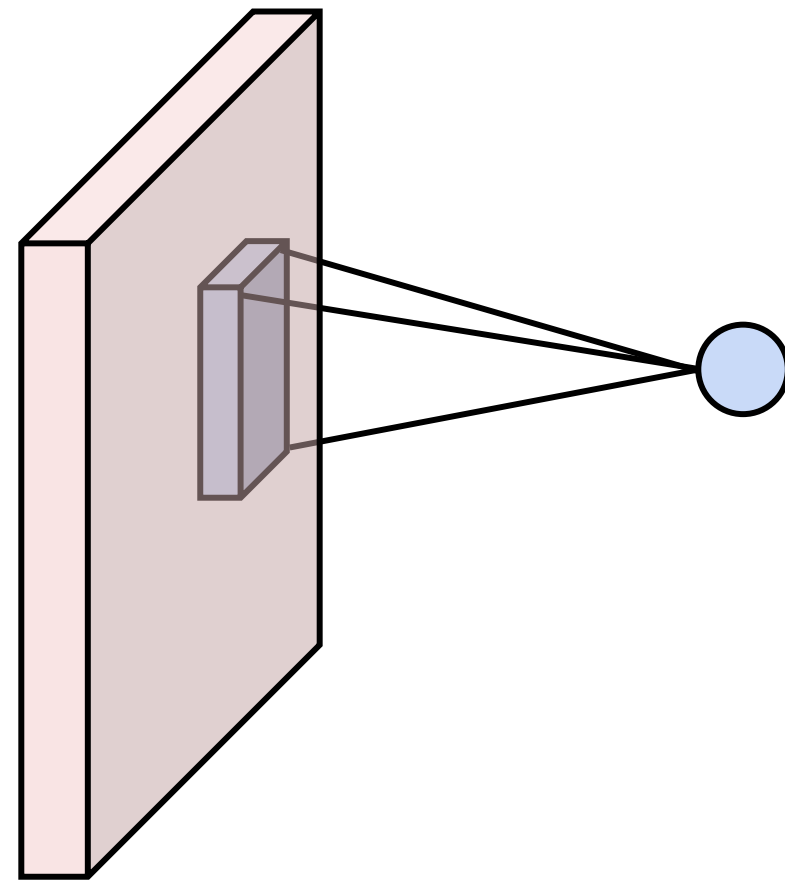
Max pooling with 2x2
kernel size and stride 2



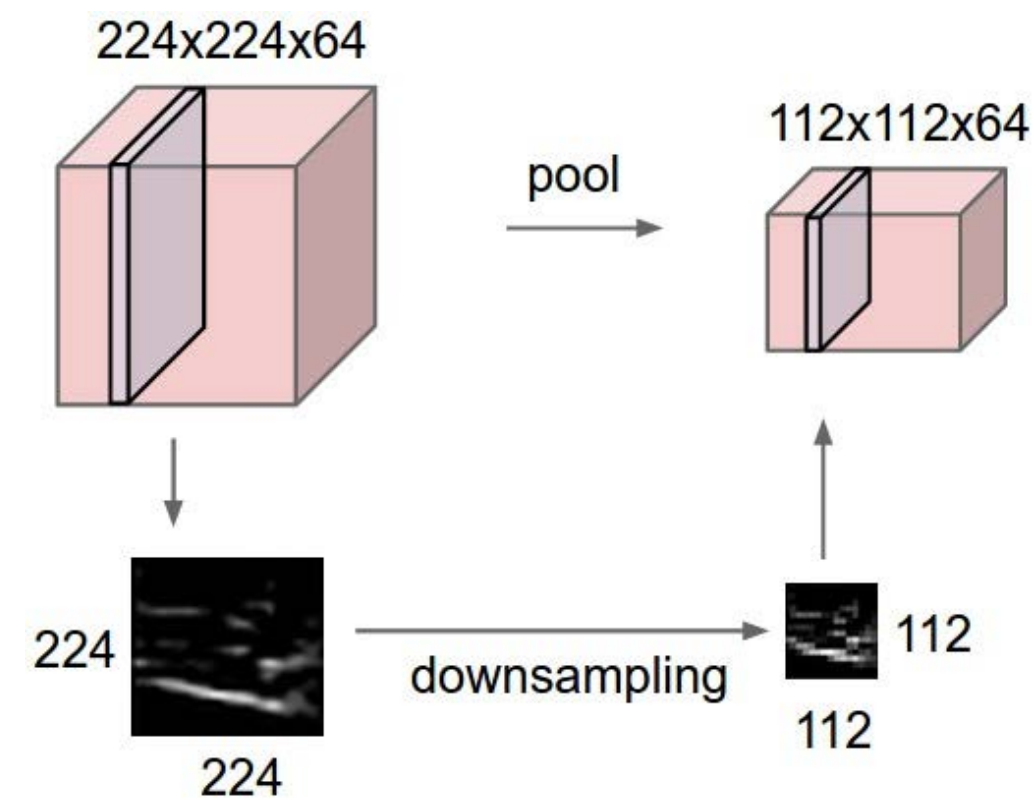
Introduces **invariance** to
small spatial shifts
No learnable parameters!

Components of a Convolutional Network

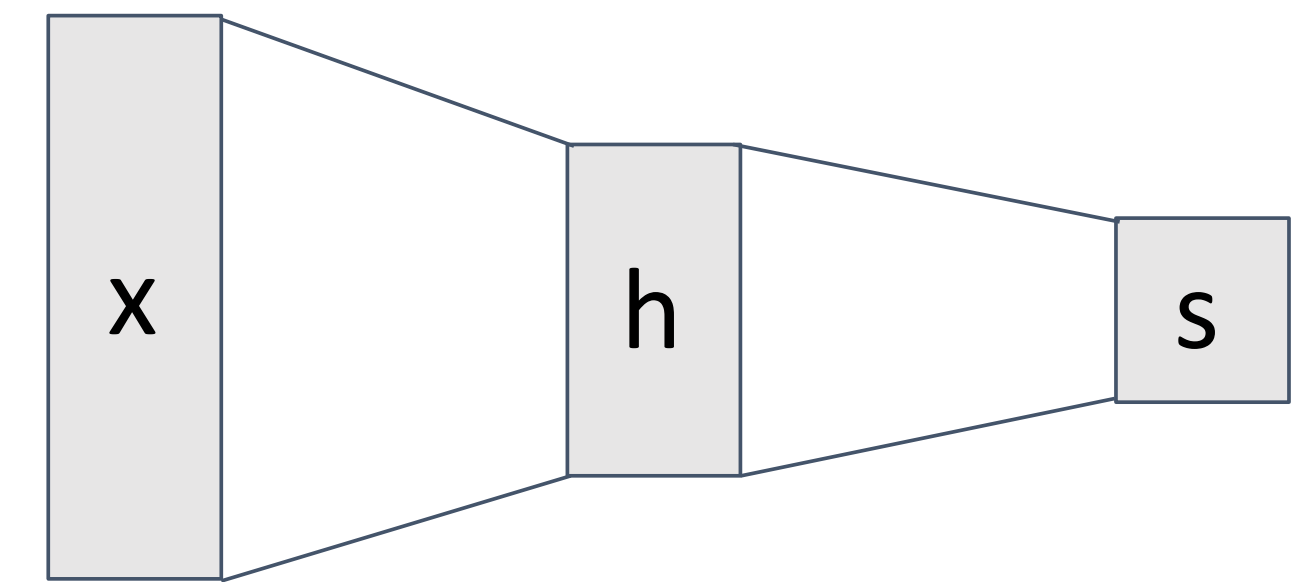
Convolution Layers



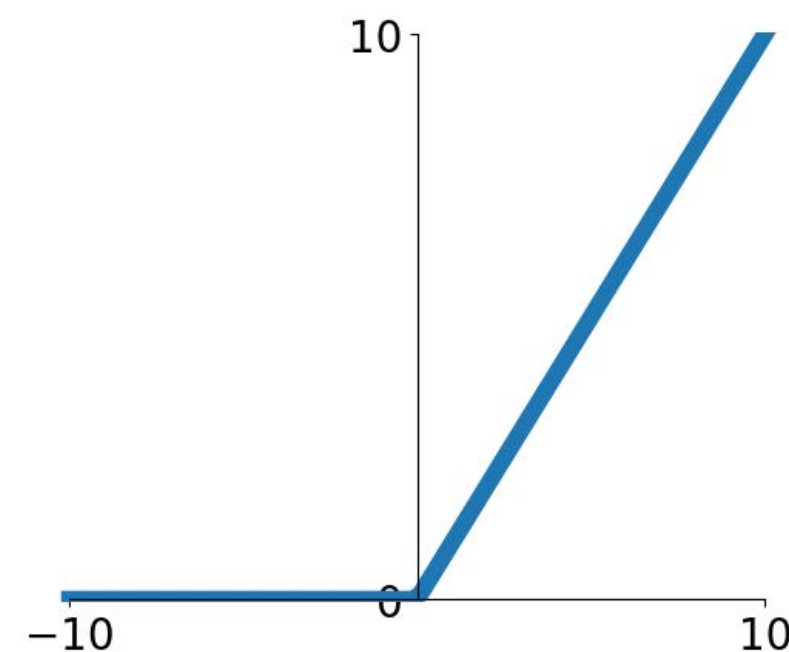
Pooling Layers



Fully-Connected Layers



Activation Function



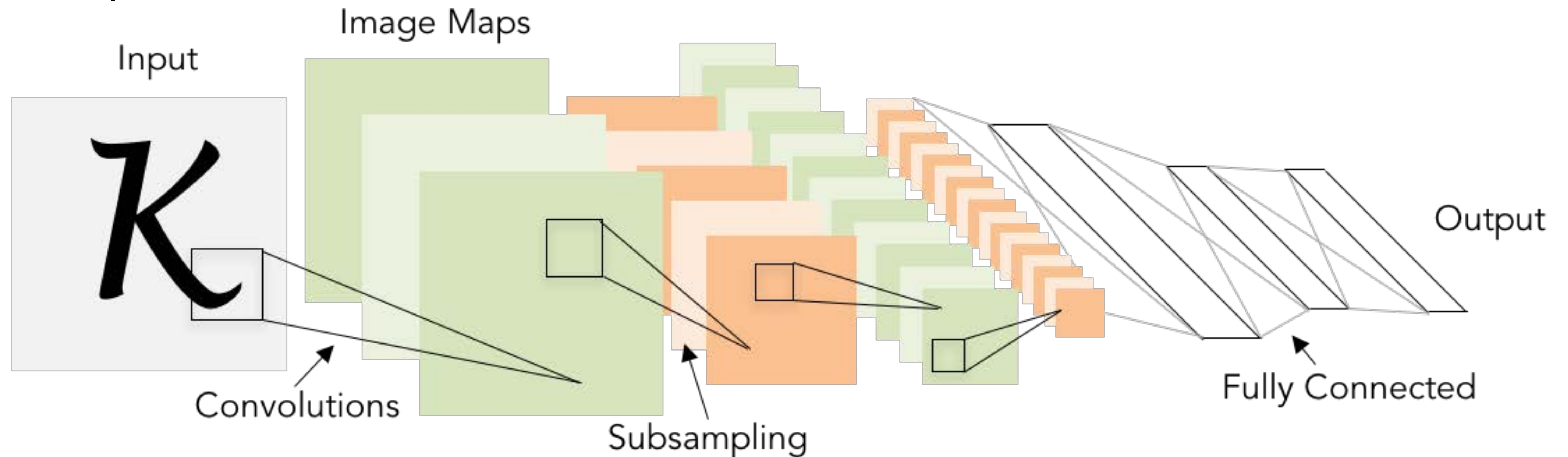
Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Convolutional Networks

Classic architecture: [Conv, ReLU, Pool] x N, flatten, [FC, ReLU] x N, FC

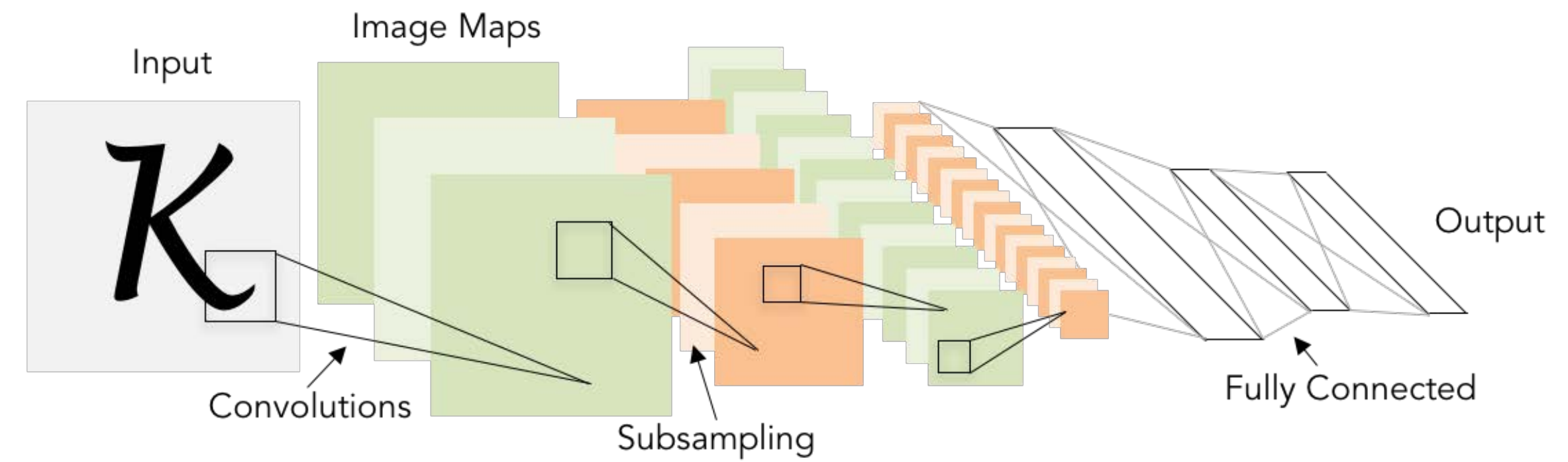
Example: LeNet-5



Lecun et al, "Gradient-based learning applied to document recognition", 1998

Example: LeNet-5

Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv ($C_{out}=20, K=5, P=2, S=1$)	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool($K=2, S=2$)	20 x 14 x 14	
Conv ($C_{out}=50, K=5, P=2, S=1$)	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	
MaxPool($K=2, S=2$)	50 x 7 x 7	
Flatten	2450	
Linear (2450 -> 500)	500	2450 x 500
ReLU	500	
Linear (500 -> 10)	10	500 x 10



As we go through the network:

Spatial size **decreases**
(using pooling or strided conv)

Number of channels **increases**
(total “volume” is preserved!)

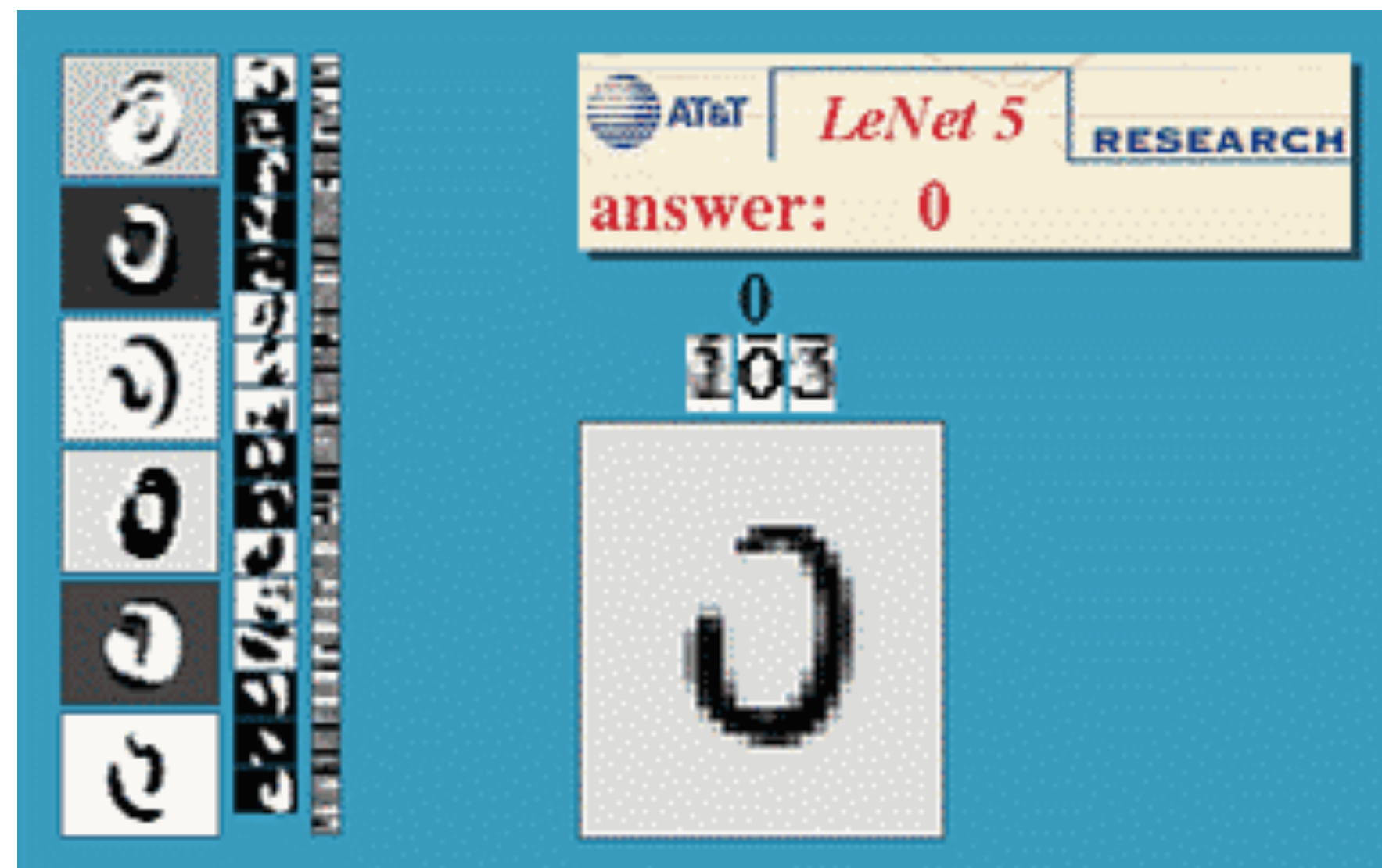
Optical Character Recognition (**OCR**)

Technology to convert **scanned documents to text**

(comes with any scanner now days)



Yann
LeCun

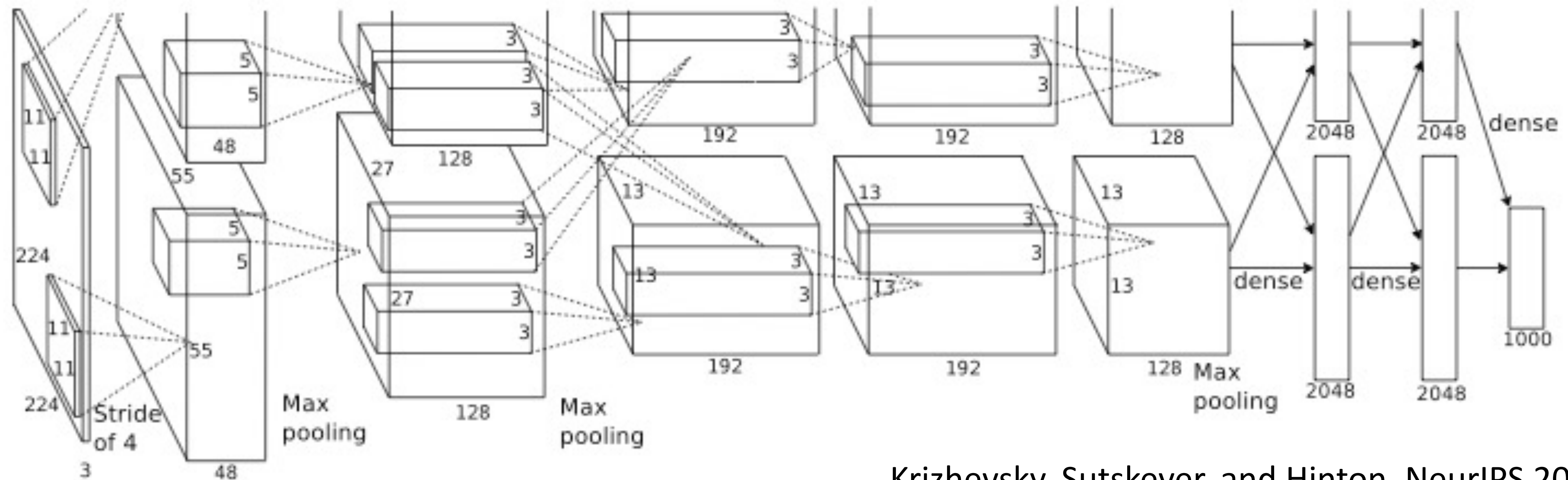


Digit recognition, AT&T labs
<http://www.research.att.com/~yann/>

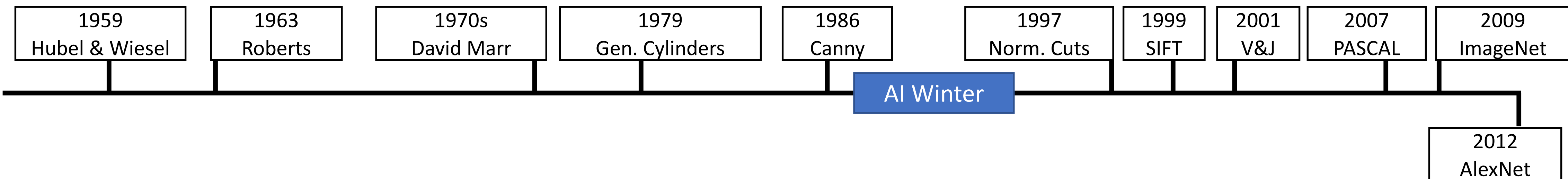


License plate readers
http://en.wikipedia.org/wiki/Automatic_number_plate_recognition

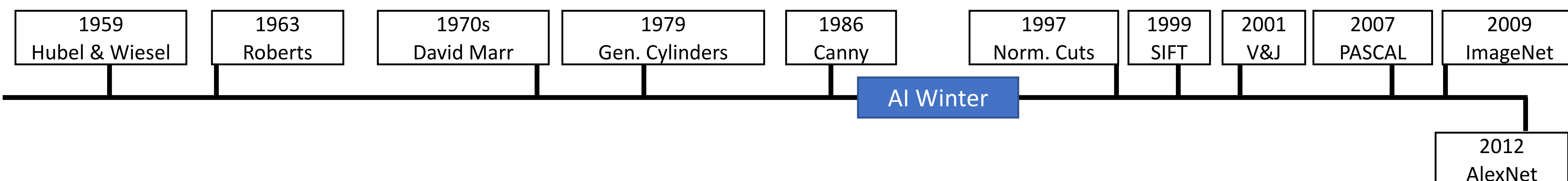
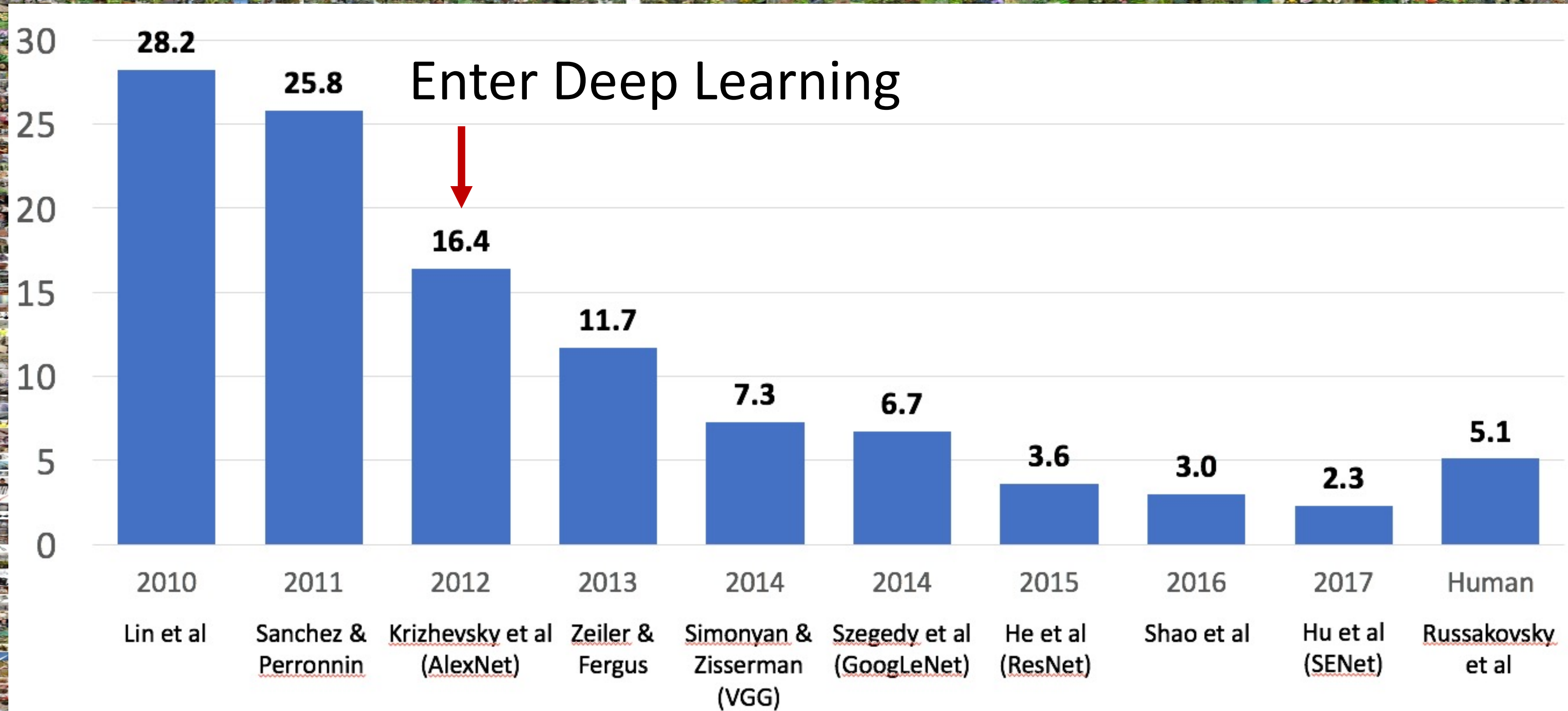
AlexNet: Deep Learning Goes Mainstream



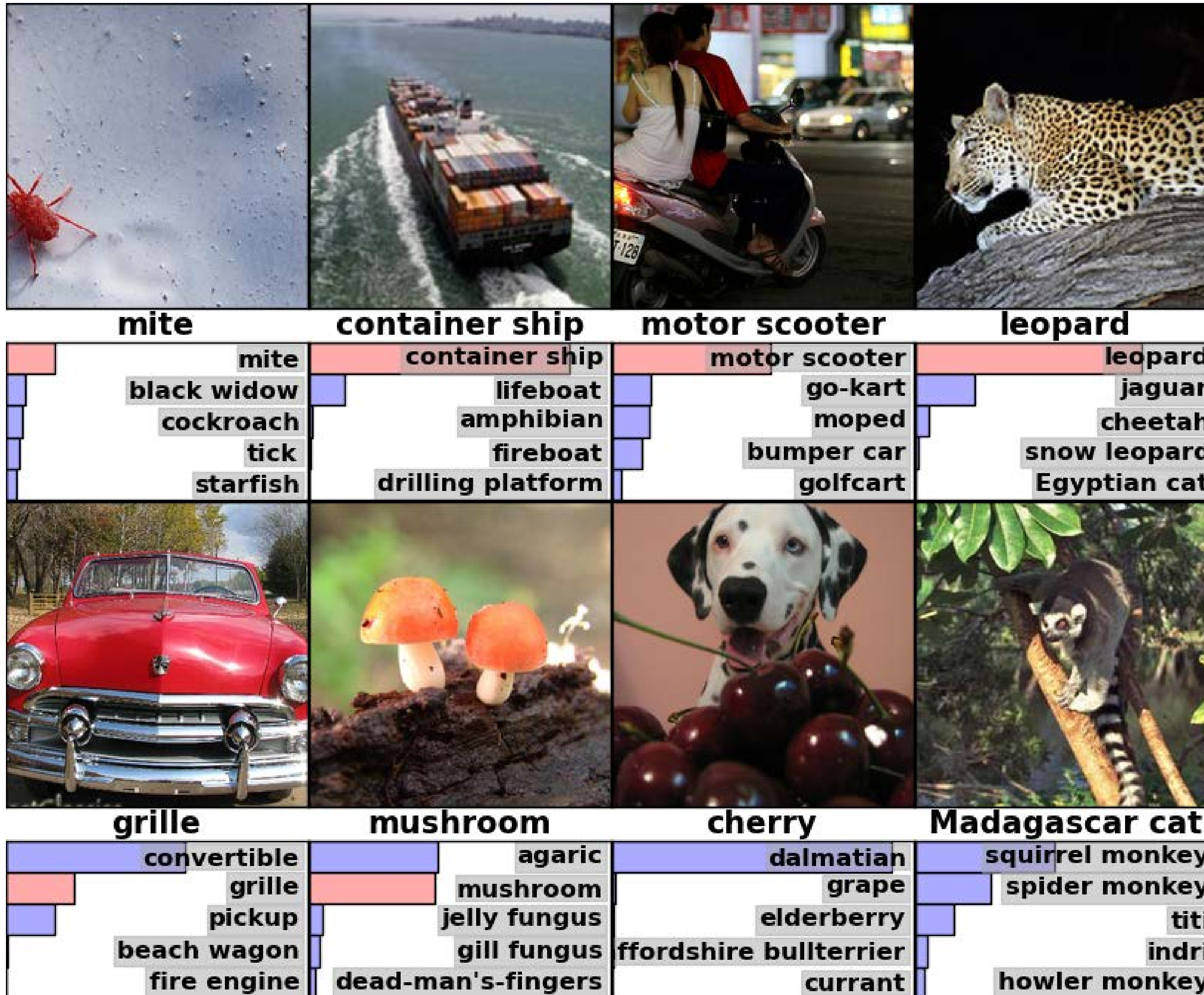
Krizhevsky, Sutskever, and Hinton, NeurIPS 2012



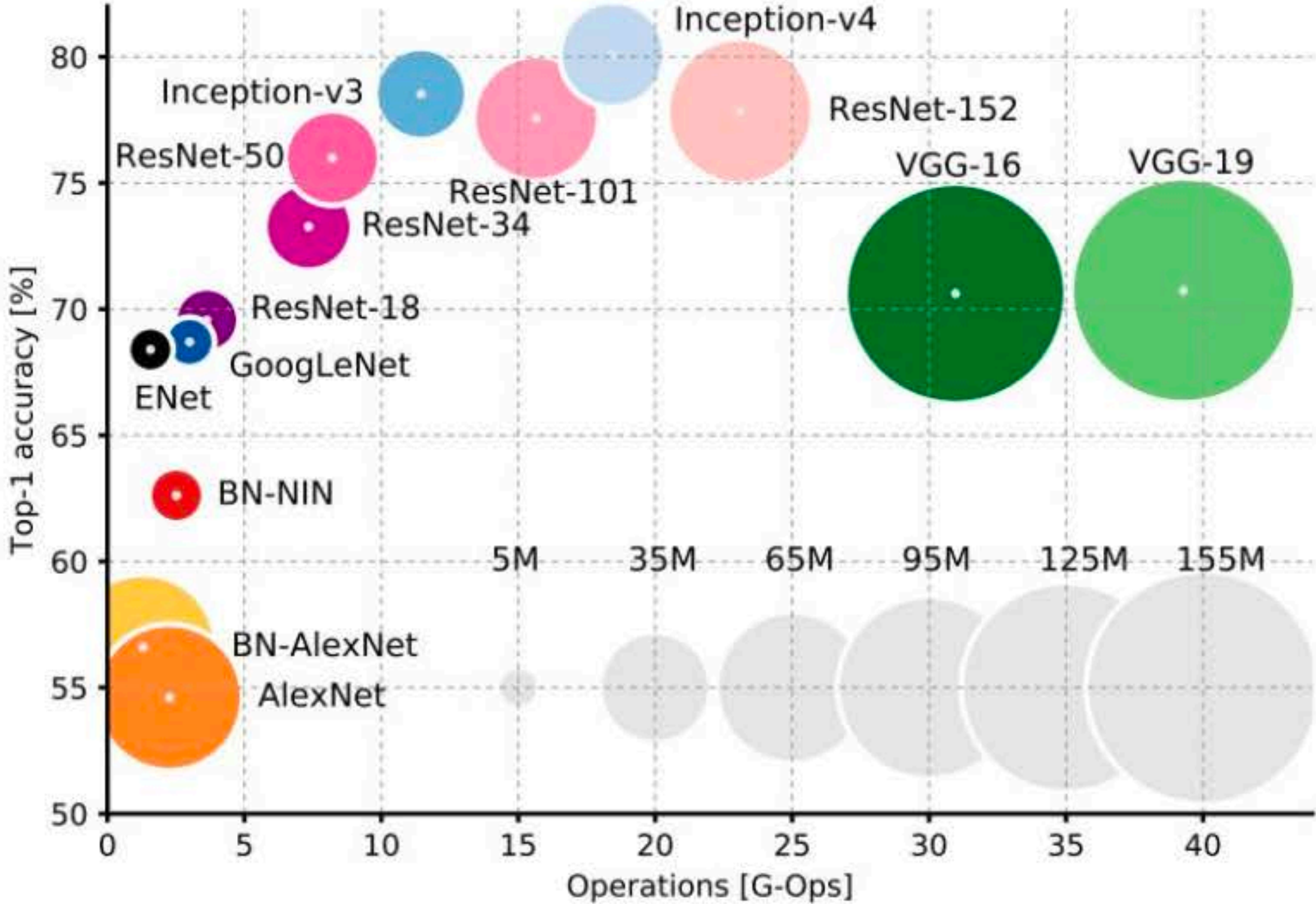
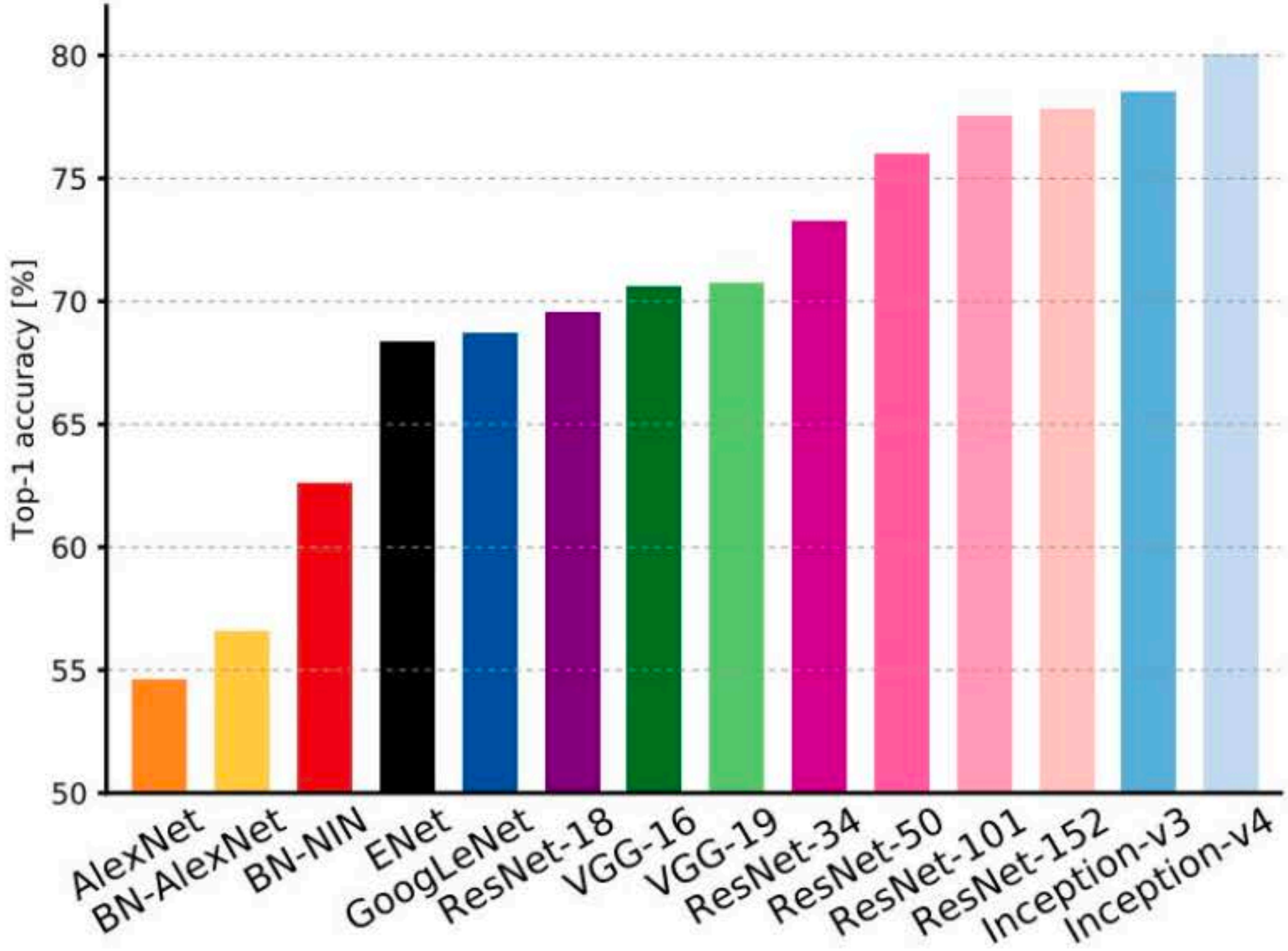
IMAGENET Large Scale Visual Recognition Challenge



AlexNet on ImageNet



Comparing Complexity



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

* adopted from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

Summary

The parameters of a neural network are learned using **backpropagation**, which computes gradients via recursive application of the chain rule

A **convolutional neural network** assumes inputs are images, and constrains the network architecture to reduce the number of parameters

A **convolutional layer** applies a set of learnable filters

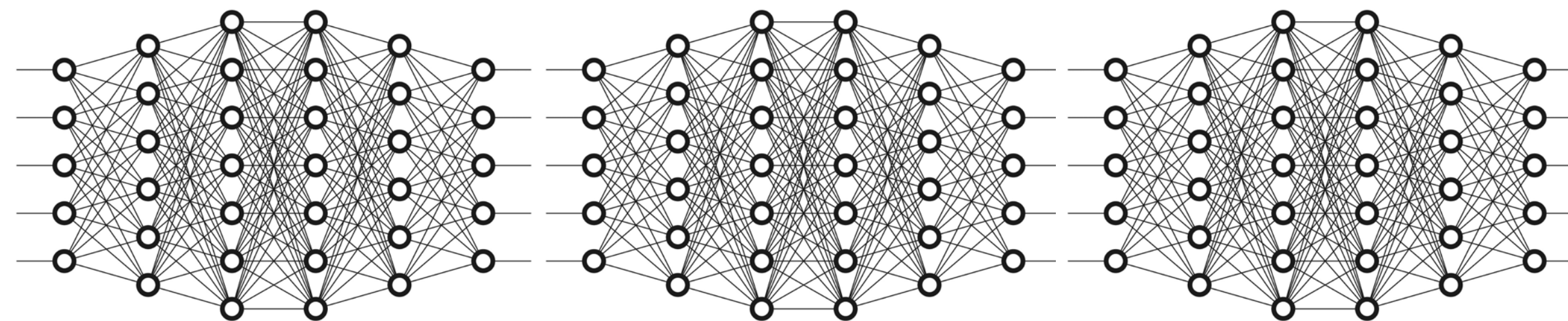
A **pooling layer** performs spatial downsampling

A **fully-connected** layer is the same as in a regular neural network

Convolutional neural networks can be seen as learning a hierarchy of filters



CPSC 425: Computer Vision



Lecture 22: Neural Networks 3

Menu for Today

Topics:

- **Neural Networks** part 3
- Weight **Initialization**
- **Normalization**
- Preventing **Overfitting**

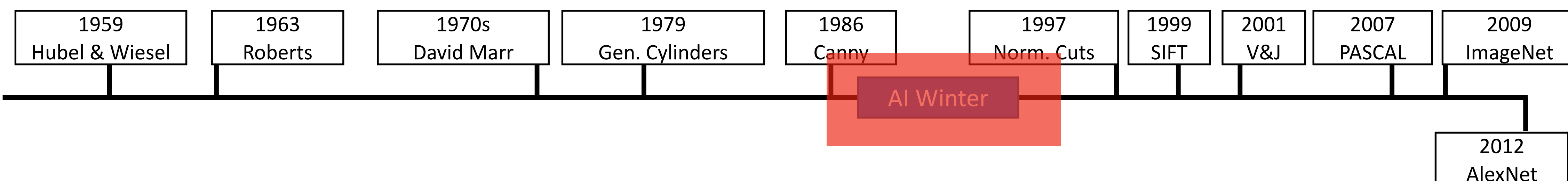
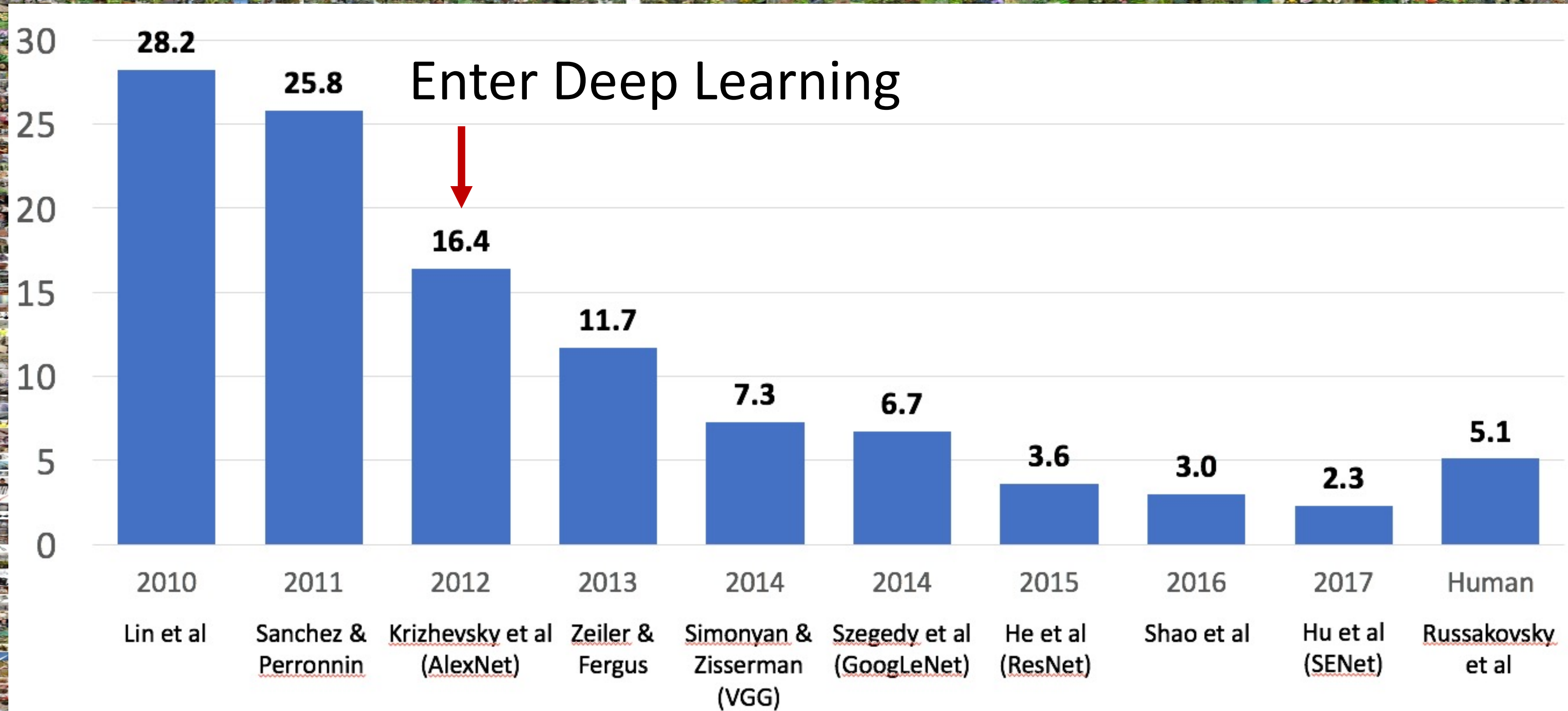
Readings:

- **Today's** Lecture: Szeliski 5.1.3, 5.3-5.4, Justin Johnson Michigan EECS 498/598

Reminders:

- **Quiz 6**: Open Apr 10th, due **Apr 11th!**
- **Assignment 6**: Deep Learning due **April 12th!**
- **Final**: April 16th

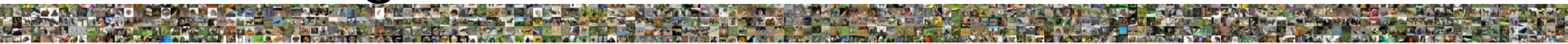
IMAGENET Large Scale Visual Recognition Challenge





So why now?

