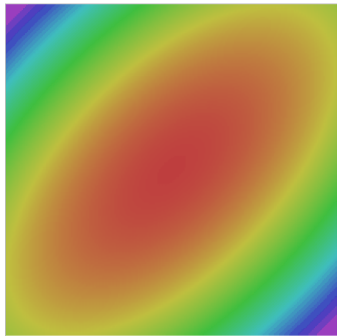
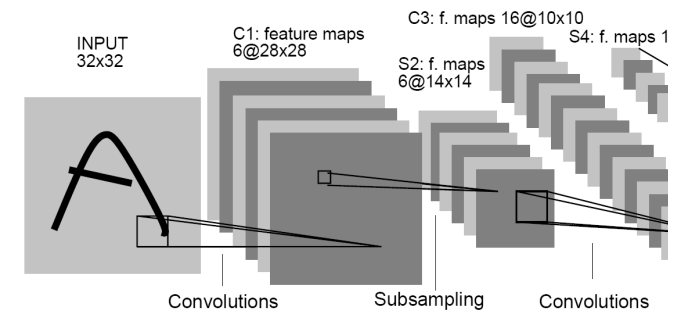


# Convolutional Neural Networks

## Computer Vision (UW EE/CSE 576)



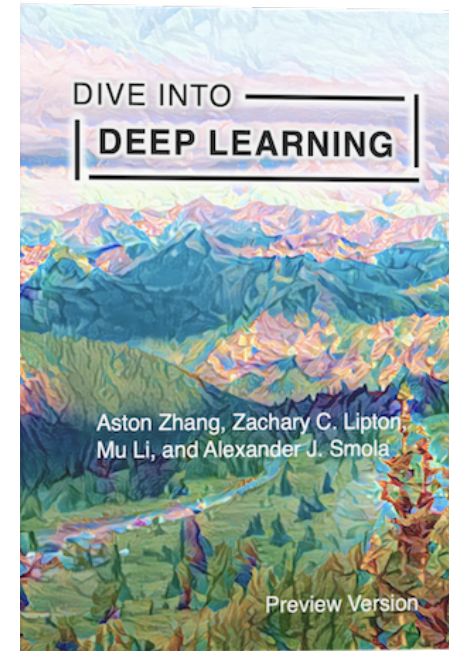
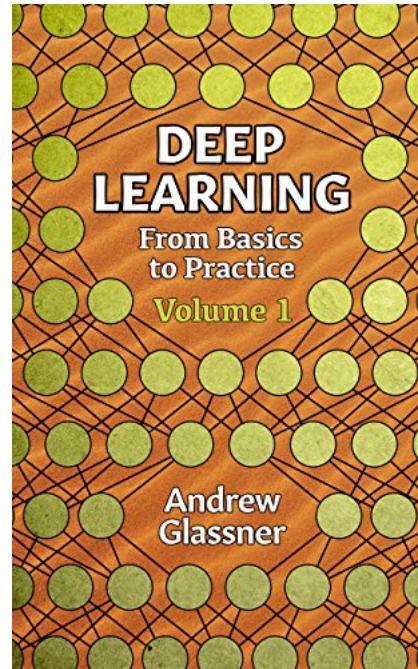
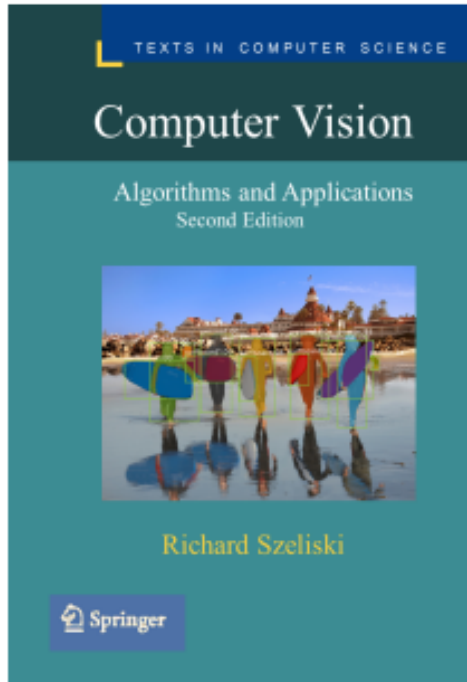
Richard Szeliski  
Facebook & UW  
Lecture 8 – Apr 23, 2020



# Class calendar

Date	Topic	Slides	Reading	Homework
April 9	Filters and convolutions	<a href="#">Google Slides</a>	<a href="#">Szeliski</a> , Chapter 3	HW1 due, <a href="#">HW2</a> assigned
April 14	Interpolation and Optimization	<a href="#">pdf</a> , <a href="#">pptx</a>	<a href="#">Szeliski</a> , Chapter 4	
April 16	Machine Learning	<a href="#">pdf</a> , <a href="#">pptx</a>	<a href="#">Szeliski</a> , Chapter 5.1-5.2	
April 21	Deep Neural Networks		<a href="#">Szeliski</a> , Chapter 5.3	
<b>April 23</b>	<b>Convolutional Neural Networks</b>		<a href="#">Szeliski</a> , Chapter 5.4	<b>HW2 due, HW3 assigned</b>
April 28	Network Architectures		<a href="#">Szeliski</a> , Chapter 5.4-5.5	
April 30	Object Detection		<a href="#">Szeliski</a> , Chapter 6.3	
May 5	Detection and Instance Segmentation		<a href="#">Szeliski</a> , Chapter 6.4	
May 7	Edges, features, matching, RANSAC		<a href="#">Szeliski</a> , Chapter 7.1-7.2, 8.1-8.2	HW3 due, HW4 assigned

# References



<https://d2l.ai/>

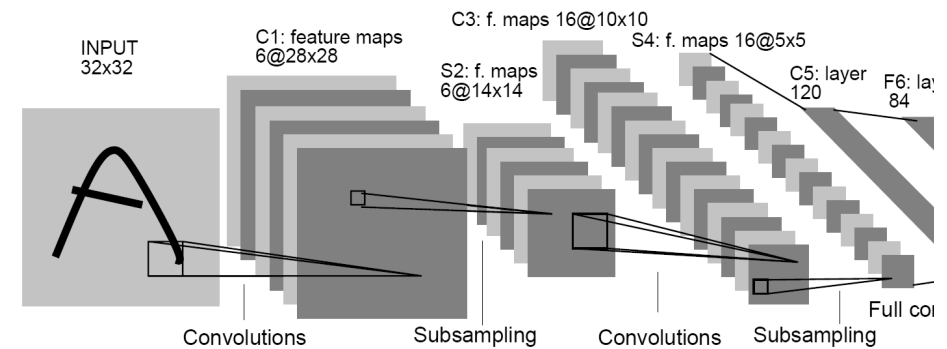
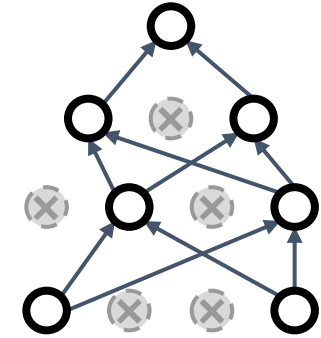
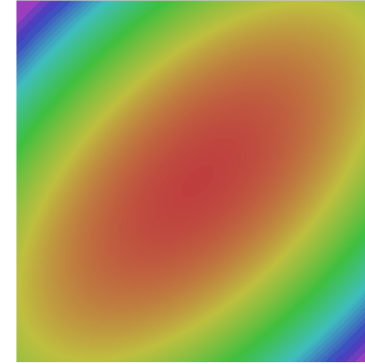
# Readings

## Chapter 5 **Deep Learning**

5.3	Deep neural networks . . . . .	272
5.3.1	Weights and layers . . . . .	274
5.3.2	Activation functions . . . . .	276
5.3.3	Regularization and normalization . . . . .	278
5.3.4	Loss functions . . . . .	283
5.3.5	Backpropagation . . . . .	285
5.3.6	Training and optimization . . . . .	289
5.4	Convolutional neural networks . . . . .	291
5.4.1	Pooling and unpooling . . . . .	295
5.4.2	<i>Application: Digit classification</i> . . . . .	297
5.4.3	Model zoos . . . . .	297
5.4.4	Visualizing weights and activations . . . . .	303
5.4.5	Adversarial examples . . . . .	306
5.4.6	Pre-training and fine-tuning networks . . . . .	306
5.5	More complex networks . . . . .	309

# Convolutional neural networks++

- Training and optimization
- More regularization (dropout, ...)
- Convolutional neural networks
- Pooling
- Batch normalization



# As before, I'm borrowing slides from



UNIVERSITY OF MICHIGAN

EECS 498-007 / 598-005  
Deep Learning for Computer Vision  
Fall 2019

## Course Description

Computer Vision has become ubiquitous in our society, with applications in search, image understanding, apps, mapping, medicine, drones, and self-driving cars. Core to many of these applications are visual recognition tasks such as image classification and object detection. Recent developments in neural network approaches have greatly advanced the performance of these state-of-the-art visual recognition systems. This course is a deep dive into details of neural-network based deep learning methods for computer vision. During this course, students will learn to implement, train and debug their own neural networks and gain a detailed understanding of cutting-edge research in computer vision. We will cover learning algorithms, neural network architectures, and practical engineering tricks for training and fine-tuning networks for visual recognition tasks.

### Instructor



Justin Johnson

### Graduate Student Instructors



Yunseok Jang



Kibok Lee



Luowei Zhou

Lecture 13	Thursday February 20	<b>Intro to Machine Learning</b> Image warping / blending Supervised vs Unsupervised learning Train / Test splits Linear Regression Regularization	<a href="#">[slides (pdf)]</a> <a href="#">[slides (pptx)]</a>
Lecture 14	Tuesday February 25	<b>Linear Models</b> Cross-Validation K-Nearest Neighbors SVM loss Cross-Entropy loss	<a href="#">[slides (pdf)]</a> <a href="#">[slides (pptx)]</a> <a href="#">[CS231n Linear Classification]</a>
Lecture 15	Thursday February 27	<b>Optimization</b> Stochastic Gradient Descent SGD + Momentum	<a href="#">[slides (pdf)]</a> <a href="#">[slides (pptx)]</a> <a href="#">[CS231n Optimization]</a>
Lecture 16	Tuesday March 10	<b>Neural Networks</b> Overfitting / Underfitting Bias / Variance tradeoff Fully-connected neural networks Biological neurons	<a href="#">[slides (pdf)]</a> <a href="#">[slides (pptx)]</a> <a href="#">[CS231n Neural Networks]</a>

Lecture 17	Tuesday March 17	<b>Backpropagation</b> Computational Graphs Backpropagation Matrix multiplication example	<a href="#">[video (from EECS 498/598)]</a> <a href="#">[slides (from EECS 498/598)]</a> <a href="#">[231n Backpropagation]</a> <a href="#">[Backprop for Matrix Multiply]</a> <a href="#">[Olah on Backprop]</a> <a href="#">[Nielsen on Backprop]</a>
Lecture 18	Thursday March 19	<b>Convolutional Networks</b> Convolution Pooling Batch Normalization	<a href="#">[video (from EECS 498/598)]</a> <a href="#">[slides (from EECS 498/598)]</a> <a href="#">[231n ConvNets]</a> <a href="#">[Goodfellow, Chapter 9]</a>
Lecture 19	Tuesday March 24	<b>CNN Architectures</b> AlexNet, VGG, ResNet Size vs Accuracy Neural Architecture Search	<a href="#">[video (from EECS 498/598)]</a> <a href="#">[slides (from EECS 498/598)]</a> <a href="#">[AlexNet paper]</a> <a href="#">[VGG paper]</a> <a href="#">[GoogLeNet paper]</a> <a href="#">[ResNet paper]</a>
Lecture 20	Thursday March 26	<b>Training Neural Networks I</b> Activation Functions Data preprocessing Weight initialization Data Augmentation Regularization	<a href="#">[video (from EECS 498/598)]</a> <a href="#">[slides (from EECS 498/598)]</a> <a href="#">[231n Training I]</a>
Lecture 21	Tuesday March 31	<b>Training Neural Networks II</b> Learning rate schedules Hyperparameter optimization Model ensembles Transfer learning Large-batch training	<a href="#">[video (EECS 498/598)]</a> <a href="#">[slides (from EECS 498/598)]</a> <a href="#">[231n Training II]</a> <a href="#">[Karpathy "Recipe for Training"]</a>

# Lecture 4: Optimization



# Loss Functions quantify preferences

- We have some dataset of  $(x, y)$
- We have a **score function**:
- We have a **loss function**:

Q: How do we find the best  $W$ ?

$$s = f(x; W) = Wx$$

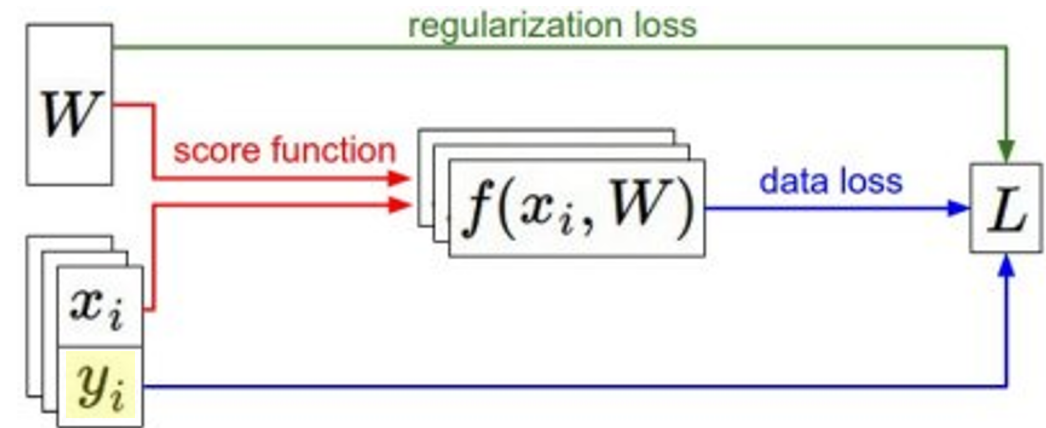
Linear classifier

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) \text{ Softmax}$$

SVM

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + R(W) \text{ Full loss}$$



# Follow the slope

In 1-dimension, the **derivative** of a function gives the slope:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

In multiple dimensions, the **gradient** is the vector of (partial derivatives) along each dimension

The slope in any direction is the **dot product** of the **direction** with the gradient  
The direction of steepest descent is the **negative gradient**

# Loss is a function of $W$ : Analytic Gradient

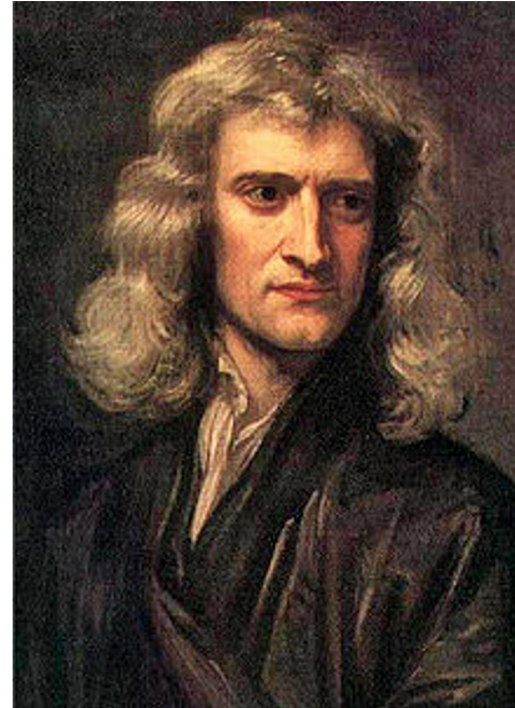
$$L = \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$s = f(x; W) = Wx$$

want  $\nabla_W L$

Use calculus to compute an  
**analytic gradient**



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



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

# Computing Gradients

- **Numeric gradient:** approximate, slow, easy to write
- **Analytic gradient:** exact, fast, error-prone

What's the difference?

Which one is better?

# Computing Gradients

- **Numeric gradient:** approximate, slow, easy to write
- **Analytic gradient:** exact, fast, error-prone

In practice: Always use analytic gradient, but check implementation with numerical gradient. This is called a **gradient check**.

```
def grad_check_sparse(f, x, analytic_grad, num_checks=10, h=1e-7):  
    """  
    sample a few random elements and only return numerical  
    in this dimensions.  
    """
```

# Gradient Descent

Iteratively step in the direction of  
the negative gradient  
(direction of local steepest descent)

```
# Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
    dw = compute_gradient(loss_fn, data, w)
    w -= learning_rate * dw
```

## Hyperparameters:

- Weight initialization method
- Number of steps
- Learning rate

$$w \leftarrow w - \alpha g \quad \text{or}$$
$$w_{t+i} = w_t - \alpha_t g_t$$

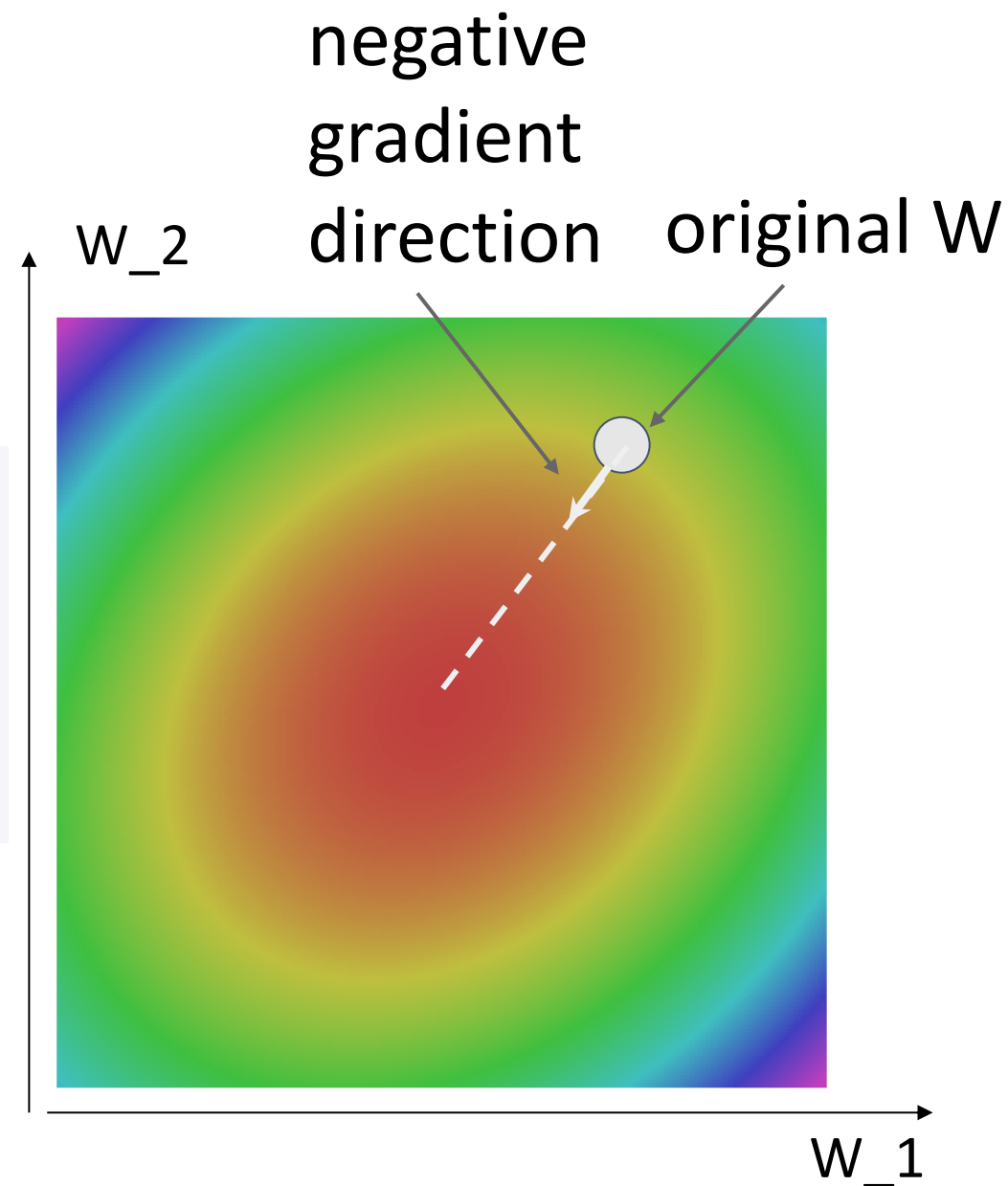
# Gradient Descent

Iteratively step in the direction of the negative gradient  
(direction of local steepest descent)

```
# Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
    dw = compute_gradient(loss_fn, data, w)
    w -= learning_rate * dw
```

## Hyperparameters:

- Weight initialization method
- Number of steps
- Learning rate



# Batch Gradient Descent

Full sum expensive  
when N is large!

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$



# [Minibatch] Stochastic Gradient Descent (SGD)

Full sum expensive  
when N is large!

Approximate sum using  
a **minibatch** of examples  
32 / 64 / 128 common

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

```
# Stochastic gradient descent
w = initialize_weights()
for t in range(num_steps):
    minibatch = sample_data(data, batch_size)
    dw = compute_gradient(loss_fn, minibatch, w)
    w -= learning_rate * dw
```

## Hyperparameters:

- Weight initialization
- Number of steps
- Learning rate
- Batch size
- Data sampling

# Stochastic Gradient Descent (SGD)

$$L(W) = \mathbb{E}_{(x,y) \sim p_{data}} [L(x, y, W)] + \lambda R(W)$$
$$\approx \frac{1}{N} \sum_{i=1}^N L(x_i, y_i, W) + \lambda R(W)$$

Think of loss as an expectation over the full **data distribution**  $p_{data}$

Approximate expectation via sampling

# Stochastic Gradient Descent (SGD)

Think of loss as an expectation over the full **data distribution**  $p_{\text{data}}$

$$L(W) = \mathbb{E}_{(x,y) \sim p_{\text{data}}} [L(x, y, W)] + \lambda R(W)$$

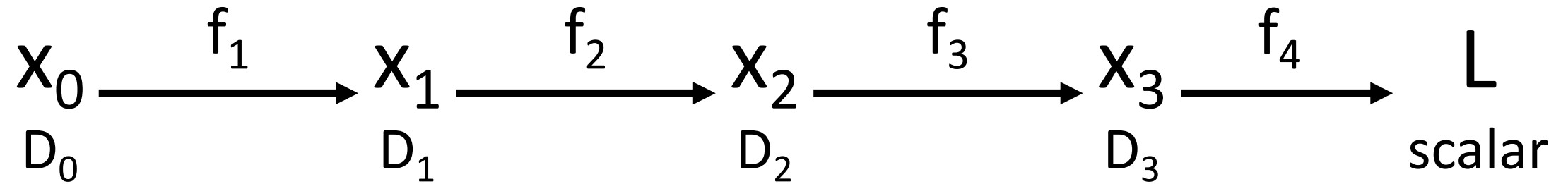
$$\approx \frac{1}{N} \sum_{i=1}^N L(x_i, y_i, W) + \lambda R(W)$$

Approximate expectation via sampling

$$\nabla_W L(W) = \nabla_W \mathbb{E}_{(x,y) \sim p_{\text{data}}} [L(x, y, W)] + \lambda \nabla_W R(W)$$

$$\approx \sum_{i=1}^N \nabla_W L_W(x_i, y_i, W) + \nabla_W R(W)$$

# Recall: Reverse-Mode Automatic Differentiation



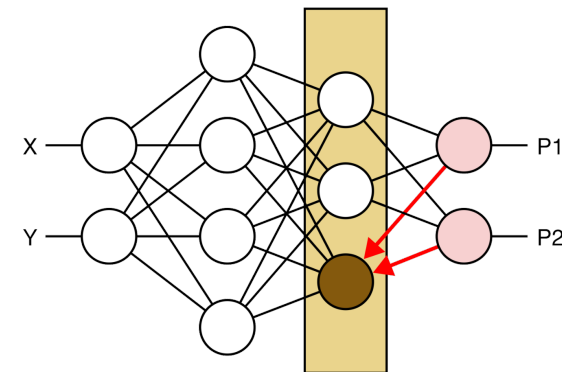
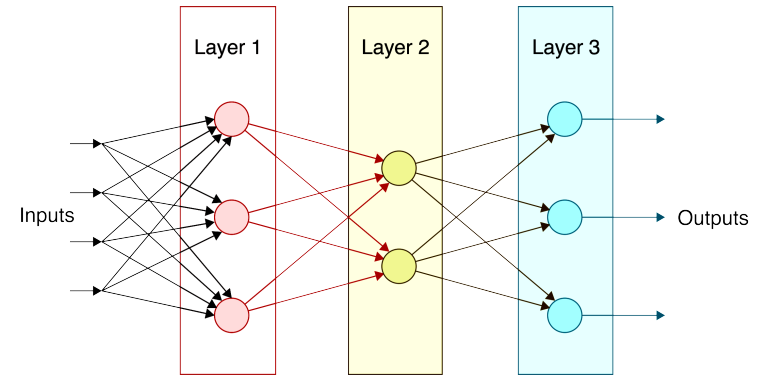
Matrix multiplication is **associative**: we can compute products in any order  
Computing products right-to-left avoids matrix-matrix products; only needs matrix-vector

Chain rule  $\frac{\partial L}{\partial x_0} = \left( \frac{\partial x_1}{\partial x_0} \right) \left( \frac{\partial x_2}{\partial x_1} \right) \left( \frac{\partial x_3}{\partial x_2} \right) \left( \frac{\partial L}{\partial x_3} \right)$

$D_0 \times D_1 \quad D_1 \times D_2 \quad D_2 \times D_3 \quad D_3$

# Mini-batch evaluation with matrices (HW3)

- DNNs are described as passing **vectors** between layers
- Why not pass all samples in a mini-batch as a **matrix**?
- What used to be column vectors are now rows
- Need to adjust weight-vector multiplies
$$\mathbf{s} = \mathbf{W} \mathbf{x}$$
becomes
$$\mathbf{S} = \mathbf{X} \mathbf{W}^T$$
- Need to adjust gradients (Jacobians) as well



# Homework 3: Neural Networks in C++

## What you'll be implementing

We will be training a fully-connected neural network for this assignment.

- `src/activation.cpp`: you will implement the forward and backward passes for several activation functions.
- `src/classifier.cpp`: you will implement gradient computation and parameter updates using algorithms we discussed in class.

You'll be training on two datasets, one is MNIST which is a digit-recognition dataset. The other is a simple visual recognition dataset called CIFAR.

## 1. Implementing Neural Networks

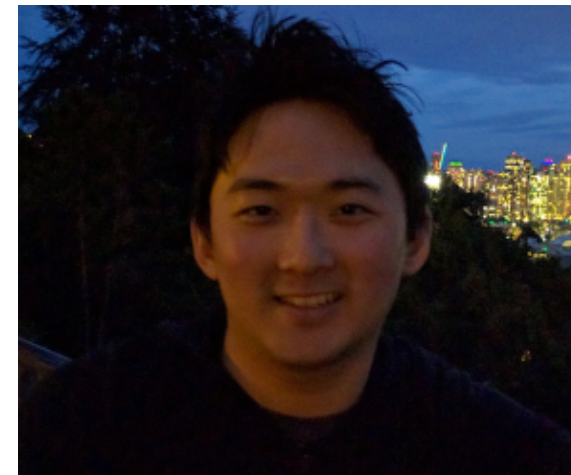
### 1.1 Activation Functions

An important part of machine learning, be it linear classifiers or neural networks, is the activation function you use.

We will be implementing the following activation functions:

- Linear:  $f(x) = x$
- Logistic:  $f(x) = 1 / (1 + \exp(-x))$
- tanh:  $f(x) = \tanh(x)$
- ReLU:  $f(x) = \max(0, x)$
- Leaky ReLU:  $f(x) = 0.01*x$  if  $x < 0$  else  $x$
- Softmax: [https://en.wikipedia.org/wiki/Softmax\\_function](https://en.wikipedia.org/wiki/Softmax_function)

- Quick overview by Keunhong Park

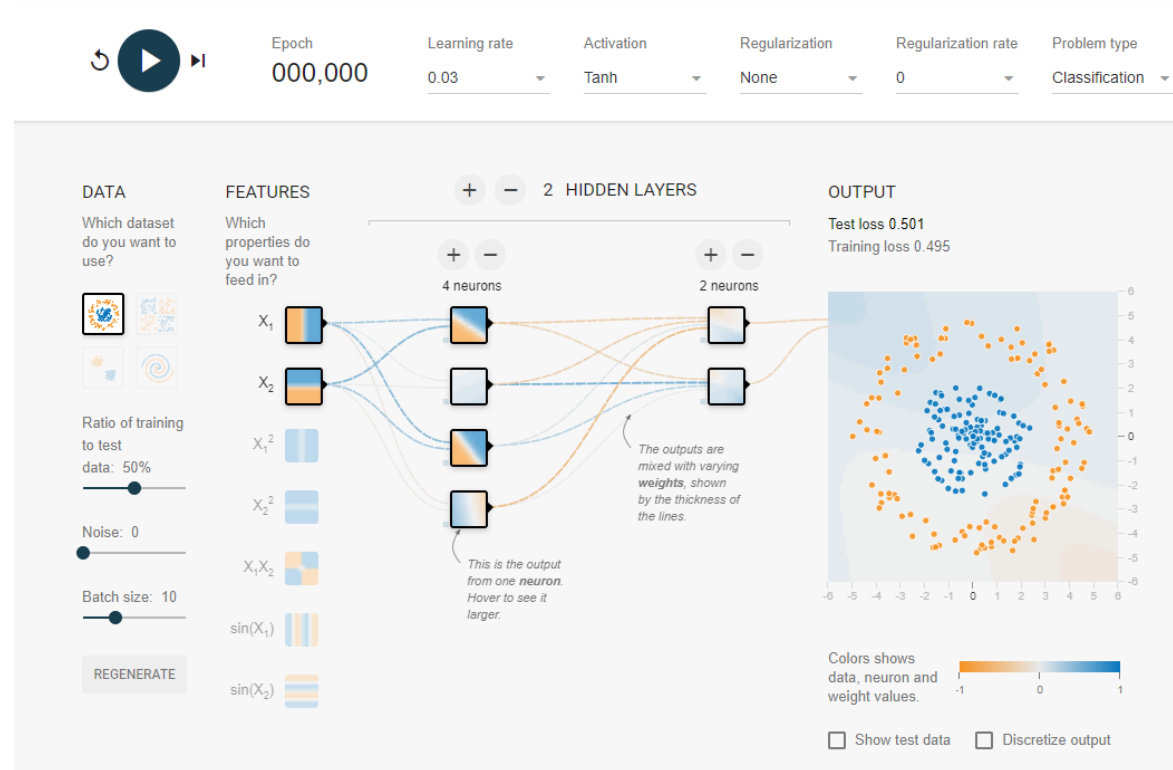


**Keunhong Park (TA)** <sup>CS</sup>

Ph.D Student

[View full profile](#)

# Interactive Web demo

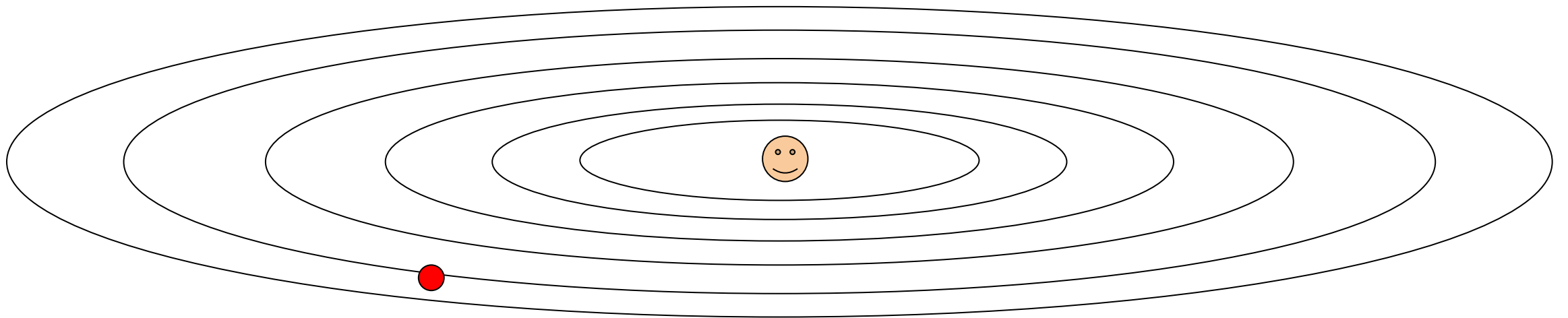


*(Thanks, Doruk)*

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

# Problems with SGD

What if loss changes quickly in one direction and slowly in another?  
What does gradient descent do?



Loss function has high **condition number**: ratio of largest to smallest singular value of the Hessian matrix is large

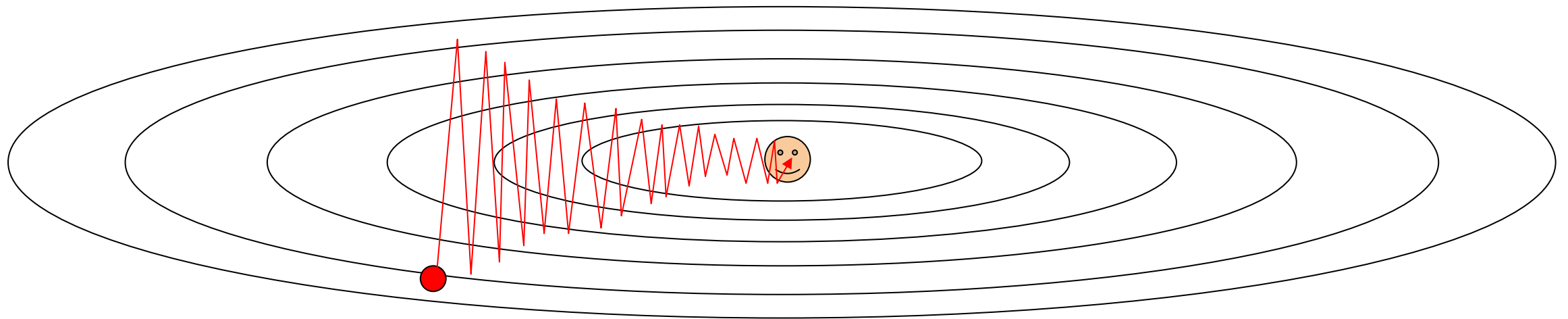


# Problems with SGD

What if loss changes quickly in one direction and slowly in another?

What does gradient descent do?

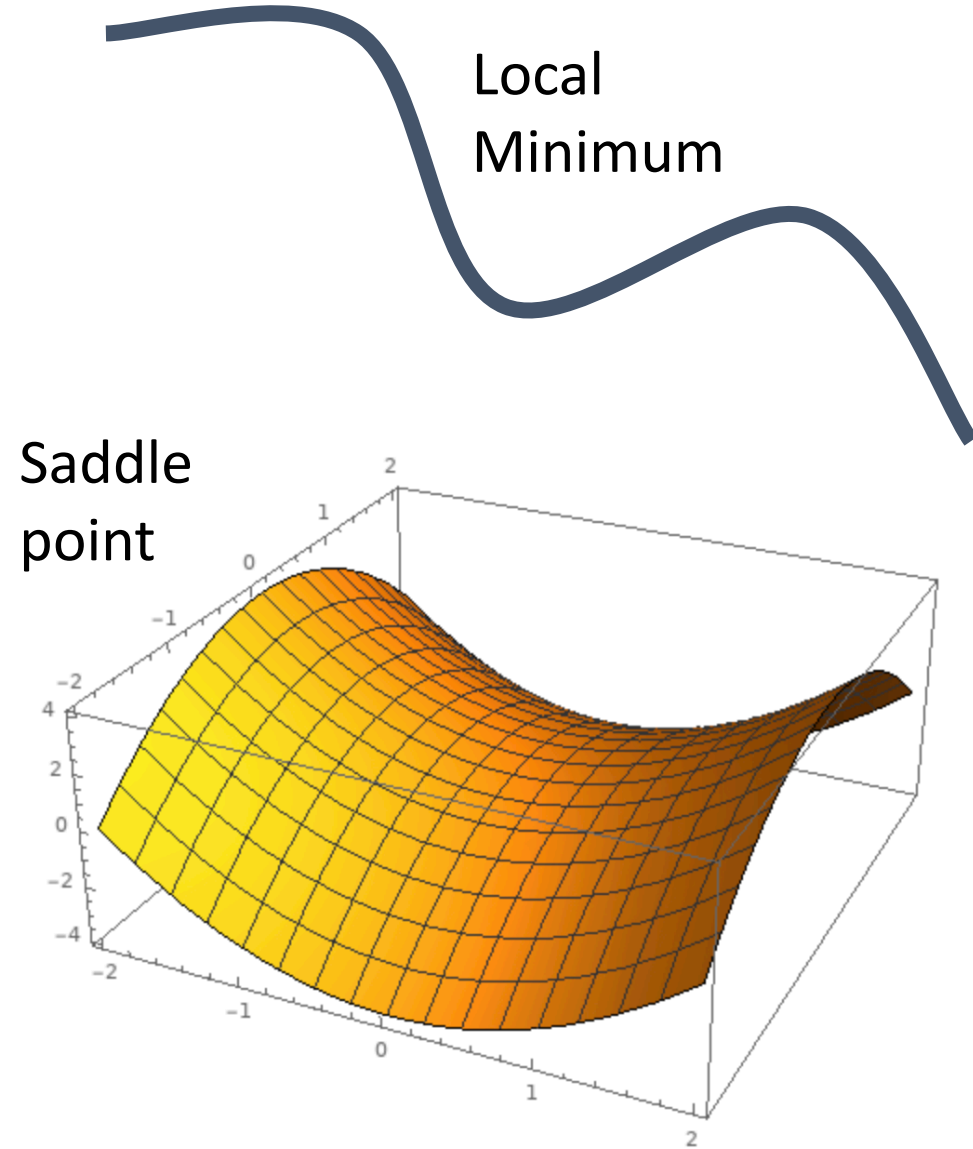
Very slow progress along shallow dimension, jitter along steep direction



Loss function has high **condition number**: ratio of largest to smallest singular value of the Hessian matrix is large

# Problems with SGD

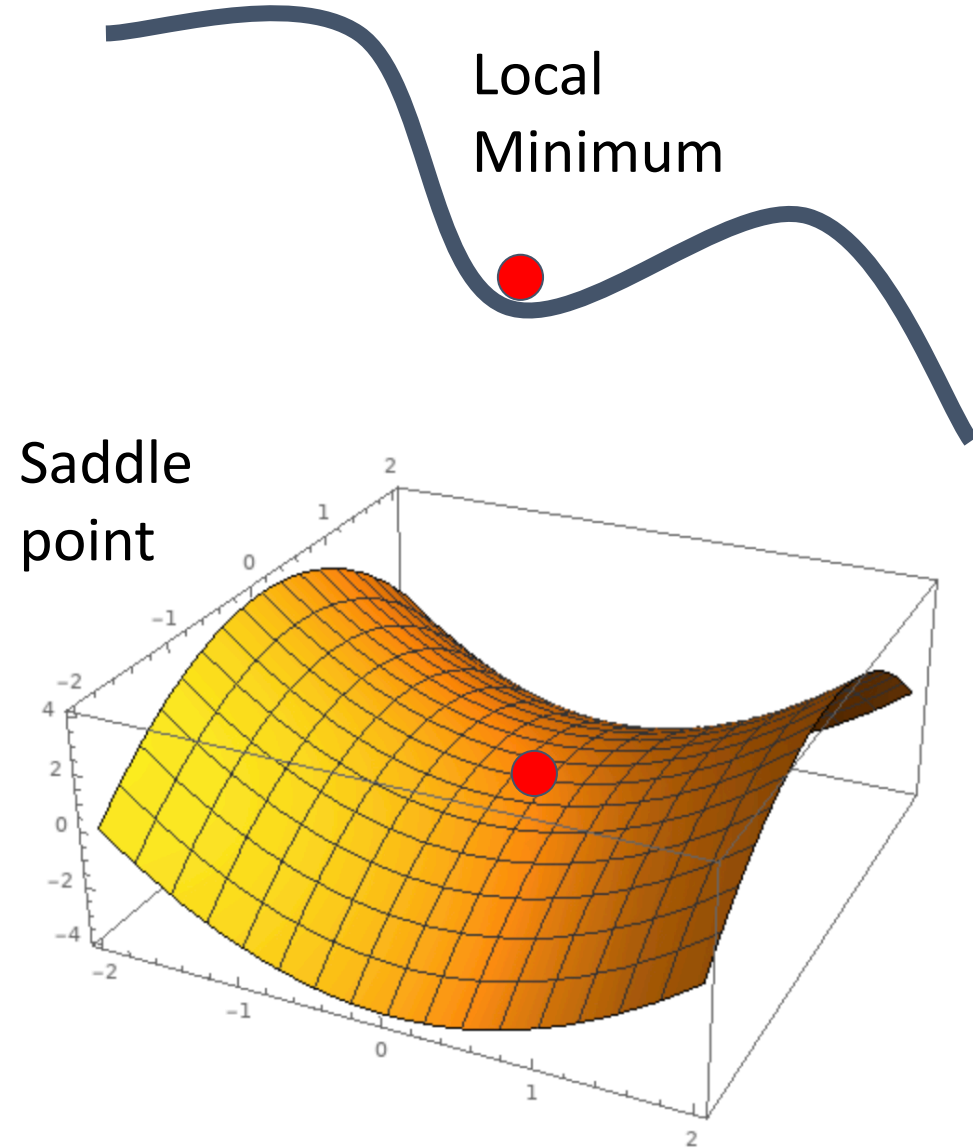
What if the loss function has a **local minimum** or **saddle point**?



