

# Convolutional Neural Networks III

CS7150, Spring 2025

Prof. Huaizu Jiang

Northeastern University

Recap

# Padding: 1D Case

Signal



Filter



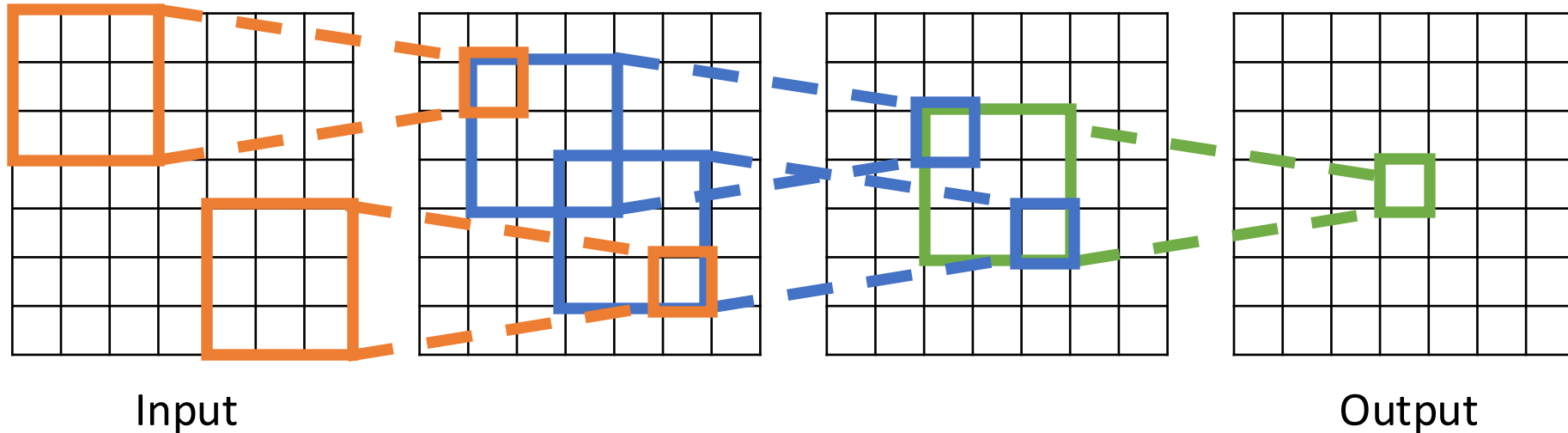
Output



Input: 7 vs output: 7

# Receptive Fields

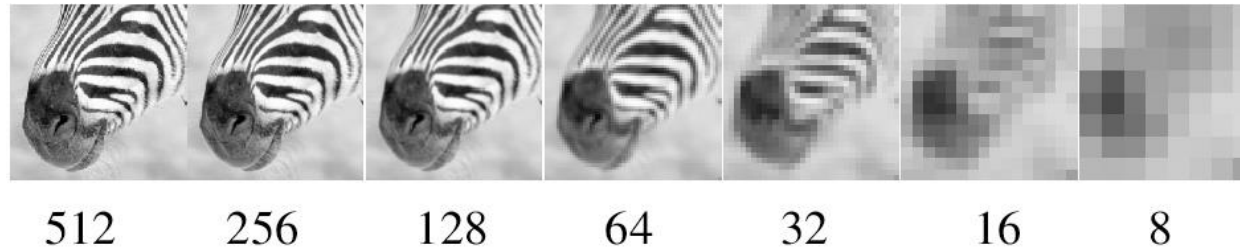
Each successive convolution adds  $K - 1$  to the receptive field size  
With  $L$  layers the receptive field size is  $1 + L * (K - 1)$



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

**Solution:** Downsample inside the network

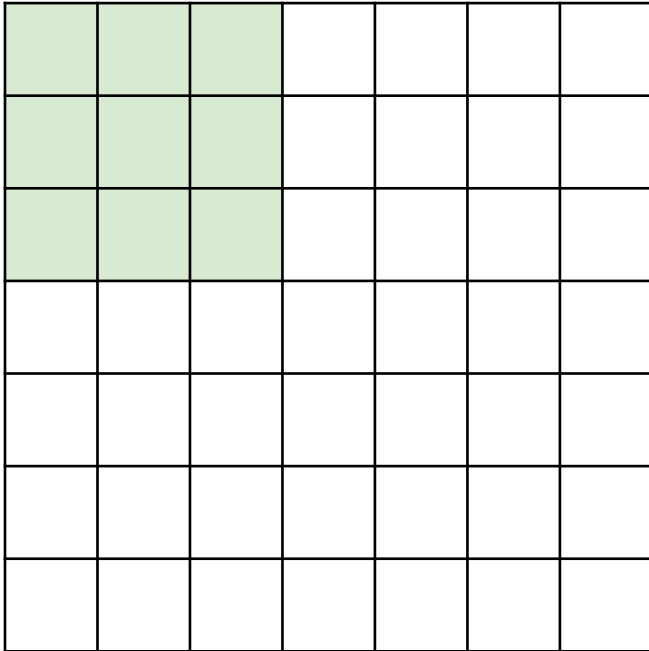
# Idea: Image/Feature Pyramid



50 px



# Strided Convolution

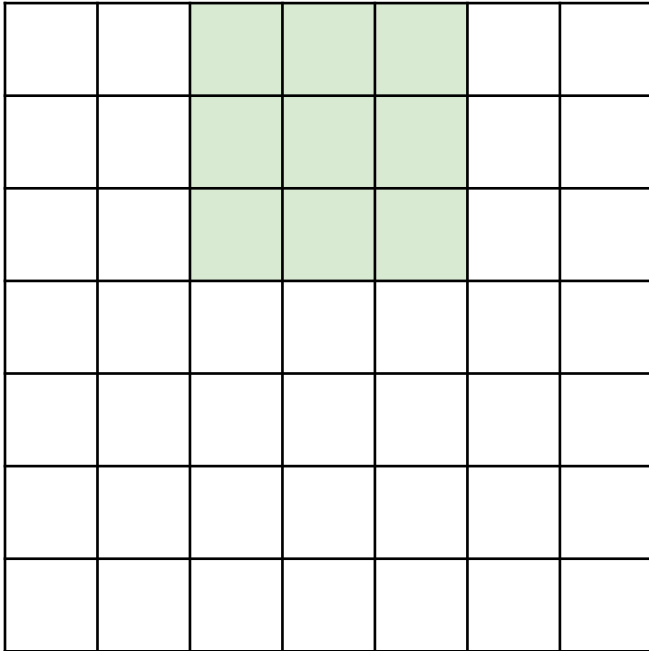


Input: 7x7

Filter: 3x3

Stride: 2

# Strided Convolution

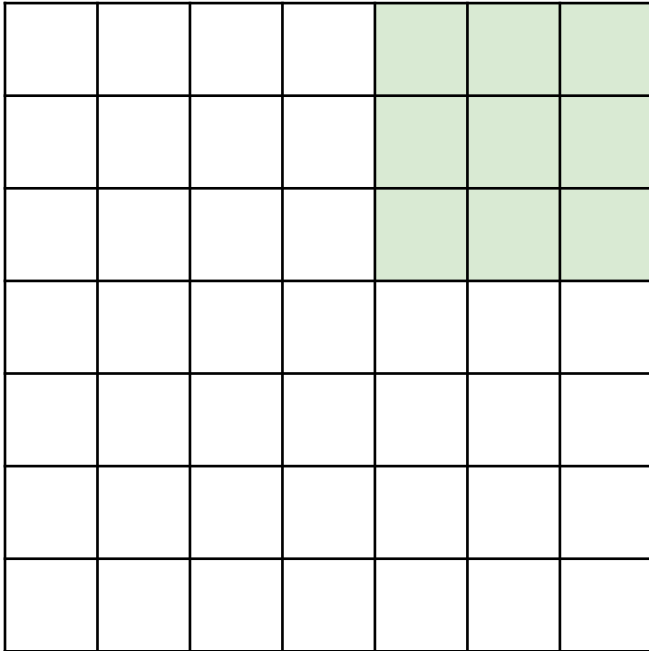


Input: 7x7

Filter: 3x3

Stride: 2

# Strided Convolution



Input: 7x7

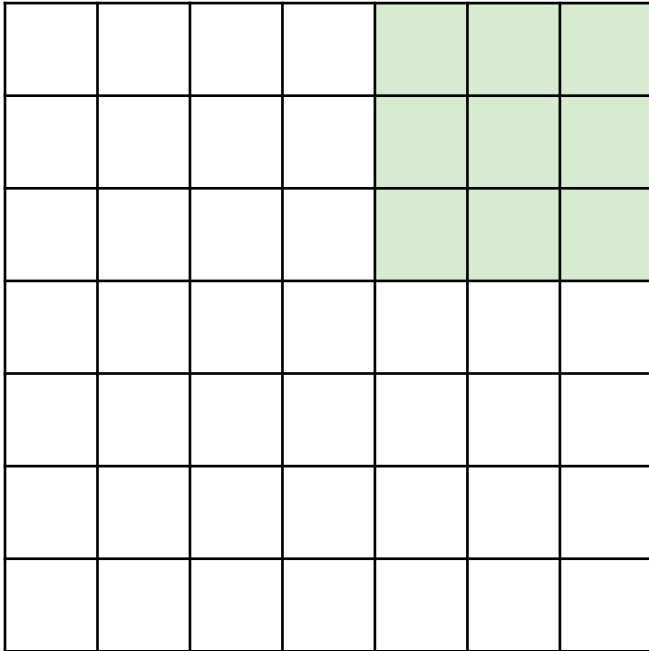
Filter: 3x3

Stride: 2

Output: 3x3



# Strided Convolution



Input: 7x7

Filter: 3x3

Output: 3x3

Stride: 2

In general:

Input:  $W$

Filter:  $K$

Padding:  $P$

Stride:  $S$

Output:  $\text{Ceil}((W - K + 1 + 2P) / S)$

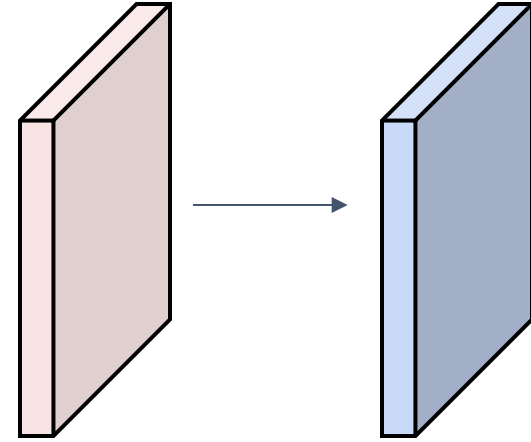
# Today's Class

- More about convolution layer
- Pooling layer
- Batch normalization layer

# Convolution Example

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

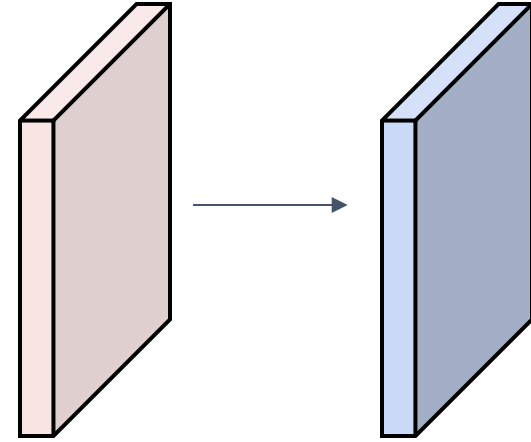
Output volume size: ?