Rise of large datasets



IMAGENET

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

Rise of large datasets

LAION-5B: A NEW ERA OF OPEN LARGE-SCALE MULTI-MODAL DATASETS | LAION

https://laion.ai/blog/laion-5b/

LAION

- Projects
- Team
- Blog
- Notes
- Press
- About
- FAQ
- Donations
- Privacy Policy
- Dataset Requests
- Impressum

LAION-5B: A NEW ERA OF OPEN LARGE-SCALE MULTI-MODAL DATASETS

by: Romain Beaumont, 31 Mar, 2022

We present a dataset of 5,85 billion CLIP-filtered image-text pairs, 14x bigger than LAION-400M, previously the biggest openly accessible image-text dataset in the world - see also our [NeurIPS2022 paper](#)

Authors: Christoph Schuhmann, Richard Vencu, Romain Beaumont, Theo Coombes, Cade Gordon, Aarush Katta, Robert Kaczmarczyk, Jenia Jitsev

Backend url: Index:

Clip retrieval works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions Display full captions Display similarities

french cat

french cat

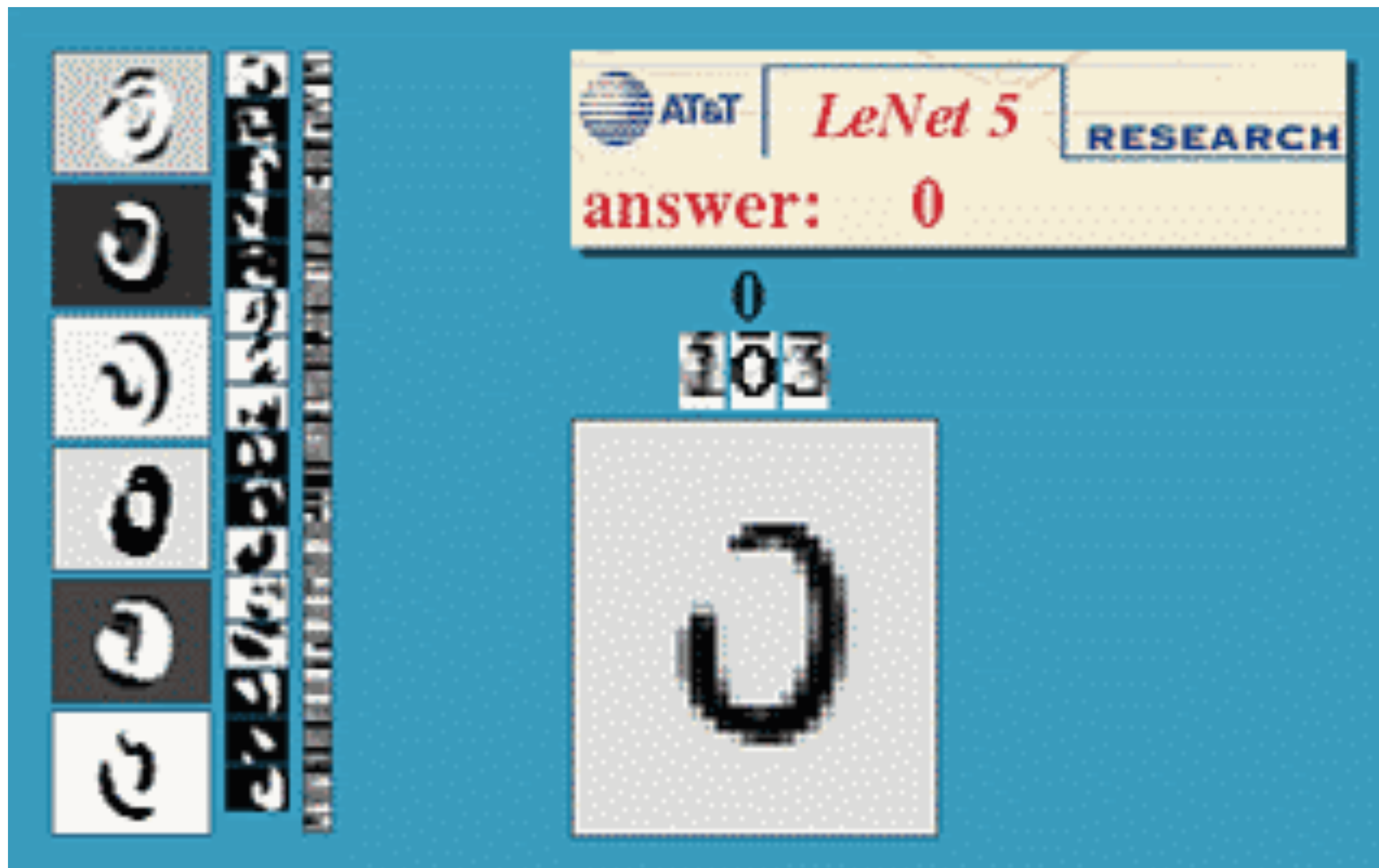
How to tell if your feline is french. He wears a b...

イケメン猫モデル「トキ・ナンタケツト」がかっこいい

Hilarious pics of funny cats! funnycatsgif.com

Clever architectures

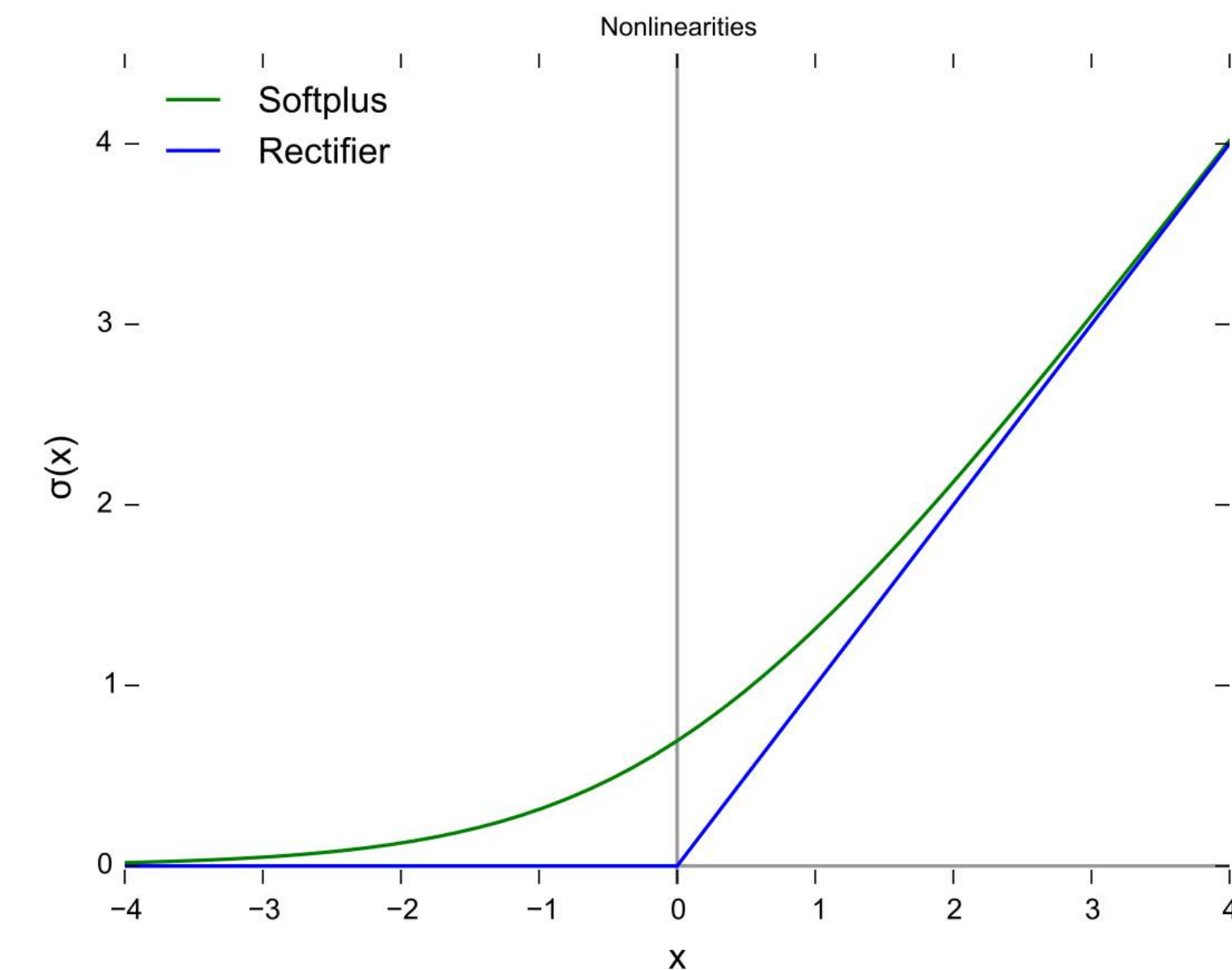
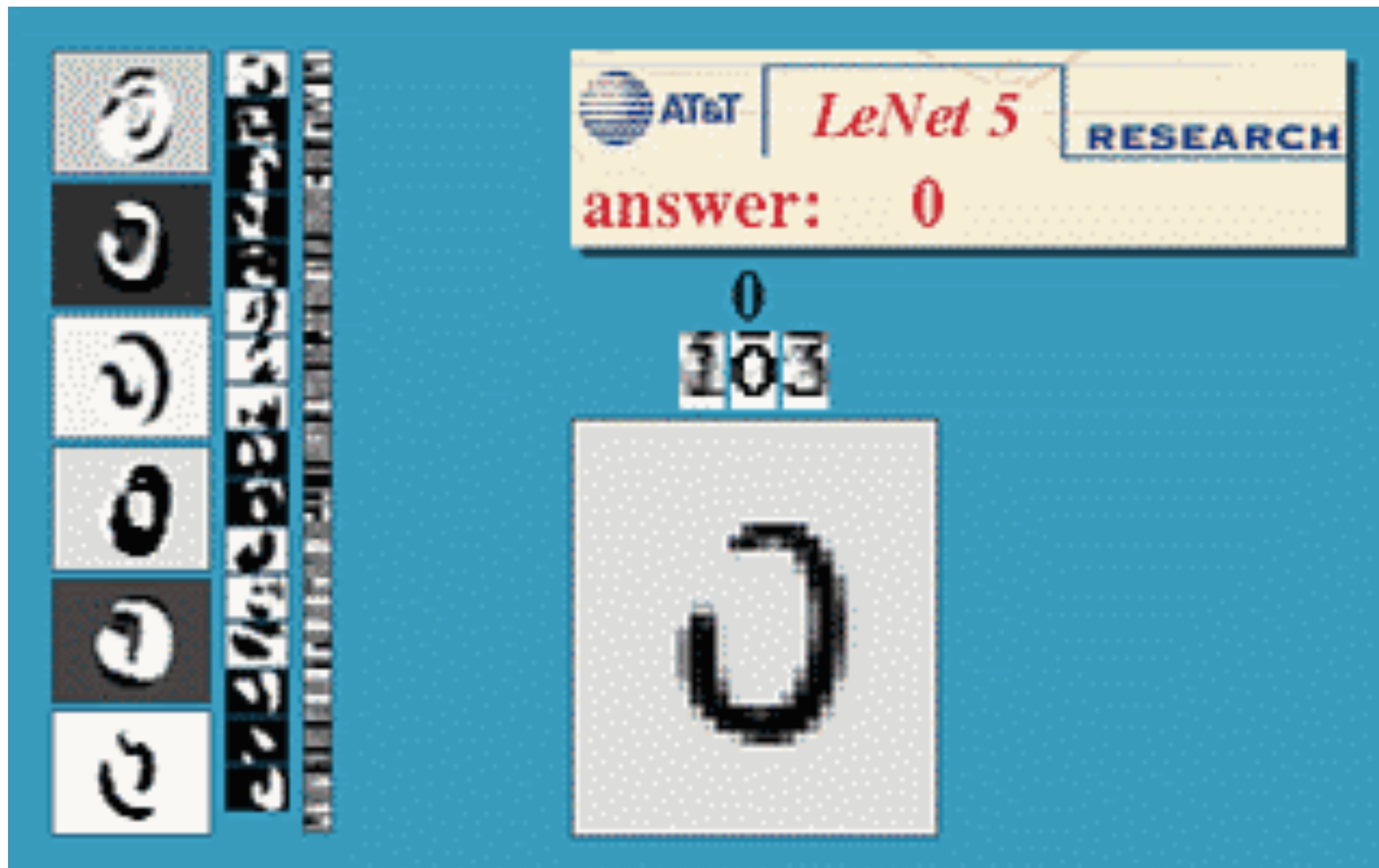
Convolutional neural networks



[Lecun, Bottou, Bengio, and Haffner, "Gradient-Based Learning Applied to Document Recognition", 1998]

Clever architectures

Convolutional neural networks



[Nair and Hinton, "Rectified Linear Units Improve Restricted Boltzmann Machines", 2010]

[Lecun, Bottou, Bengio, and Haffner, "Gradient-Based Learning Applied to Document Recognition", 1998]

Clever architectures

Transformers, Vaswani et al., 2017

Attention is all you need

Authors Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, Illia Polosukhin

Publication date 2017

Journal Advances in neural information processing systems

Volume 30

Description The dominant sequence transduction models are based on complex recurrent or convolutional neural networks in an encoder and decoder configuration. The best performing such models also connect the encoder and decoder through an attention mechanism. We propose a novel, simple network architecture based solely on an attention mechanism, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our single model with 165 million parameters, achieves 27.5 BLEU on English-to-German translation, improving over the existing best ensemble result by over 1 BLEU. On English-to-French translation, we outperform the previous single state-of-the-art with model by 0.7 BLEU, achieving a BLEU score of 41.1.

Total citations Cited by 87925

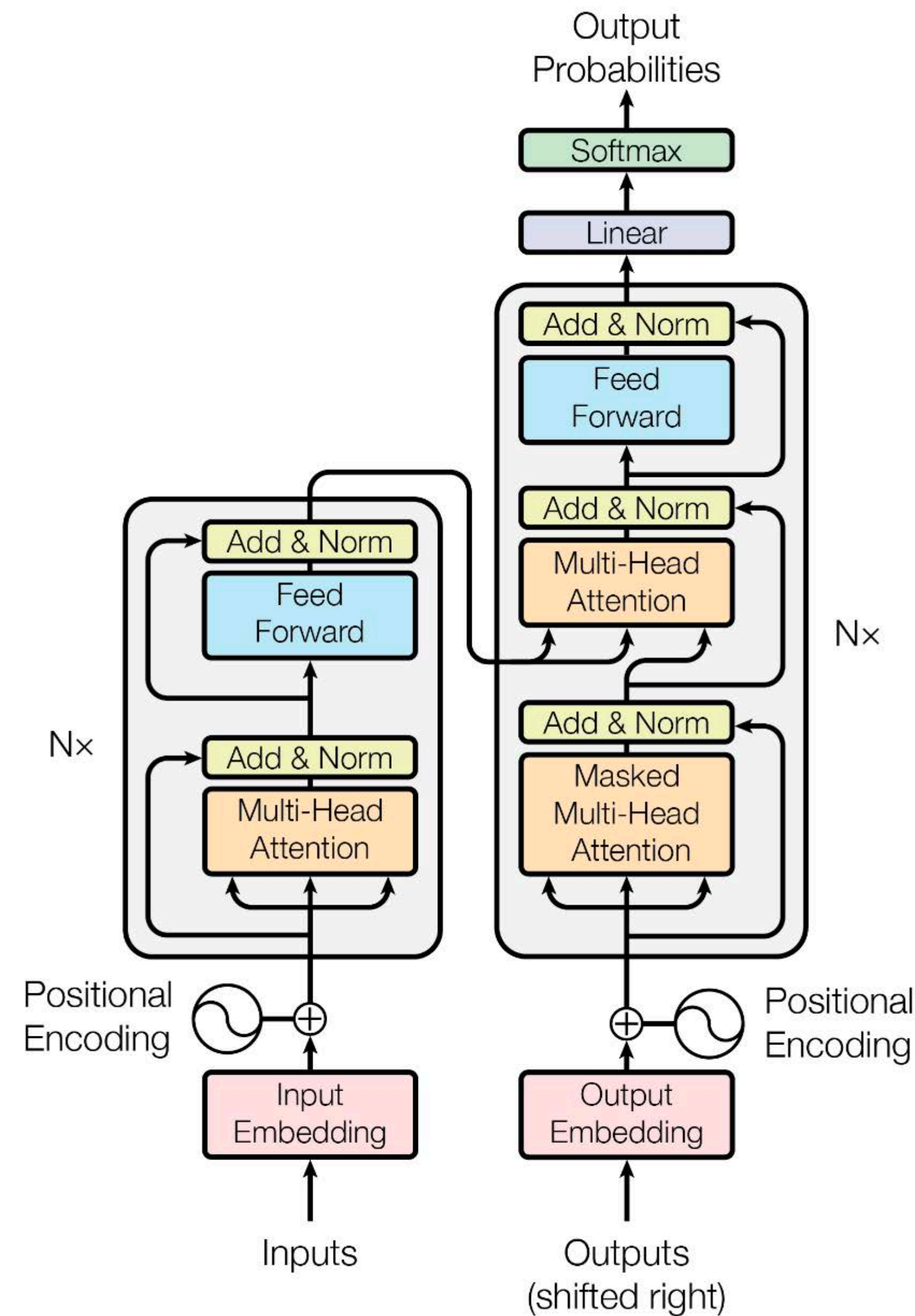
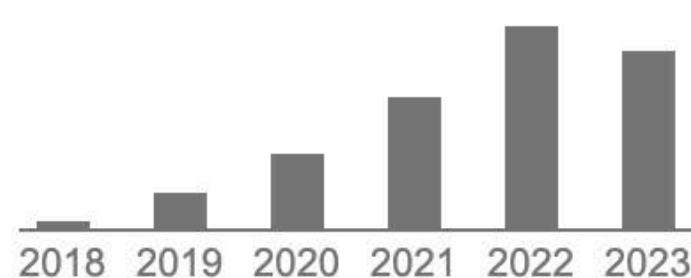
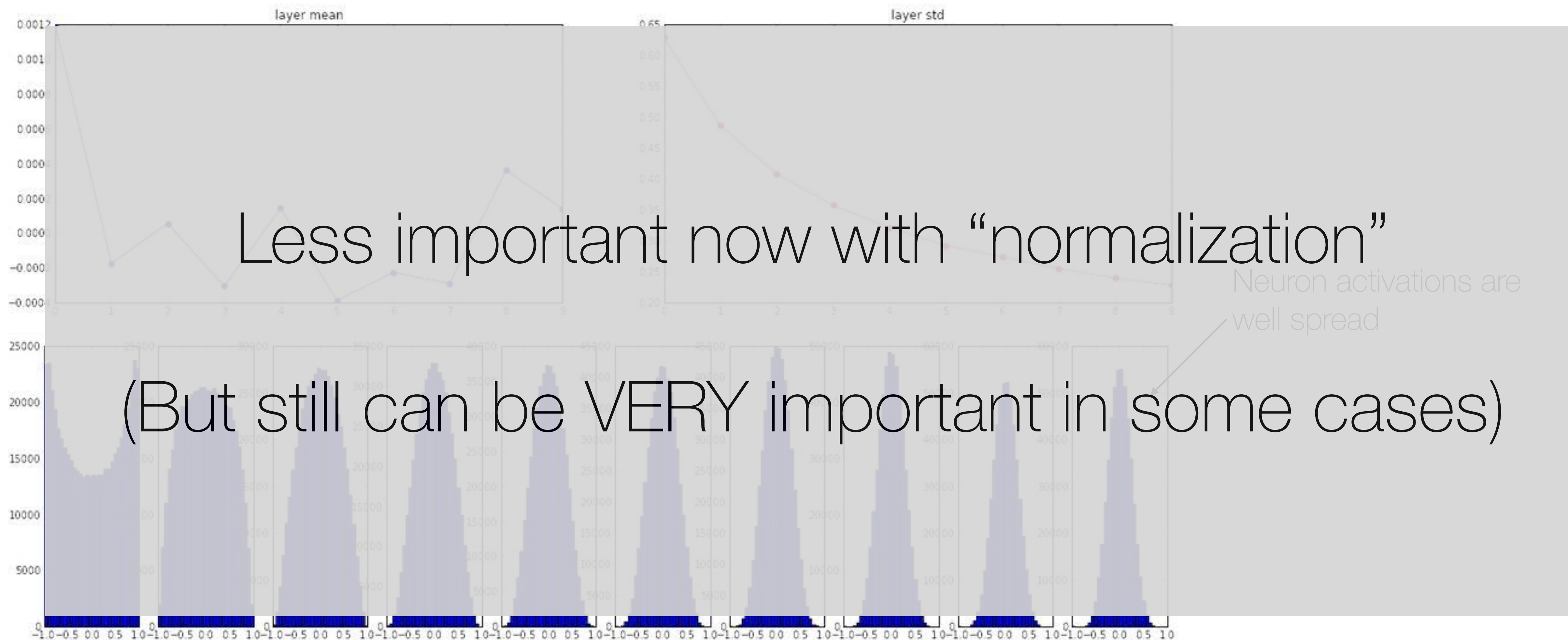


Figure 1: The Transformer - model architecture.

Proper initialization schemes

Much more important than you think

“Xavier initialization”
[Glorot et al., 2010]



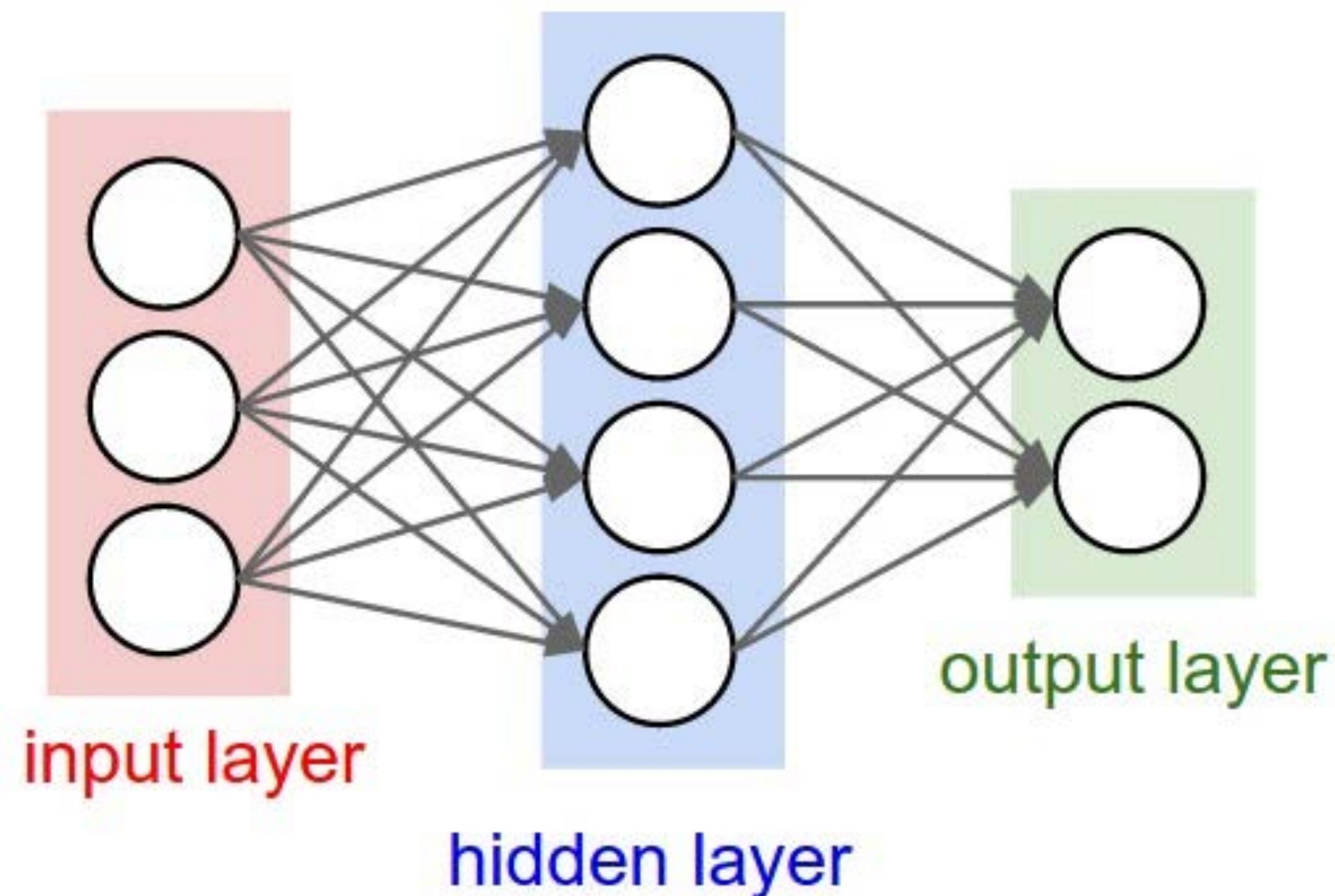


Weight initialization

Weight initialization

A look into ways things can go wrong

Q: what happens when $W=0$ init is used?



We
A loc

```
27 parser.add_argument("--img_size", type=int, default=32, help="size of each image dimension")
28 parser.add_argument("--channels", type=int, default=1, help="number of image channels")
29 parser.add_argument("--sample_interval", type=int, default=1000, help="number of image channels")
30 opt = parser.parse_args()
31 print(opt)
32
33 cuda = True if torch.cuda.is_available() else False
34
35
36 def weights_init_normal(m):
37     classname = m.__class__.__name__
38     if classname.find("Conv") != -1:
39         torch.nn.init.normal_(m.weight.data, 0.0, 0.02)
40     elif classname.find("BatchNorm") != -1:
41         torch.nn.init.normal_(m.weight.data, 1.0, 0.02)
42         torch.nn.init.constant_(m.bias.data, 0.0)
43
44
45 class Generator(nn.Module):
46     def __init__(self):
47         super(Generator, self).__init__()
48
49         self.init_size = opt.img_size // 4
50         self.ll = nn.Sequential(nn.Linear(opt.latent_dim, 128 * self.init_size * 2))
```


Weight initialization

A look into ways things can go wrong

Let's look at
some
activation
statistics

E.g. 10-layer net with
500 neurons on each
layer, using tanh non-
linearities, and initializing
as described in last slide.

```
# assume some unit gaussian 10-D input data
D = np.random.randn(1000, 500)
hidden_layer_sizes = [500]*10
nonlinearities = ['tanh']*len(hidden_layer_sizes)

act = {'relu':lambda x:np.maximum(0,x), 'tanh':lambda x:np.tanh(x)}
Hs = {}
for i in xrange(len(hidden_layer_sizes)):
    X = D if i == 0 else Hs[i-1] # input at this layer
    fan_in = X.shape[1]
    fan_out = hidden_layer_sizes[i]
    W = np.random.randn(fan_in, fan_out) * 0.01 # layer initialization

    H = np.dot(X, W) # matrix multiply
    H = act[nonlinearities[i]](H) # nonlinearity
    Hs[i] = H # cache result on this layer

# look at distributions at each layer
print 'input layer had mean %f and std %f' % (np.mean(D), np.std(D))
layer_means = [np.mean(H) for i,H in Hs.iteritems()]
layer_stds = [np.std(H) for i,H in Hs.iteritems()]
for i,H in Hs.iteritems():
    print 'hidden layer %d had mean %f and std %f' % (i+1, layer_means[i], layer_stds[i])

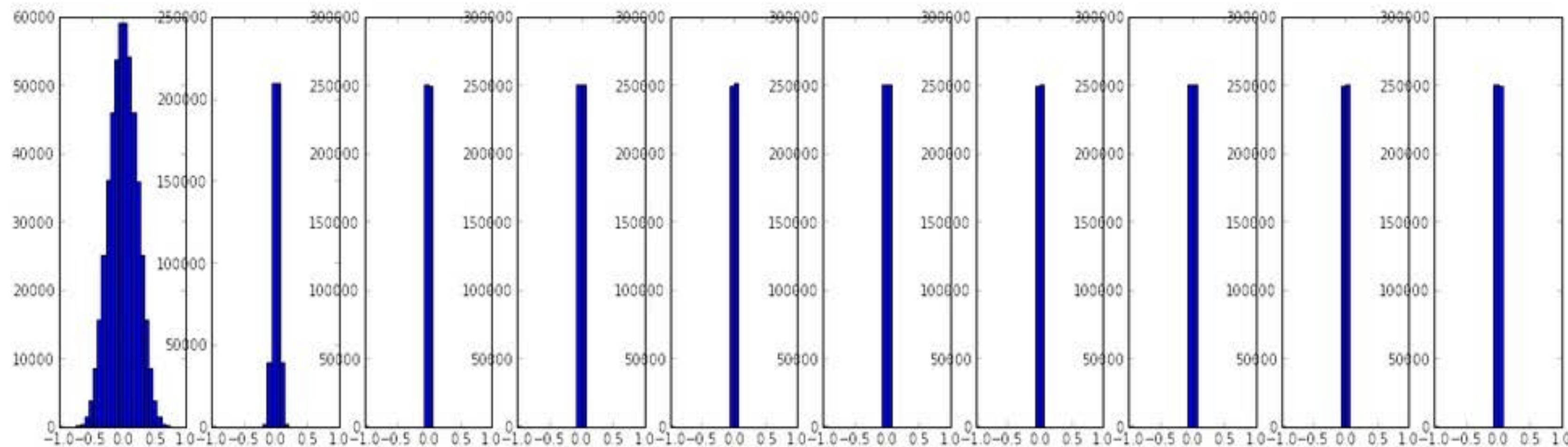
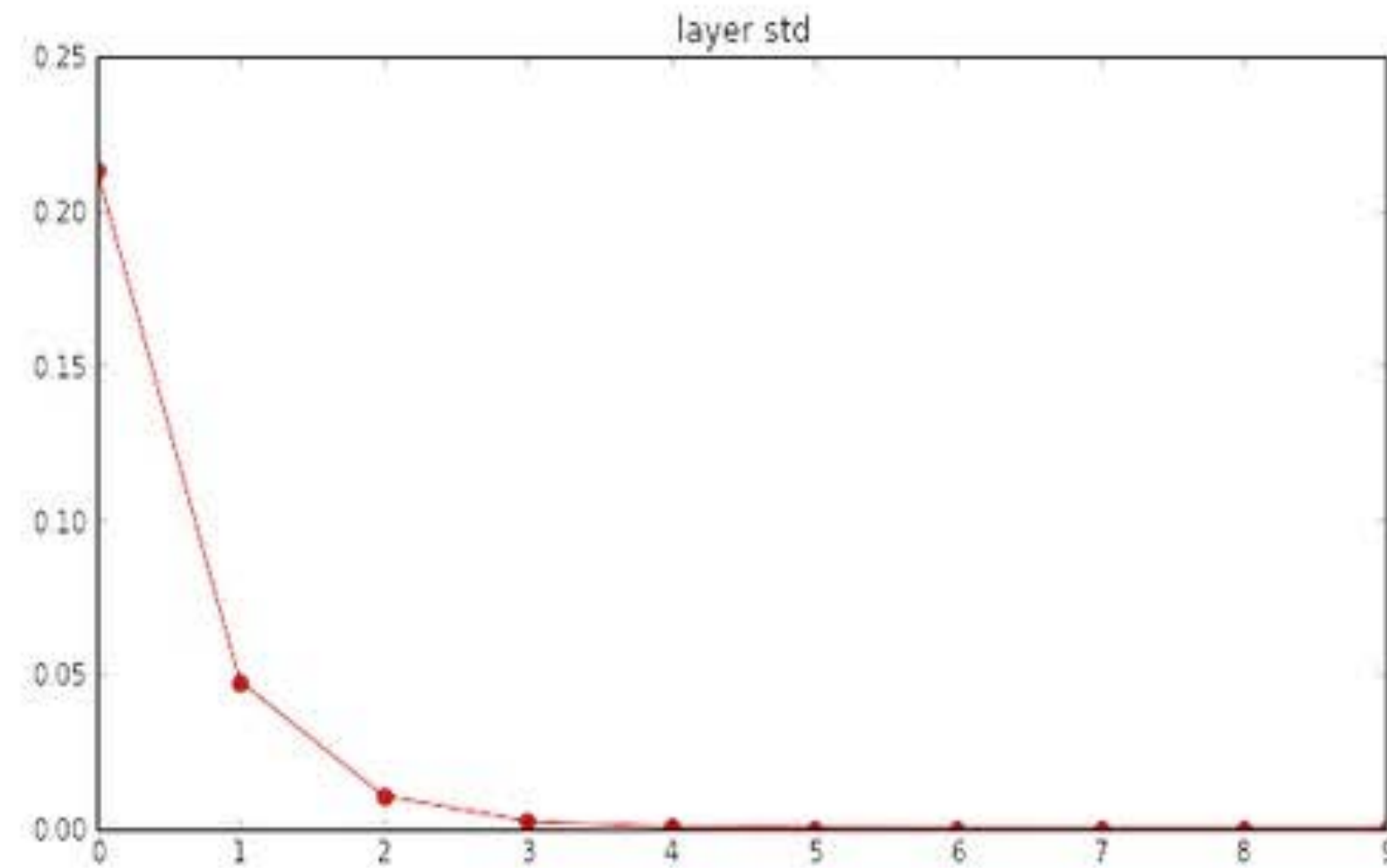
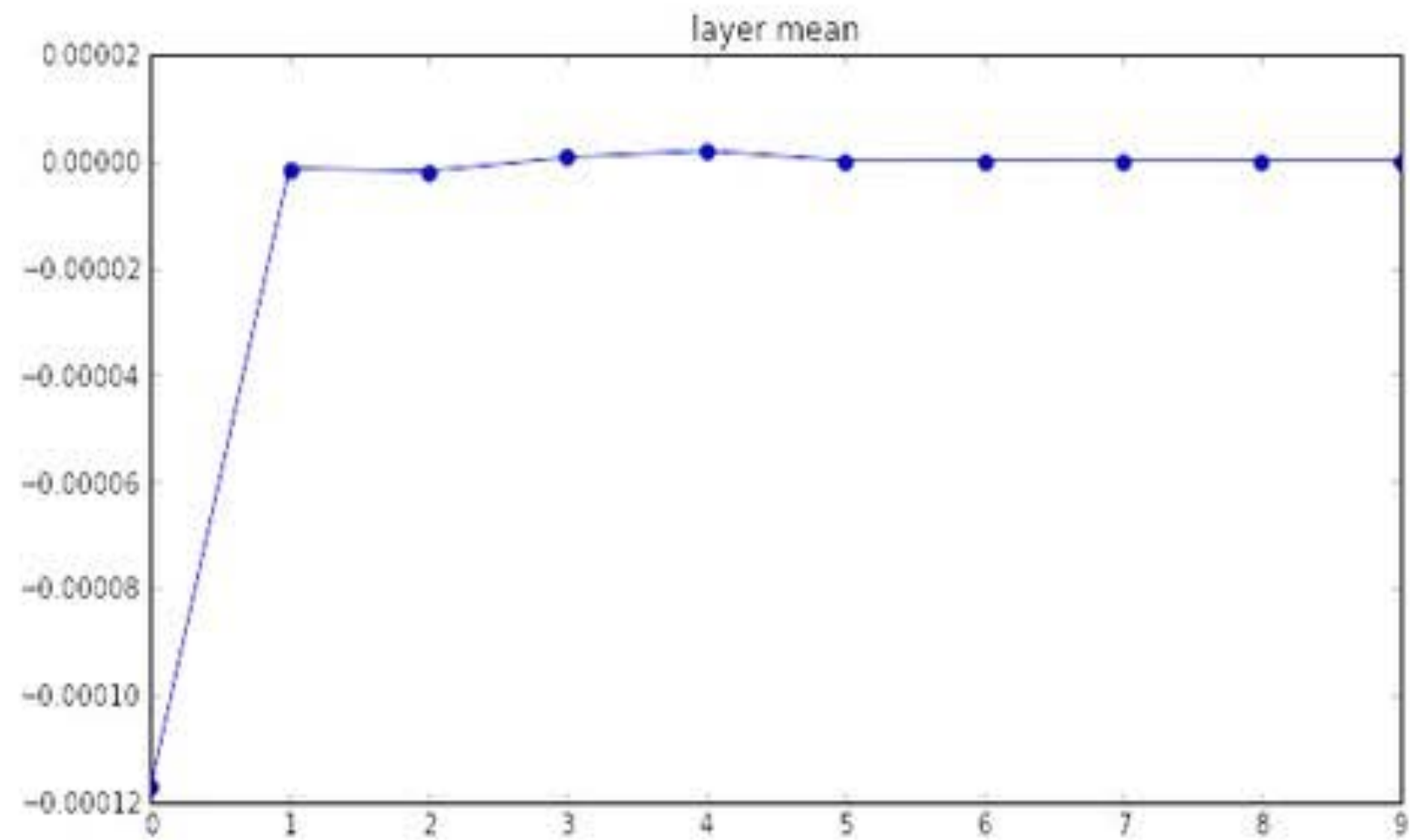
# plot the means and standard deviations
plt.figure()
plt.subplot(121)
plt.plot(Hs.keys(), layer_means, 'ob-')
plt.title('layer mean')
plt.subplot(122)
plt.plot(Hs.keys(), layer_stds, 'or-')
plt.title('layer std')

# plot the raw distributions
plt.figure()
for i,H in Hs.iteritems():
    plt.subplot(1,len(Hs),i+1)
    plt.hist(H.ravel(), 30, range=(-1,1))
```

