

Lecture 7: Convolutional Networks

Reminder: A2

Due Monday, September 30, 11:59pm (Even if you enrolled late!)

Your submission must pass the validation script

Slight schedule change

Content originally planned for today got split into two lectures

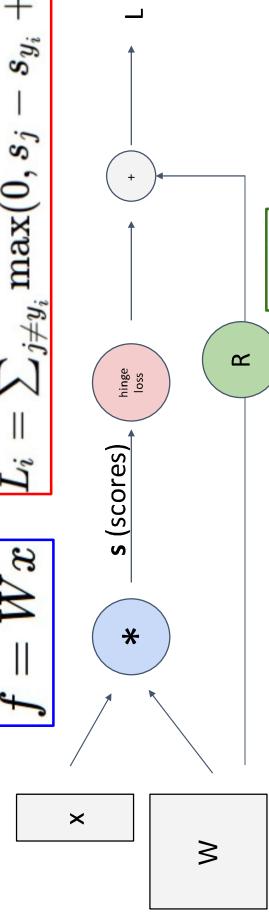
Pushes the schedule back a bit:

- A4 Due Date: Friday 11/1 -> Friday 11/8
- A5 Due Date: Friday 11/15 -> Friday 11/22
- A6 Due Date: Still Friday 12/6

Last Time: Backpropagation

Represent complex expressions as **computational graphs**

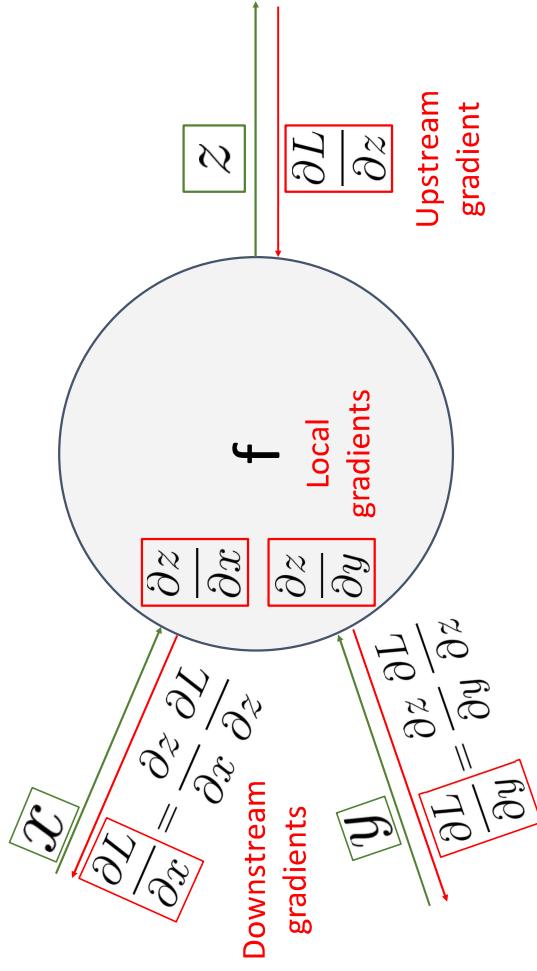
$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$



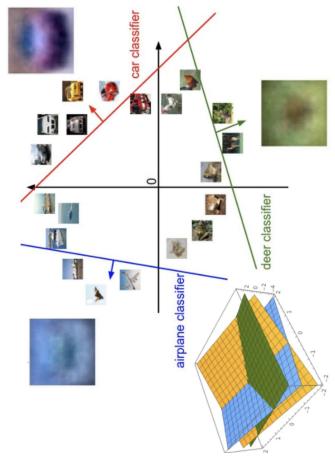
Forward pass computes outputs

Backward pass computes gradients

During the backward pass, each node in the graph receives **upstream gradients** and multiplies them by **local gradients** to compute **downstream gradients**

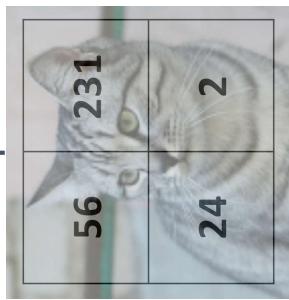


$$f(x, W) = Wx$$



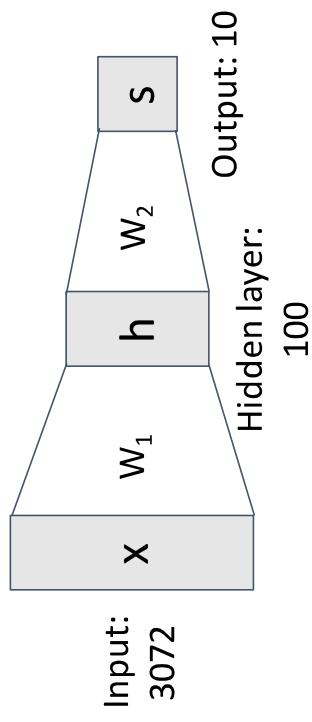
Stretch pixels into column

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



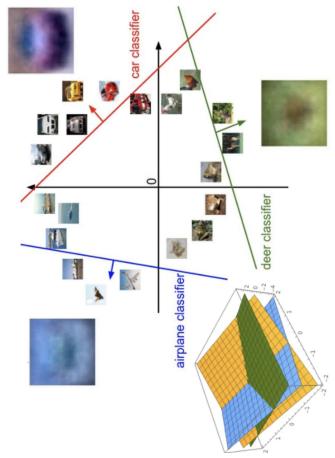
Input image
(2, 2)

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

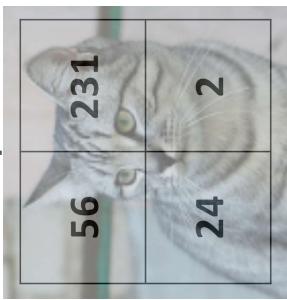


(4,)

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



Stretch pixels into column

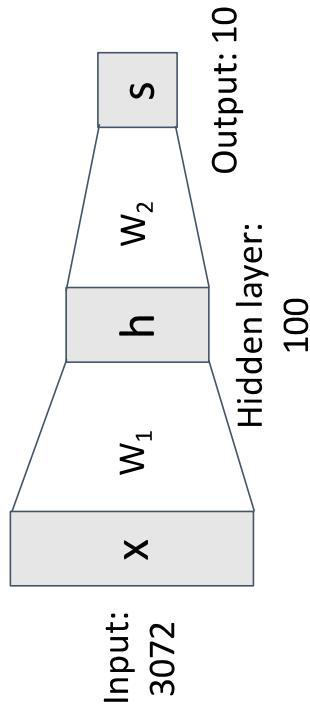


Problem: So far our classifiers don't respect the spatial structure of images!

Input image
(2, 2)

Solution: Define new computational nodes (4,) that operate on images!

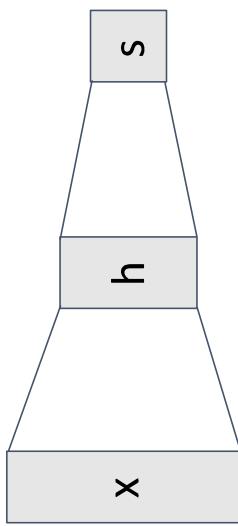
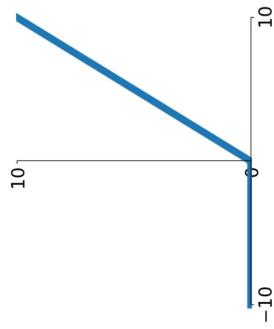
$$f = W_2 \max(0, W_1 x)$$



Components of a Full-Connected Network

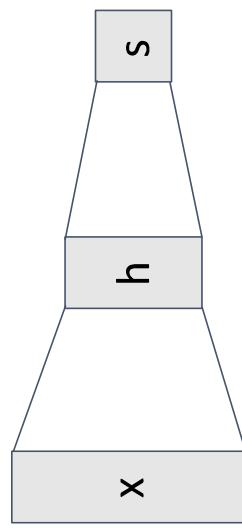
Fully-Connected Layers

Activation Function

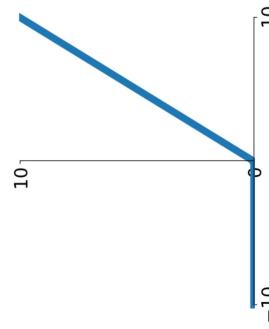


Components of a Convolutional Network

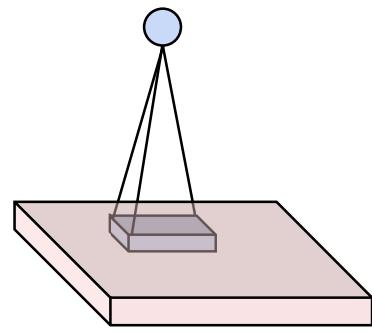
Fully-Connected Layers



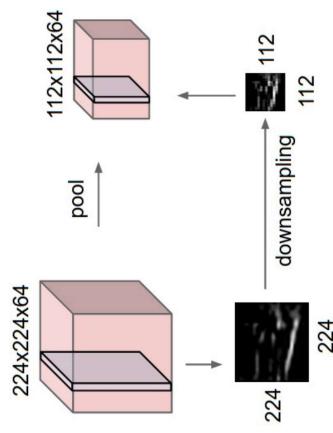
Activation Function



Convolution Layers



Pooling Layers

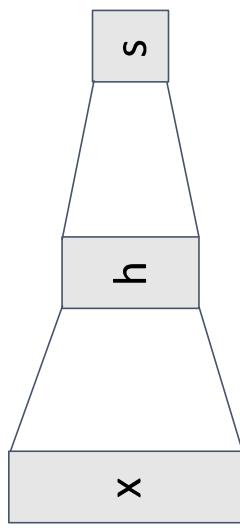


Normalization

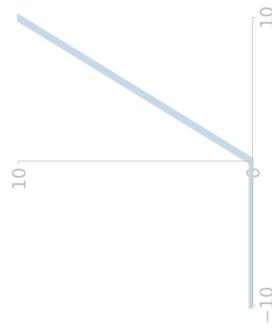
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Components of a Convolutional Network

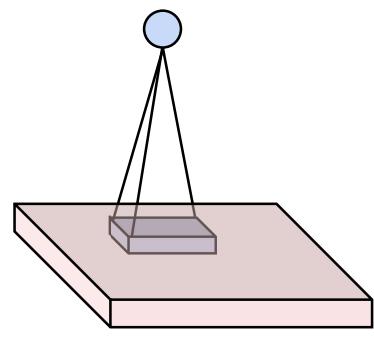
Fully-Connected Layers



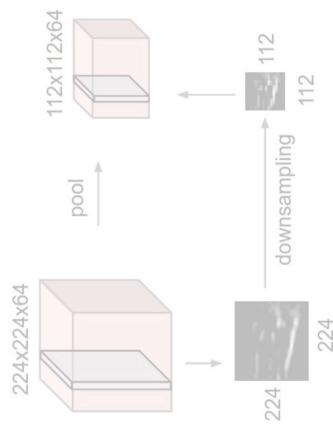
Activation Function



Convolution Layers



Pooling Layers

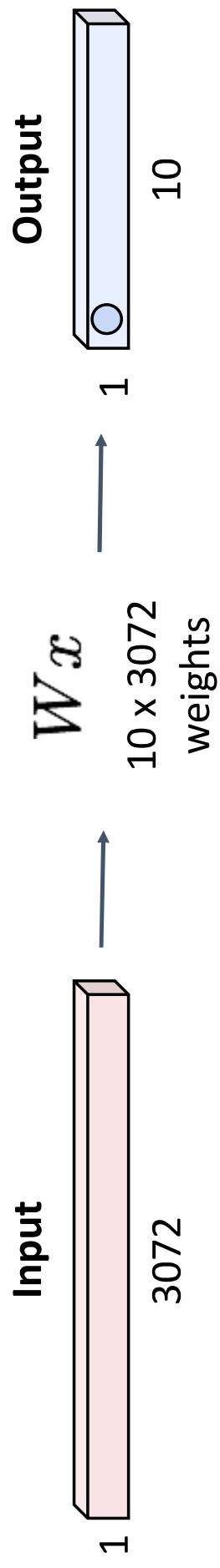


Normalization

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

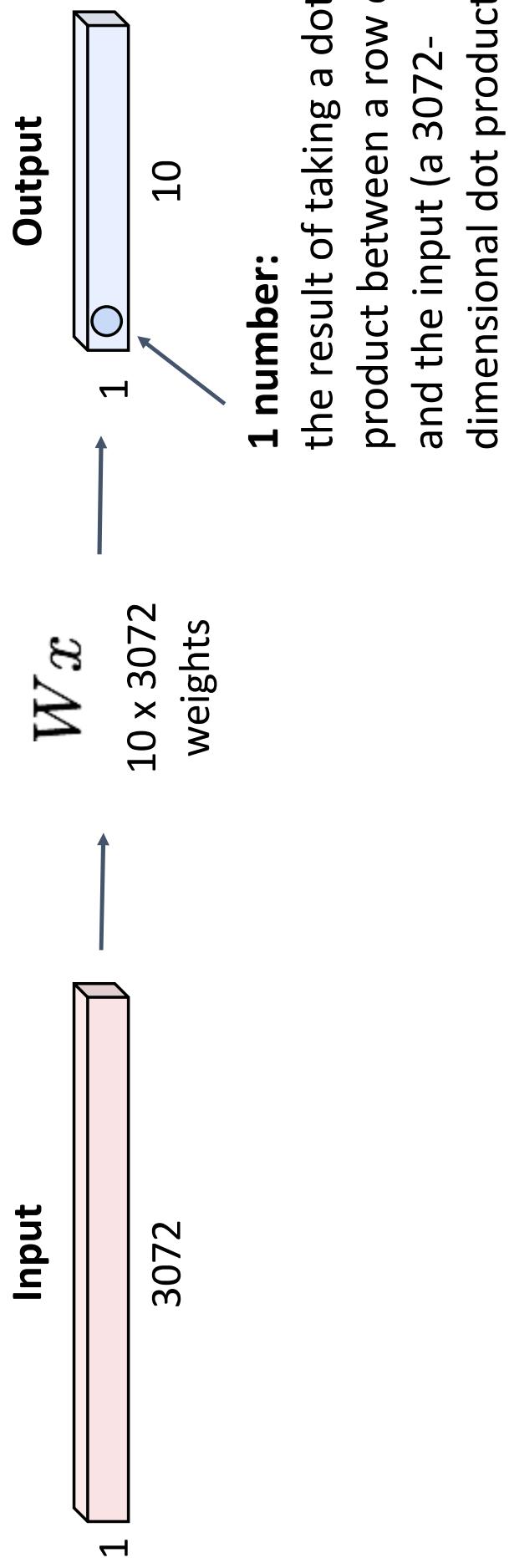
Fully-Connected Layer

32x32x3 image -> stretch to 3072×1



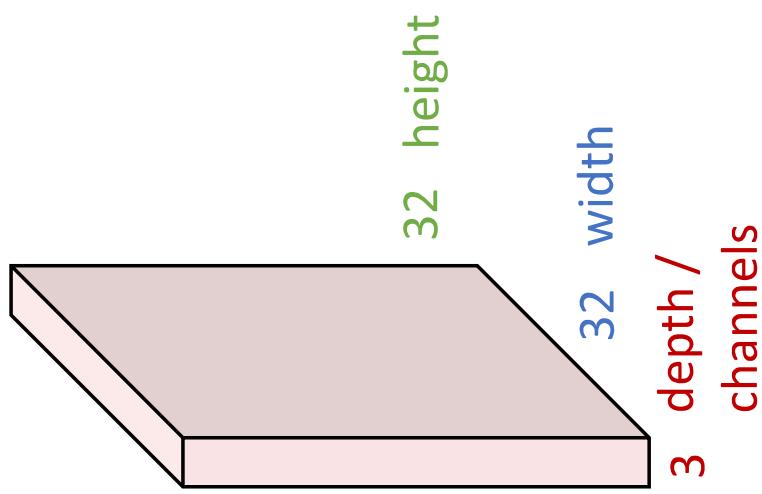
Fully-Connected Layer

32x32x3 image -> stretch to 3072×1



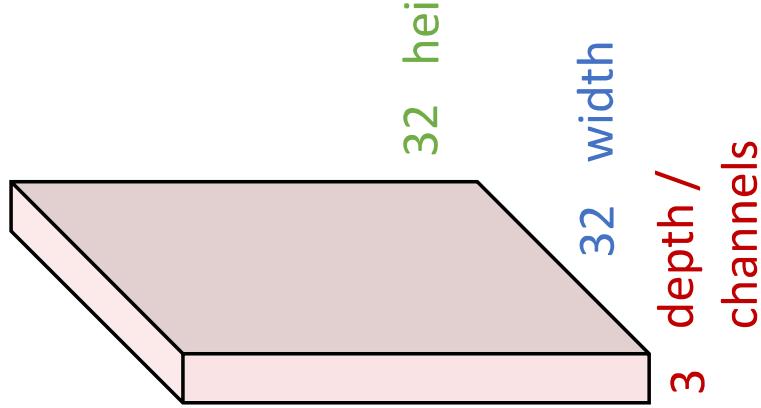
Convolution Layer

3x32x32 image: preserve spatial structure

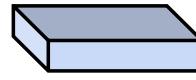


Convolution Layer

3x32x32 image



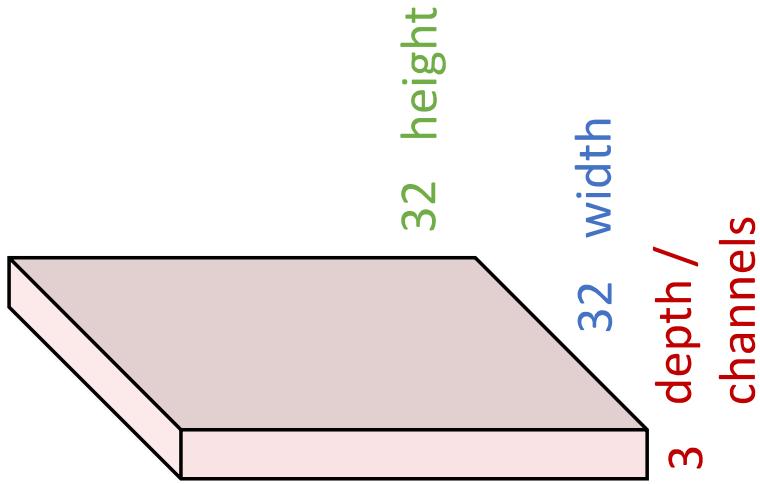
3x5x5 filter



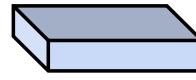
Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Convolution Layer