# Convolution Example

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

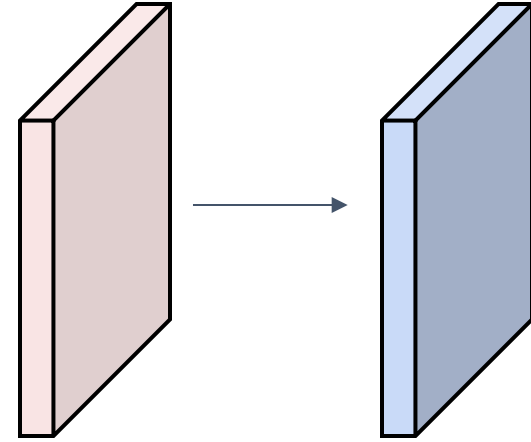
Output volume size:  
 $\text{Ceil}((32 + 2 * 2 - 5 + 1) / 1) = 32$  spatially, so  
10 x 32 x 32



# 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

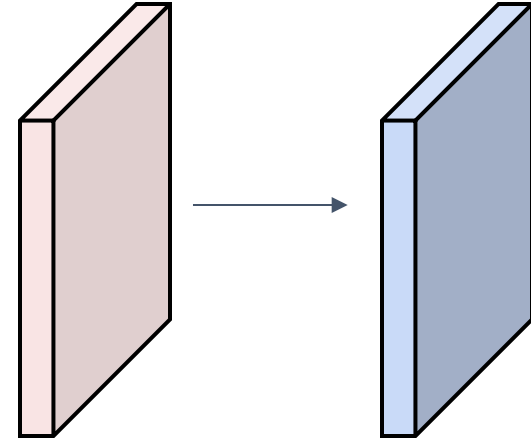
Input volume: 3 x 32 x 32  
10 5x5 filters with stride 1, pad 2

Output volume size: 10 x 32 x 32

Number of learnable parameters: 760

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

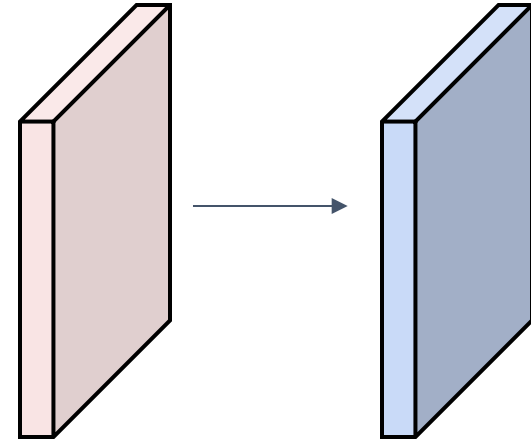
10 filters, so total is  $10 * 76 = 760$



# 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

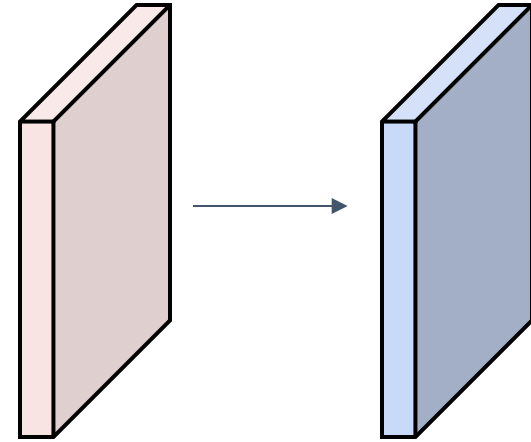
Input volume: **3** x 32 x 32  
10 **5x5** filters with stride 1, pad 2

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 =  $75 * 10240 = 768K$





# 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:

- $H' = \text{Ceil}((H - K + 2P + 1) / S)$
- $W' = \text{Ceil}((W - K + 2P + 1) / S)$

# 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:

- $H' = \text{Ceil}((H - K + 2P + 1) / S)$
- $W' = \text{Ceil}((W - K + 2P + 1) / S)$

**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)

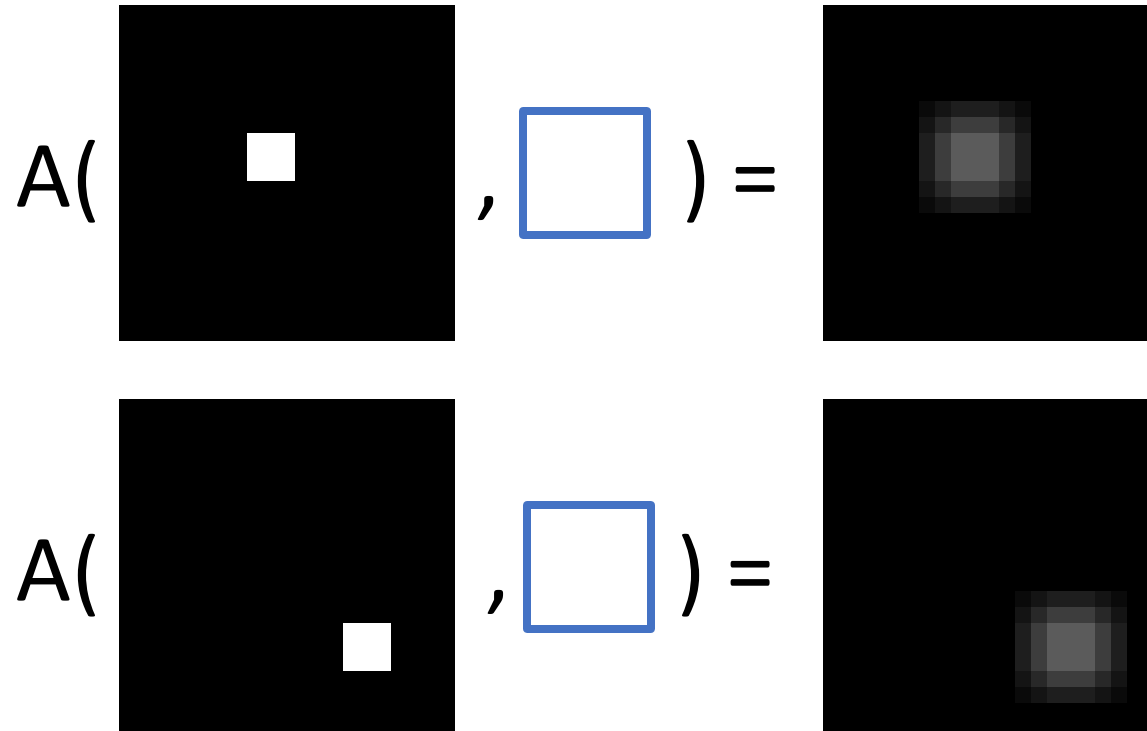
$K = 3, P = 1, S = 2$  (Downsample by 2)

# Properties: Shift-Invariant

Assume:  $I$  image,  $f$  filter

**Shift-invariant:**  $\text{shift}(\text{apply}(I, f)) = \text{apply}(\text{shift}(I, f))$

Intuitively: only depends on filter neighborhood



# Forward pass of a convolution layer

```
x_padded = pad_input(x)
H_out, W_out = compute_output_dimension()
N, C, H, W = x.shape
F, C, HH, WW = w.shape
# shape of (F, C'), where  $C' = C * HH * WW$ 
w_reshape = reshape_to_F_Cp(w).T
for i in range(W_out):
    for j in range(H_out):
        startx, starty = get_top_left_position()
        # shape of (N, C, HH, WW)
        x_data = get_patch(x_padded, startx, starty, HH, WW)
        # shape of (N, C'), where  $C' = C * HH * WW$ 
        x_data_reshape = reshape_to_N_Cp(x_data)
        # shape of each location is (N, F)
        out[:, :, j, i] = dot_product(x_data_reshape, w_reshape) + b
```

# Gradients of weights in convolution

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

input

1	2	3
4	5	6
7	8	9

conv kernel

# Gradients of weights in convolution

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

input

1	2	3
4	5	6
7	8	9

conv kernel

Gradients for this particular weight

# Gradients of weights in convolution

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

input

1	2	3
4	5	6
7	8	9

conv kernel

Gradients for this particular weight

Sum the gradients **from different patches** (determined by the padding and stride).

The same idea applies to the bias term.

# Gradients of input in convolution

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

input

1	2	3
4	5	6
7	8	9

conv kernel

Gradients for this particular input

(why we ever need to compute the gradient for the input?)