# Problems with SGD

What if the loss function has a **local minimum** or **saddle point**?

Zero gradient, gradient descent gets stuck

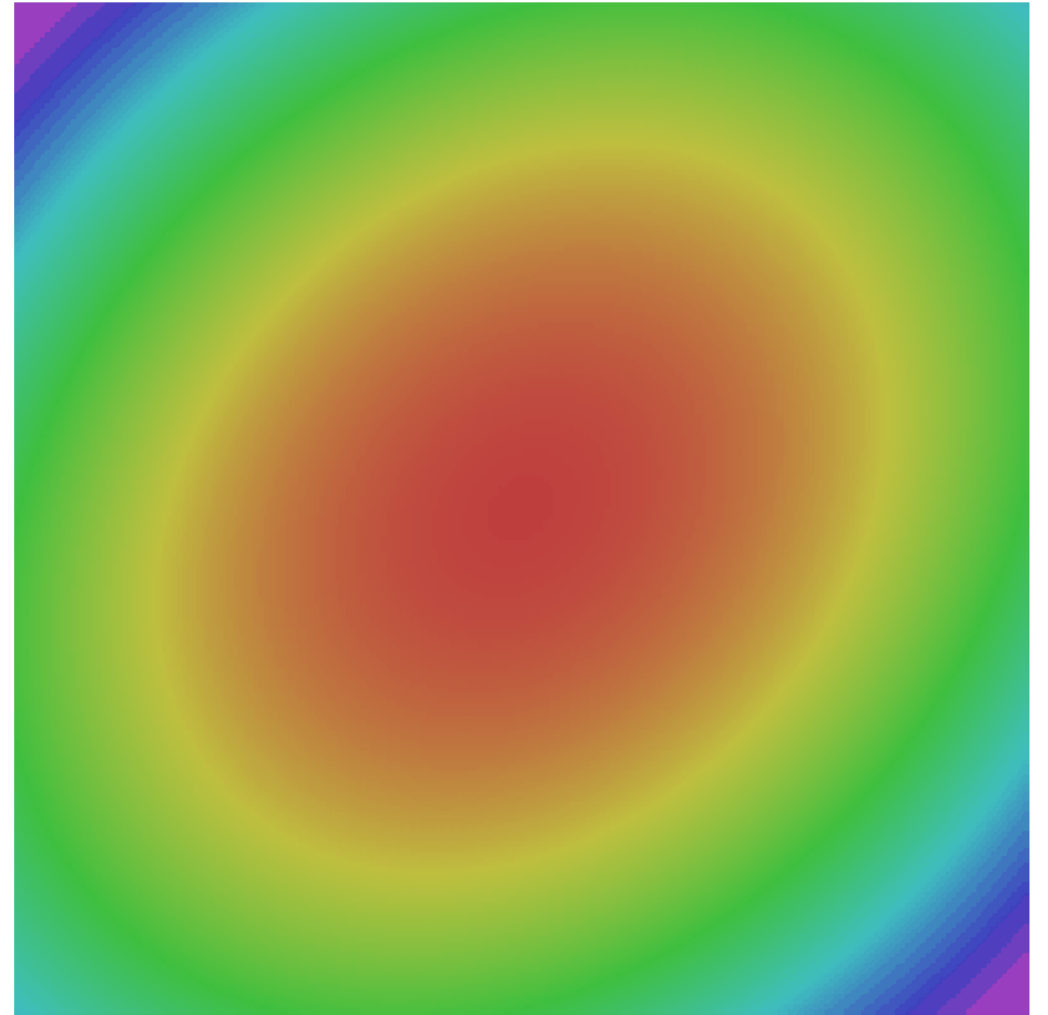


# Problems with SGD

Our gradients come from minibatches  
so they can be noisy!

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W)$$



# SGD

## SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
for t in range(num_steps):  
    dw = compute_gradient(w)  
    w -= learning_rate * dw
```

# SGD + Momentum

## SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
for t in range(num_steps):  
    dw = compute_gradient(w)  
    w -= learning_rate * dw
```

## SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
v = 0  
for t in range(num_steps):  
    dw = compute_gradient(w)  
    v = rho * v + dw  
    w -= learning_rate * v
```

- Build up “velocity” as a running mean of gradients
- Rho gives “friction”; typically rho=0.9 or 0.99

# SGD + Momentum

## SGD+Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v - learning_rate * dw
    w += v
```

## SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

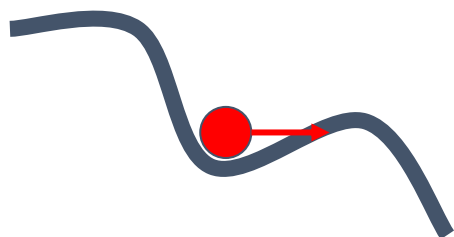
$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```

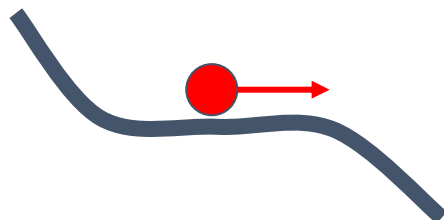
You may see SGD+Momentum formulated different ways, but they are equivalent - give same sequence of  $x$

# SGD + Momentum

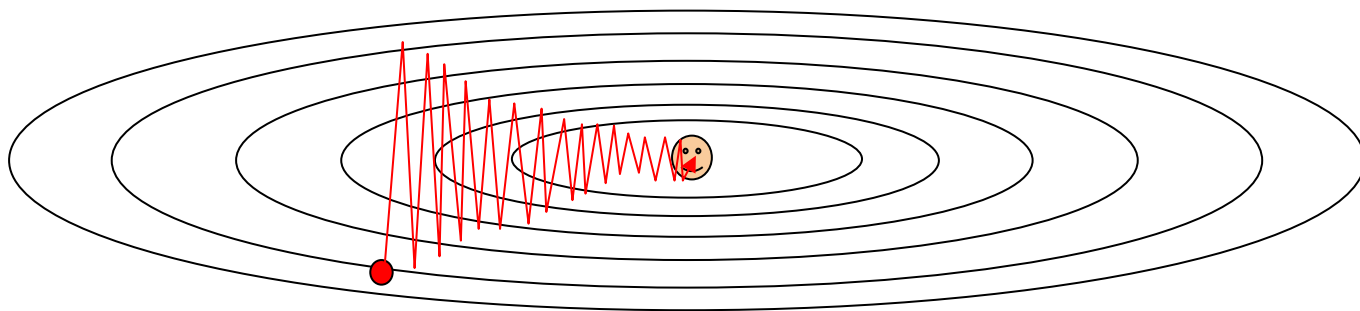
Local Minima



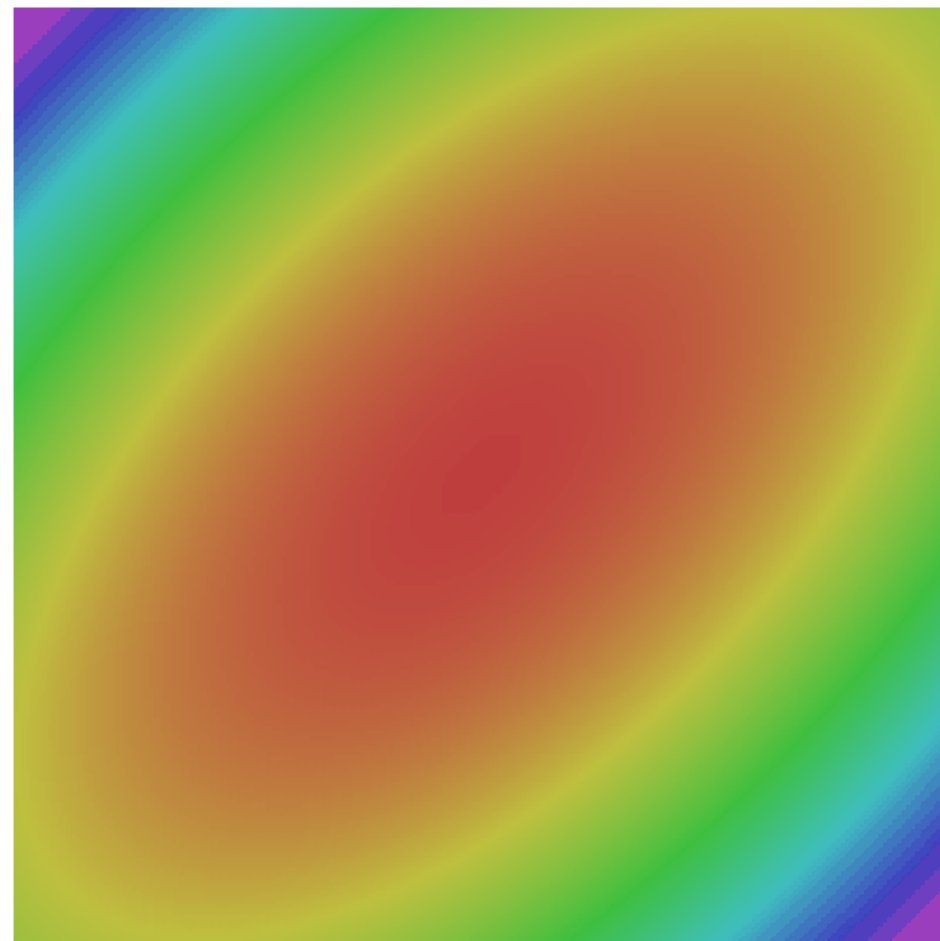
Saddle points



Poor Conditioning



# Gradient Noise



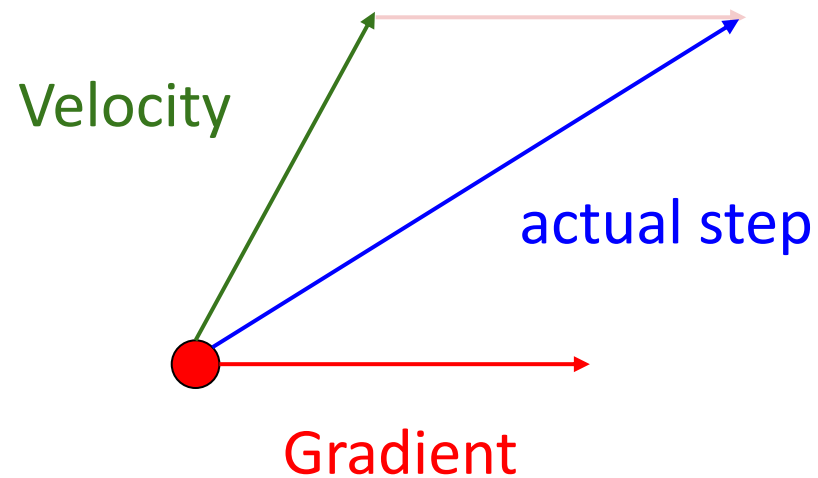
— SGD — SGD+Momentum

Sutskever et al, "On the importance of initialization and momentum in deep learning", ICML 2013



# SGD + Momentum

Momentum update:



Combine gradient at current point with velocity to get step used to update weights

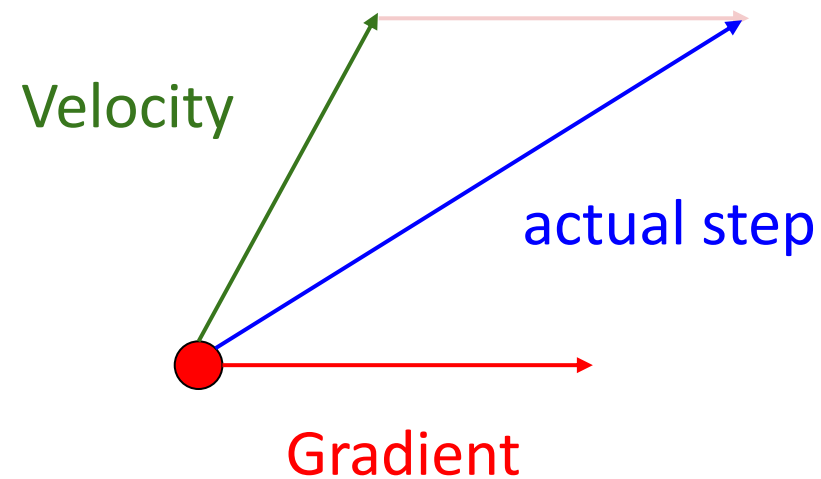
Nesterov, "A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ ", 1983

Nesterov, "Introductory lectures on convex optimization: a basic course", 2004

Sutskever et al, "On the importance of initialization and momentum in deep learning", ICML 2013

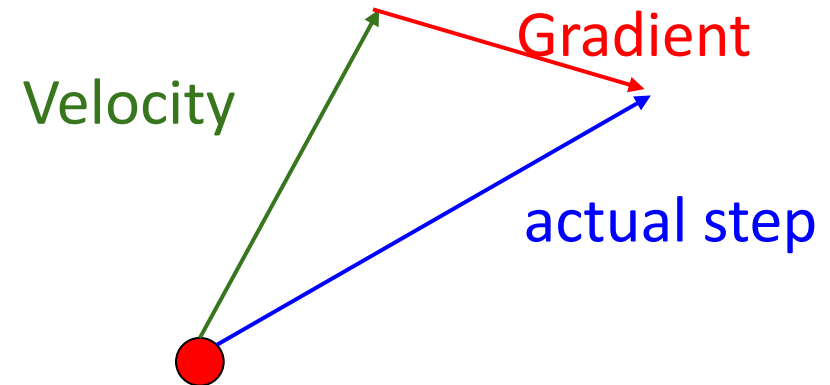
# Nesterov Momentum

Momentum update:



Combine gradient at current point with velocity to get step used to update weights

Nesterov Momentum



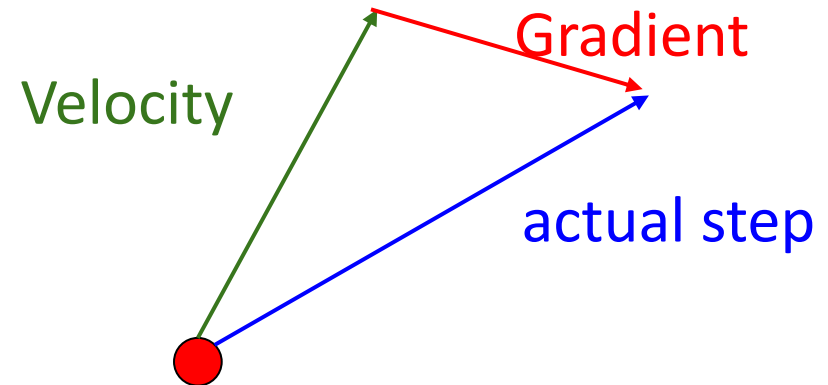
“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

Nesterov, “A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ ”, 1983  
Nesterov, “Introductory lectures on convex optimization: a basic course”, 2004  
Sutskever et al, “On the importance of initialization and momentum in deep learning”, ICML 2013

# Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$



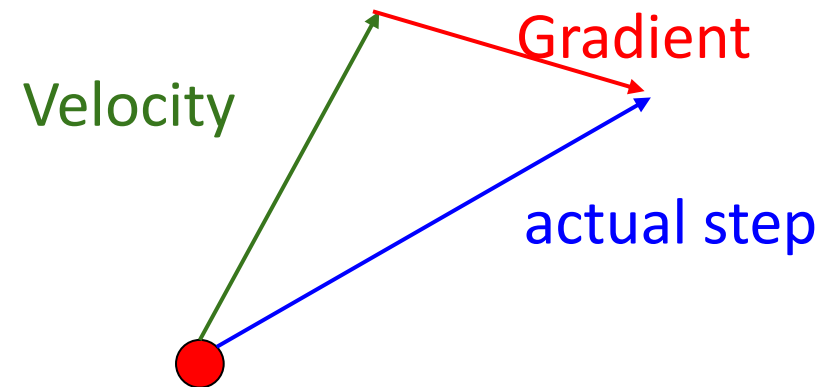
“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

# Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

Annoying, usually we want update in terms of  $x_t, \nabla f(x_t)$



“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

# Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

Change of variables  $\tilde{x}_t = x_t + \rho v_t$   
and rearrange:

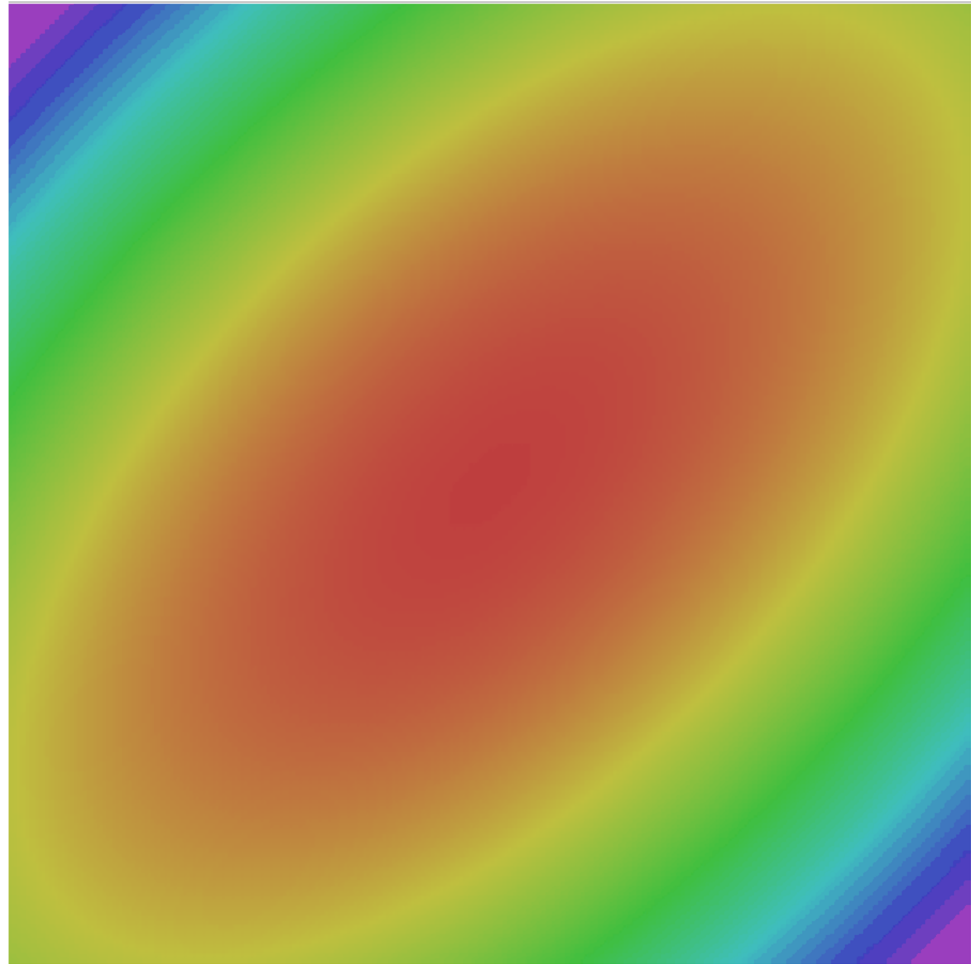
$$v_{t+1} = \rho v_t - \alpha \nabla f(\tilde{x}_t)$$

$$\begin{aligned}\tilde{x}_{t+1} &= \tilde{x}_t - \rho v_t + (1 + \rho)v_{t+1} \\ &= \tilde{x}_t + v_{t+1} + \rho(v_{t+1} - v_t)\end{aligned}$$

Annoying, usually we want  
update in terms of  $x_t, \nabla f(x_t)$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    old_v = v
    v = rho * v - learning_rate * dw
    w -= rho * old_v - (1 + rho) * v
```

# Nesterov Momentum



- SGD
- SGD+Momentum
- Nesterov

# AdaGrad

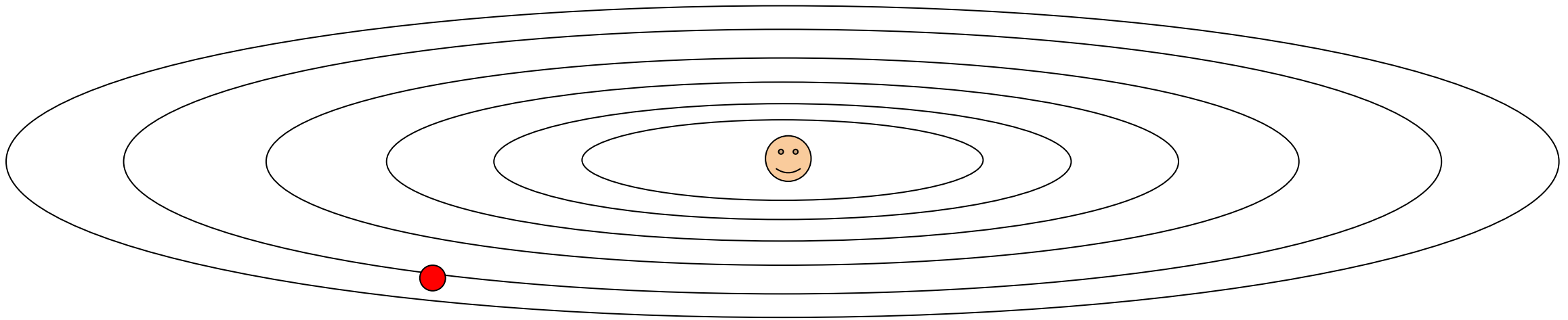
```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

Added element-wise scaling of the gradient based on the historical sum of squares in each dimension

“Per-parameter learning rates”  
or “adaptive learning rates”

# AdaGrad

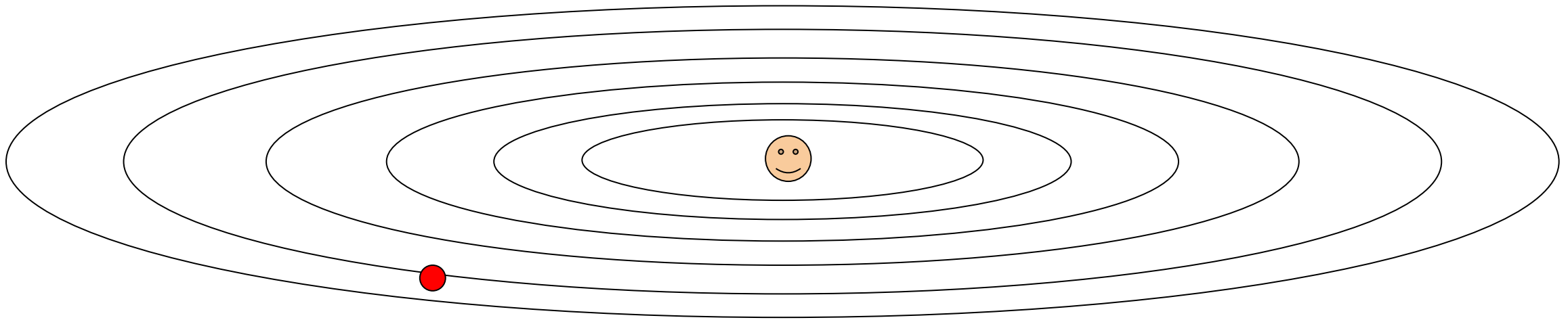
```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```





# AdaGrad

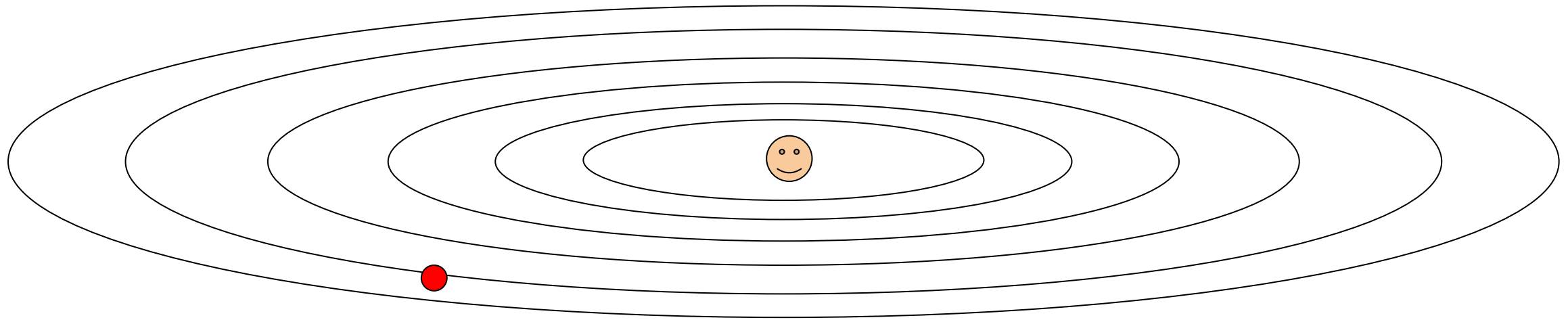
```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```



Q: What happens with AdaGrad?

# AdaGrad

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```



Q: What happens with AdaGrad?

Progress along “steep” directions is damped;  
progress along “flat” directions is accelerated

# RMSProp: “Leaky Adagrad”

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

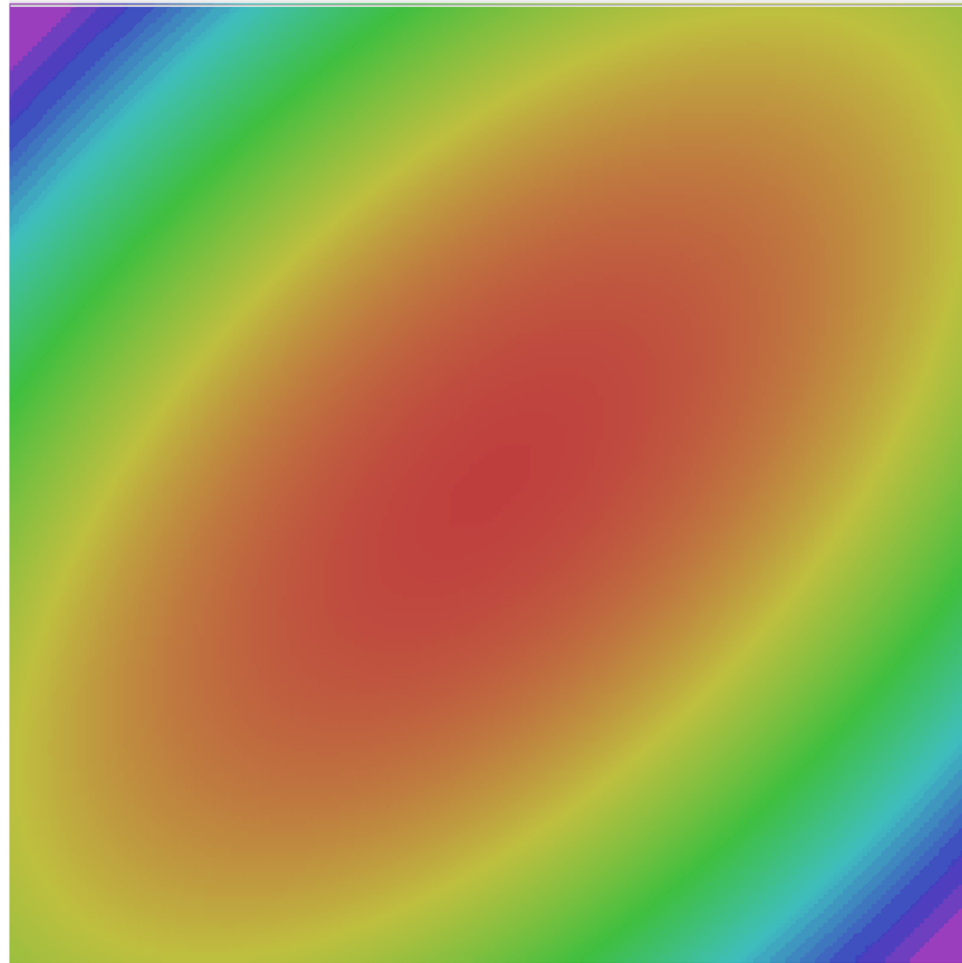
AdaGrad

↓

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

RMSProp

# RMSProp



- SGD
- SGD+Momentum
- RMSProp

# Adam (almost): RMSProp + Momentum

```
moment1 = 0
moment2 = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

# Adam (almost): RMSProp + Momentum

```
moment1 = 0
moment2 = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

Adam

Momentum

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```

SGD+Momentum

# Adam (almost): RMSProp + Momentum

```
moment1 = 0
moment2 = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

Adam

Momentum

AdaGrad / RMSProp

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

RMSProp

# Adam: Very Common in Practice!

for input to the CNN; each colored pixel in the image yields a 7D one-hot vector. Following common practice, the network is trained end-to-end using stochastic gradient descent with the Adam optimizer [22]. We anneal the learning rate to 0 using a half cosine schedule without restarts [28].

Bakhtin, van der Maaten, Johnson, Gustafson, and Girshick, NeurIPS 2019

We train all models using Adam [23] with learning rate  $10^{-4}$  and batch size 32 for 1 million iterations; training takes about 3 days on a single Tesla P100. For each mini-batch we first update  $f$ , then update  $D_{img}$  and  $D_{obj}$ .

Johnson, Gupta, and Fei-Fei, CVPR 2018

ganized into three residual blocks. We train for 25 epochs using Adam [27] with learning rate  $10^{-4}$  and 32 images per batch on 8 Tesla V100 GPUs. We set the cubify thresh-

Gkioxari, Malik, and Johnson, ICCV 2019

sampled with each bit drawn uniformly at random. For gradient descent, we use Adam [29] with a learning rate of  $10^{-3}$  and default hyperparameters. All models are trained with batch size 12. Models are trained for 200 epochs, or 400 epochs if being trained on multiple noise layers.

Zhu, Kaplan, Johnson, and Fei-Fei, ECCV 2018

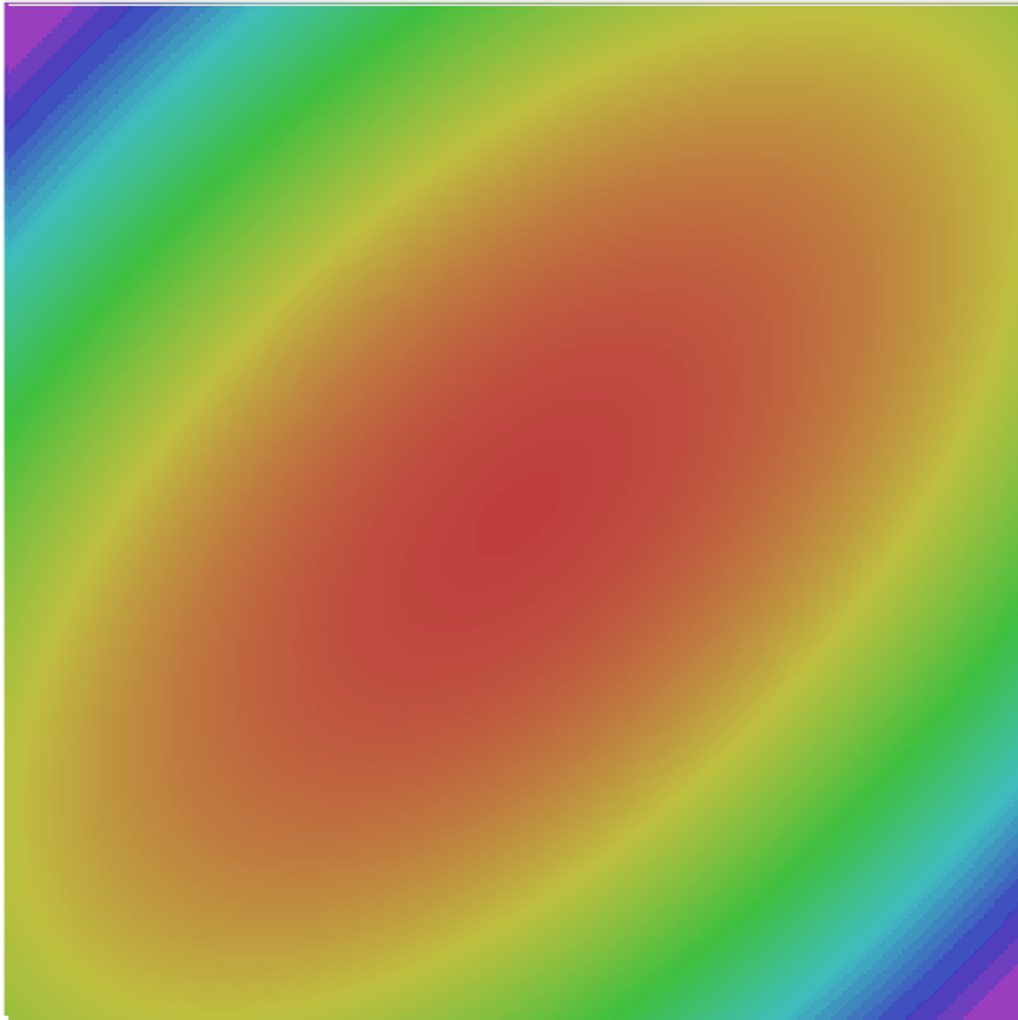
16 dimensional vectors. We iteratively train the Generator and Discriminator with a batch size of 64 for 200 epochs using Adam [22] with an initial learning rate of 0.001.

Gupta, Johnson, et al, CVPR 2018

Adam with  $\beta_1 = 0.9$ ,  
 $\beta_2 = 0.999$ , and  $\text{learning\_rate} = 1e-3, 5e-4, 1e-4$   
is a great starting point for many models!



# Adam



- SGD
- SGD+Momentum
- RMSProp
- Adam

# Optimization Algorithm Comparison

Algorithm	Tracks first moments (Momentum)	Tracks second moments (Adaptive learning rates)	Leaky second moments	Bias correction for moment estimates
SGD	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>
SGD+Momentum	✓	<b>X</b>	<b>X</b>	<b>X</b>
Nesterov	✓	<b>X</b>	<b>X</b>	<b>X</b>
AdaGrad	<b>X</b>	✓	<b>X</b>	<b>X</b>
RMSProp	<b>X</b>	✓	✓	<b>X</b>
Adam	✓	✓	✓	✓

In practice:

- **Adam** is a good default choice in many cases  
**SGD+Momentum** can outperform Adam but may require more tuning
- If you can afford to do full batch updates then try out **L-BFGS** (and don't forget to disable all sources of noise)

# Lecture 10: Training Neural Networks (Part 1)

# Overview

## **1. One time setup**

Activation functions, data preprocessing, weight initialization, regularization

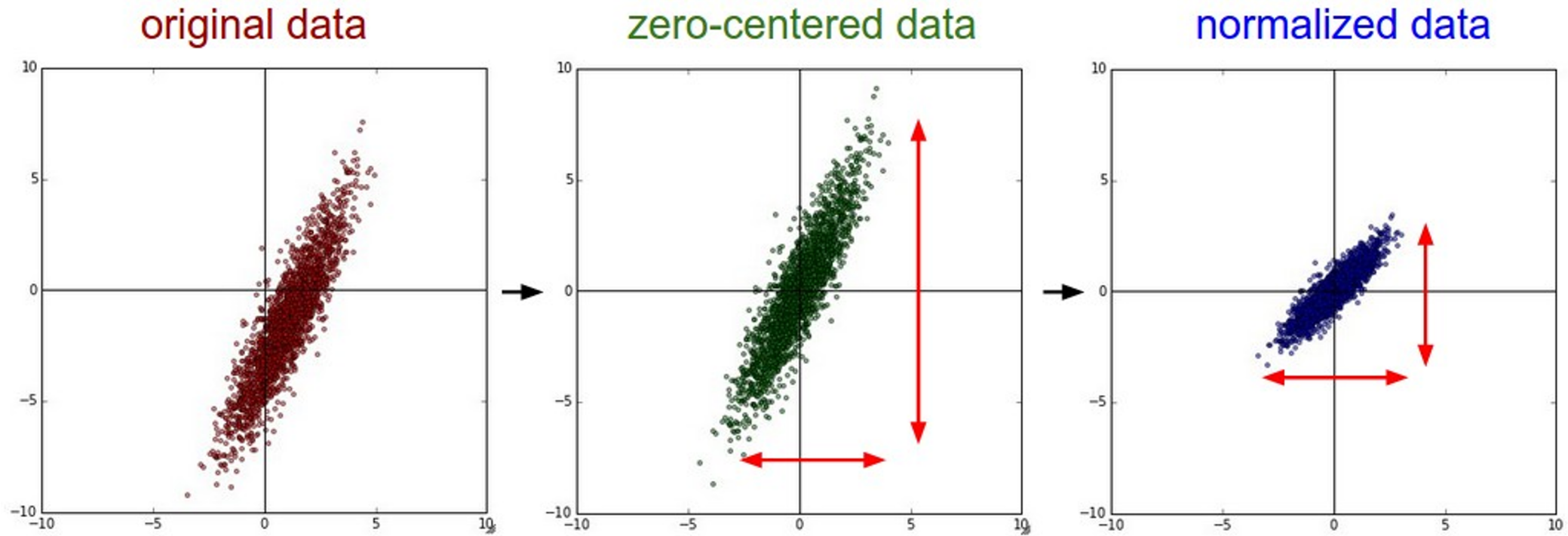
## **2. Training dynamics**

Learning rate schedules; large-batch training; hyperparameter optimization

## **3. After training**

Model ensembles, transfer learning

# Snapshot: Data Preprocessing

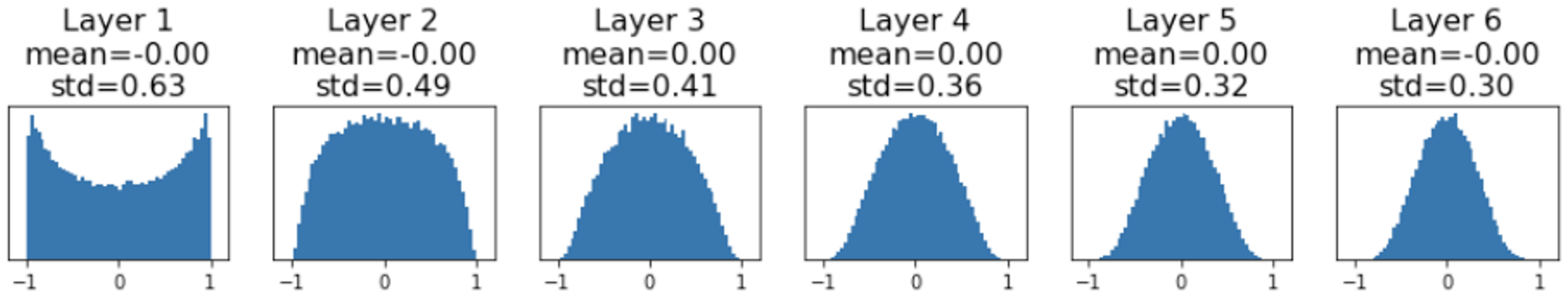


# Snapshot: Weight Initialization

```
dims = [4096] * 7
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

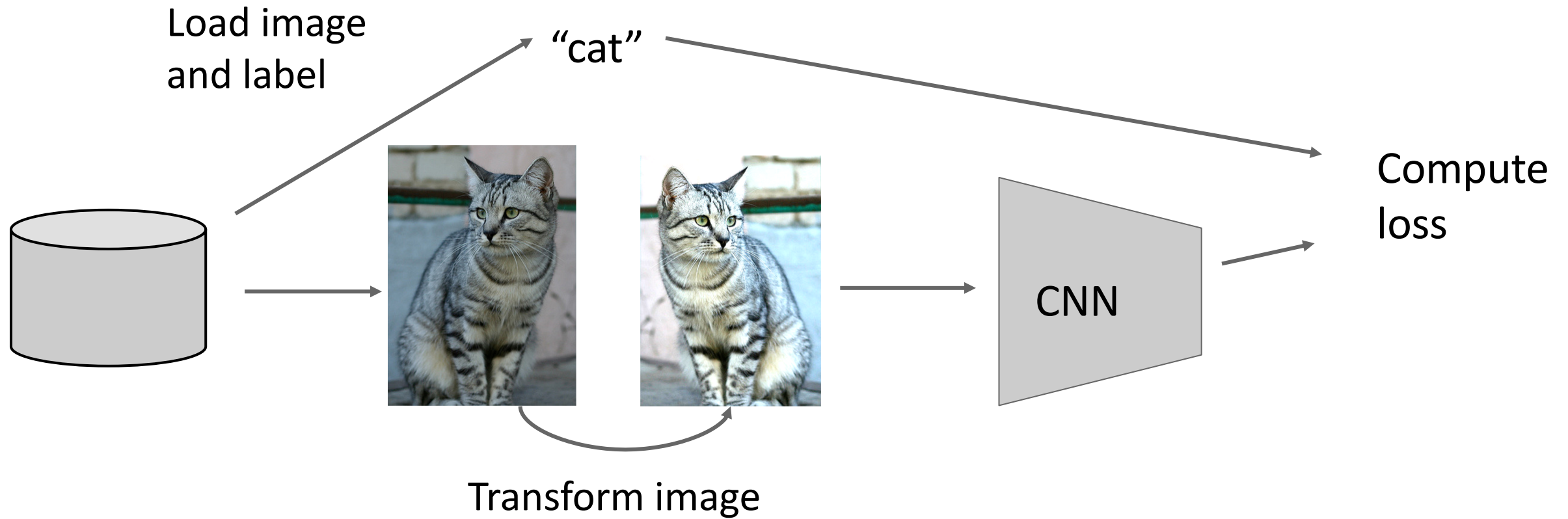
“Xavier” initialization:  
std =  $1/\sqrt{D_{in}}$

“Just right”: Activations are nicely scaled for all layers!