$3 \times 32 \times 32$ image
Filters always extend the full depth of the input volume



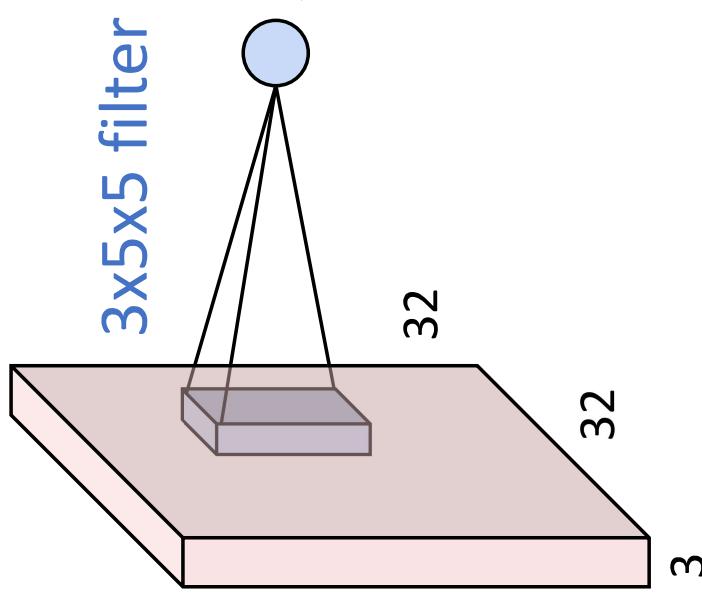
$3 \times 5 \times 5$ filter



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

Convolution Layer

3x32x32 image



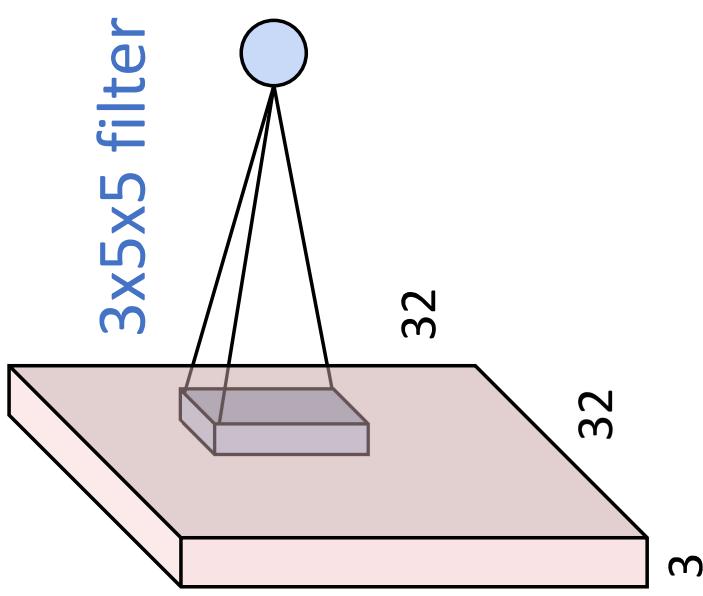
1 number:

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

$$w^T x + b$$

Convolution Layer

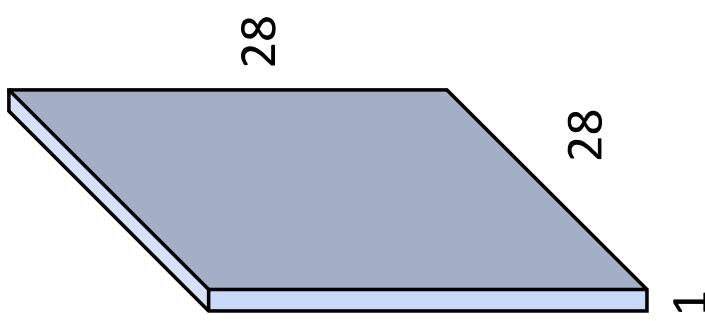
3x32x32 image



convolve (slide) over
all spatial locations

1x28x28

activation map

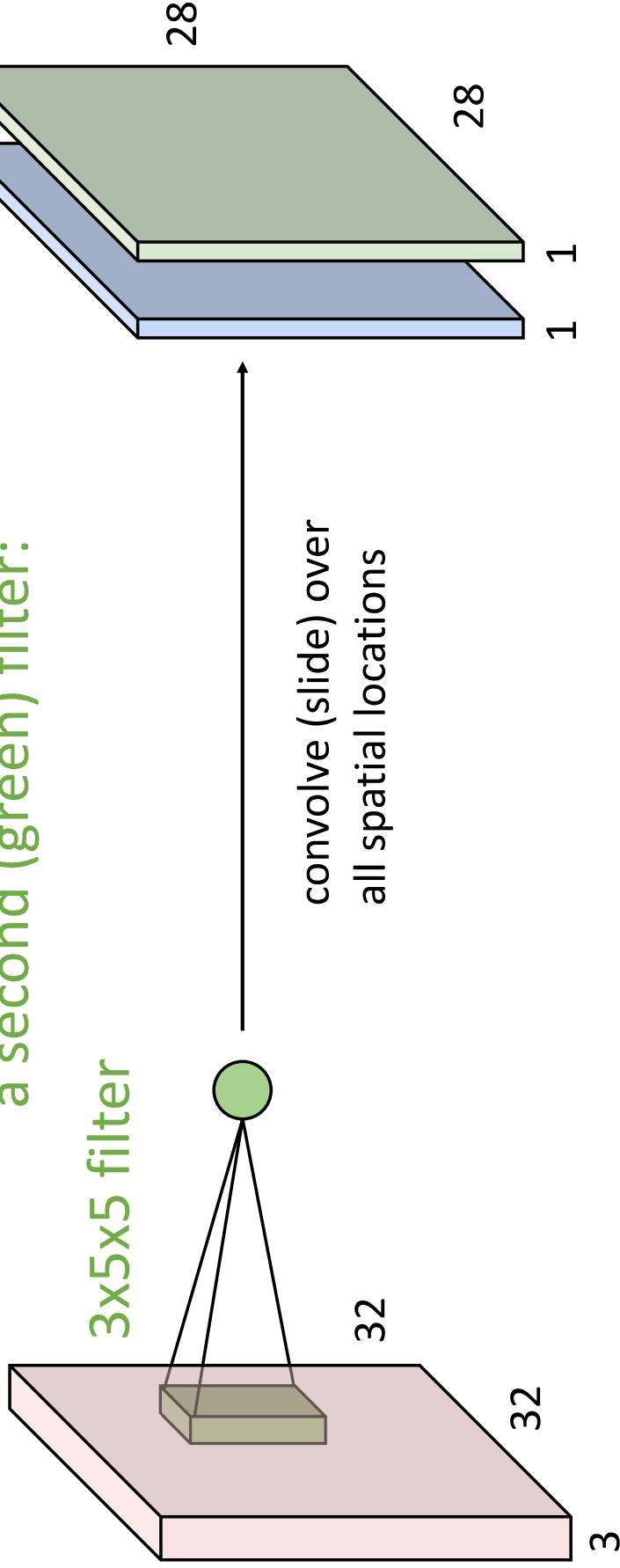


Convolution Layer

3x32x32 image

Consider repeating with
a second (green) filter:

3x5x5 filter



two 1x28x28 activation map

activation map

Convolution Layer

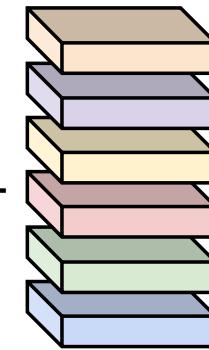
3x32x32 image

Consider 6 filters,
each $3 \times 5 \times 5$

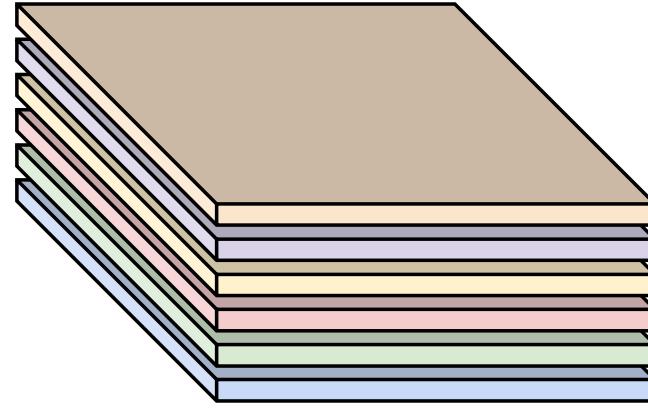


Convolution
Layer

6x3x5x5
filters



6 activation maps,
each $1 \times 28 \times 28$



Stack activations to get a
 $6 \times 28 \times 28$ output image!

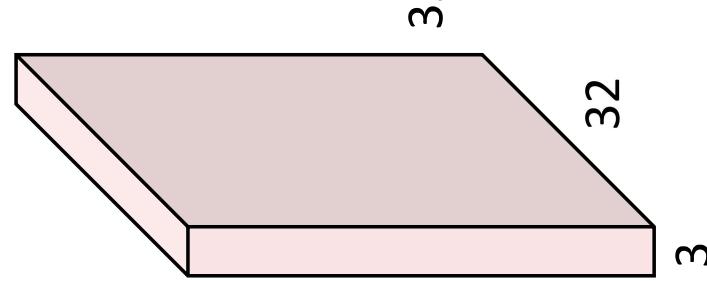
Convolution Layer

3x32x32 image

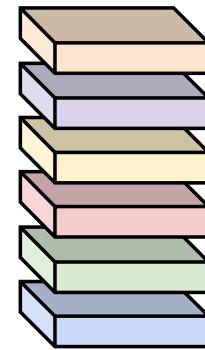
Also 6-dim bias vector:



Convolution
Layer



6x3x5x5
filters

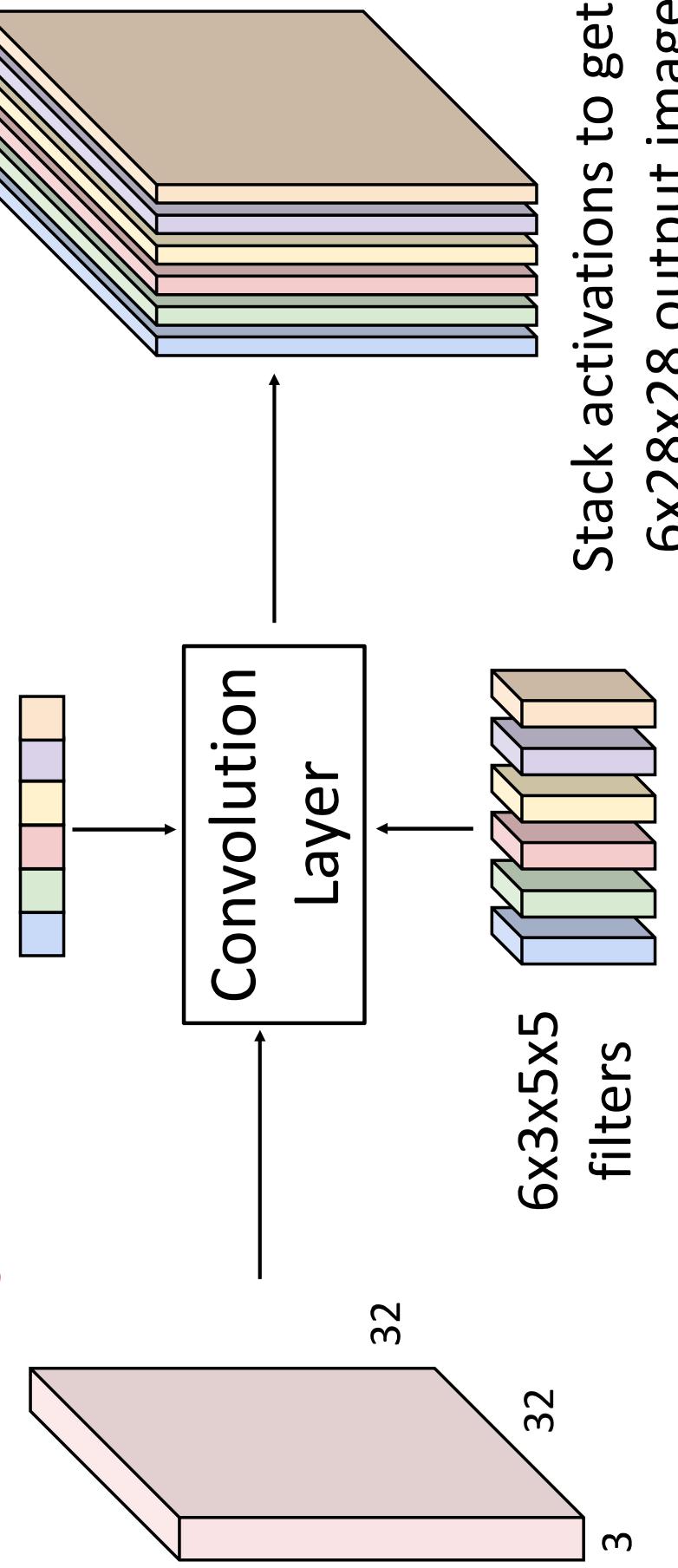


6 activation maps,
each 1x28x28

Convolution Layer

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

3x32x32 image
Also 6-dim bias vector:



Convolution Layer

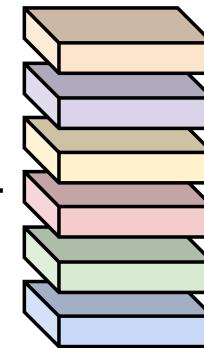
2x3x32x32

Batch of images

Also 6-dim bias vector:



Convolution
Layer



6x3x5x5
filters

32

32

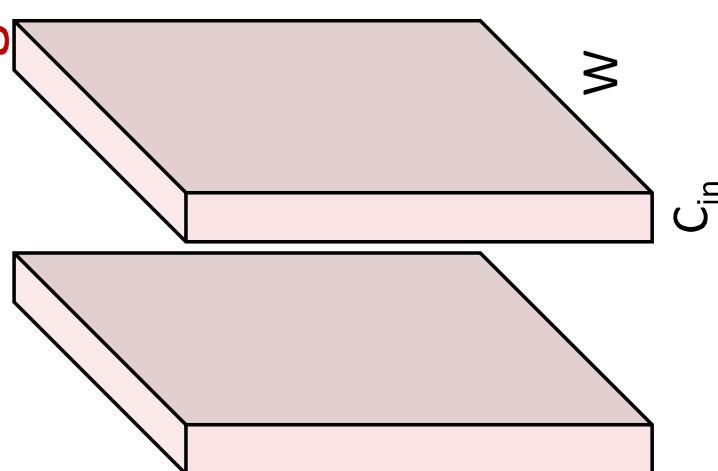
3

2x6x28x28
Batch of outputs

Convolution Layer

$N \times C_{in} \times H \times W$

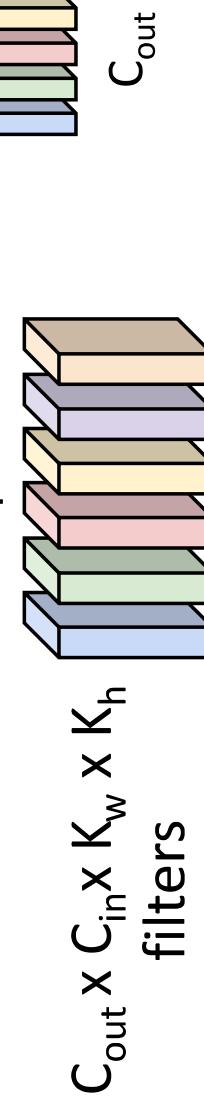
Batch of images



Also C_{out} -dim bias vector:

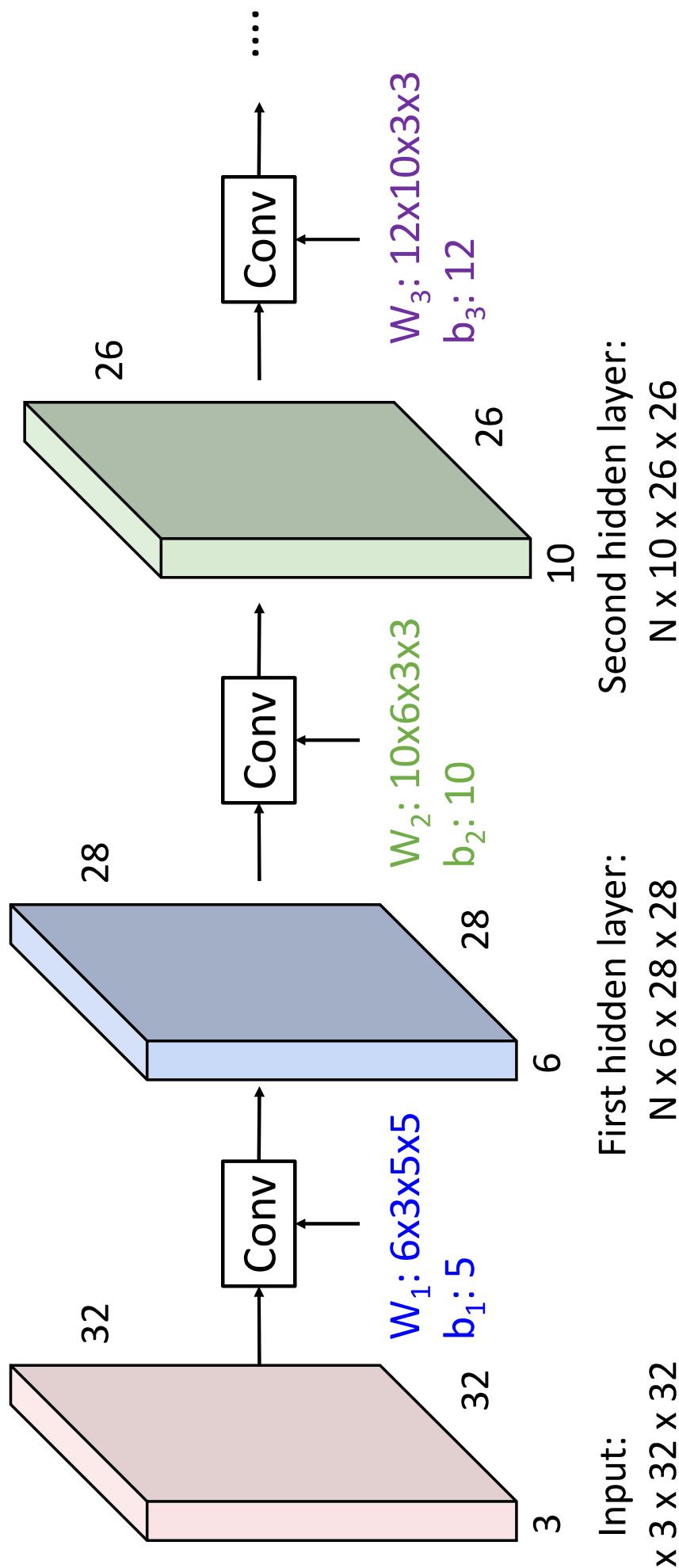


Convolution
Layer



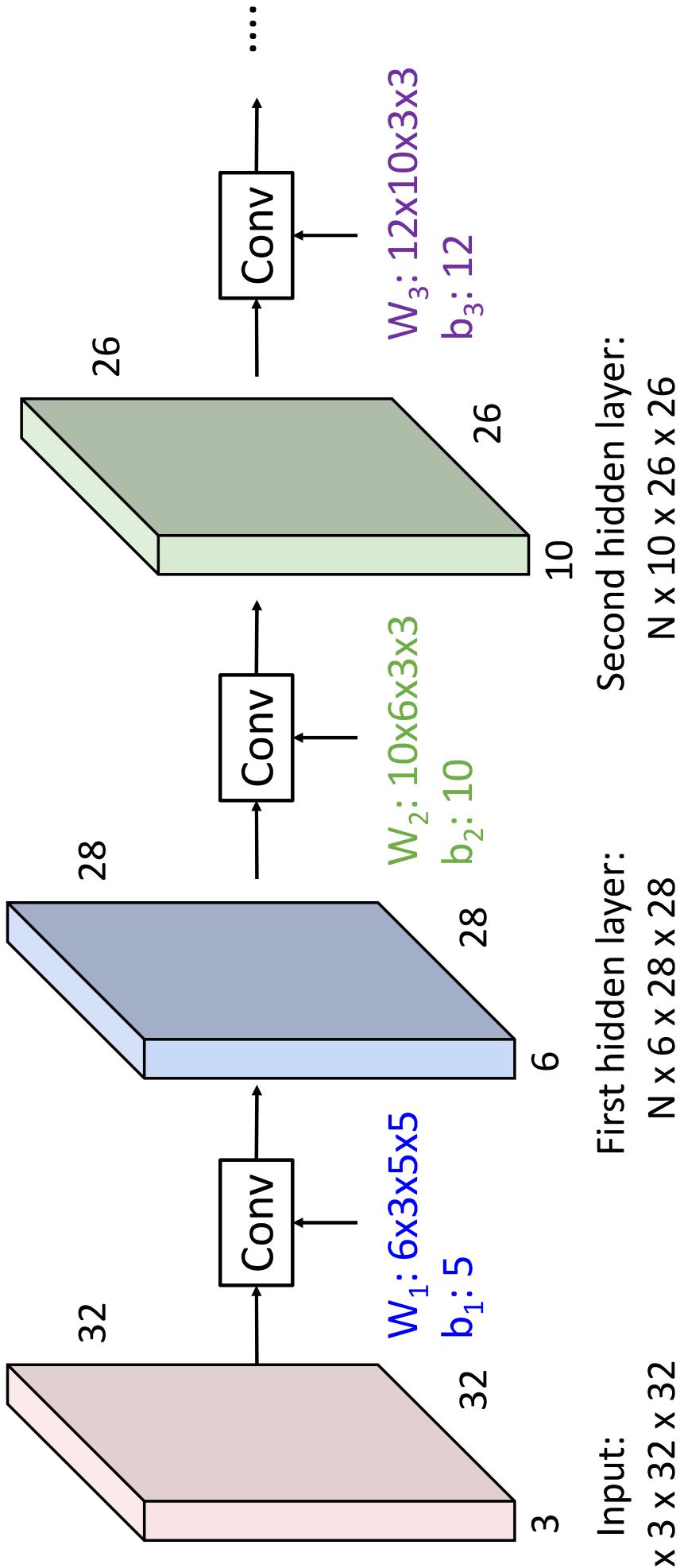
$N \times C_{out} \times H' \times W'$
Batch of outputs

Stacking Convolutions



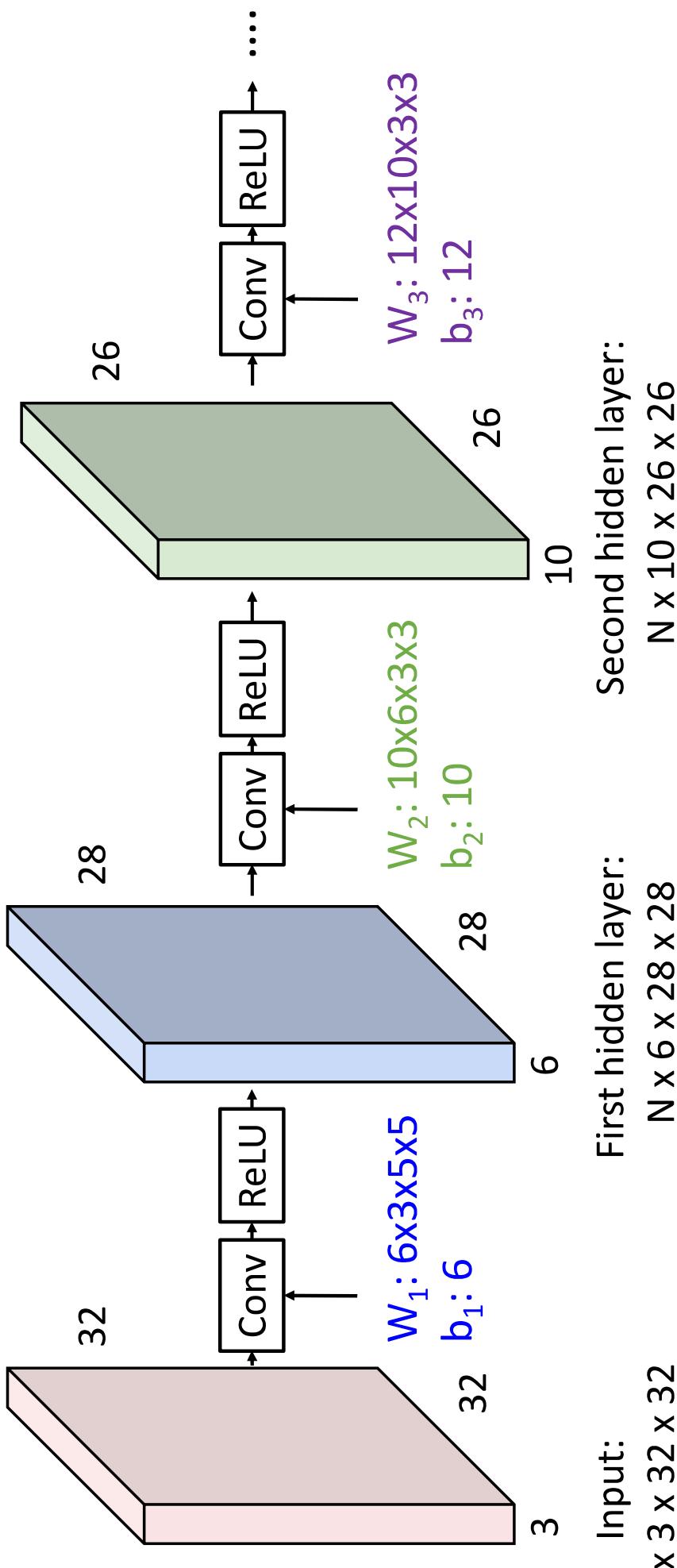
Stacking Convolutions

Q: What happens if we stack
two convolution layers?

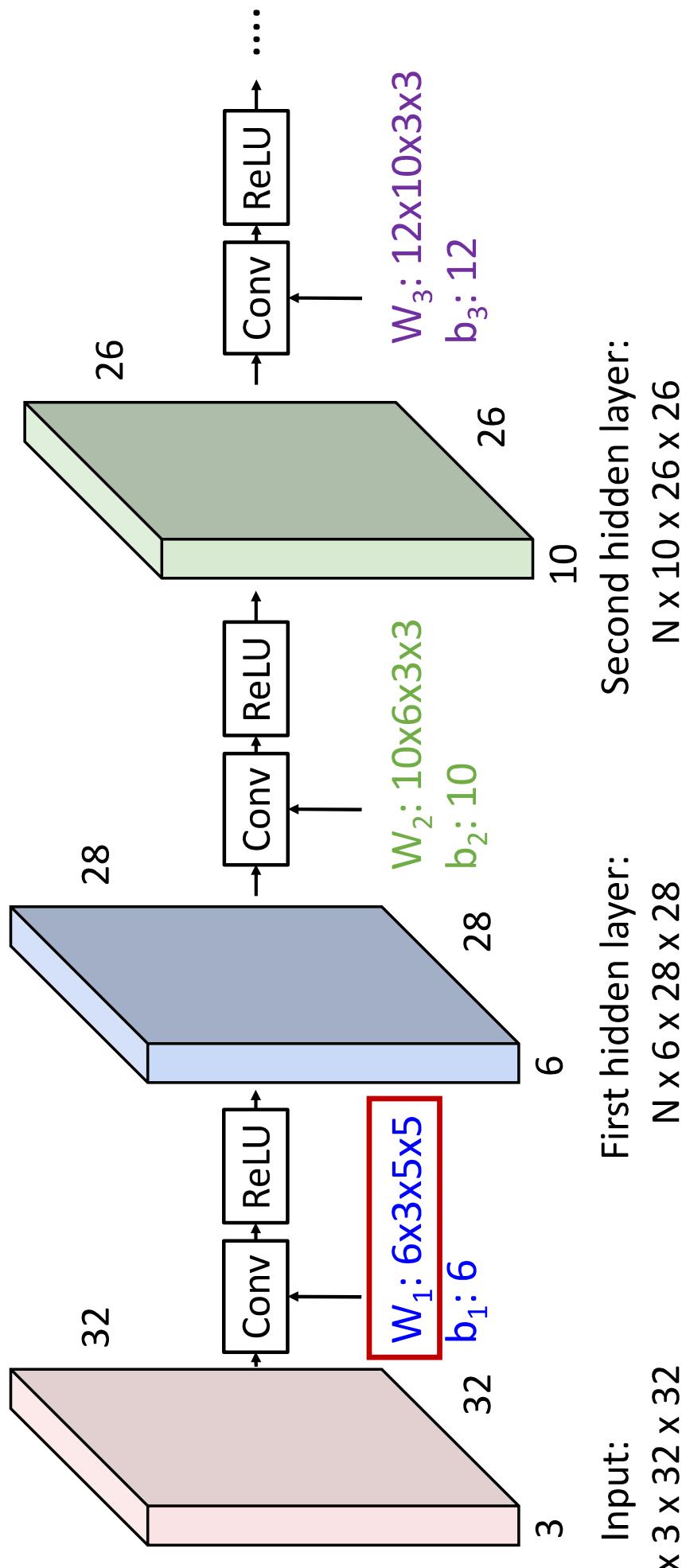


Stacking Convolutions

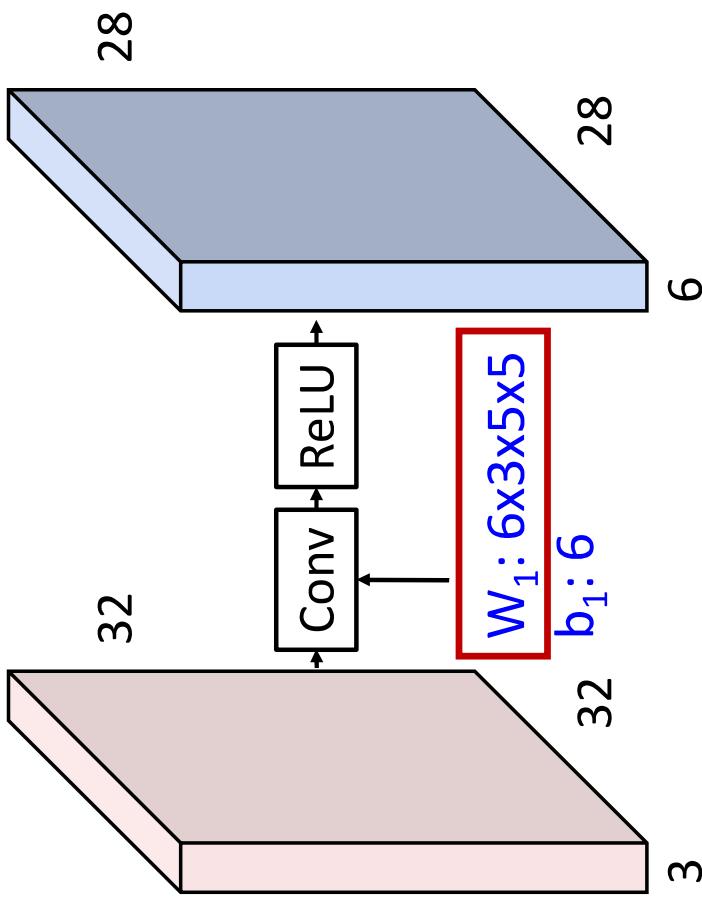
Q: What happens if we stack (Recall $y=W_2W_1x$ is
two convolution layers?
A: We get another convolution!



What do convolutional filters learn?



What do convolutional filters learn?

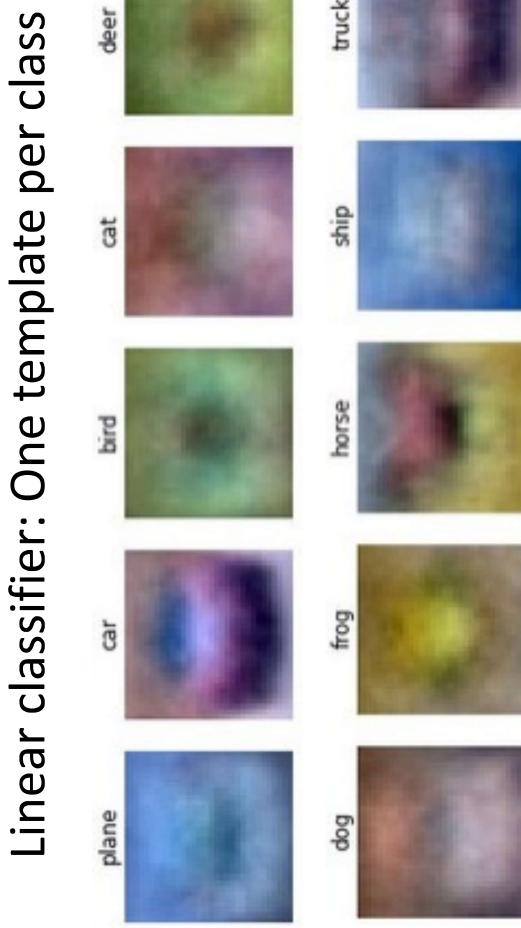


Input:

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

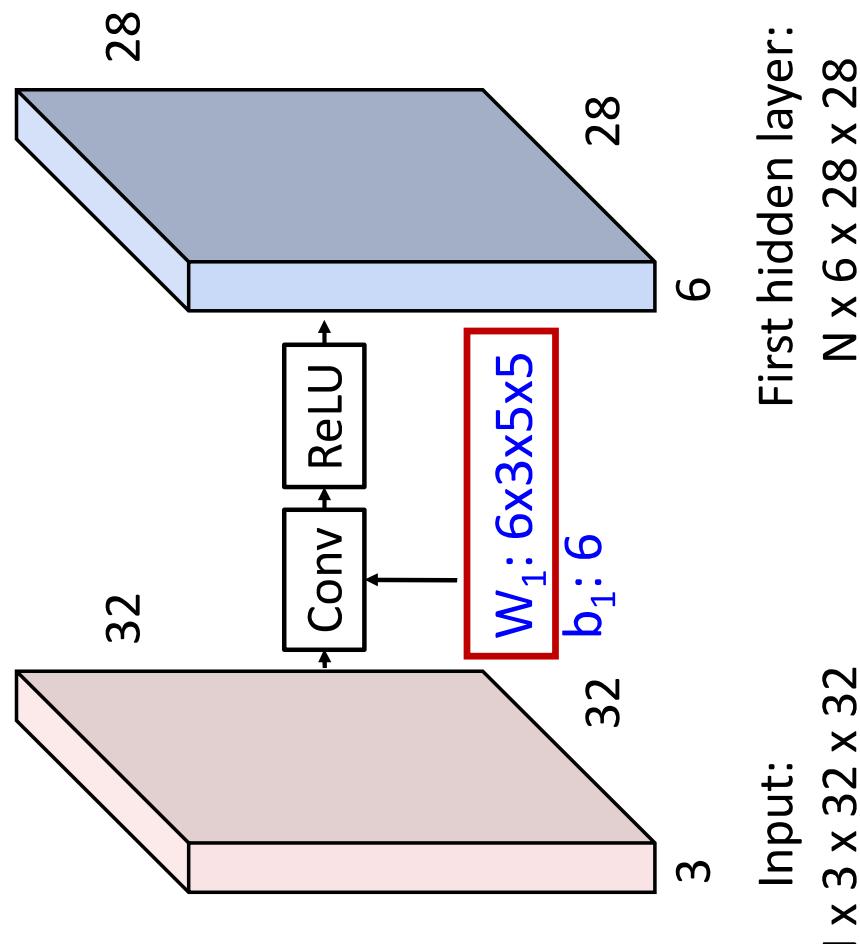
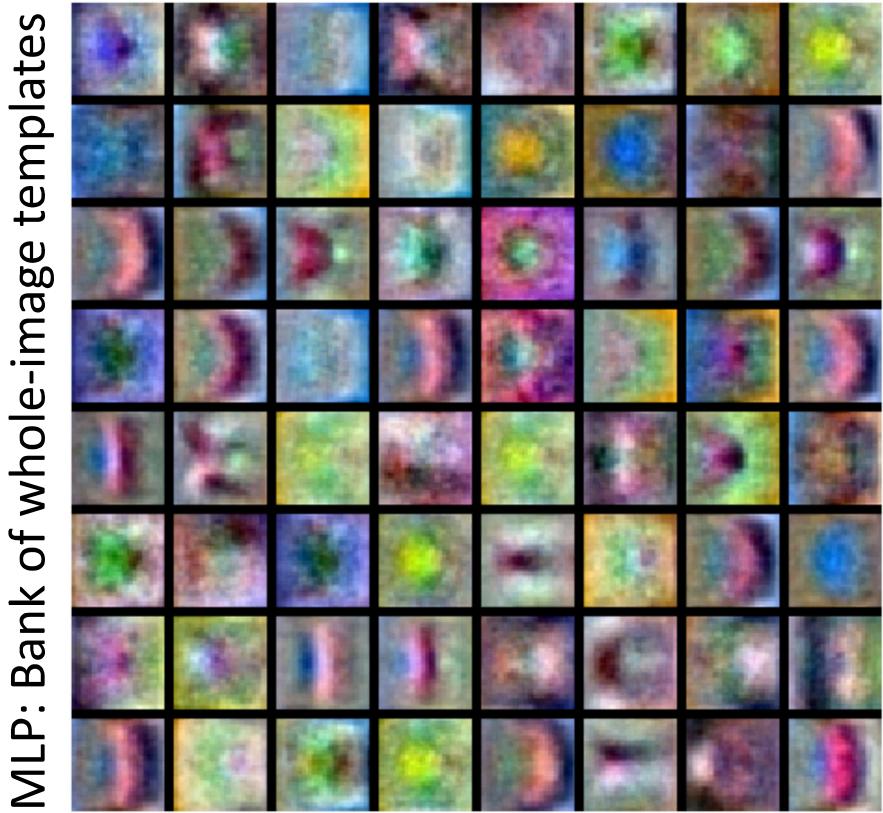
First hidden layer:

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



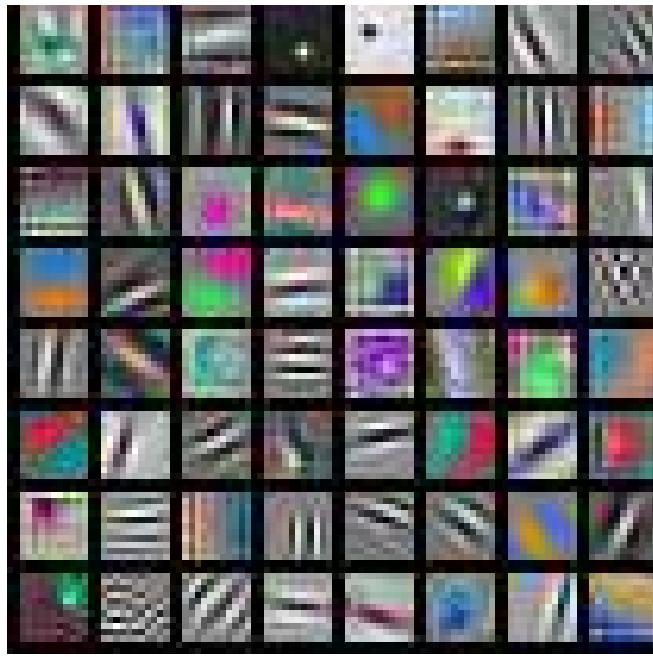
Linear classifier: One template per class

What do convolutional filters learn?

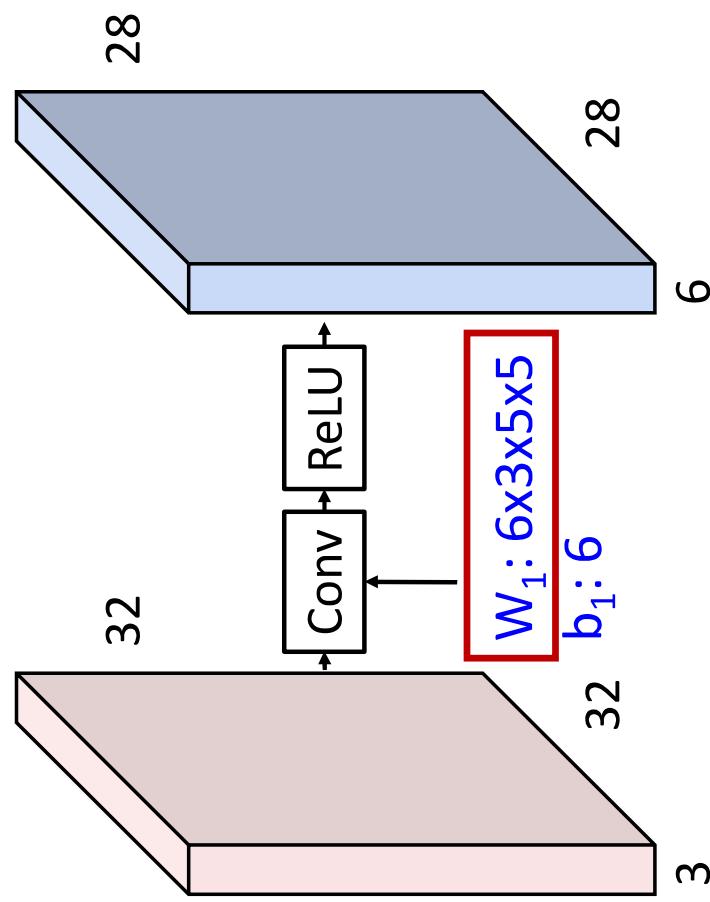


What do convolutional filters learn?

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



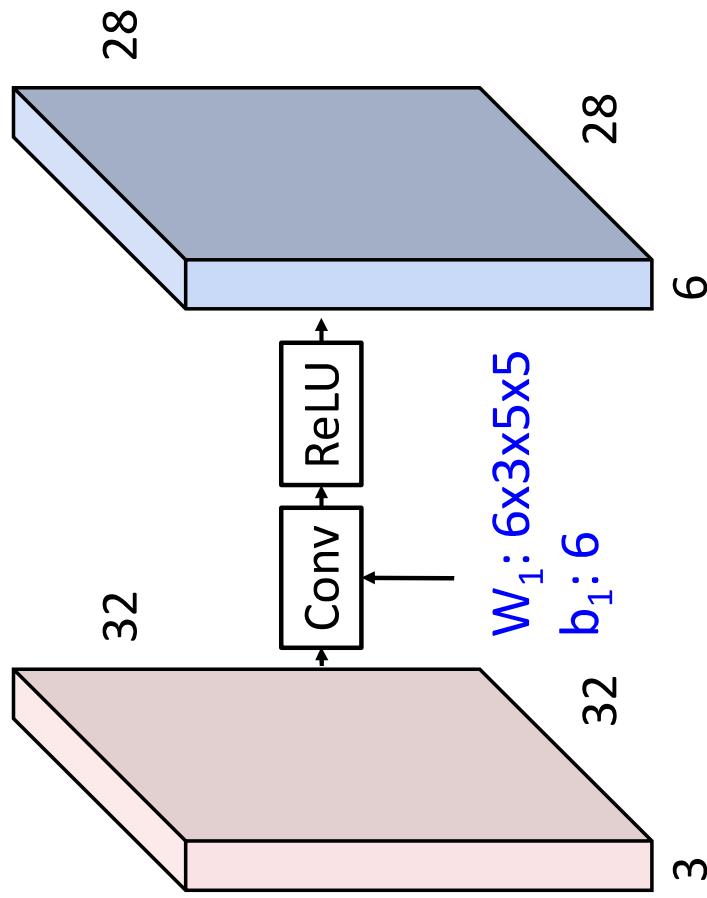
Input:
 $N \times 3 \times 32 \times 32$