# Gradients of input in convolution

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

input

1	2	3
4	5	6
7	8	9

conv kernel

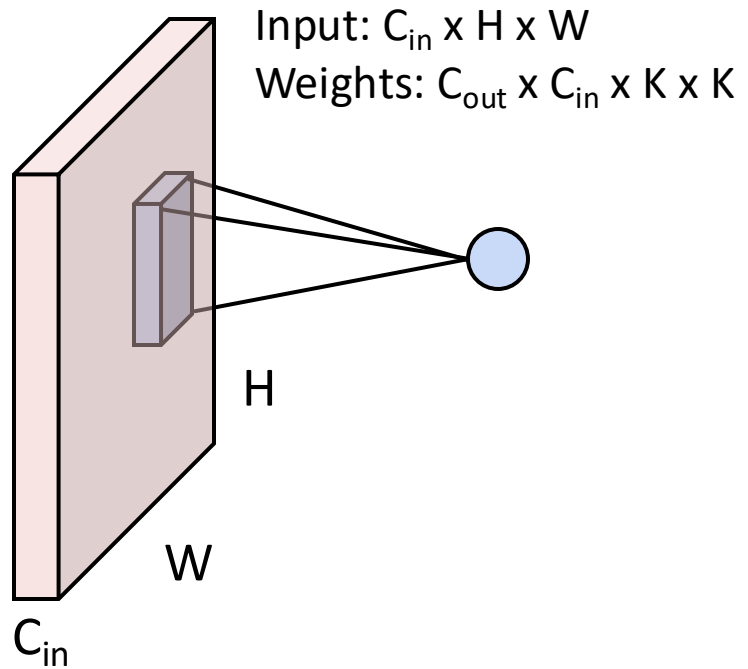
10	11	12
13	14	15
16	17	18

Gradients for this particular input

Sum the gradients from different convolution kernels.