Glorot and Bengio, “Understanding the difficulty of training deep feedforward neural networks”, AISTAT 2010

# Snapshot: Data Augmentation



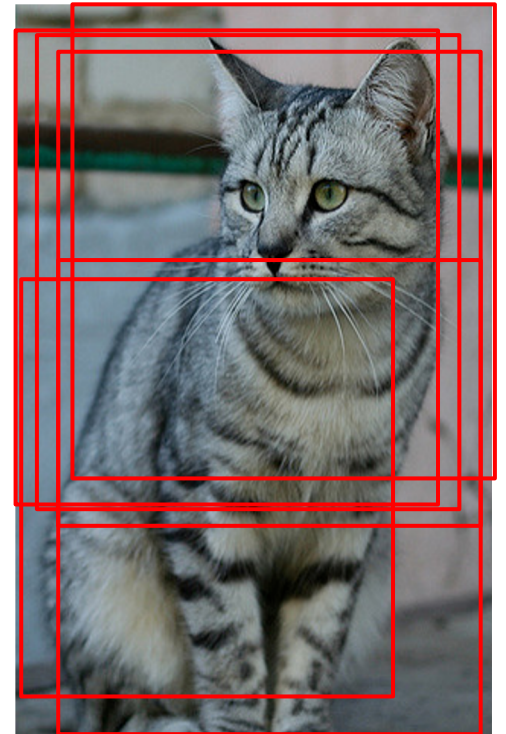


# Data Augmentation: Random Crops and Scales

**Training:** sample random crops / scales

ResNet:

1. Pick random  $L$  in range  $[256, 480]$
2. Resize training image, short side =  $L$
3. Sample random  $224 \times 224$  patch



# Regularization

**Training:** Add randomness

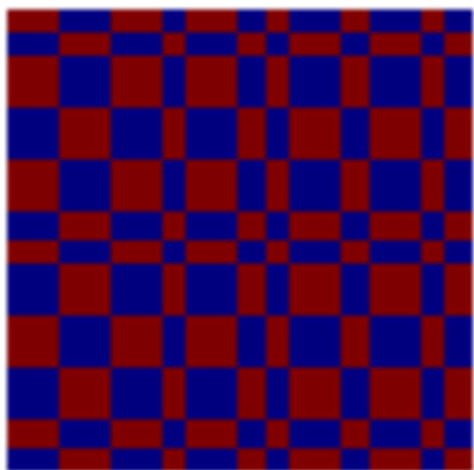
**Testing:** Marginalize out randomness

## Examples:

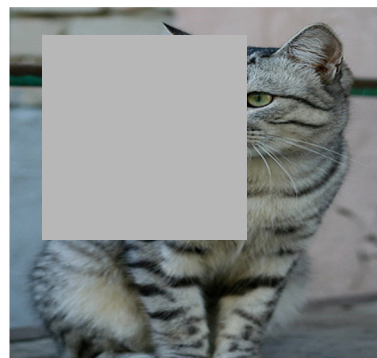
Batch Normalization

Data Augmentation

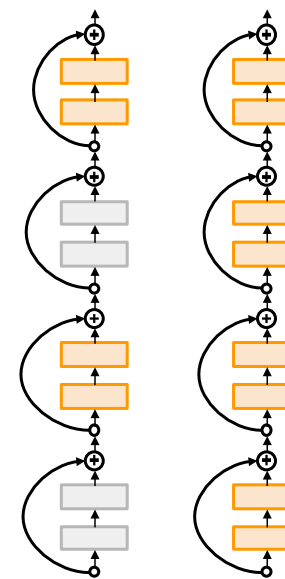
Fractional pooling



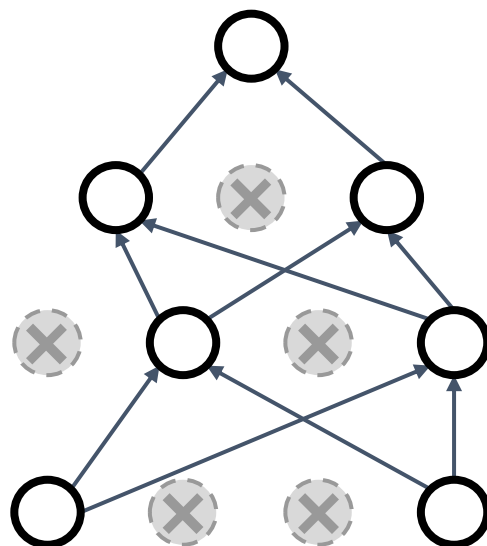
## Cutout



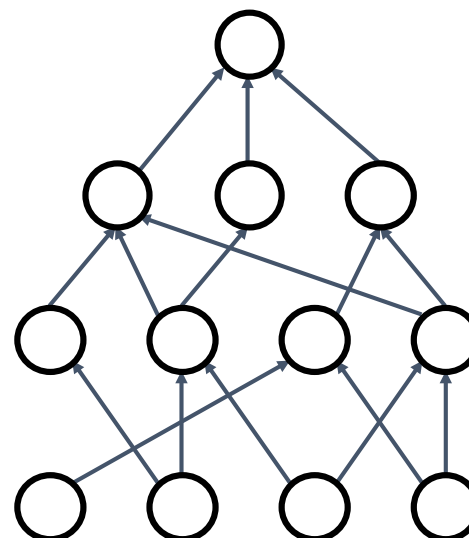
## Stochastic Depth



## Dropout



## DropConnect



## Mixup



(Old style) regularization: Add term to the loss

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

In common use:

**L2 regularization**

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

L1 regularization

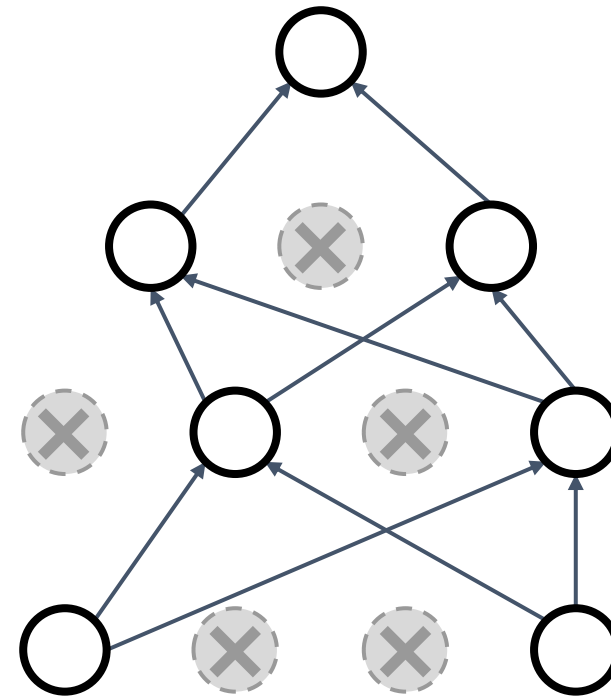
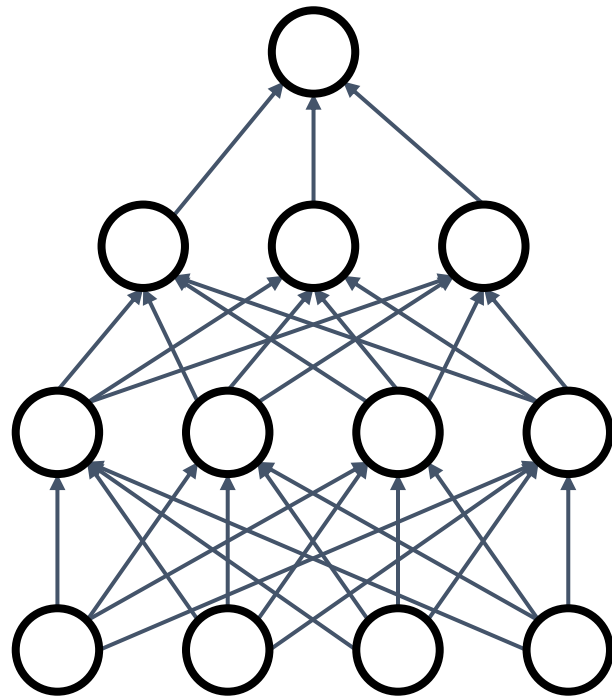
$$R(W) = \sum_k \sum_l |W_{k,l}|$$

Elastic net (L1 + L2)

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

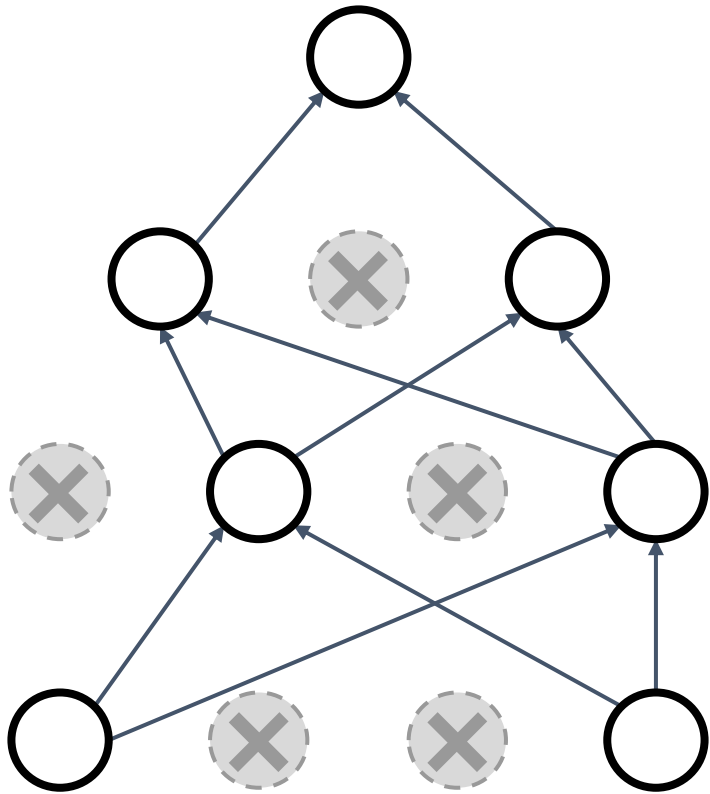
# Regularization: Dropout

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



Srivastava et al, "Dropout: A simple way to prevent neural networks from overfitting", JMLR 2014

# Regularization: Dropout



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...

# Dropout: Test Time

```
def predict(X):  
    # ensembled forward pass  
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations  
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations  
    out = np.dot(W3, H2) + b3
```

At test time all neurons are active always  
=> We must scale the activations so that for each neuron:  
output at test time = expected output at training time

# More common: “Inverted dropout”

```
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = (np.random.rand(*H1.shape) < p) / p # first dropout mask. Notice /p!
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = (np.random.rand(*H2.shape) < p) / p # second dropout mask. Notice /p!
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)

def predict(X):
    # ensembled forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) # no scaling necessary
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    out = np.dot(W3, H2) + b3
```

Drop and scale  
during training

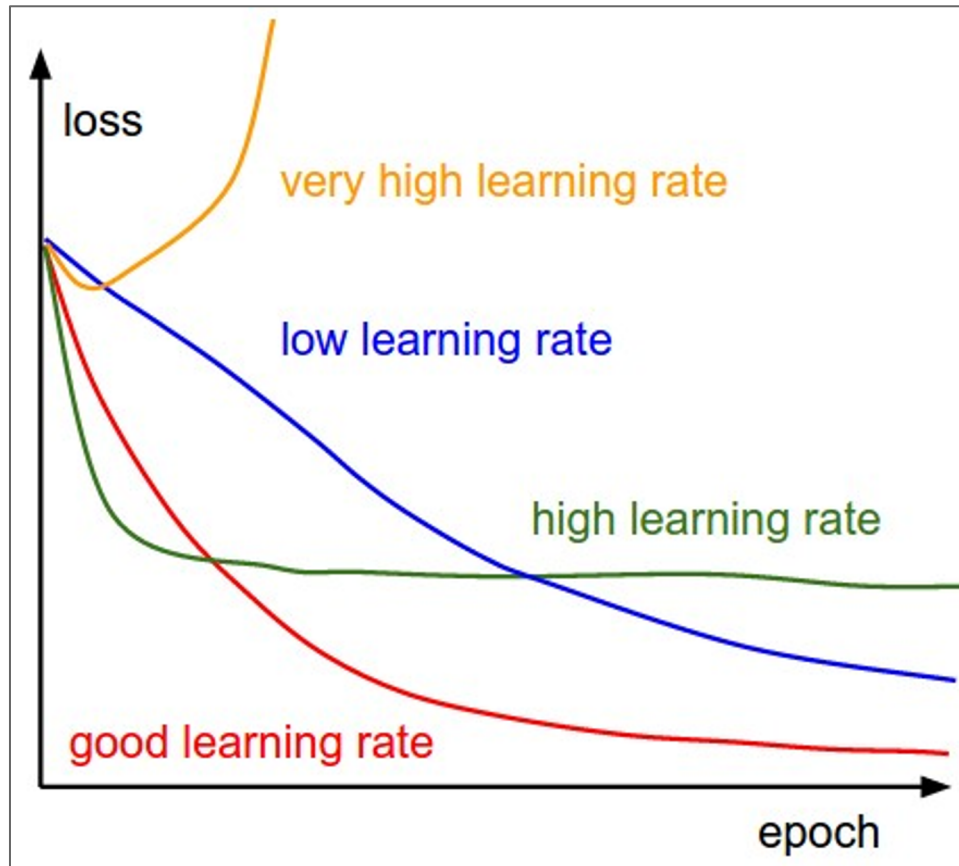
test time is unchanged!



# Learning Rate Schedules



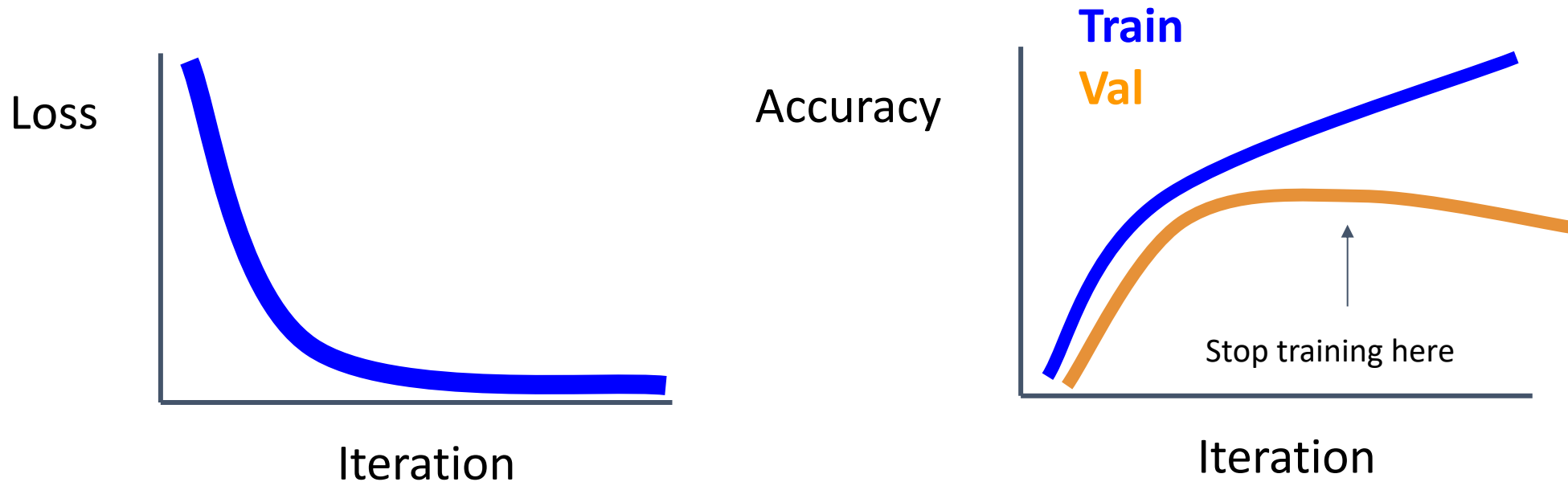
SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.



Q: Which one of these learning rates is best to use?

A: All of them! Start with large learning rate and decay over time

# How long to train? Early Stopping



Stop training the model when accuracy on the validation set decreases  
Or train for a long time, but always keep track of the model snapshot that worked best on val. **Always a good idea to do this!**

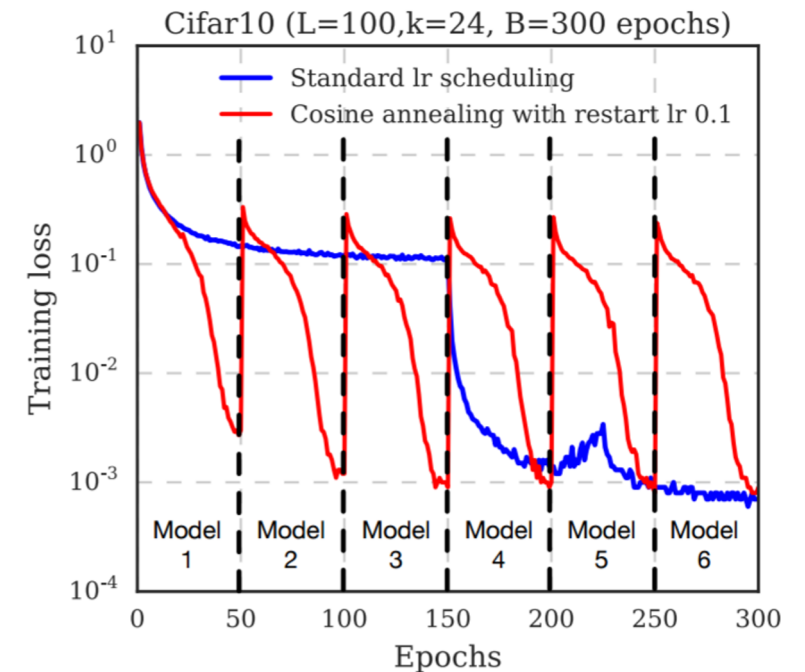
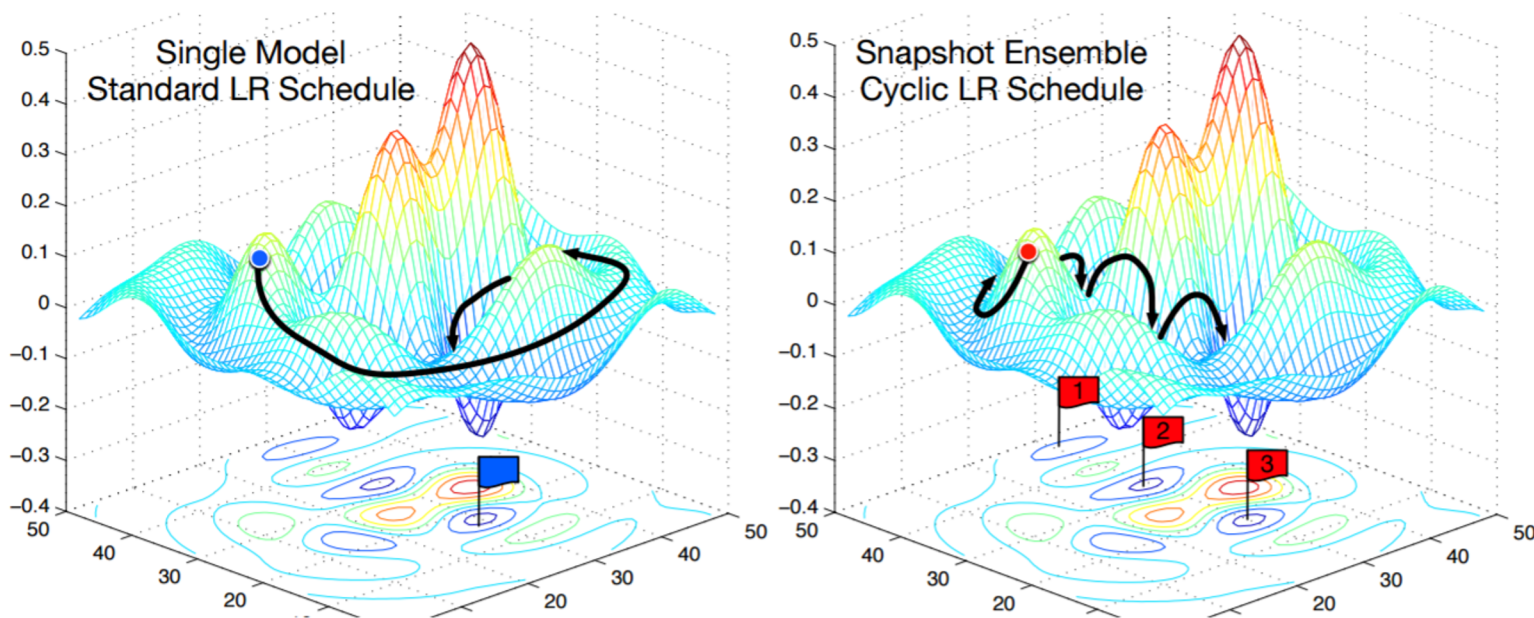
# Model Ensembles

1. Train multiple independent models
2. At test time average their results  
(Take average of predicted probability distributions, then choose argmax)

Enjoy 2% extra performance

# Model Ensembles: Tips and Tricks

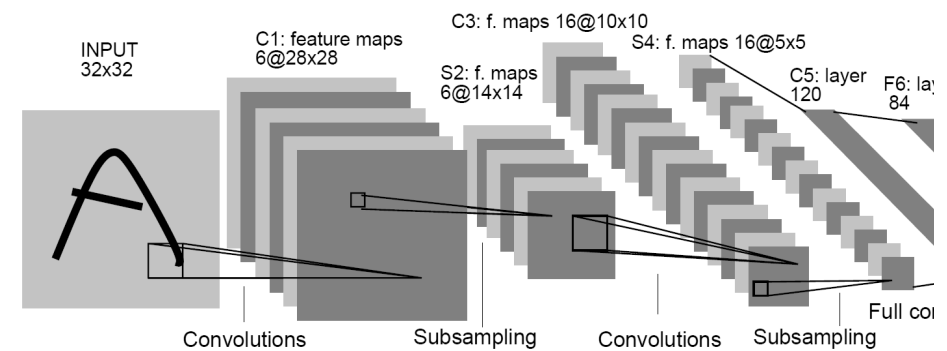
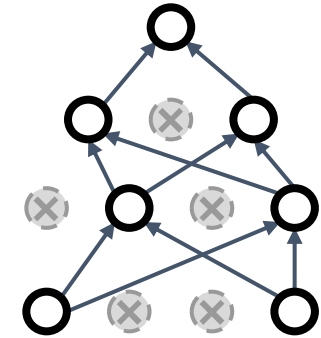
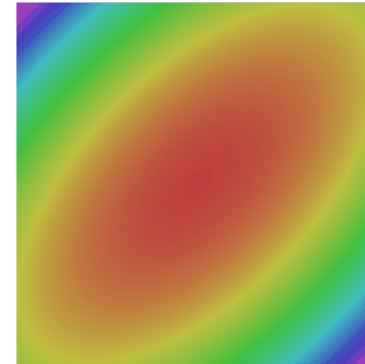
Instead of training independent models, use multiple snapshots of a single model during training!



Cyclic learning rate schedules can make this work even better!

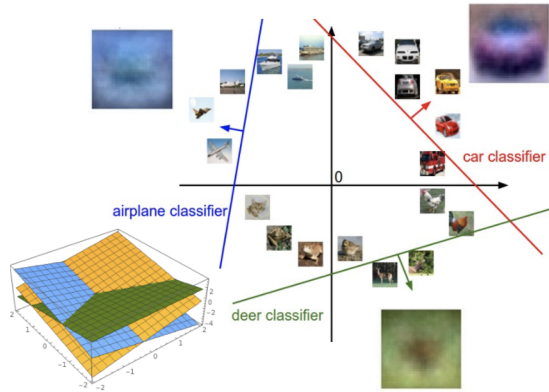
# Convolutional neural networks++

- Training and optimization
- More regularization (dropout, ...)
- Convolutional neural networks
- Pooling
- Batch normalization

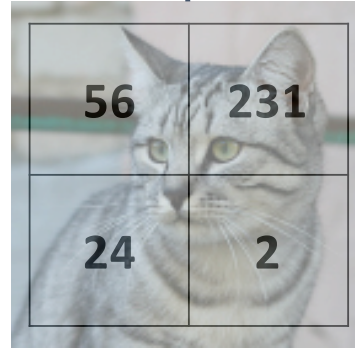


# Lecture 7: Convolutional Networks

$$f(x,W) = Wx$$



Stretch pixels into column



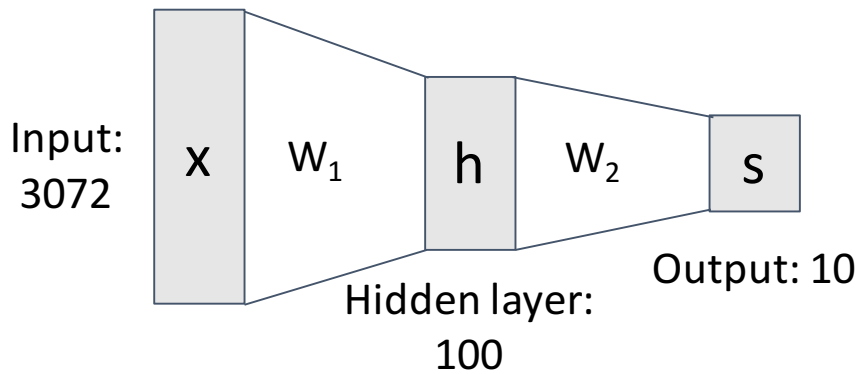
Input image  
(2, 2)

**Problem:** So far our classifiers don't respect the spatial structure of images!

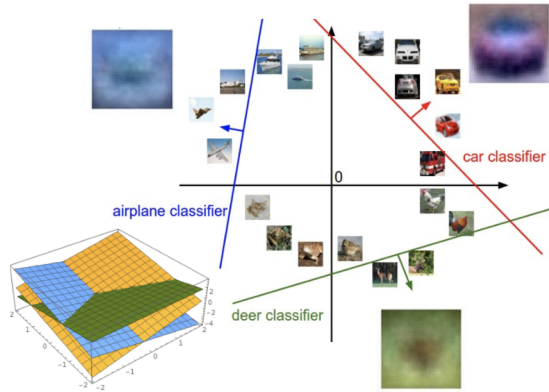
56
231
24
2

(4,)

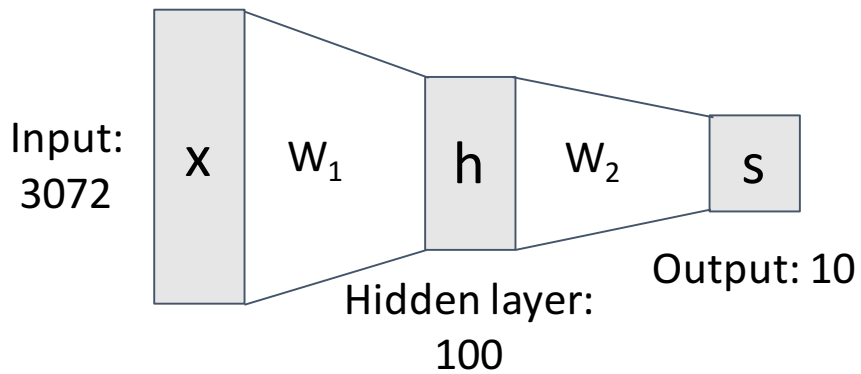
$$f = W_2 \max(0, W_1 x)$$



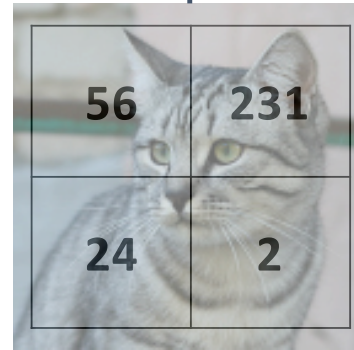
$$f(x,W) = Wx$$



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



Stretch pixels into column



Input image  
(2, 2)

**Problem:** So far our classifiers don't respect the spatial structure of images!

**Solution:** Define new computational nodes that operate on images!

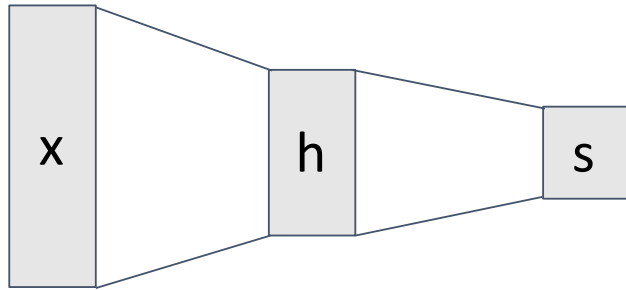
56
231
24
2

(4,)

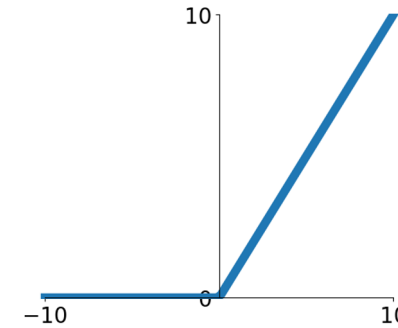


# Components of a Fully-Connected Network

Fully-Connected Layers

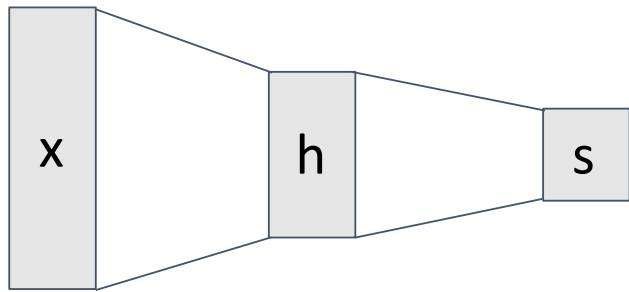


Activation Function

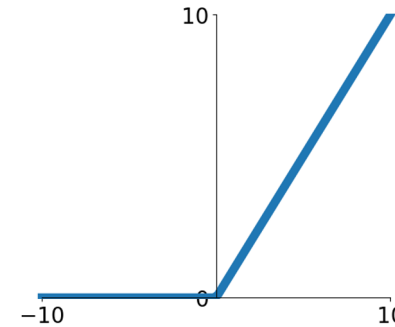


# Components of a Convolutional Network

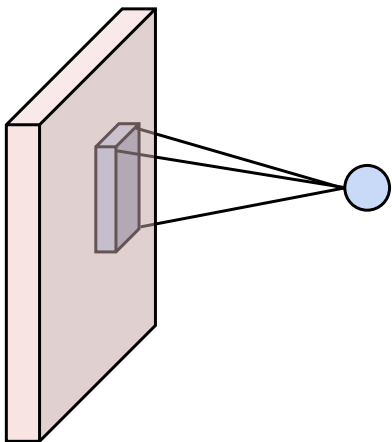
## Fully-Connected Layers



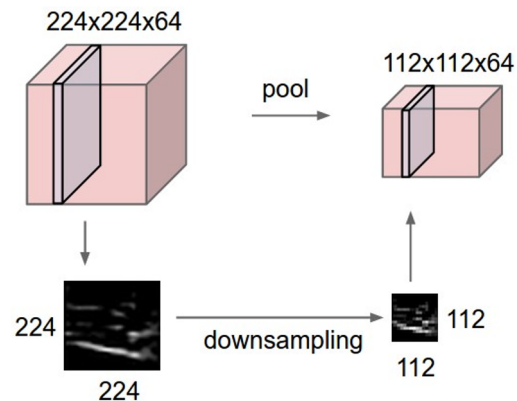
## Activation Function



## Convolution Layers



## Pooling Layers

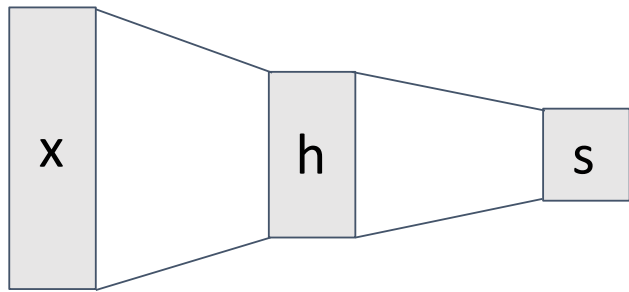


## Normalization

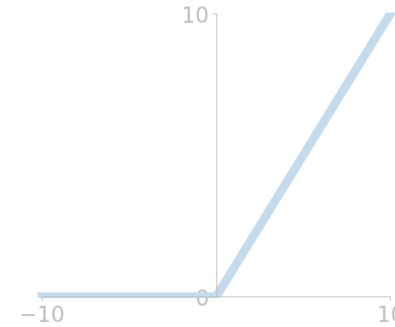
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