Init with small random numbers

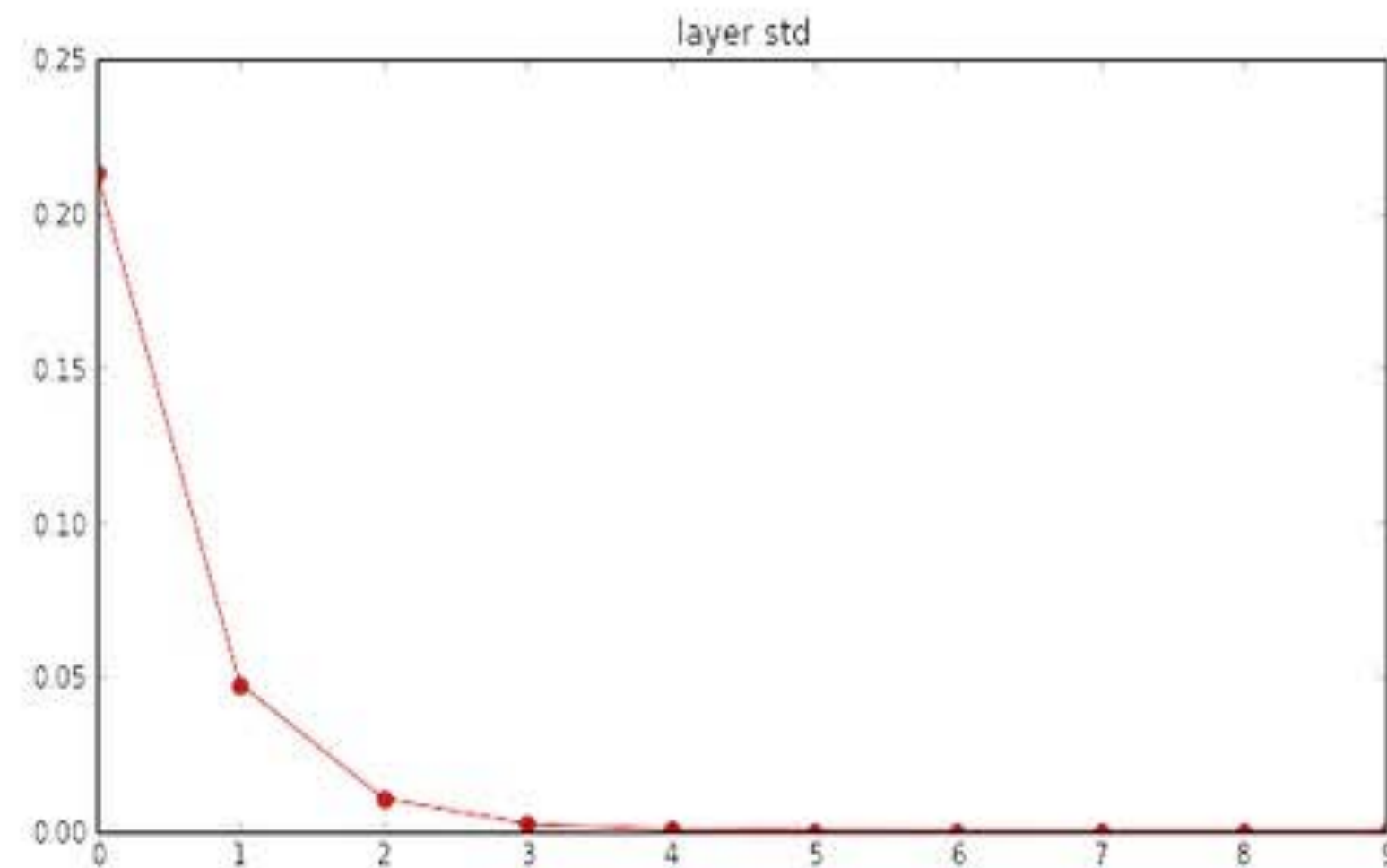
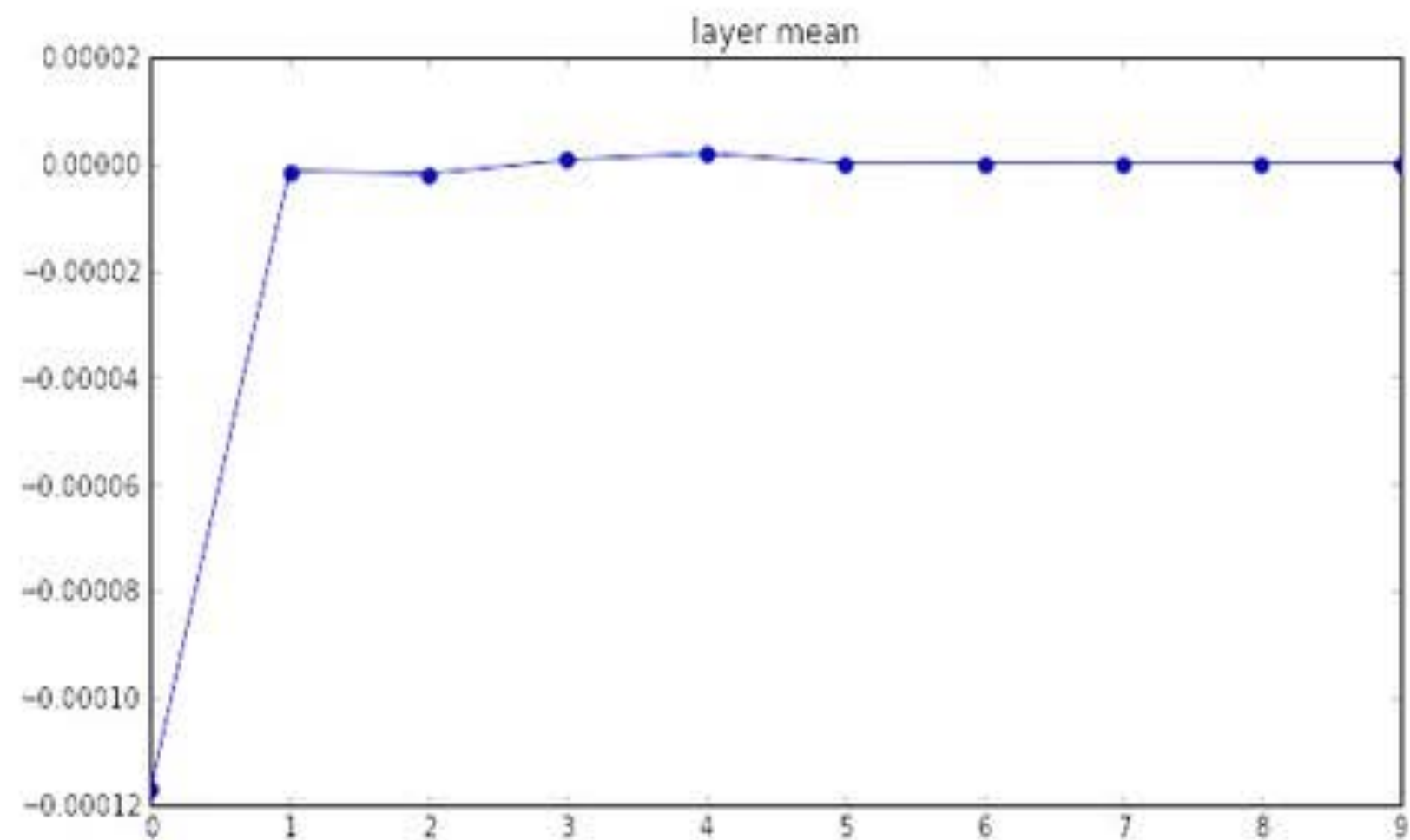
Weight initialization

A look into ways things can go wrong



Weight initialization

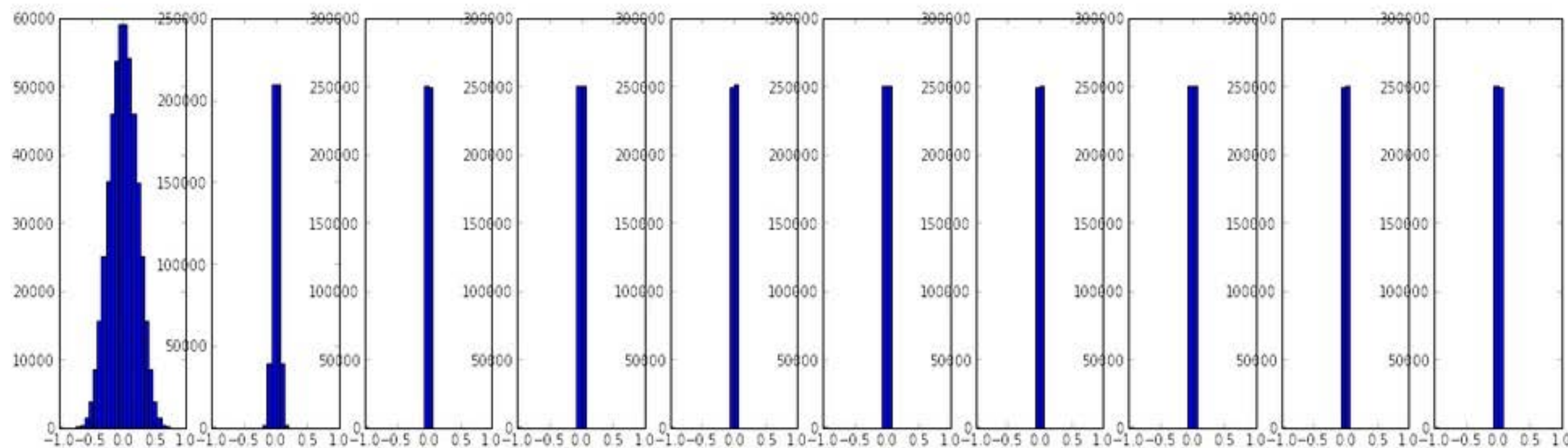
A look into ways things can go wrong



All activations become zero!

Q: think about the backward pass. What do the gradients look like?

Hint: think about backward pass for a $W \cdot X$ gate.

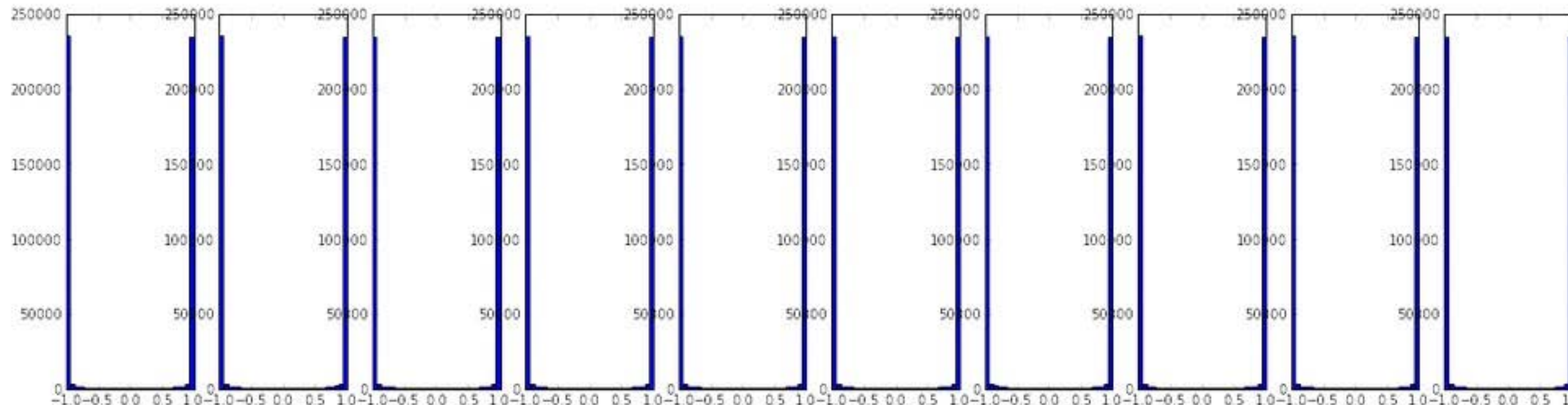
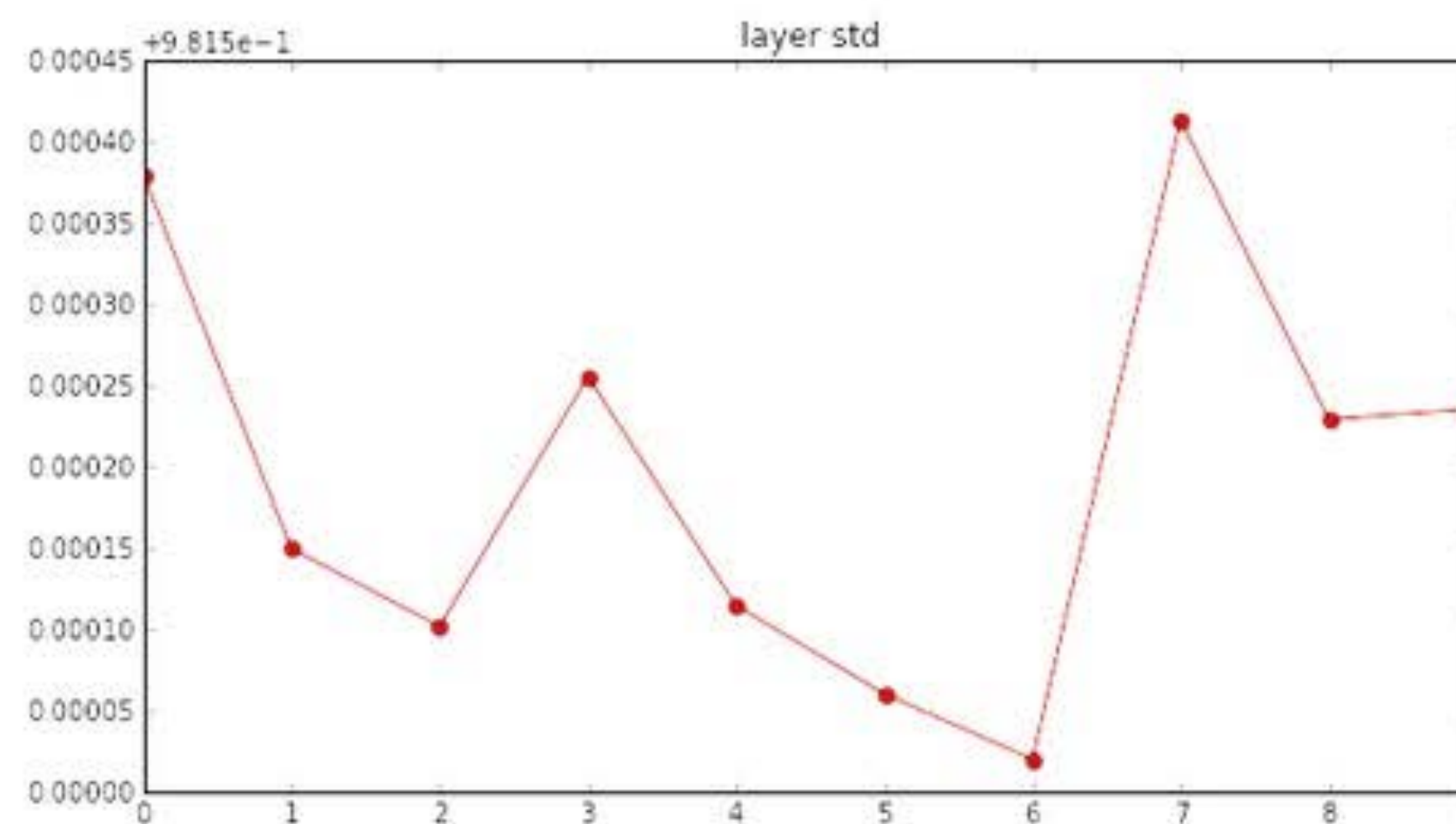
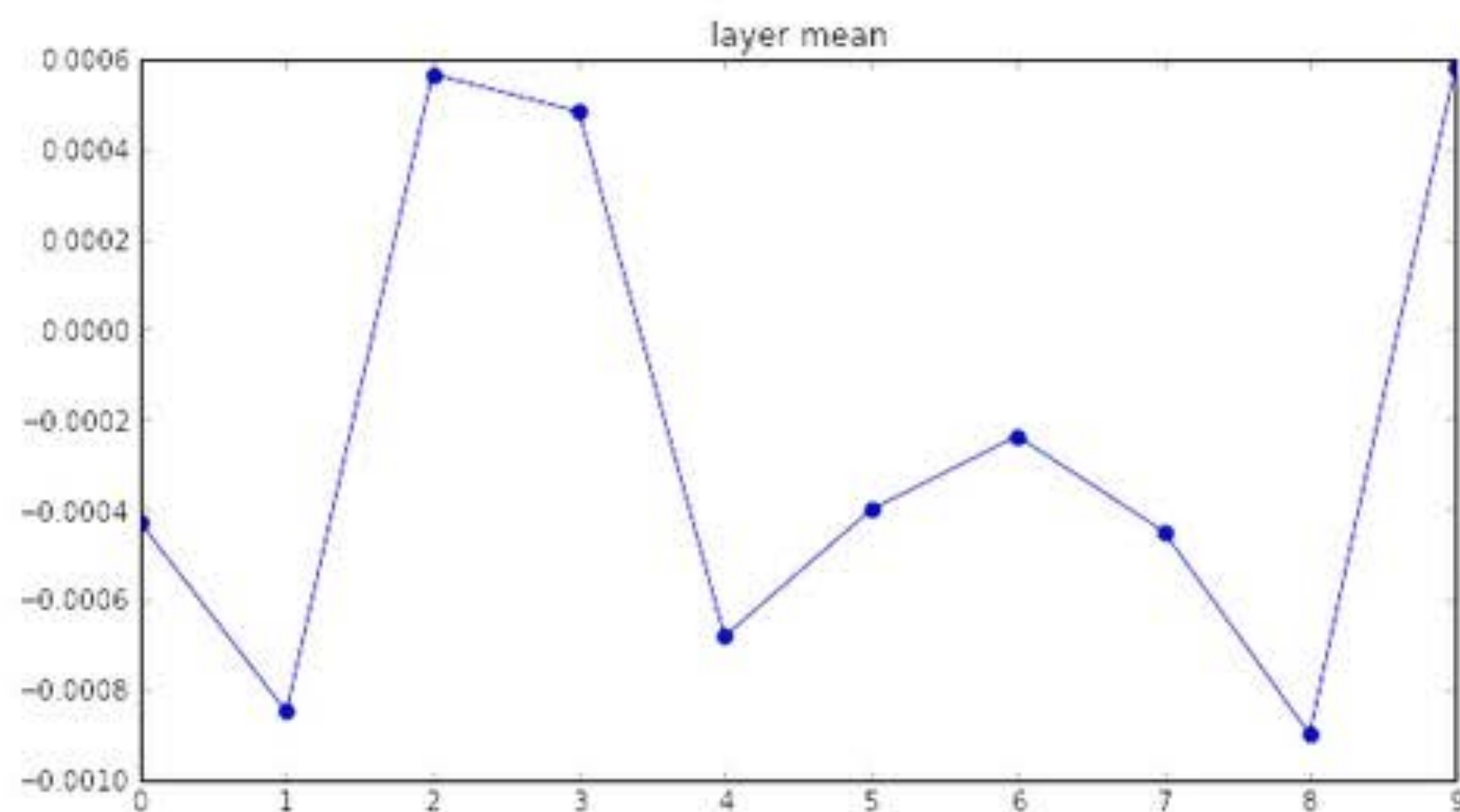


Weight initialization

A look into ways things can go wrong

```
W = np.random.randn(fan_in, fan_out) * 1.0 # layer initialization
```

*1.0 instead of *0.01



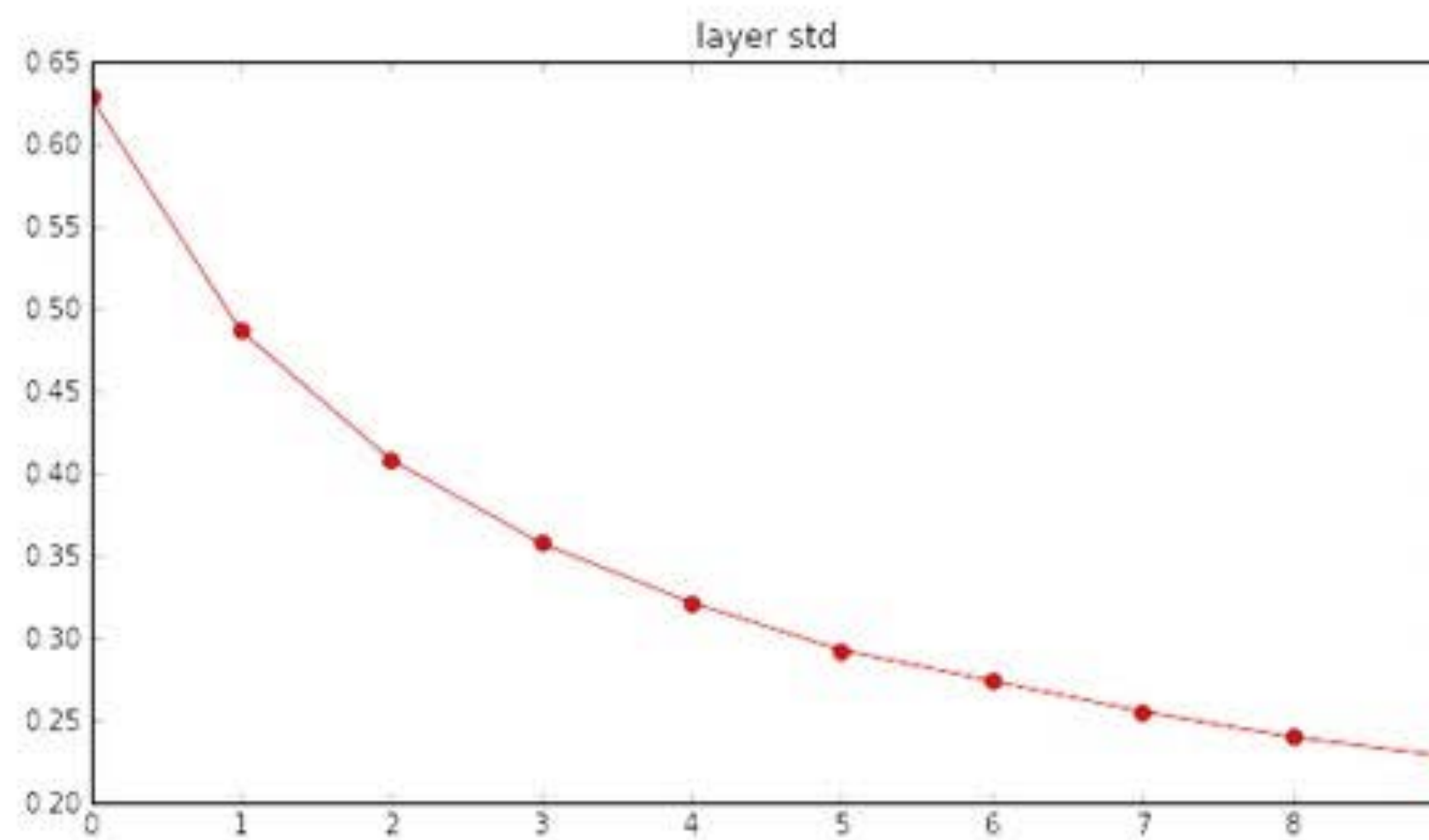
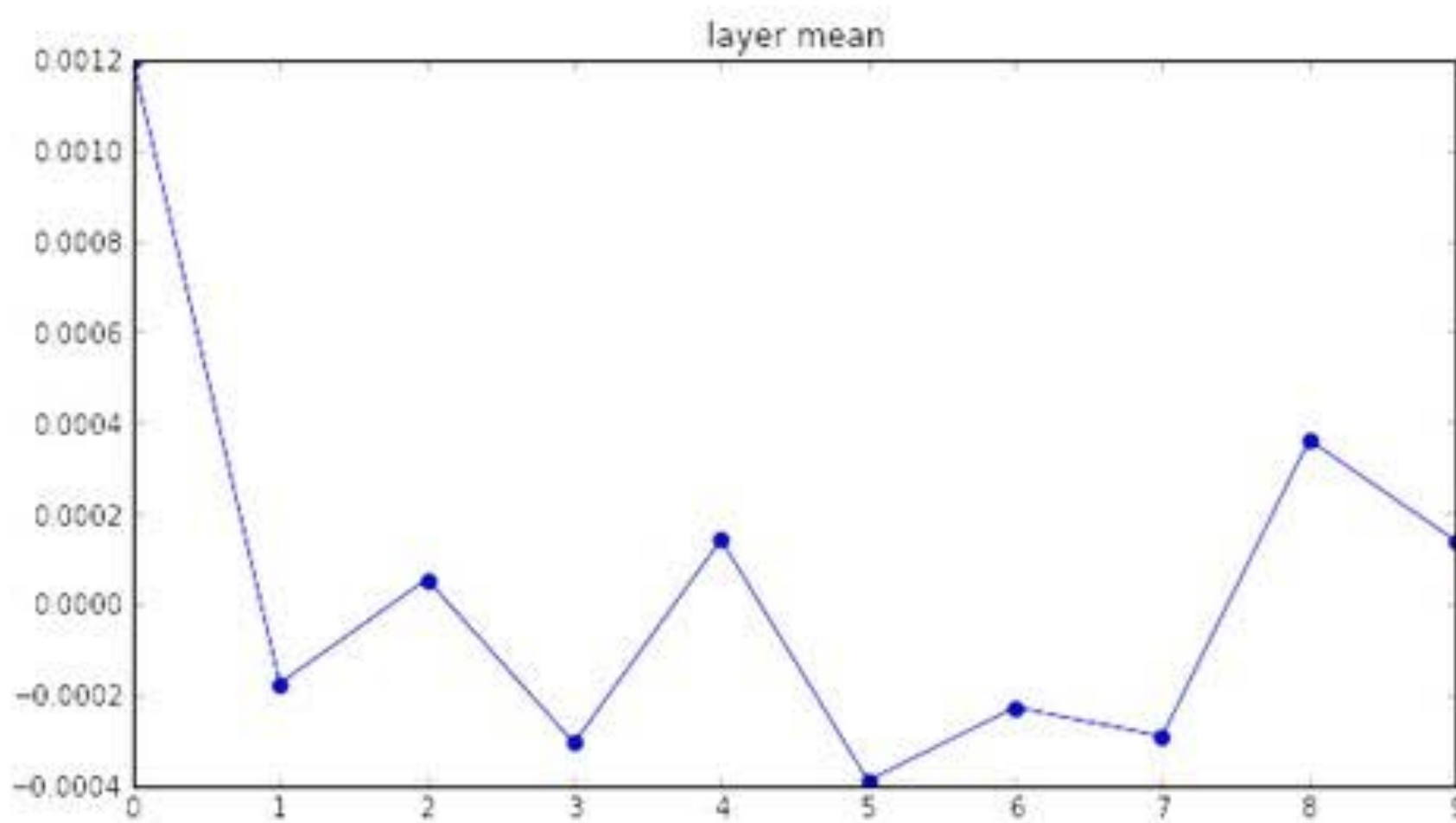
Almost all neurons completely saturated, either -1 and 1. Gradients will be all zero.

Weight initialization

```
W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in) # layer initialization
```

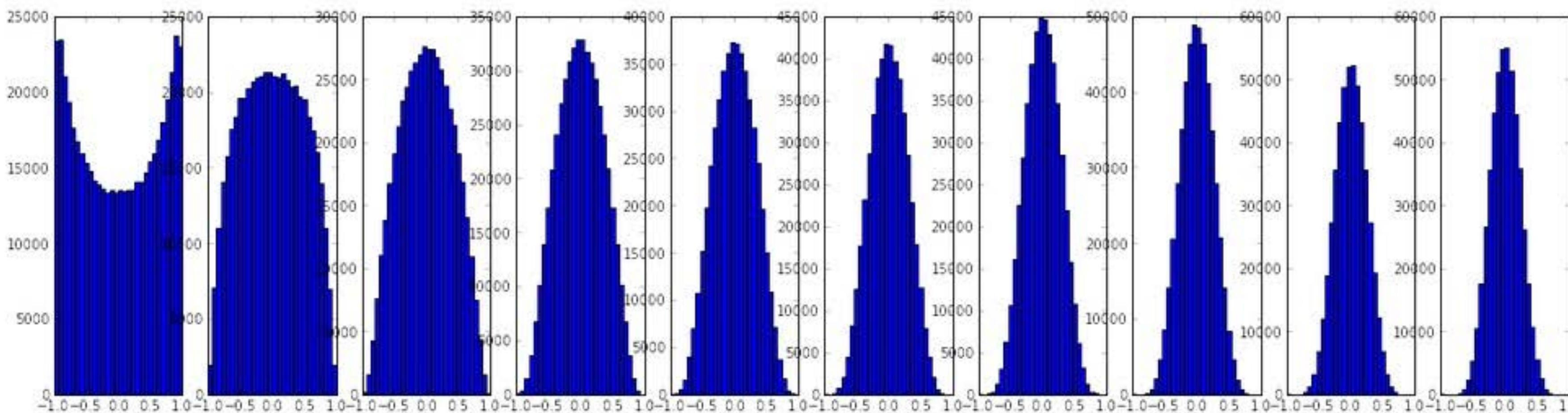
Glorot initialization [Glorot et al., 2010]

Some number according to "fan in"



Statistically motivated

Good for tanh

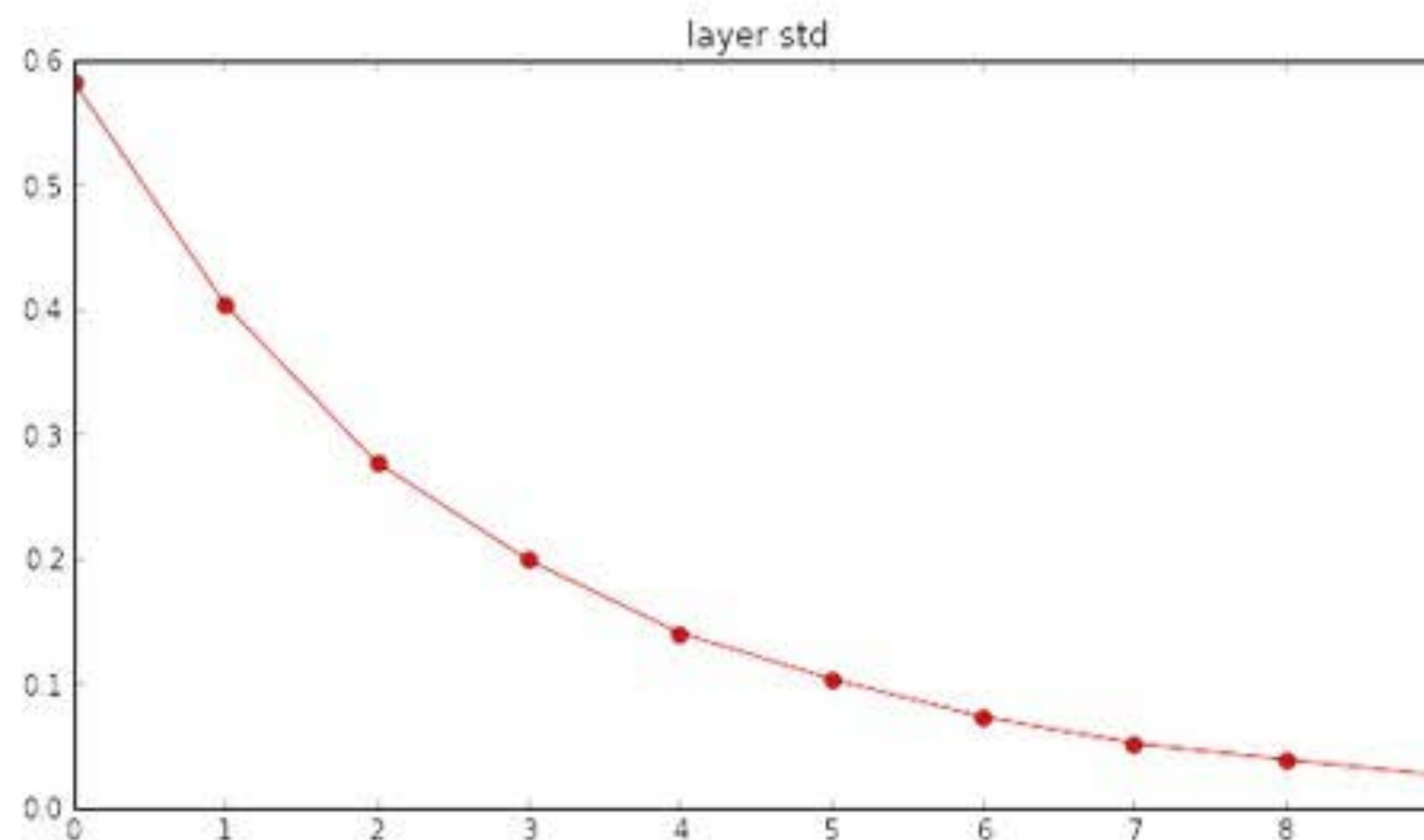
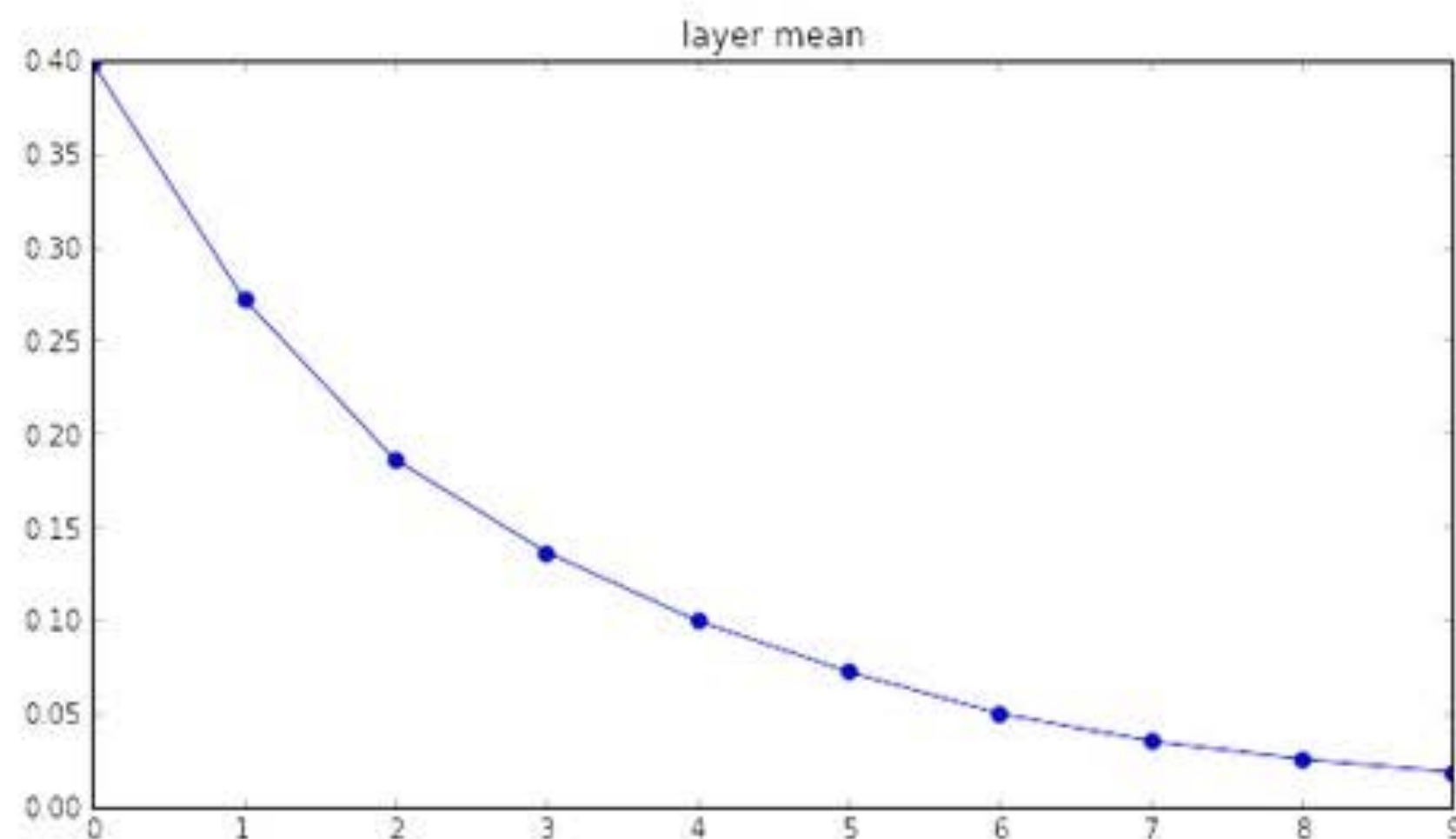


Weight initialization

```
W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in) # layer initialization
```

Glorot initialization [Glorot et al., 2010]

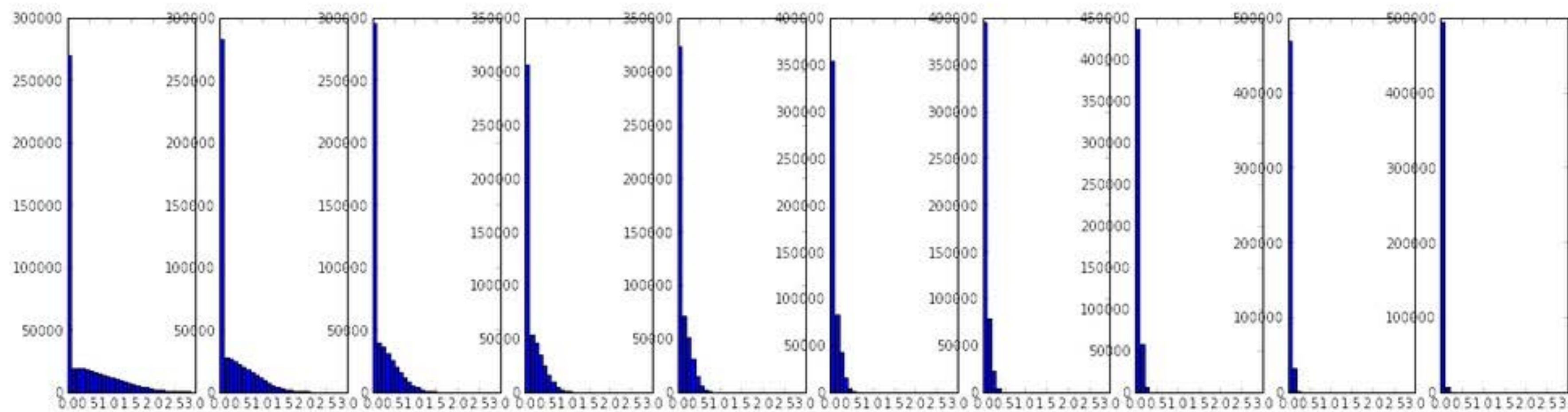
Some number according to "fan in"



Statistically motivated

Good for tanh

Not so good for ReLU

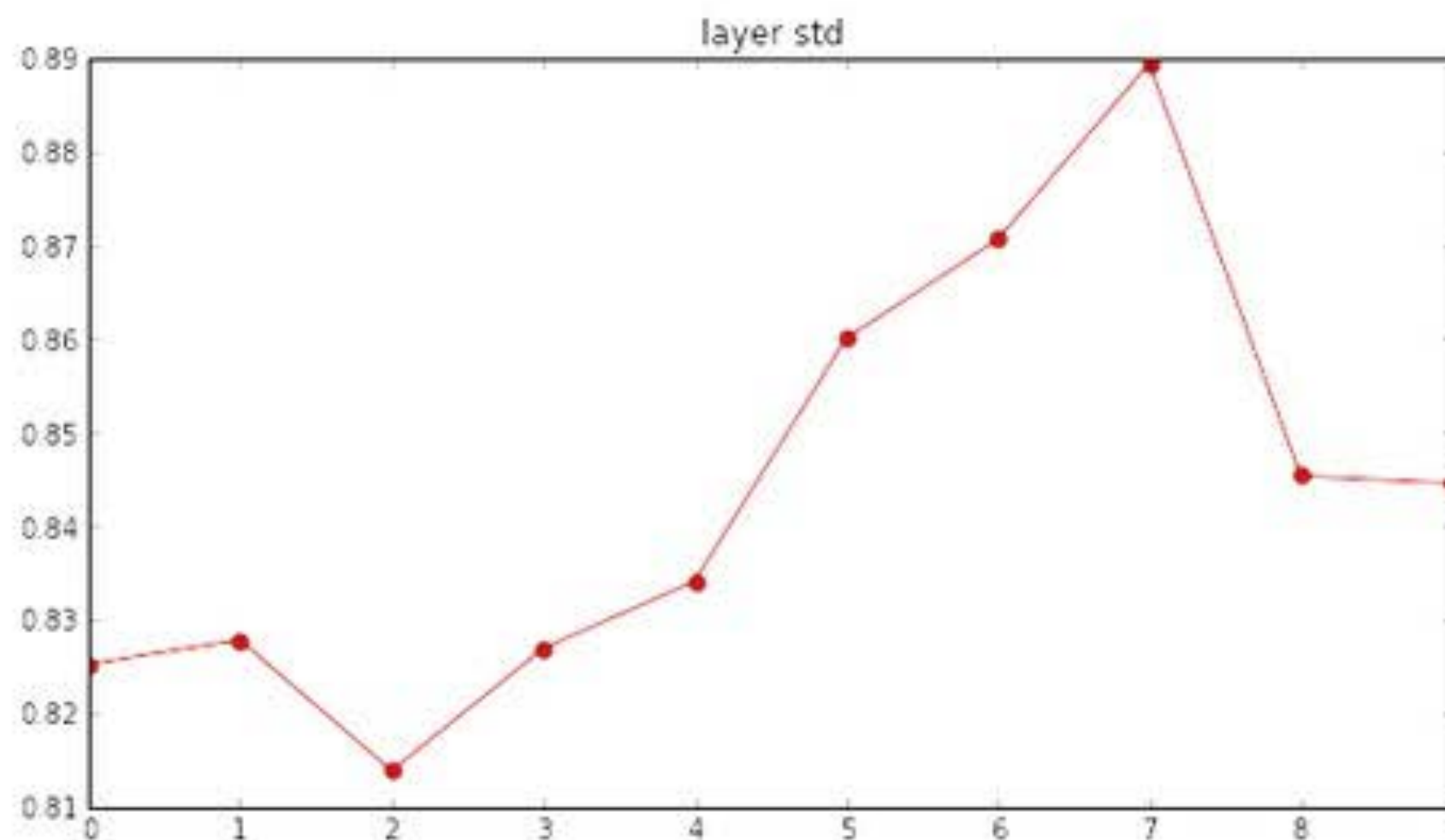
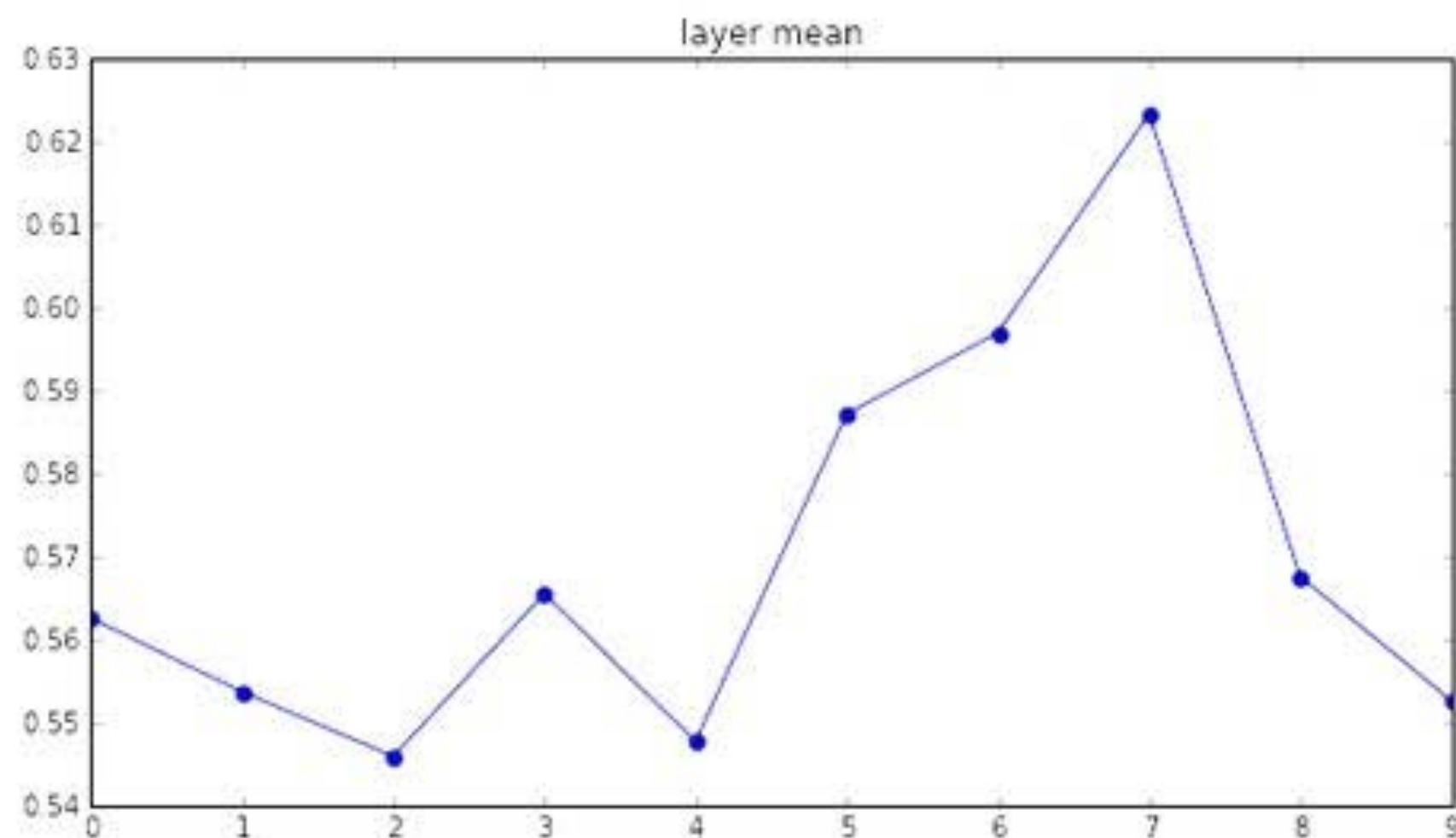