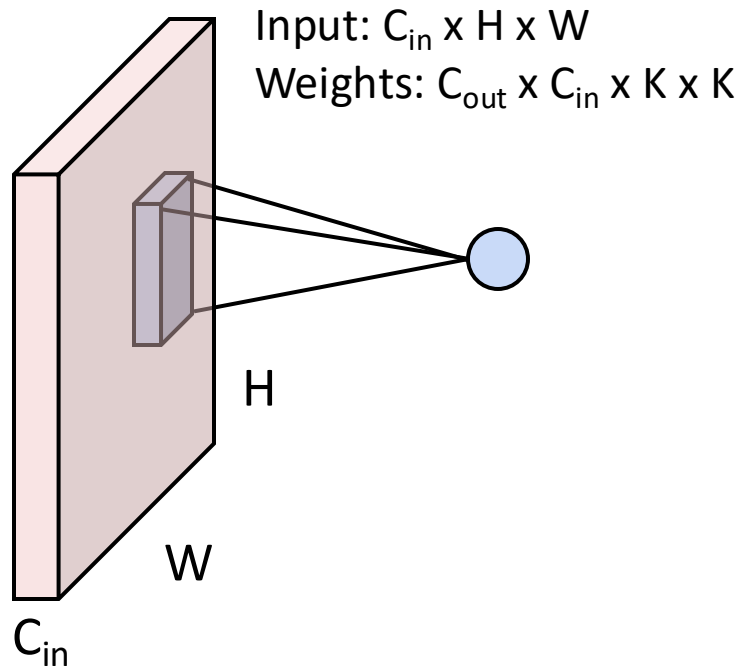
# Other types of convolution

So far: 2D Convolution

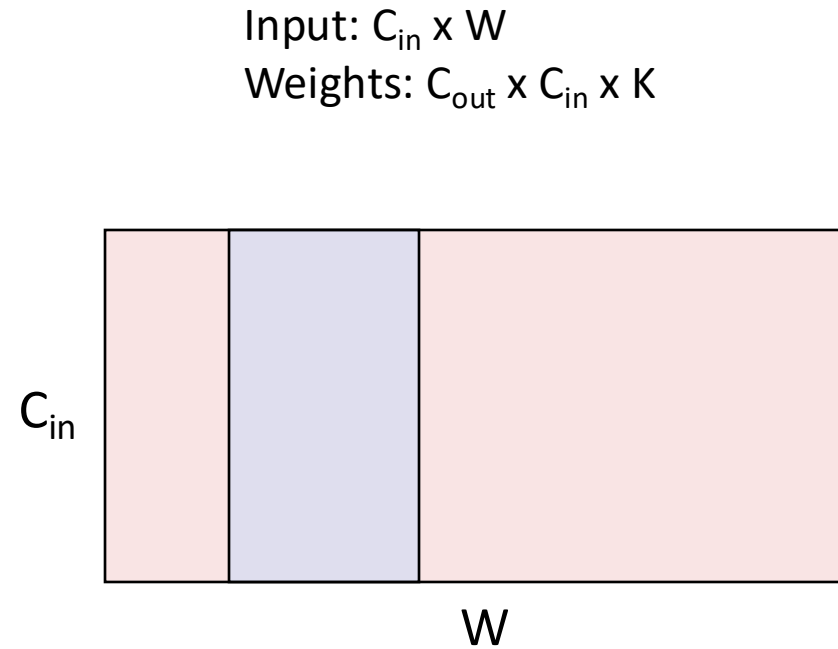


# Other types of convolution

## So far: 2D Convolution

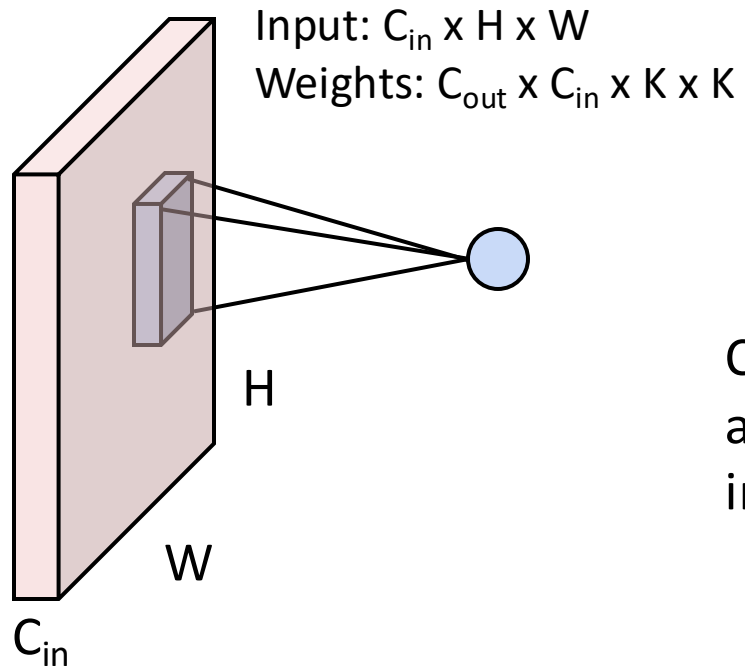


## 1D Convolution

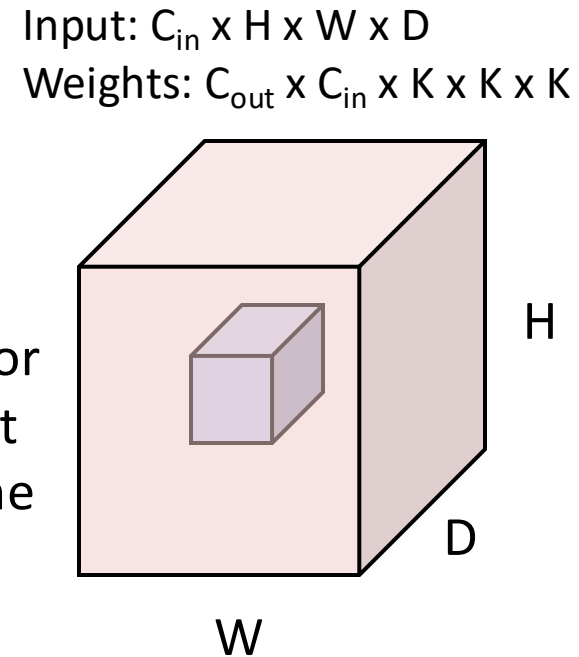


# Other types of convolution

## So far: 2D Convolution



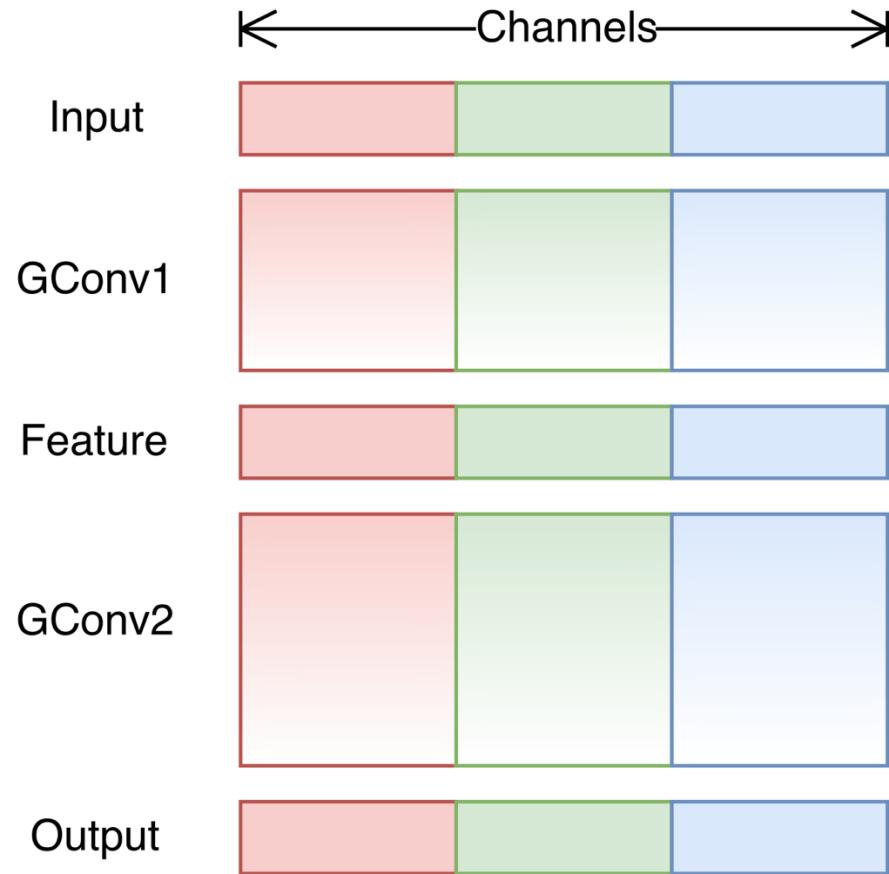
## 3D Convolution



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

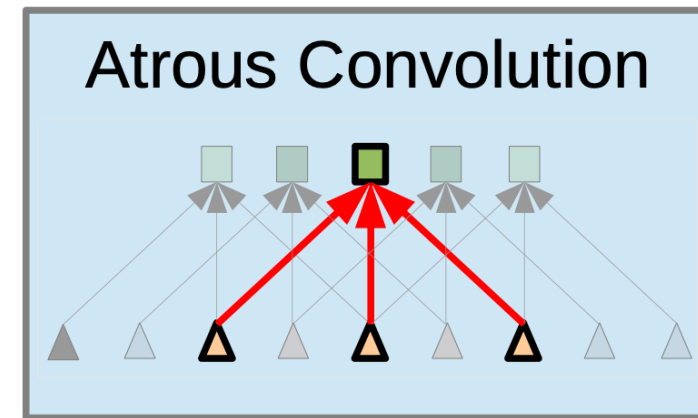
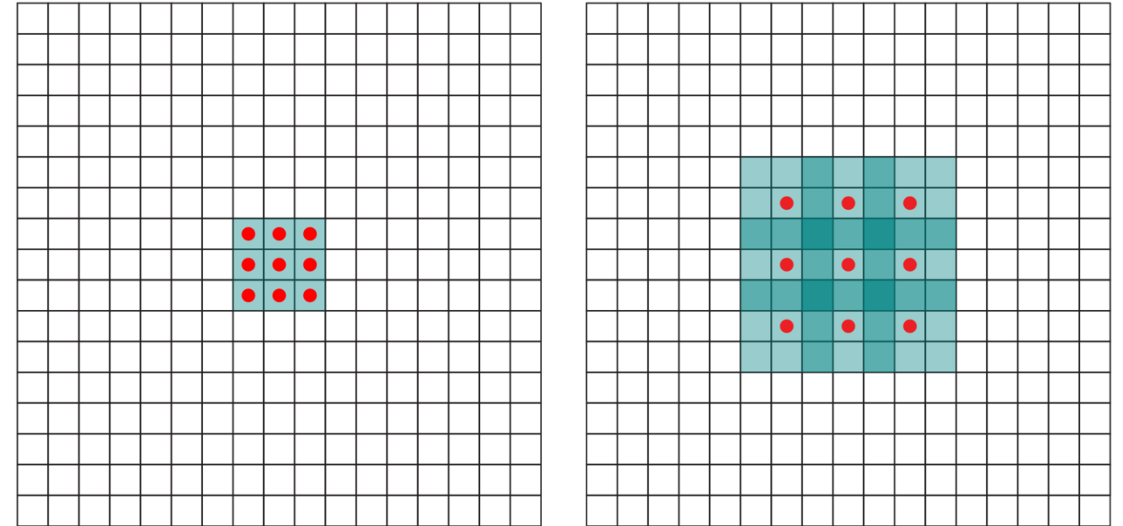
# Other types of convolution

## Group-based convolution



[Zhang et al., CVPR 2017]

## Dilated convolution



[Yu and Koltun. ICLR 2016]

[Chen et al., ICLR 2015]

# 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:

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

# PyTorch Convolution Layers

## Conv2d

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

[\[SOURCE\]](#)

## Conv1d

**CLASS** `torch.nn.Conv1d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros')`

[\[SOURCE\]](#) 

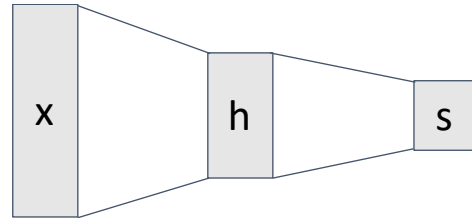
## Conv3d

**CLASS** `torch.nn.Conv3d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros')`

[\[SOURCE\]](#)

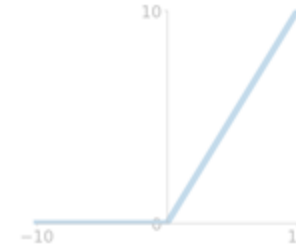
# Components of a Convolutional Network

Fully-Connected Layers



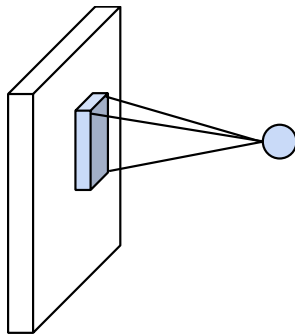
$$y = Wx + b$$

Activation Function

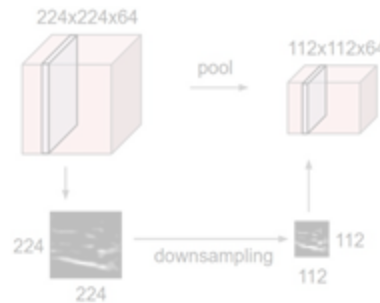


$$y = \max(0, x)$$

Convolution Layers



Pooling Layers



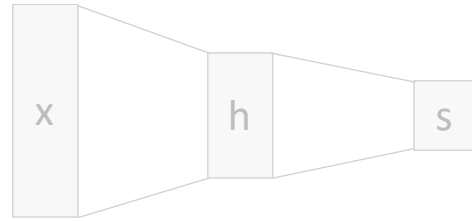
Normalization

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



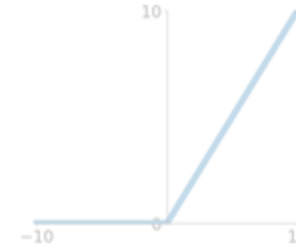
# Components of a Convolutional Network

Fully-Connected Layers



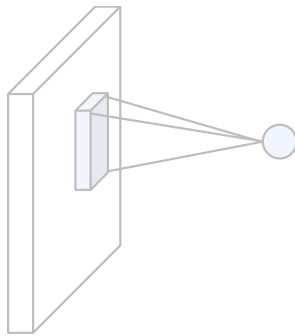
$$y = Wx + b$$

Activation Function

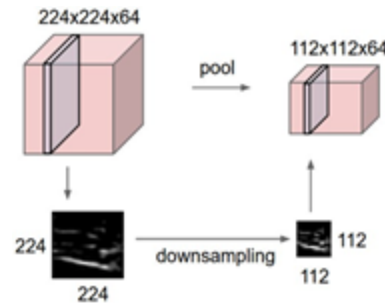


$$y = \max(0, x)$$

Convolution Layers



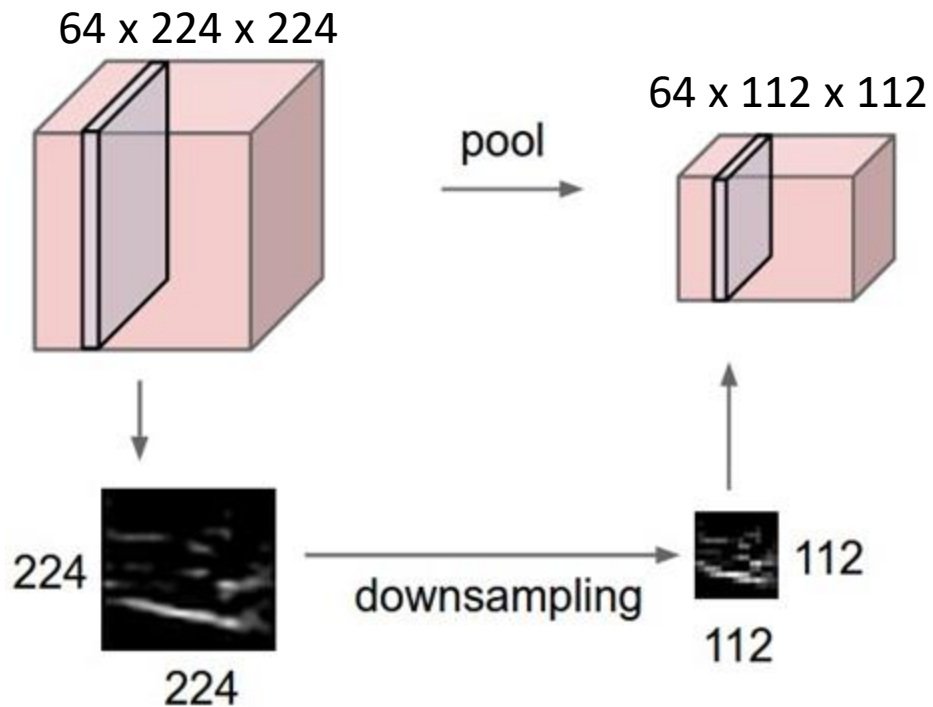
Pooling Layers



Normalization

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

# Pooling Layers: Downsampling



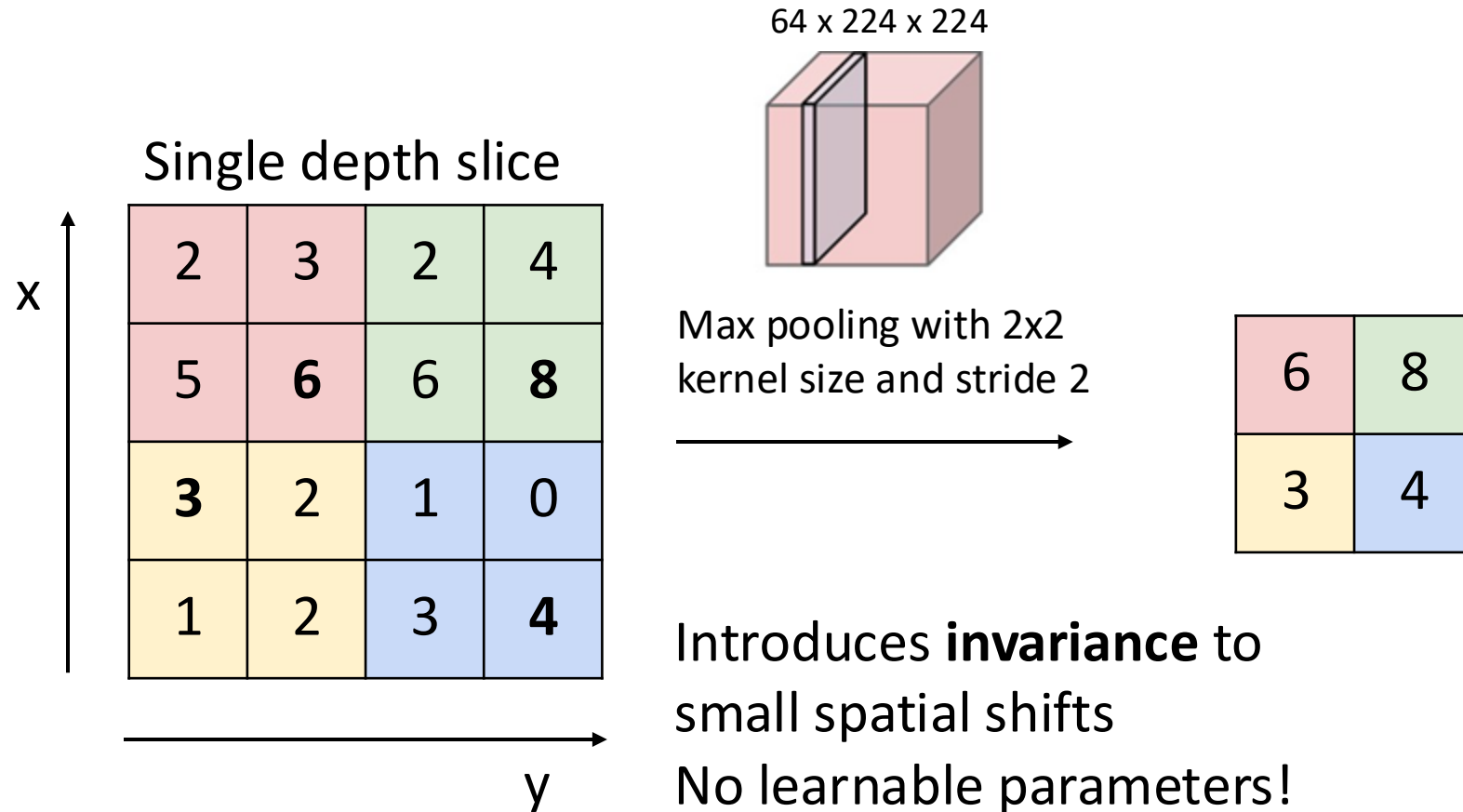
## Hyperparameters:

Kernel Size

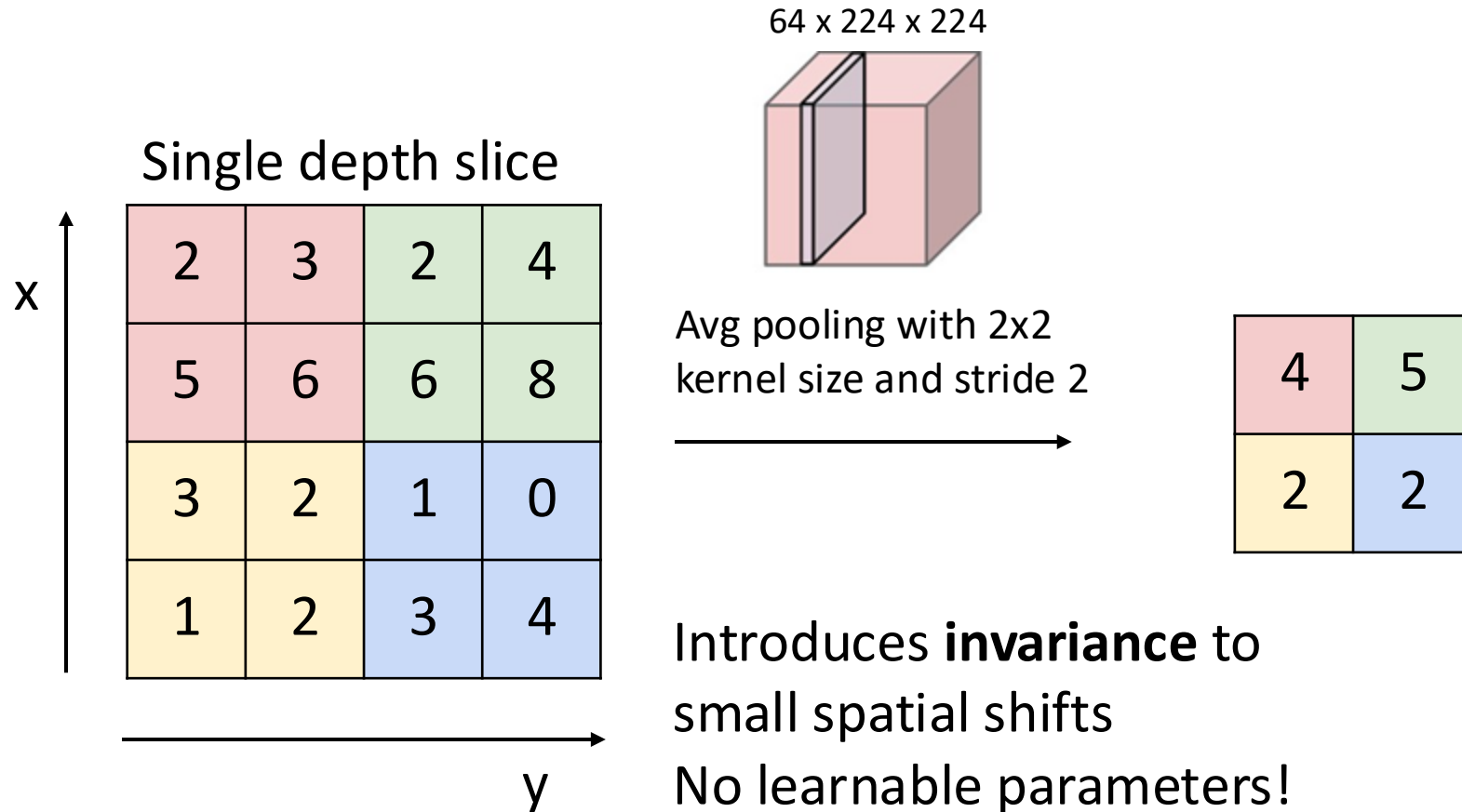
Stride

Pooling function

# Max Pooling



# Average Pooling



# 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

- $H' = \text{Ceil}((H - K + 2P + 1) / S)$
- $W' = \text{Ceil}((W - K + 2P + 1) / S)$

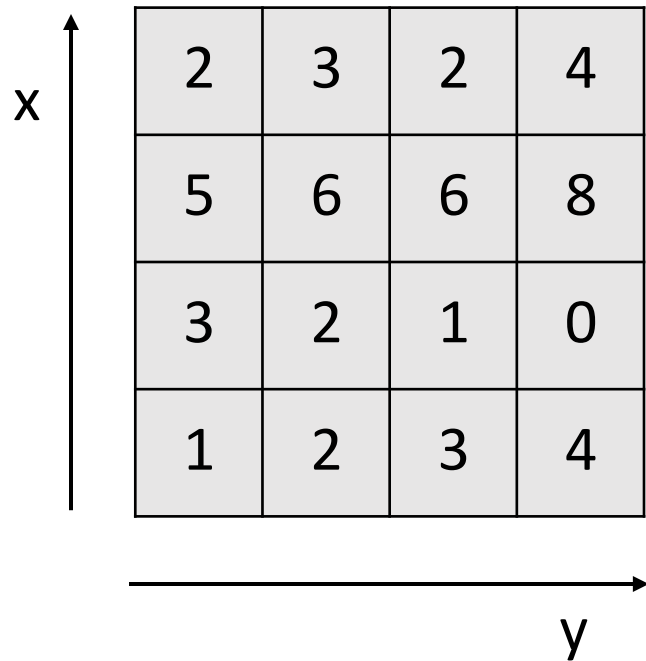
**Learnable parameters:** None!

Common settings:

max,  $K = 2$ ,  $S = 2$

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

# Global Average Pooling



A 4x4 grid of numbers with x and y axes. The x-axis is a vertical arrow on the left, and the y-axis is a horizontal arrow at the bottom. The grid contains the following values:

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



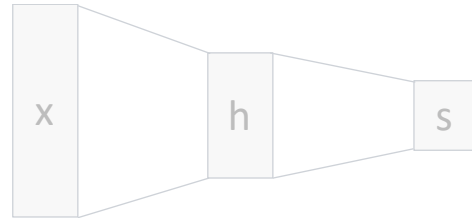
3.25
------

Average per channel (1 channel here).

Gradients?