# Components of a Convolutional Network

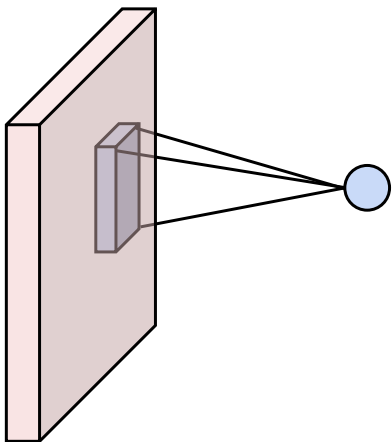
## Fully-Connected Layers



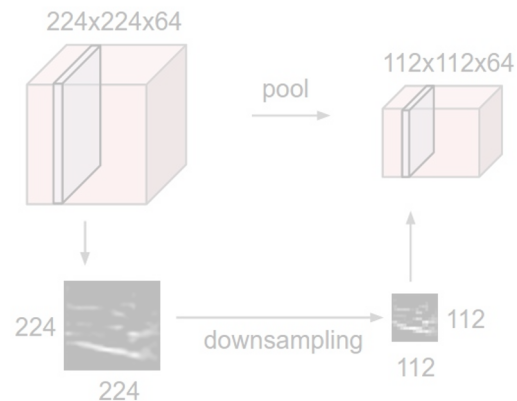
## Activation Function



## Convolution Layers



## Pooling Layers

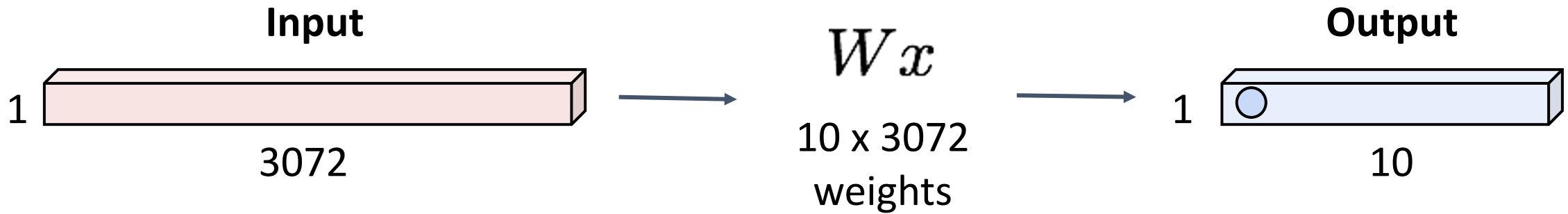


## Normalization

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

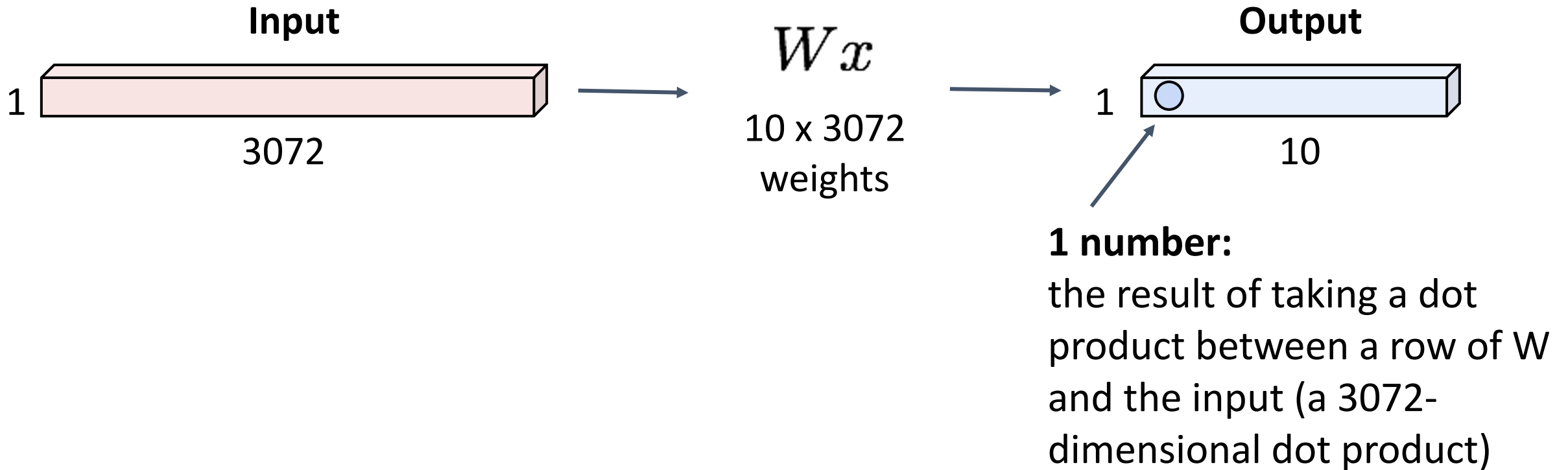
# Fully-Connected Layer

32x32x3 image -> stretch to 3072 x 1



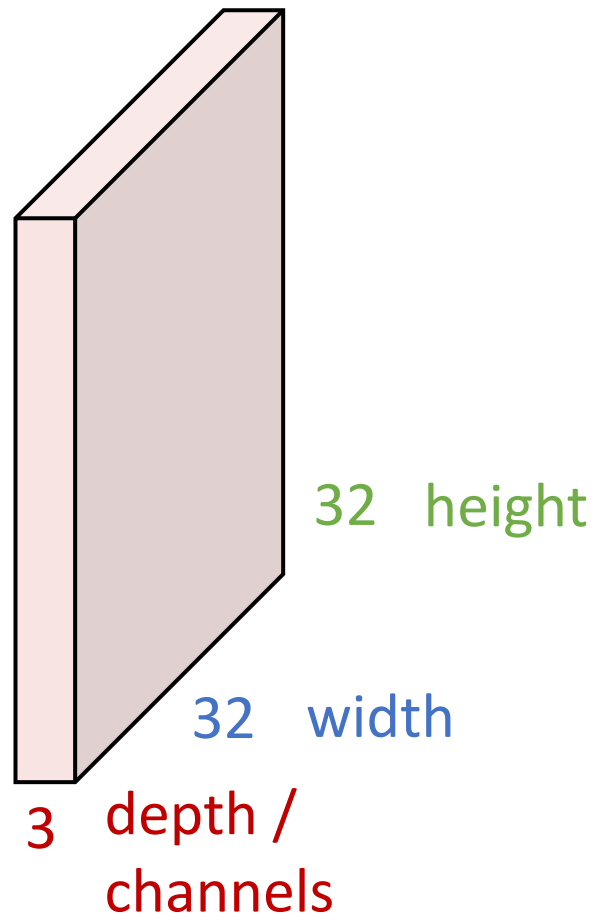
# Fully-Connected Layer

32x32x3 image -> stretch to 3072 x 1



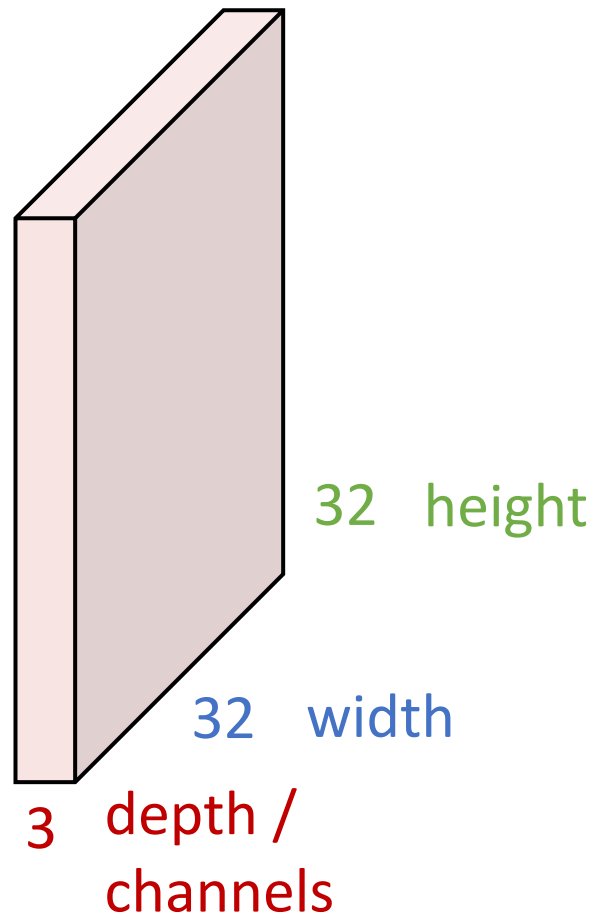
# Convolution Layer

3x32x32 image: preserve spatial structure



# Convolution Layer

3x32x32 image



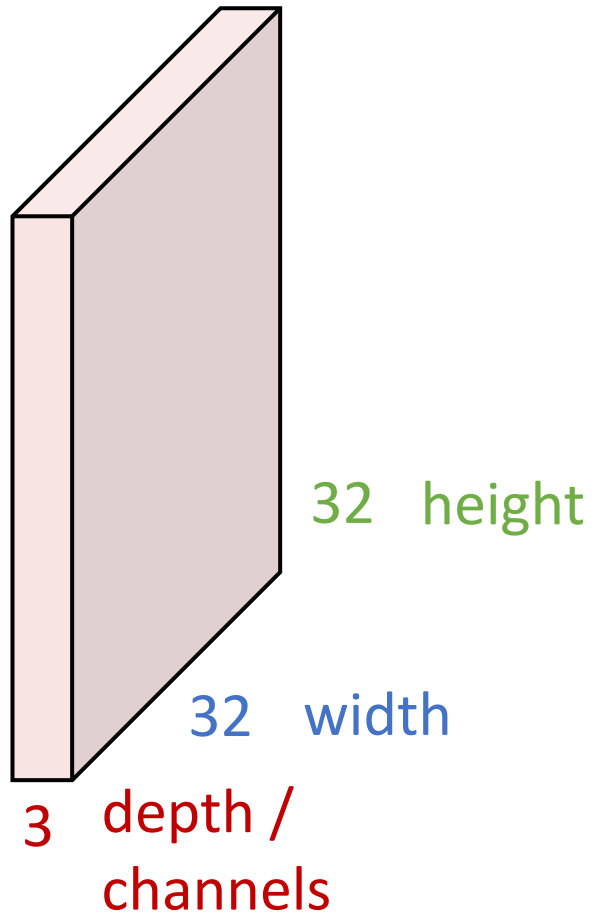
3x5x5 filter



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

# Convolution Layer

3x32x32 image



Filters (almost) always extend the full depth of the input volume

3x5x5 filter

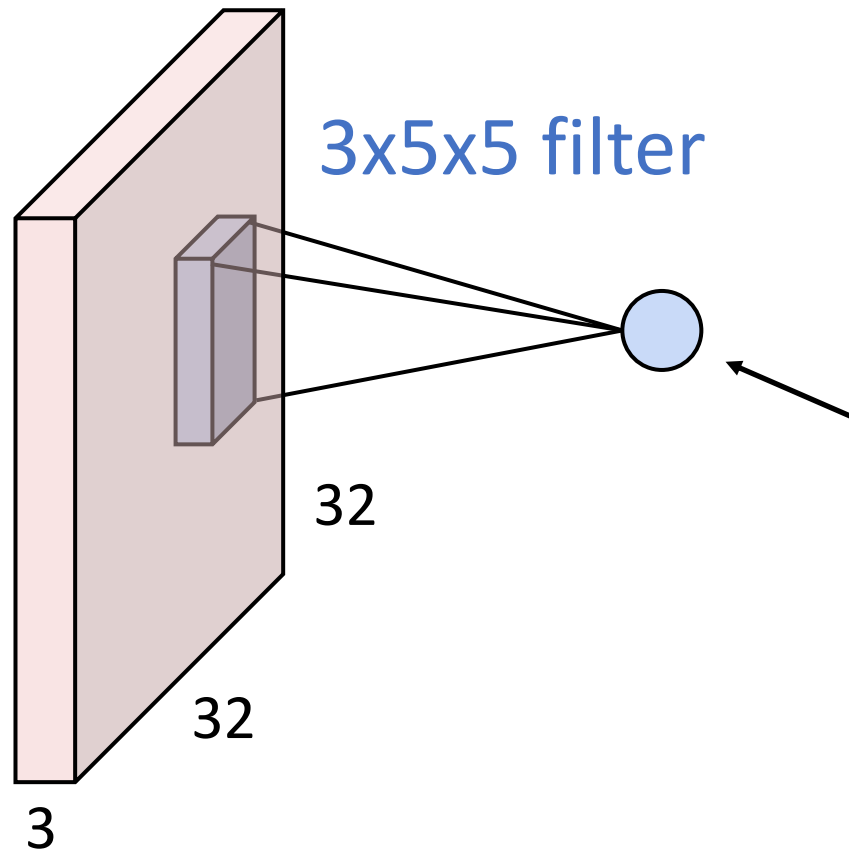


**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”



# Convolution Layer

3x32x32 image



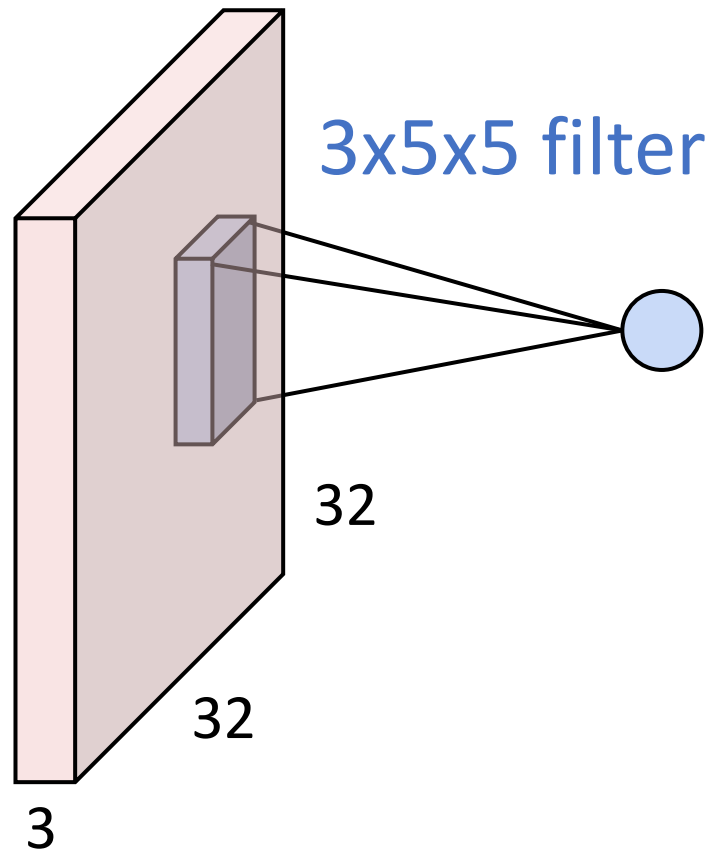
**1 number:**

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

$$w^T x + b$$

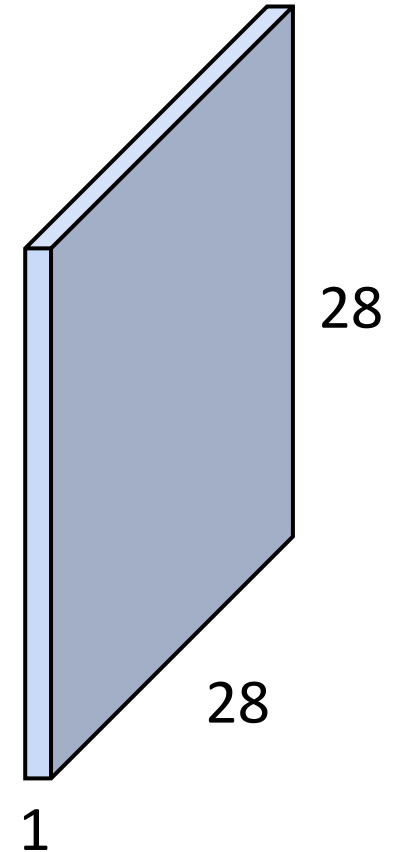
# Convolution Layer

3x32x32 image



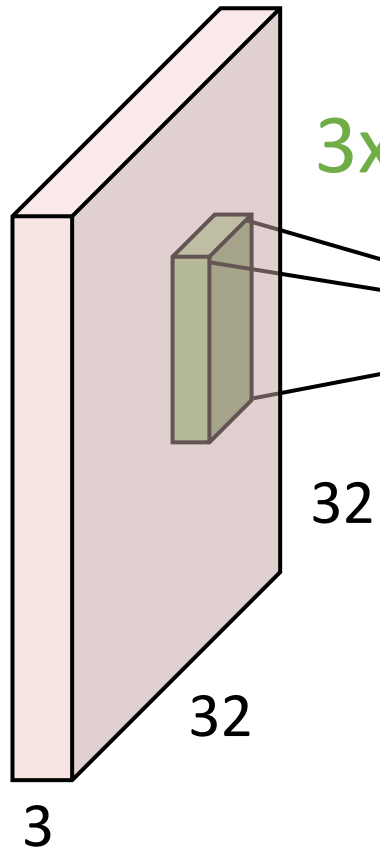
convolve (slide) over  
all spatial locations

1x28x28  
activation map

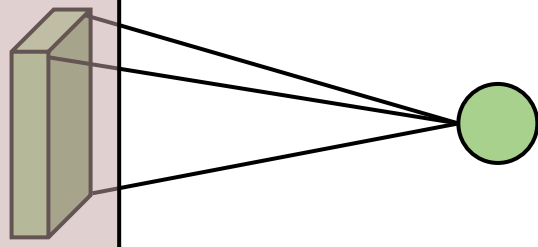


# Convolution Layer

3x32x32 image



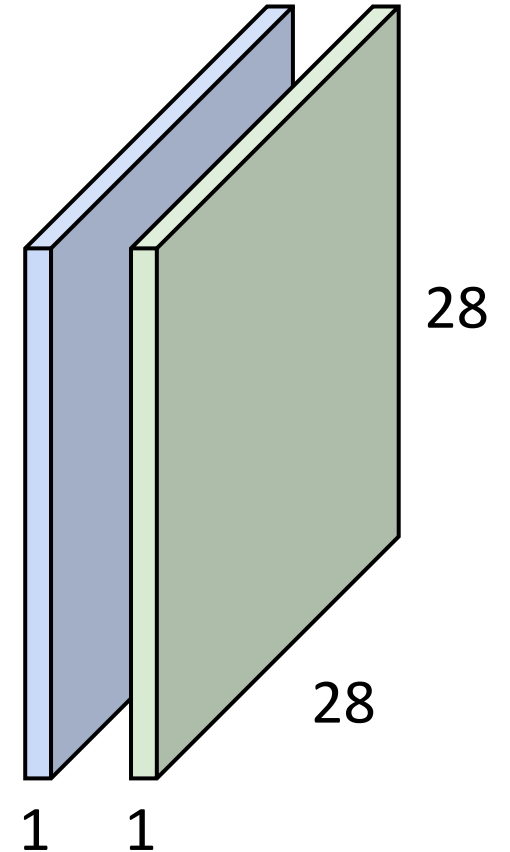
3x5x5 filter



Consider repeating with a second (green) filter:

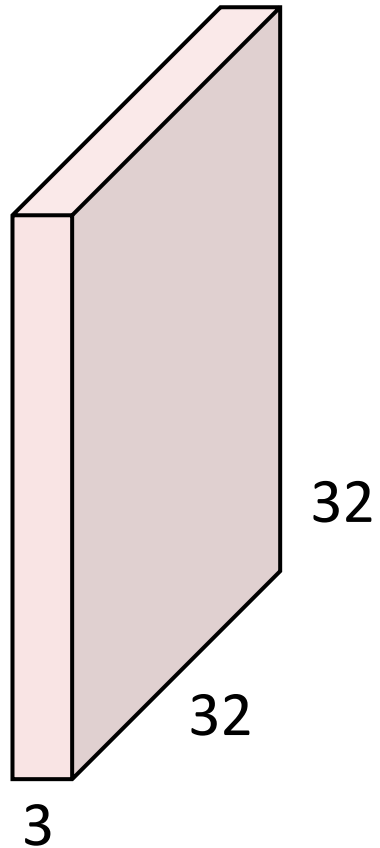
convolve (slide) over all spatial locations

two 1x28x28 activation map



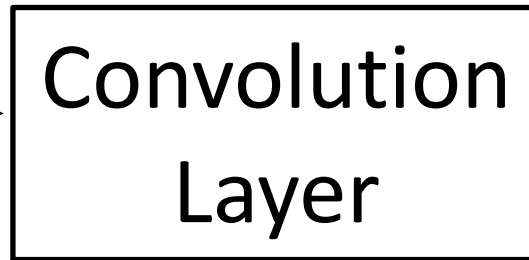
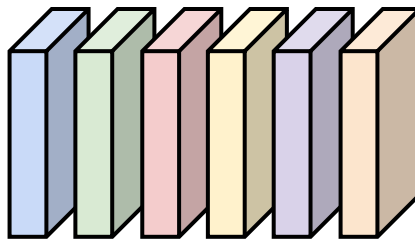
# Convolution Layer

3x32x32 image

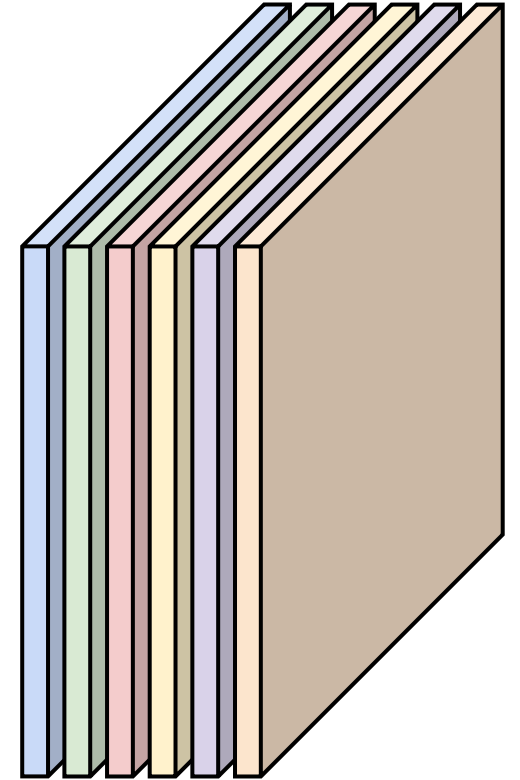


Consider 6 filters, each 3x5x5

6x3x5x5 filters



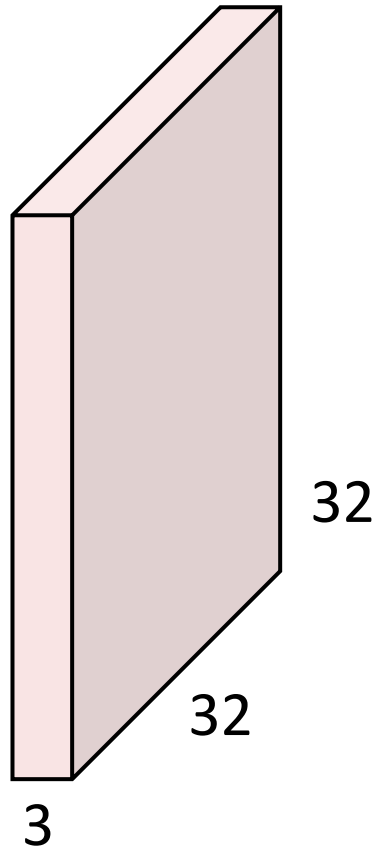
6 activation maps, each 1x28x28



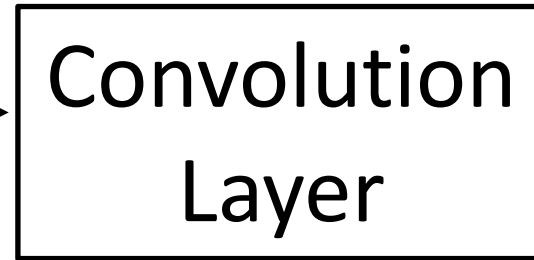
Stack activations to get a 6x28x28 output image!

# Convolution Layer

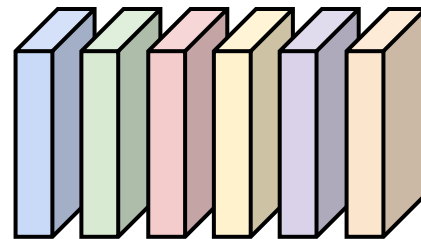
3x32x32 image



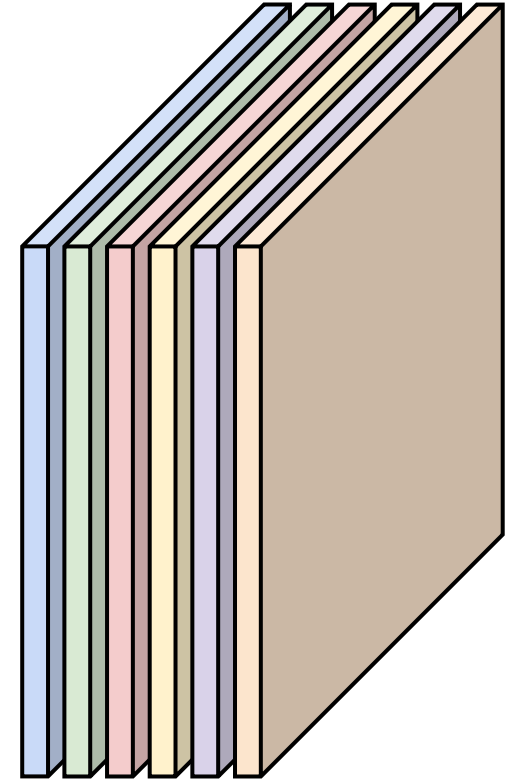
Also 6-dim bias vector:



6x3x5x5 filters



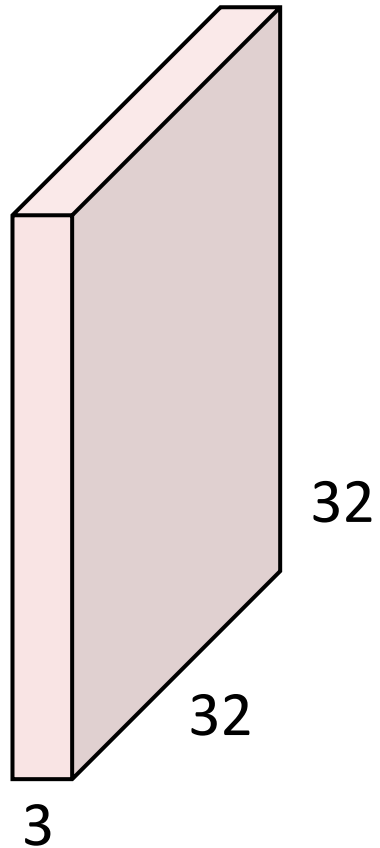
6 activation maps,  
each 1x28x28



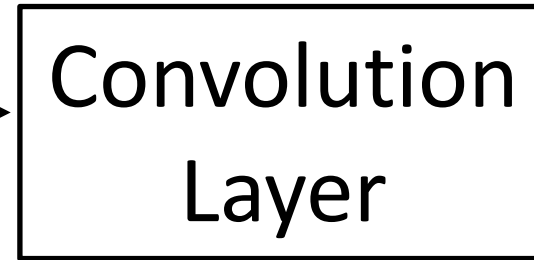
Stack activations to get a  
6x28x28 output image!

# Convolution Layer

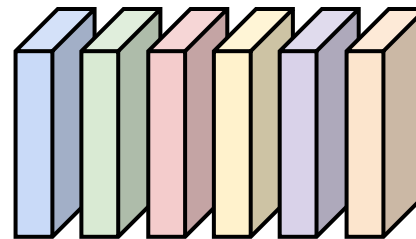
3x32x32 image



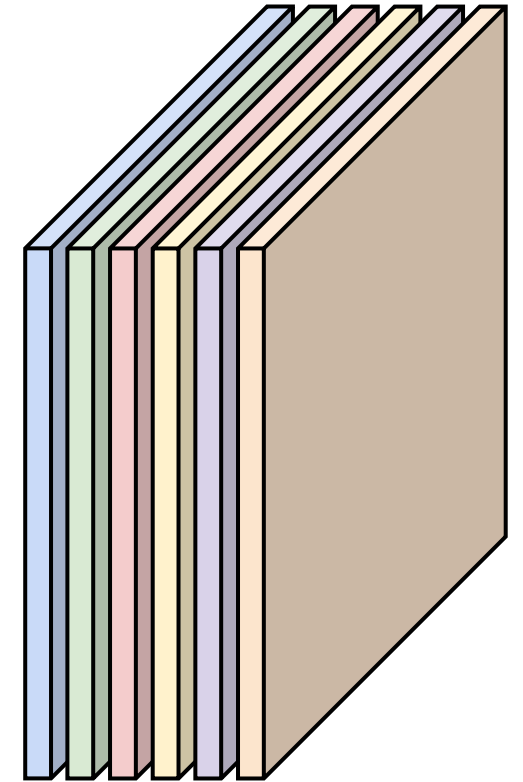
Also 6-dim bias vector:



6x3x5x5 filters



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

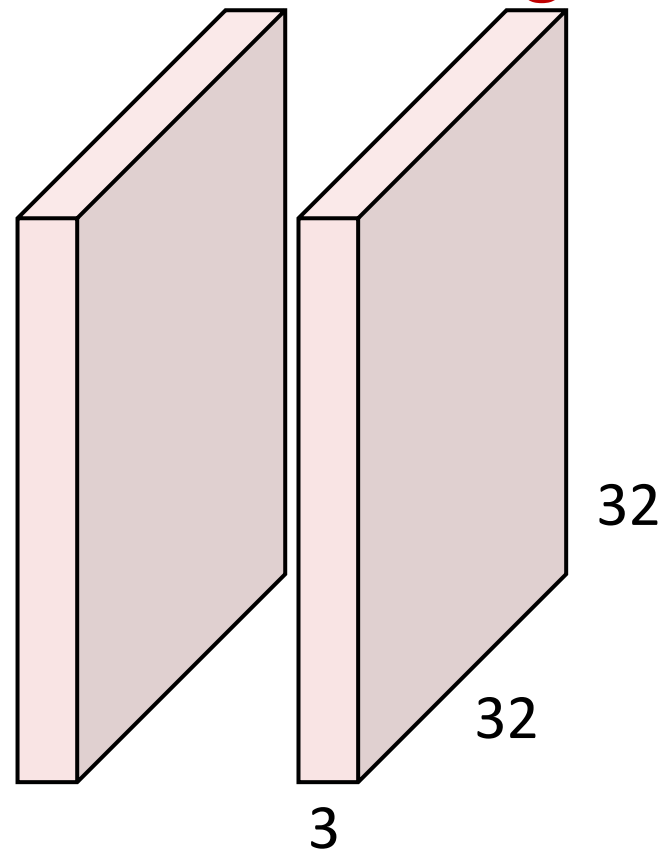


Stack activations to get a 6x28x28 output image!

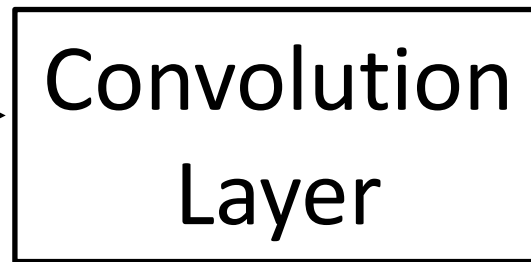
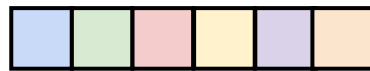
# Convolution Layer

2x3x32x32

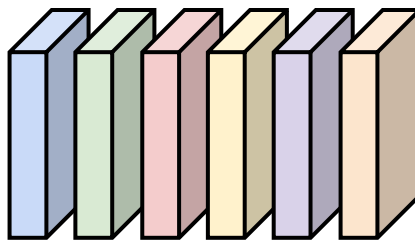
Batch of images



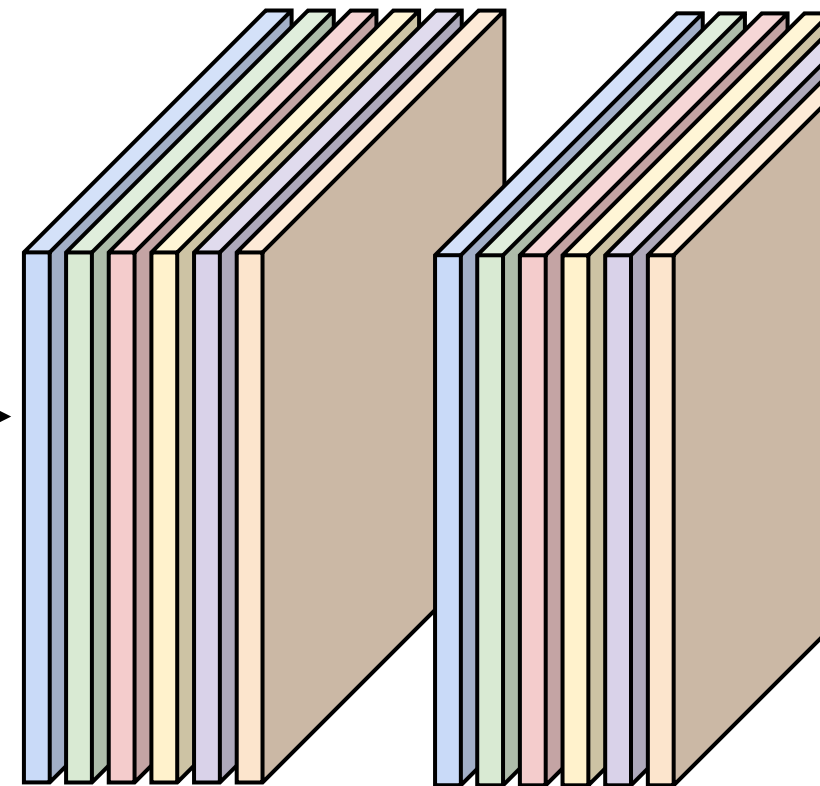
Also 6-dim bias vector:



6x3x5x5  
filters

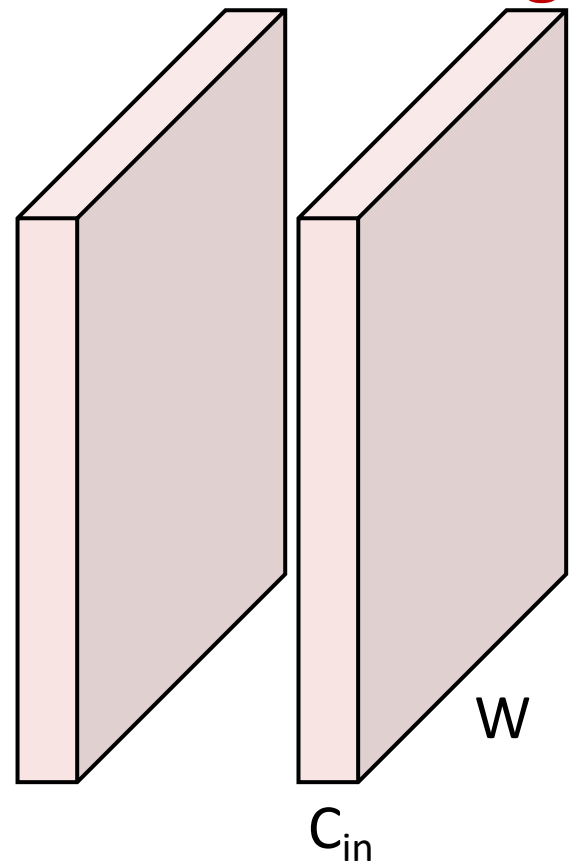


2x6x28x28  
Batch of outputs

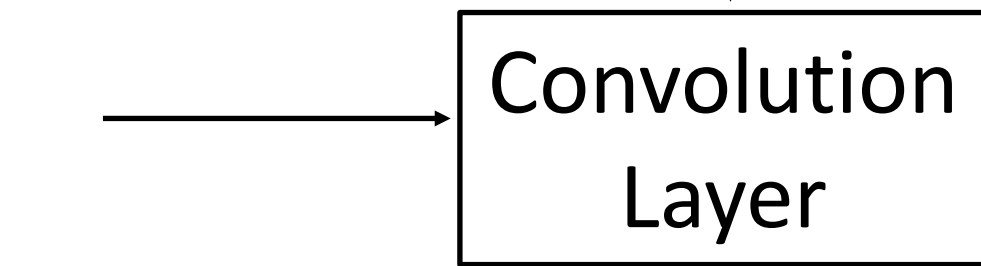
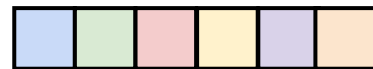


# Convolution Layer

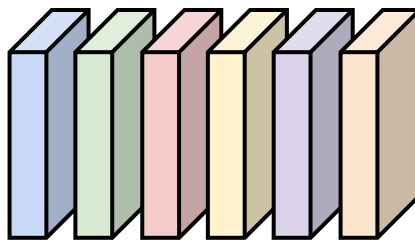
$N \times C_{in} \times H \times W$   
Batch of images



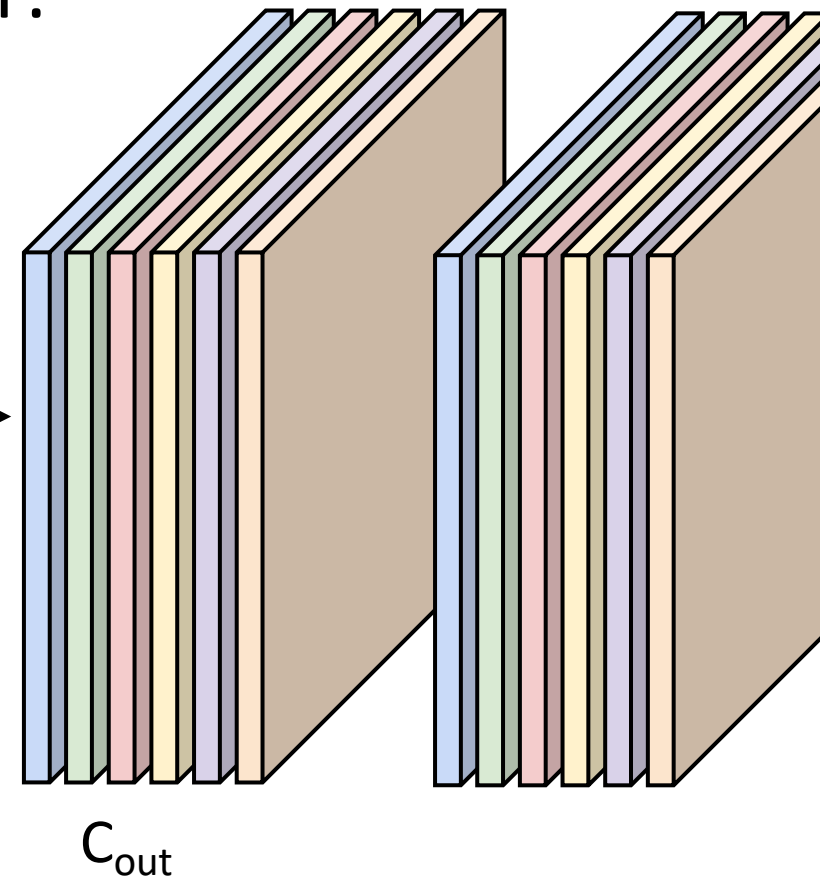
Also  $C_{out}$ -dim bias vector:



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

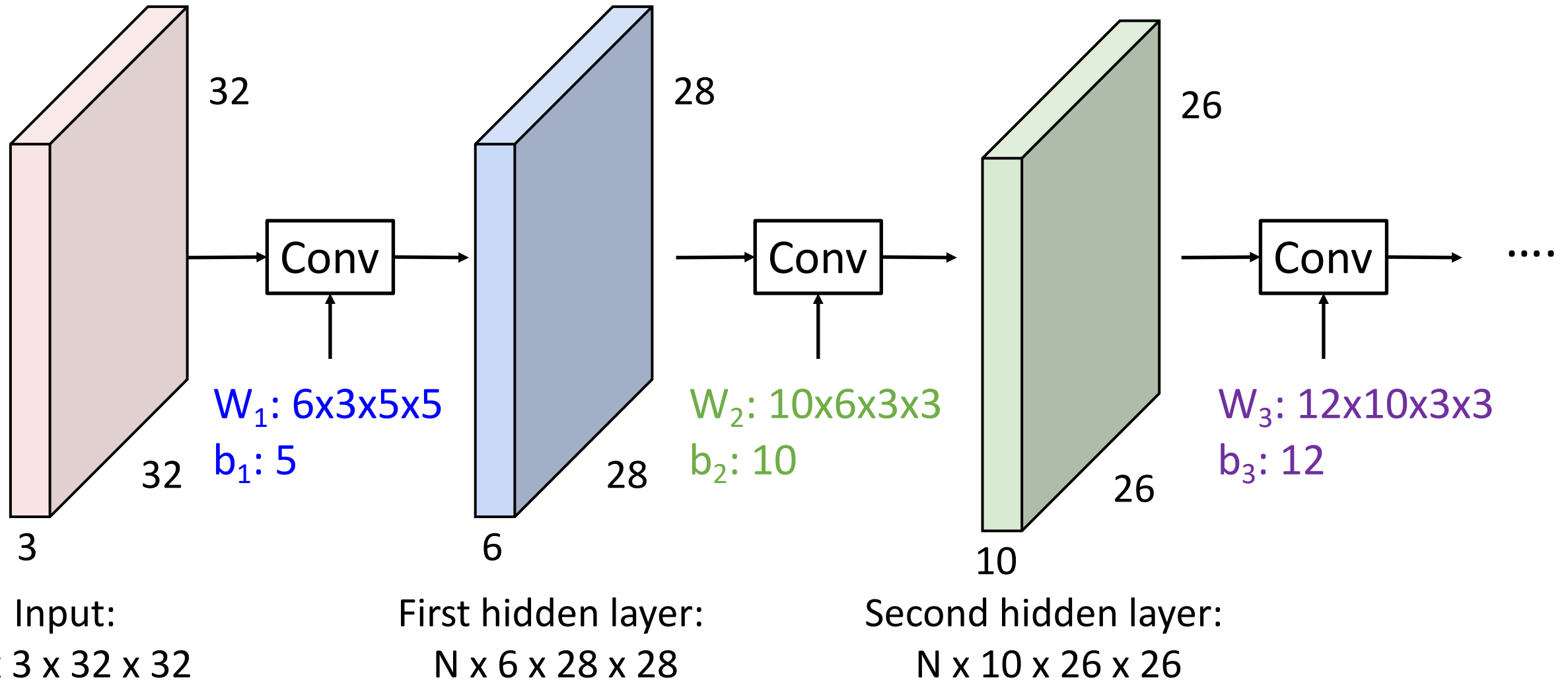


$N \times C_{out} \times H' \times W'$   
Batch of outputs





# Stacking Convolutions

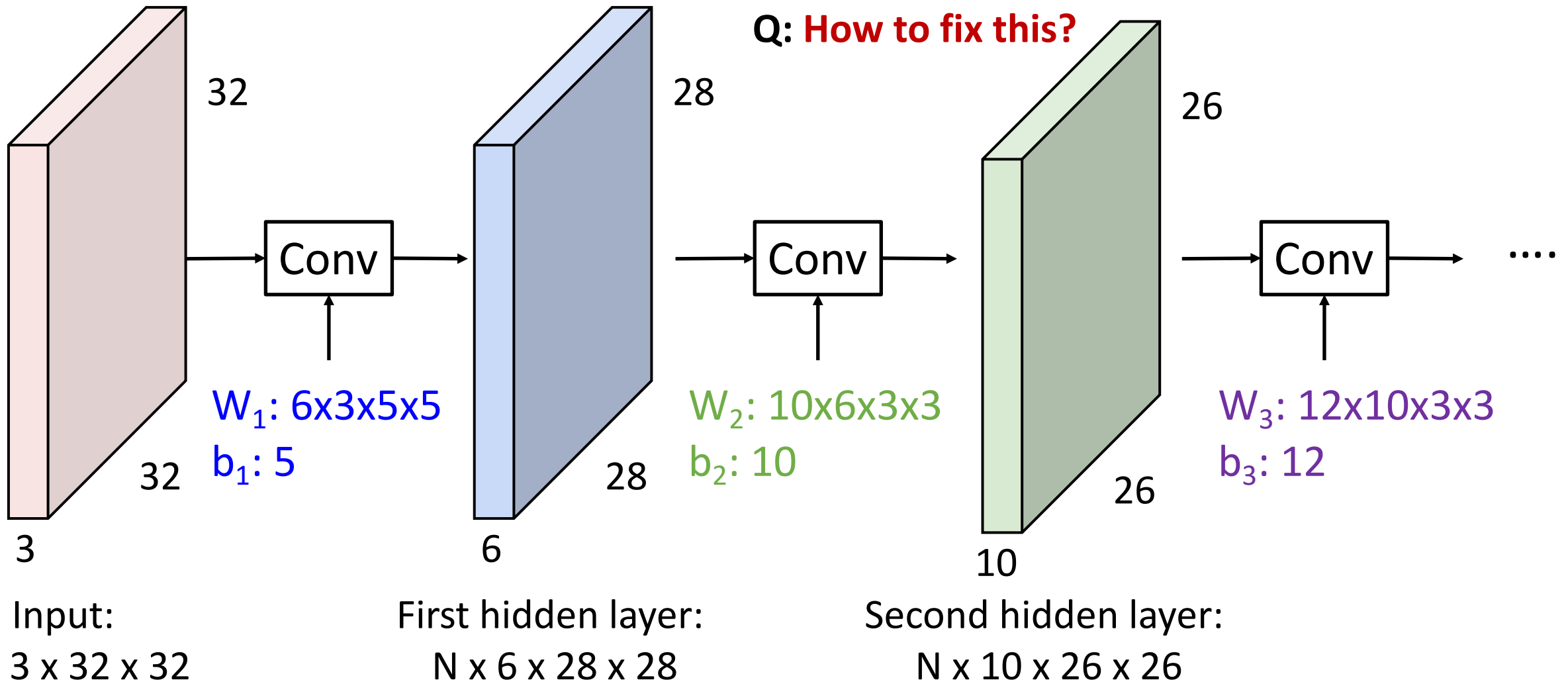


# Stacking Convolutions

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

A: We get another convolution!