Weight initialization

He initialization [He et al., 2015]

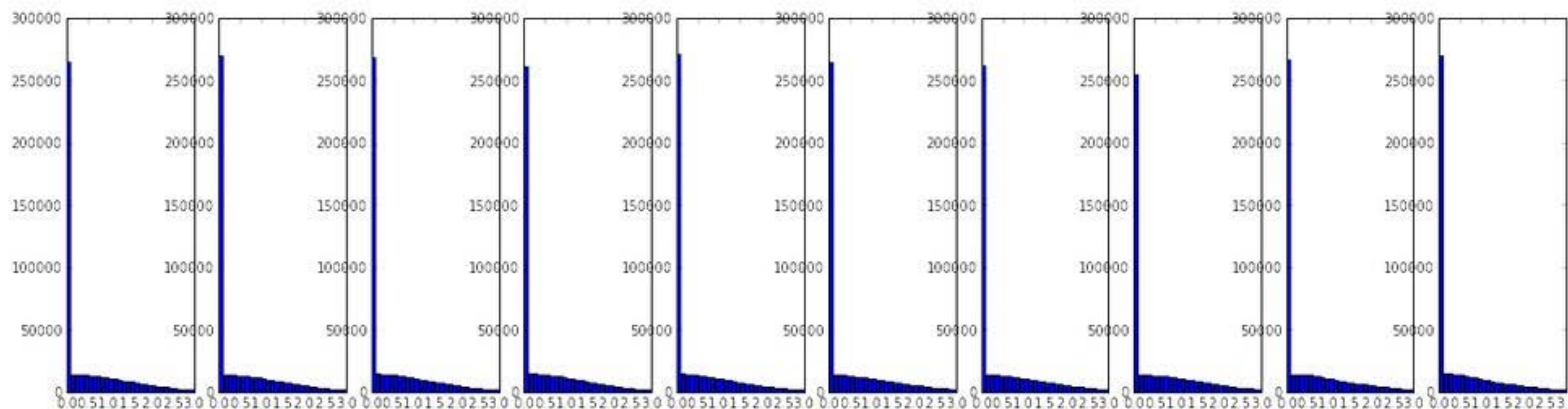
```
W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in/2) # layer initialization
```

magic number 2



Statistically motivated

Good for ReLU

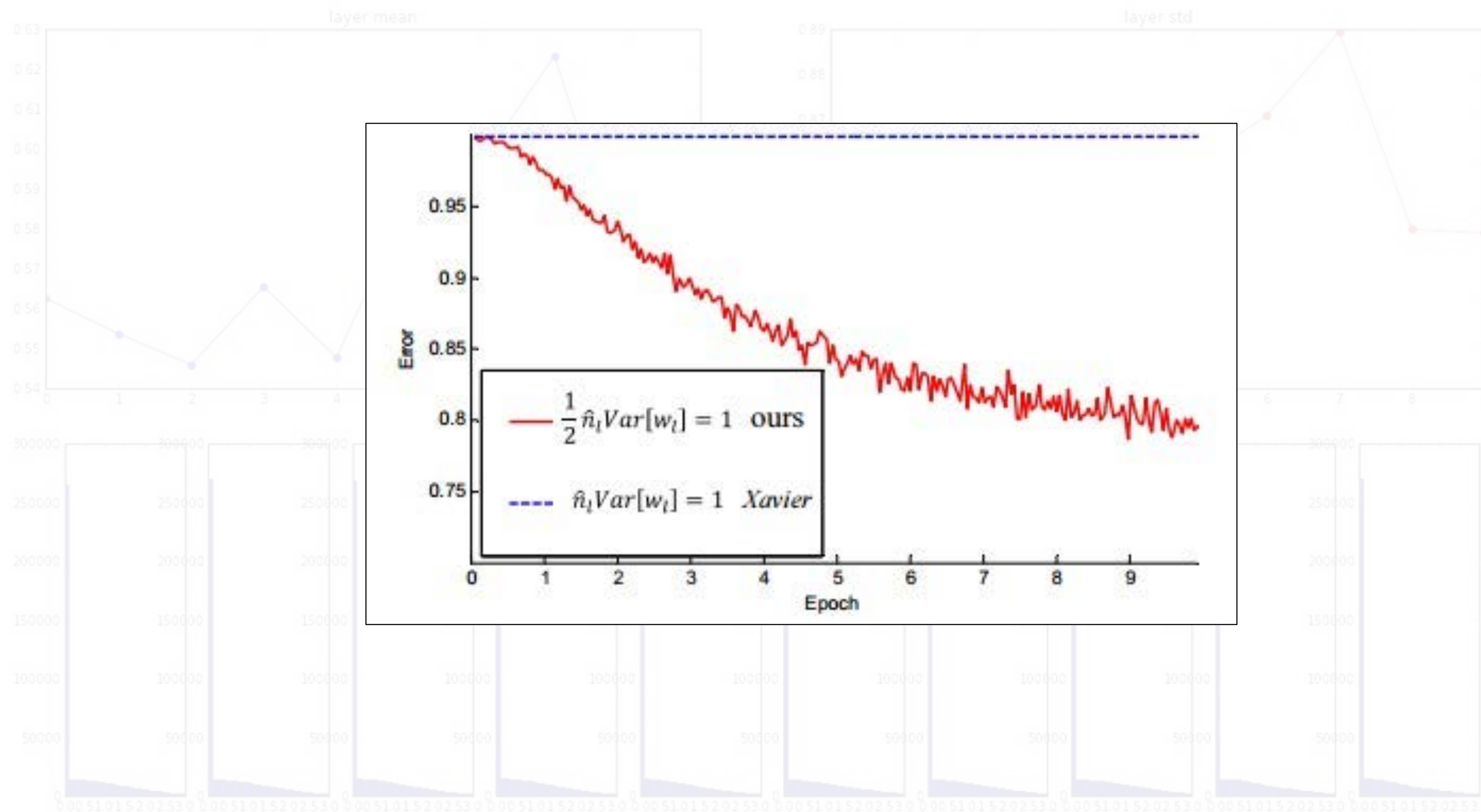


Weight initialization

He initialization [He et al., 2015]

```
W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in/2) # layer initialization
```

magic number 2

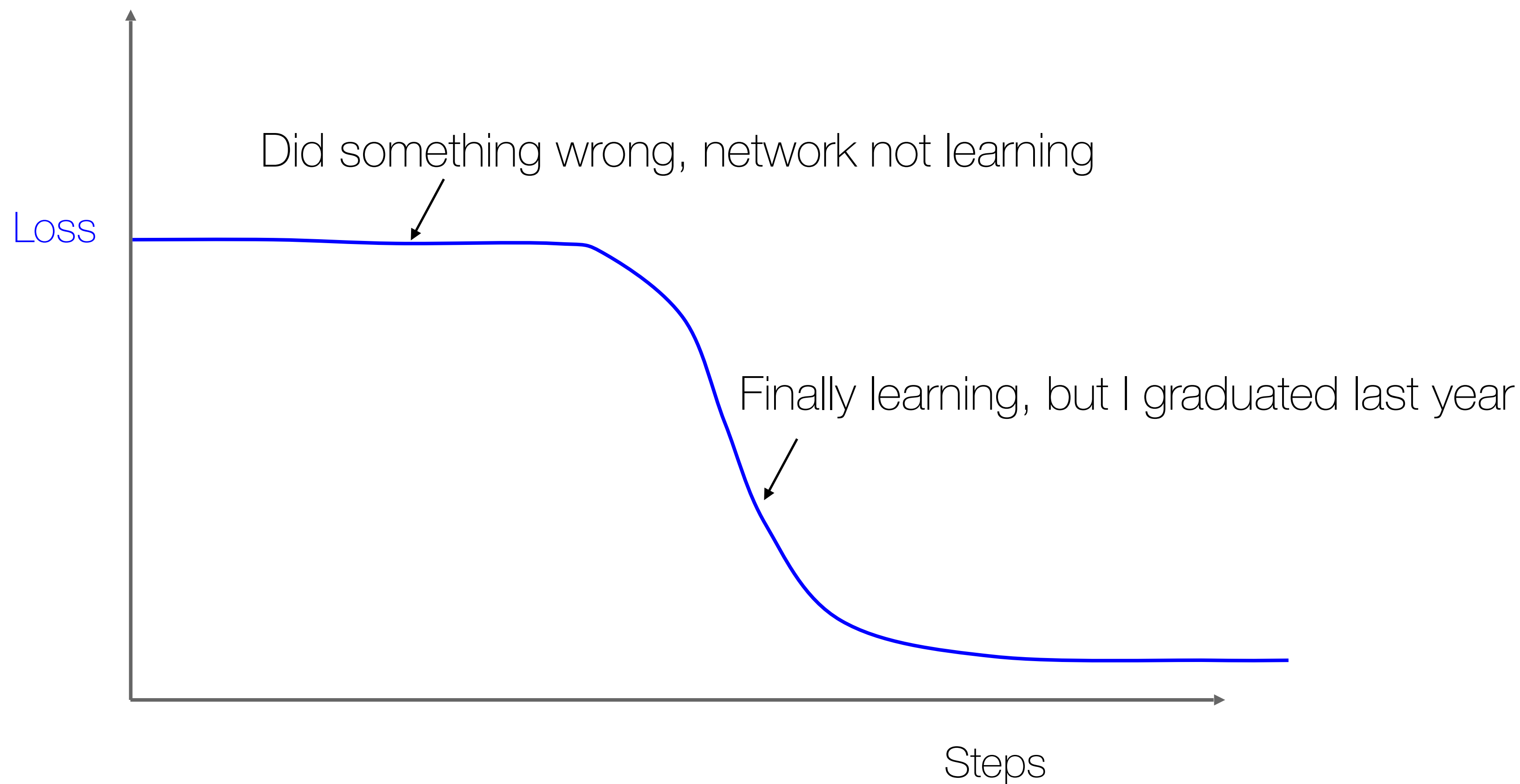


Statistically motivated

Good for ReLU

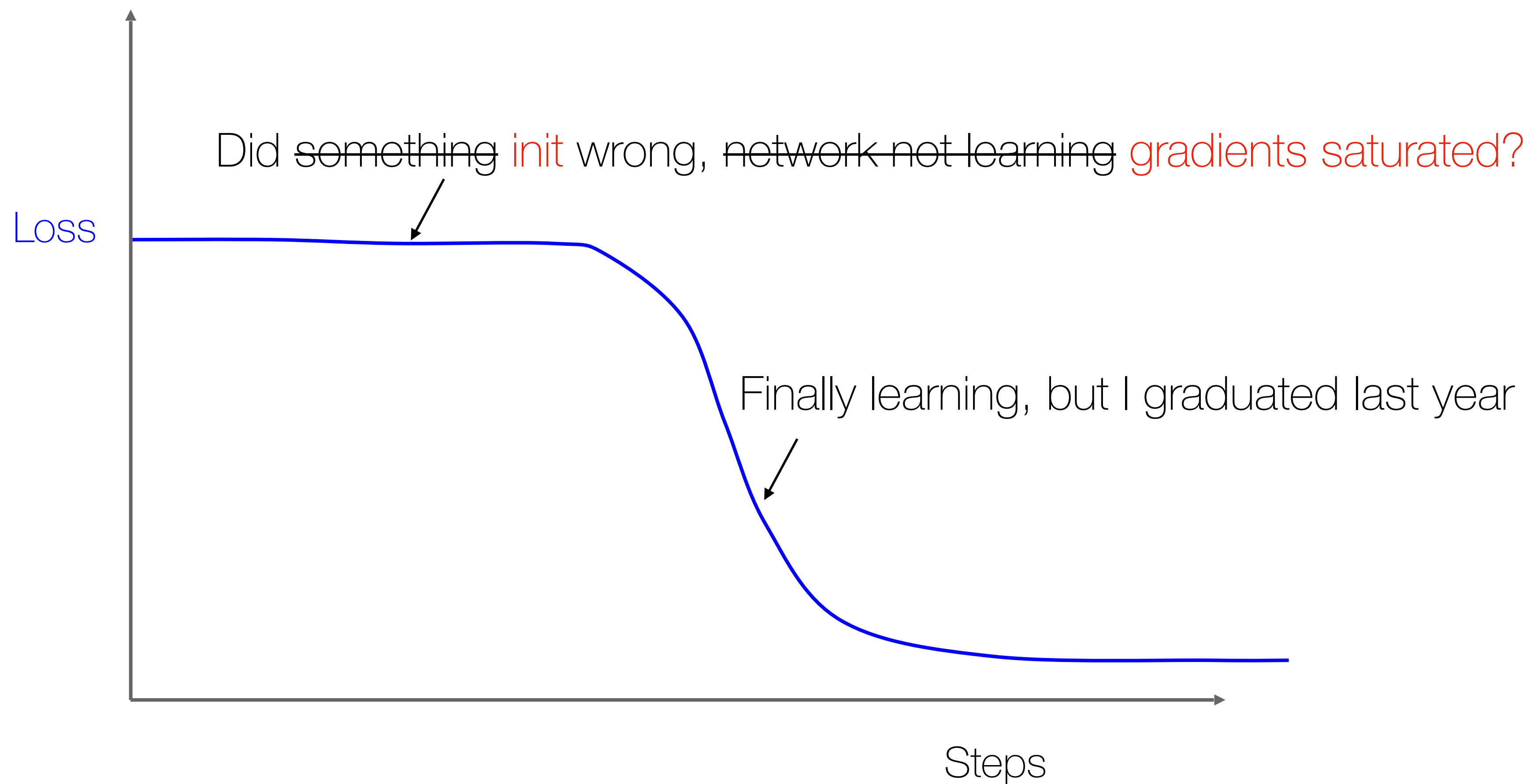
Recall: But it is never that easy

A typical sad loss curve



Recall: But it is never that easy

A typical sad loss curve



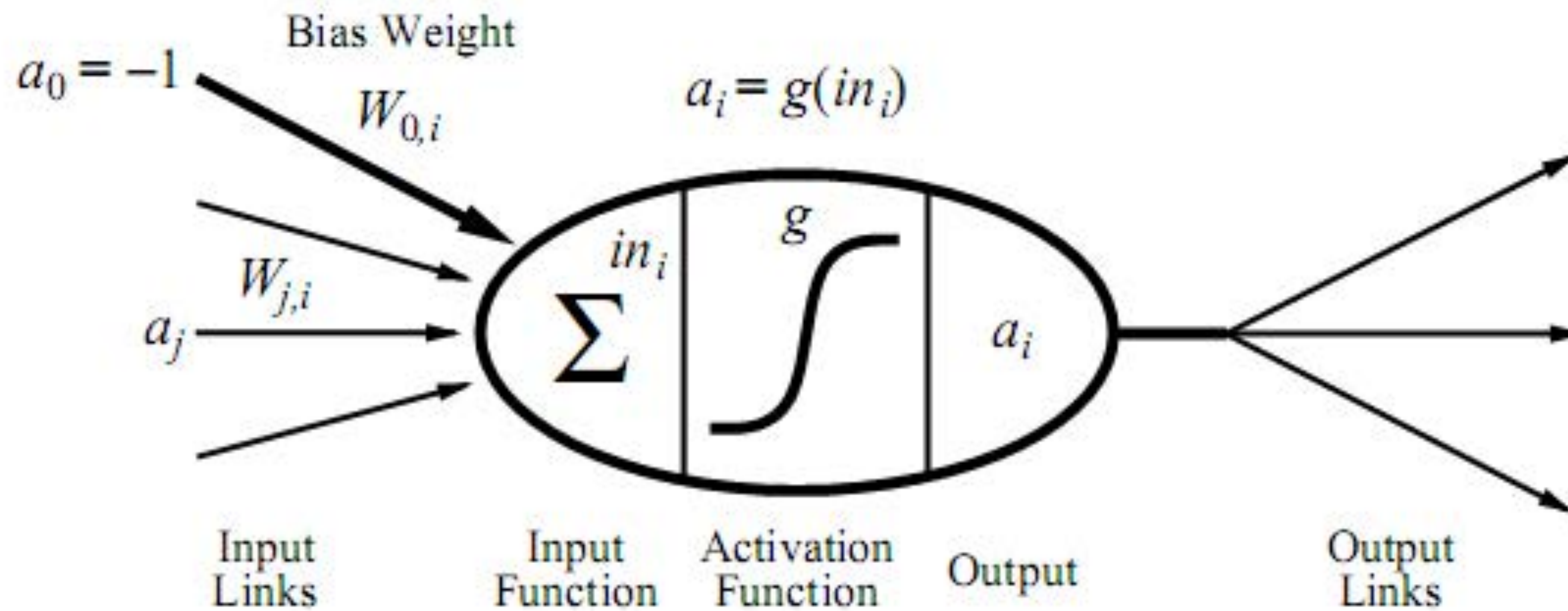


Normalization

Batch normalization [Ioffe and Szegedy, 2015]

Recall...

$$a_i \leftarrow g(in_i) = g(\sum_j W_{j,i} a_j)$$

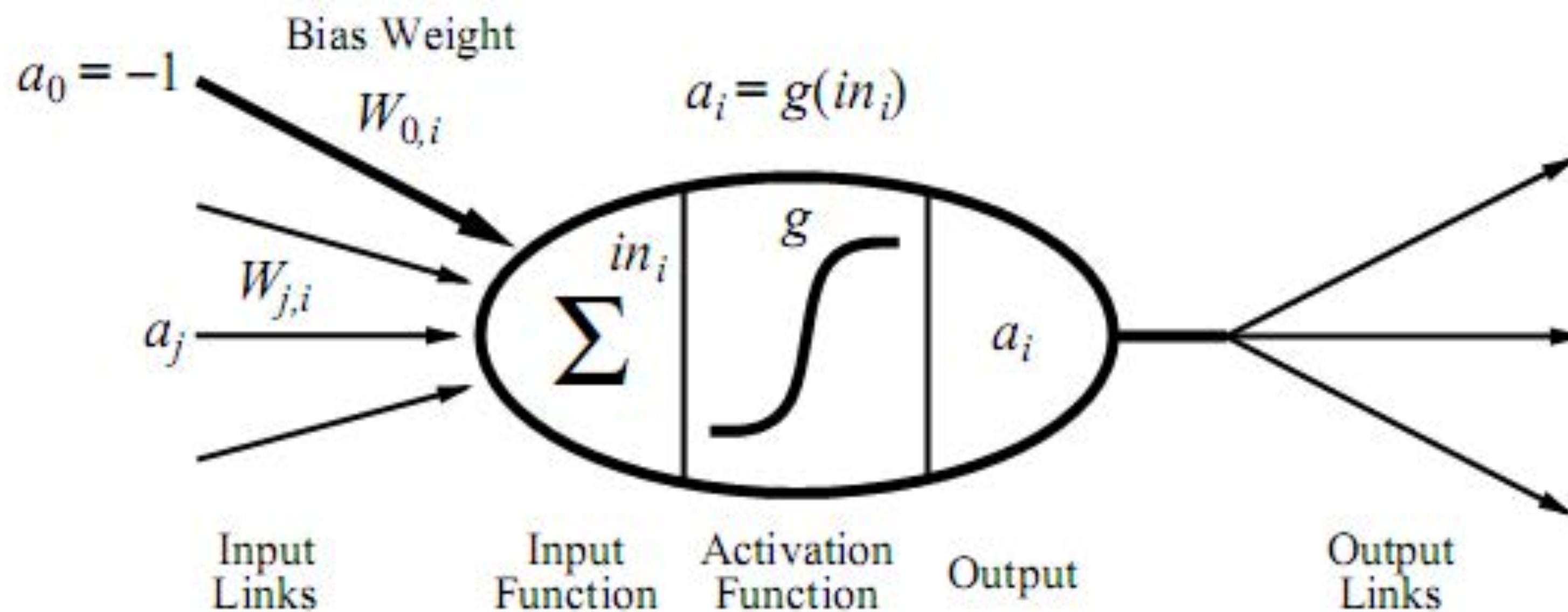


Batch normalization [Ioffe and Szegedy, 2015]

Recall...

Linear operations should cancel out

$$a_i \leftarrow g(in_i) = g(\sum_j W_{j,i} a_j)$$



Batch normalization [Ioffe and Szegedy, 2015]

Forcing a zero-mean and unit standard deviation

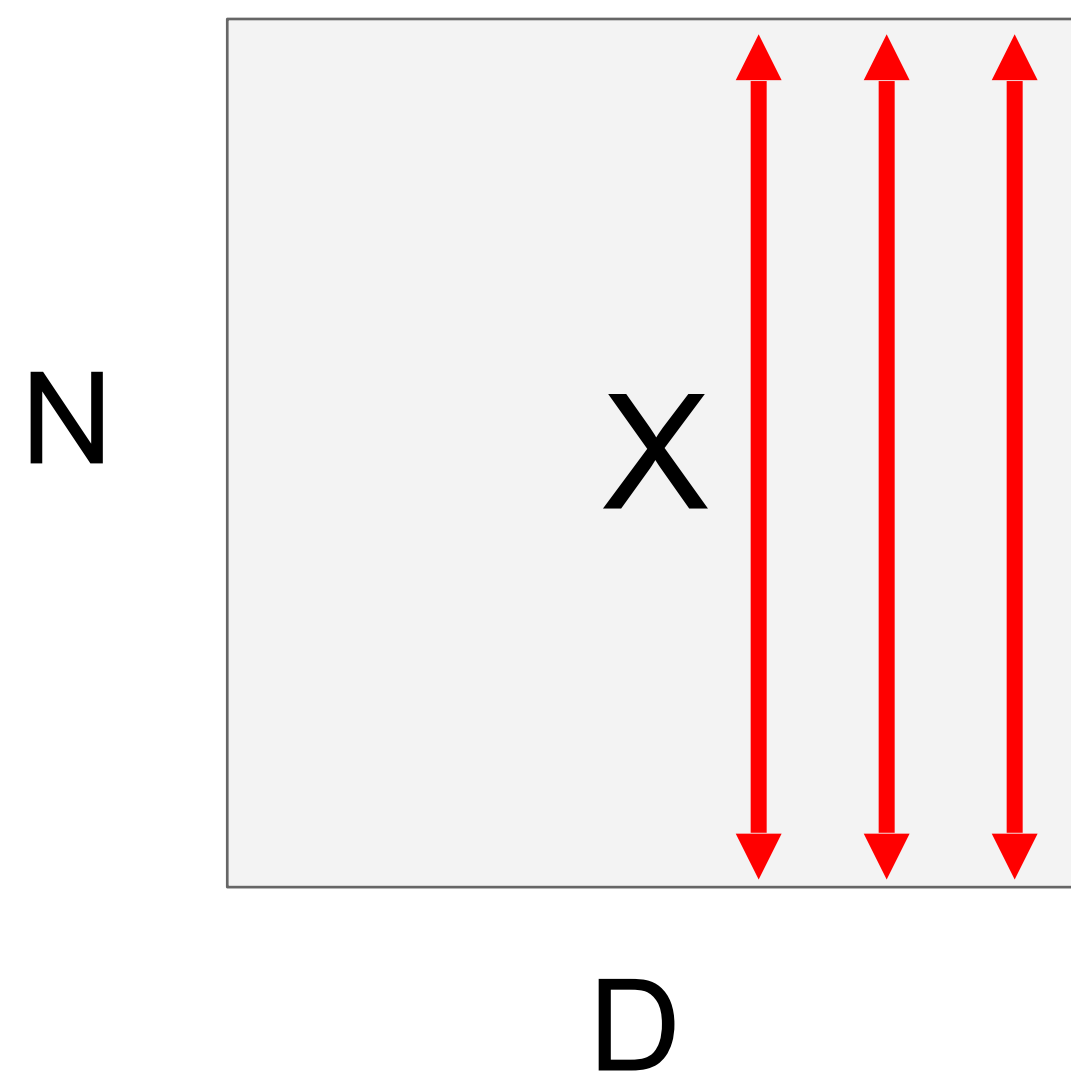
consider a batch of activations at some layer. To make each dimension unit gaussian, apply:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbf{E}[x^{(k)}]}{\sqrt{\mathbf{Var}[x^{(k)}]}}$$

this is a linear differentiable function...

Batch normalization [Ioffe and Szegedy, 2015]

Forcing a zero-mean and unit standard deviation



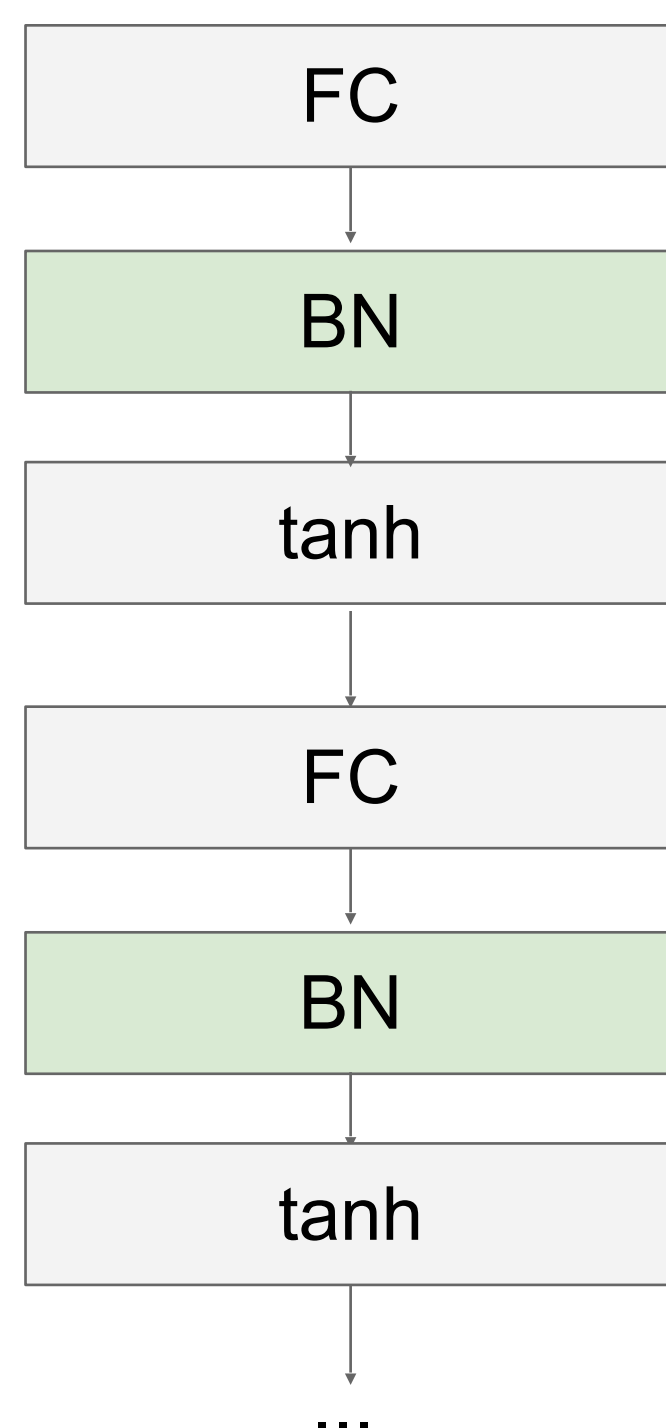
1. Compute the empirical mean and variance independently for each dimension.

2. Normalize

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Batch normalization [Ioffe and Szegedy, 2015]

Forcing a zero-mean and unit standard deviation



Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Batch normalization [Ioffe and Szegedy, 2015]

Introducing learnable scale / shift

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbf{E}[x^{(k)}]}{\sqrt{\mathbf{Var}[x^{(k)}]}}$$

And then allow the network to squash the range if it wants to:

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

Note, the network can learn:

$$\gamma^{(k)} = \sqrt{\mathbf{Var}[x^{(k)}]}$$

$$\beta^{(k)} = \mathbf{E}[x^{(k)}]$$

to recover the identity mapping.

Batch normalization [Ioffe and Szegedy, 2015]

Introducing learnable scale / shift

IMPORTANT: At test time, we don't have these — use training time stats

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbf{E}[x^{(k)}]}{\sqrt{\mathbf{Var}[x^{(k)}]}}$$

And then allow the network to squash the range if it wants to:

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

Note, the network can learn:

$$\gamma^{(k)} = \sqrt{\mathbf{Var}[x^{(k)}]}$$

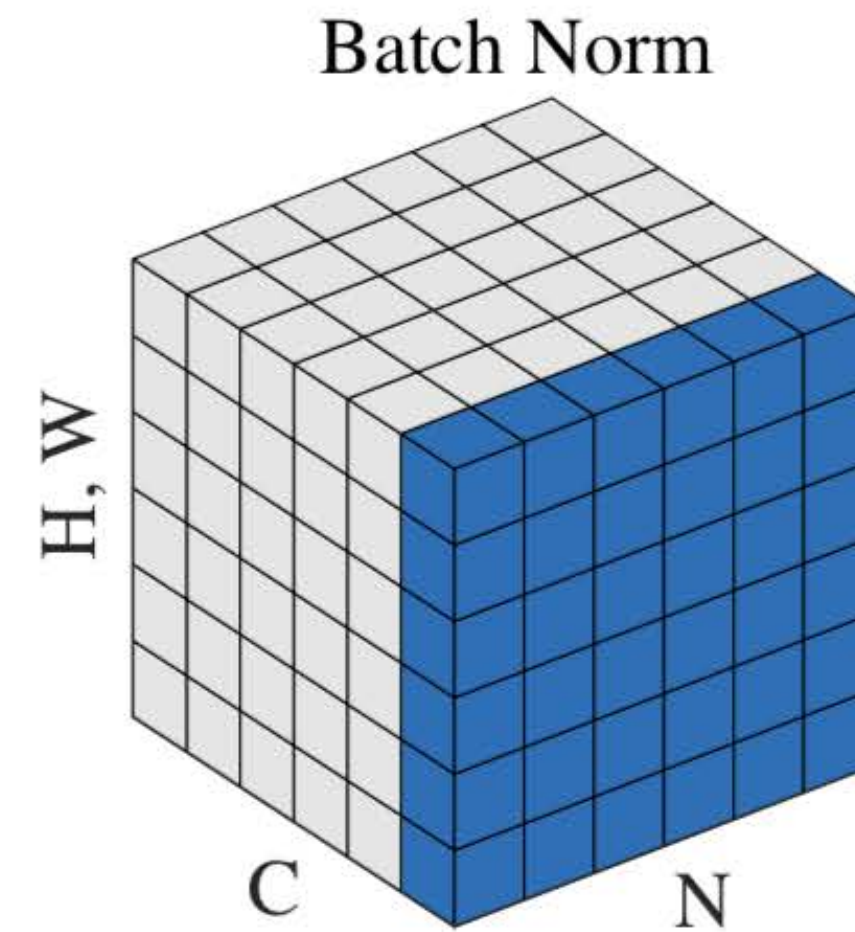
$$\beta^{(k)} = \mathbf{E}[x^{(k)}]$$

to recover the identity mapping.

Other normalization techniques

Batch Normalization

Skipped in class
(outside of scope)



Skipped in class
(outside of scope)

Other normalization techniques

Batch Normalization

Batch Normalization for **fully-connected** networks

$$\begin{array}{l}
 \mathbf{x} : \mathbf{N} \times \mathbf{D} \\
 \text{Normalize} \quad \downarrow \\
 \boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{1} \times \mathbf{D} \\
 \boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{D} \\
 \mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}
 \end{array}$$

Batch Normalization for **convolutional** networks
(Spatial Batchnorm, BatchNorm2D)

$$\begin{array}{l}
 \mathbf{x} : \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W} \\
 \text{Normalize} \quad \downarrow \quad \downarrow \quad \downarrow \\
 \boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\
 \boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\
 \mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}
 \end{array}$$

Skipped in class
(outside of scope)

Other normalization techniques

Batch Normalization

Batch Normalization for **fully-connected** networks

$\mathbf{x} : \mathbf{N} \times \mathbf{D}$

Normalize \downarrow

$\boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{1} \times \mathbf{D}$

$\boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{D}$

$\mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}$

Batch Normalization for **convolutional** networks
(Spatial Batchnorm, BatchNorm2D)

$\mathbf{x} : \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W}$

Normalize $\downarrow \downarrow \downarrow$

$\boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1}$

$\boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1}$

$\mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}$

This is why train/test needs to be different

Skipped in class
(outside of scope)

Other normalization techniques

Batch Normalization

Batch Normalization for **fully-connected** networks

$\mathbf{x} : \mathbf{N} \times \mathbf{D}$

Normalize

$\boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{1} \times \mathbf{D}$

$\boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{D}$

$\mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}$

Batch Normalization for **convolutional** networks
(Spatial Batchnorm, BatchNorm2D)

$\mathbf{x} : \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W}$

Normalize

$\boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1}$

$\boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1}$

$\mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}$

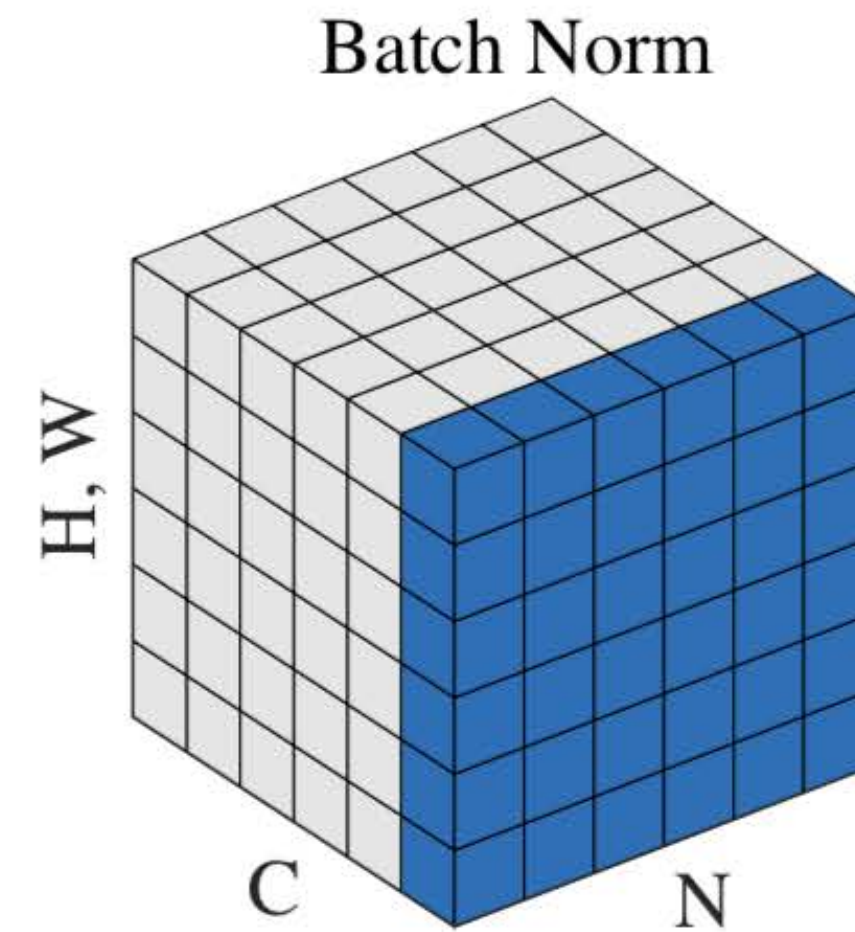
Always watch out when implementing!!!