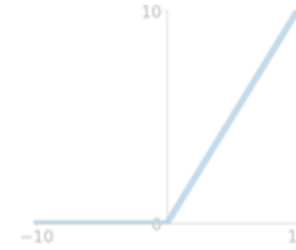
# Components of a Convolutional Network

Fully-Connected Layers



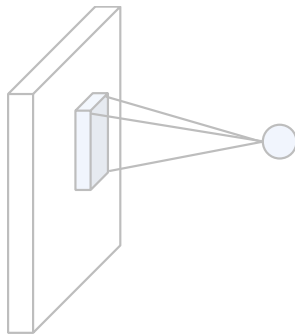
$$y = Wx + b$$

Activation Function

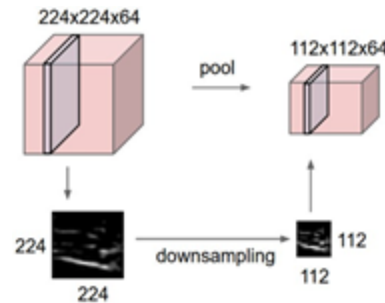


$$y = \max(0, x)$$

Convolution Layers



Pooling Layers

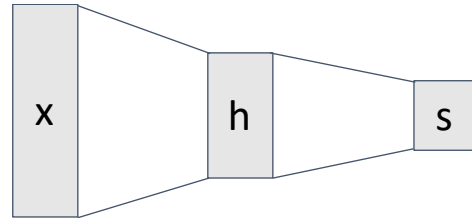


Normalization

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

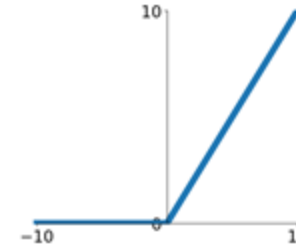
# Components of a Convolutional Network

Fully-Connected Layers



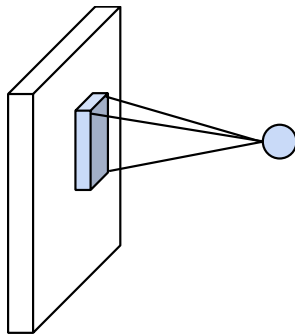
$$y = Wx + b$$

Activation Function

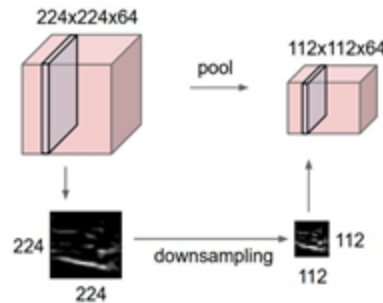


$$y = \max(0, x)$$

Convolution Layers



Pooling Layers



Normalization

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

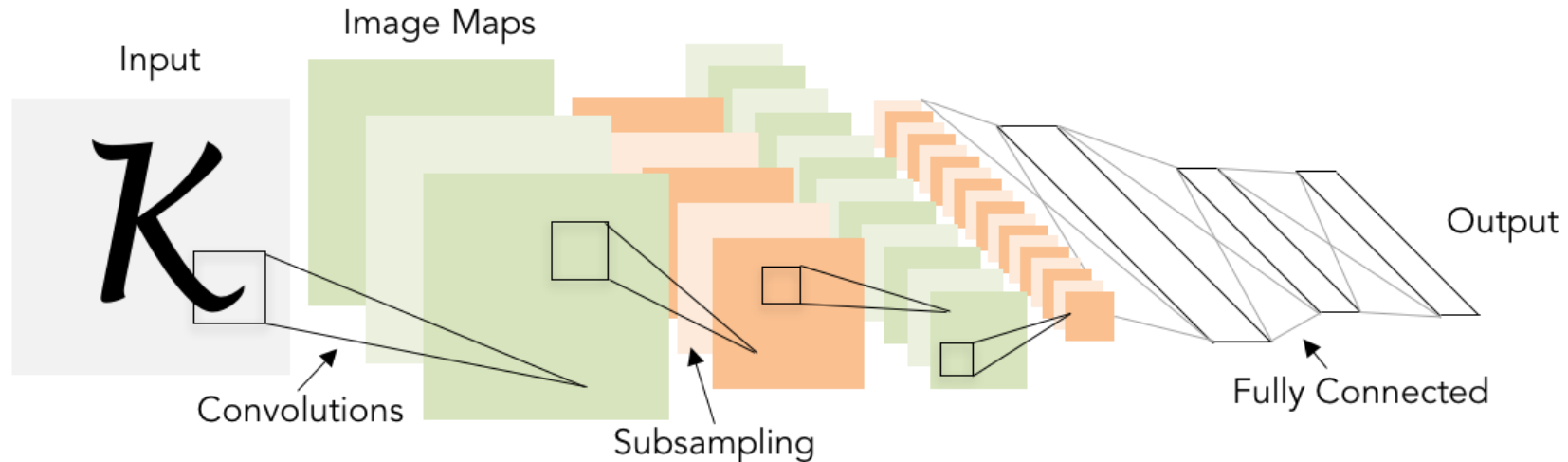


# Convolutional Networks

Classic architecture:

[Conv, ReLU, Pool] x N, flatten, [FC, ReLU] x N, FC

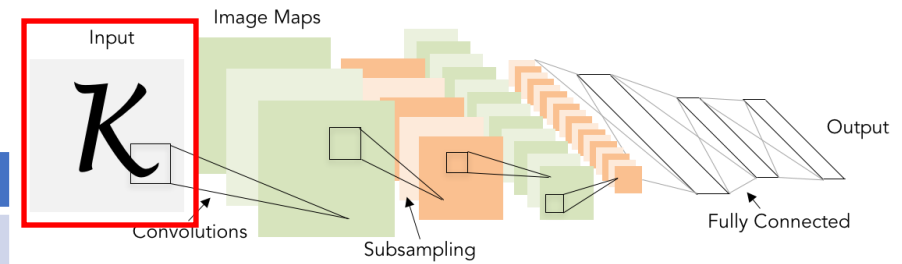
Example: LeNet-5



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

# Example: LeNet-5

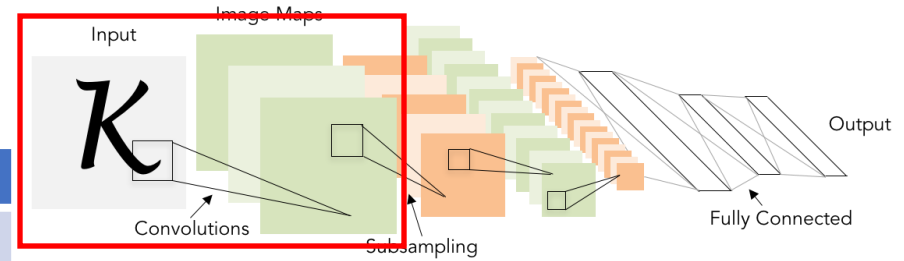
Layer	Output Size	Weight Size
Input	1 x 28 x 28	



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

# Example: LeNet-5

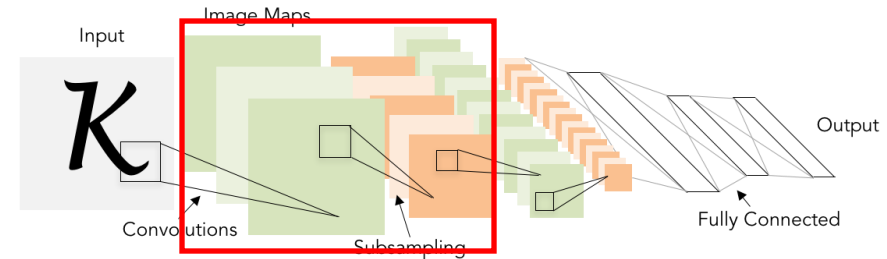
Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv ( $C_{out}=20$ , $K=5$ , $P=2$ , $S=1$ )	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	



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

# Example: LeNet-5

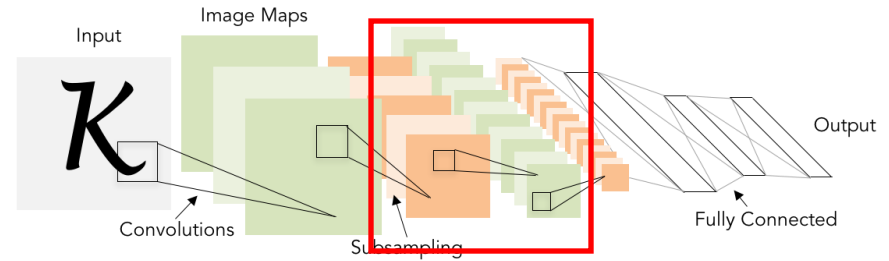
Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv ( $C_{out}=20$ , $K=5$ , $P=2$ , $S=1$ )	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool( $K=2$ , $S=2$ )	20 x 14 x 14	



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

# Example: LeNet-5

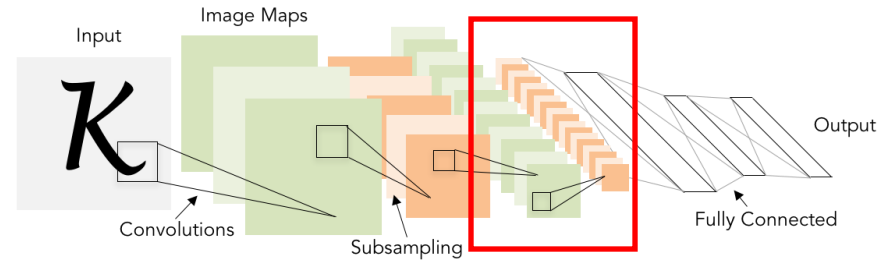
Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv ( $C_{out}=20$ , $K=5$ , $P=2$ , $S=1$ )	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool( $K=2$ , $S=2$ )	20 x 14 x 14	
Conv ( $C_{out}=50$ , $K=5$ , $P=2$ , $S=1$ )	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	



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

# Example: LeNet-5

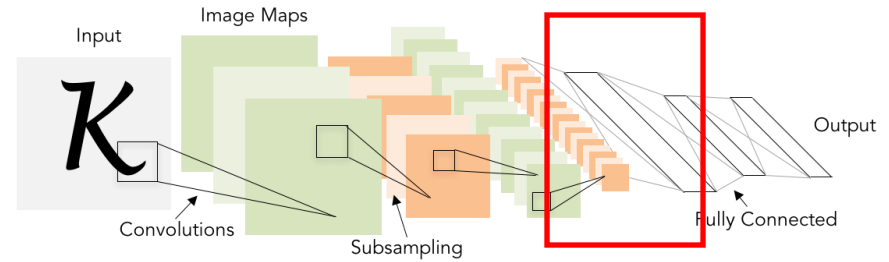
Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv ( $C_{\text{out}}=20$ , $K=5$ , $P=2$ , $S=1$ )	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool( $K=2$ , $S=2$ )	20 x 14 x 14	
Conv ( $C_{\text{out}}=50$ , $K=5$ , $P=2$ , $S=1$ )	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	
MaxPool( $K=2$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



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

# Example: Lenet-5

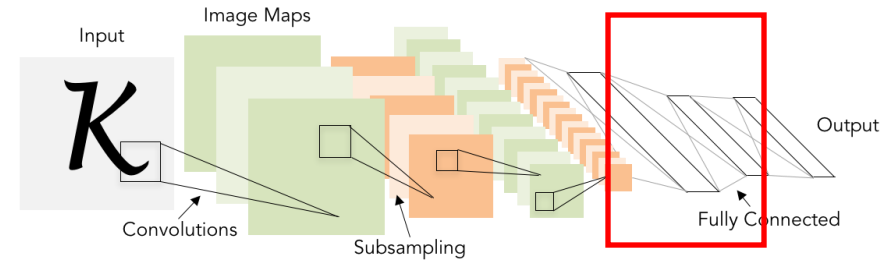
Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv ( $C_{out}=20$ , $K=5$ , $P=2$ , $S=1$ )	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool( $K=2$ , $S=2$ )	20 x 14 x 14	
Conv ( $C_{out}=50$ , $K=5$ , $P=2$ , $S=1$ )	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	
MaxPool( $K=2$ , $S=2$ )	50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