First hidden layer:
 $N \times 6 \times 28 \times 28$

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

A closer look at spatial dimensions

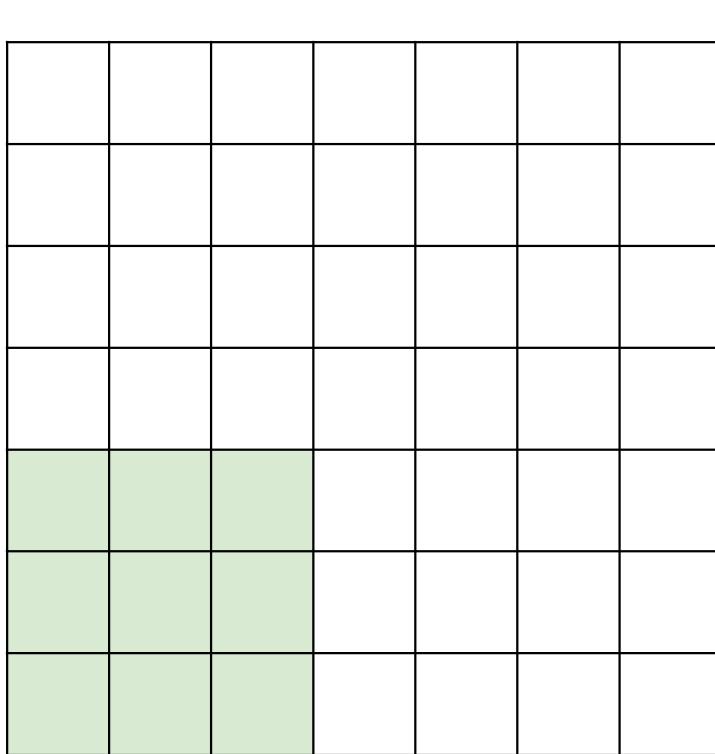


Input: $N \times 3 \times 32 \times 32$

First hidden layer:
 $N \times 6 \times 28 \times 28$

A closer look at spatial dimensions

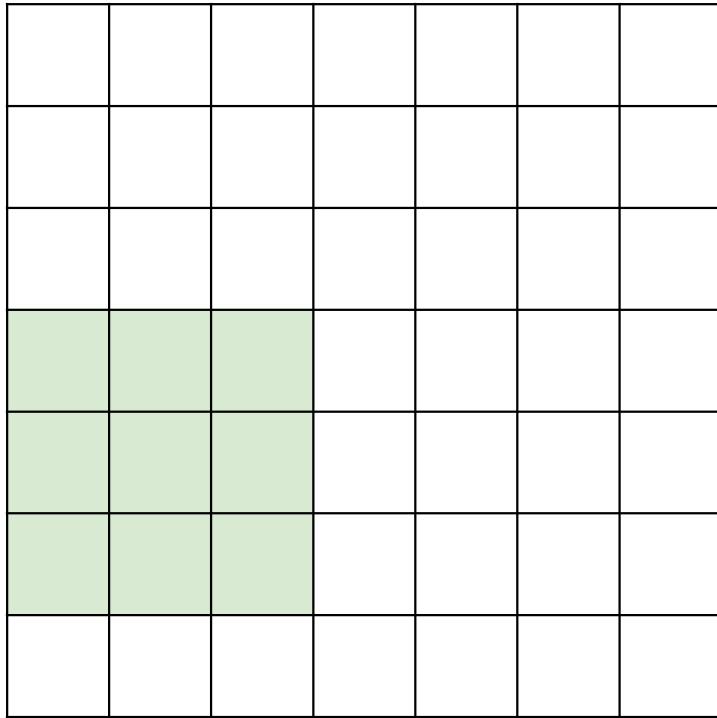
Input: 7x7
Filter: 3x3



7

A closer look at spatial dimensions

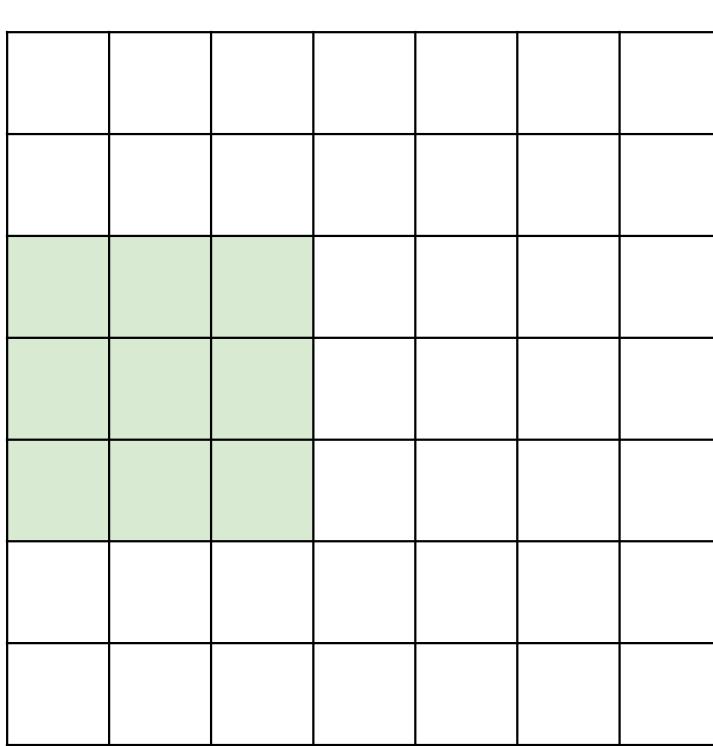
Input: 7x7
Filter: 3x3



7

A closer look at spatial dimensions

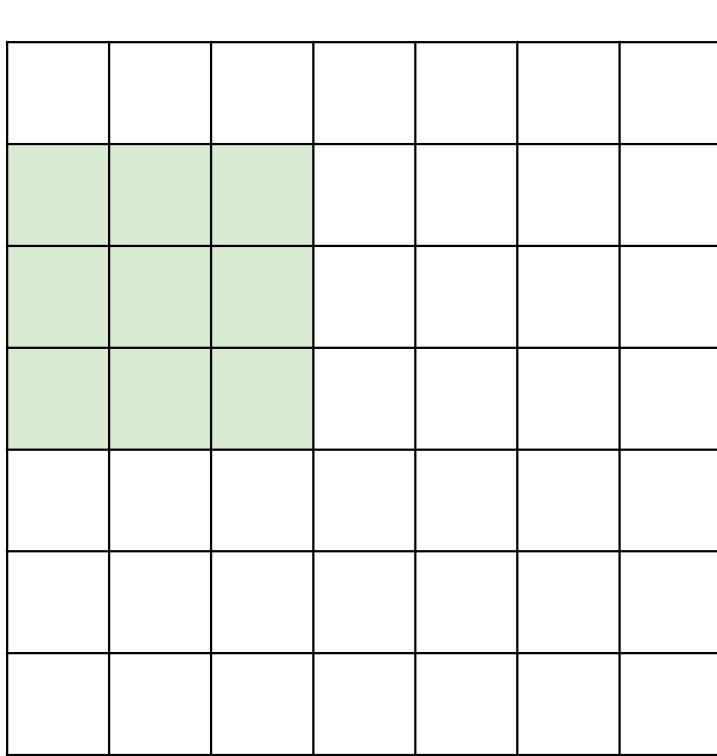
Input: 7x7
Filter: 3x3



7

A closer look at spatial dimensions

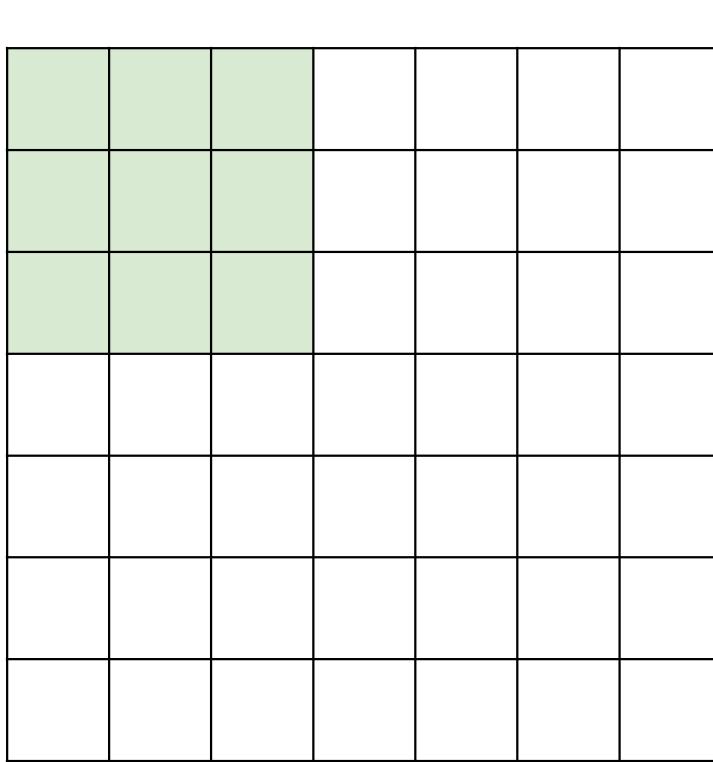
Input: 7x7
Filter: 3x3



7

A closer look at spatial dimensions

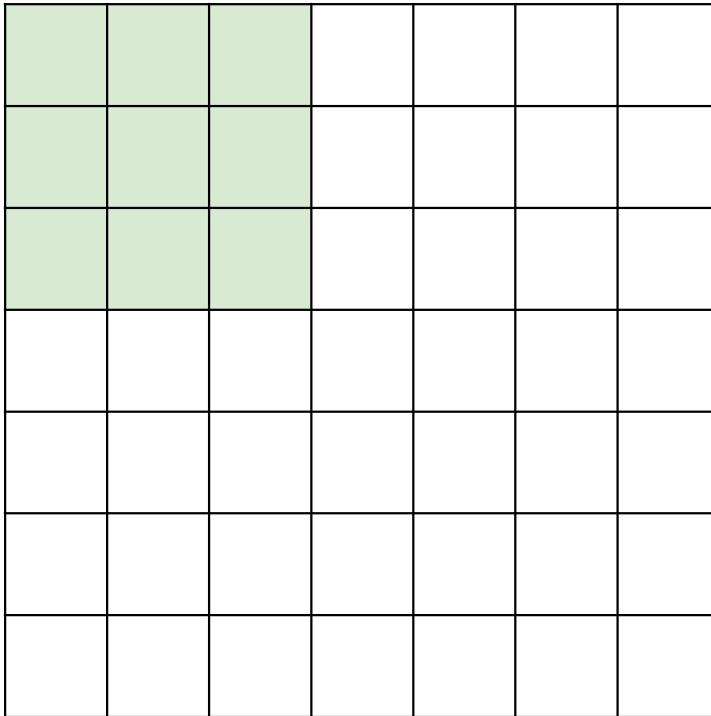
Input: 7x7
Filter: 3x3
Output: 5x5



7

A closer look at spatial dimensions

Input: 7x7
Filter: 3x3
Output: 5x5



In general:
Input: W
Filter: K
Output: $W - K + 1$

Problem: Feature maps “shrink” with each layer!

7

A closer look at spatial dimensions

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								

Input: 7x7

Filter: 3x3

Output: 5x5

In general:

Input: W

Filter: K

Padding: P

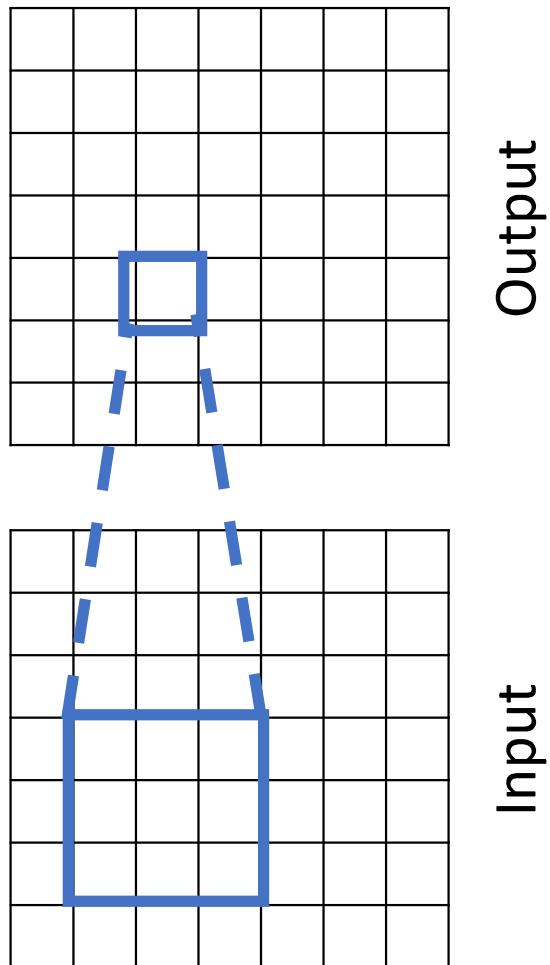
Output: $W - K + 1 + 2P$

Very common:

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

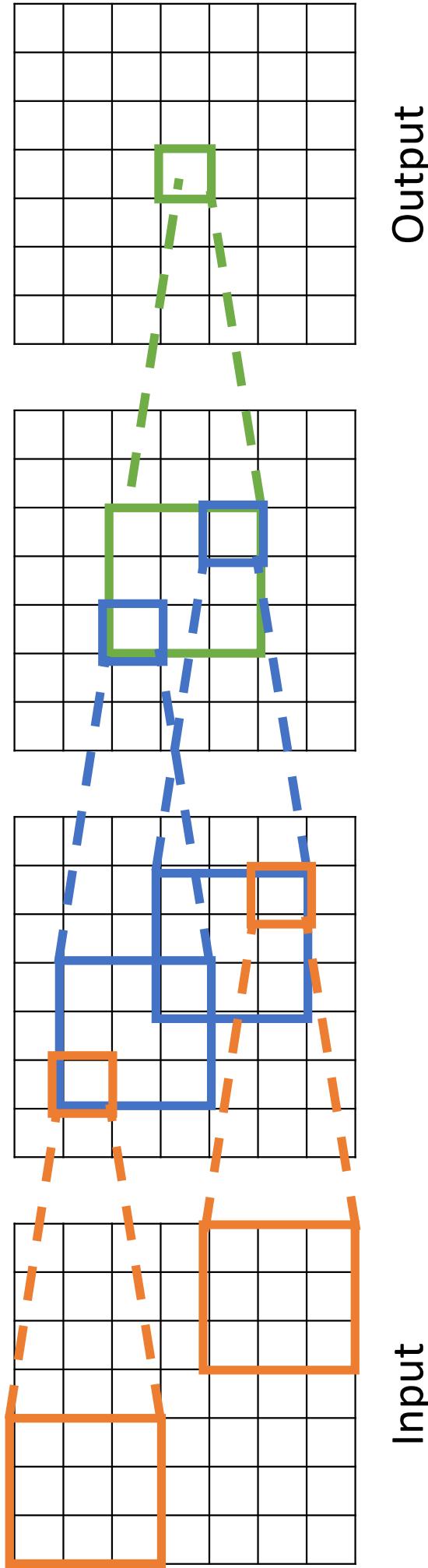
Receptive Fields

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



Receptive Fields

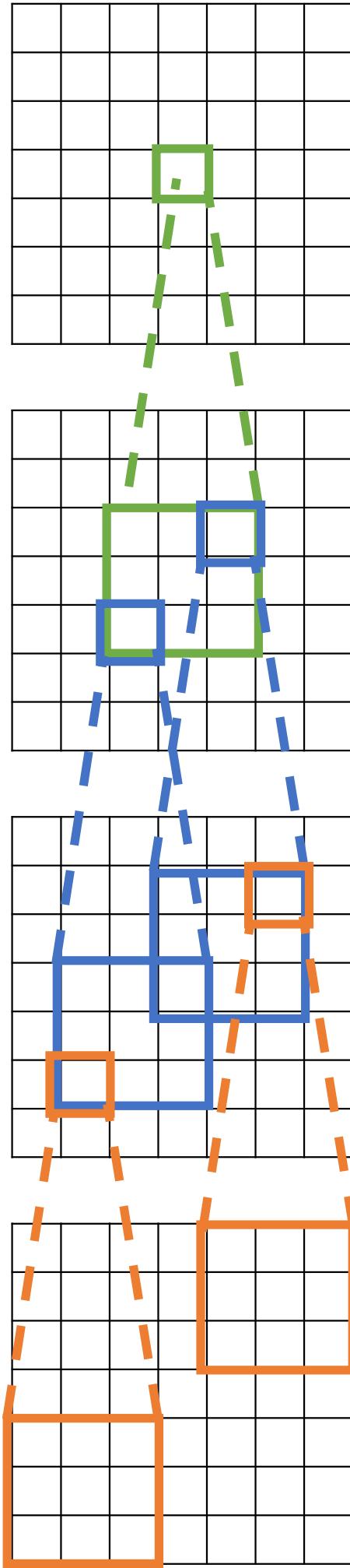
Each successive convolution adds $K - 1$ to the receptive field size
With L layers the receptive field size is $1 + L * (K - 1)$



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 $1 + L * (K - 1)$

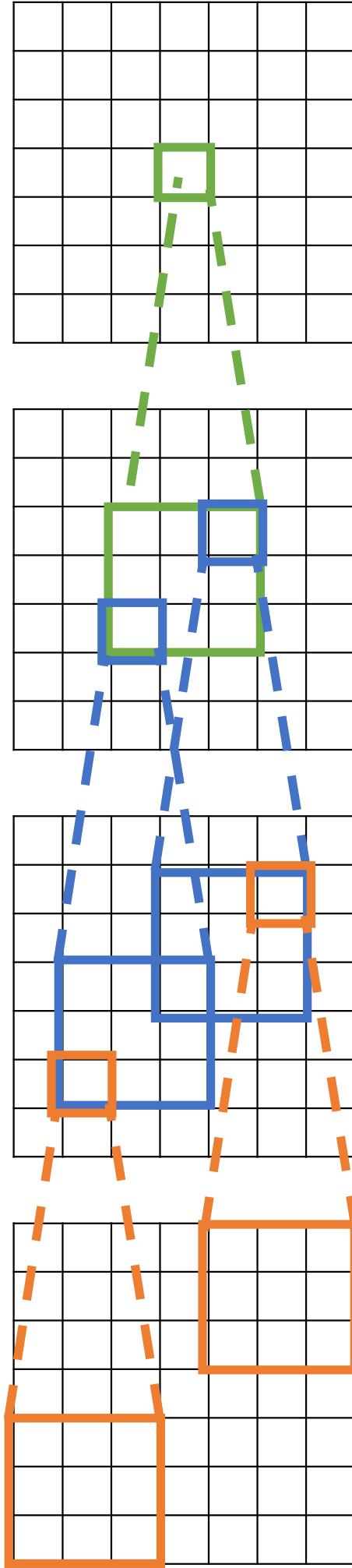


Input

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

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

Output

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

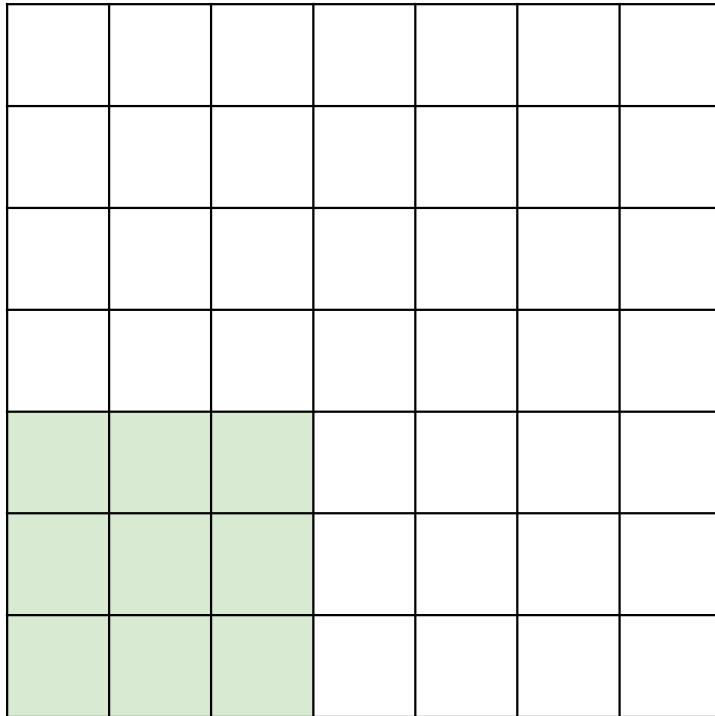
Solution: Downsample inside the network

Strided Convolution

Input: 7x7

Filter: 3x3

Stride: 2

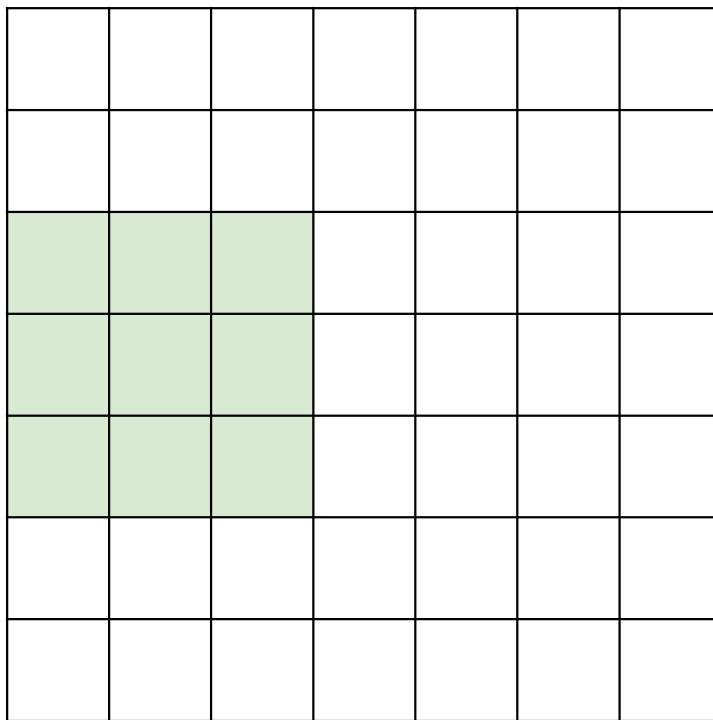


Strided Convolution

Input: 7x7

Filter: 3x3

Stride: 2

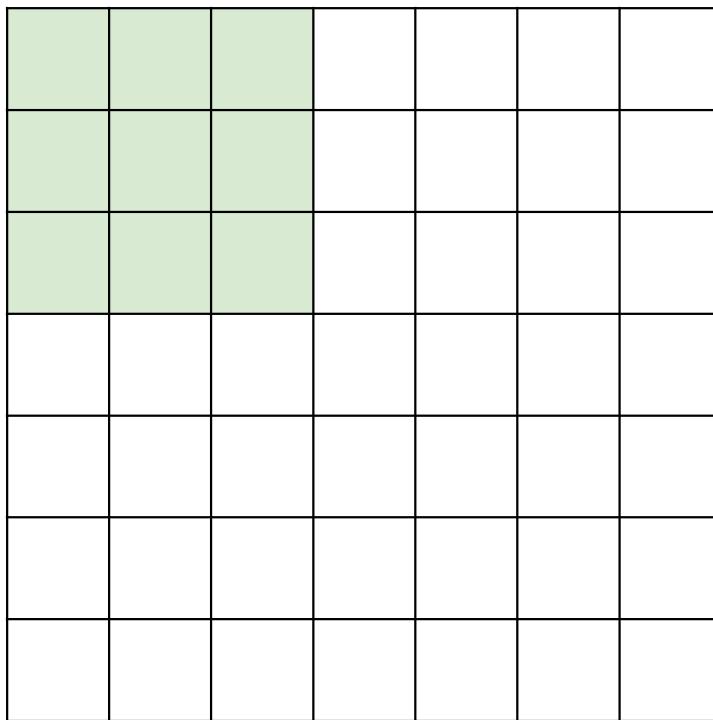


Strided Convolution

Input: 7×7

Filter: 3×3

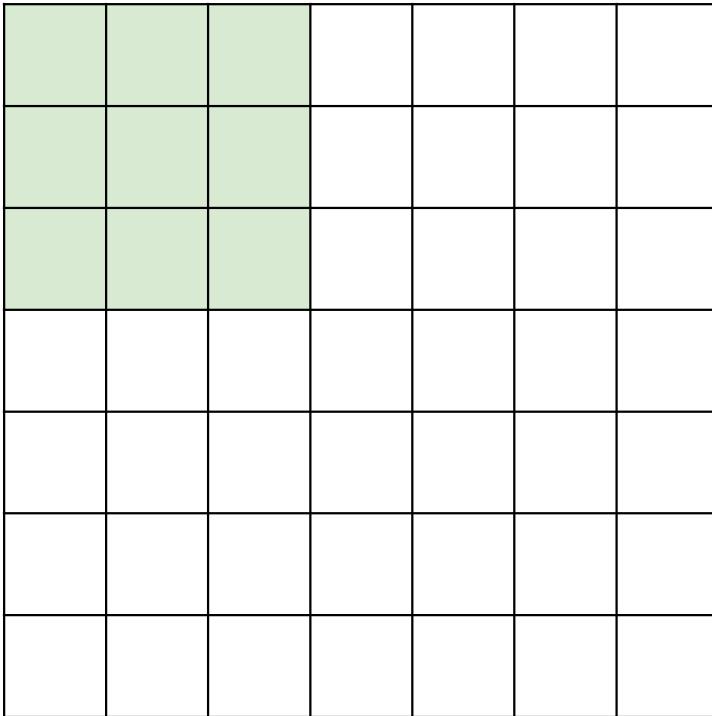
Stride: 2



Output: 3×3

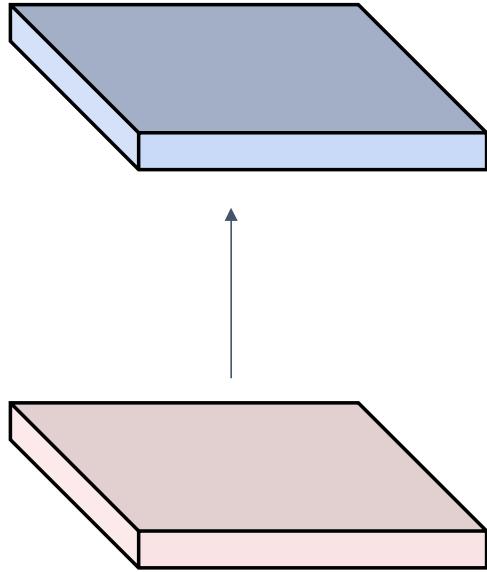
Strided Convolution

Input: 7x7
Filter: 3x3
Stride: 2



In general:
Input: W
Filter: K
Padding: P
Stride: S
Output: $(W - K + 2P) / S + 1$

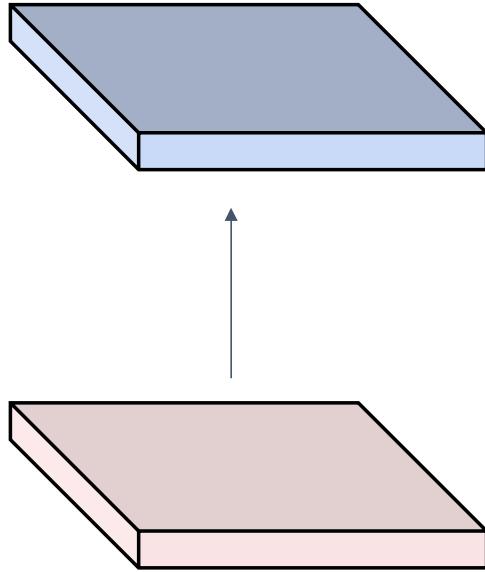
Convolution Example



Input volume: $3 \times 32 \times 32$
10 5x5 filters with stride 1, pad 2

Output volume size: ?