This is why train/test needs to be different

Other normalization techniques

Batch Normalization

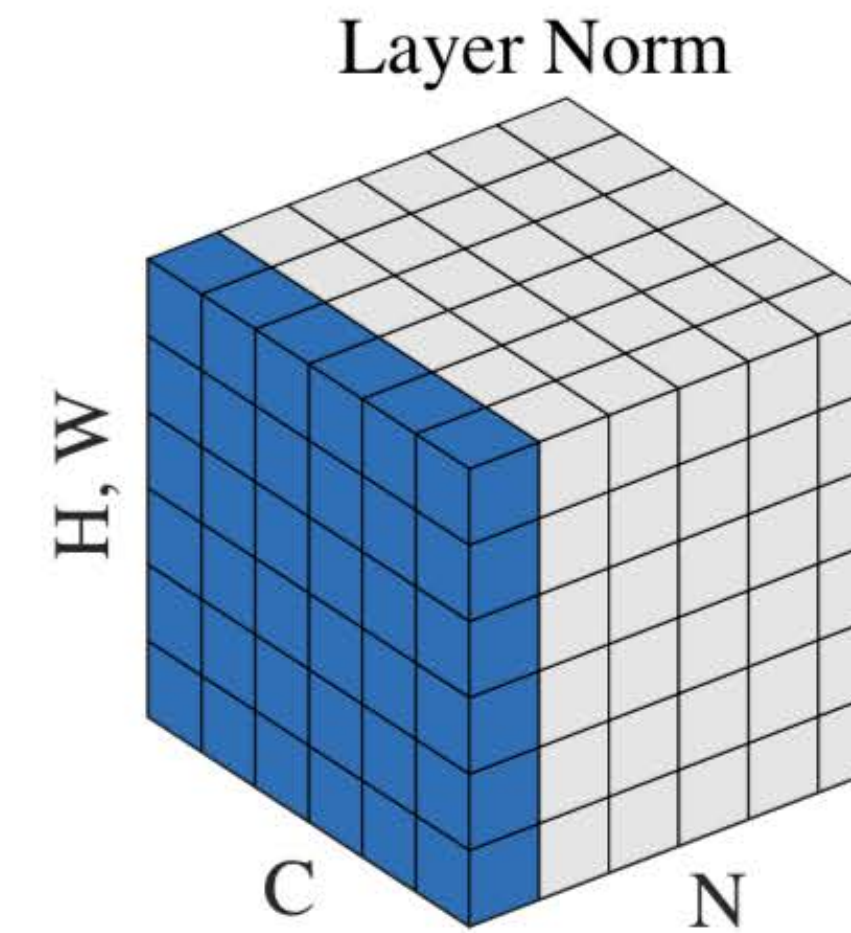
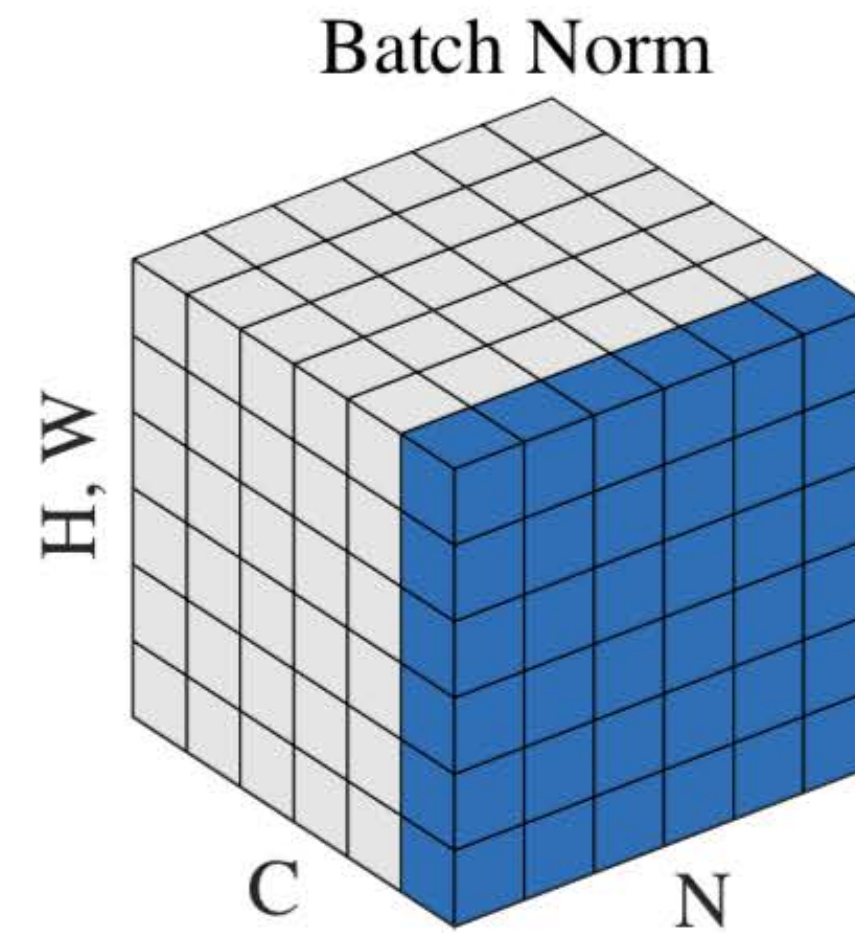
Skipped in class
(outside of scope)



Skipped in class
(outside of scope)

Other normalization techniques

Layer Normalization



Skipped in class
(outside of scope)

Other normalization techniques

Layer Normalization

Batch Normalization for **fully-connected** networks

$$\begin{aligned}
 & \mathbf{x} : \mathbf{N} \times \mathbf{D} \\
 \text{Normalize} & \quad \downarrow \\
 & \boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{1} \times \mathbf{D} \\
 & \boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{D} \\
 & \mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}
 \end{aligned}$$

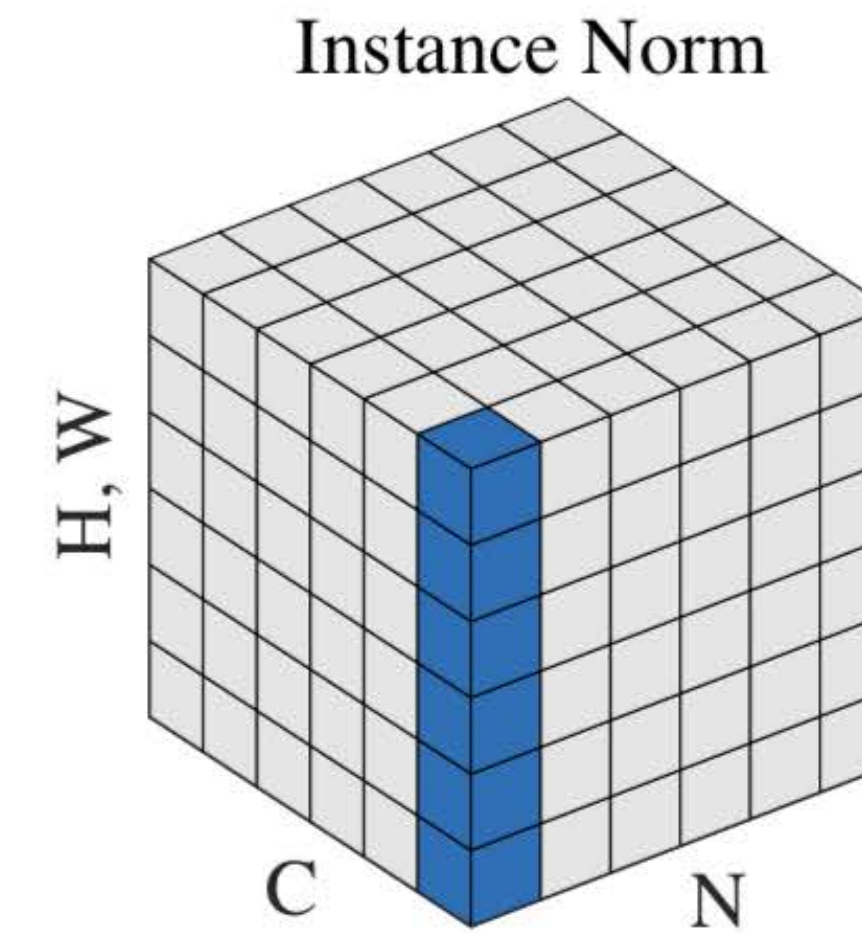
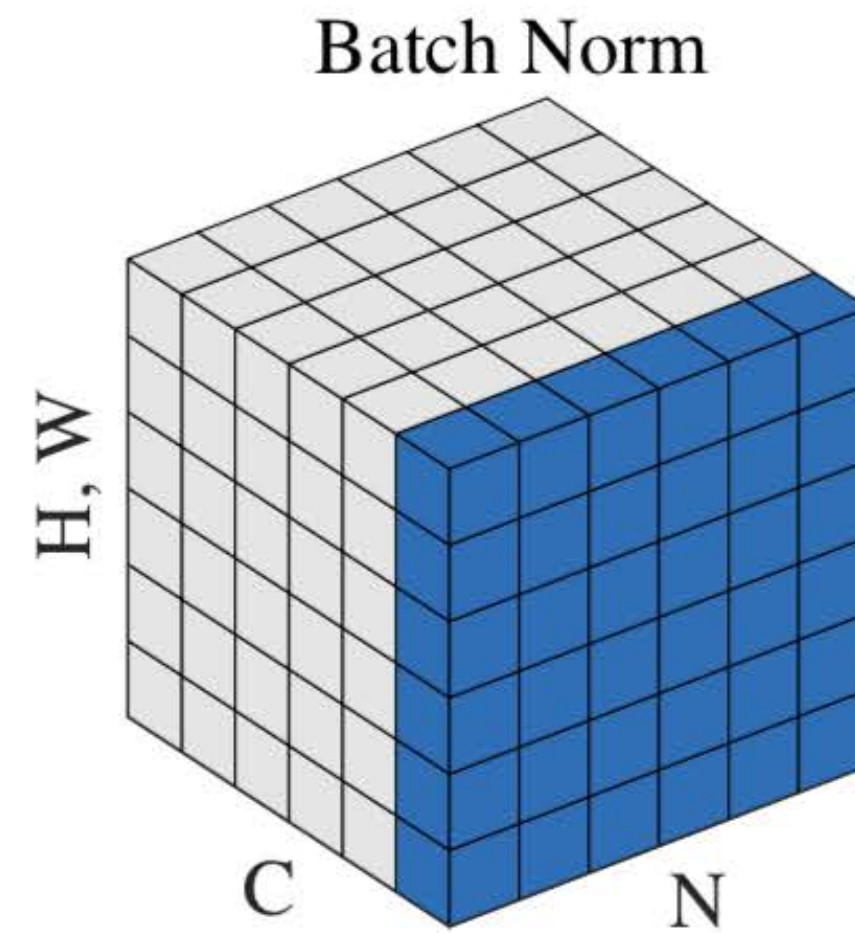
Layer Normalization for fully-connected networks
Same behavior at train and test!
Can be used in recurrent networks

$$\begin{aligned}
 & \mathbf{x} : \mathbf{N} \times \mathbf{D} \\
 \text{Normalize} & \quad \downarrow \\
 & \boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{N} \times \mathbf{1} \\
 & \boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{D} \\
 & \mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}
 \end{aligned}$$

Skipped in class
(outside of scope)

Other normalization techniques

Instance Normalization



Skipped in class
(outside of scope)

Other normalization techniques

Instance Normalization

Batch Normalization for convolutional networks

$$\begin{array}{l}
 \mathbf{x} : \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W} \\
 \text{Normalize} \quad \downarrow \quad \downarrow \quad \downarrow \\
 \boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\
 \boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\
 \mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}
 \end{array}$$

Instance Normalization for convolutional networks
Same behavior at train / test!

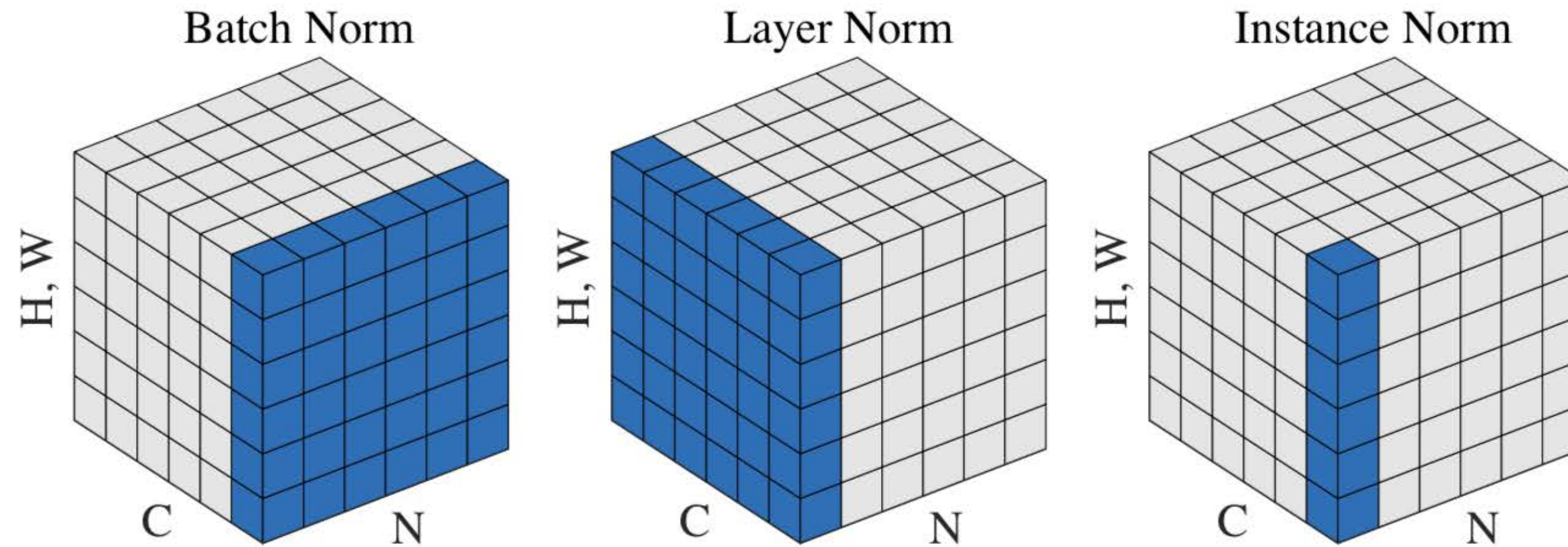
$$\begin{array}{l}
 \mathbf{x} : \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W} \\
 \text{Normalize} \quad \quad \downarrow \quad \downarrow \\
 \boldsymbol{\mu}, \boldsymbol{\sigma} : \mathbf{N} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\
 \boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\
 \mathbf{y} = \boldsymbol{\gamma} (\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}
 \end{array}$$

Ulyanov et al, Improved Texture Networks: Maximizing Quality and Diversity in Feed-forward Stylization and Texture Synthesis, CVPR 2017

Skipped in class
(outside of scope)

Other normalization techniques

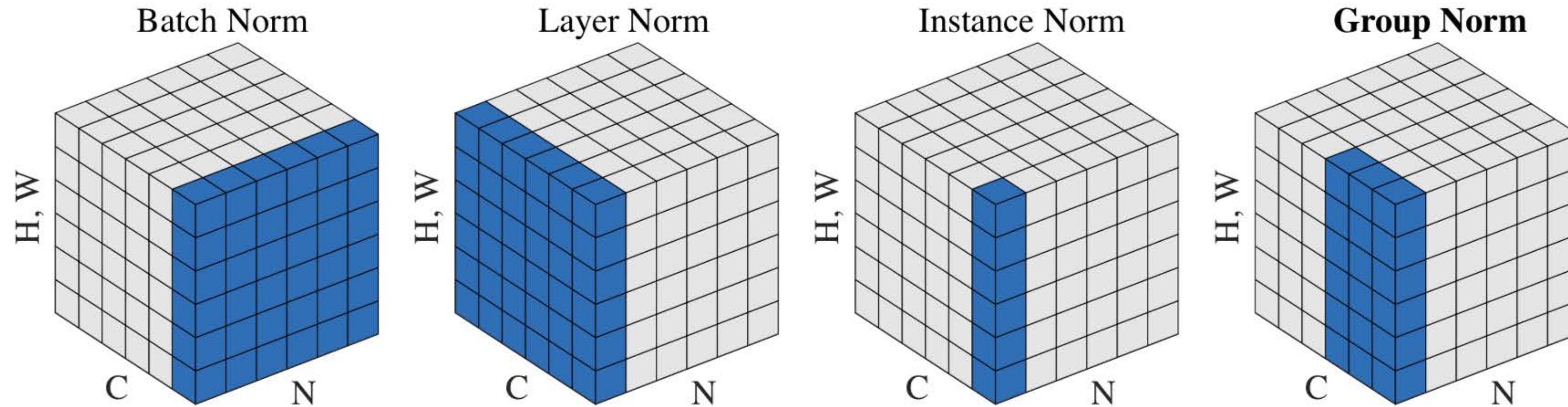
Group Normalization



Skipped in class
(outside of scope)

Other normalization techniques

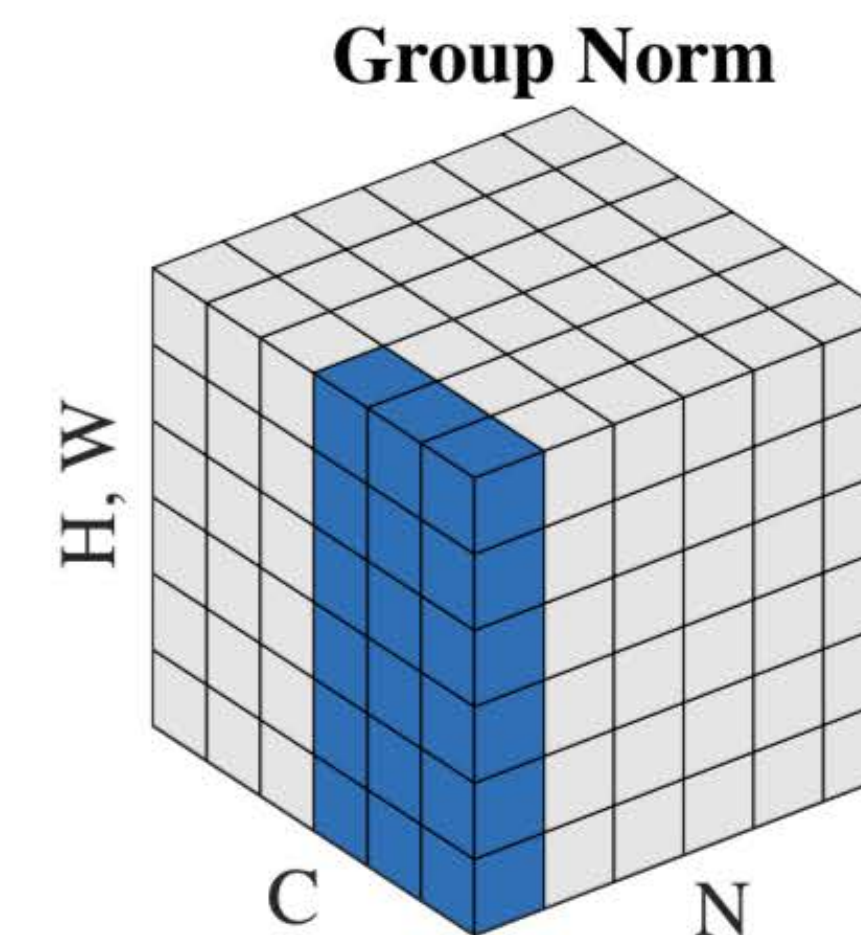
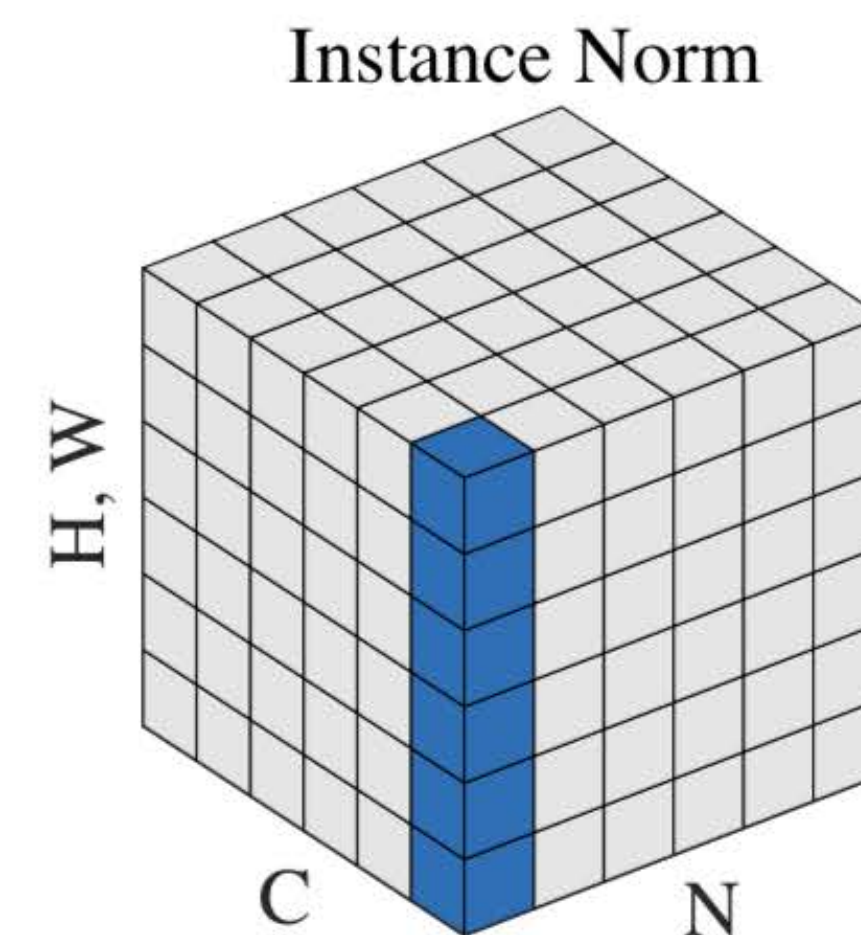
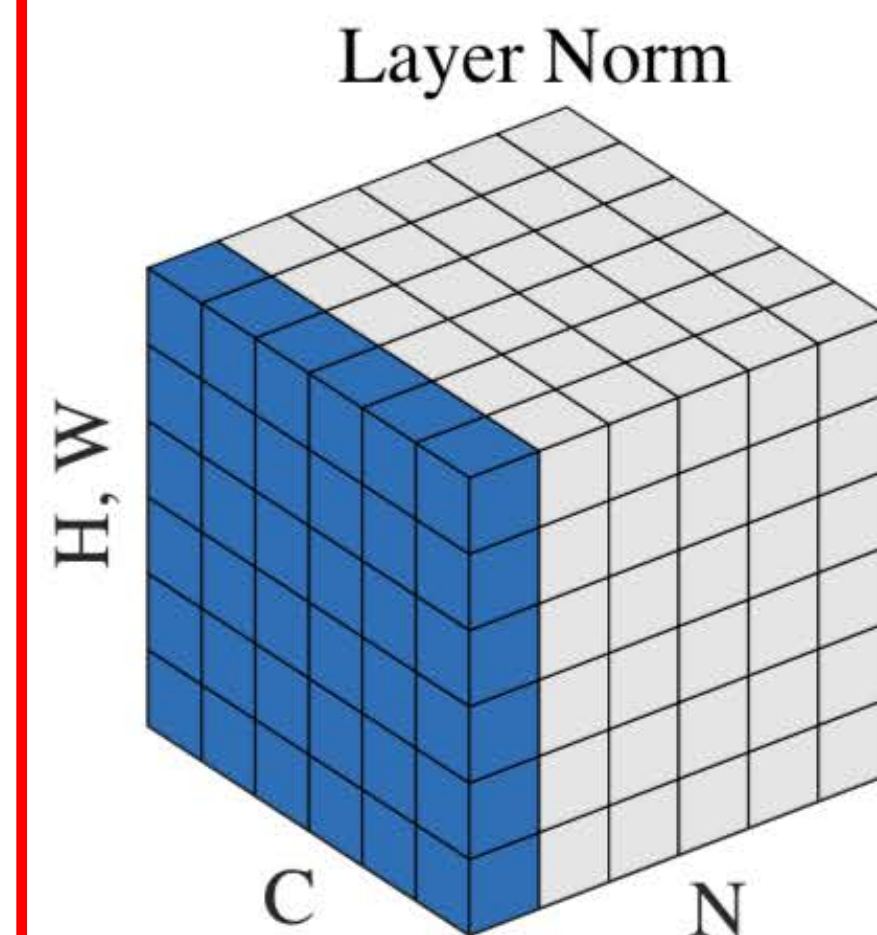
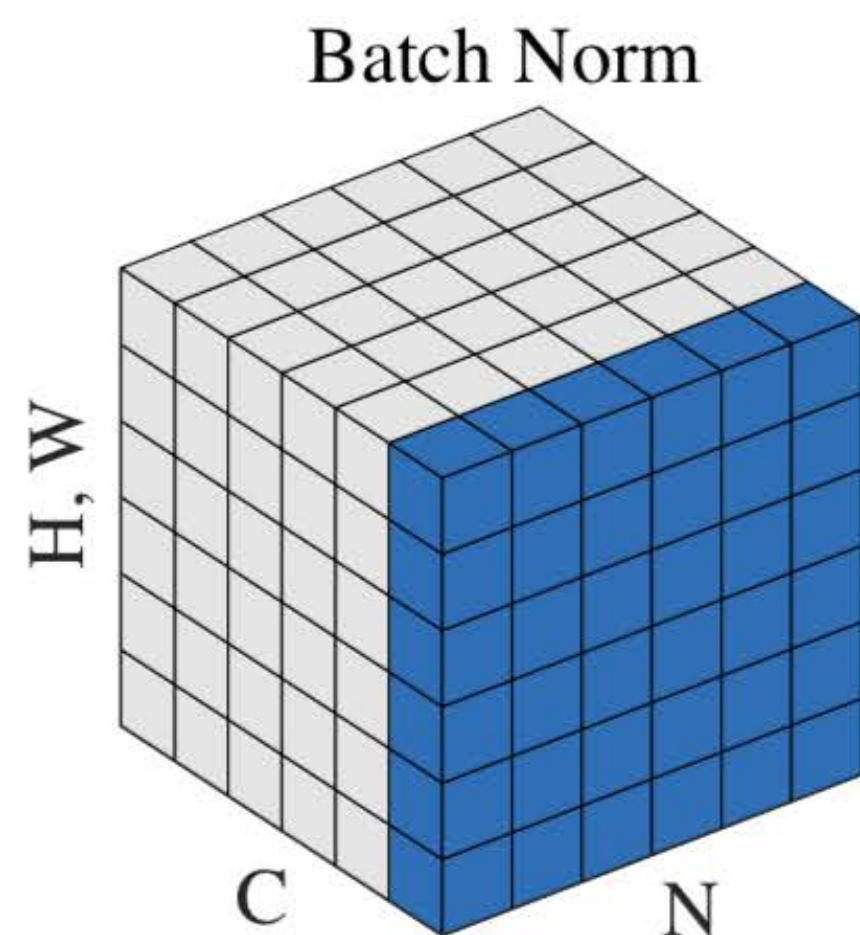
Group Normalization



Skipped in class
(outside of scope)

Other normalization techniques

Group Normalization



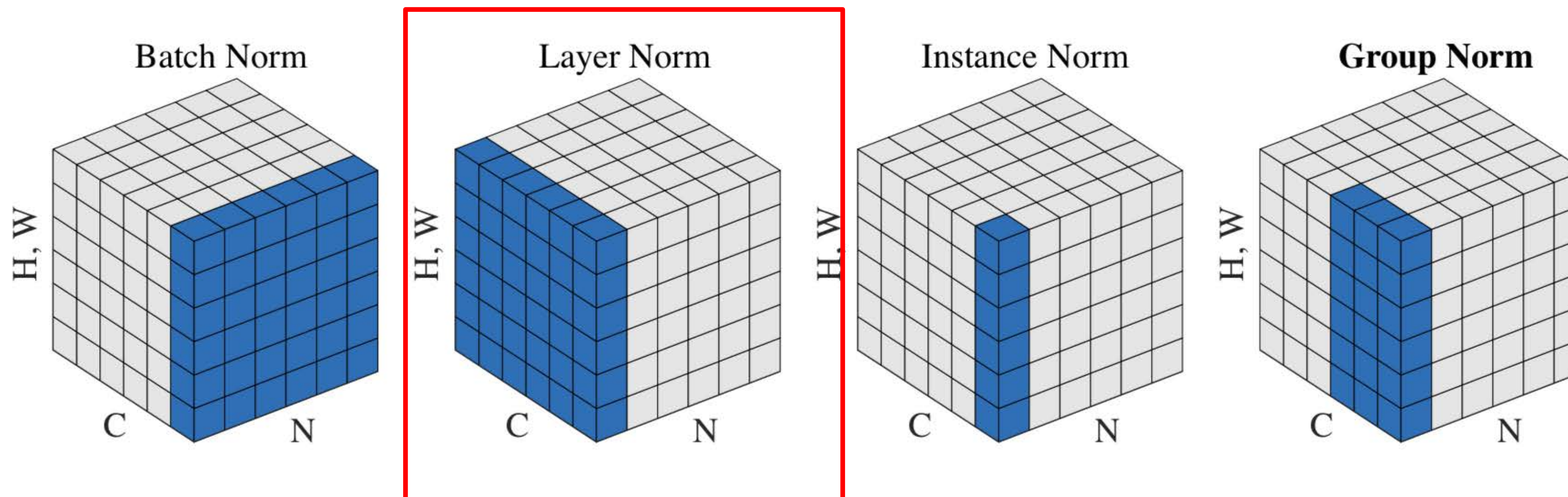
No train/test-time differences.

Much preferred in my opinion.

Skipped in class
(outside of scope)

Other normalization techniques

Group Normalization

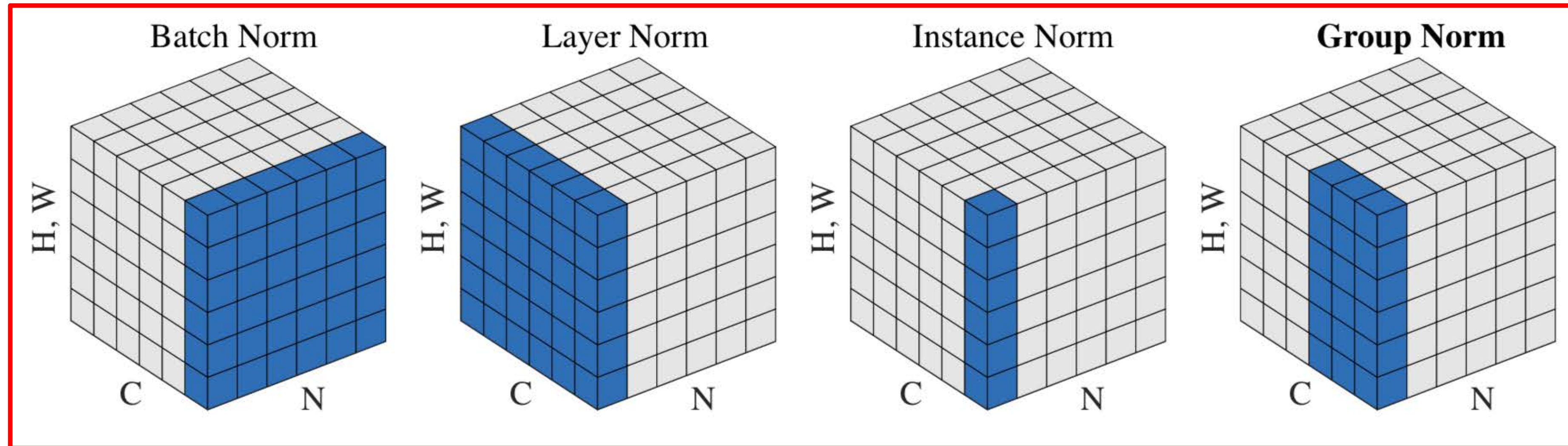


Can be implemented using PyTorch's Group norm.

Skipped in class
(outside of scope)

Other normalization techniques

Group Normalization

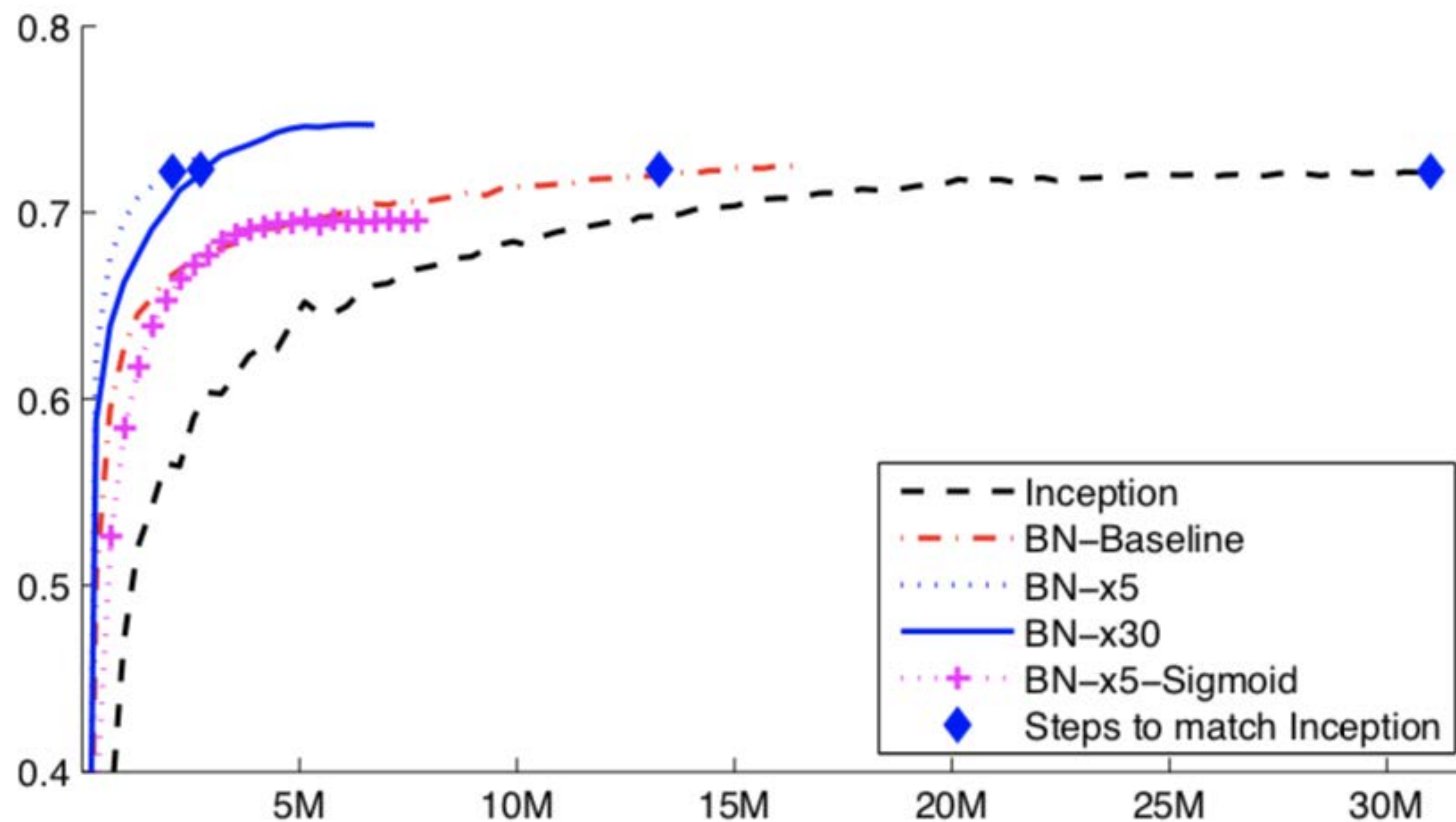


Choice of normalization should be data dependent



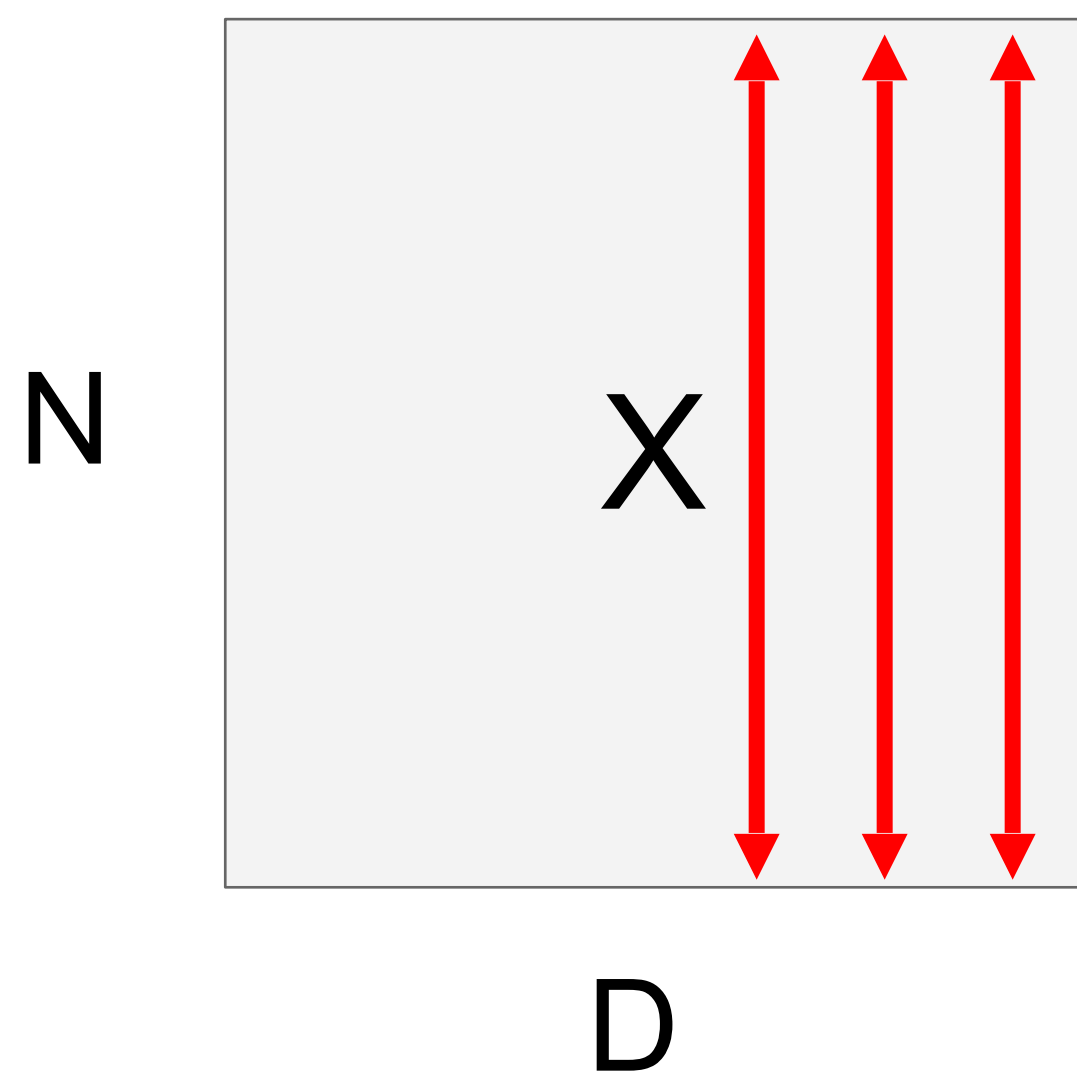
By the way . . . with normalization
something else also happens

Batch normalization



Batch normalization

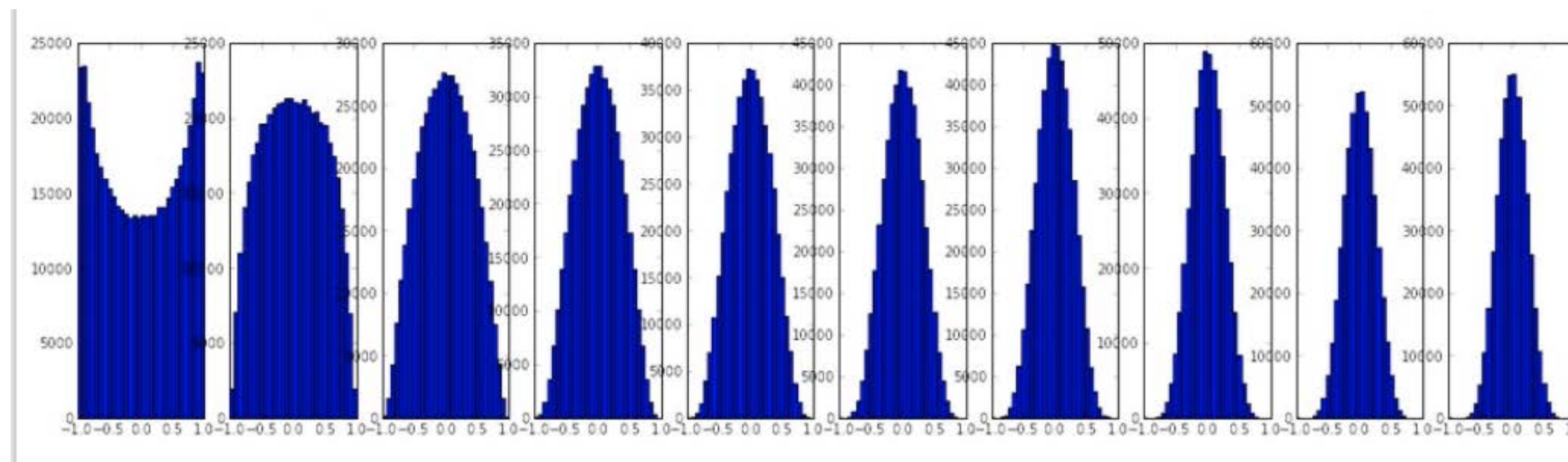
Recall...



1. compute the empirical mean and variance independently for each dimension.