Q: **How to fix this?**



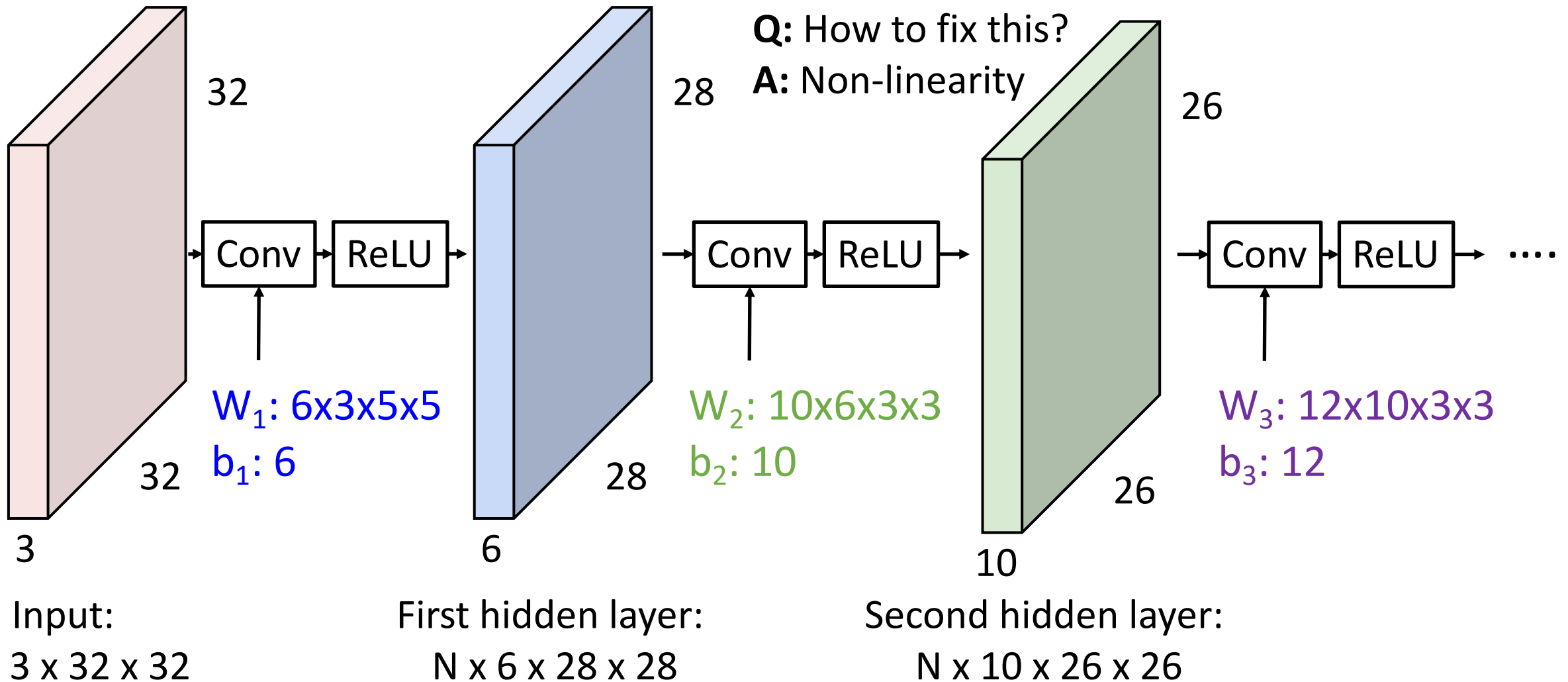
# Stacking Convolutions

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

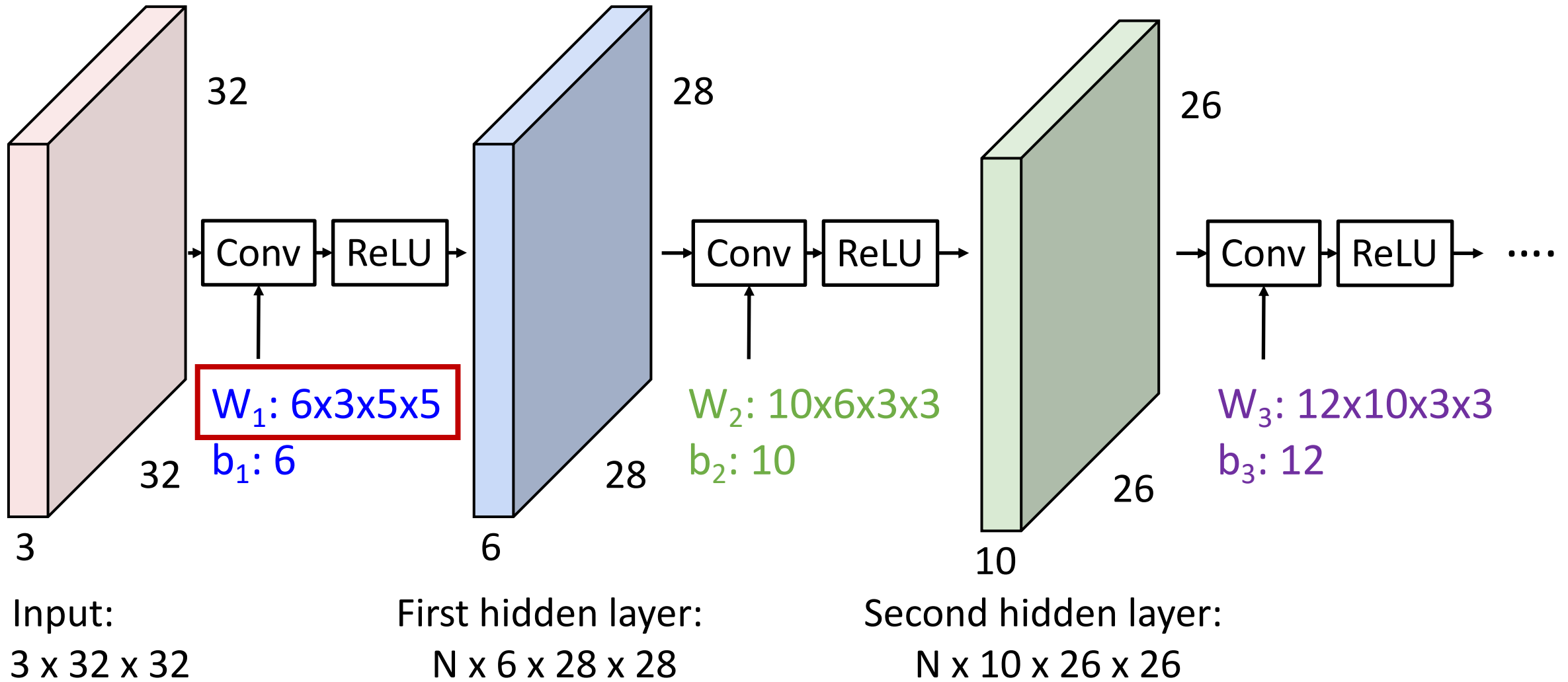
**A:** We get another convolution!

**Q:** How to fix this?

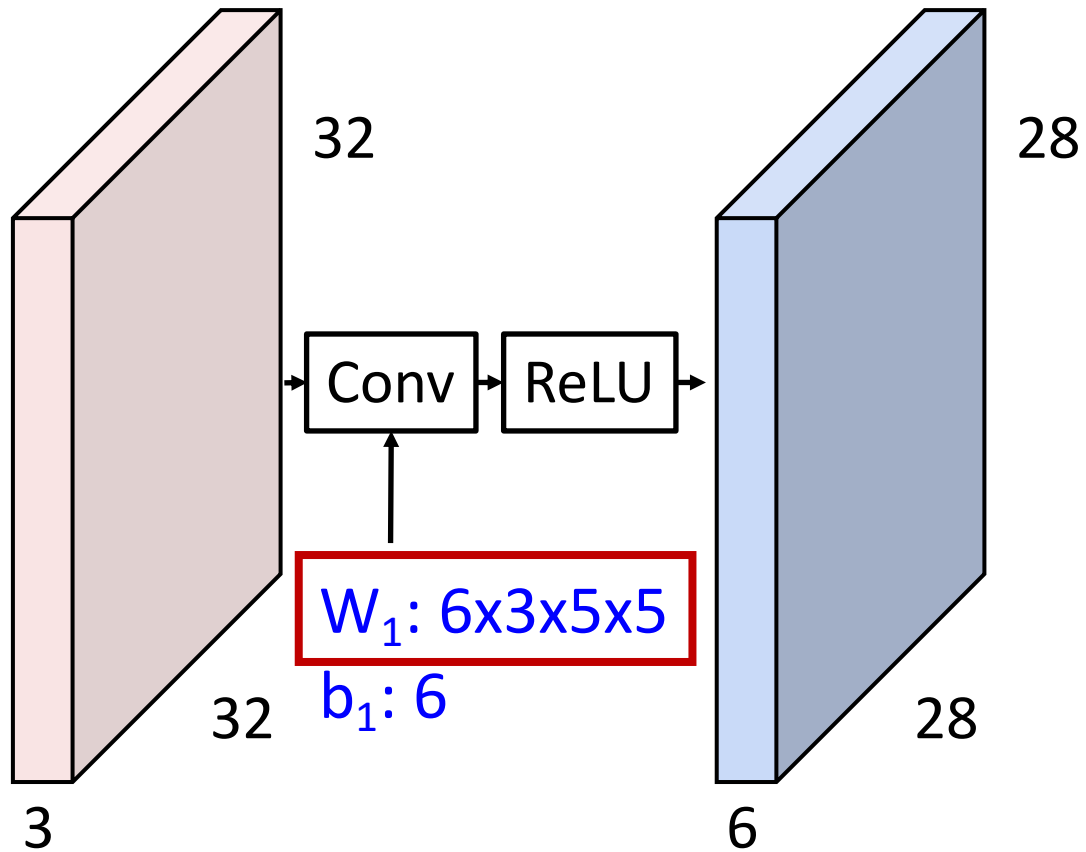
**A:** Non-linearity



# What do convolutional filters learn?



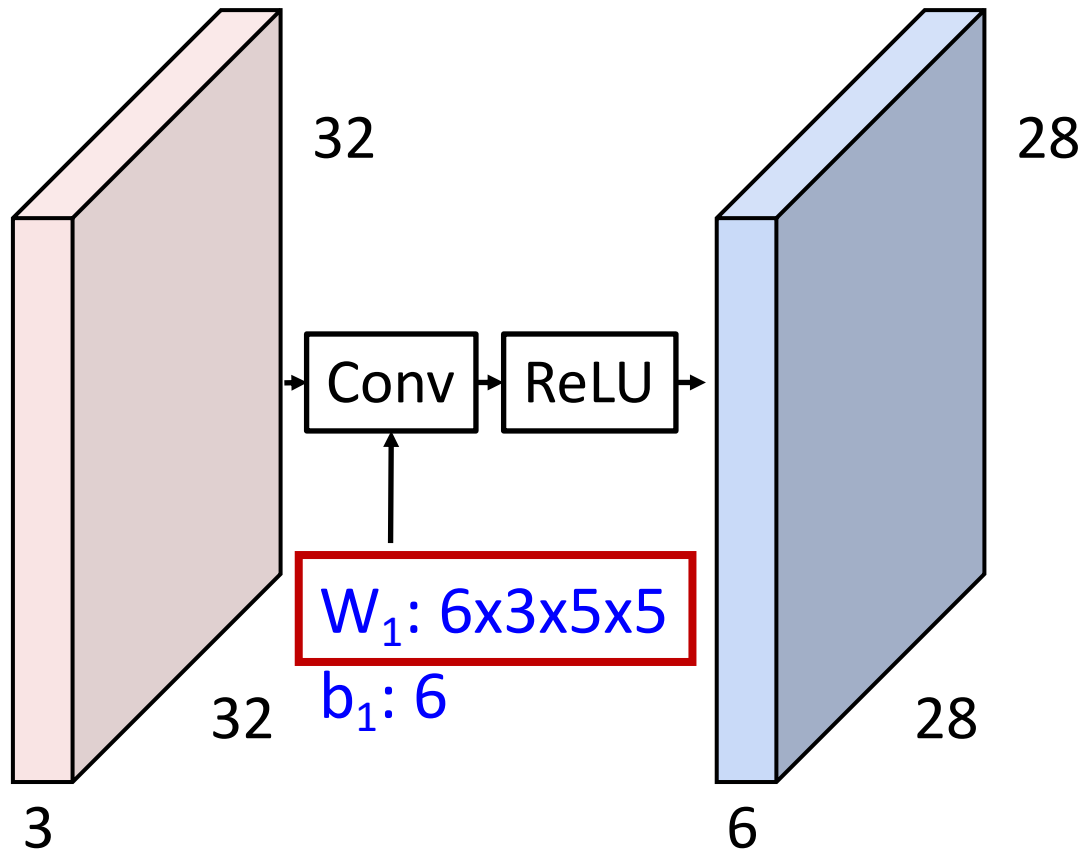
# What do convolutional filters learn?



Linear classifier: One template per class



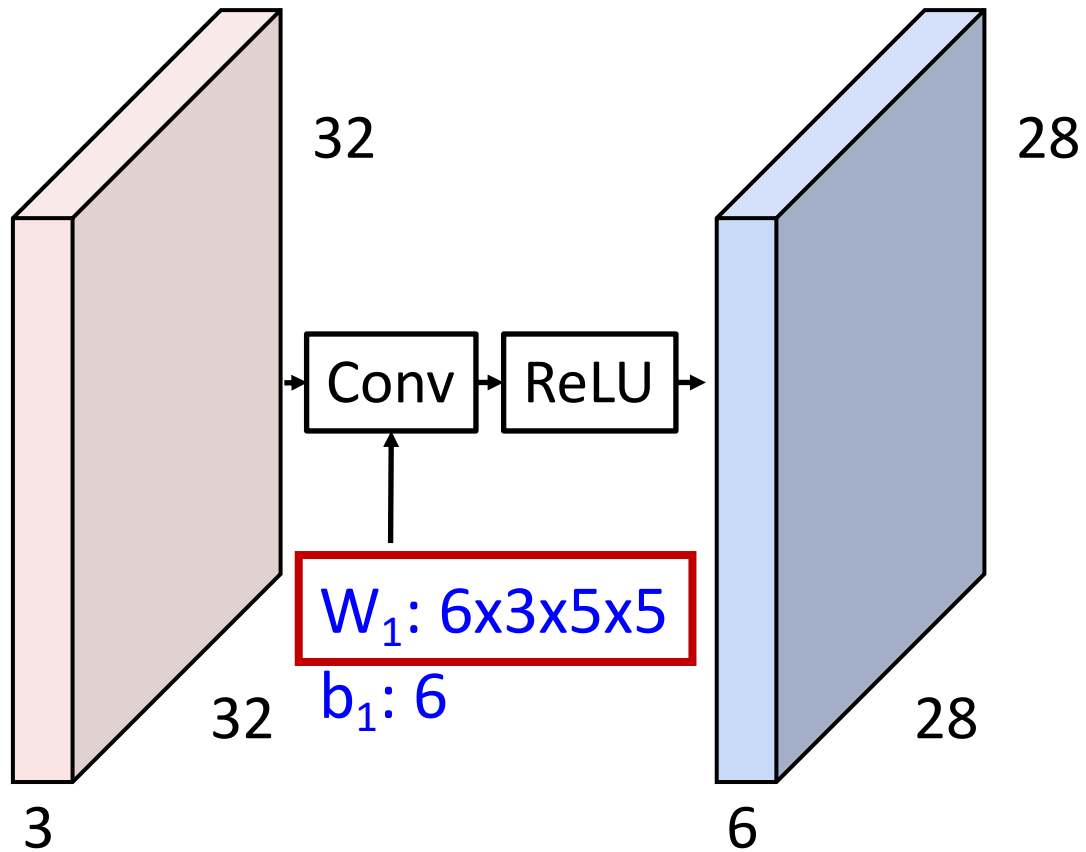
# What do convolutional filters learn?



MLP: Bank of whole-image templates



# What do convolutional filters learn?



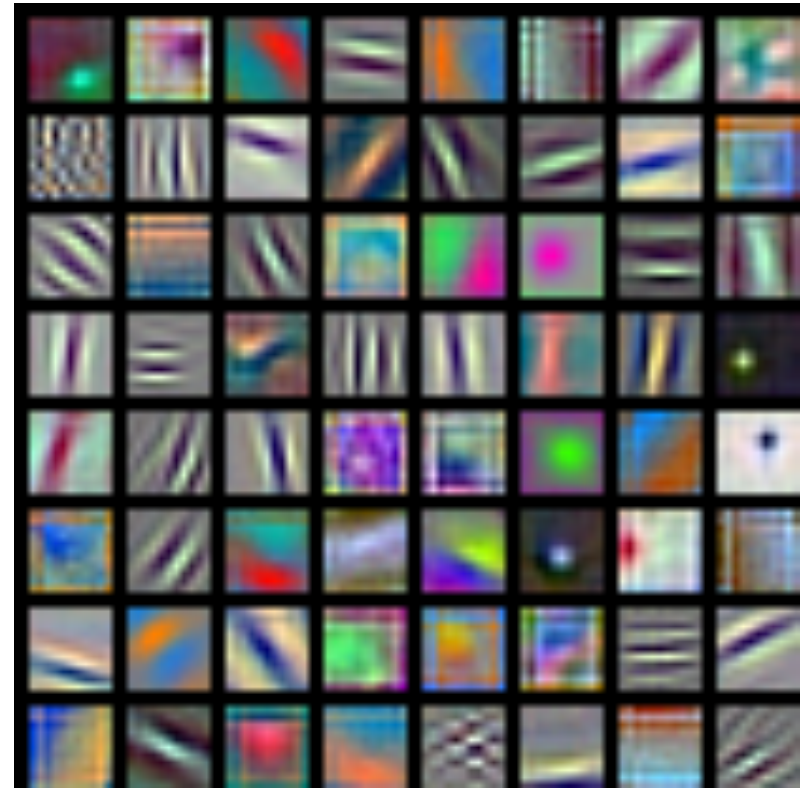
Input:

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

First hidden layer:

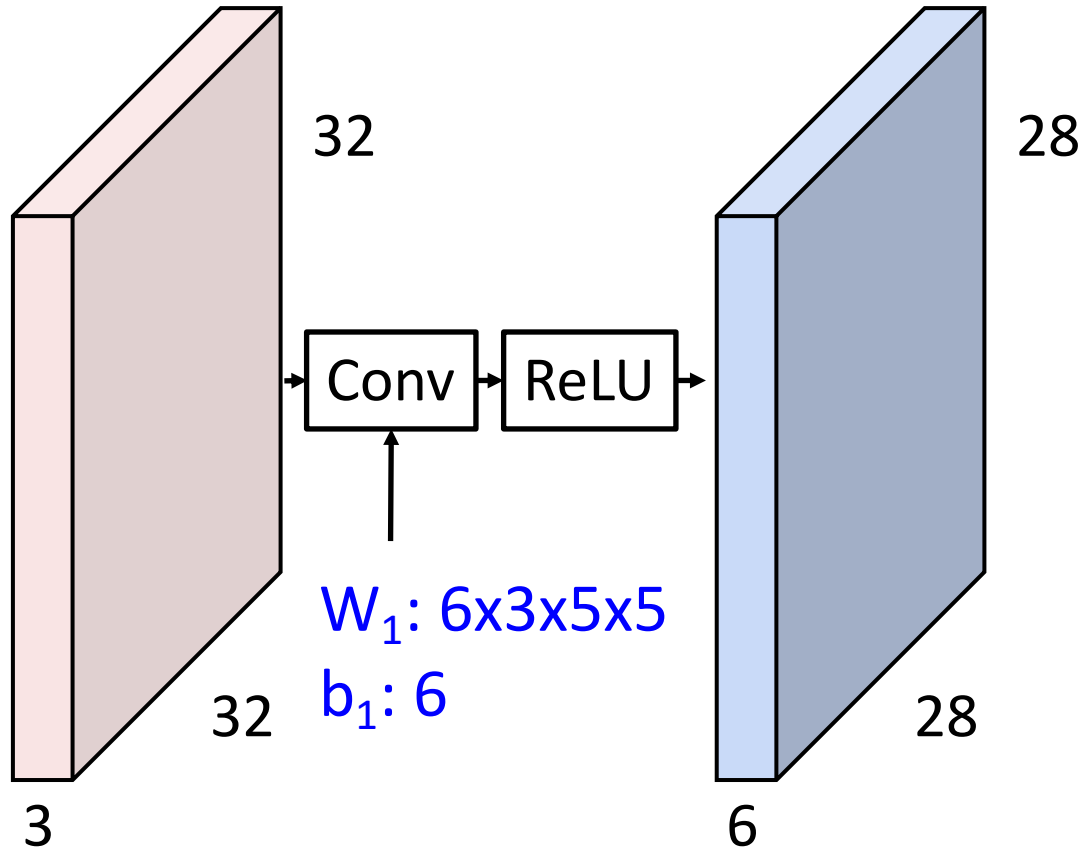
$N \times 6 \times 28 \times 28$

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



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

# A closer look at spatial dimensions



Input:

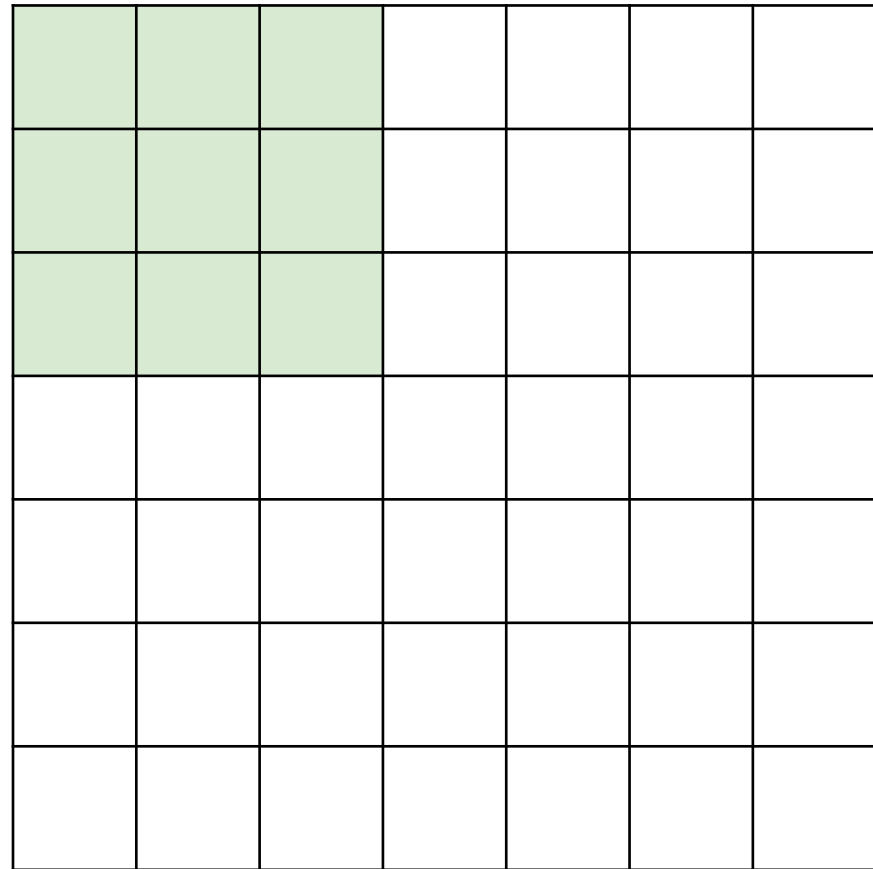
$N \times 3 \times 32 \times 32$

First hidden layer:

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



# A closer look at spatial dimensions



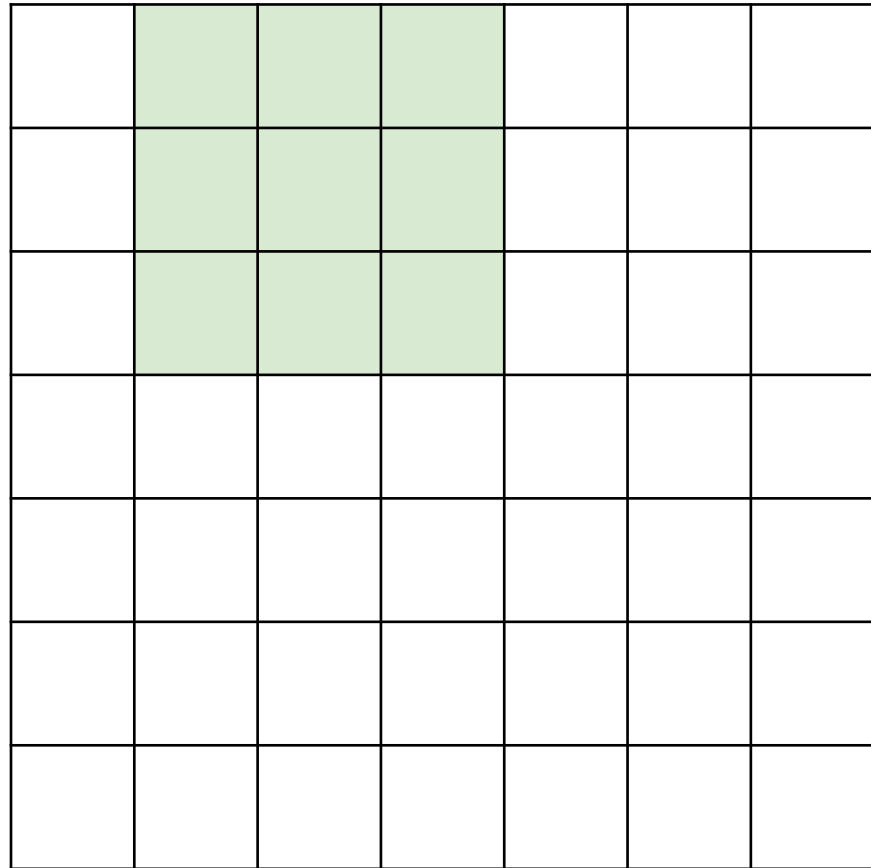
Input: 7x7

Filter: 3x3

7

7

# A closer look at spatial dimensions



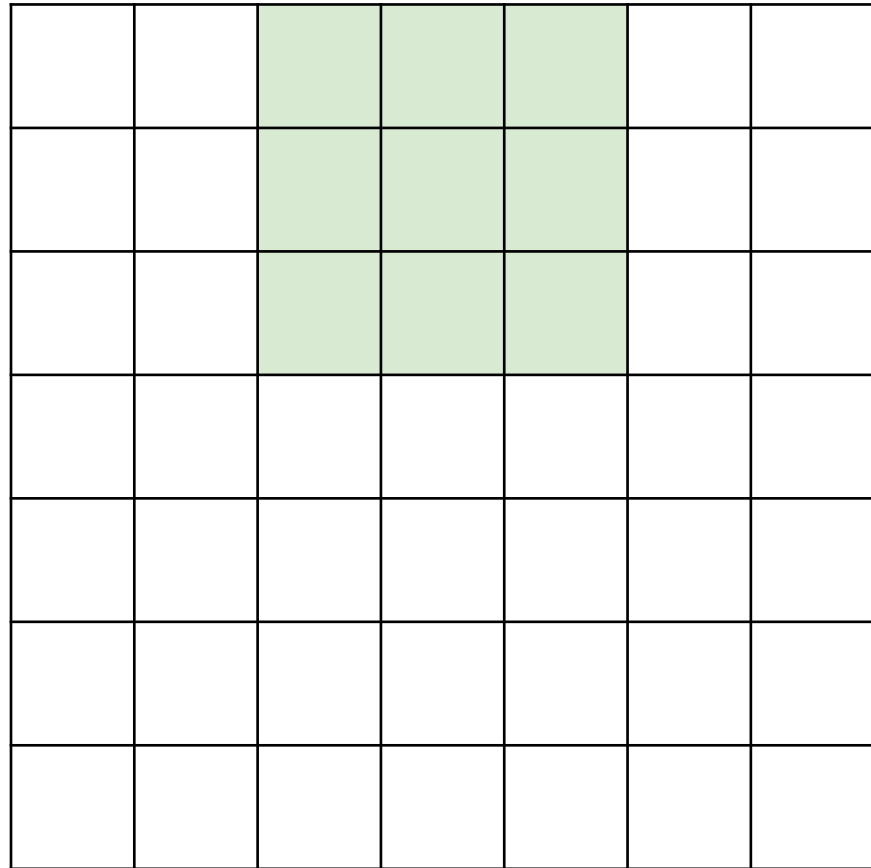
Input: 7x7

Filter: 3x3

7

7

# A closer look at spatial dimensions



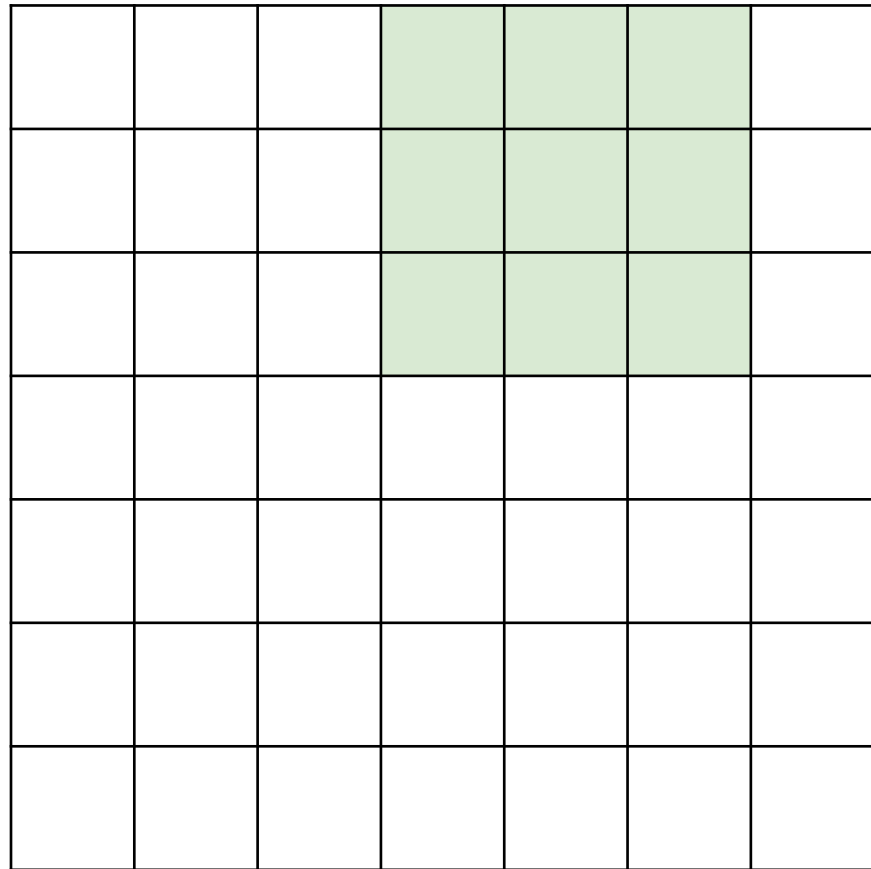
7

7

Input: 7x7

Filter: 3x3

# A closer look at spatial dimensions



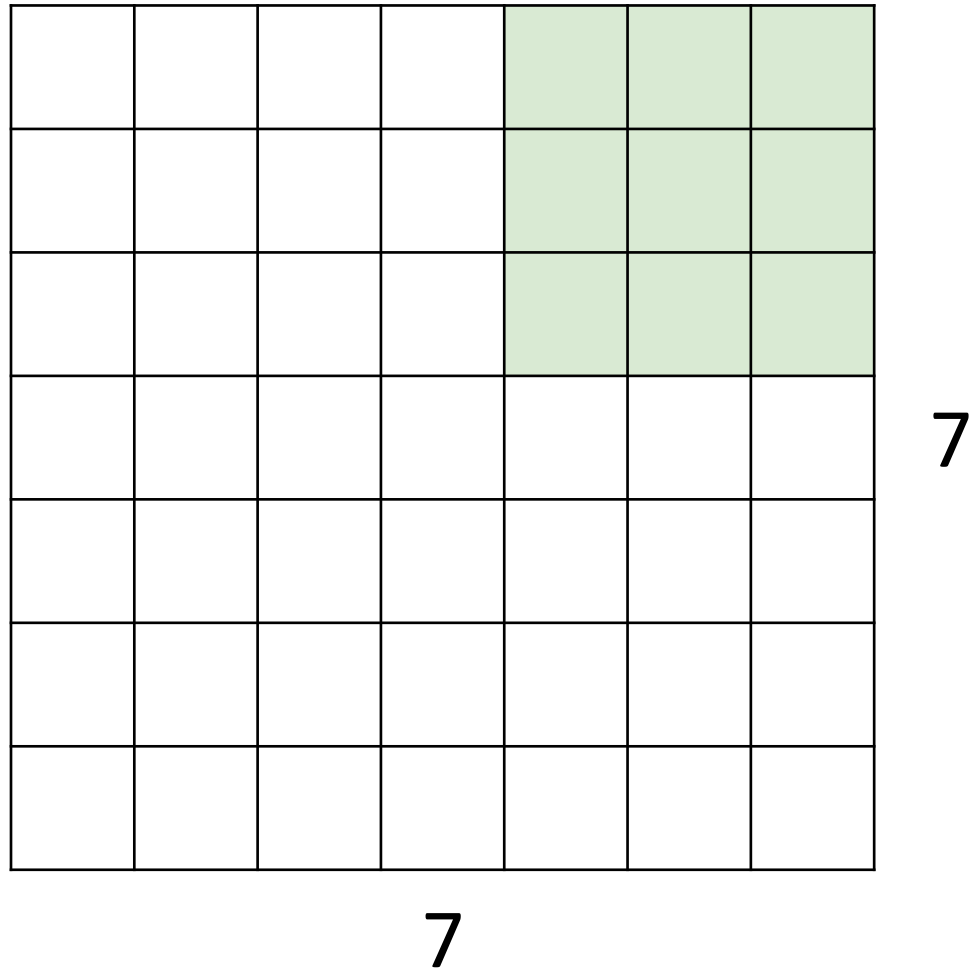
Input: 7x7

Filter: 3x3

7

7

# A closer look at spatial dimensions

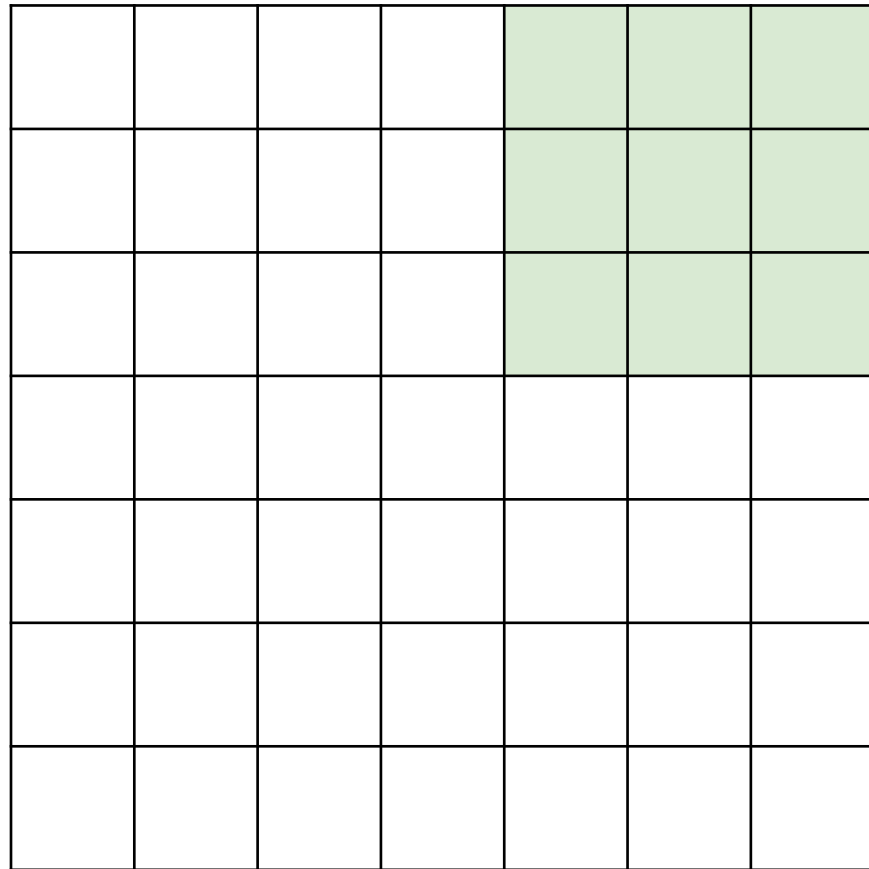


Input: 7x7

Filter: 3x3

Output: 5x5

# A closer look at spatial dimensions



Input: 7x7

Filter: 3x3

Output: 5x5

In general:      **Problem: Feature maps “shrink” with each layer!**  
Input:  $W$   
Filter:  $K$   
Output:  $W - K + 1$

# A closer look at spatial dimensions

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input: 7x7

Filter: 3x3

Output: 5x5

In general:

Input:  $W$

Filter:  $K$

Output:  $W - K + 1$

Problem: Feature maps “shrink” with each layer!