Convolution Example

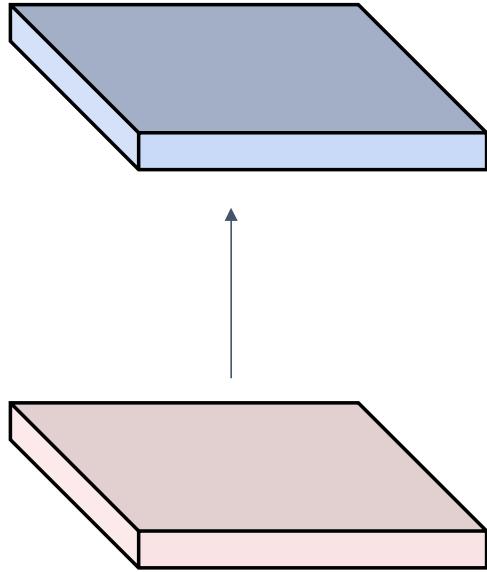


Input volume: $3 \times 32 \times 32$
10 5x5 filters with stride **1**, pad **2**

Output volume size:

$(32+2*2-5)/1+1 = 32$ spatially, so
10 $\times 32 \times 32$

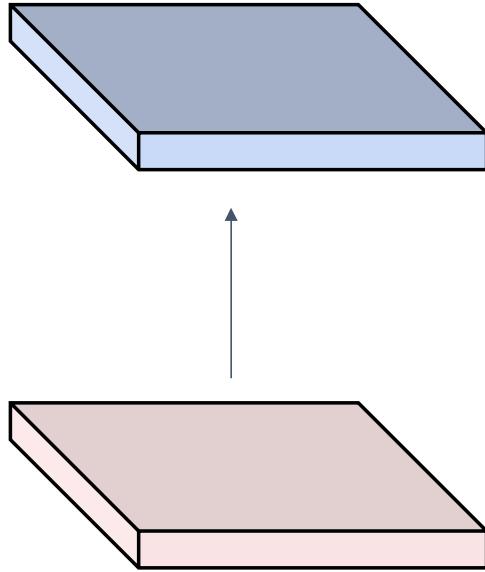
Convolution Example



Input volume: $3 \times 32 \times 32$
10 5x5 filters with stride 1, pad 2

Output volume size: $10 \times 32 \times 32$
Number of learnable parameters: ?

Convolution Example



Input volume: **3** × 32 × 32

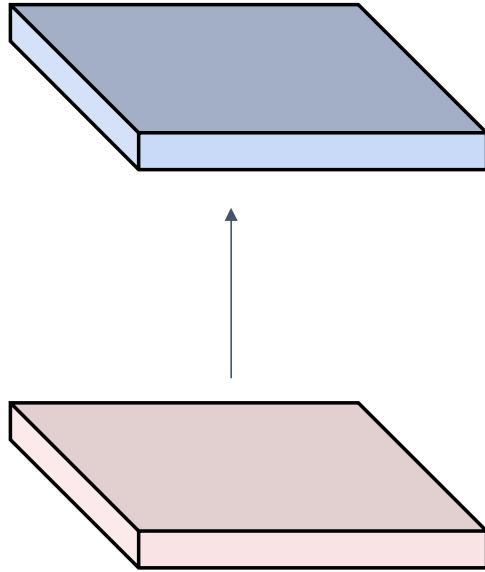
10 5x5 filters with stride 1, pad 2

Output volume size: 10 × 32 × 32

Number of learnable parameters: **760**

Parameters per filter: **3 * 5 * 5 + 1** (for bias) = **76**
10 filters, so total is **10 * 76 = 760**

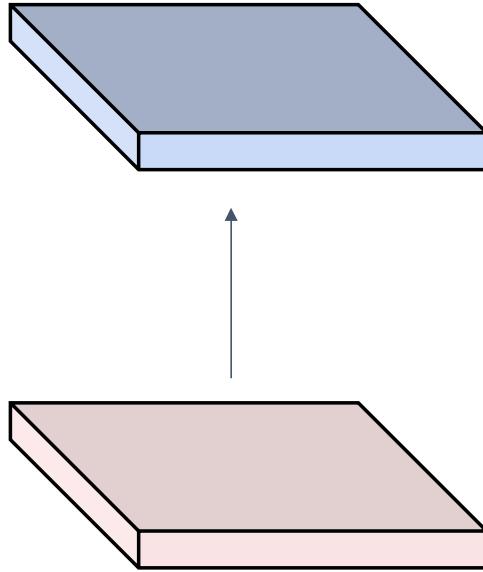
Convolution Example



Input volume: $3 \times 32 \times 32$
10 5x5 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

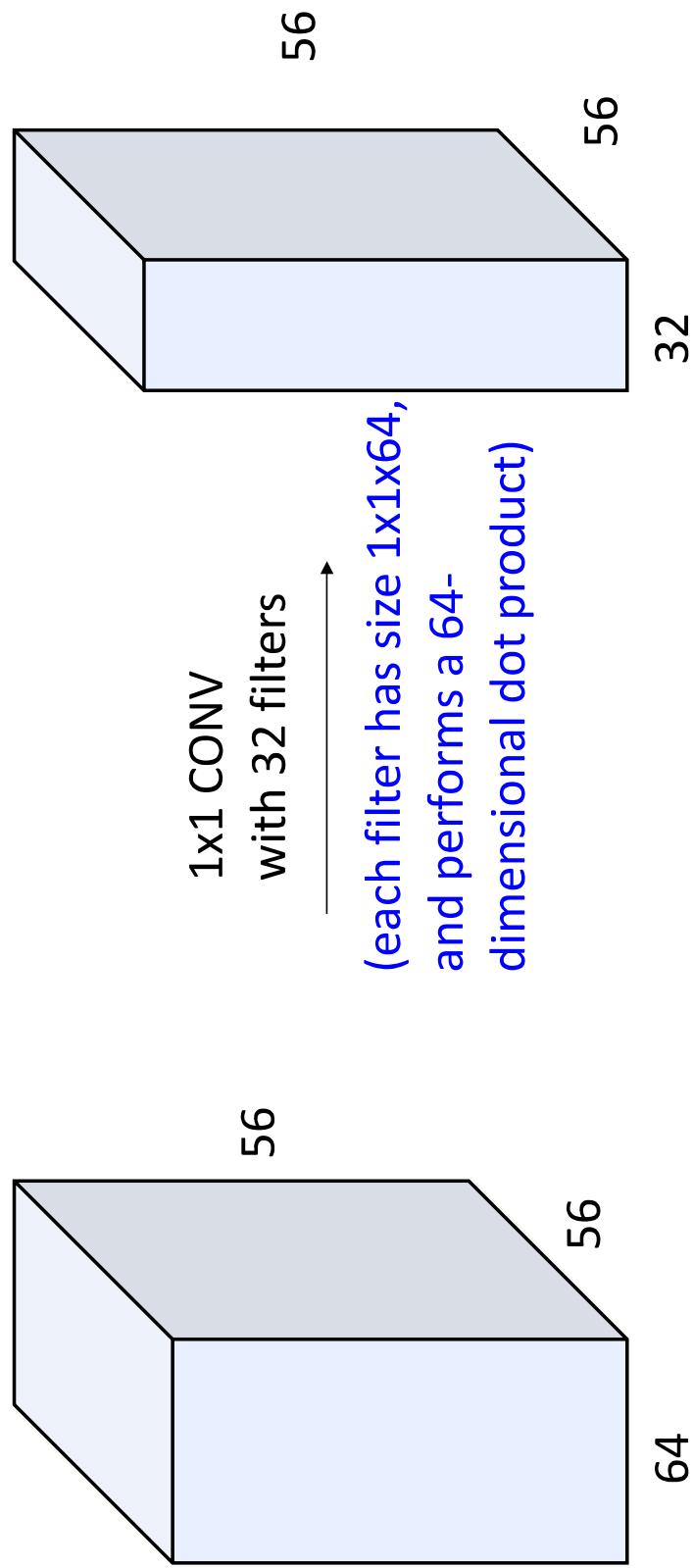
Output volume size: **10 x 32 x 32**

Number of learnable parameters: 760

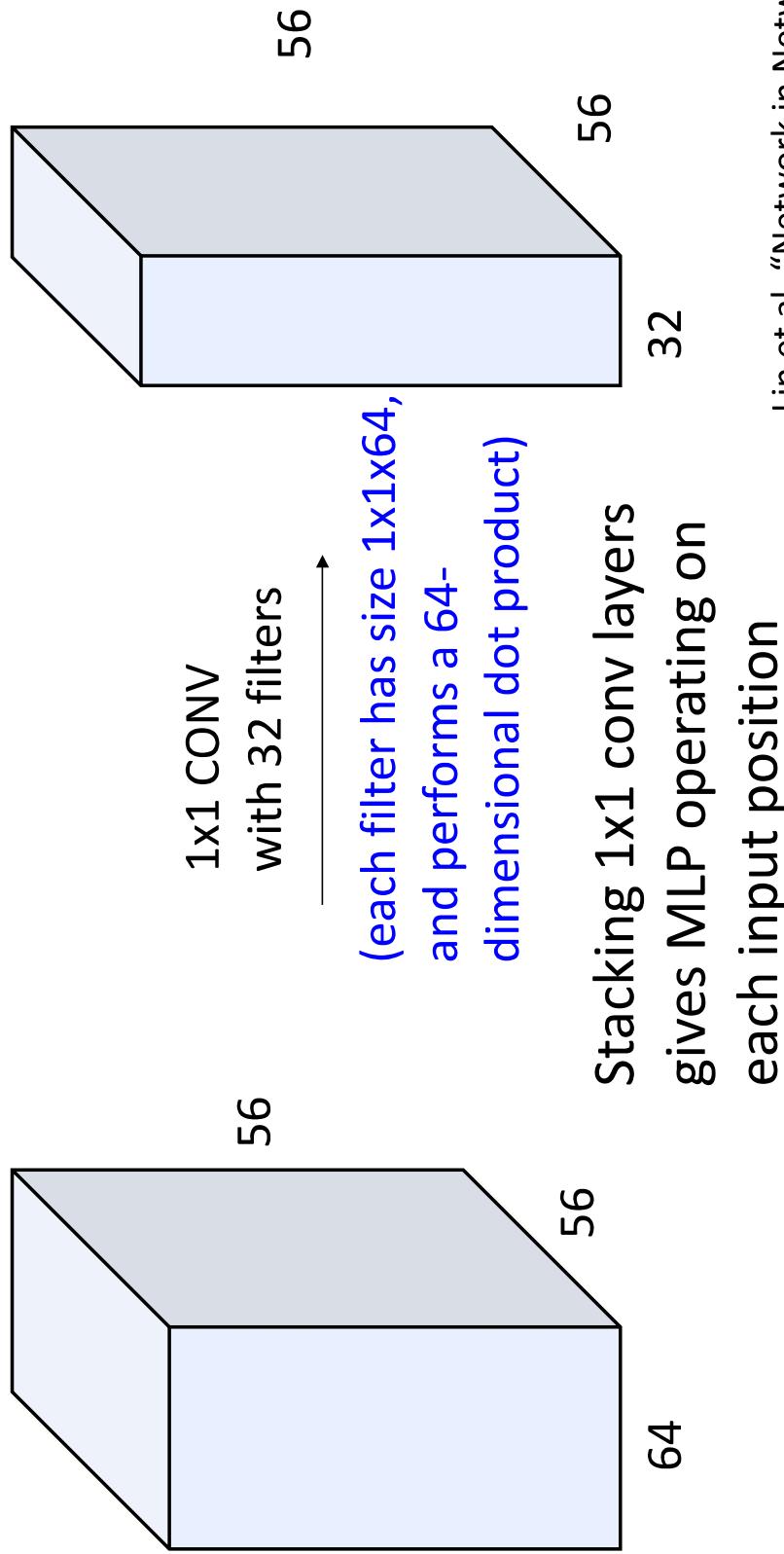
Number of multiply-add operations: **768,000**

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

Example: 1×1 Convolution



Example: 1×1 Convolution



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

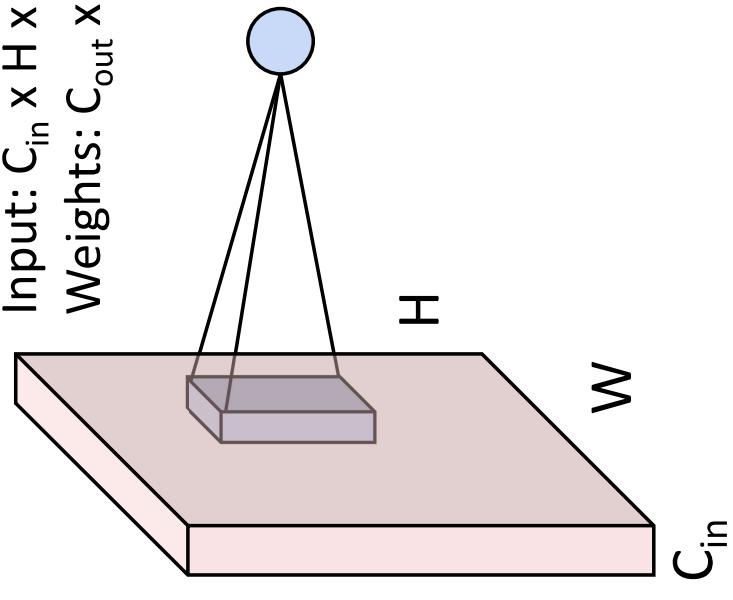
Hyperparameters:

- **Kernel size:** $K_H \times K_W$
 - **Number filters:** C_{out}
 - **Padding:** P
 - **Stride:** S
- 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)
 - $K = 3, P = 1, S = 1$ (3×3 conv)
 - $K = 5, P = 2, S = 1$ (5×5 conv)
 - $K = 1, P = 0, S = 1$ (1×1 conv)
 - $K = 3, P = 1, S = 2$ (Downsample by 2)
- 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:
- $H' = (H - K + 2P) / S + 1$
 - $W' = (W - K + 2P) / S + 1$