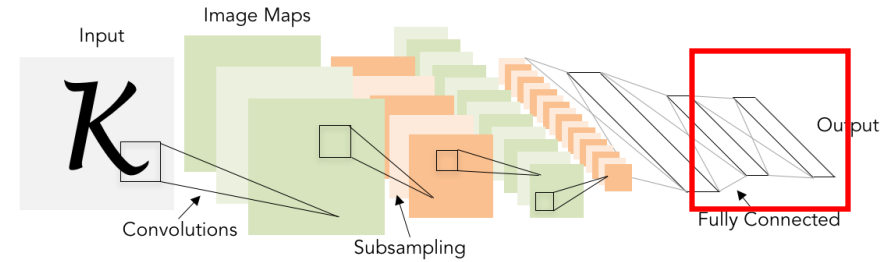


Lecun et al, "Gradient-based learning applied to document recognition", 1998



# Example: Lenet-5

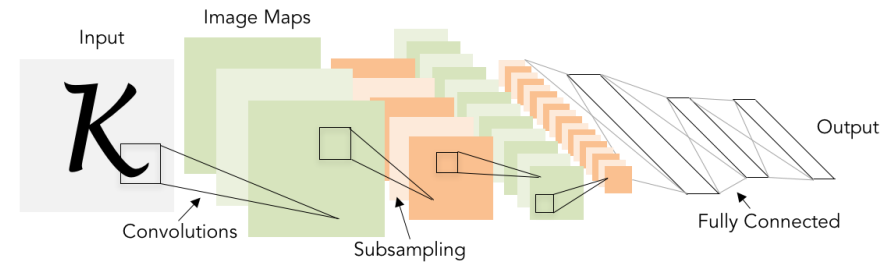
Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv ( $C_{out}=20$ , $K=5$ , $P=2$ , $S=1$ )	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool( $K=2$ , $S=2$ )	20 x 14 x 14	
Conv ( $C_{out}=50$ , $K=5$ , $P=2$ , $S=1$ )	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	
MaxPool( $K=2$ , $S=2$ )	50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



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 ( $C_{\text{out}}=20$ , $K=5$ , $P=2$ , $S=1$ )	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 ( $C_{\text{out}}=50$ , $K=5$ , $P=2$ , $S=1$ )	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



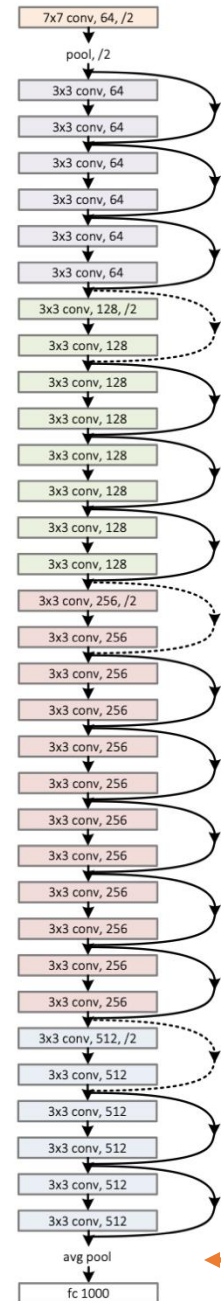
As we go through the network:

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

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

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

# ResNet



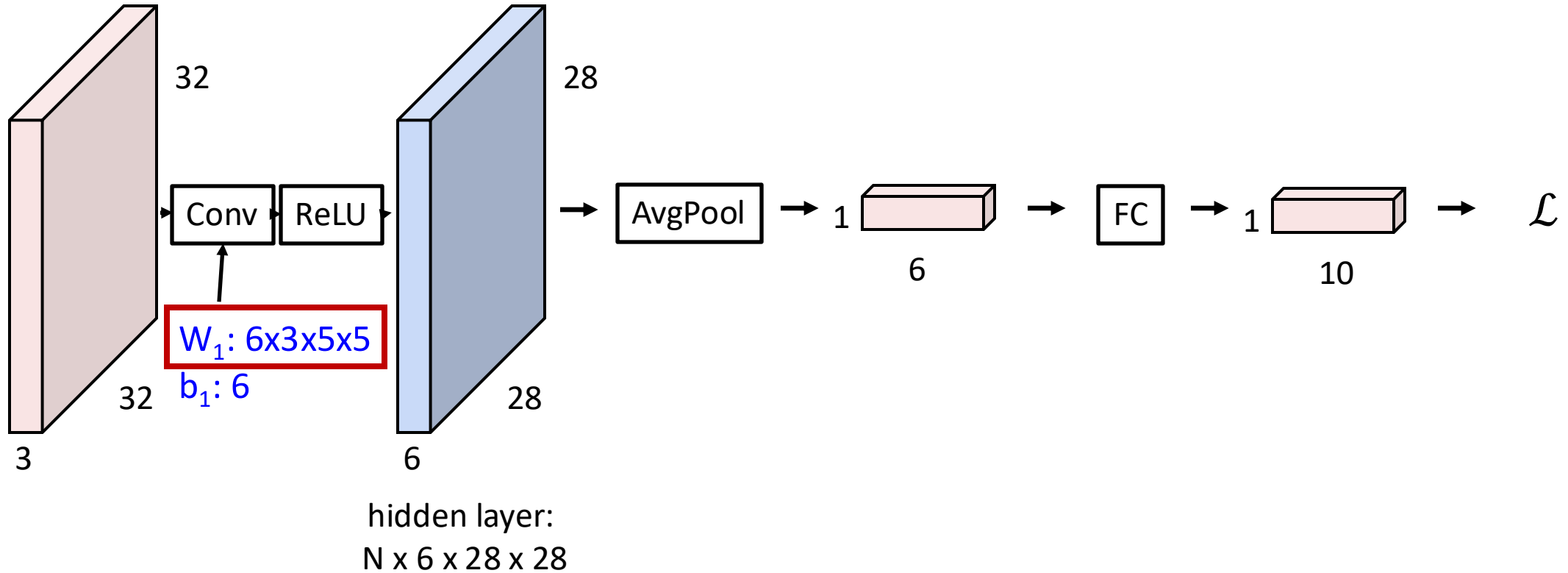
Why is it a better idea than flattening the feature map before the classifier (a fully-connect layer)?

Number of parameters.

Spatial dimension of the input.

global average pooling

# Backpropagation for CNNs

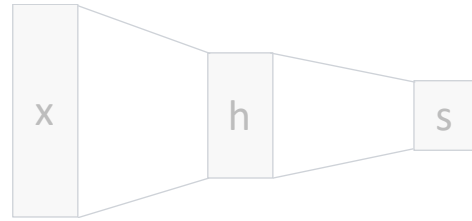


What does its computational graph look like?

Problem: Deep Networks  
very hard to train!

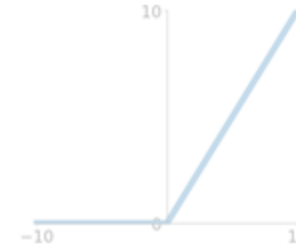
# Components of a Convolutional Network

Fully-Connected Layers



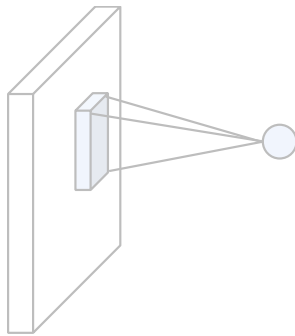
$$y = Wx + b$$

Activation Function

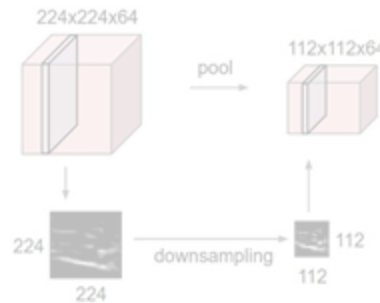


$$y = \max(0, x)$$

Convolution Layers



Pooling Layers



Normalization

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

# Batch Normalization

**Idea:** “Normalize” the outputs of each layer so they have zero mean and unit variance

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

# Batch Normalization

**Idea:** “Normalize” the outputs of each layer so they have zero mean and unit variance

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

We can normalize a batch of activations like this:

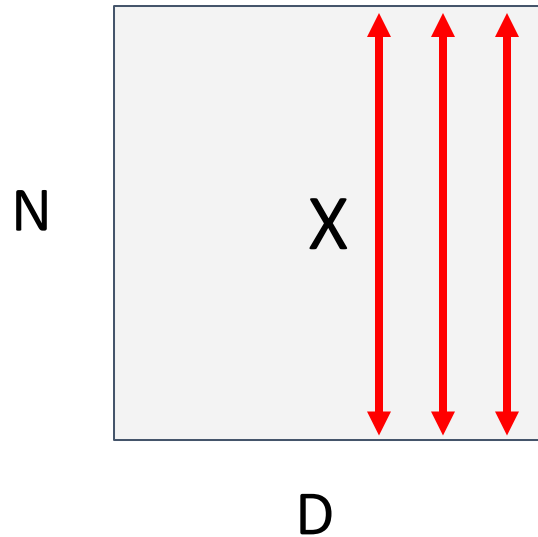
$$\hat{x} = \frac{x - E[x]}{\sqrt{Var[x]}}$$

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



# Batch Normalization

**Input:**  $x \in \mathbb{R}^{N \times D}$



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

Per-channel  
mean, shape is D

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

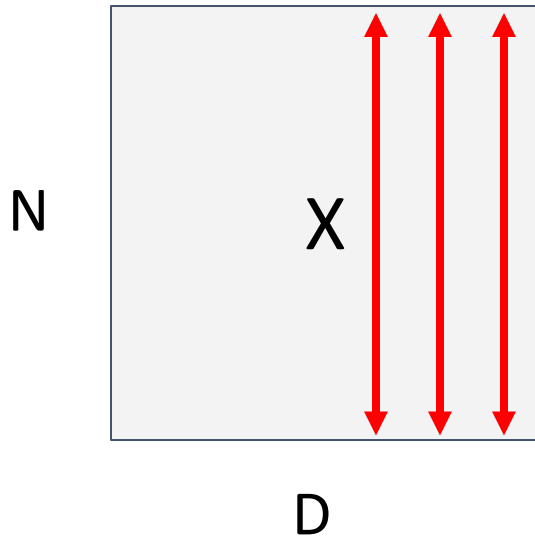
Per-channel  
std, shape is D

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

Normalized x,  
Shape is N x D

# Batch Normalization

**Input:**  $x \in \mathbb{R}^{N \times D}$



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

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 restrictive?