Solution: **padding**

Add zeros around the input

# A closer look at spatial dimensions

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input: 7x7

Filter: 3x3

Output: 5x5

In general:

Input:  $W$

Filter:  $K$

Padding:  $P$

Output:  $W - K + 1 + 2P$

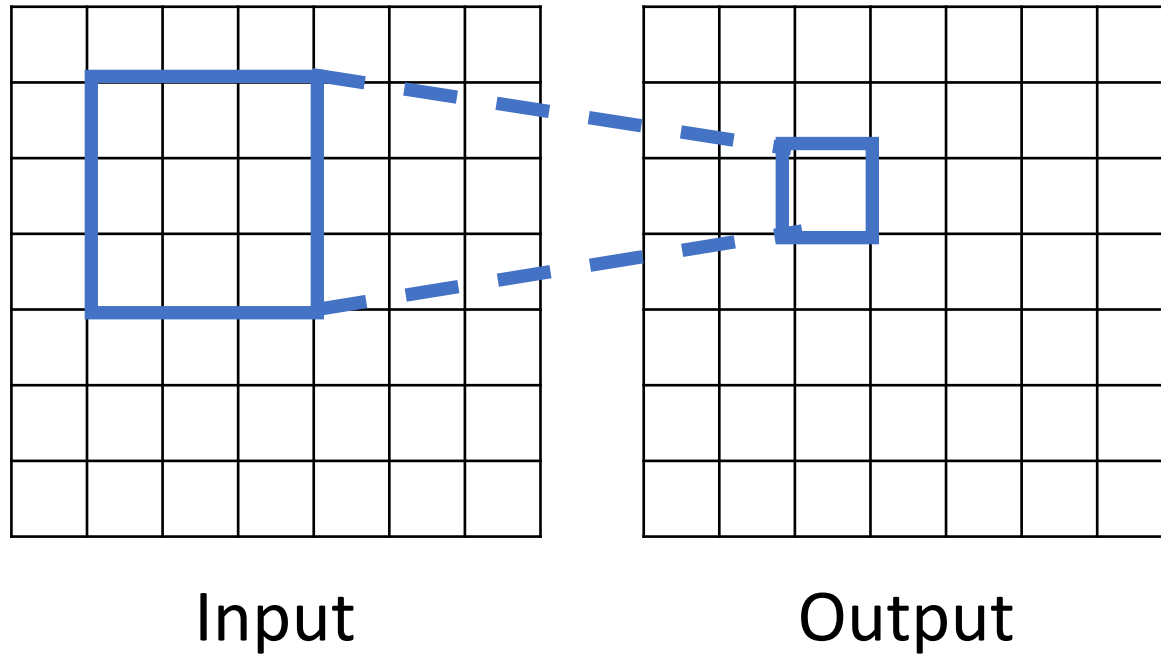
Very common:

Set  $P = (K - 1) / 2$  to  
make output have  
same size as input!



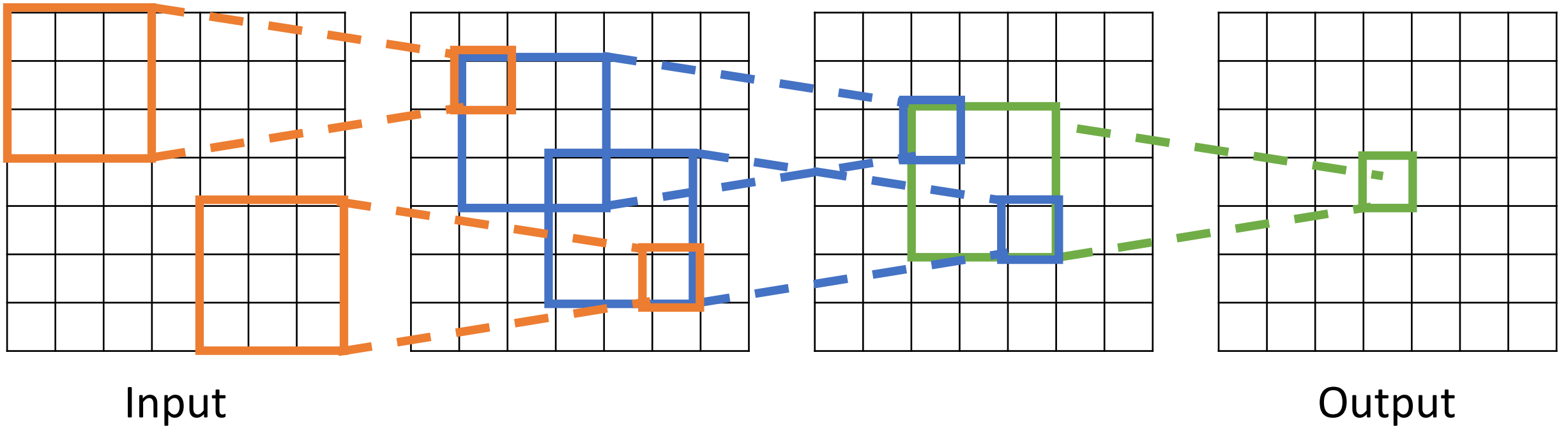
# Receptive Fields

For convolution with kernel size  $K$ , each element in the output depends on a  $K \times K$  **receptive field** in the input



# Receptive Fields

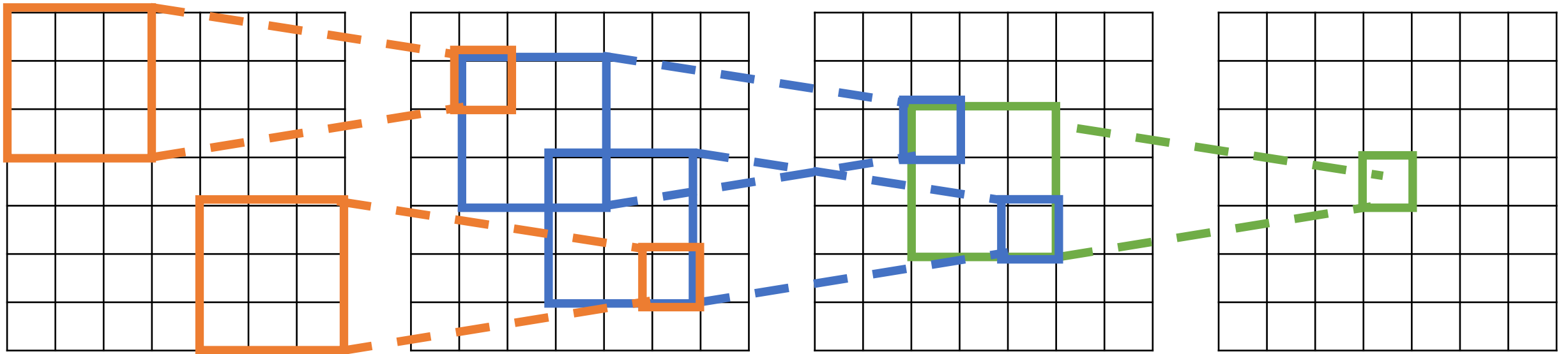
Each successive convolution adds  $K - 1$  to the receptive field size  
With  $L$  layers the receptive field size is  $1 + L * (K - 1)$



Be careful – “receptive field in the input” vs “receptive field in the previous layer”  
Hopefully clear from context!

# Receptive Fields

Each successive convolution adds  $K - 1$  to the receptive field size  
With  $L$  layers the receptive field size is  $1 + L * (K - 1)$



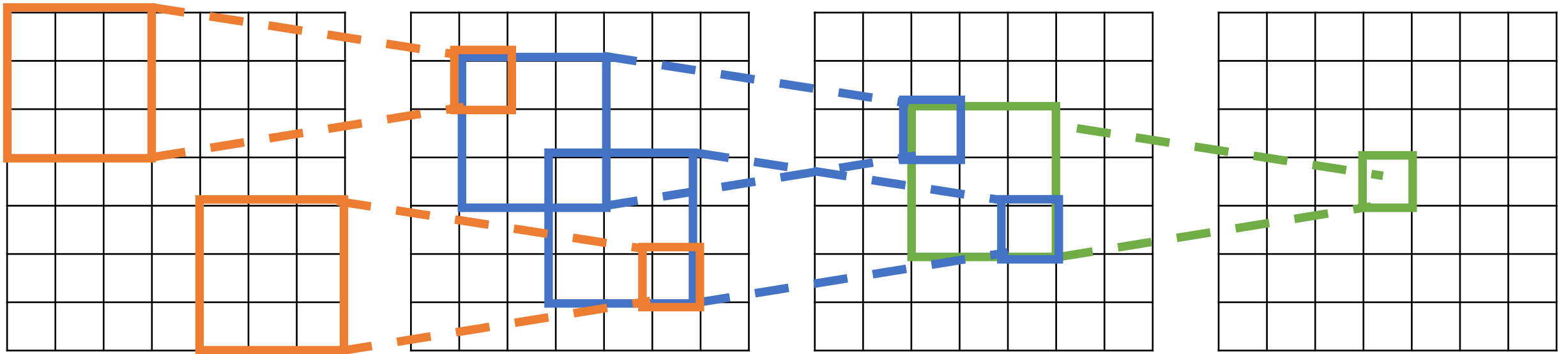
Input

Problem: For large images we need many layers for each output to “see” the whole image

Output

# Receptive Fields

Each successive convolution adds  $K - 1$  to the receptive field size  
With  $L$  layers the receptive field size is  $1 + L * (K - 1)$



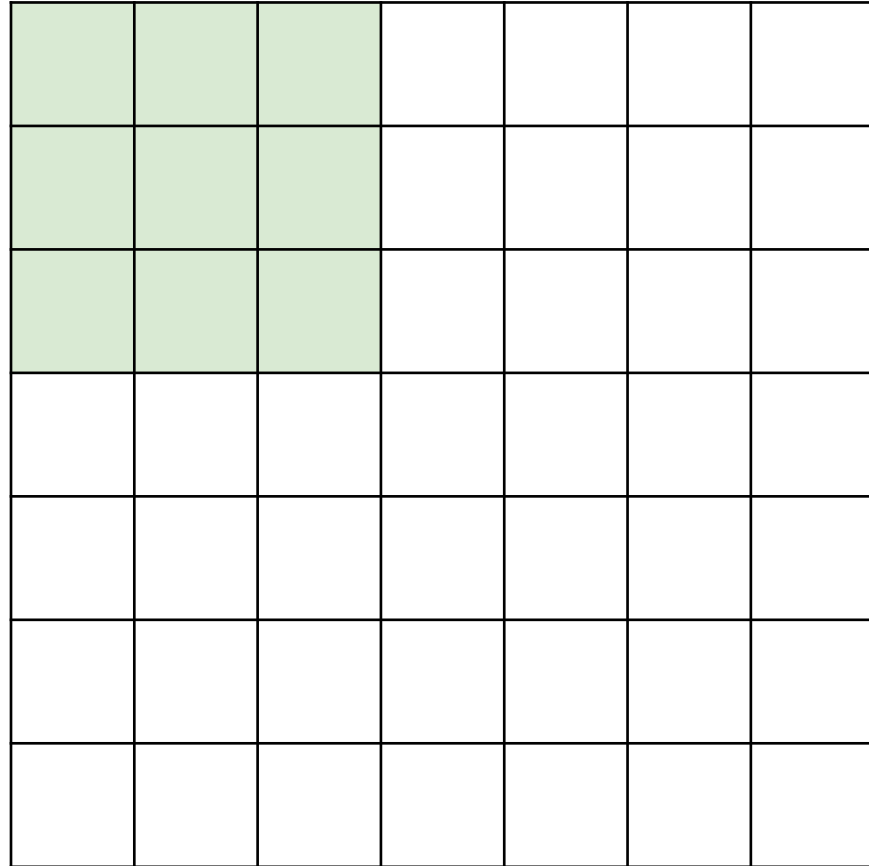
Input

Problem: For large images we need many layers for each output to “see” the whole image

Solution: Downsample inside the network

Output

# Strided Convolution

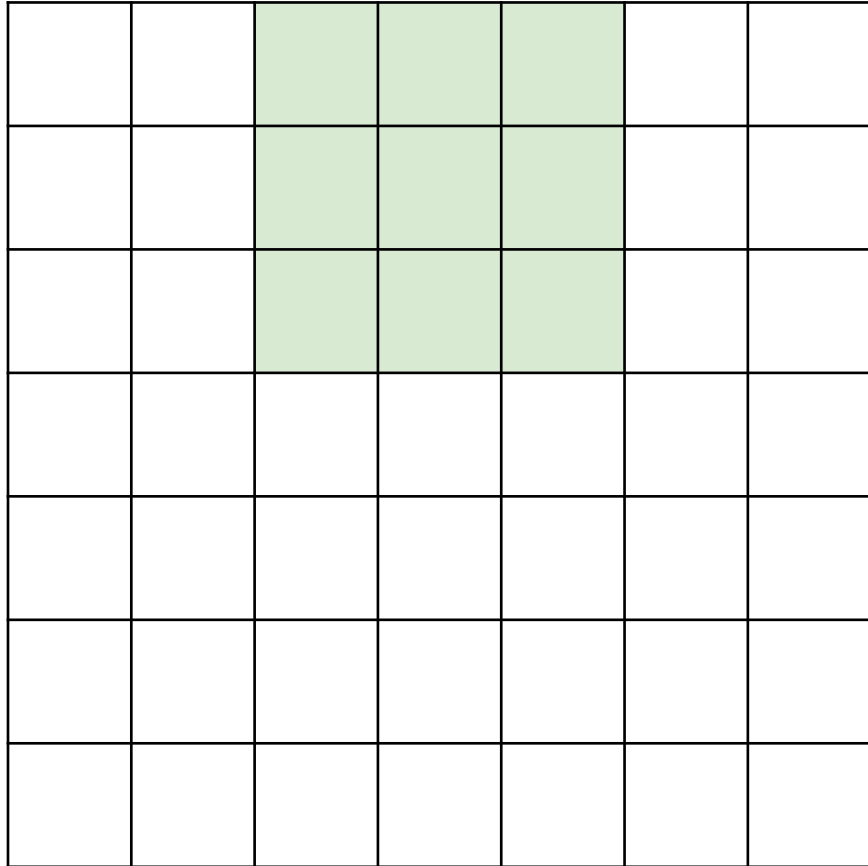


Input: 7x7

Filter: 3x3

Stride: 2

# Strided Convolution

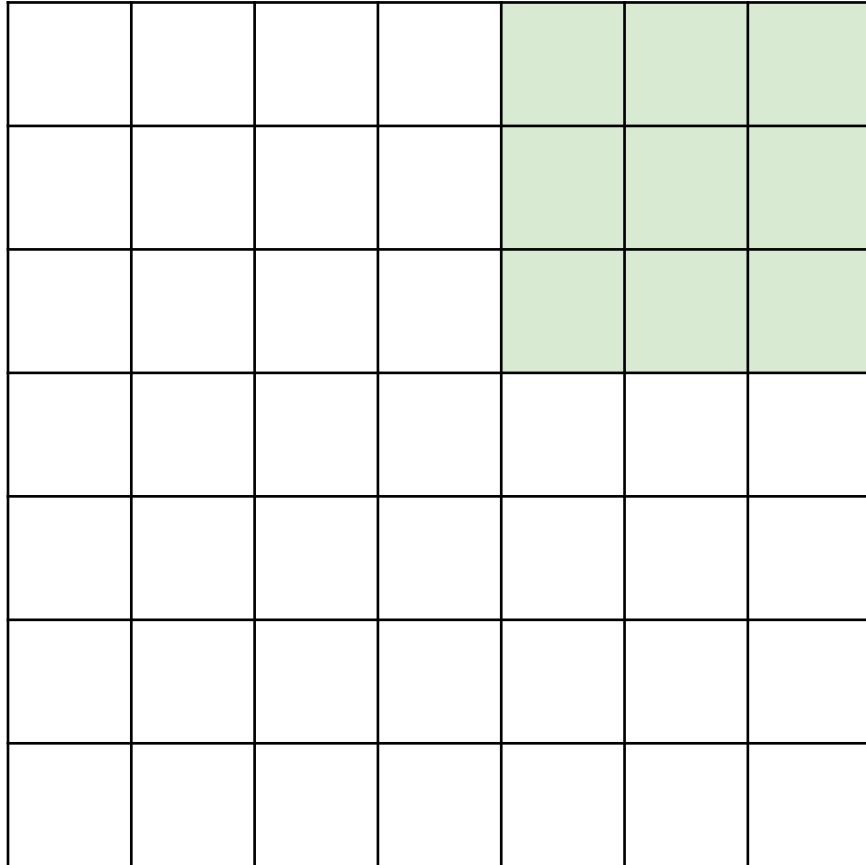


Input: 7x7

Filter: 3x3

Stride: 2

# Strided Convolution



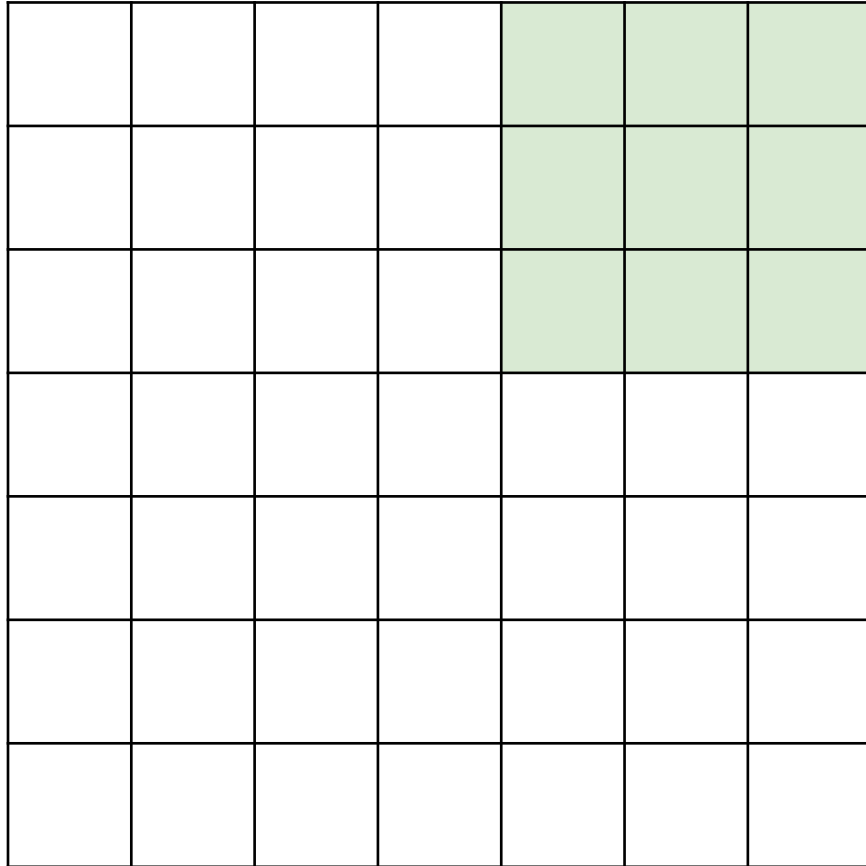
Input: 7x7

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:  $S$

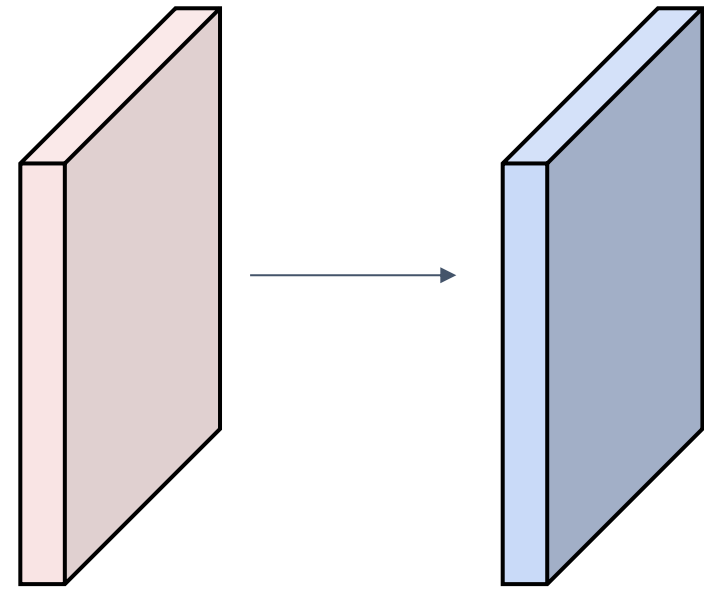
Output:  $(W - K + 2P) / S + 1$



# Convolution Example

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

Output volume size: ?



# Convolution Example

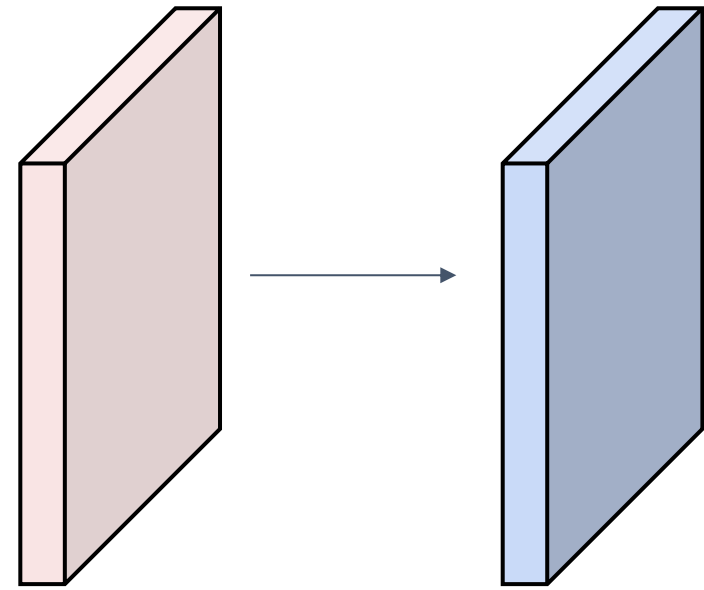
Input volume: 3 x 32 x 32

10 5x5 filters with stride 1, pad 2

Output volume size:

$(32 + 2 * 2 - 5) / 1 + 1 = 32$  spatially, so

10 x 32 x 32



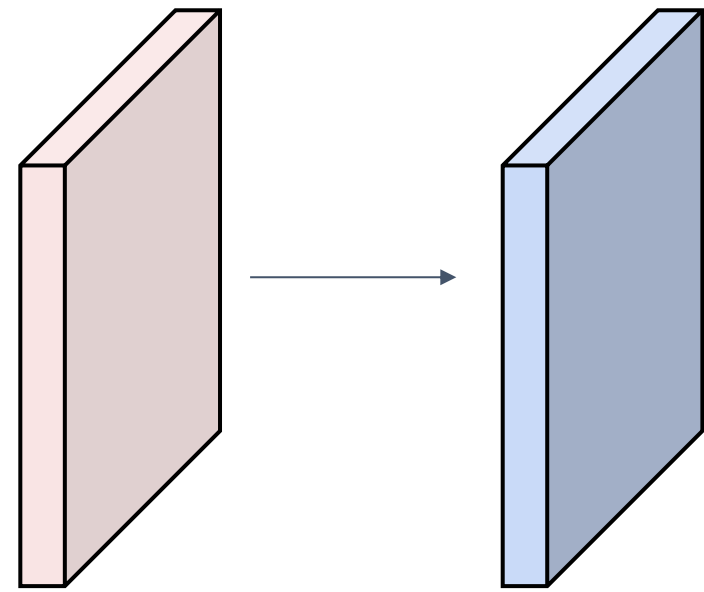
# Convolution Example

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

10  $5 \times 5$  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

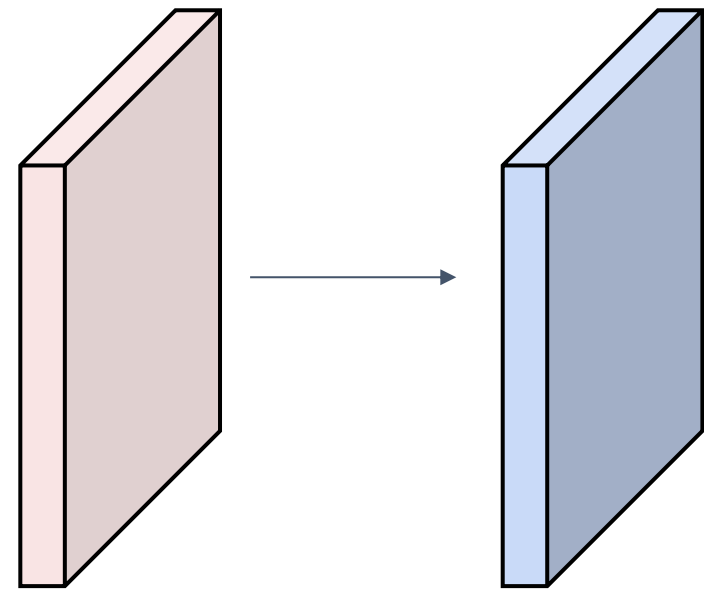
**10** **5x5** filters with stride 1, pad 2

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

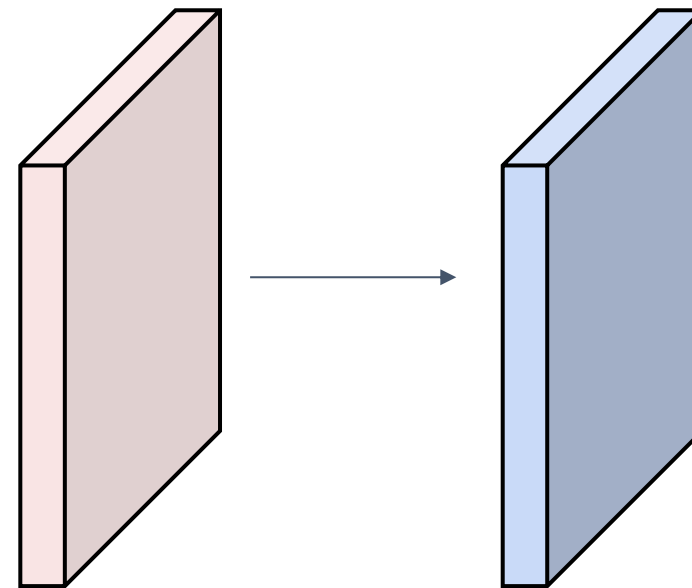
Input volume:  $3 \times 32 \times 32$

10  $5 \times 5$  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

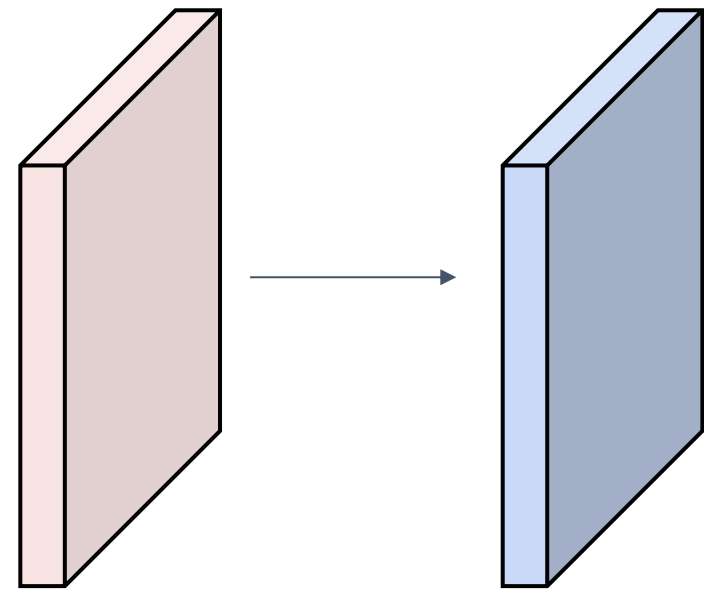
Input volume: **3** x 32 x 32  
10 **5x5** filters with stride 1, pad 2

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

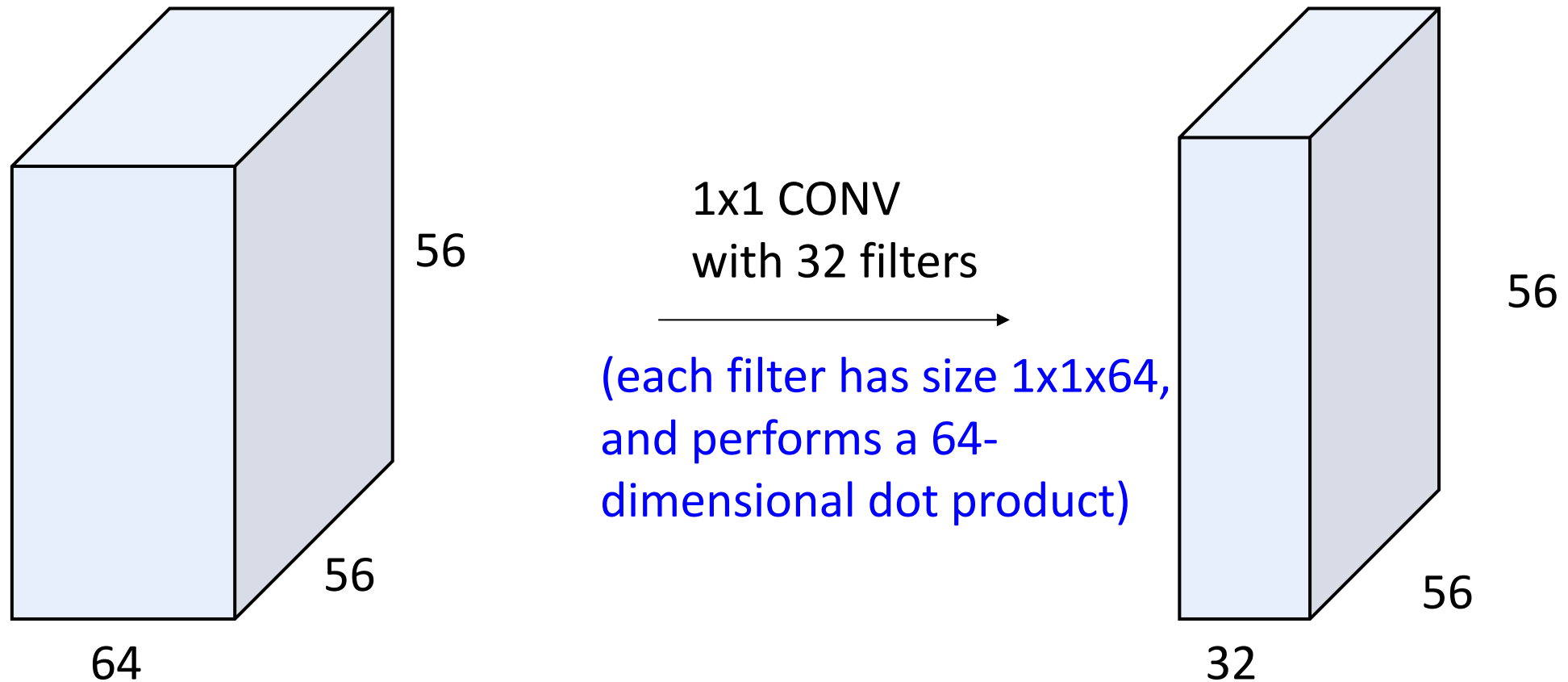
Number of learnable parameters: 760

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

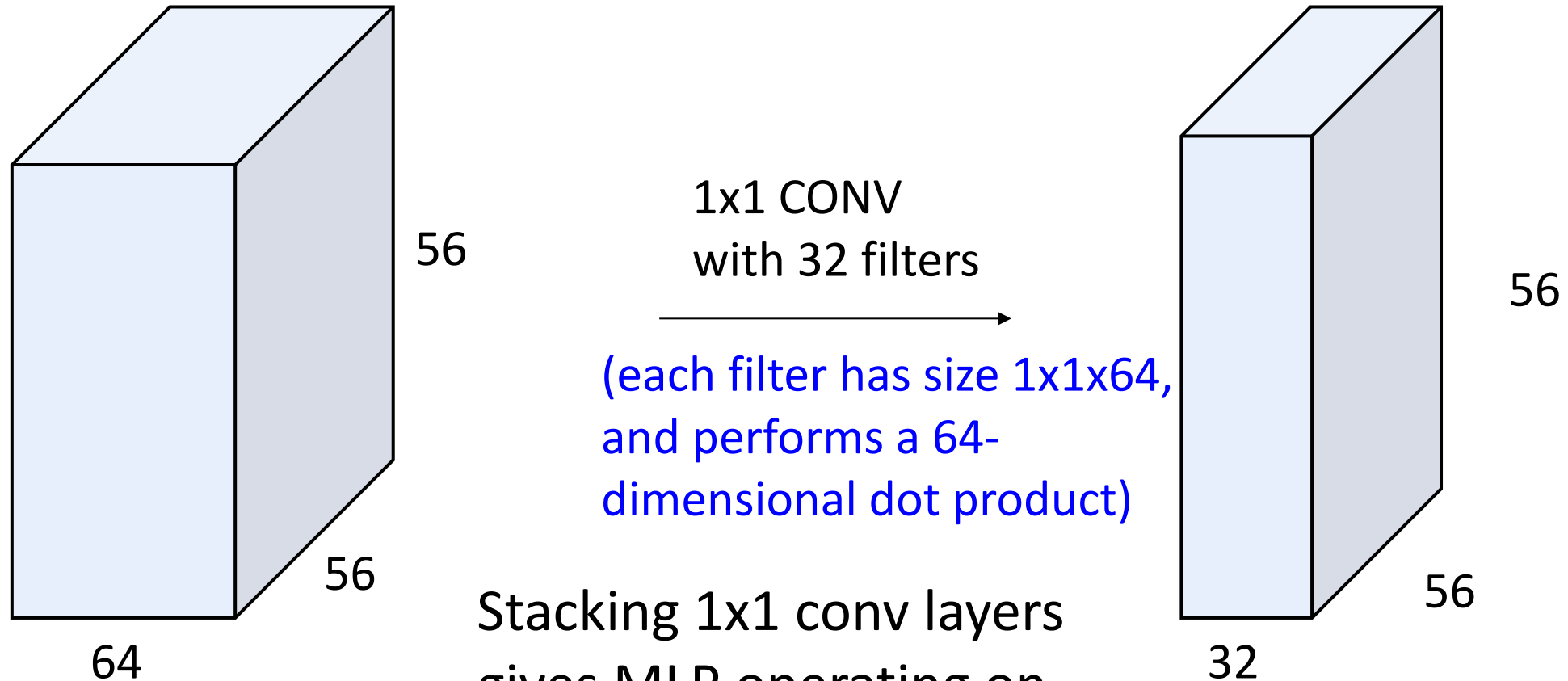
**10\*32\*32** = 10,240 outputs; each output is the inner product of two **3x5x5** tensors (75 elems); total =  $75 * 10240 = 768K$



# Example: 1x1 Convolution



# Example: 1x1 Convolution



Stacking 1x1 conv layers  
gives MLP operating on  
each input position

Lin et al, "Network in Network", ICLR 2014



# Convolution Summary

**Input:**  $C_{in} \times H \times W$

**Hyperparameters:**

- **Kernel size:**  $K_H \times K_W$
- **Number filters:**  $C_{out}$
- **Padding:**  $P$
- **Stride:**  $S$

**Weight matrix:**  $C_{out} \times C_{in} \times K_H \times K_W$

giving  $C_{out}$  filters of size  $C_{in} \times K_H \times K_W$

**Bias vector:**  $C_{out}$

**Output size:**  $C_{out} \times H' \times W'$  where:

- $H' = (H - K + 2P) / S + 1$
- $W' = (W - K + 2P) / S + 1$

# Convolution Summary

**Input:**  $C_{in} \times H \times W$

**Hyperparameters:**

- **Kernel size:**  $K_H \times K_W$
- **Number filters:**  $C_{out}$
- **Padding:**  $P$
- **Stride:**  $S$

**Weight matrix:**  $C_{out} \times C_{in} \times K_H \times K_W$   
giving  $C_{out}$  filters of size  $C_{in} \times K_H \times K_W$

**Bias vector:**  $C_{out}$

**Output size:**  $C_{out} \times H' \times W'$  where:

- $H' = (H - K + 2P) / S + 1$
- $W' = (W - K + 2P) / S + 1$

Common settings:

$K_H = K_W$  (Small square filters)

$P = (K - 1) / 2$  ("Same" padding)

$C_{in}, C_{out} = 32, 64, 128, 256$  (powers of 2)

$K = 3, P = 1, S = 1$  (3x3 conv)

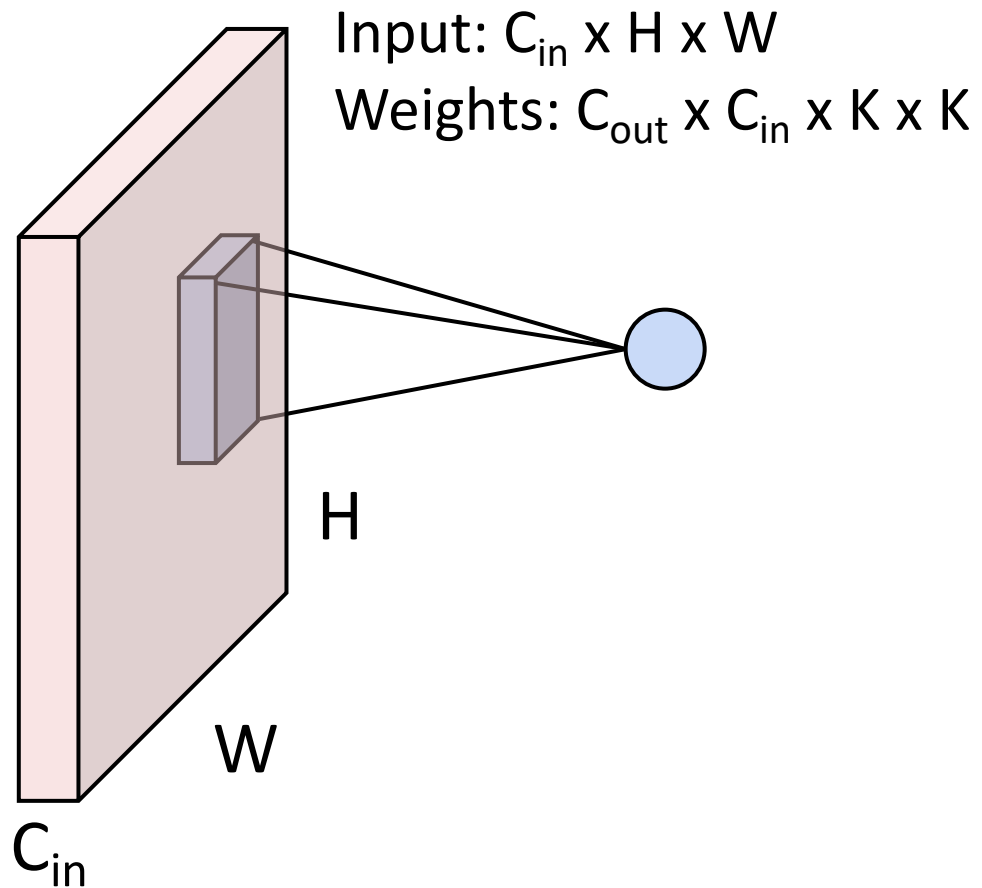
$K = 5, P = 2, S = 1$  (5x5 conv)

$K = 1, P = 0, S = 1$  (1x1 conv)

$K = 3, P = 1, S = 2$  (Downsample by 2)