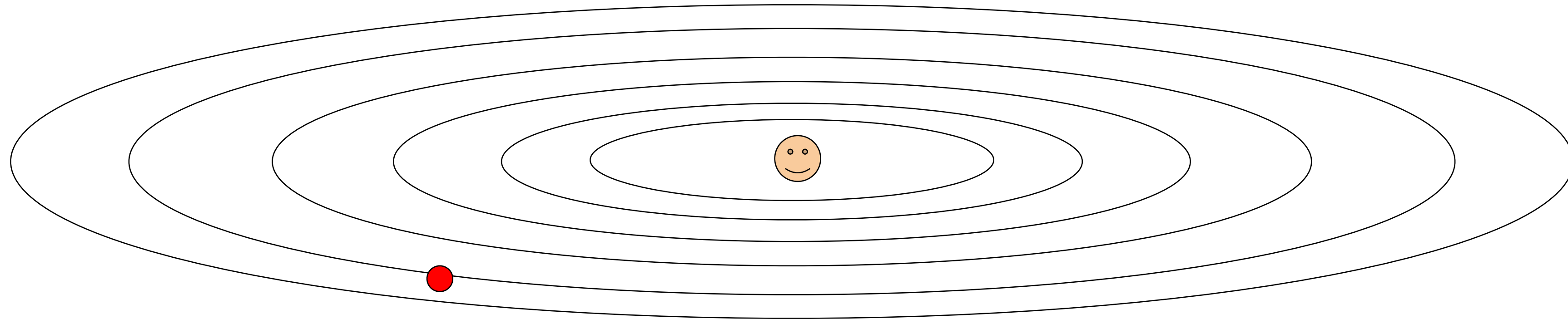
2. Normalize

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

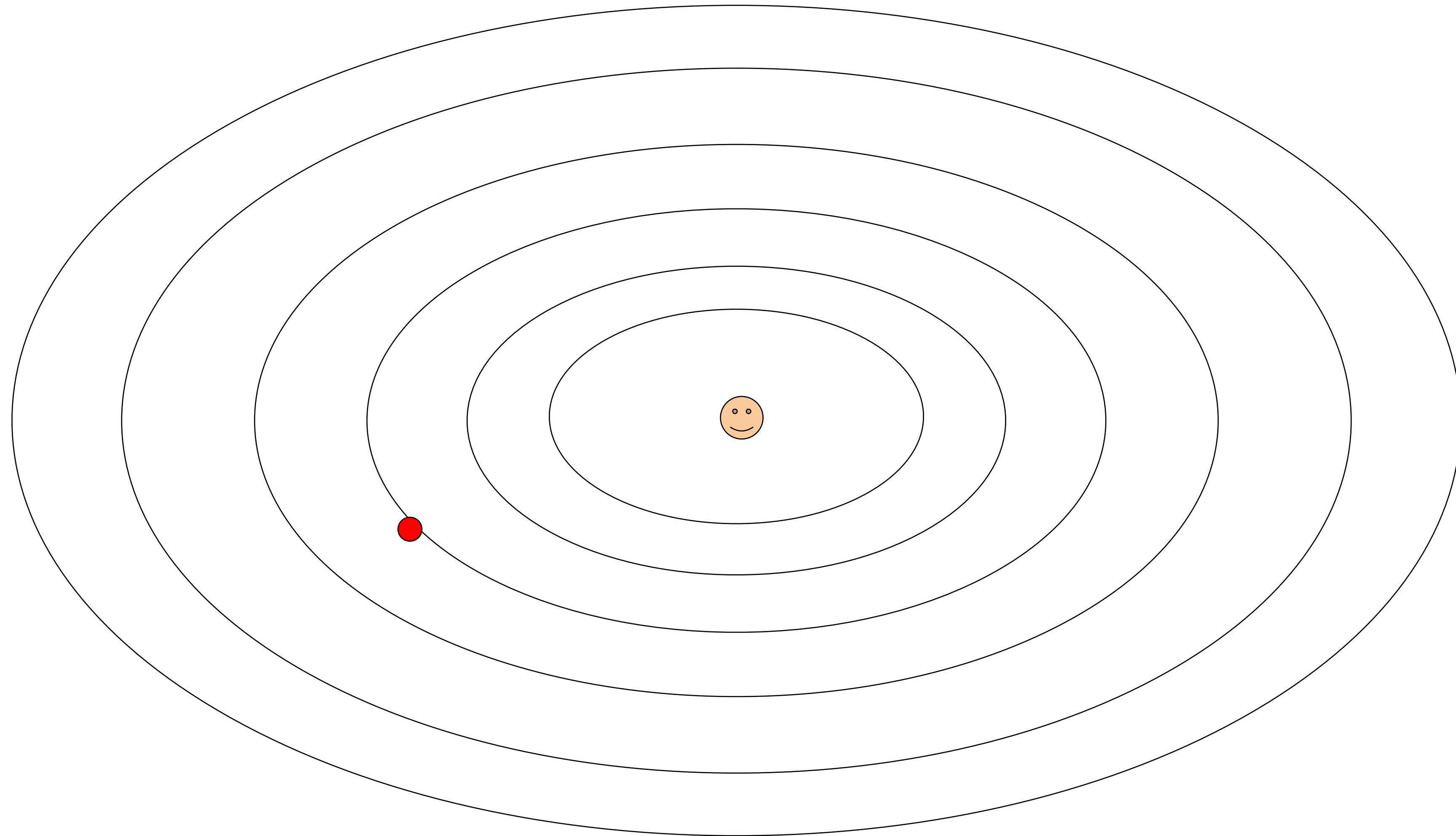


Batch normalization

This imbalance between
dimensions is the problem

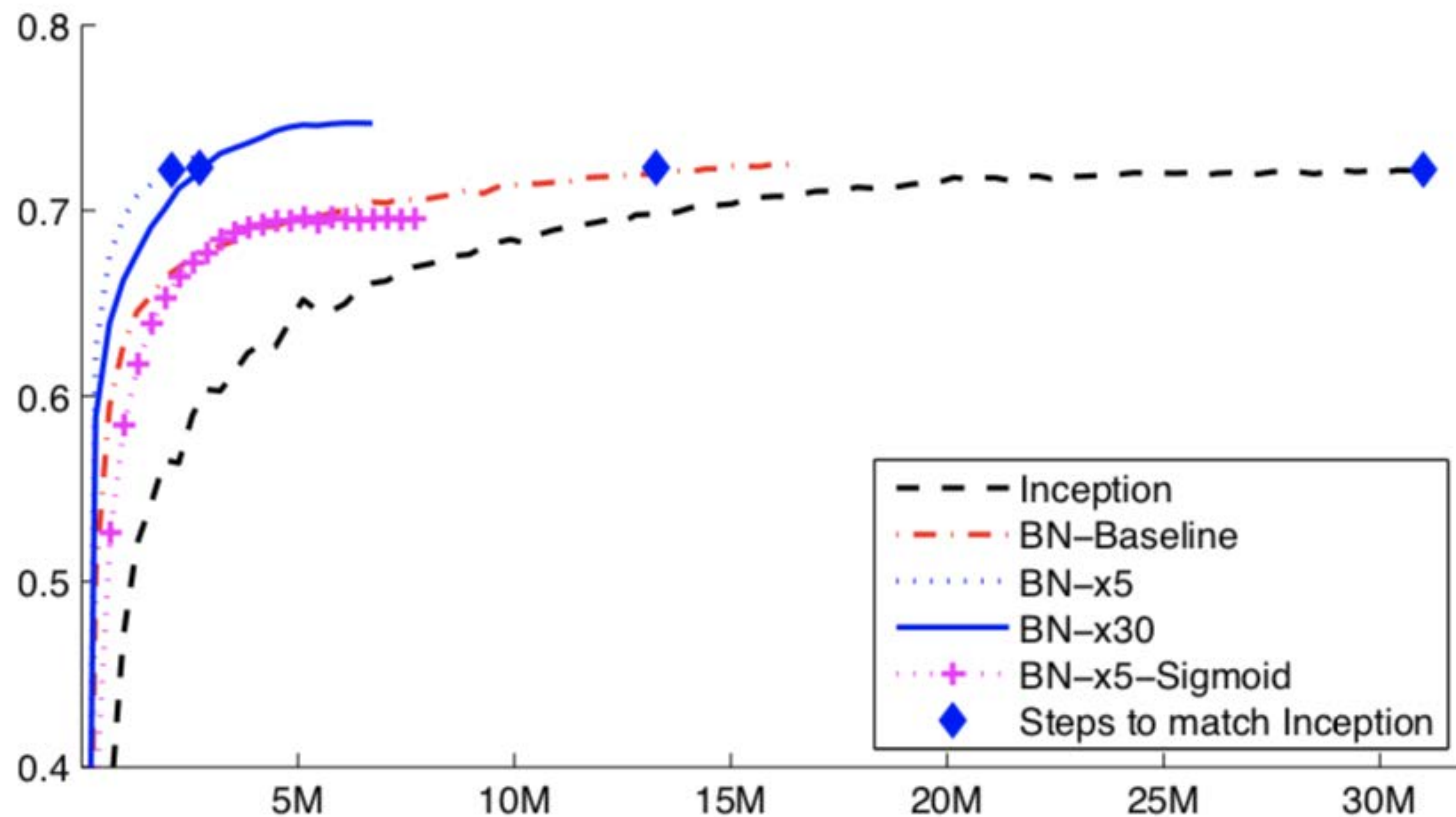


Batch normalization



Let's artificially make it like this!

Batch normalization

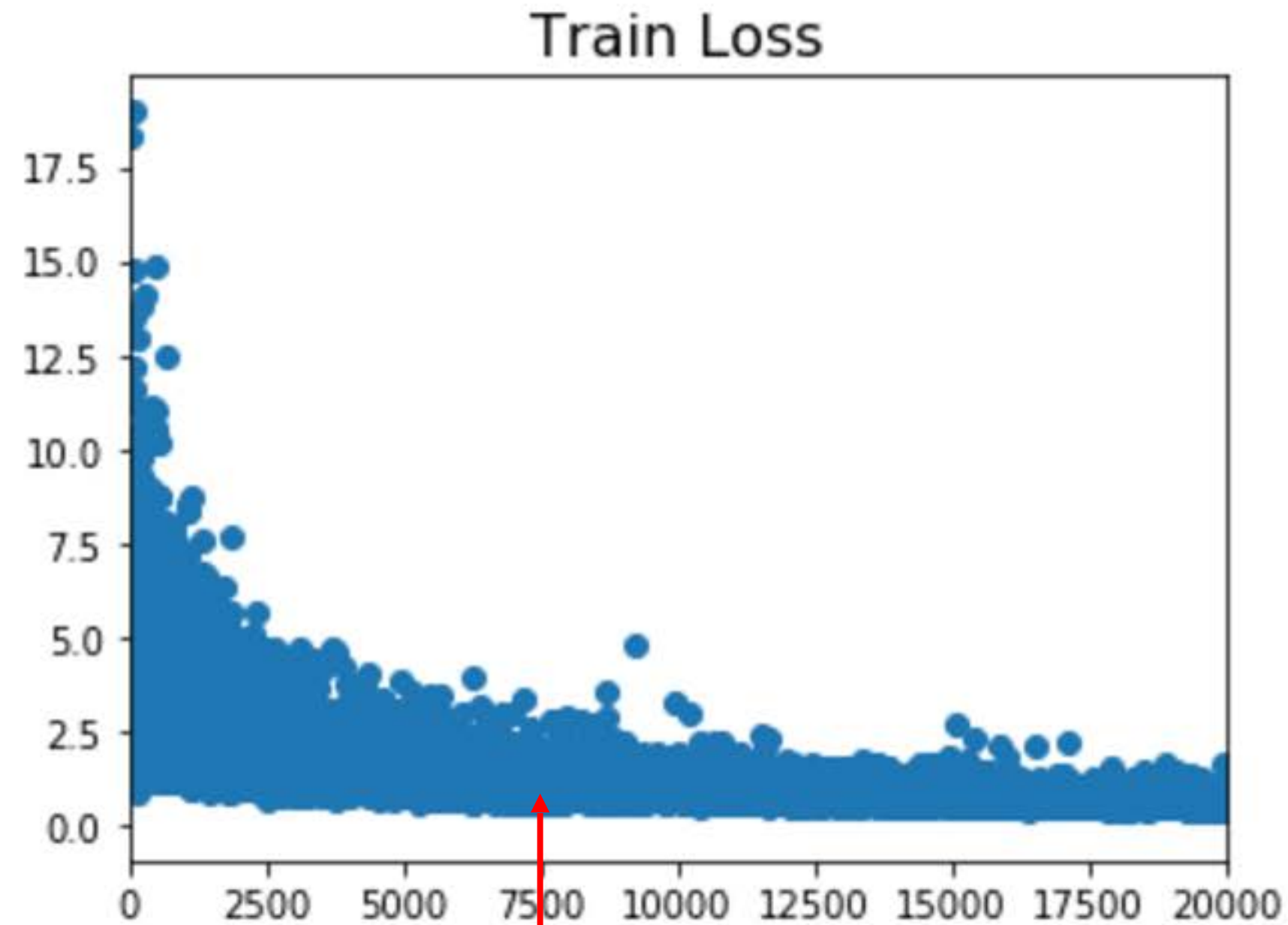




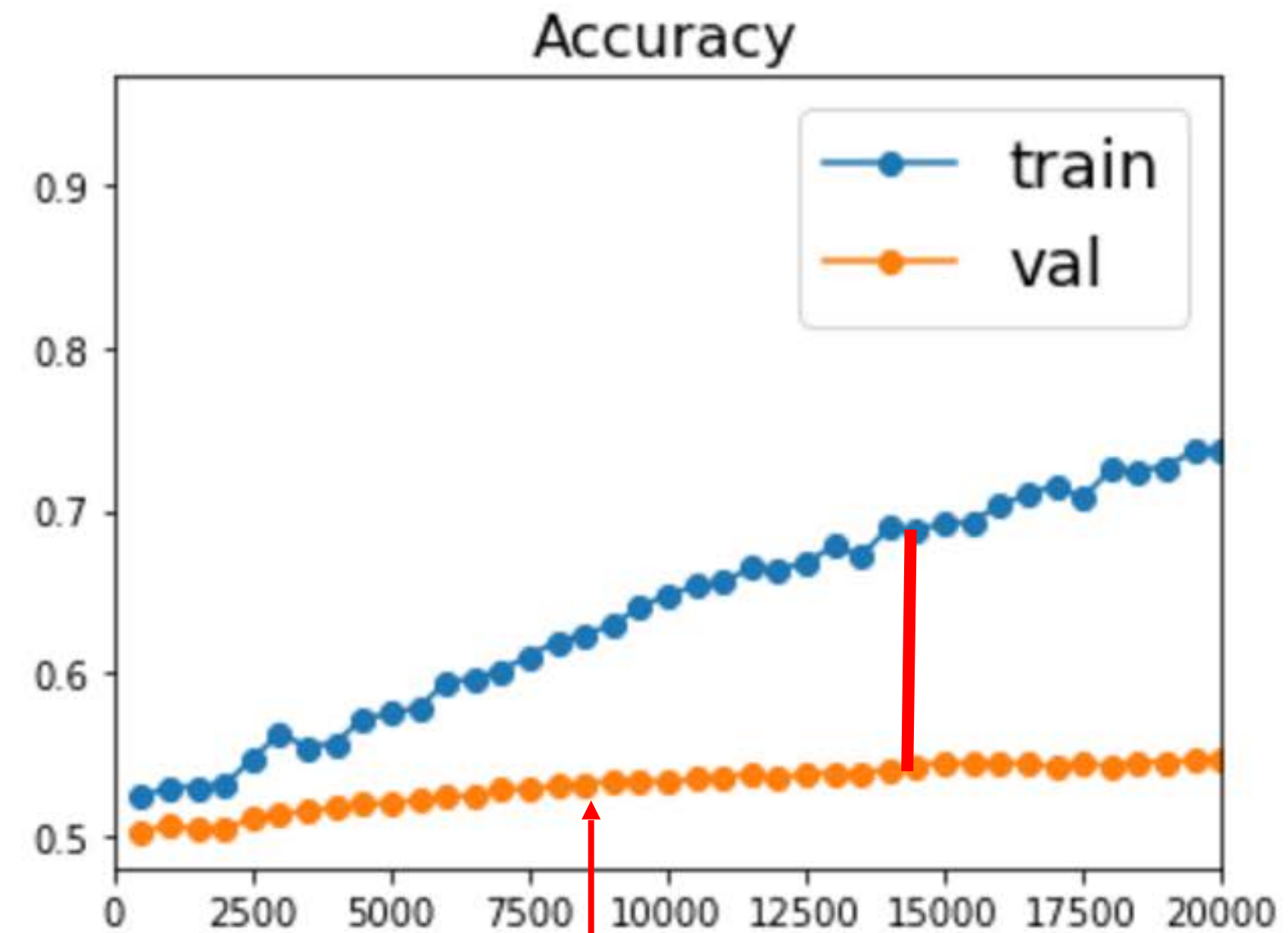
Preventing overfitting

Beyond training loss

Recall the other problem



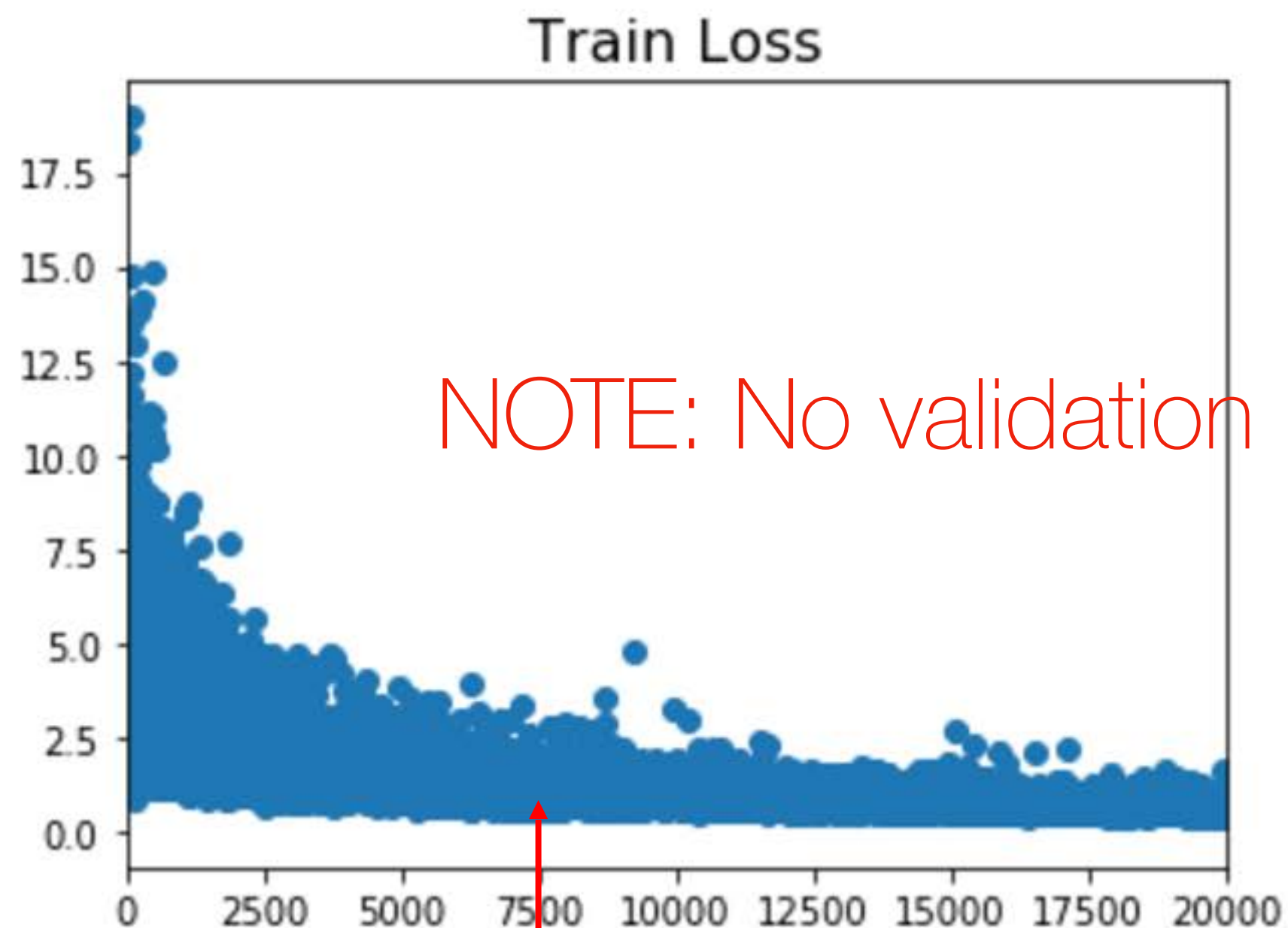
Better optimization algorithms help reduce training loss



But we really care about error on new data - how to reduce the gap?

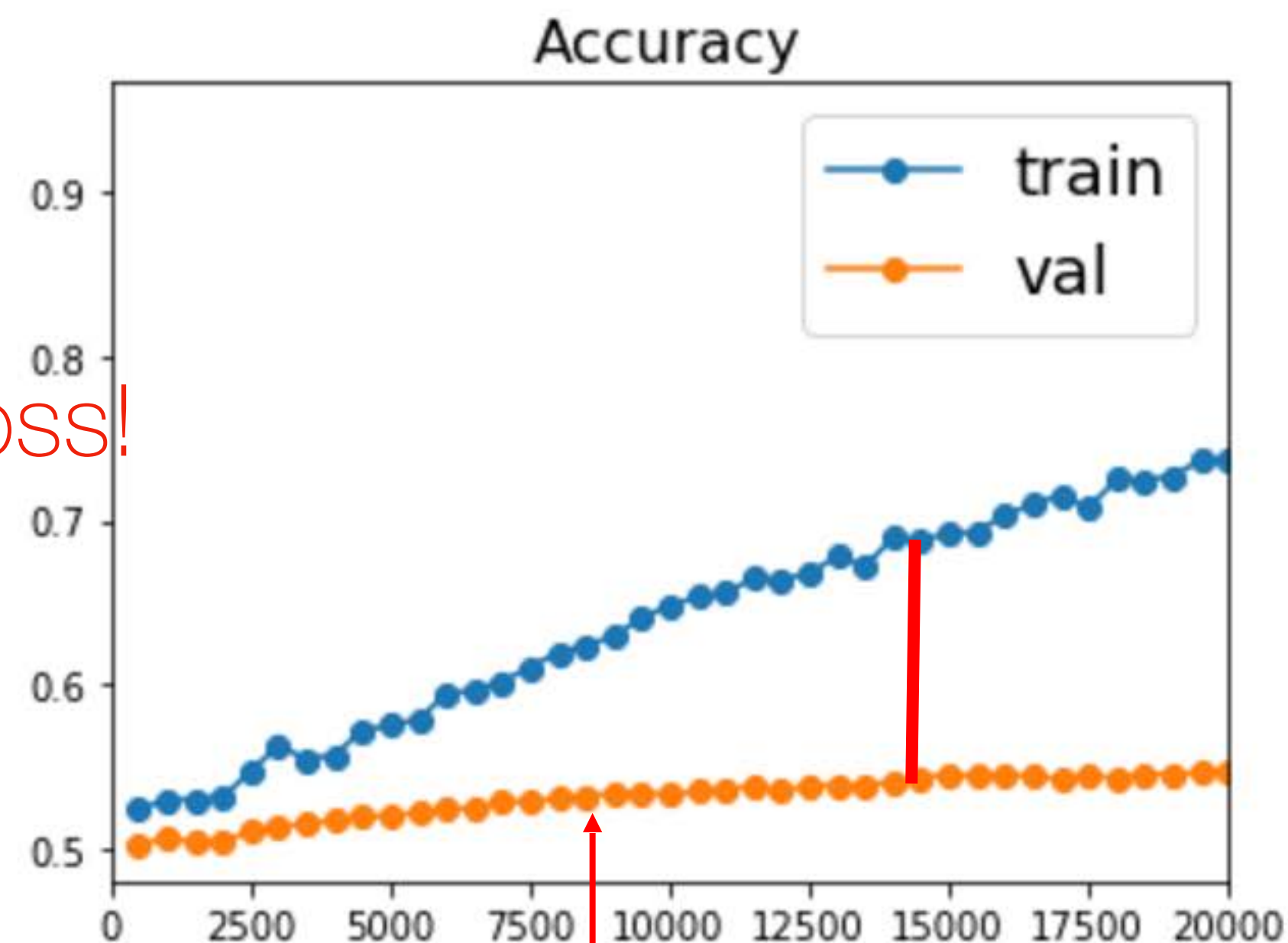
Beyond training loss

Recall the other problem



NOTE: No validation loss!

Better optimization algorithms help reduce training loss



But we really care about error on new data - how to reduce the gap?

A typical approach to overfitting

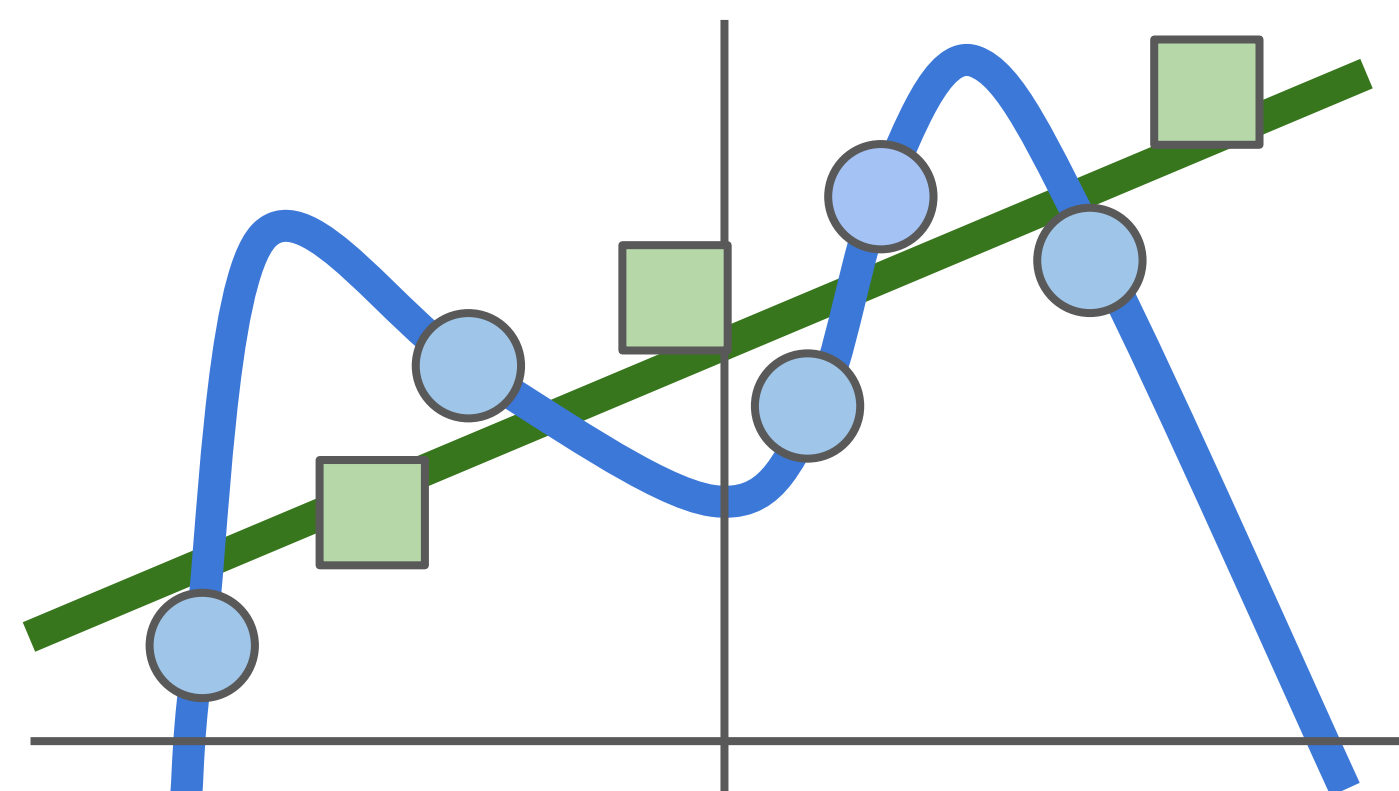
Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss}} + \underbrace{\lambda R(W)}_{\text{Regularization}}$$

Data loss: Model predictions should match training data

Regularization: Model should be “simple”, so it works on test data

$$\|W\|_2^2$$



Occam's Razor:

“Among competing hypotheses, the simplest is the best”

William of Ockham, 1285 - 1347

Common regularizers

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$$

In common use:

L2 regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2 \quad (\text{Weight decay})$$

L1 regularization

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

Elastic net (L1 + L2)

$$R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

Common regularizers

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$$

In common use:

L2 regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2 \quad (\text{Weight decay})$$

L1 regularization

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

Elastic net (L1 + L2)

$$R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

Common regularizers

My personal warning against L2

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$$

In common use:

L2 regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2 \quad (\text{Weight decay})$$

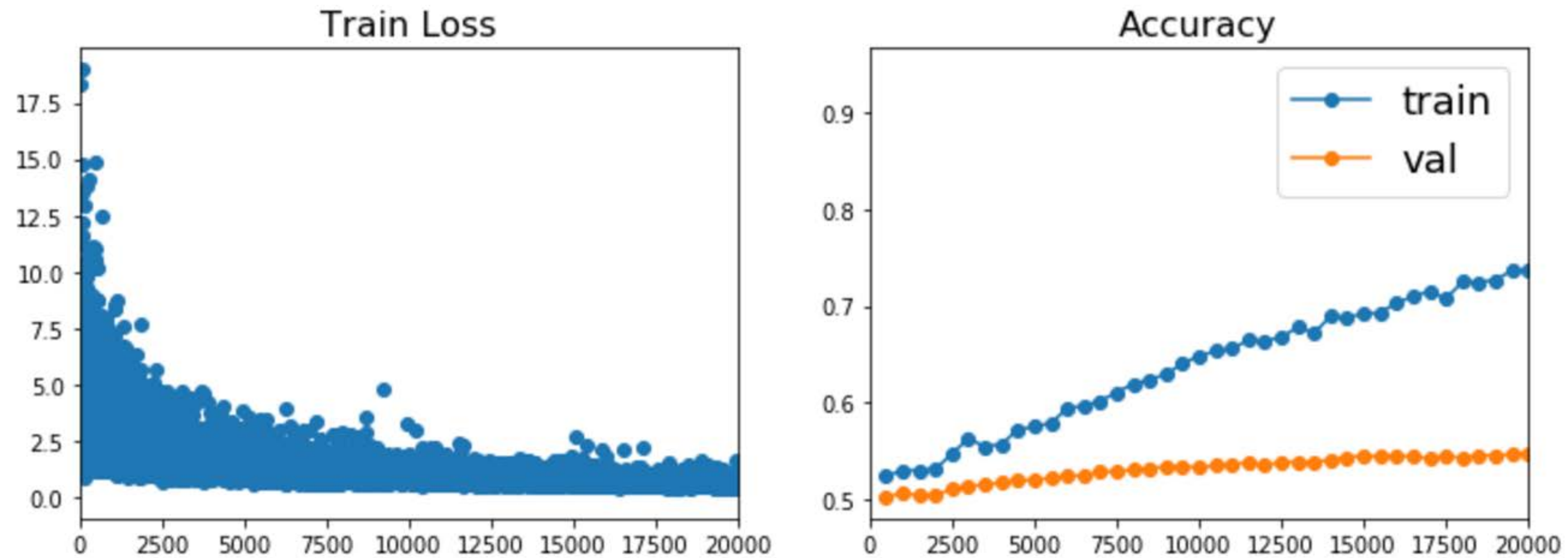
L1 regularization

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

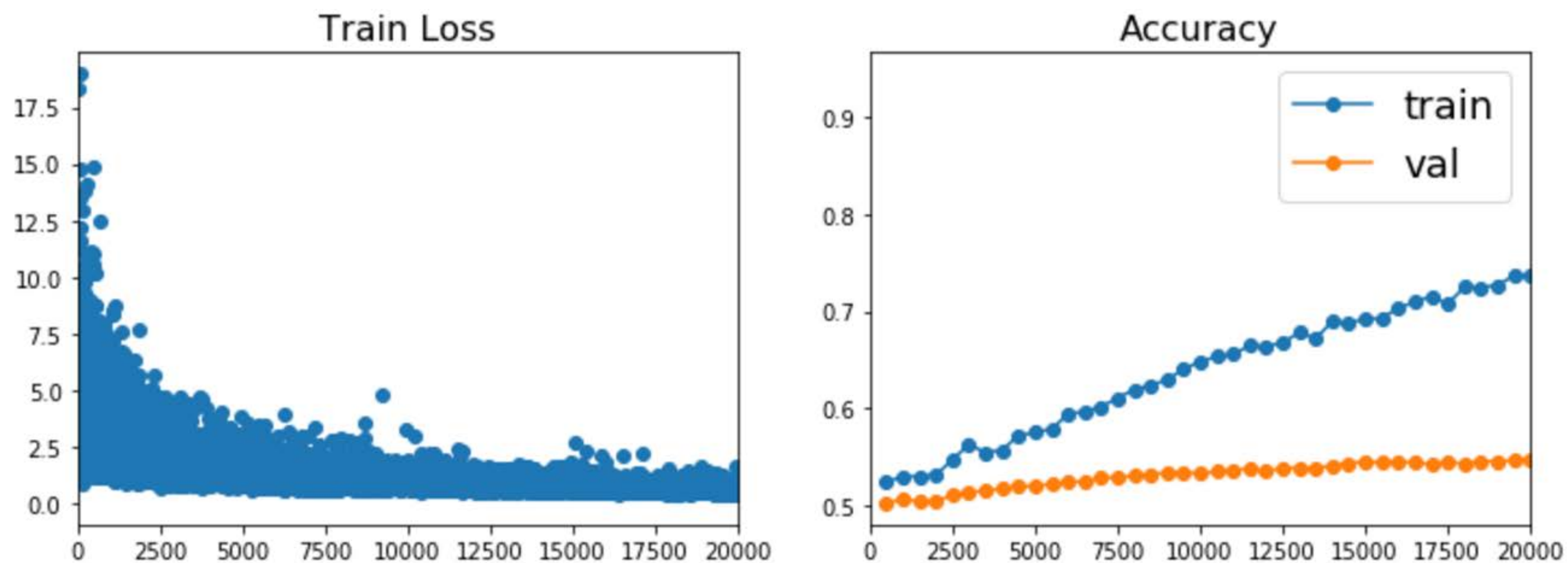
Ela

[Laarhoven, 2017](#), “However, we show that L2 regularization has **no regularizing effect when combined with normalization**. Instead, regularization has an influence on the scale of weights, and thereby on the effective learning rate.”

Why does this happen in the first place?



Why does this happen in the first place?

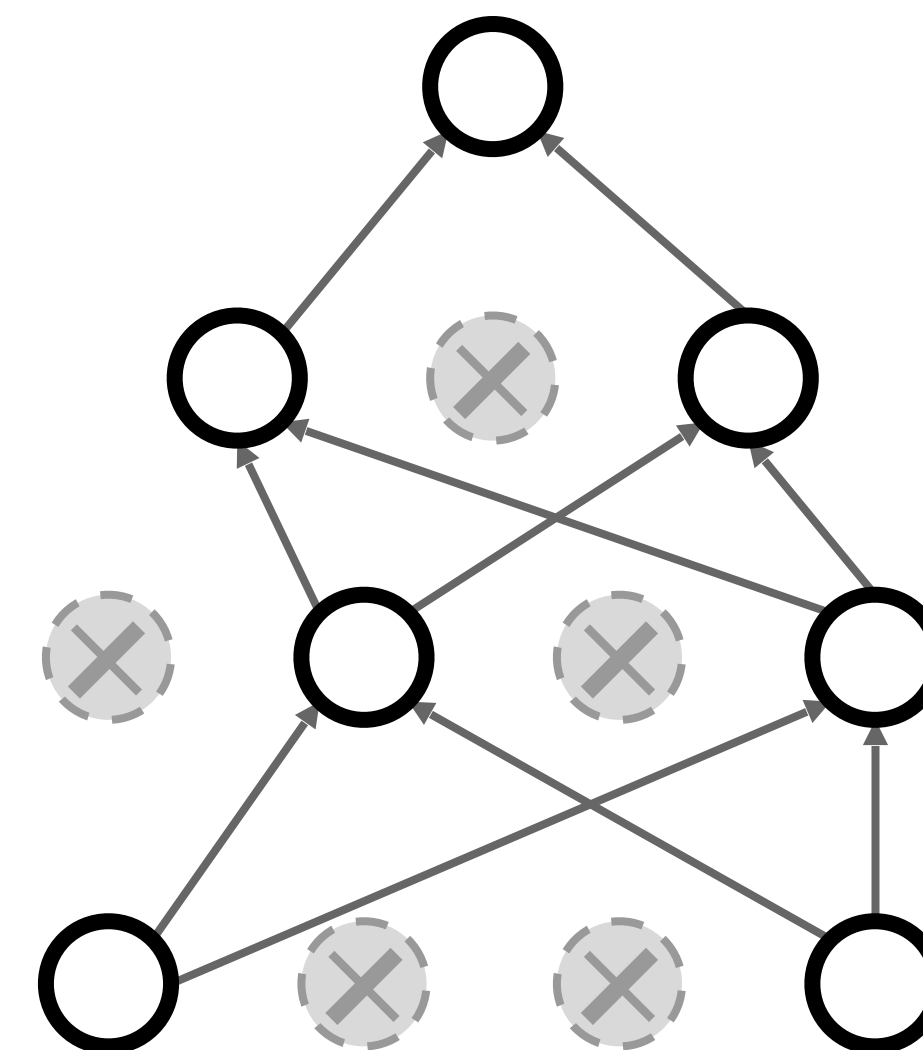
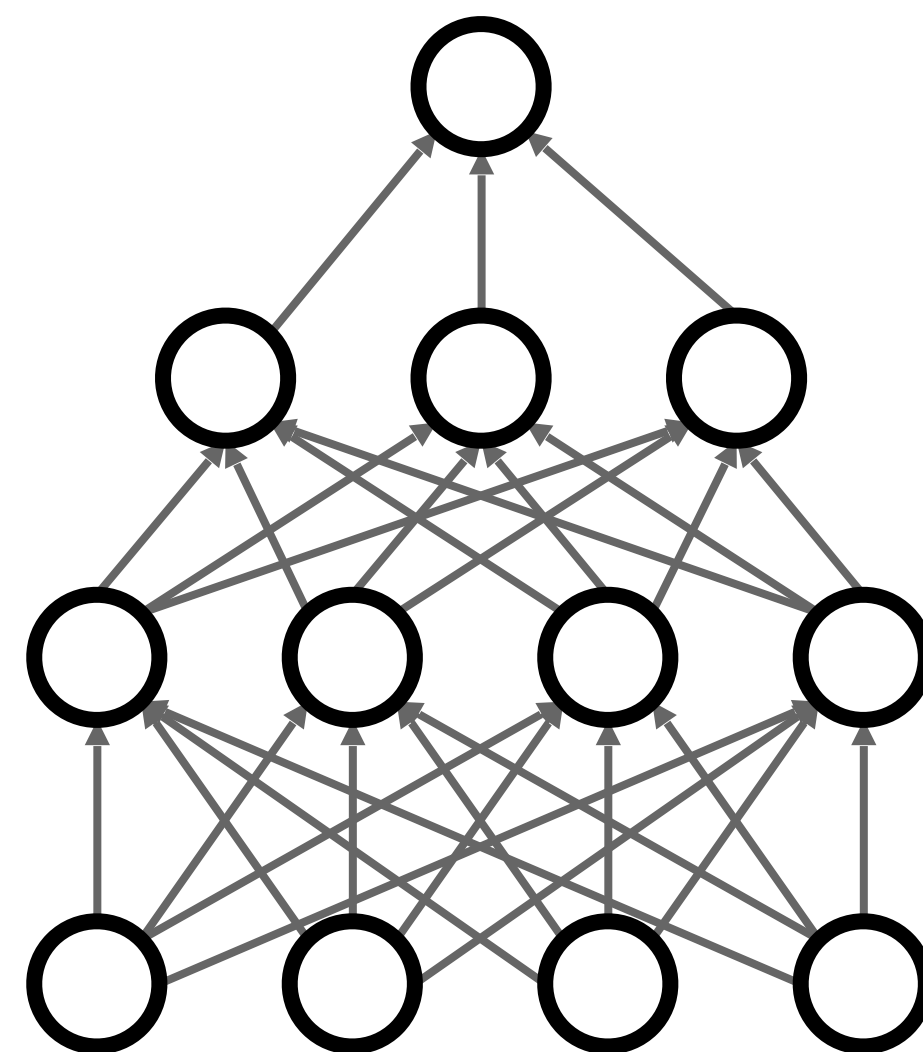


Can we somehow encode
uncertainty in data?

Regularization: Dropout

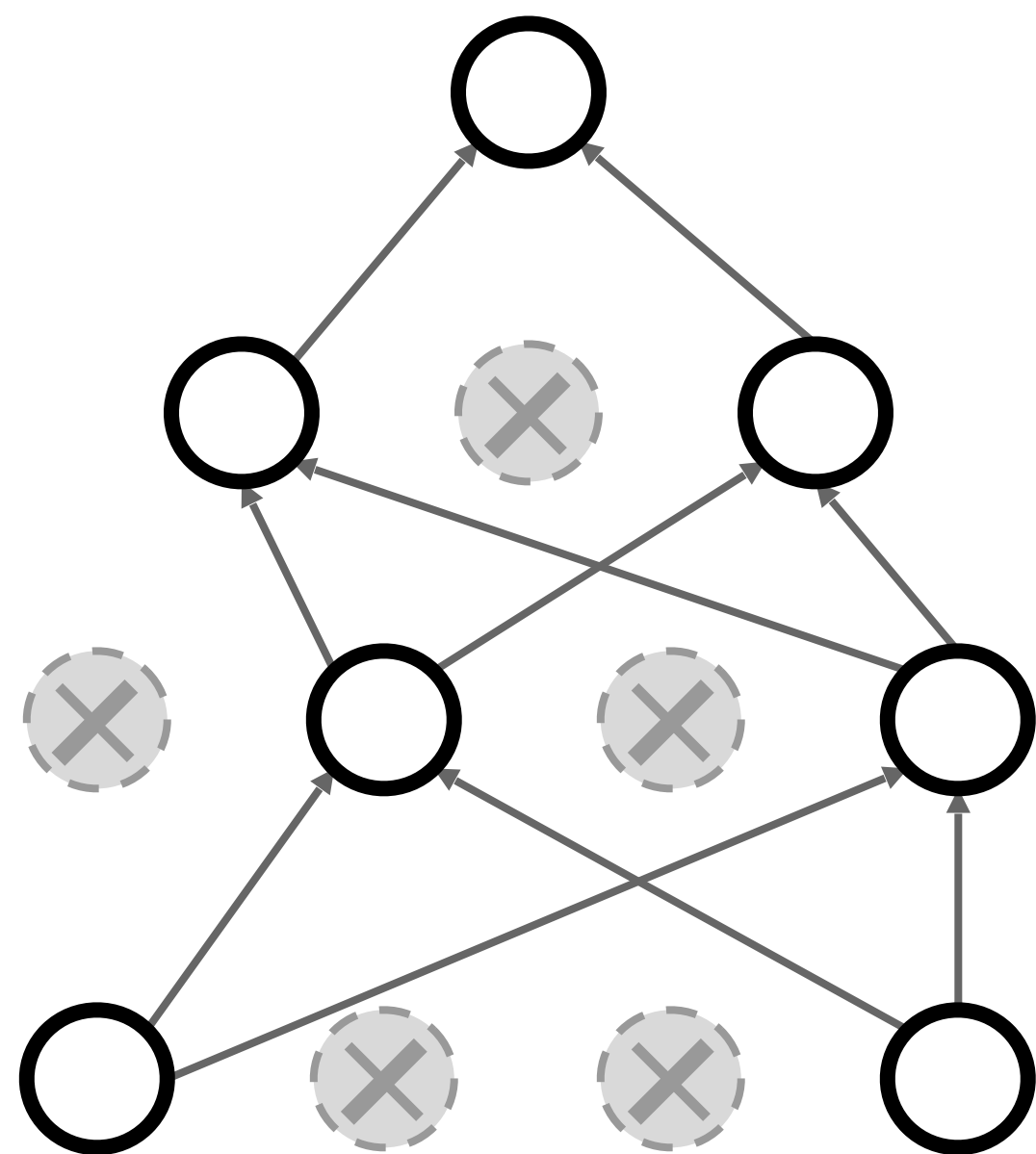
Making it impossible to trust the data 100%

In each forward pass, randomly set some neurons to zero
Probability of dropping is a hyperparameter; 0.5 is common



Regularization: Dropout

Making it impossible to trust the data 100%



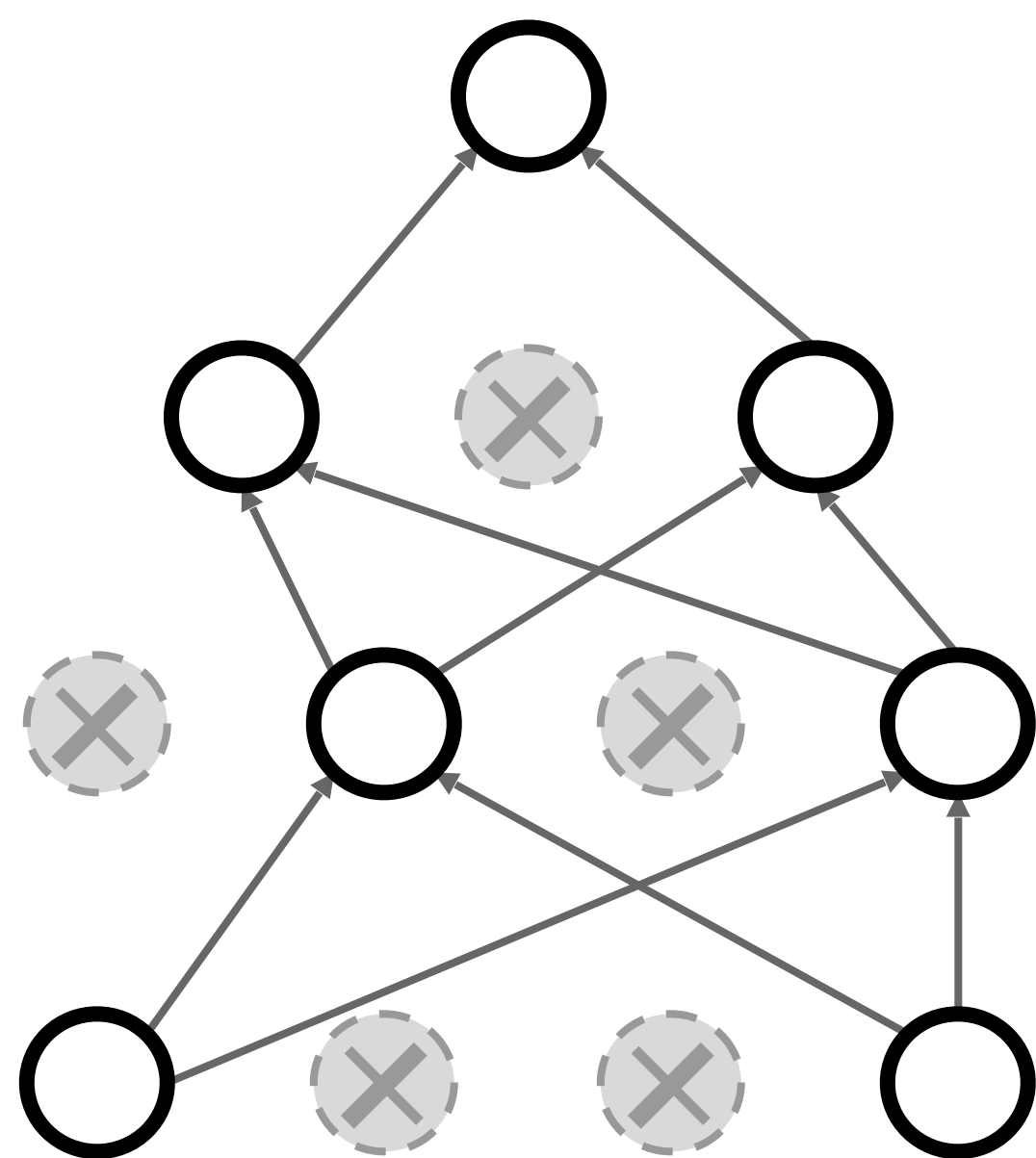
Forces the network to have a redundant representation;
Prevents co-adaptation of features



Skipped in class
(outside of scope)

Regularization: Dropout

Making it impossible to trust the data 100%



Another interpretation:

Dropout is training a large **ensemble** of models (that share parameters).