# Batch Normalization

**Input:**  $x \in \mathbb{R}^{N \times D}$

**Learnable scale and shift parameters:**

$$\gamma, \beta \in \mathbb{R}^D$$

Learning  $\gamma = \sigma$ ,  $\beta = \mu$   
will recover the identity  
function (in expectation)

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

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

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output,  
Shape is N x D

# Batch Normalization

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

**Input:**  $x \in \mathbb{R}^{N \times D}$

**Learnable scale and  
shift parameters:**

$$\gamma, \beta \in \mathbb{R}^D$$

Learning  $\gamma = \sigma$ ,  $\beta = \mu$   
will recover the identity  
function (in expectation)

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

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

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output,  
Shape is N x D

# Batch Normalization: Test-Time

**Input:**  $x \in \mathbb{R}^{N \times D}$

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

**Learnable scale and shift parameters:**

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

$$\gamma, \beta \in \mathbb{R}^D$$

Learning  $\gamma = \sigma, \beta = \mu$   
will recover the identity  
function (in expectation)

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

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:**  $x \in \mathbb{R}^{N \times D}$

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

**Learnable scale and shift parameters:**

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

$\gamma, \beta \in \mathbb{R}^D$

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

During testing batchnorm becomes a linear operator!

Can be fused with the previous fully-connected or conv layer

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

# Batch Normalization for ConvNets

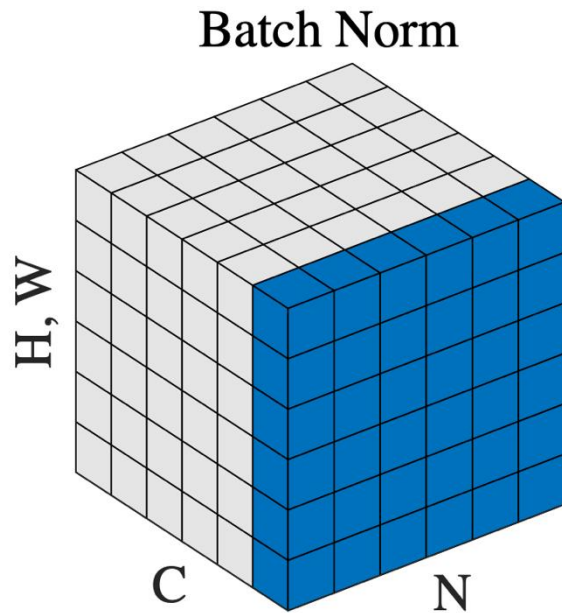
Batch Normalization for  
**fully-connected** networks

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

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

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

# Batch Normalization for ConvNets



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

$$x : N \times C \times H \times W$$

Normalizing

$$\mu, \sigma$$

$$: 1 \times C \times 1 \times 1$$

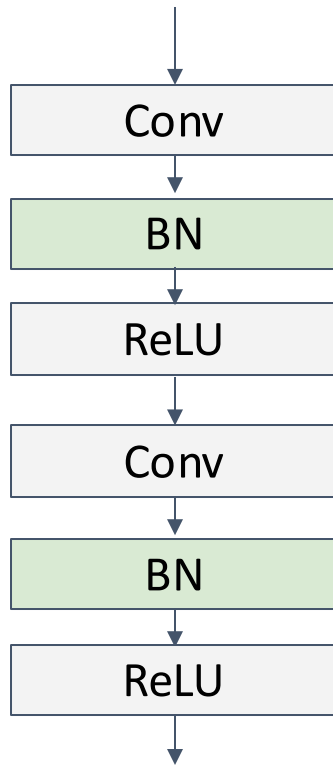
$$\gamma, \beta$$

$$: 1 \times C \times 1 \times 1$$

$$y = \frac{(x - \mu)}{\sigma} \gamma + \beta$$



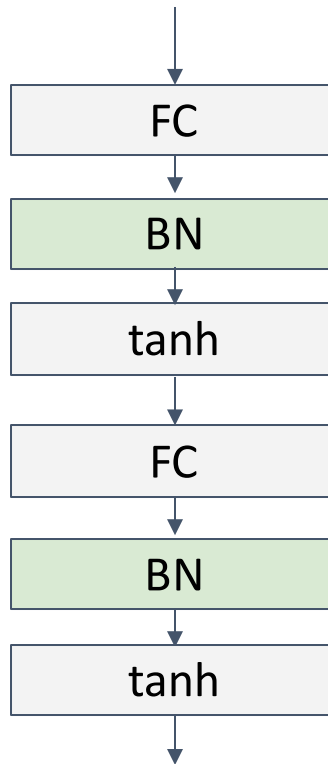
# Batch Normalization



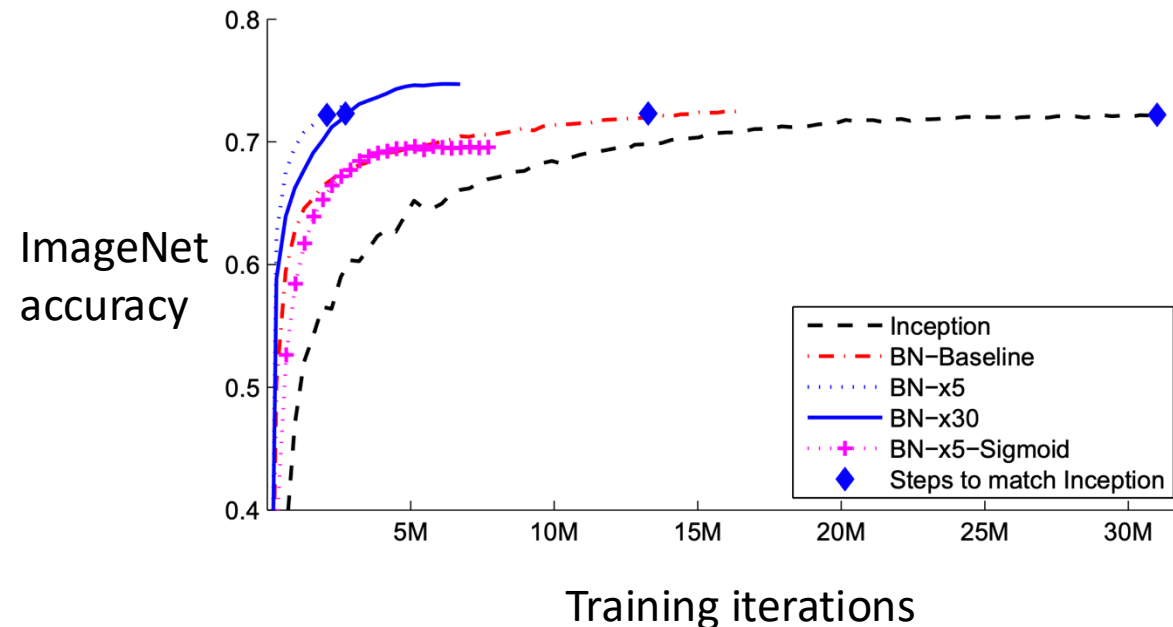
Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

$$\hat{x} = \frac{x - E[x]}{\sqrt{Var[x]}}$$

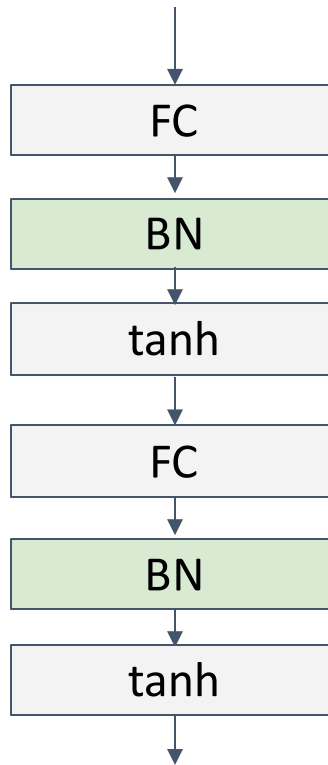
# Batch Normalization



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



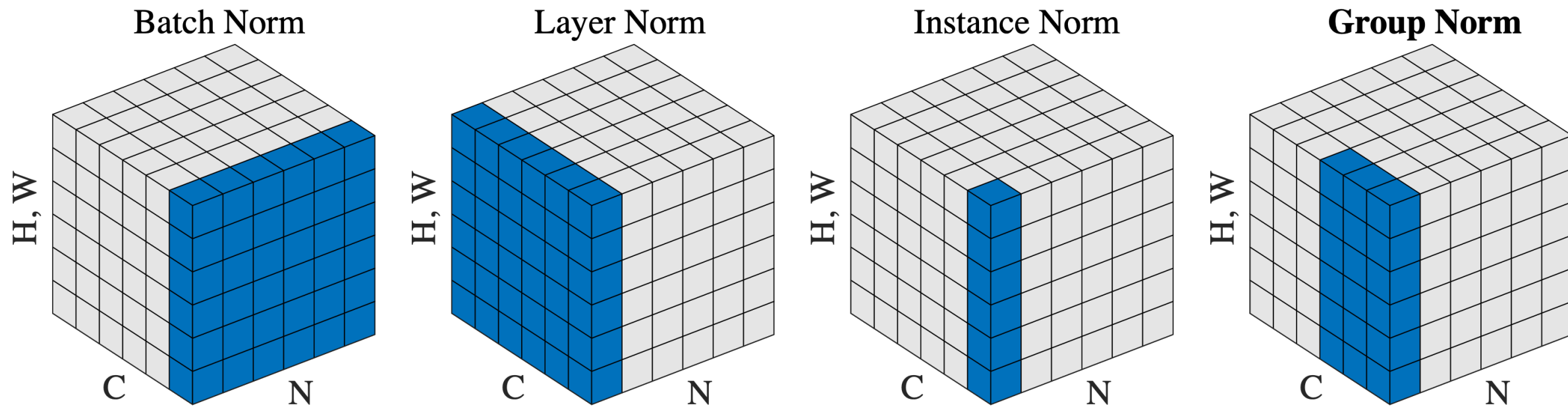
# Batch Normalization



Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Free 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!**

# Different Normalization Layers



[Wu and He. Group Normalization. ECCV 2018. Best paper honorable mention.]

# Devils in the details (mainly for PyTorch)

1. To reliably estimate BN statistics (running mean and average), you need at least 8 samples on each GPU
2. If you don't have enough samples on each GPU
  1. Synchronized BN layers
  2. Group normalization layers
3. Instance normalization layers are useful for some applications: such as style transfer, dense correspondence.

Next Class

More about  
Convolutional Neural Networks