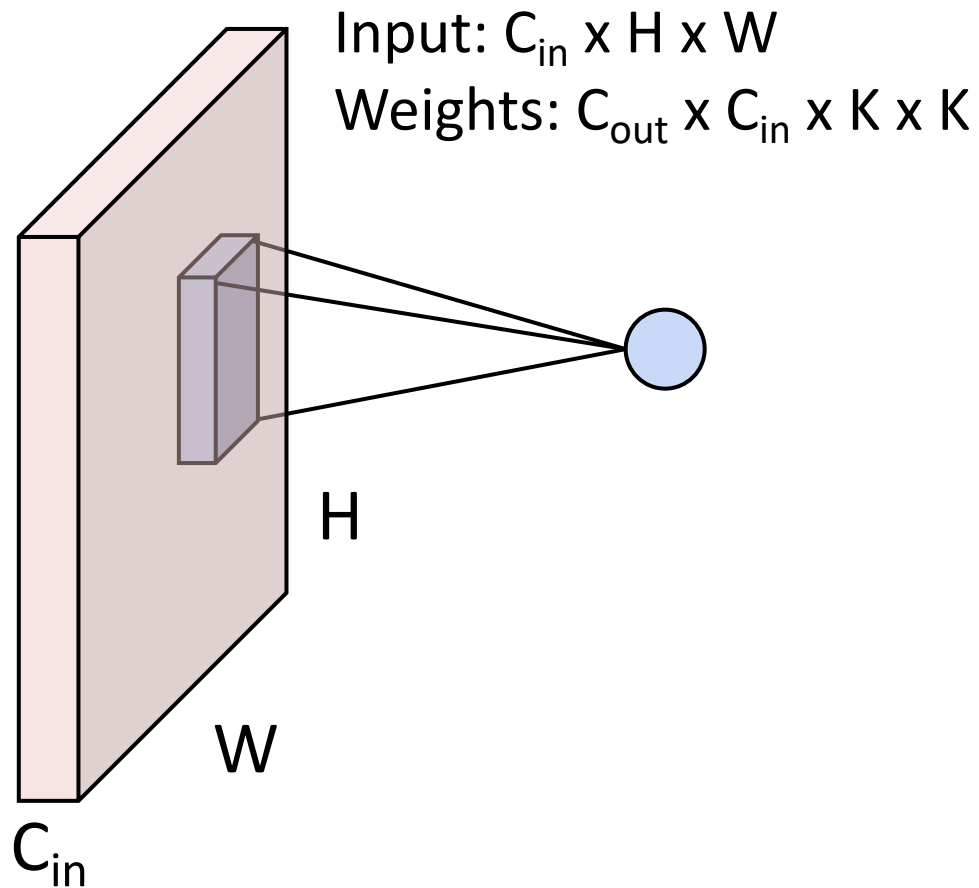
# Other types of convolution

So far: 2D Convolution

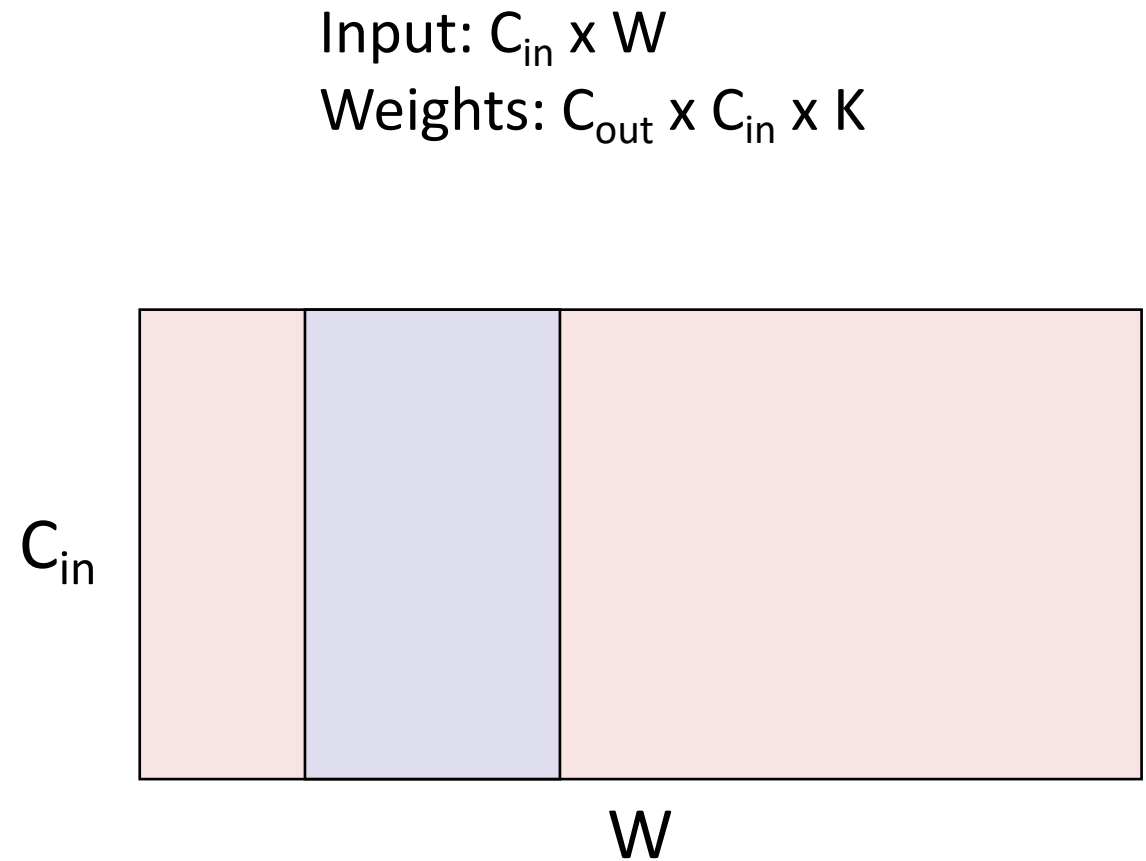


# Other types of convolution

So far: 2D Convolution

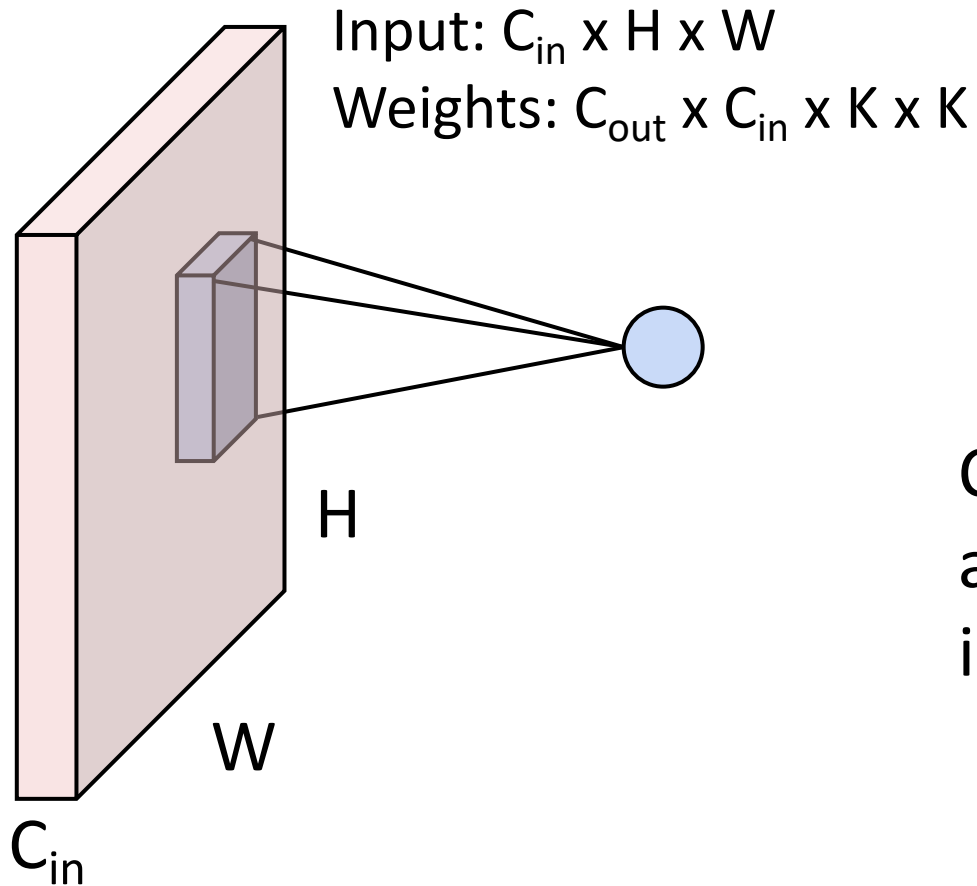


1D Convolution



# Other types of convolution

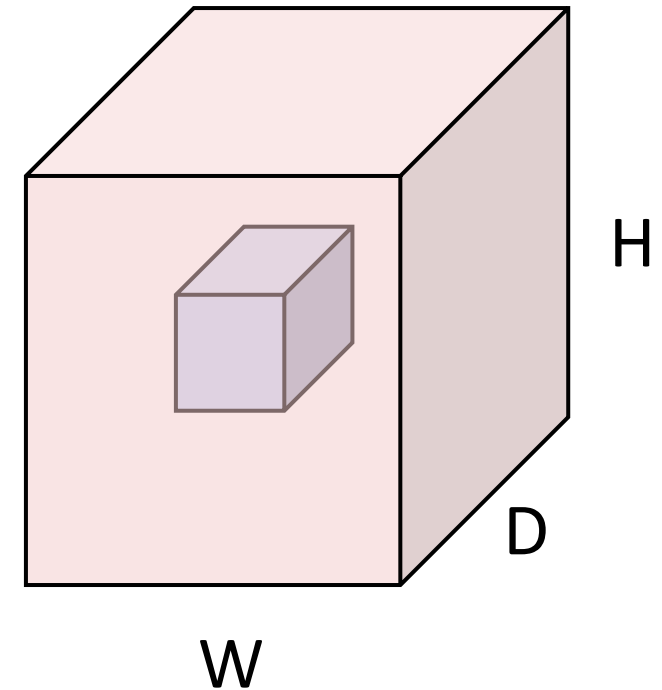
So far: 2D Convolution



$C_{in}$ -dim vector  
at each point  
in the volume

3D Convolution

Input:  $C_{in} \times H \times W \times D$   
Weights:  $C_{out} \times C_{in} \times K \times K \times K$



# PyTorch Convolution Layer

## Conv2d

**CLASS** `torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros')`

[SOURCE]

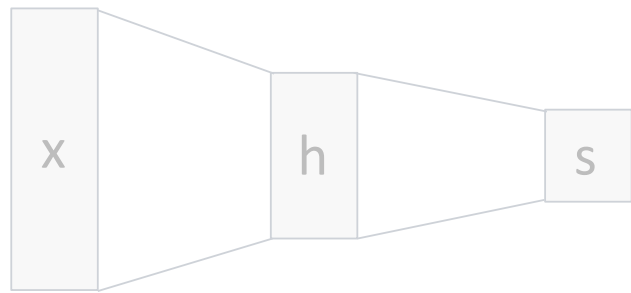
Applies a 2D convolution over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size  $(N, C_{\text{in}}, H, W)$  and output  $(N, C_{\text{out}}, H_{\text{out}}, W_{\text{out}})$  can be precisely described as:

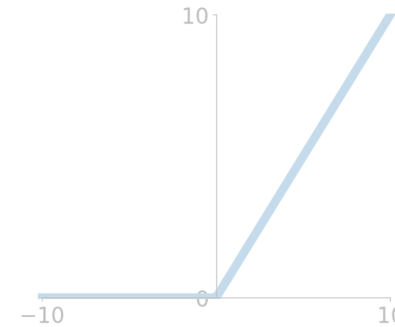
$$\text{out}(N_i, C_{\text{out}_j}) = \text{bias}(C_{\text{out}_j}) + \sum_{k=0}^{C_{\text{in}}-1} \text{weight}(C_{\text{out}_j}, k) \star \text{input}(N_i, k)$$

# Components of a Convolutional Network

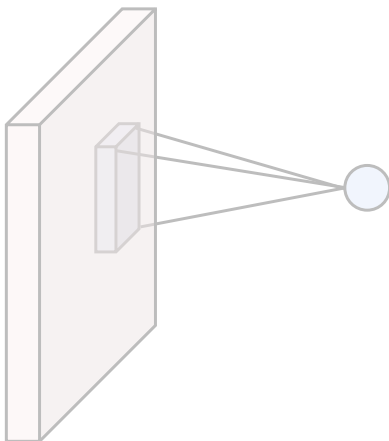
Fully-Connected Layers



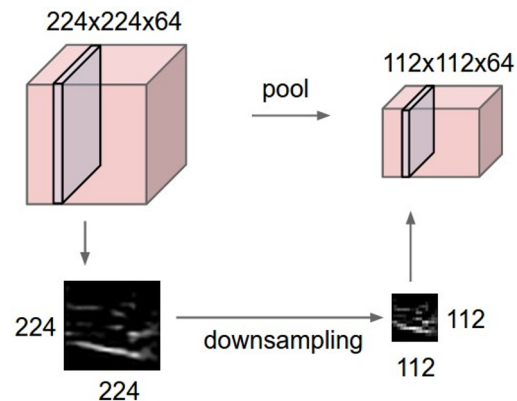
Activation Function



Convolution Layers



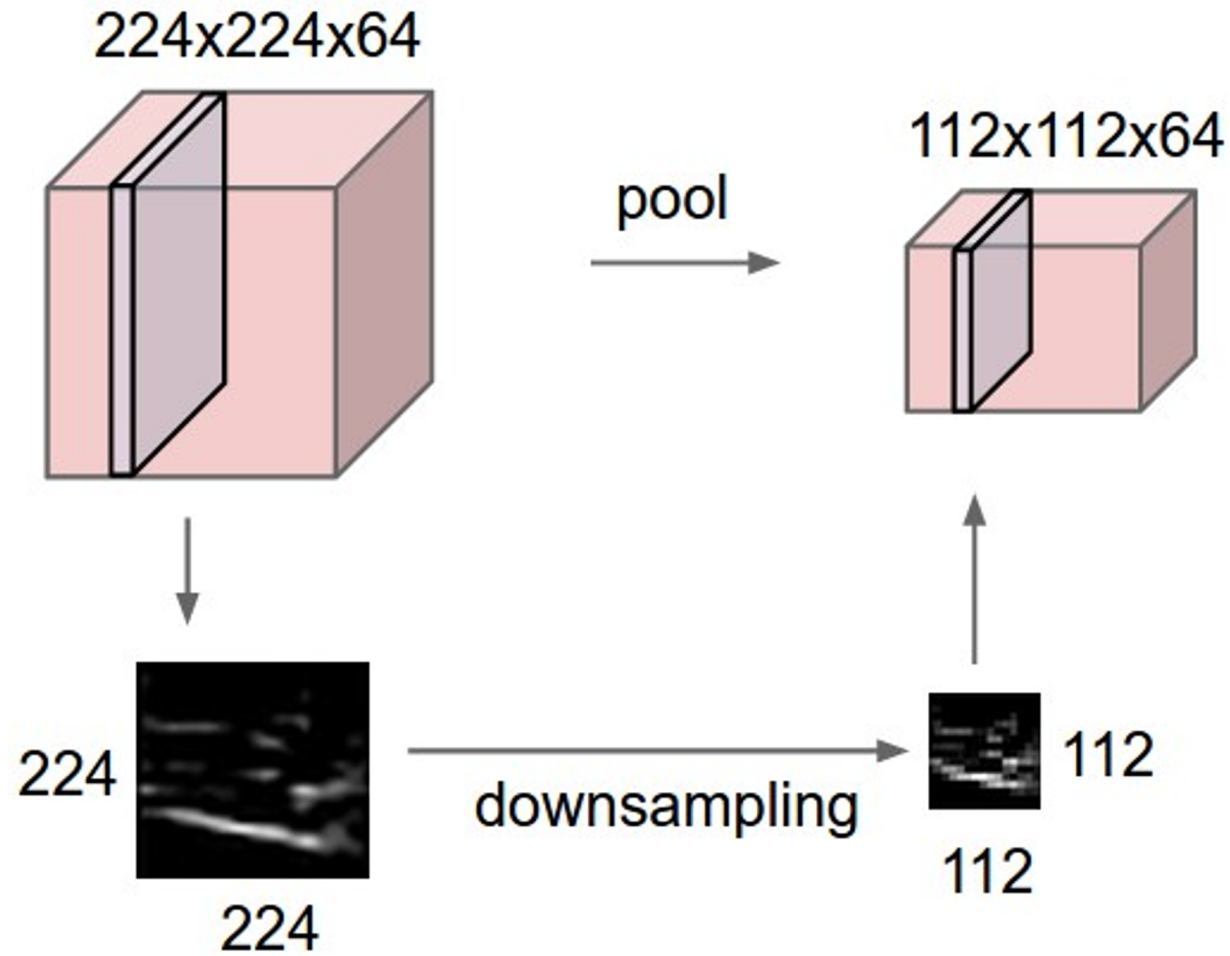
Pooling Layers



Normalization

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

# Pooling Layers: Another way to downsample

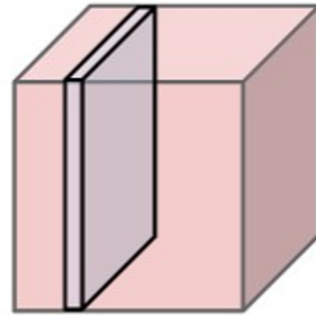


**Hyperparameters:**  
Kernel Size  
Stride  
Pooling function



# Max Pooling

224x224x64



Single depth slice

x

1	1	2	4
5	<b>6</b>	7	<b>8</b>
<b>3</b>	2	1	0
1	2	3	<b>4</b>

y

Max pooling with 2x2  
kernel size and stride 2



6	8
3	4

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

# Pooling Summary

**Input:**  $C \times H \times W$

**Hyperparameters:**

- Kernel size:  $K$
- Stride:  $S$
- Pooling function (max, avg)

**Output:**  $C \times H' \times W'$  where

- $H' = (H - K) / S + 1$
- $W' = (W - K) / S + 1$

**Learnable parameters:** None!

Common settings:

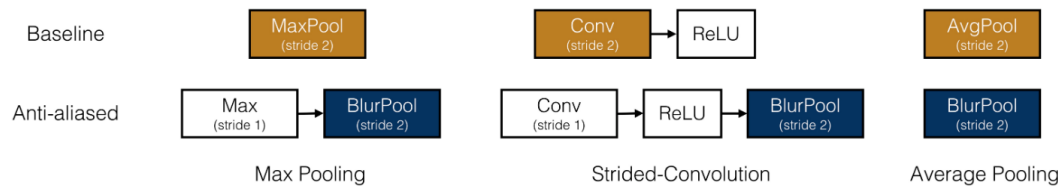
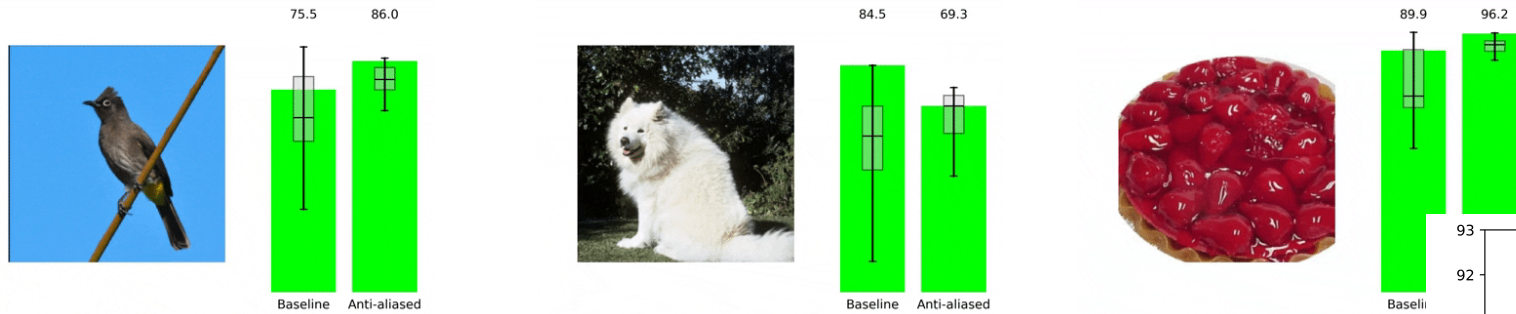
max,  $K = 2, S = 2$

max,  $K = 3, S = 2$  (AlexNet)

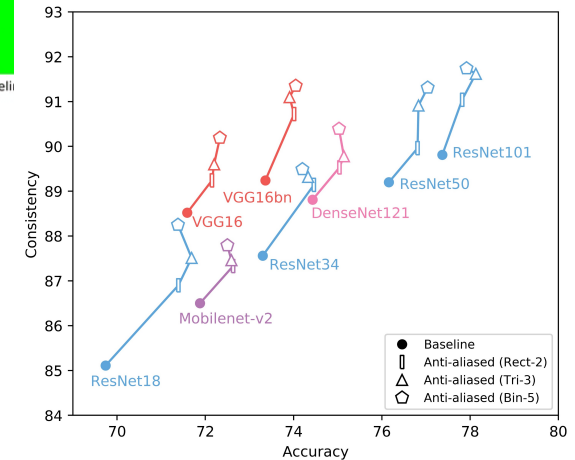
# What about shift invariance?

## Making Convolutional Networks Shift-Invariant Again

Richard Zhang



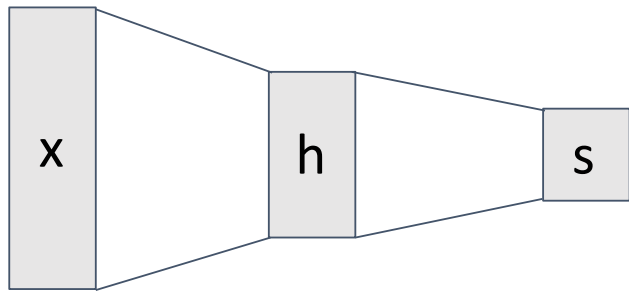
We anti-alias modern networks with classic signal processing, making them more shift-invariant. Predominant downsampling methods ignore the Nyquist sampling theorem. We make the following replacements: MaxPool→MaxBlurPool (pictured above), StridedConv→ConvBlurPool, and AvgPool→BlurPool.



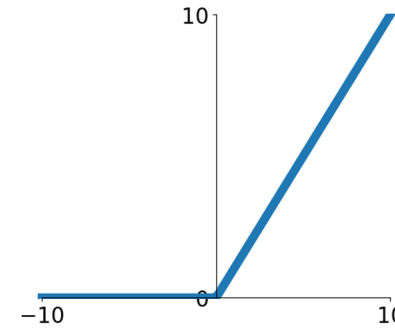
<https://richzhang.github.io/antialiased-cnns/>

# Components of a Convolutional Network

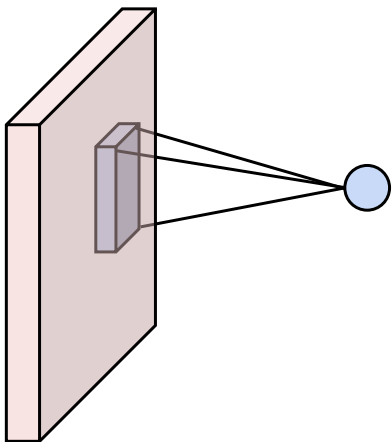
## Fully-Connected Layers



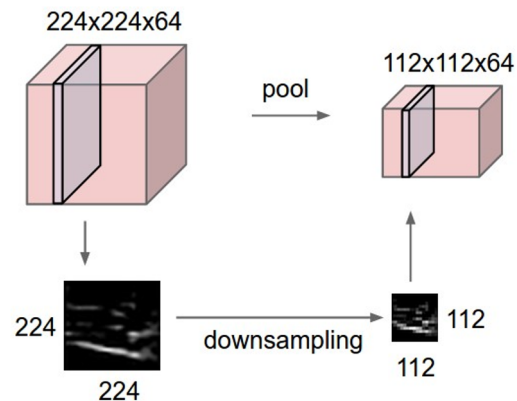
## Activation Function



## Convolution Layers



## Pooling Layers



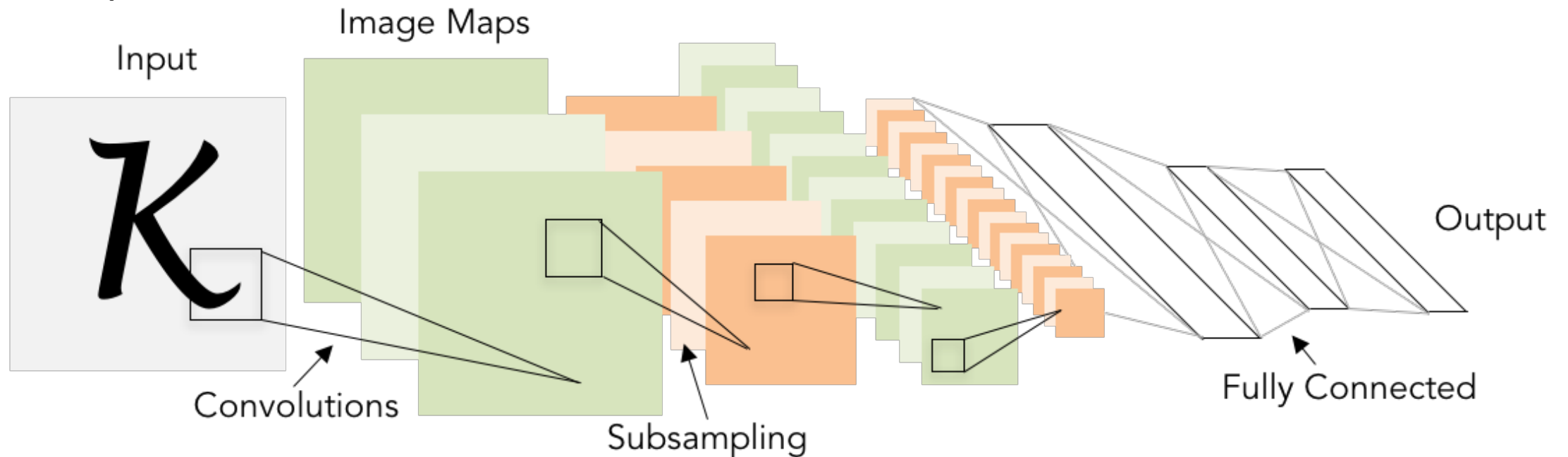
## Normalization

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

# Convolutional Networks

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

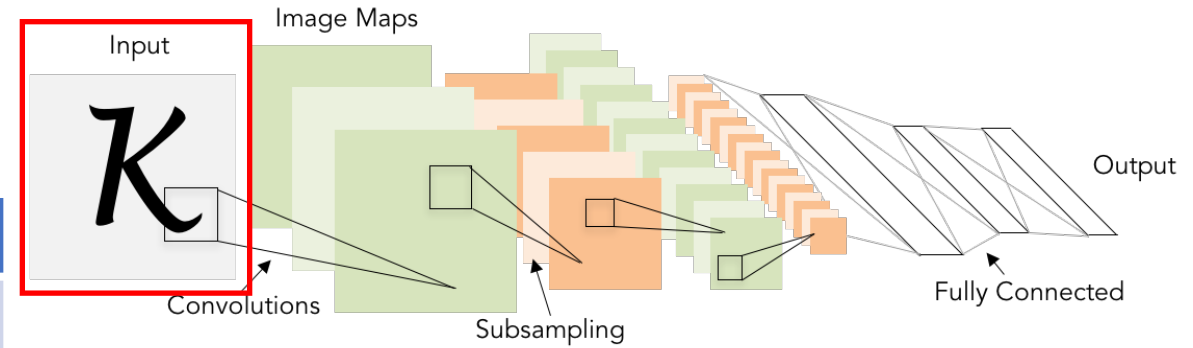
Example: LeNet-5



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

# Example: LeNet-5

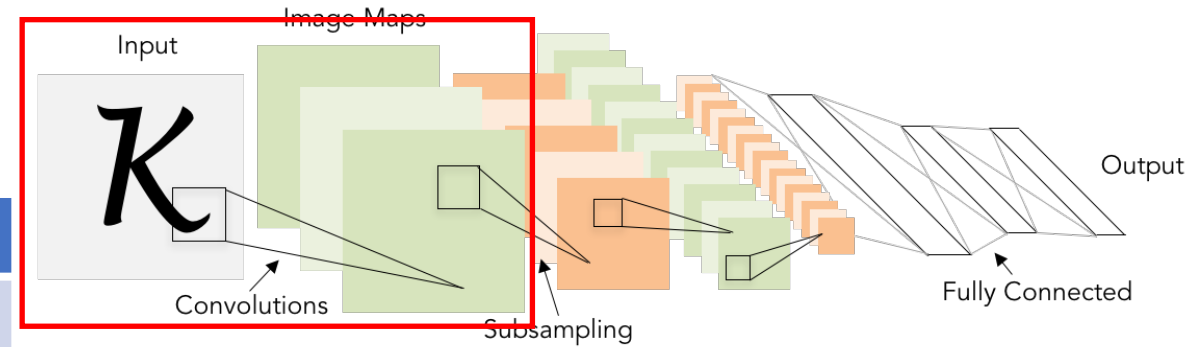
Layer	Output Size	Weight Size
Input	1 x 28 x 28	



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	



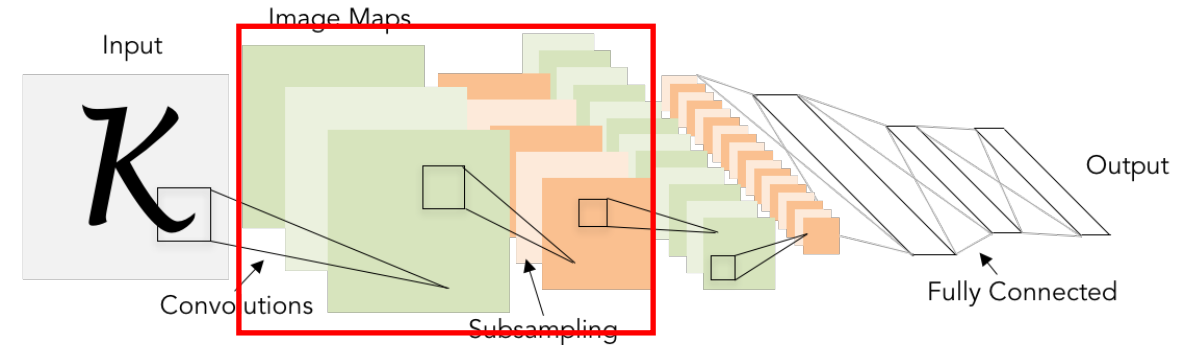
\* Original paper:  $C_{out} = 6$

\*\* Original paper: sigmoid

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	



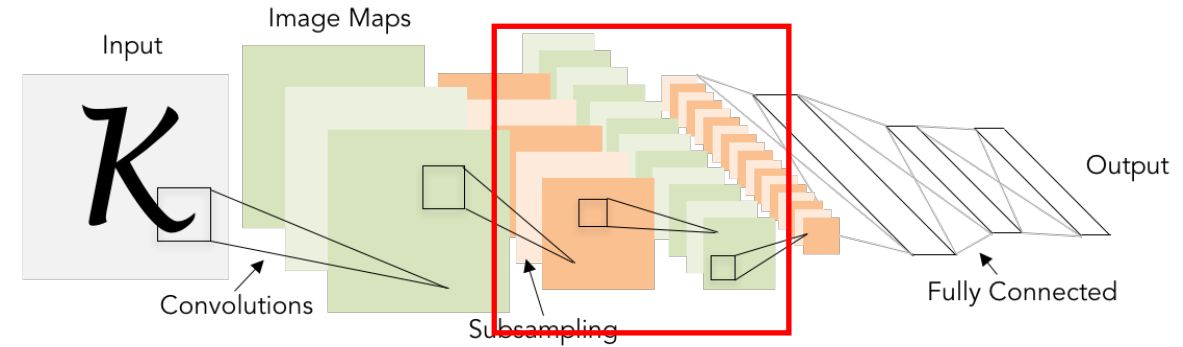
\* 2x2 strided convolution

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	

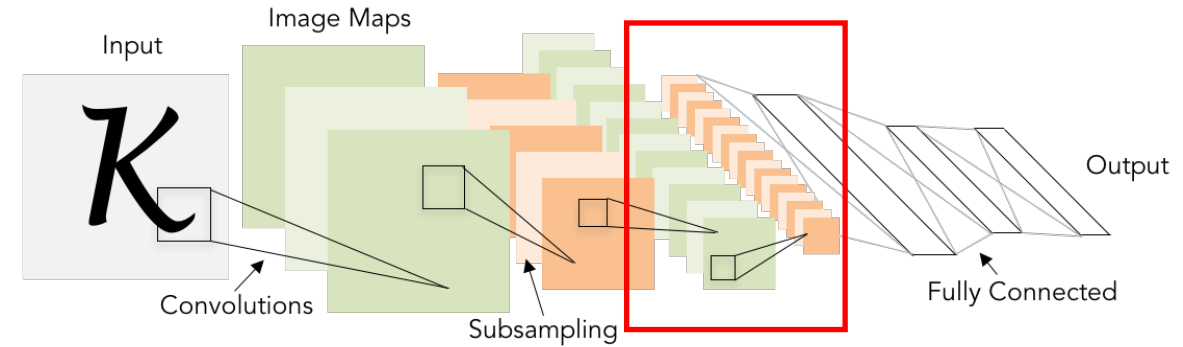


\* Original paper:  $C_{out} = 16$ , grouped convolutions

\*\* Original paper: sigmoid

# 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	

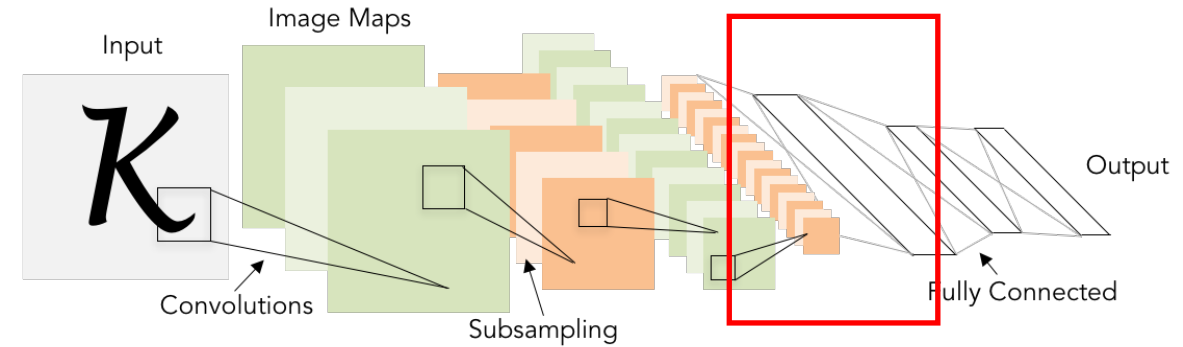


\* 2x2 strided convolution

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

# Example: LeNet-5

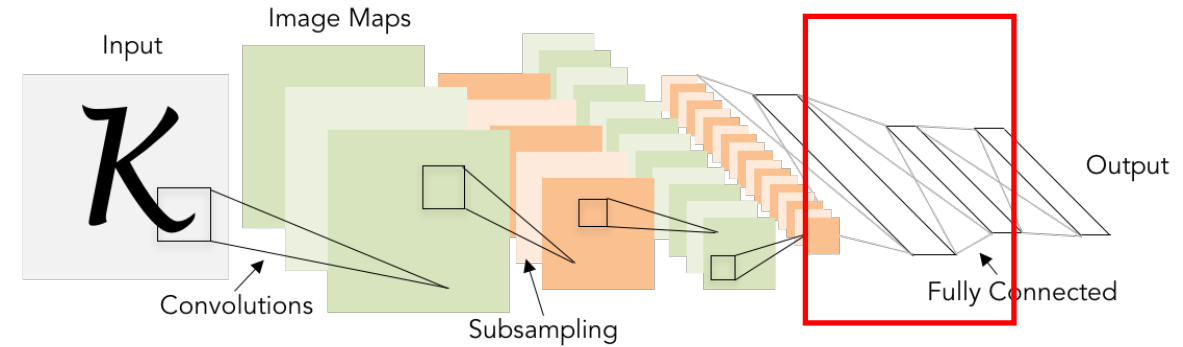
Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv ( $C_{out}=20$ , $K=5$ , $P=2$ , $S=1$ )	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool( $K=2$ , $S=2$ )	20 x 14 x 14	
Conv ( $C_{out}=50$ , $K=5$ , $P=2$ , $S=1$ )	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	
MaxPool( $K=2$ , $S=2$ )	50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



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	

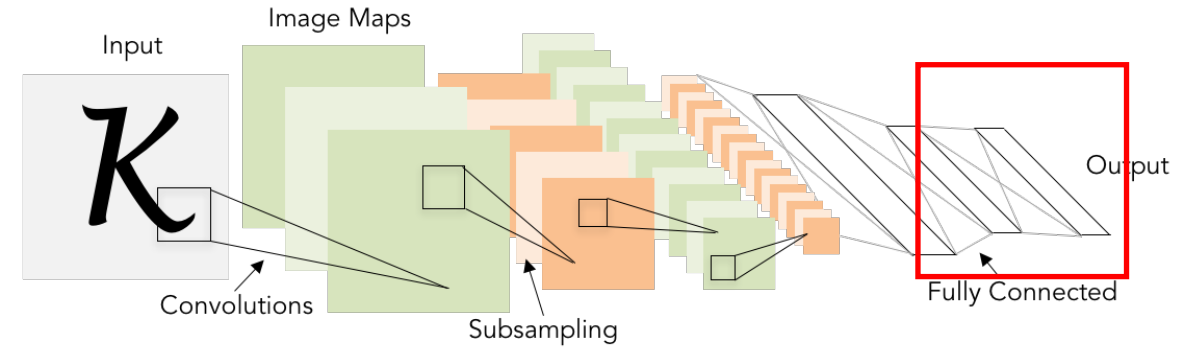


\* Original paper has different 1x1 convolutions, sigmoid non-linearities

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

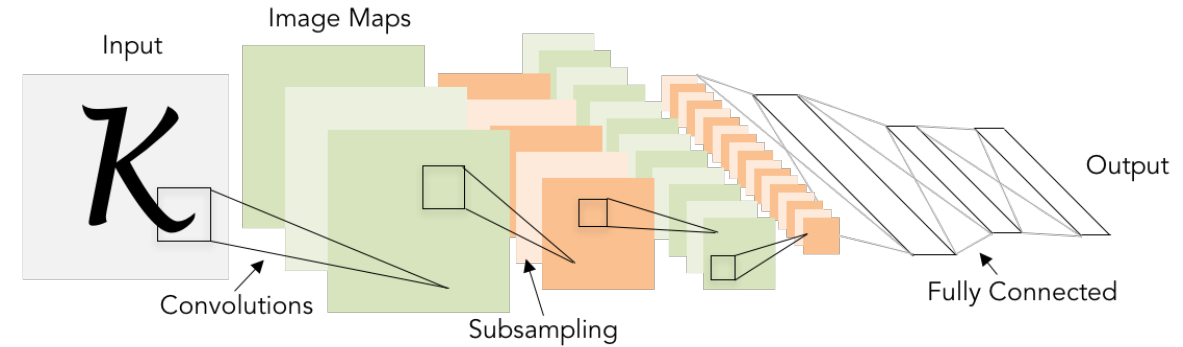


\* Original paper uses RBF (radial basis function) kernels instead of a softmax

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:

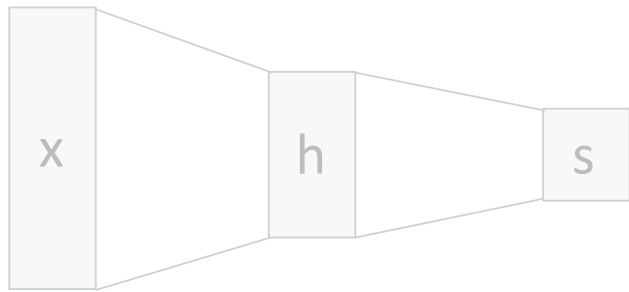
Spatial size **decreases**  
(using pooling or strided conv)

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

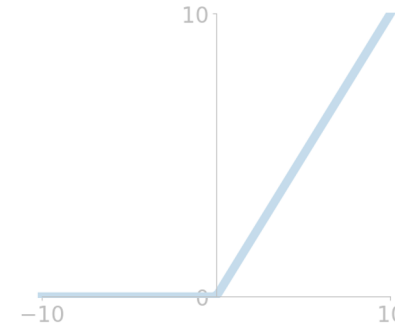
Problem: Deep Networks very hard to train!