Each binary mask is one model

An FC layer with 4096 units has $2^{4096} \sim 10^{1233}$ possible masks!

Only $\sim 10^{82}$ atoms in the universe...

Skipped in class

Regularization: Dropout at **test** time (outside of scope)

Again the train / test gap

Dropout makes our output random!

Output (label) Input (image)

$$\boxed{y} = f_W(\boxed{x}, \boxed{z})$$

Random mask

Want to “average out” the randomness at test-time

$$y = f(x) = E_z[f(x, z)] = \int p(z) f(x, z) dz$$

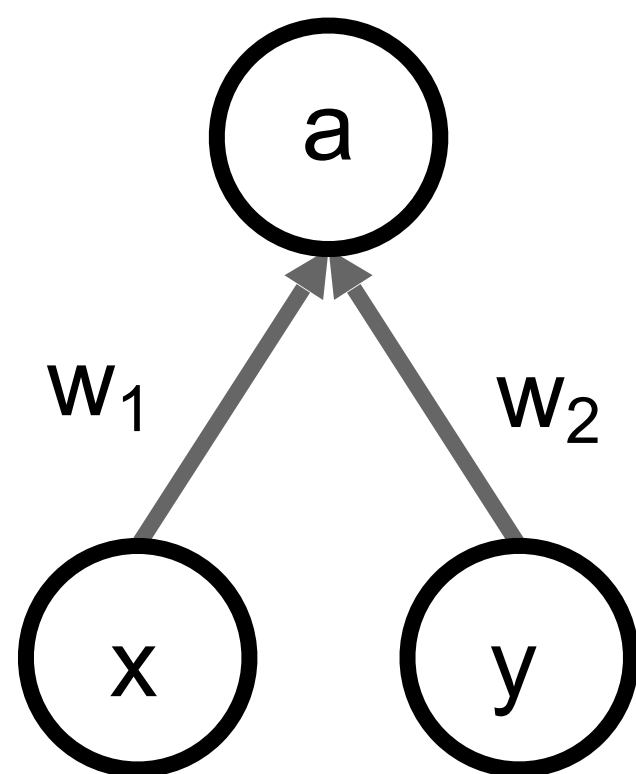
But this integral seems hard ...

Regularization: Dropout at **test** time (outside of scope)

An approximate solution

Want to approximate the integral $y = f(x) = E_z [f(x, z)] = \int p(z) f(x, z) dz$

Consider a single neuron.



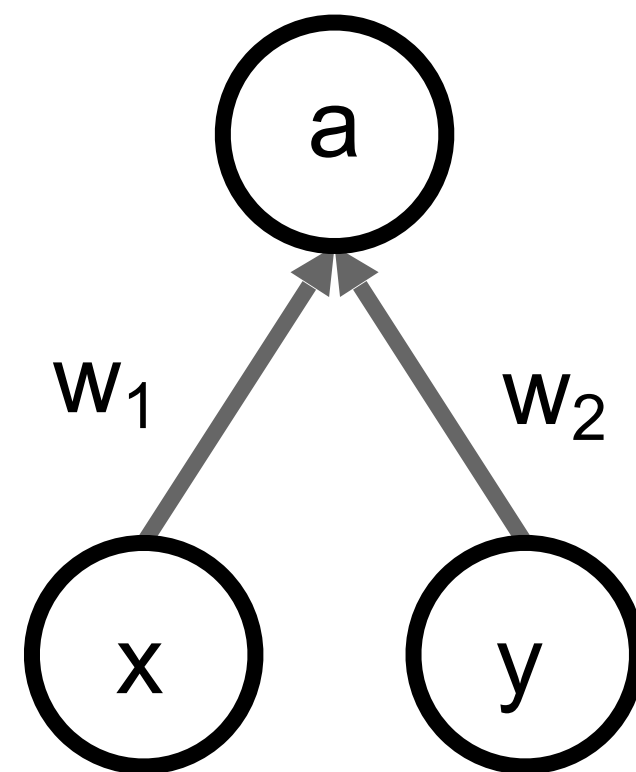
Regularization: Dropout at **test** time (outside of scope)

An approximate solution

Want to approximate the integral $y = f(x) = E_z [f(x, z)] = \int p(z) f(x, z) dz$

Consider a single neuron.

At test time we have: $E[a] = w_1x + w_2y$

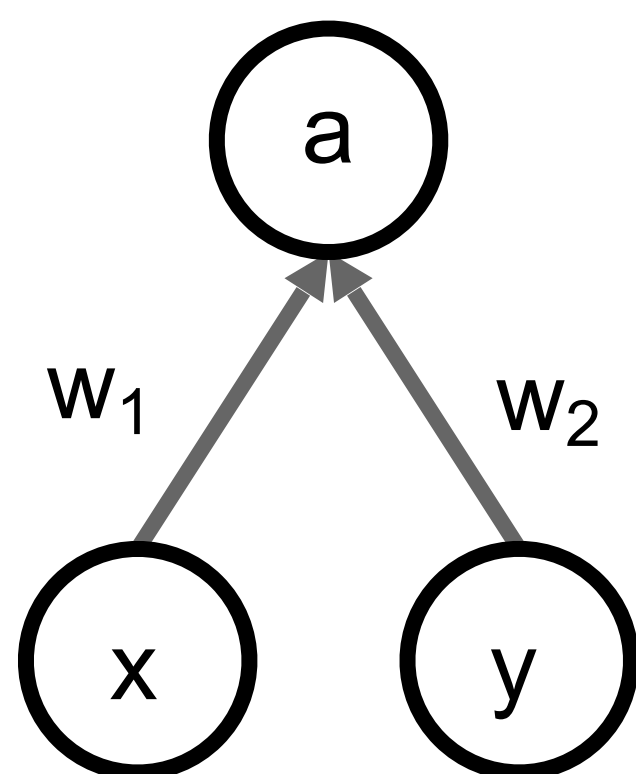


Regularization: Dropout at **test** time (outside of scope)

An approximate solution

Want to approximate the integral $y = f(x) = E_z [f(x, z)] = \int p(z) f(x, z) dz$

Consider a single neuron.



At test time we have: $E[a] = w_1x + w_2y$

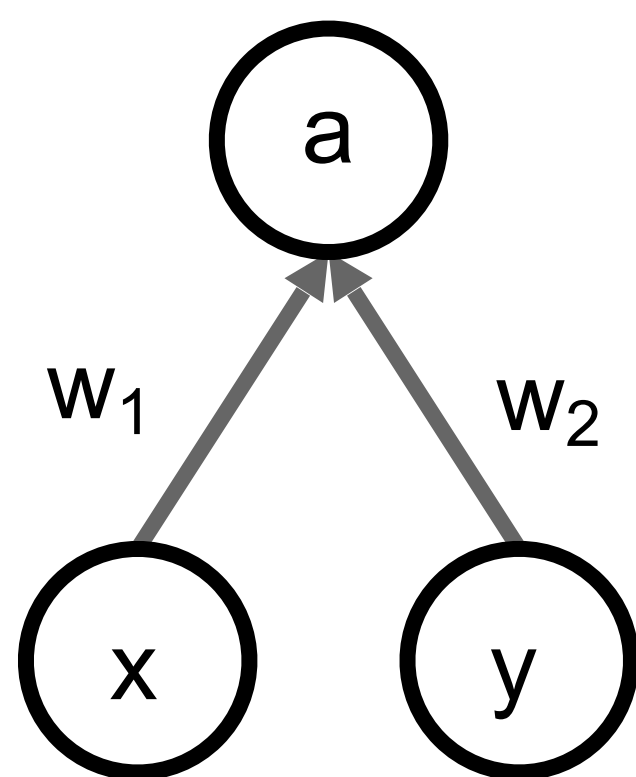
During training we have:
$$E[a] = \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) = \frac{1}{2}(w_1x + w_2y)$$

Regularization: Dropout at **test** time (outside of scope)

An approximate solution

Want to approximate the integral $y = f(x) = E_z [f(x, z)] = \int p(z) f(x, z) dz$

Consider a single neuron.



At test time we have: $E[a] = w_1x + w_2y$

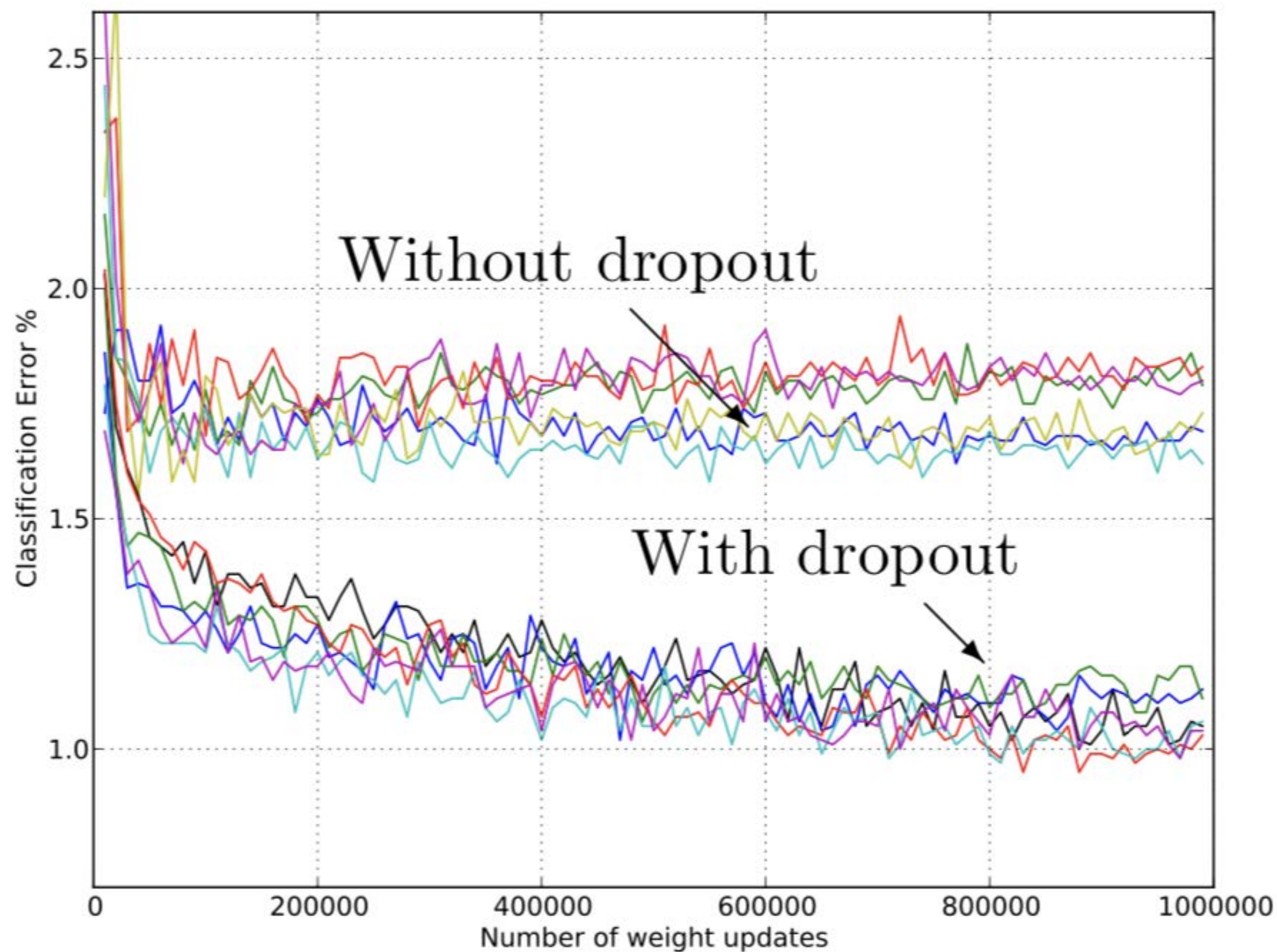
During training we have: $E[a] = \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) = \frac{1}{2}(w_1x + w_2y)$

At test time, **multiply** by dropout probability

Skipped in class
(outside of scope)

Regularization: Dropout

How good is it?



Regularization: A common pattern (outside of scope)

Training: Add some kind of randomness

$$y = f_W(x, z)$$

Testing: Average out randomness
(sometimes approximate)

$$y = f(x) = E_z [f(x, z)] = \int p(z) f(x, z) dz$$

Skipped in class

Regularization: A common pattern (outside of scope)

Training: Add some kind of randomness

$$y = f_W(x, z)$$

Testing: Average out randomness
(sometimes approximate)

$$y = f(x) = E_z [f(x, z)] = \int p(z) f(x, z) dz$$

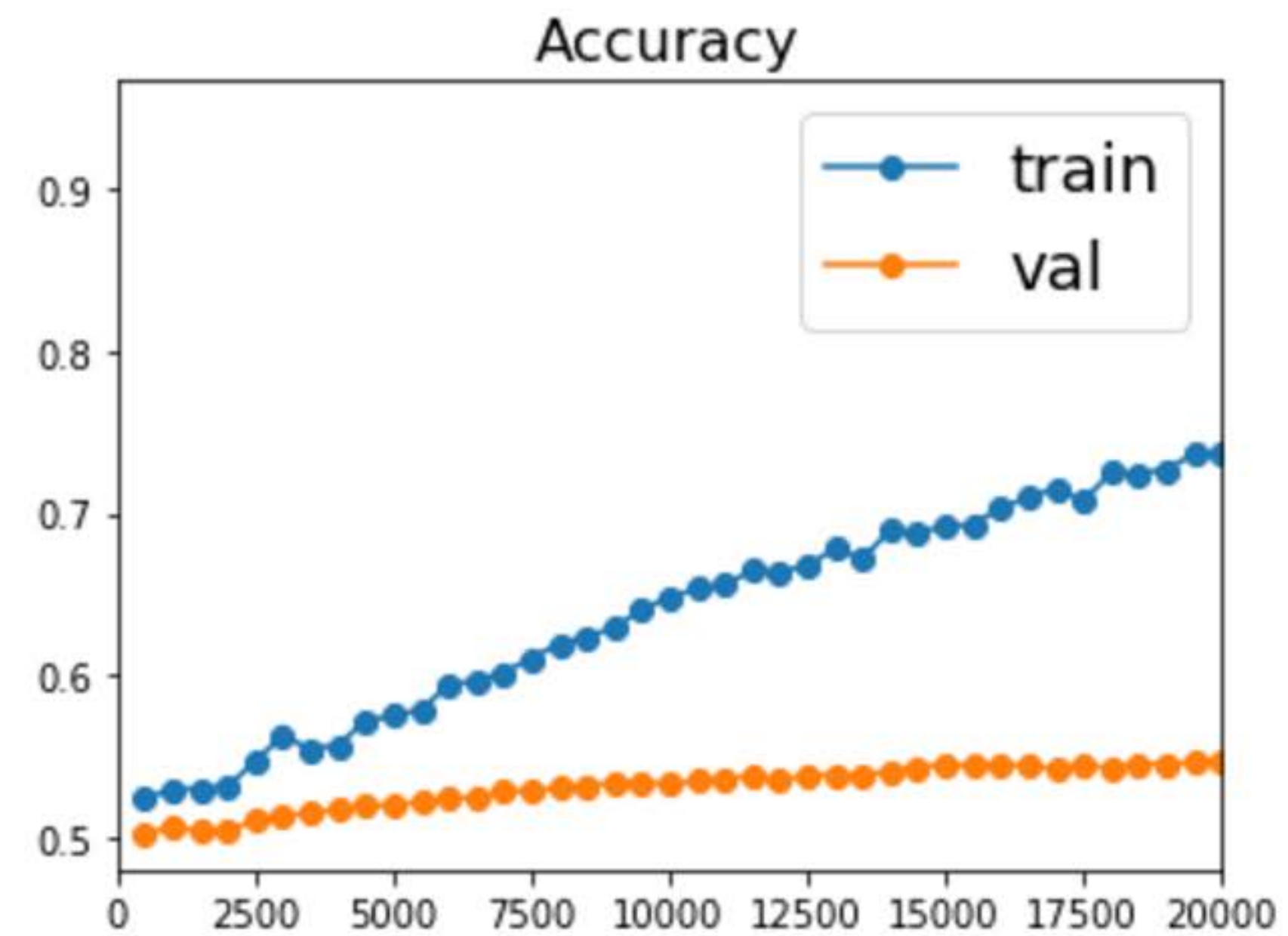
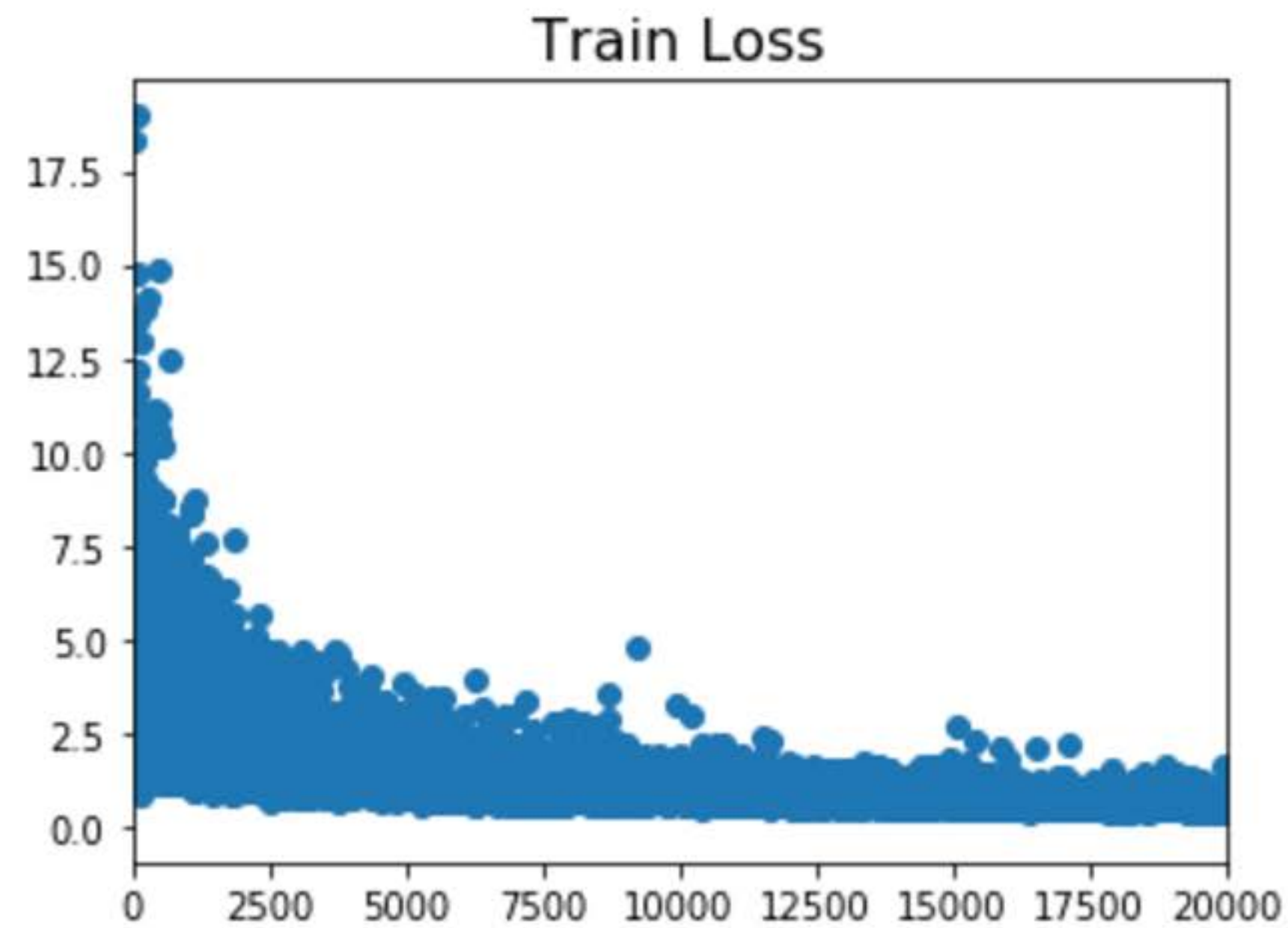
Example: Batch Normalization

Training: Normalize using stats from random minibatches

Testing: Use fixed stats to normalize

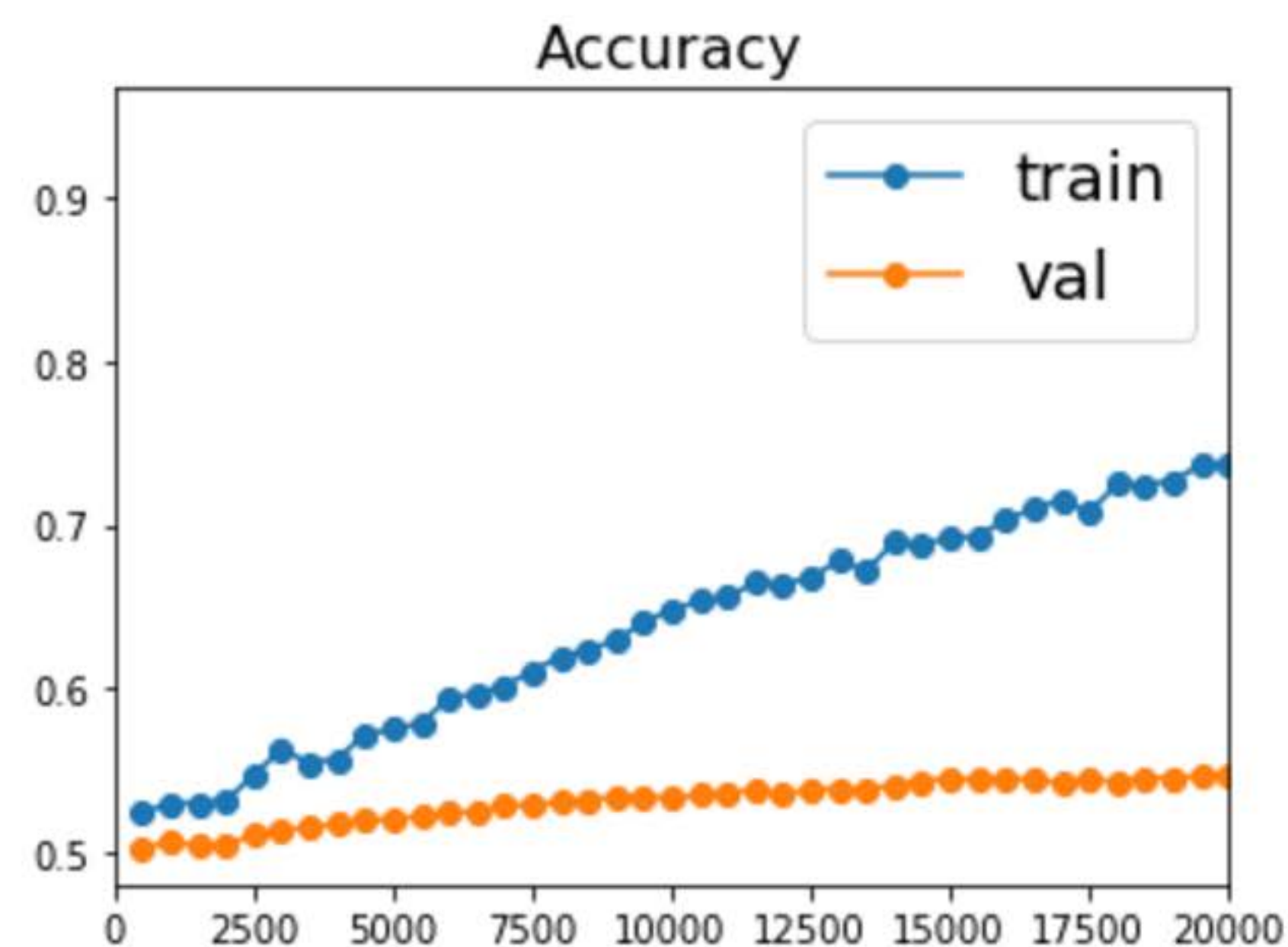
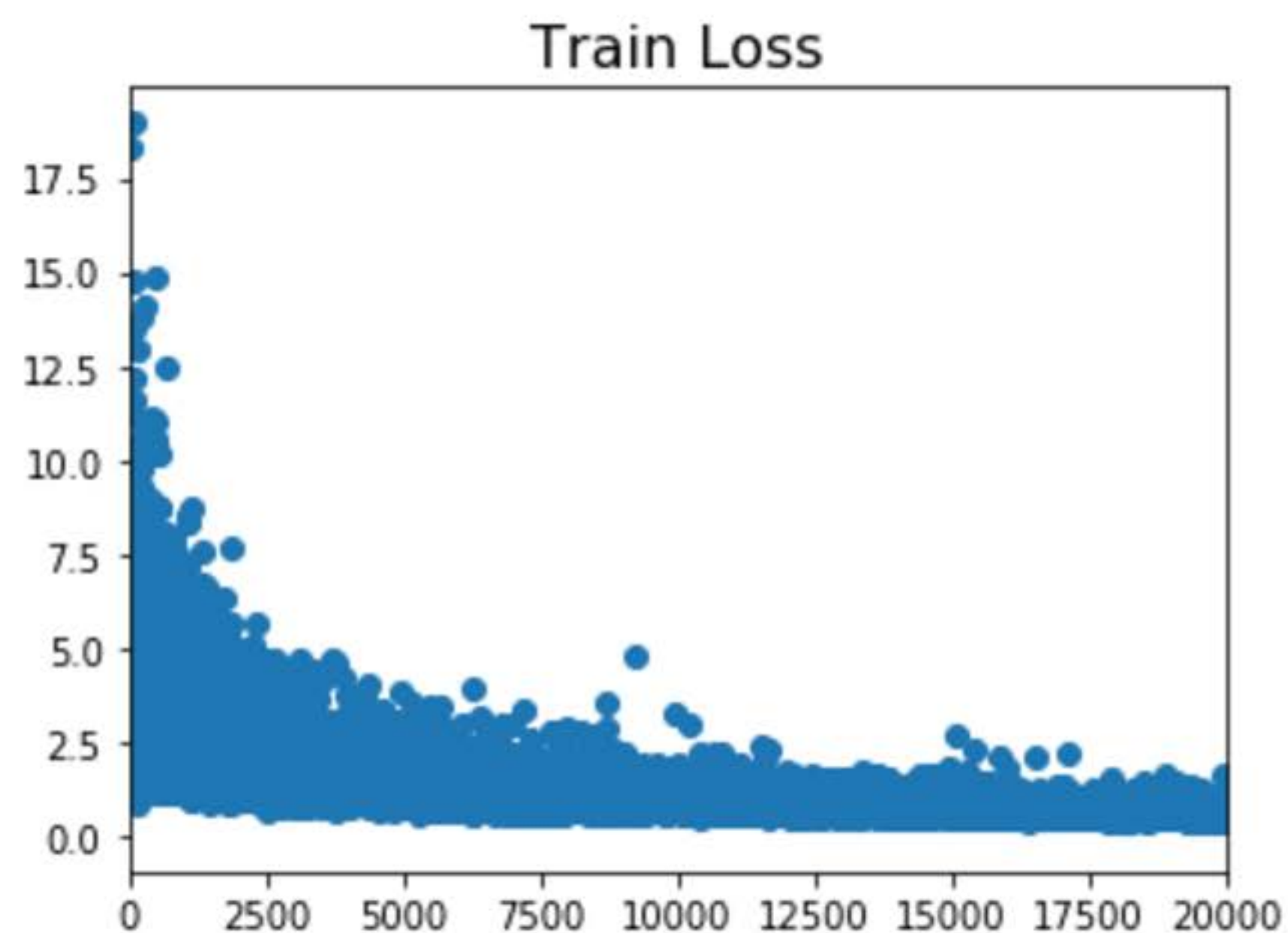
Skipped in class

Why does this happen in the first place? (outside of scope)



Skipped in class

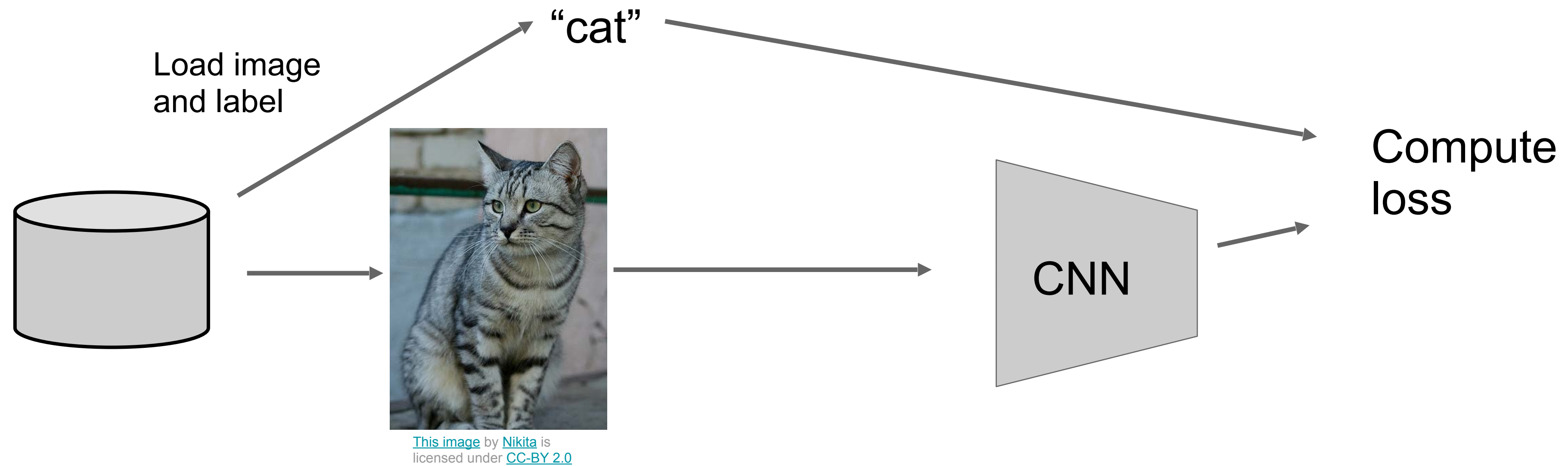
Why does this happen in the first place? (outside of scope)



How can we have more data?

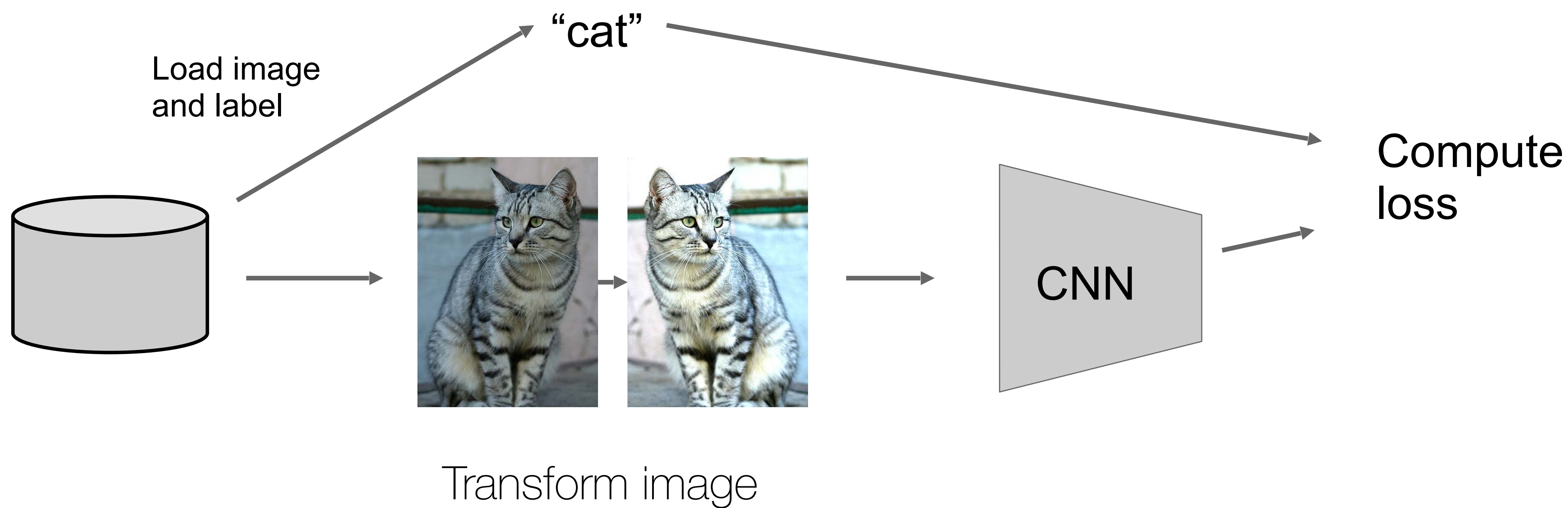
Skipped in class

Regularization: Data augmentation (outside of scope)

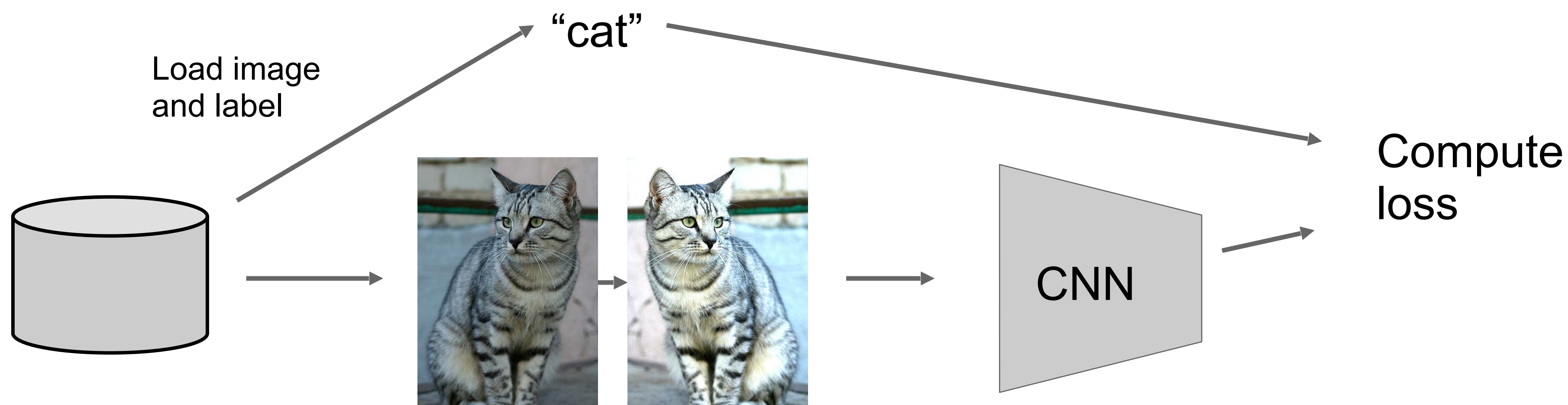


Skipped in class

Regularization: Data augmentation (outside of scope)



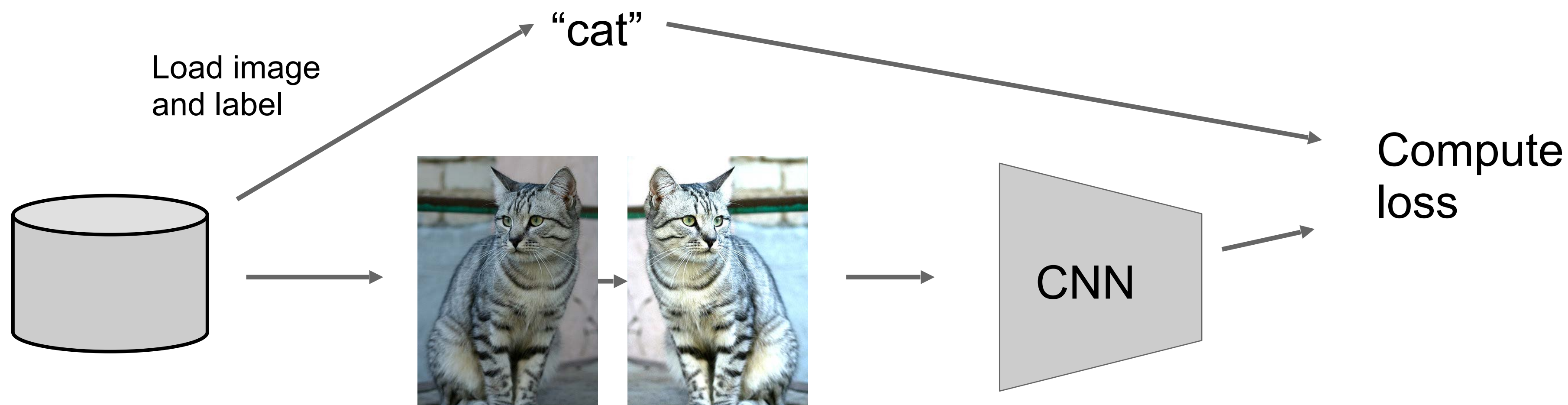
Regularization: Data augmentation (outside of scope)



Transform image

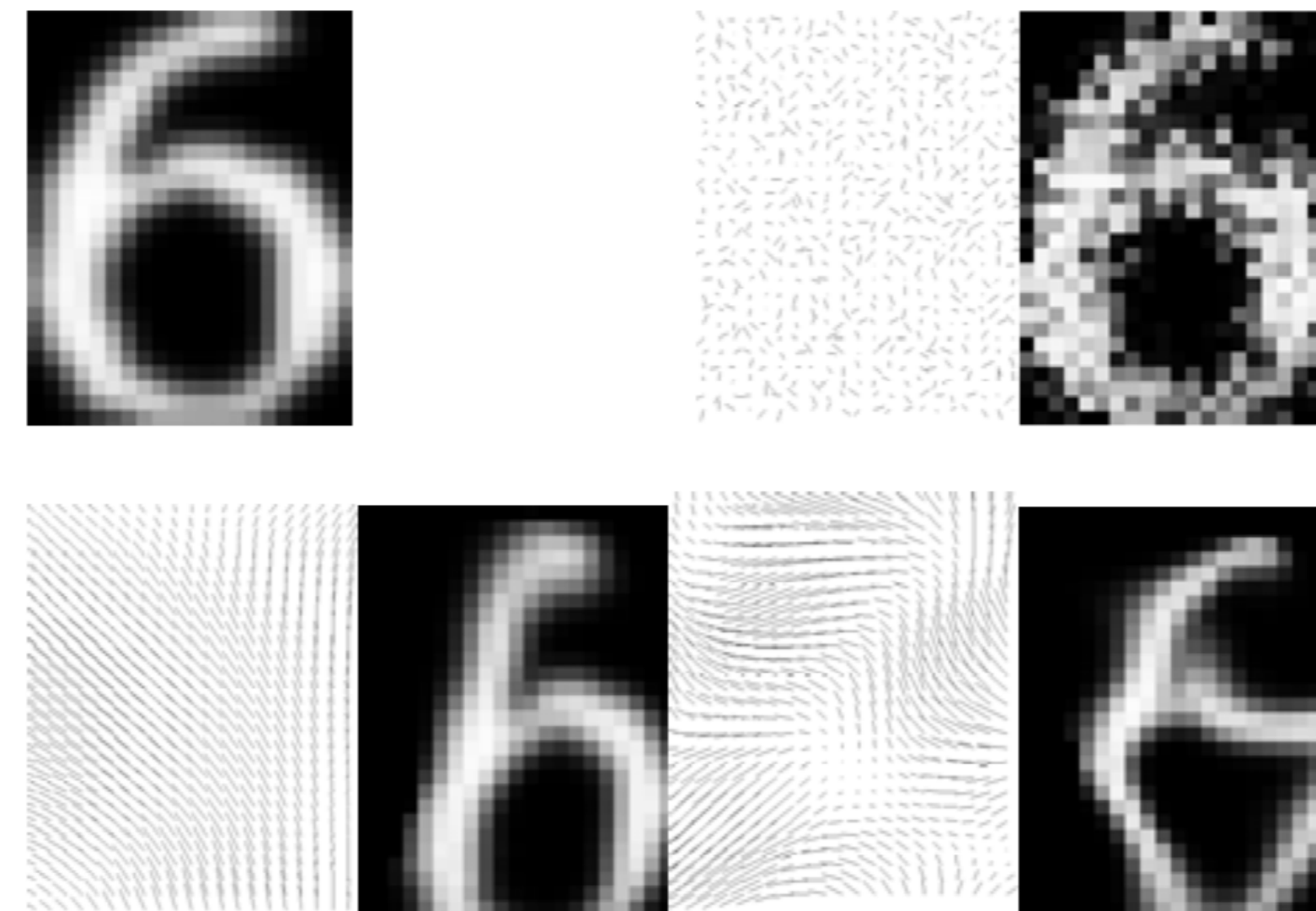
- Horizontal / vertical flips
- Color / brightness
- Rotations / scaling
- Elastic transformation

Regularization: Data augmentation (outside of scope)



Transform image

- Horizontal / vertical flips
- Color / brightness
- Rotations / scaling
- Elastic transformation

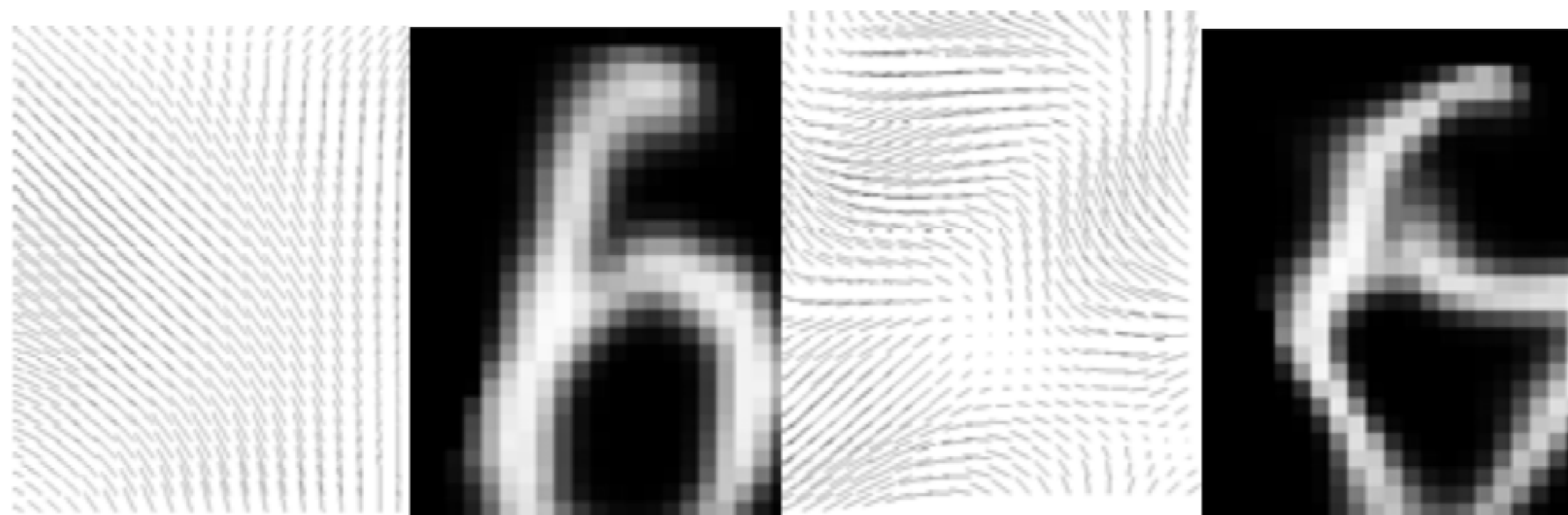


Regularization: Data augmentation (outside of scope)

Elastic deformations



1. Create random displacement field with uniform distribution



2. Smooth the displacement field with a Gaussian

Figures copyright IEEE, 2003. Reproduced for educational purposes.