# Components of a Convolutional Network

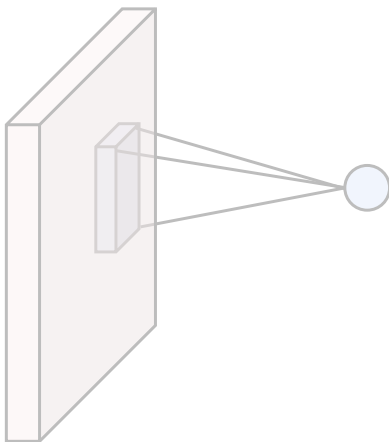
Fully-Connected Layers



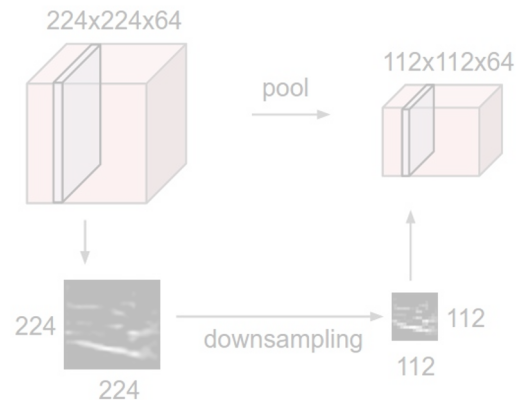
Activation Function



Convolution Layers



Pooling Layers



Normalization

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



# Batch Normalization

Idea: “Normalize” the outputs of a layer so they have zero mean and unit variance. **Why?**

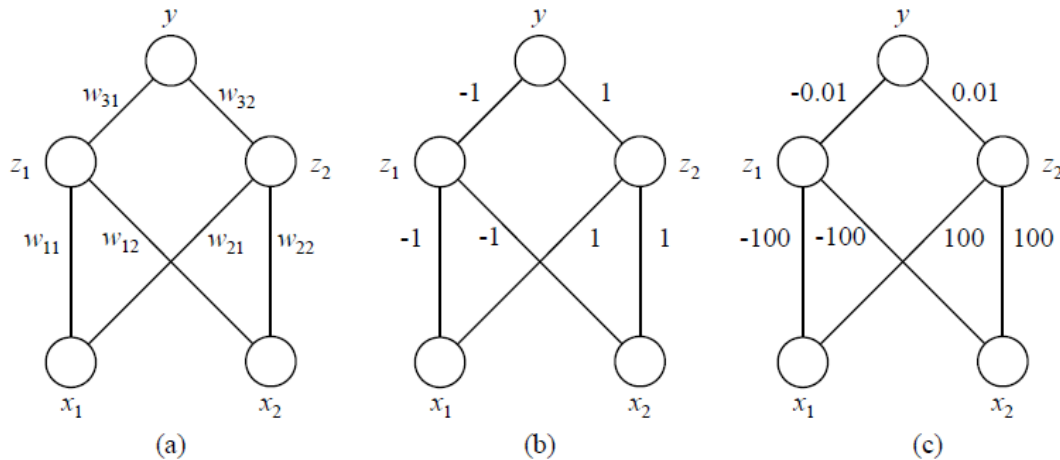
Why? Helps reduce “internal covariate shift”, improves optimization

We can normalize a batch of activations like this:

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

This is a **differentiable function**, so we can use it as an operator in our networks and backprop through it!

# Activation and weight scaling

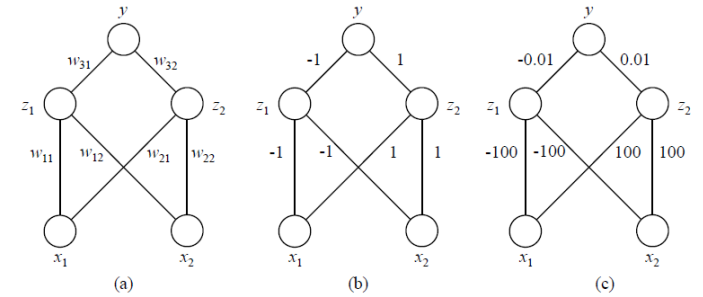


**Figure 5.53** Simple two hidden unit network with a ReLU activation function and no bias parameters for regressing the function  $y = |x_1 + 1.1x_2|$ : (a) can you guess a set of weights would fit this function? (b) a reasonable set of starting weights; (c) a poorly scaled set of weights.

2. Starting with the weights shown in column b, compute the activations for the hidden and final units as well as the regression loss for the four input values  $(x_1, x_2) \in \{-1, 0, 1\} \times \{-1, 0, 1\}$ .
3. Now compute the gradients of the squared loss with respect to all six weights using the backpropagation chain rule equations (5.78–5.81) and sum them up across the training samples to get a final gradient.

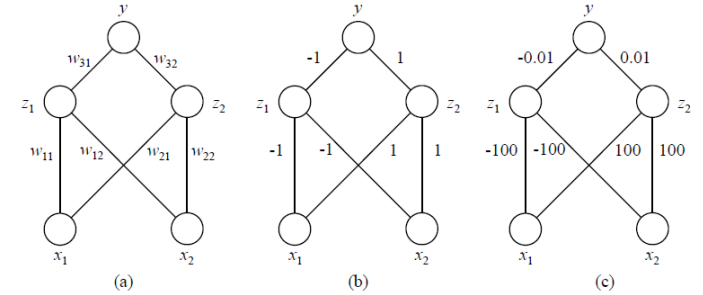
4. What step size should you take in the gradient direction, and what would your update squared loss become?
5. Repeat this exercise for the initial weights in column (c) of Figure 5.52.
6. Given this new set of weights, how much worse is your error decrease, and how many iterations would you expect it to take to achieve a reasonable solution?
7. Would batch normalization help in this case?

# Activation and weight scaling



x1	x2	w11	w12	w21	w22	z1	z2	w31	w32	y	target	L2 loss
-1	-1	-1	-1	1	1	2	0	1	1	2	2.1	0.01
-1	0					1	0			1	1	0
-1	1					0	0			0	0.1	0.01
0	-1					1	0			1	1.1	0.01
0	0					0	0			0	0	0
0	1					0	1			1	1.1	0.01
1	-1					0	0			0	0.1	0.01
1	0					0	1			1	1	0
1	1					0	2			2	2.1	0.01
												0.06
x1	x2	d w11	d w12	d w21	d w22	d z1	d z2	d w31	d w32	g y		
-1	-1	0.1	0.1	0	0	-0.1	-0.1	-0.2	0	-0.1		
-1	0	0	0	0	0	0	0	0	0	0		
-1	1	0	0	0	0	0	0	0	0	-0.1		
0	-1	0	0.1	0	0	-0.1	-0.1	-0.1	0	-0.1		
0	0	0	0	0	0	0	0	0	0	0		
0	1	0	0	0	-0.1	-0.1	-0.1	0	-0.1	-0.1		
1	-1	0	0	0	0	0	0	0	0	-0.1		
1	0	0	0	0	0	0	0	0	0	0		
1	1	0	0	-0.1	-0.1	-0.1	-0.1	0	-0.2	-0.1		
grad-:		-0.1	-0.2	0.1	0.2			0.3	0.3			

# Activation and weight scaling



x1	x2	w11	w12	w21	w22	z1	z2	w31	w32	y	target	L2 loss
-1	-1	-1	-1	1	1	2	0	1	1	2	2.1	0.01
-1	0											
-1	1											
0	-1											
0	0											
0	1											
1	-1											
1	0											
1	1											

x1	x2	w11	w12	w21	w22	z1	z2	w31	w32	y	target	L2 loss
-1	-1	-100	-100	100	100	200	0	0.01	0.01	2	2.1	0.01
-1	0					100	0			1	1	0
-1	1					0	0			0	0.1	0.01
0	-1					100	0			1	1.1	0.01
0	0					0	0			0	0	0
0	1					0	100			1	1.1	0.01
1	-1					0	0			0	0.1	0.01
1	0					0	100			1	1	0
1	1					0	200			2	2.1	0.01

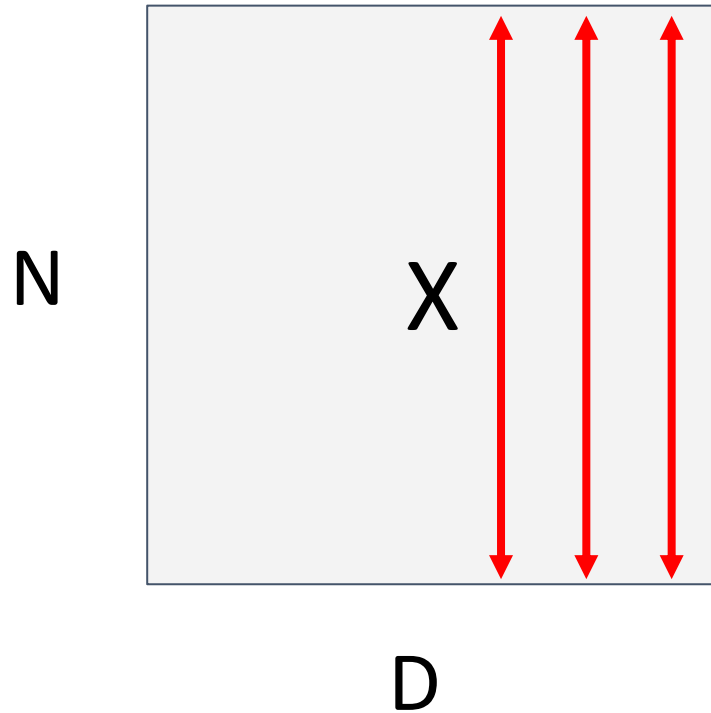
x1	x2	d w11	d w12	d w21	d w22	d z1
-1	-1	0.1	0.1	0	0	
-1	0	0	0	0	0	
-1	1	0	0	0	0	
0	-1	0	0.1	0	0	
0	0	0	0	0	0	
0	1	0	0	0	-0.1	
1	-1	0	0	0	0	
1	0	0	0	0	0	
1	1	0	0	-0.1	-0.1	
grad-		-0.1	-0.2	0.1	0.2	

x1	x2	d w11	d w12	d w21	d w22	d z1	d z2	d w31	d w32	g y
-1	-1	0.001	0.001	0	0	-0.001	-0.001	-20	0	-0.1
-1	0	0	0	0	0	0	0	0	0	0
-1	1	0	0	0	0	0	0	0	0	-0.1
0	-1	0	0.001	0	0	-0.001	-0.001	-10	0	-0.1
0	0	0	0	0	0	0	0	0	0	0
0	1	0	0	0	-0.001	-0.001	-0.001	0	-10	-0.1
1	-1	0	0	0	0	0	0	0	0	-0.1
1	0	0	0	0	0	0	0	0	0	0
1	1	0	0	-0.001	-0.001	-0.001	-0.001	0	-20	-0.1
grad-		-0.001	-0.002	0.001	0.002			30	30	

# Batch Normalization

**Input:**  $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean, shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

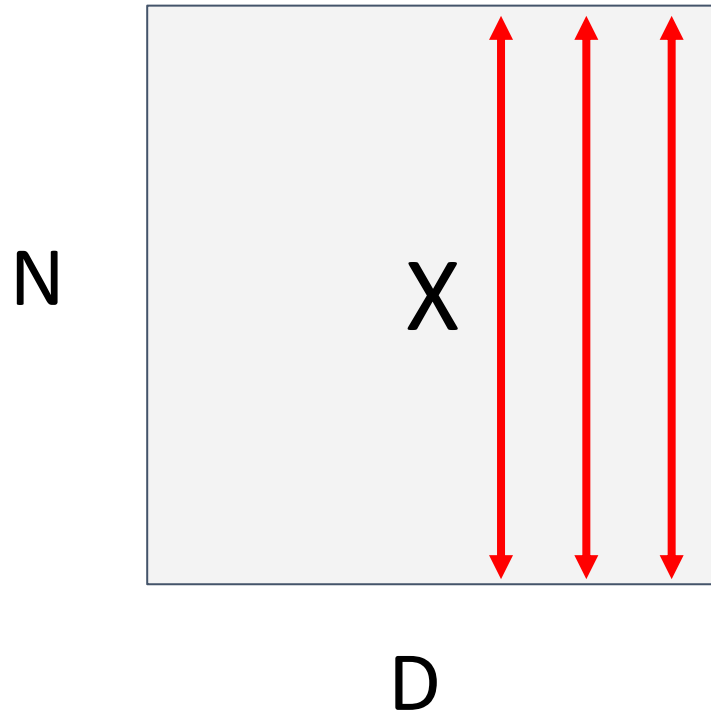
Per-channel std, shape is D

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

Normalized x,  
Shape is N x D

# Batch Normalization

**Input:**  $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean, shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel std, shape is D

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

Normalized x,  
Shape is N x D

**Problem: What if zero-mean, unit variance is too hard of a constraint?**

# Batch Normalization

**Input:**  $x : N \times D$

**Learnable scale and shift parameters:**

$$\gamma, \beta : D$$

Learning  $\gamma = \sigma$ ,  
 $\beta = \mu$  will recover the identity function!

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \quad \text{Per-channel mean, shape is } D$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \quad \text{Per-channel std, shape is } D$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}} \quad \text{Normalized } x, \text{ Shape is } N \times D$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j \quad \text{Output, Shape is } N \times D$$

# Batch Normalization: Test-Time

**Problem:** Estimates depend on minibatch; can't do this at test-time!

**Input:**  $x : N \times D$

**Learnable scale and shift parameters:**

$$\gamma, \beta : D$$

Learning  $\gamma = \sigma$ ,  
 $\beta = \mu$  will recover the identity function!

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \quad \text{Per-channel mean, shape is } D$$
$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \quad \text{Per-channel std, shape is } D$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}} \quad \text{Normalized } x, \text{ Shape is } N \times D$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j \quad \text{Output, Shape is } N \times D$$



# Batch Normalization: Test-Time

**Input:**  $x : N \times D$

**Learnable scale and shift parameters:**

$$\gamma, \beta : D$$

Learning  $\gamma = \sigma$ ,  
 $\beta = \mu$  will recover the  
identity function!

$\mu_j =$  (Running) average of values seen during training  
Per-channel mean, shape is D

$\sigma_j^2 =$  (Running) average of values seen during training  
Per-channel std, shape is D

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

$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$  Output,  
Shape is N x D

# Batch Normalization: Test-Time

**Input:**  $x : N \times D$

**Learnable scale and shift parameters:**

$\gamma, \beta : D$

During testing batchnorm becomes a linear operator!  
Can be fused with the previous fully-connected or conv layer

$\mu_j =$  (Running) average of values seen during training  
Per-channel mean, shape is D

$\sigma_j^2 =$  (Running) average of values seen during training  
Per-channel std, shape is D

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

$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$  Output,  
Shape is N x D

# Batch Normalization for ConvNets

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)

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



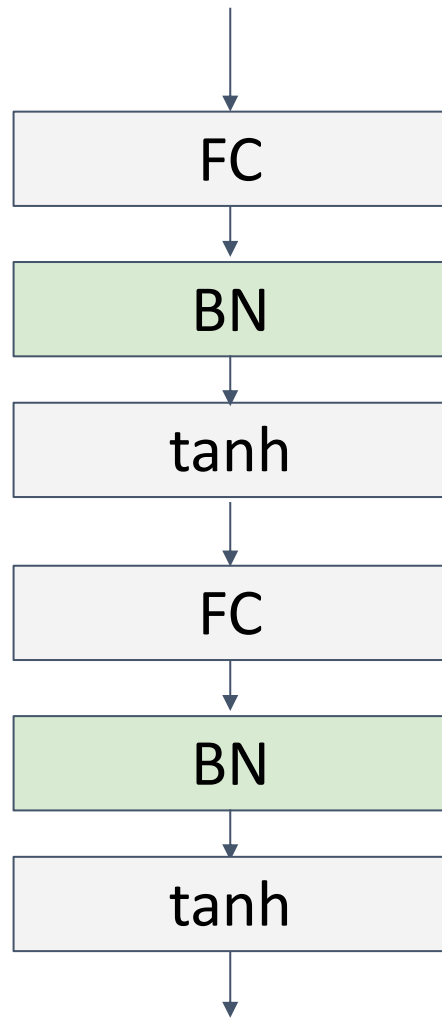
Normalize

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

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

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

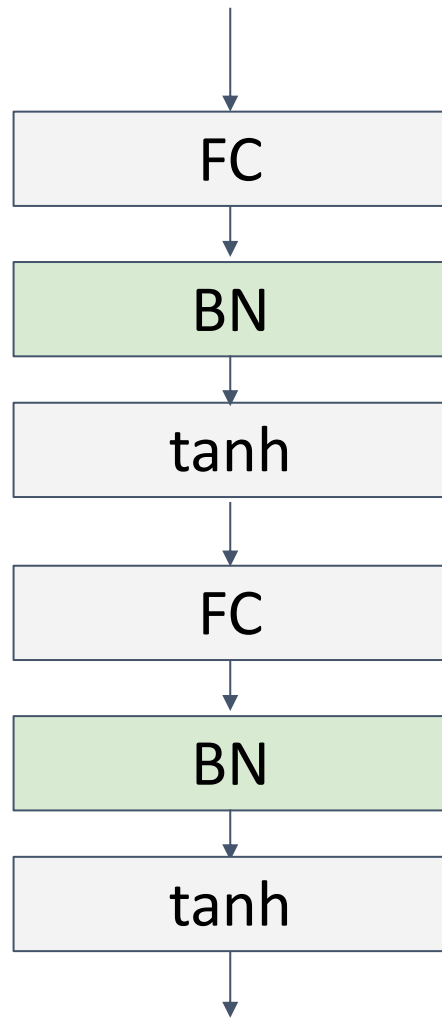
# Batch Normalization



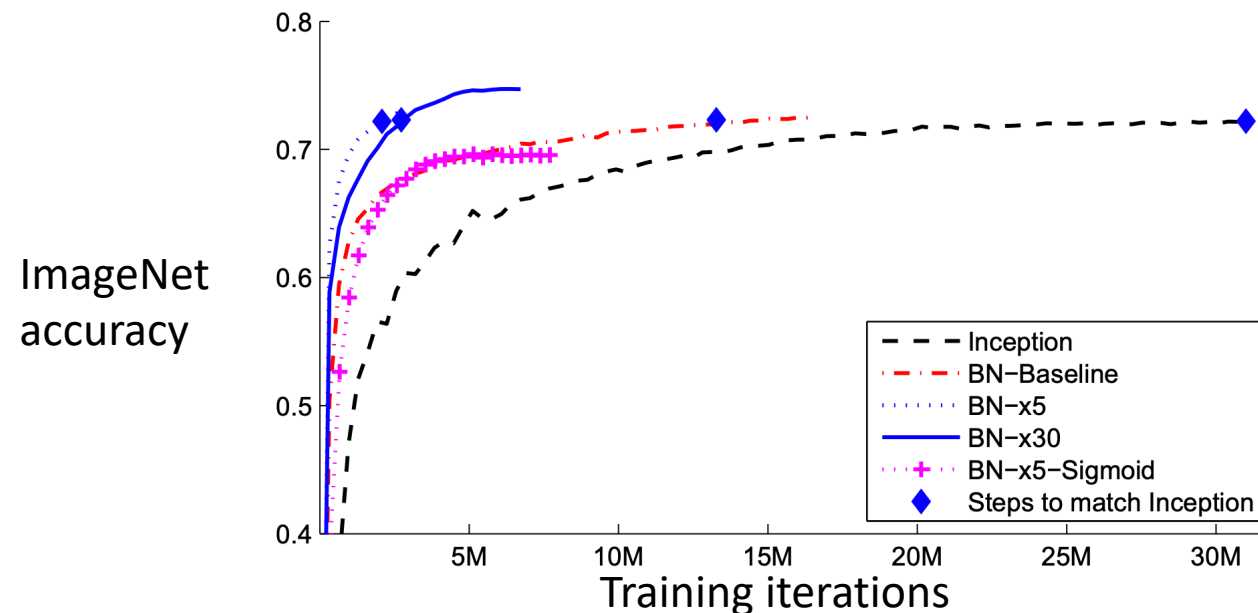
Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

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

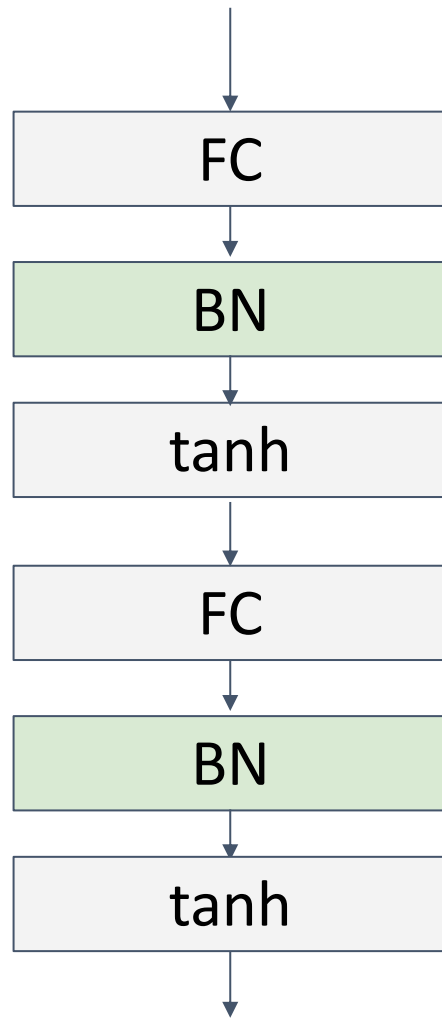
# Batch Normalization



- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!



# Batch Normalization



- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!
- **Not well-understood theoretically (yet)**
- **Behaves differently during training and testing: this is a very common source of bugs!**

# Layer Normalization

**Batch Normalization** for fully-connected networks

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

Normalize



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

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

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

**Layer Normalization** for fully-connected networks

Same behavior at train and test!

Used in RNNs, Transformers

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

Normalize



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

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

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

# Instance Normalization

**Batch Normalization** for convolutional networks

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

Normalize



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

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

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

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

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

Normalize



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

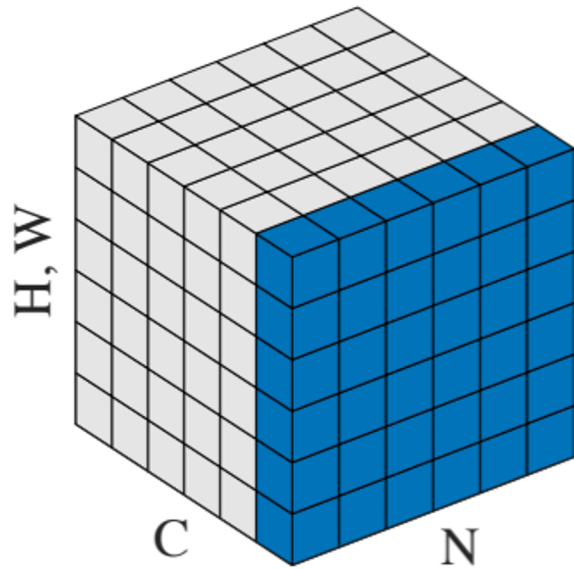
$$\boldsymbol{\gamma}, \boldsymbol{\beta} : \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1}$$

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

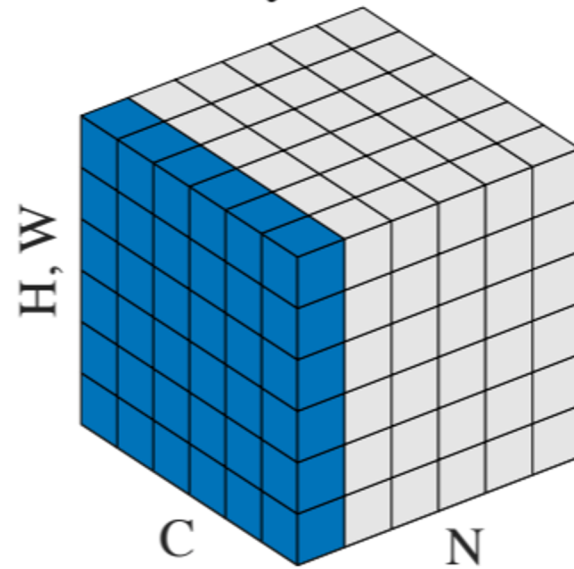


# Comparison of Normalization Layers

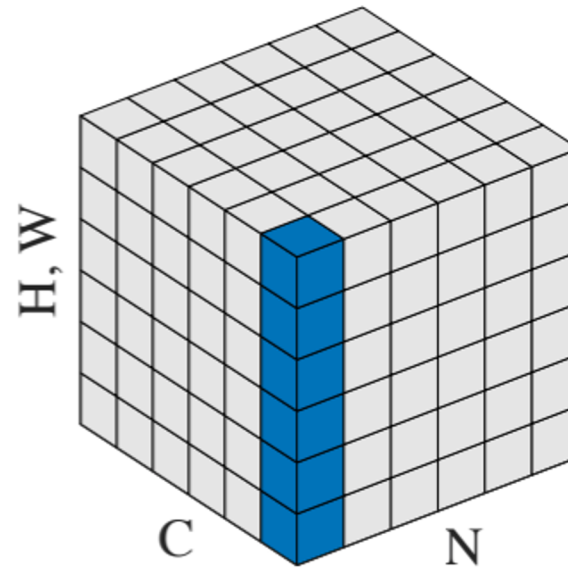
Batch Norm



Layer Norm



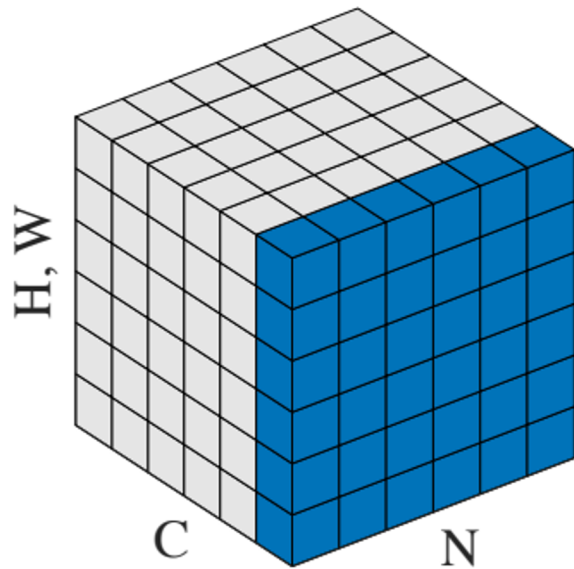
Instance Norm



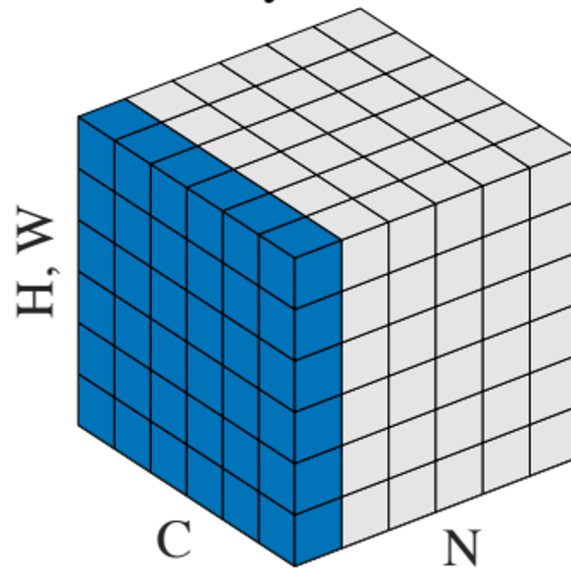
Wu and He, "Group Normalization", ECCV 2018

# Group Normalization

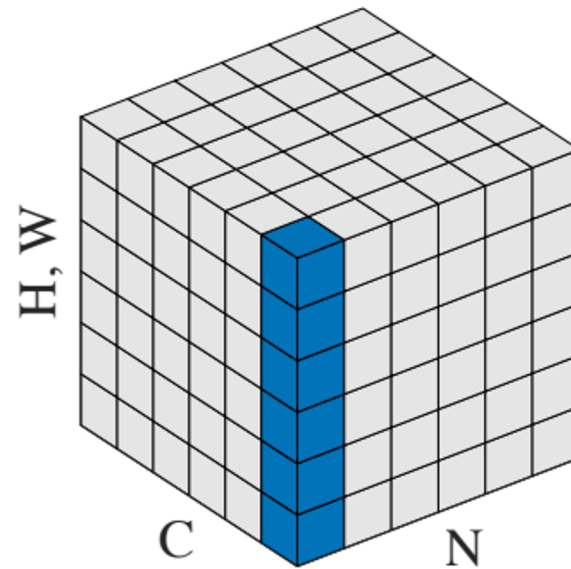
Batch Norm



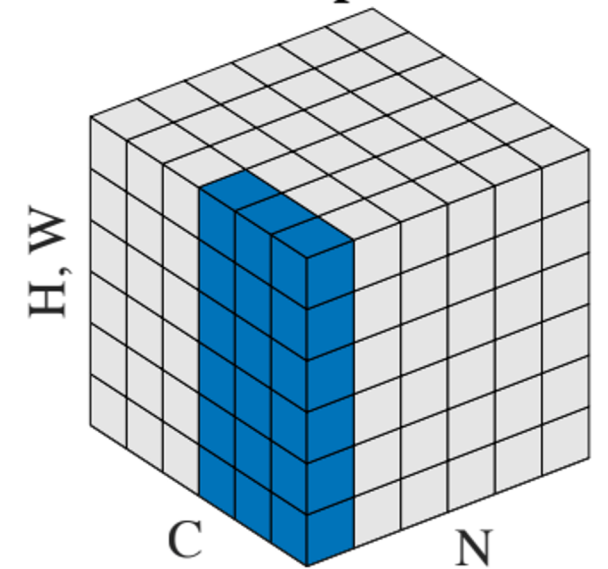
Layer Norm



Instance Norm



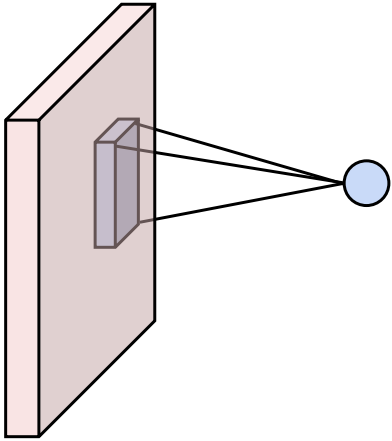
**Group Norm**



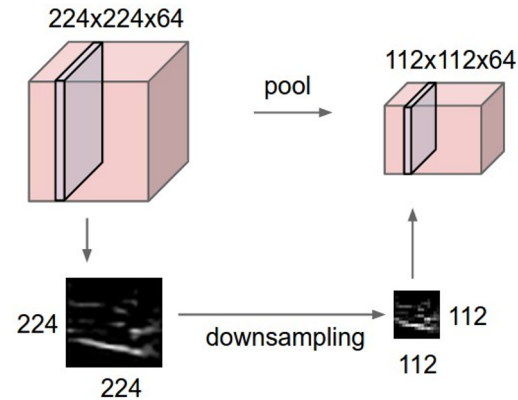
Wu and He, "Group Normalization", ECCV 2018

# Components of a Convolutional Network

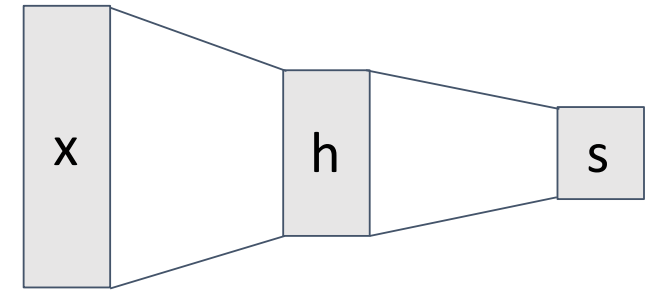
## Convolution Layers



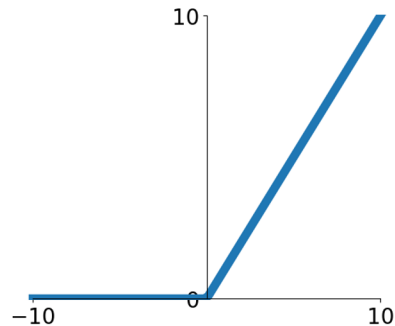
## Pooling Layers



## Fully-Connected Layers



## Activation Function

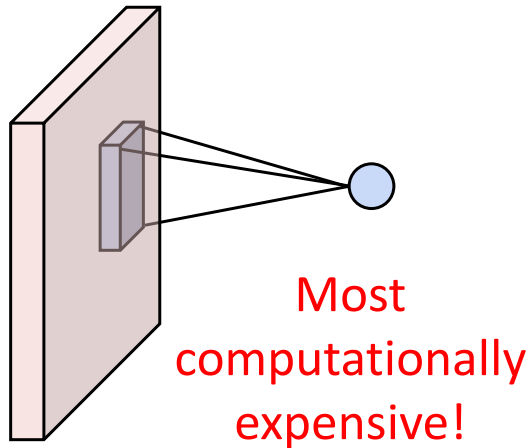


## Normalization

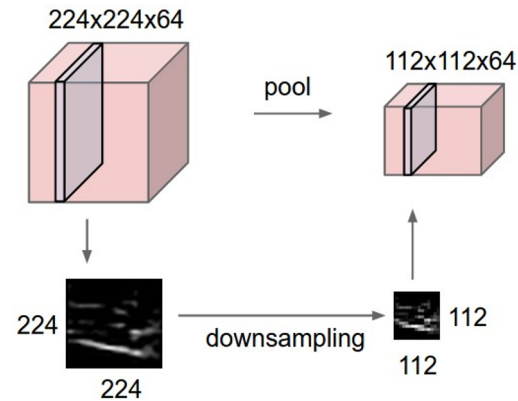
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

# Components of a Convolutional Network

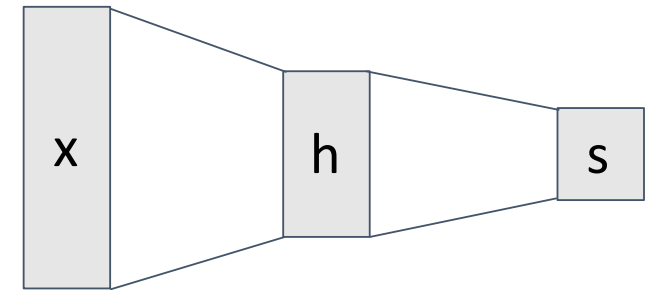
## Convolution Layers



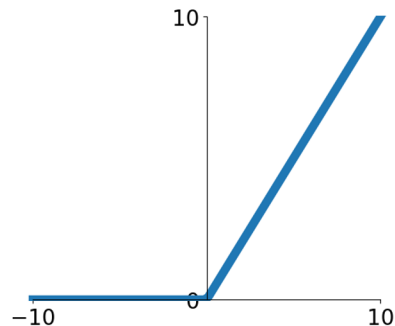
## Pooling Layers



## Fully-Connected Layers



## Activation Function

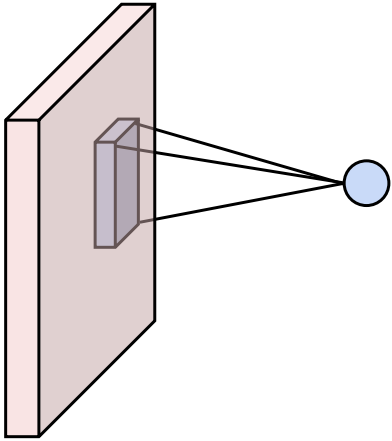


## Normalization

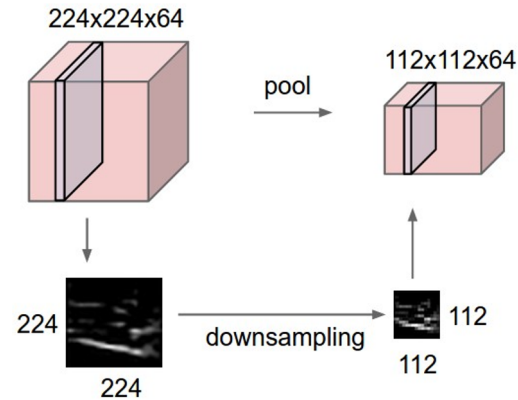
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

# Summary: Components of a Convolutional Network

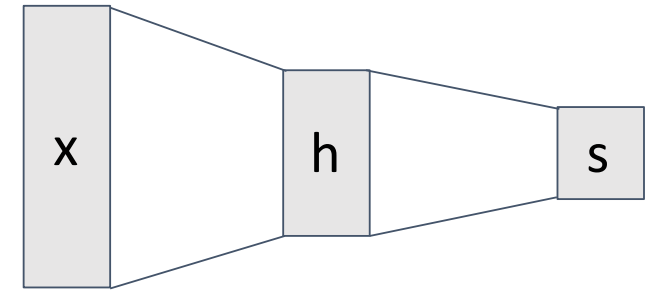
## Convolution Layers



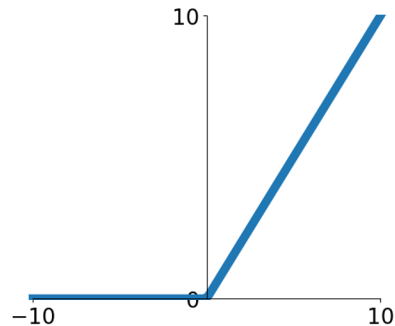
## Pooling Layers



## Fully-Connected Layers



## Activation Function

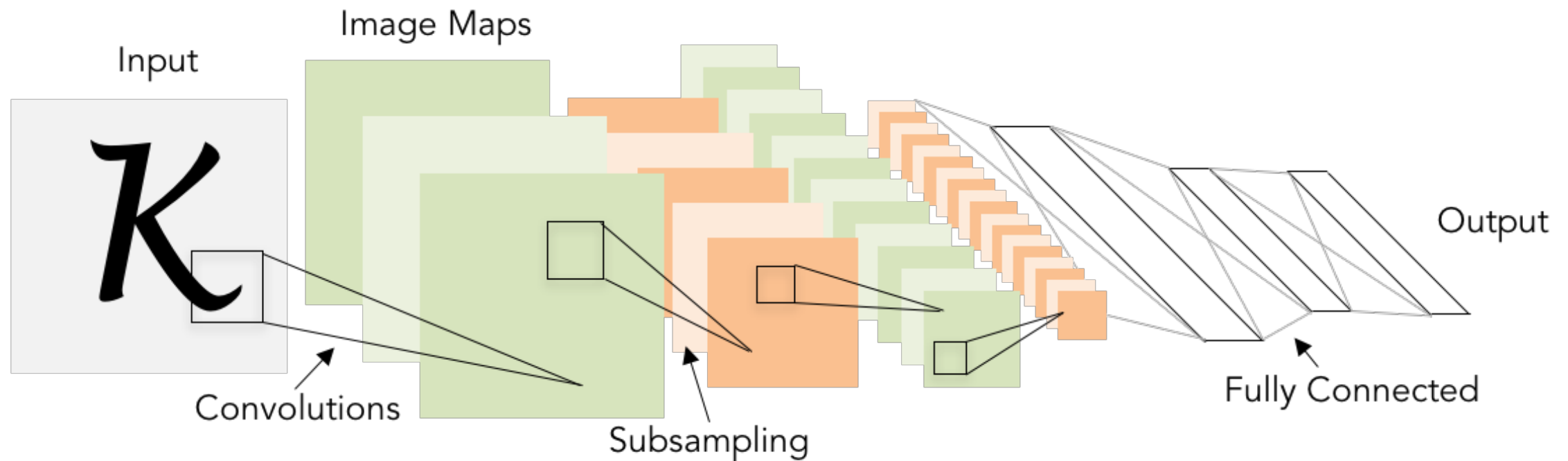


## Normalization

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

# Summary: Components of a Convolutional Network

**Problem:** What is the right way to combine all these components?



# Convolutional neural networks++

- Training and optimization
- More regularization (dropout, ...)
- Convolutional neural networks
- Pooling
- Batch normalization
- CNN architectures