Regularization: Data augmentation (outside of scope)

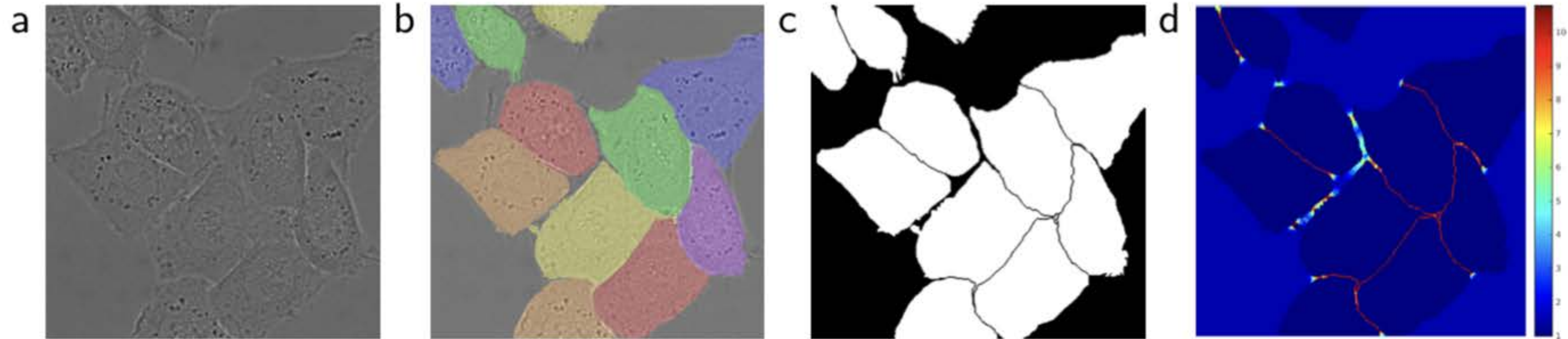
Elastic deformations

Algorithm	Distortion	Error	Ref.
2 layer MLP (MSE)	affine	1.6%	[3]
SVM	affine	1.4%	[9]
Tangent dist.	affine+thick	1.1%	[3]
Lenet5 (MSE)	affine	0.8%	[3]
Boost. Lenet4 MSE	affine	0.7%	[3]
Virtual SVM	affine	0.6%	[9]
2 layer MLP (CE)	none	1.6%	this paper
2 layer MLP (CE)	affine	1.1%	this paper
2 layer MLP (MSE)	elastic	0.9%	this paper
2 layer MLP (CE)	elastic	0.7%	this paper
Simple conv (CE)	affine	0.6%	this paper
Simple conv (CE)	elastic	0.4%	this paper

Table 1. Comparison between various algorithms.

Regularization: Data augmentation (outside of scope)

Elastic deformations

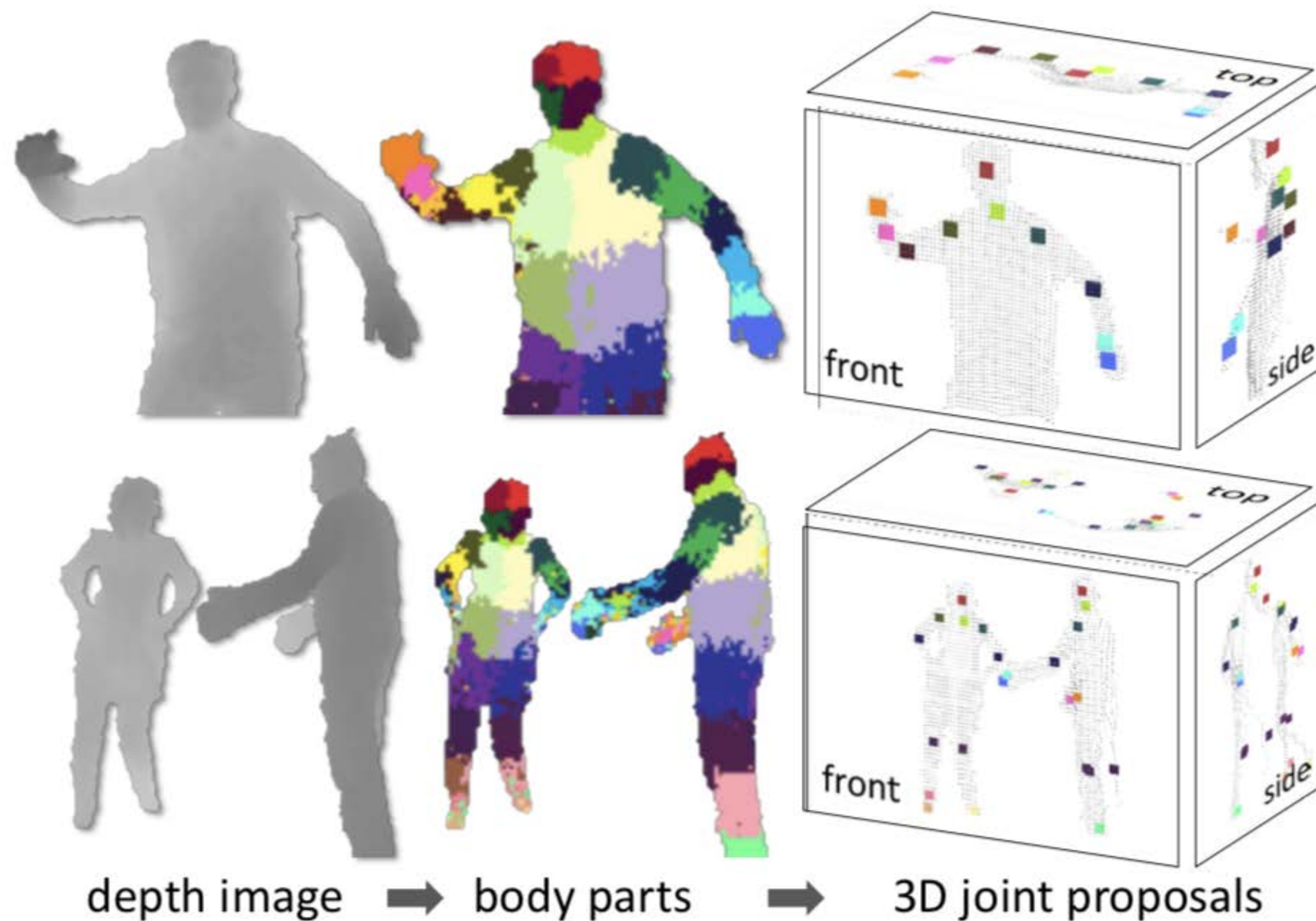


Name	PhC-U373	DIC-HeLa
IMCB-SG (2014)	0.2669	0.2935
KTH-SE (2014)	0.7953	0.4607
HOUS-US (2014)	0.5323	-
second-best 2015	0.83	0.46
u-net (2015)	0.9203	0.7756

Skipped in class

Regularization: Data augmentation (outside of scope)

Synthetic data



Skipped in class

Regularization: Data augmentation (outside of scope)

Synthetic data + generative models



Unlabeled Real Images

Simulated images



Synthetic

Refined

Unlabeled Real Images

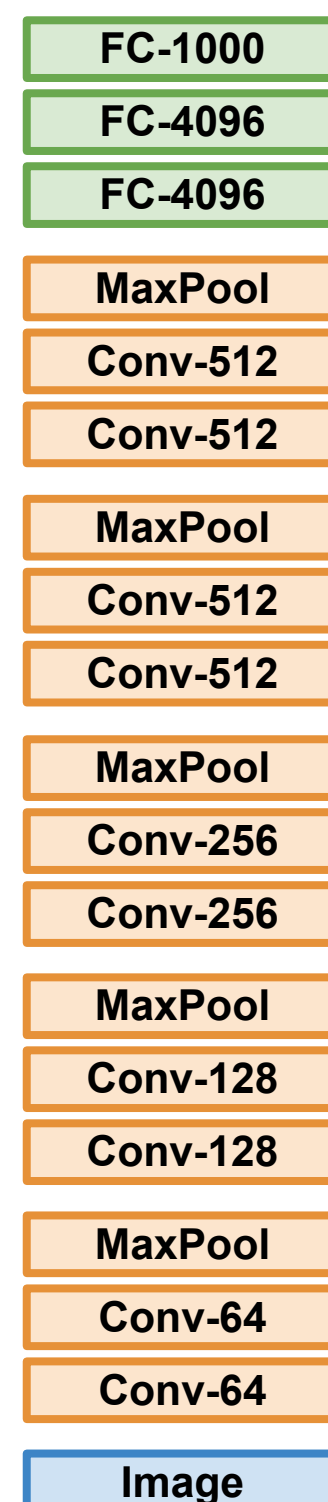
Simulated images

Skipped in class
(outside of scope)

Douvan et al. "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", IJCVML 2014
Razavian et al. "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

Using pretrained networks

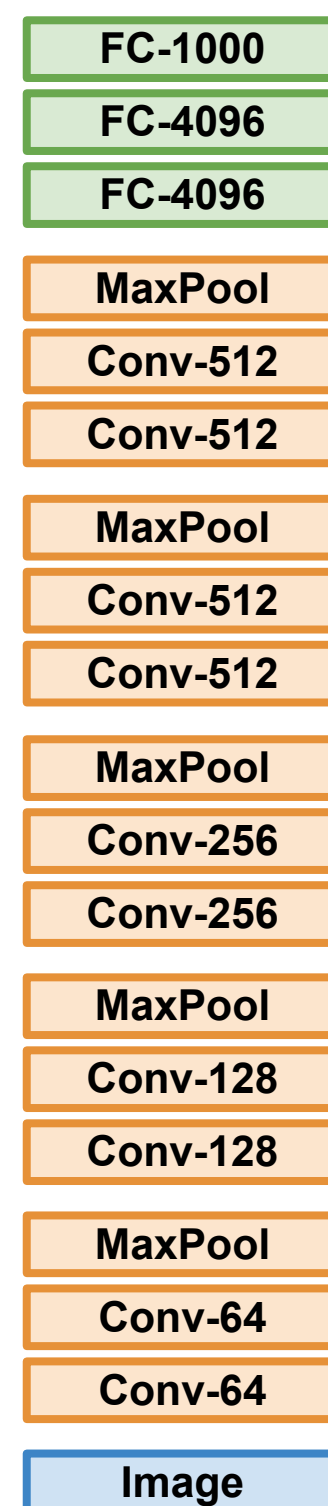
1. Train on Imagenet



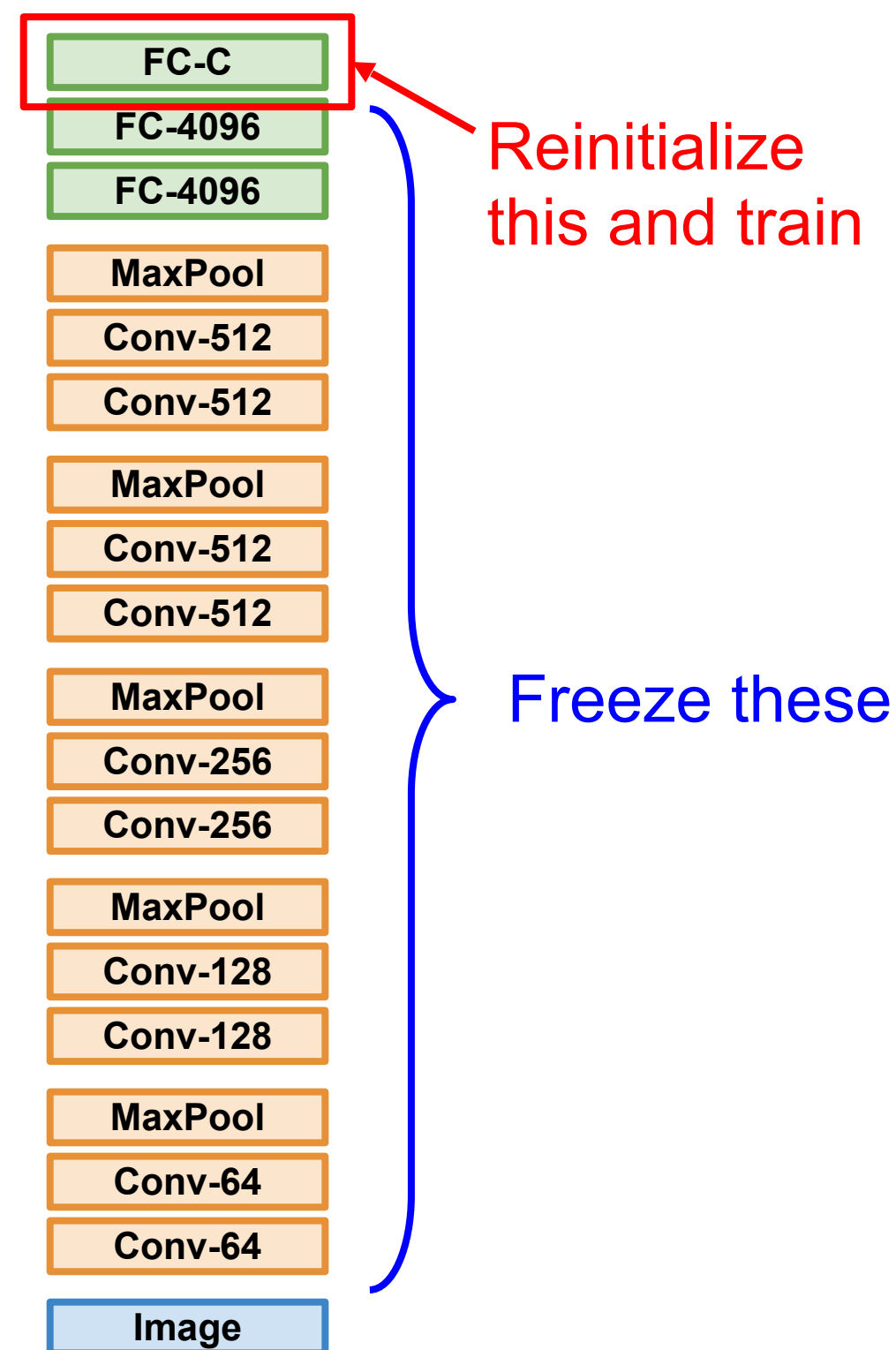
Skipped in class
(outside of scope)

Using pretrained networks

1. Train on Imagenet



2. Small Dataset (C classes)



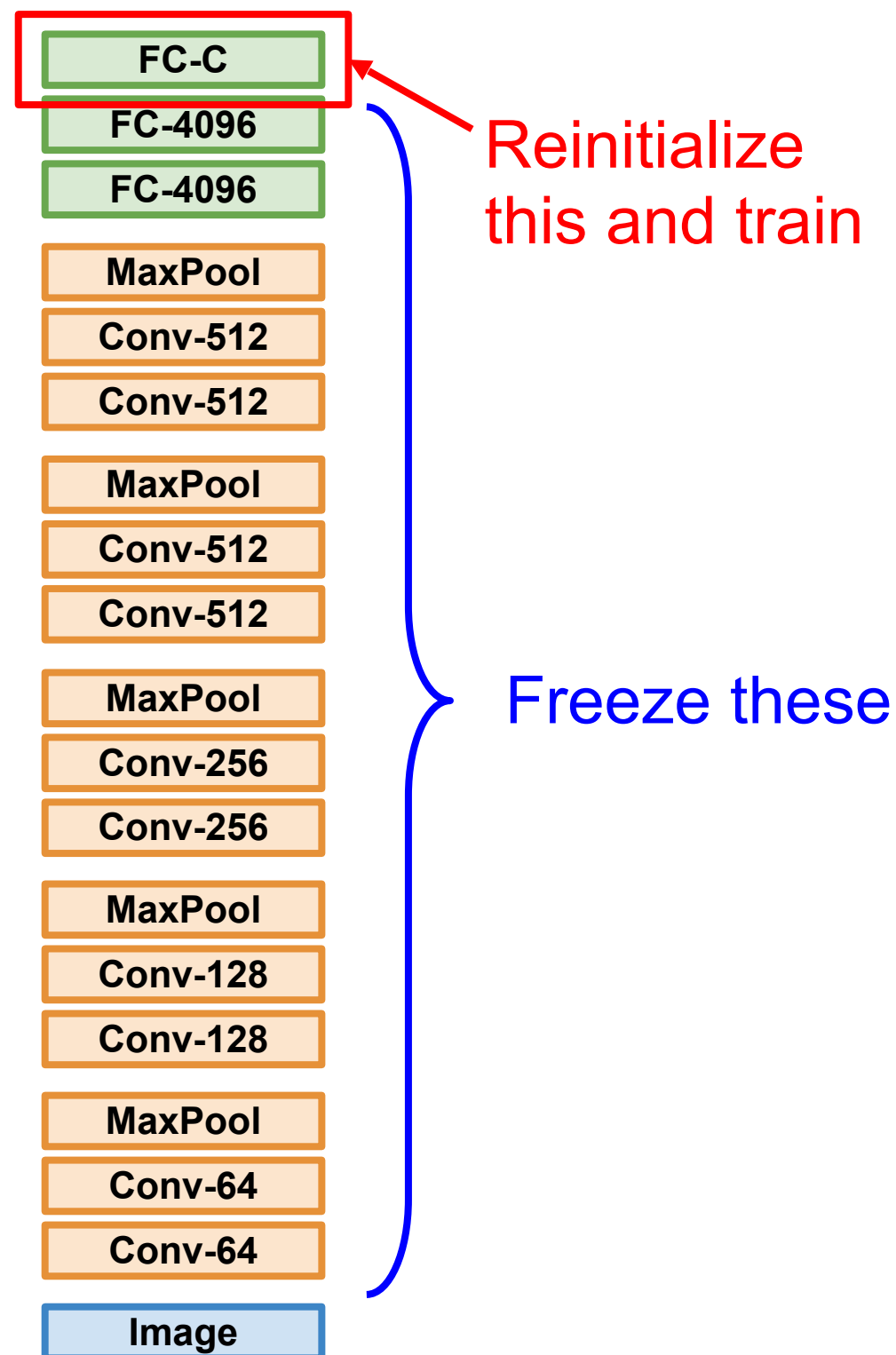
Skipped in class (outside of scope)

Using pretrained networks

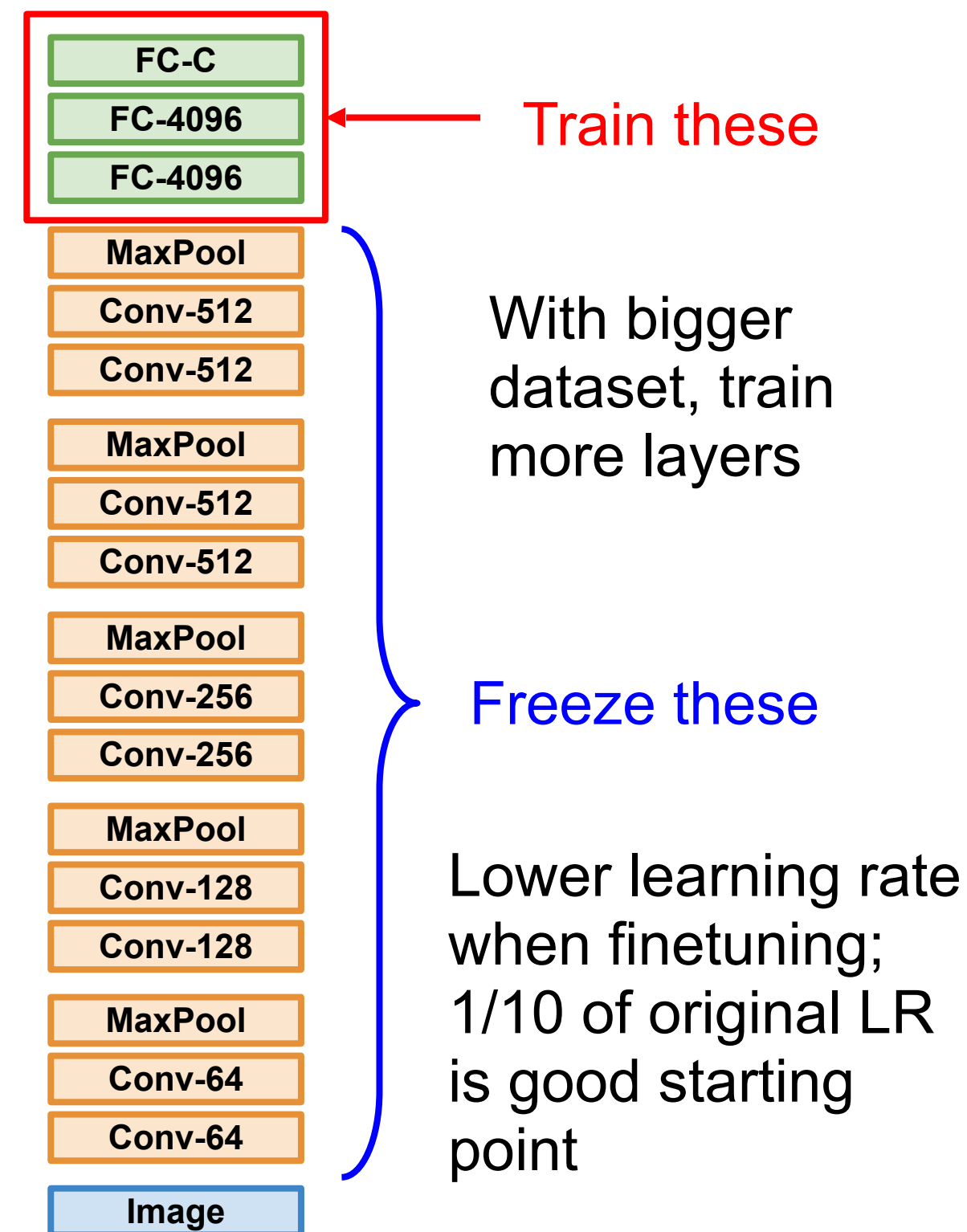
1. Train on Imagenet



2. Small Dataset (C classes)



3. Bigger dataset



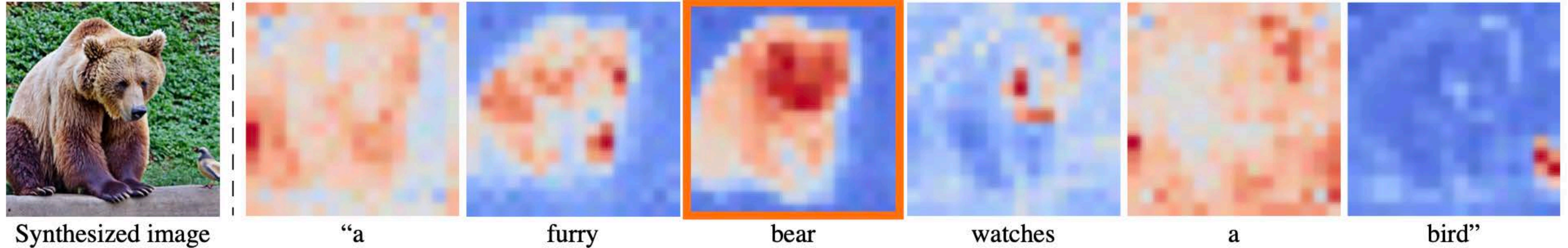


Large generative models (Skipped in class (outside of scope))



Skipped in class
(outside of scope)

Fishing information within SD



Average cross-attention maps across all timestamps

Cross-attention maps for individual timestamps

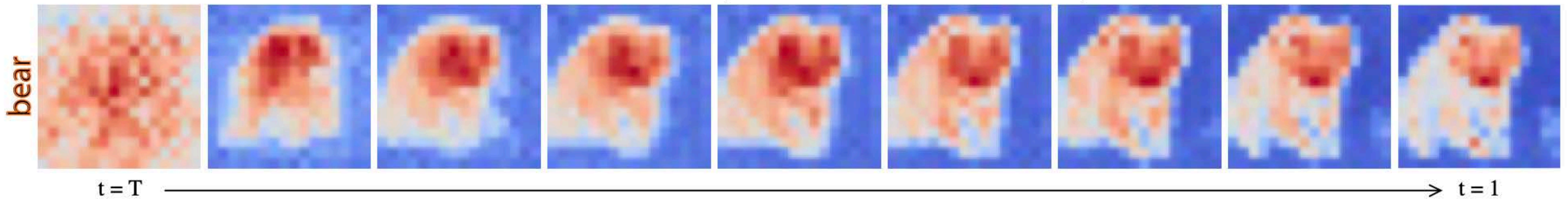
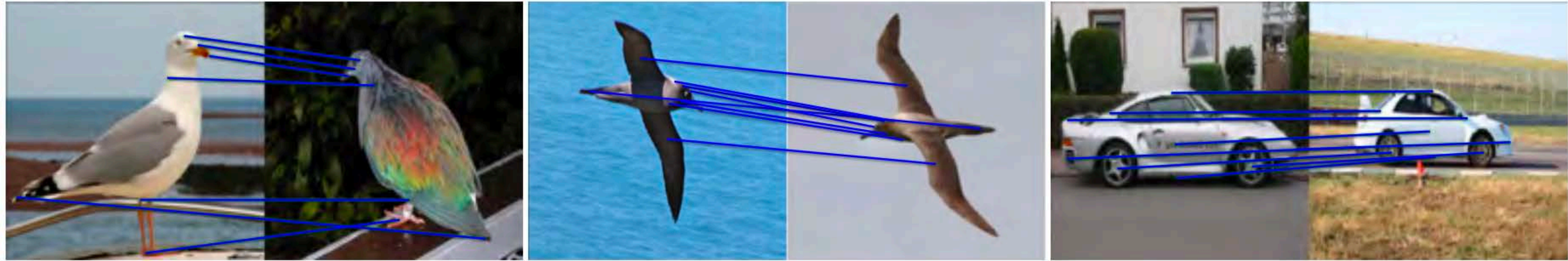


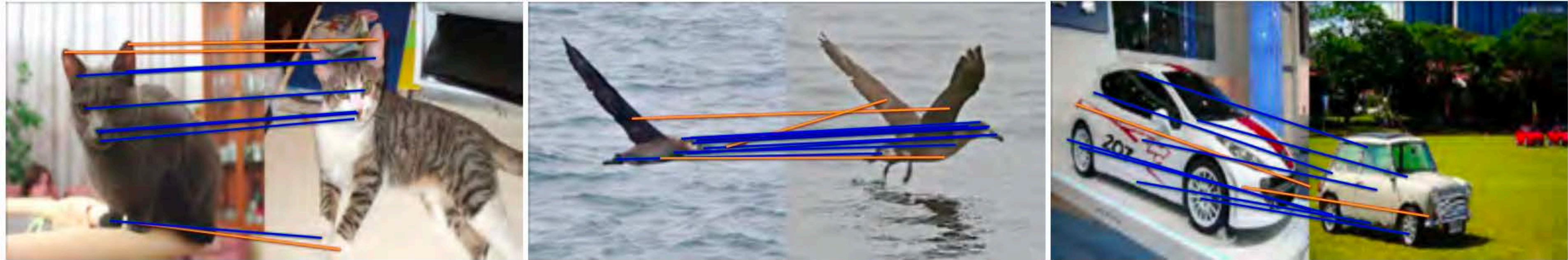
Image from [Hertz et al., ICLR, 2023]

Correspondences from SD (Skipped in class outside of scope)

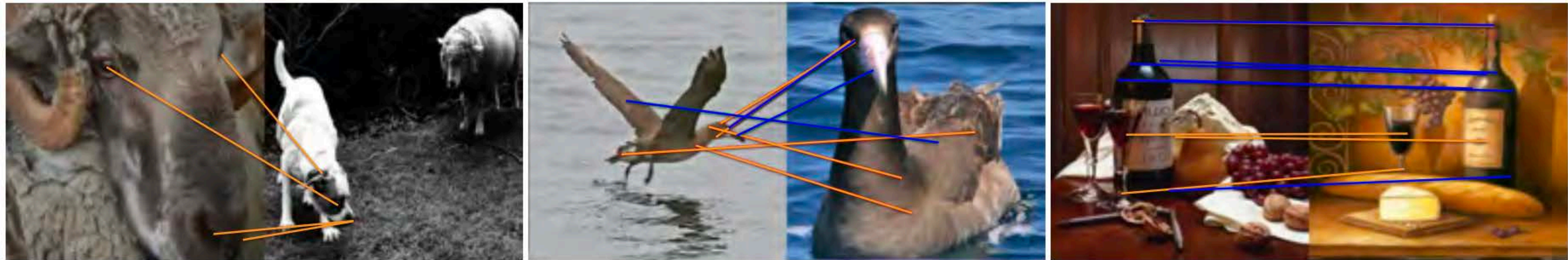
Successful



Mixed



Failure



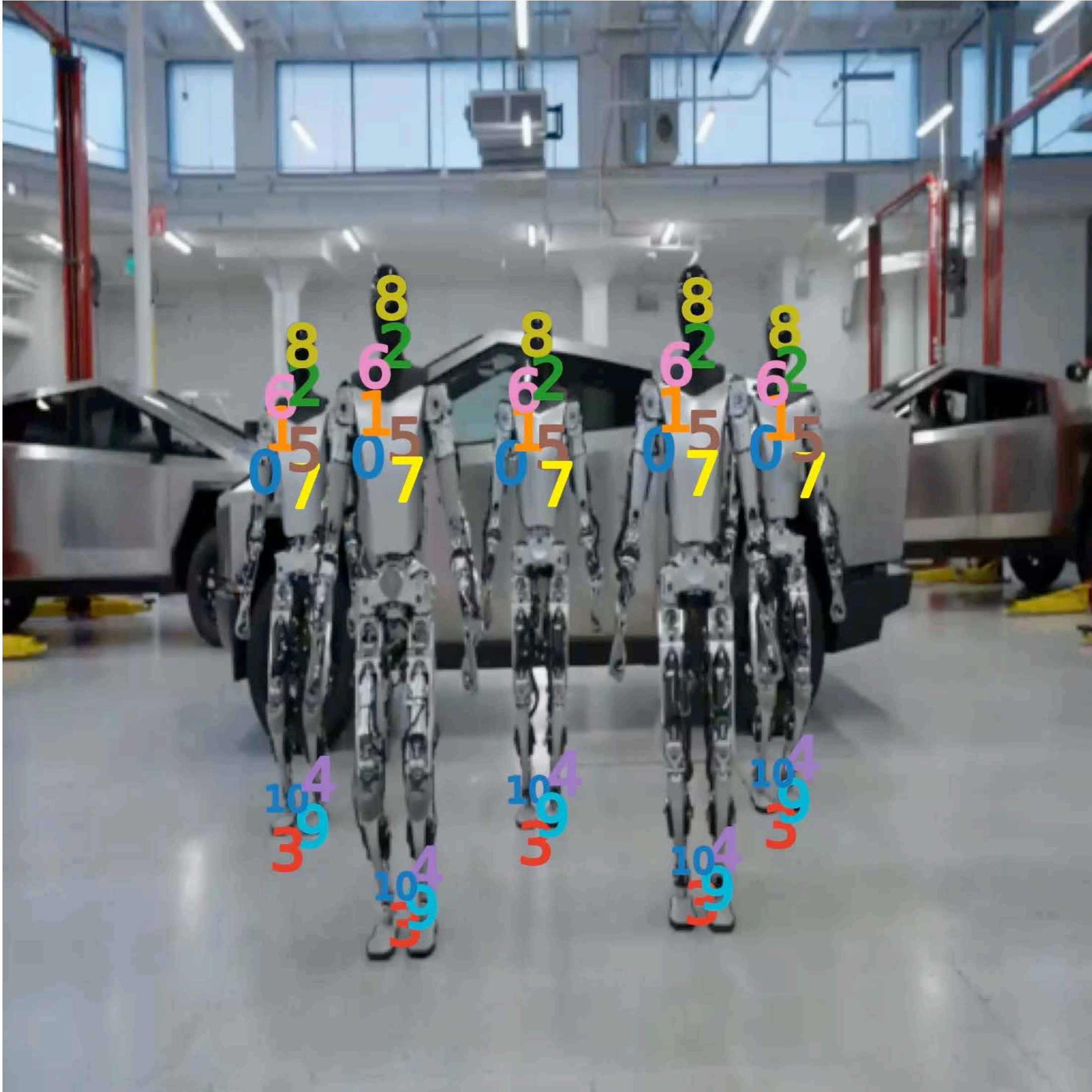
Spair-71k

CUB-200

PF-Willow



Keypoints from SD (Skipped in class outside of scope)



Text-to-3D from SD **Skipped in class** **(outside of scope)**

Input

Multi-view images



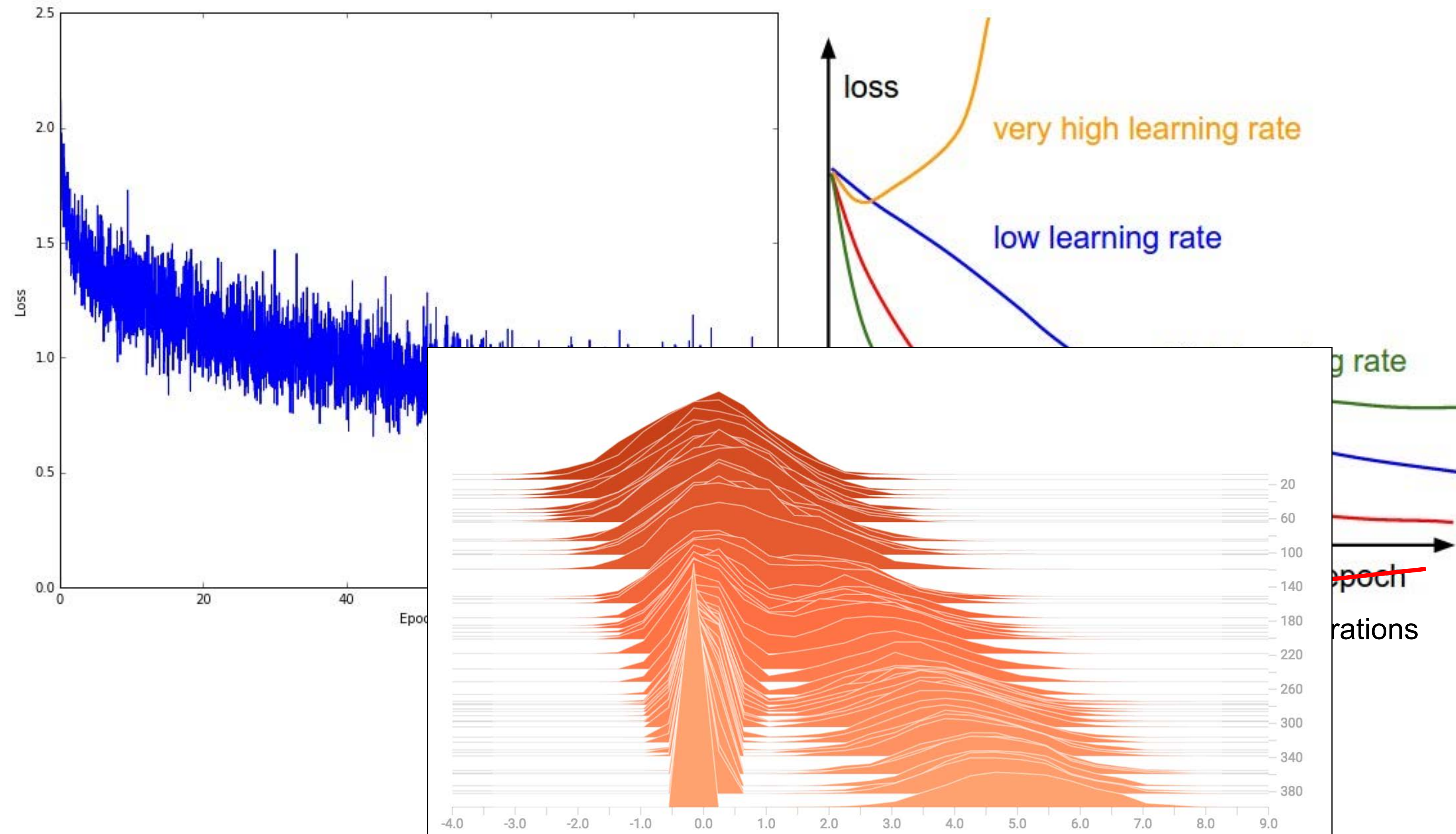
“A pepperoni pizza with arms and legs”



“A cute squirrel”

Skipped in class

Visualize VISUALIZE VISUALIZE (outside of scope)



More on Neural Networks

Skipped in class
(outside of scope)

Lots more to learn! A good place to start is

Justin Johnson, University of Michigan, EECS 498/598, e.g.,

<https://web.eecs.umich.edu/~justincj/teaching/eecs498/WI2022/>

Training Neural Nets: **Clever** Hans

Skipped in class
(outside of scope)



Hans could get 89% of the math questions right

Training Neural Nets: **Clever** Hans

Skipped in class
(outside of scope)



Hans could get 89% of the math questions right