Other types of convolution

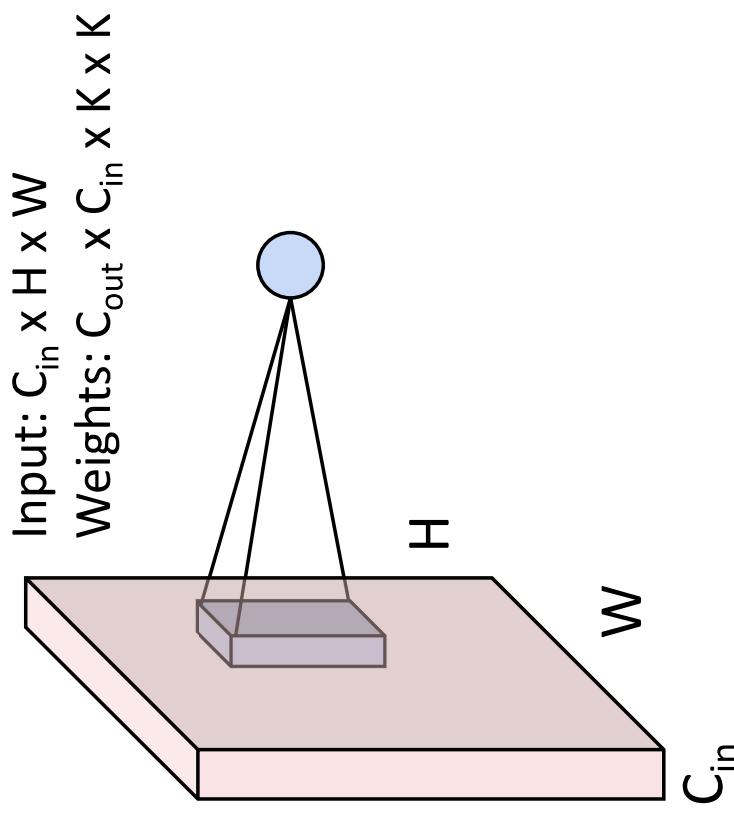
So far: 2D Convolution

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



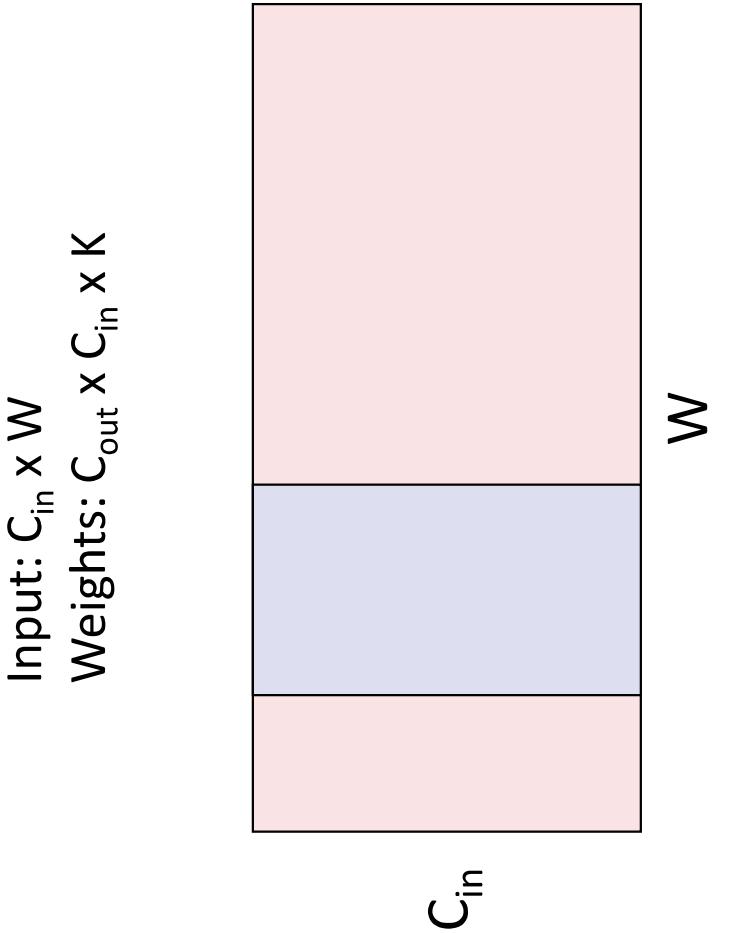
Other types of convolution

So far: 2D Convolution



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

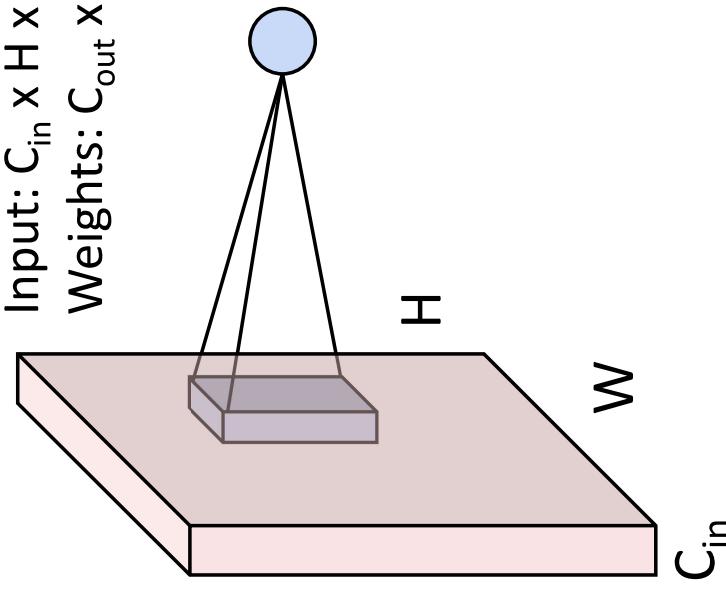
1D Convolution



Other types of convolution

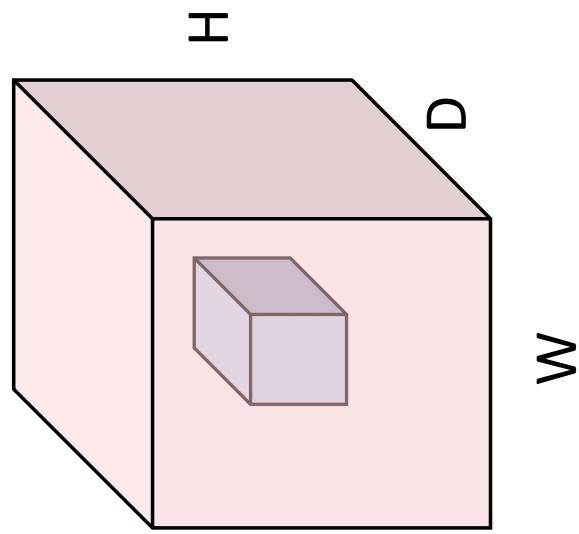
So far: 2D Convolution

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



3D Convolution

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



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

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')`

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_{in}, H, W) and output $(N, C_{\text{out}}, H_{\text{out}}, W_{\text{out}})$ can be precisely described as:

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

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’)
```

Conv1d

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

Conv3d

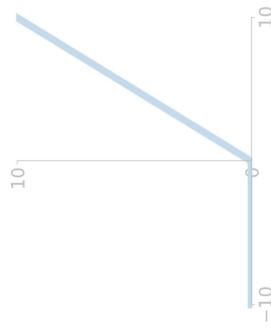
```
CLASS torch.nn.Conv3d(in_channels, out_channels, kernel_size, stride=1, padding=0,  
dilation=1, groups=1, bias=True, padding_mode=‘zeros’)
```

Components of a Convolutional Network

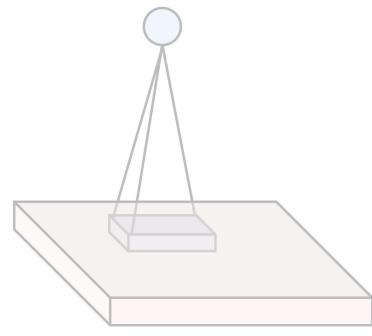
Fully-Connected Layers



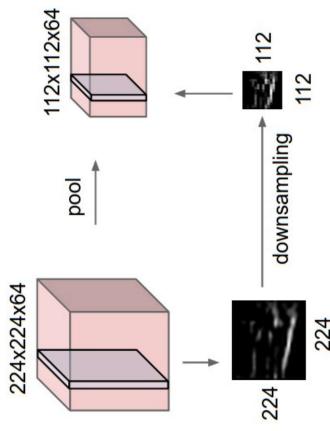
Activation Function



Convolution Layers



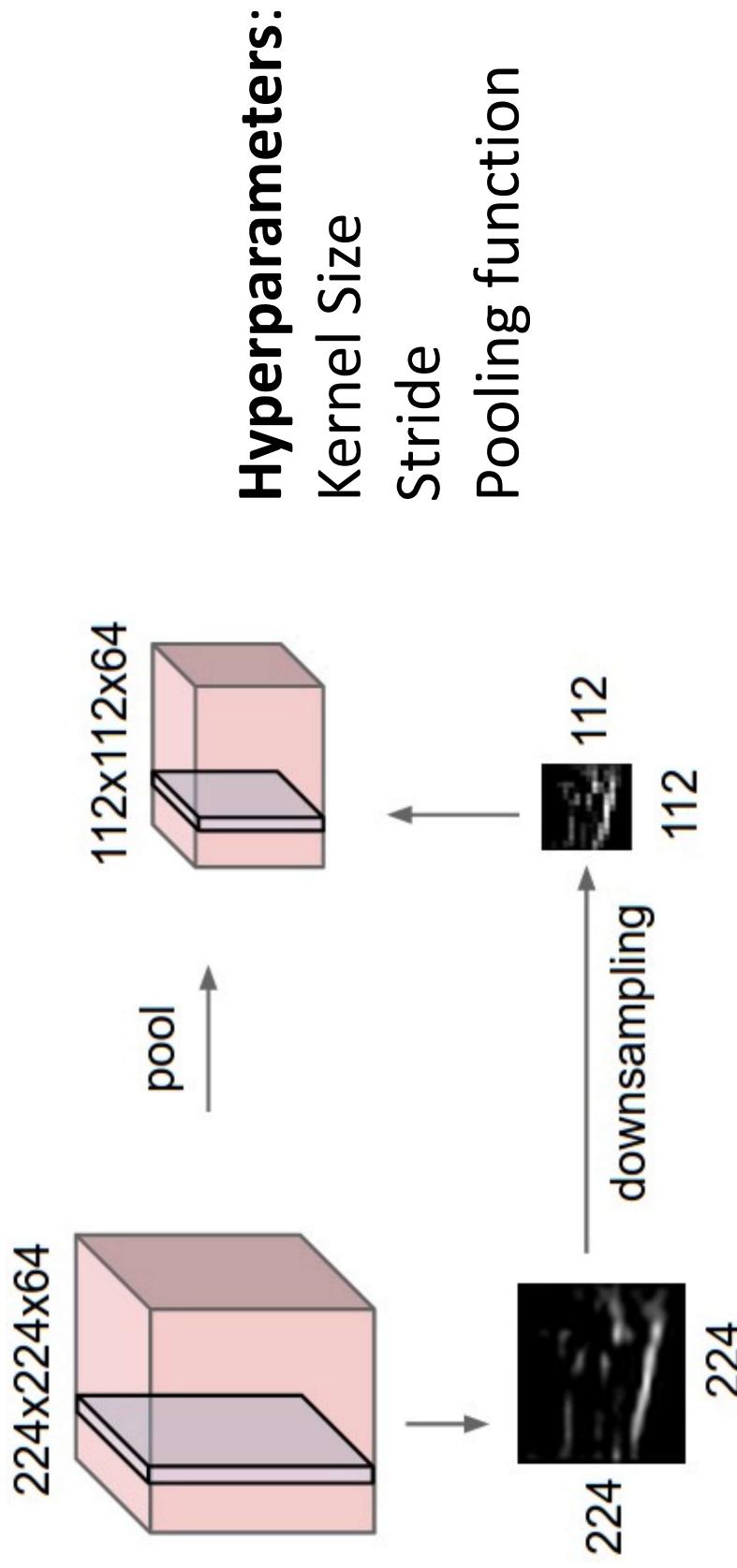
Pooling Layers



Normalization

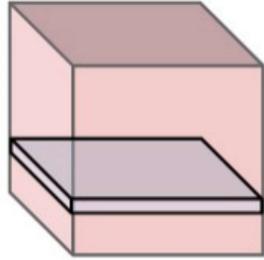
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Pooling Layers: Another way to downsample



Max Pooling

224x224x64



Single depth slice

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

x

6	8
3	4

Max pooling with 2x2
kernel size and stride 2

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

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

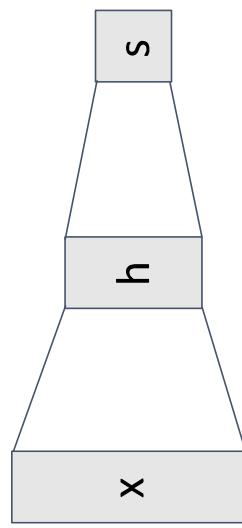
Learnable parameters: None!

Common settings:

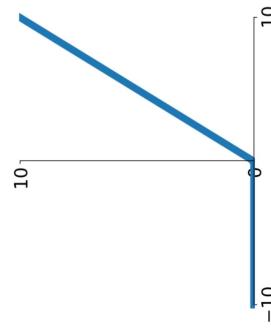
- max, $K = 2, S = 2$
- max, $K = 3, S = 2$ (AlexNet)

Components of a Convolutional Network

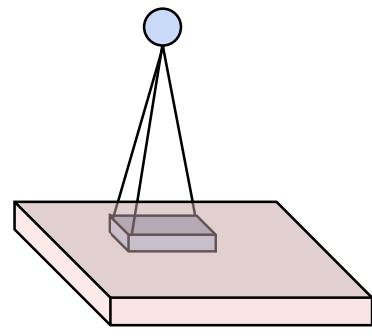
Fully-Connected Layers



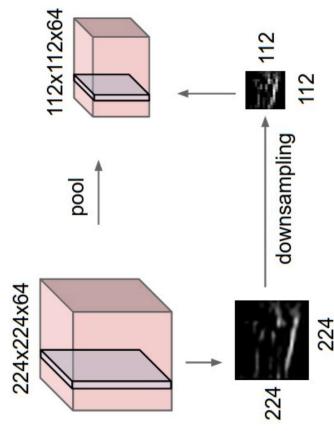
Activation Function



Convolution Layers



Pooling Layers



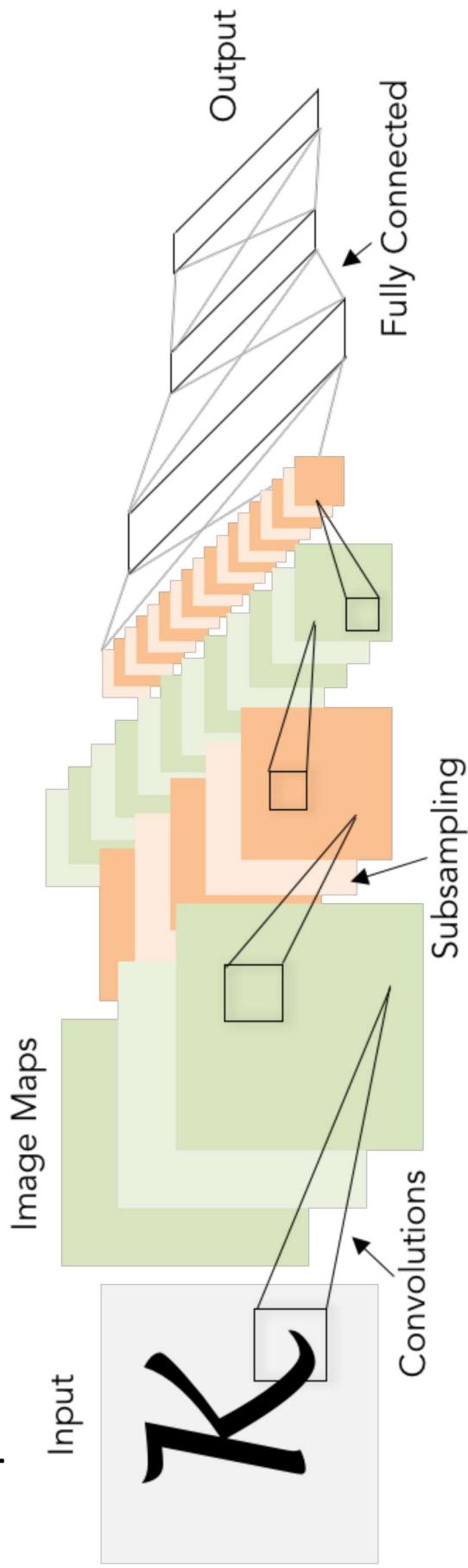
Normalization

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

Convolutional Networks

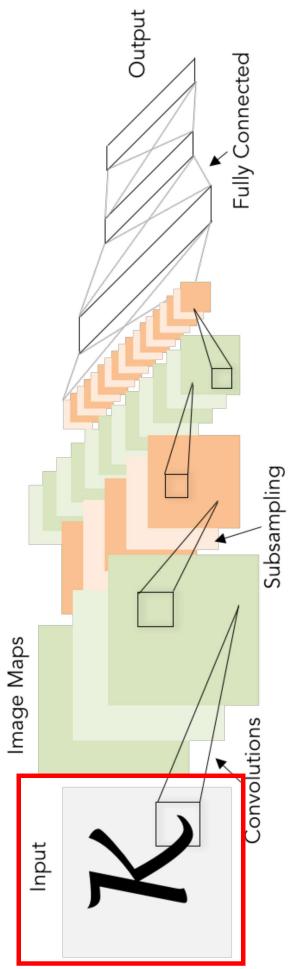
Classic architecture: [Conv, ReLU, Pool] $\times N$, flatten, [FC, ReLU] $\times N$, FC

Example: LeNet-5



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

Example: LeNet-5



Layer	Output Size	Weight Size
Input	$1 \times 28 \times 28$	

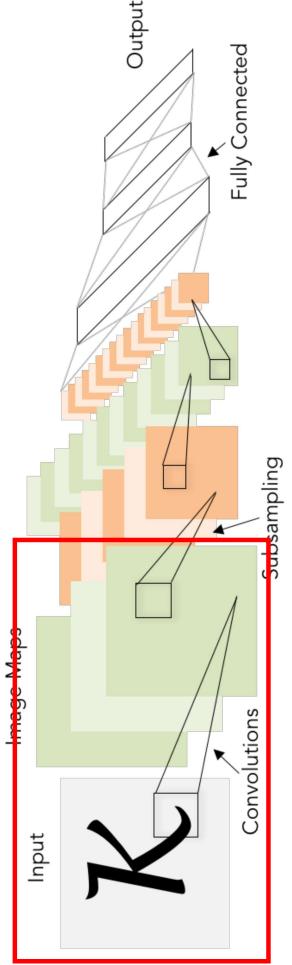
Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 68

September 24, 2019

Example: LeNet-5



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

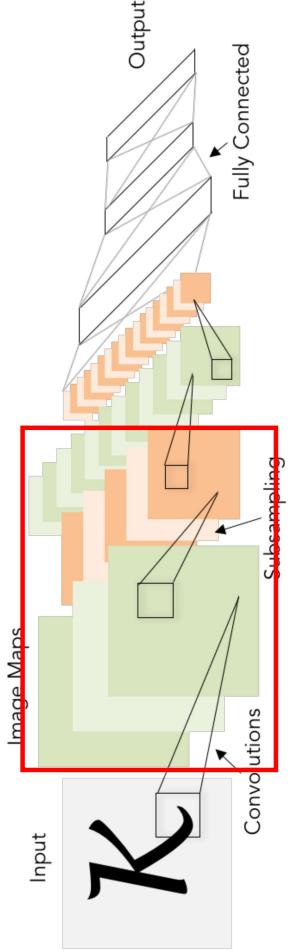
Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 69

September 24, 2019

Example: LeNet-5



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

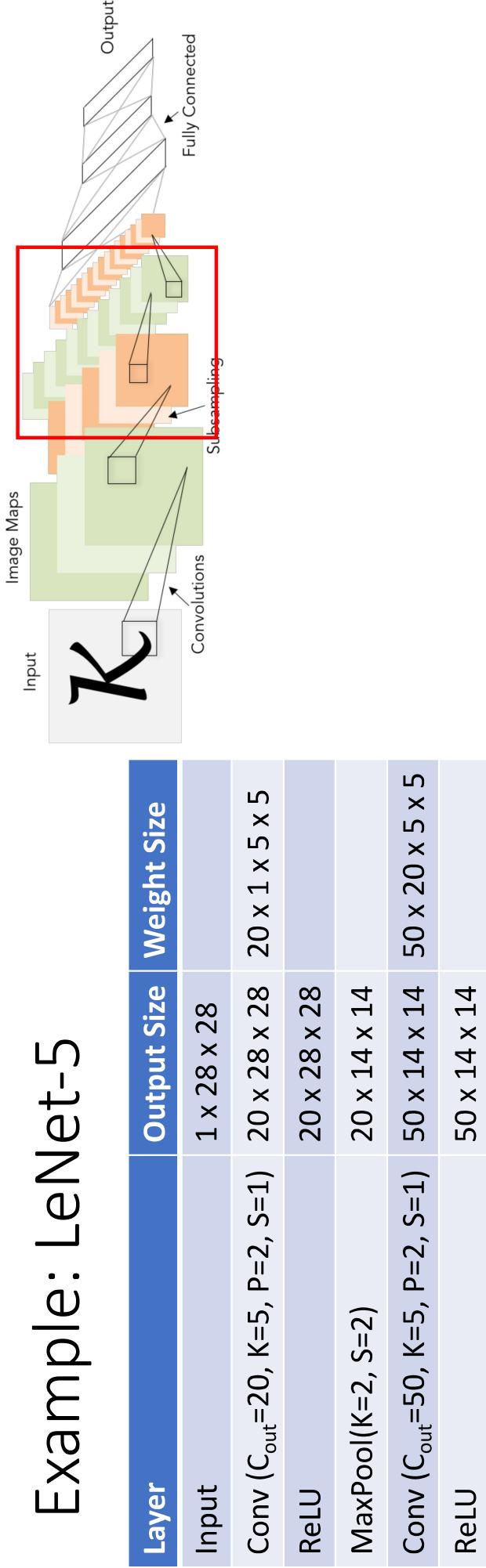
Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 70

September 24, 2019

Example: LeNet-5



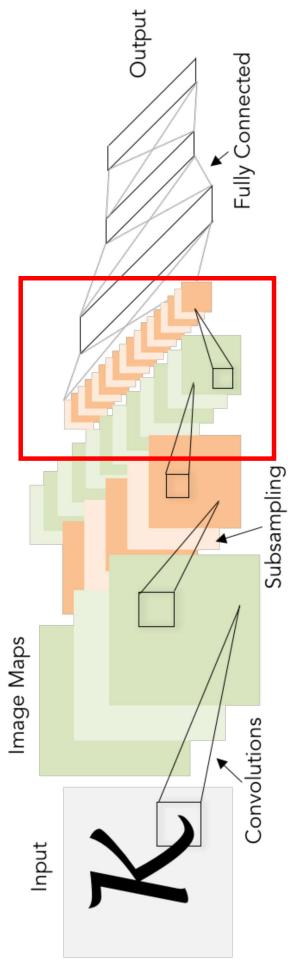
Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 71

September 24, 2019

Example: LeNet-5

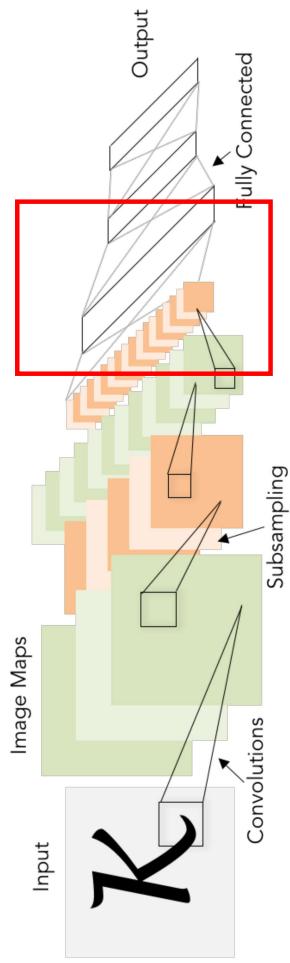


Layer	Output Size	Weight Size
Input	$1 \times 28 \times 28$	
Conv ($C_{out}=20$, $K=5$, $P=2$, $S=1$)	$20 \times 28 \times 28$	$20 \times 1 \times 5 \times 5$
ReLU	$20 \times 28 \times 28$	
MaxPool($K=2$, $S=2$)	$20 \times 14 \times 14$	
Conv ($C_{out}=50$, $K=5$, $P=2$, $S=1$)	$50 \times 14 \times 14$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

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

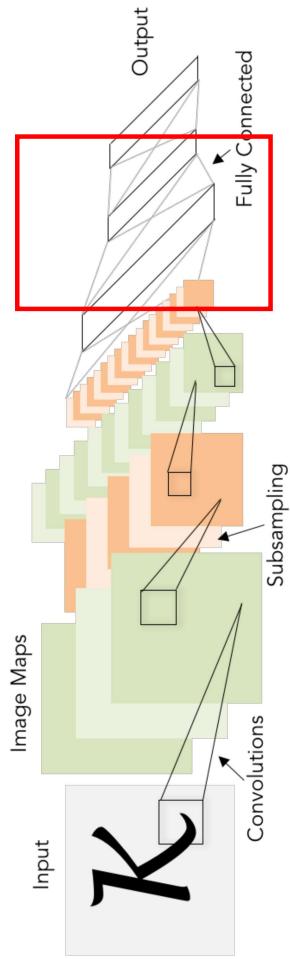
Example: LeNet-5

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



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

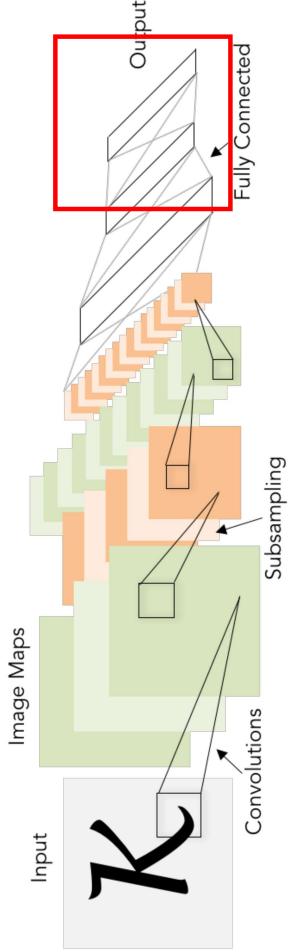
Example: LeNet-5



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

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

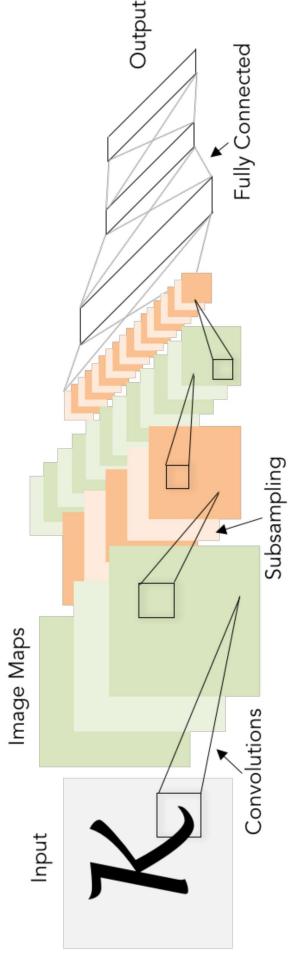
Example: LeNet-5



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

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

Example: LeNet-5



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

As we go through the network:
Spatial size decreases
(using pooling or strided conv)

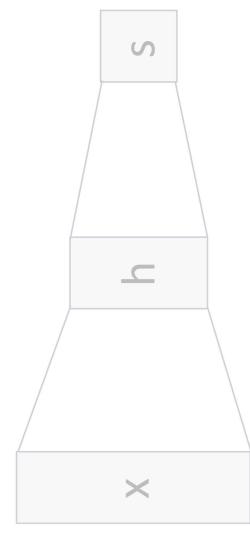
Number of channels increases
(total “volume” is preserved!)

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

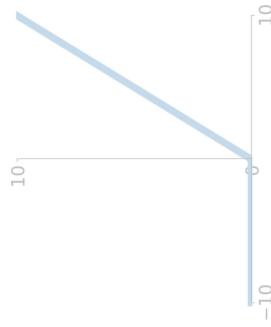
Problem: Deep Networks very hard to train!

Components of a Convolutional Network

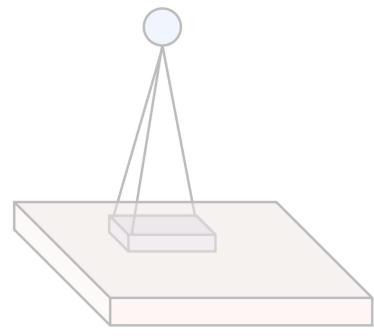
Fully-Connected Layers



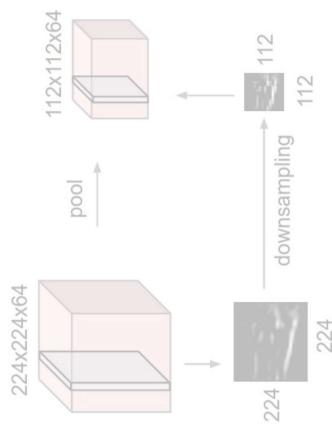
Activation Function



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? Helps reduce “internal covariate shift”, improves optimization

We can normalize a batch of activations like this:

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

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

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

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

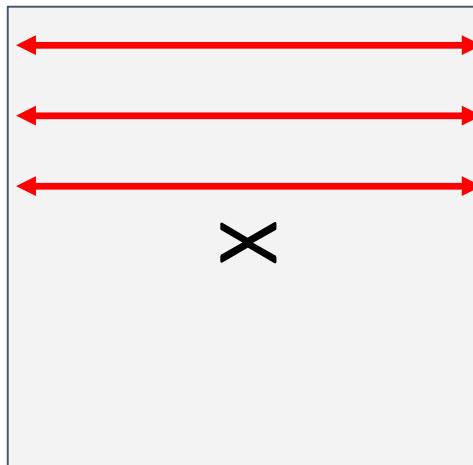
$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

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 $N \times D$

D



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

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

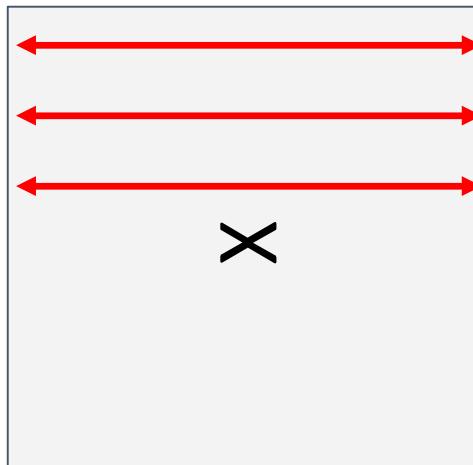
$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

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 $N \times D$

D



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

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

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

**Learnable scale and
shift parameters:**

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

Per-channel
std, shape is D

$\gamma, \beta : D$

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

Normalized x,
Shape is $N \times D$

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

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

Input: $x : N \times D$

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

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

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

Learnable scale and shift parameters:

$\gamma, \beta : D$

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

Normalized x,
Shape is $N \times D$

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

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

Output,
Shape is $N \times D$

Batch Normalization: Test-Time

Input: $x : N \times D$

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

Learnable scale and shift parameters:

$\gamma, \beta : D$

$\sigma_j^2 =$ (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 $N \times D$

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

$\beta = \mu$ will recover the identity function!

Batch Normalization: Test-Time

Input: $x : N \times D$

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

Learnable scale and shift parameters:

$\gamma, \beta : D$

$\sigma_j^2 =$ (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 $N \times D$

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 $N \times D$

Batch Normalization for ConvNets

Batch Normalization for
fully-connected networks

Batch Normalization for
convolutional networks
(Spatial Batchnorm, BatchNorm2D)

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

Normalize

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

Normalize

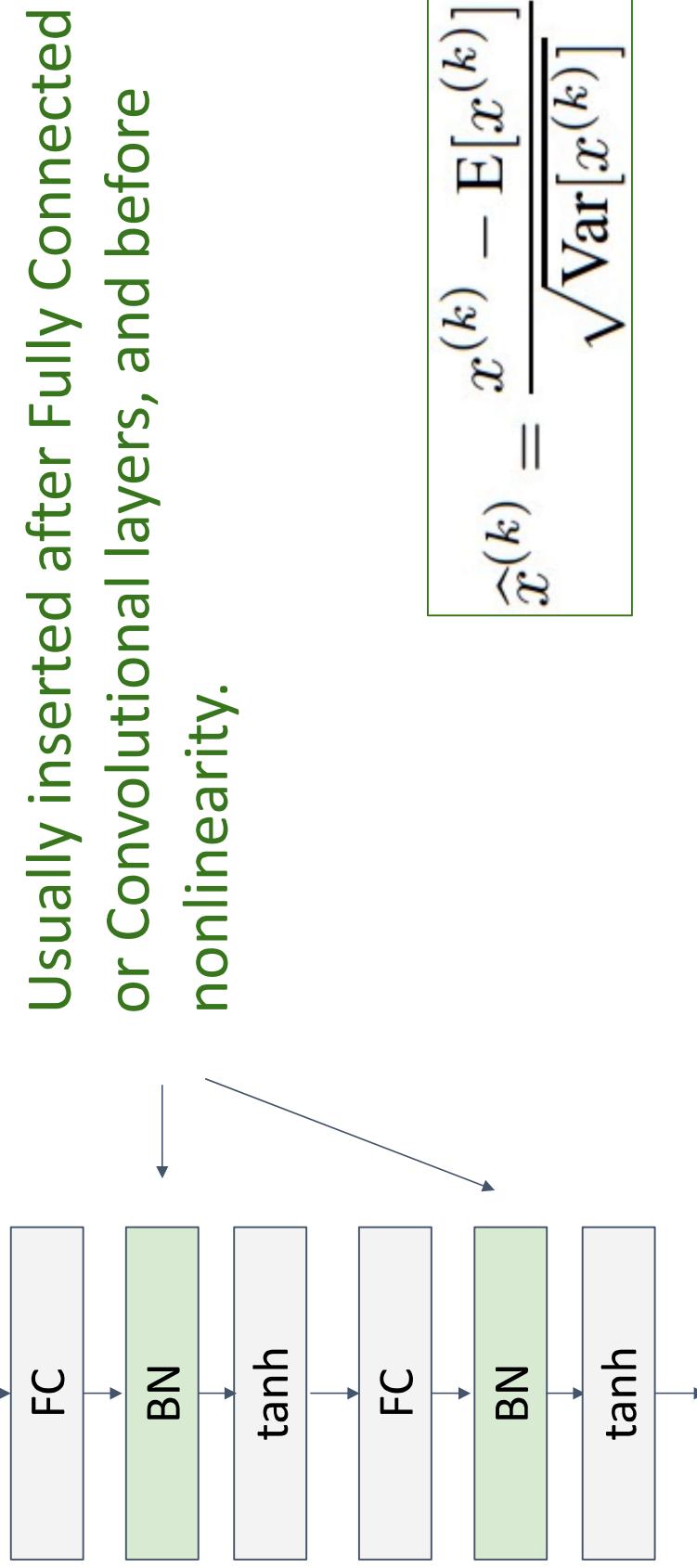
$$\boldsymbol{\mu}, \boldsymbol{\sigma} : 1 \times C \times 1 \times 1$$

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

$$y = \gamma(x - \mu) / \sigma + \beta$$

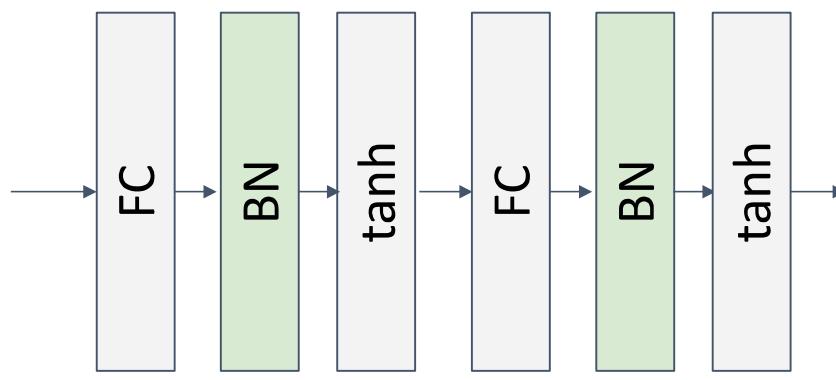
$$y = \gamma(x - \mu) / \sigma + \beta$$

Batch Normalization

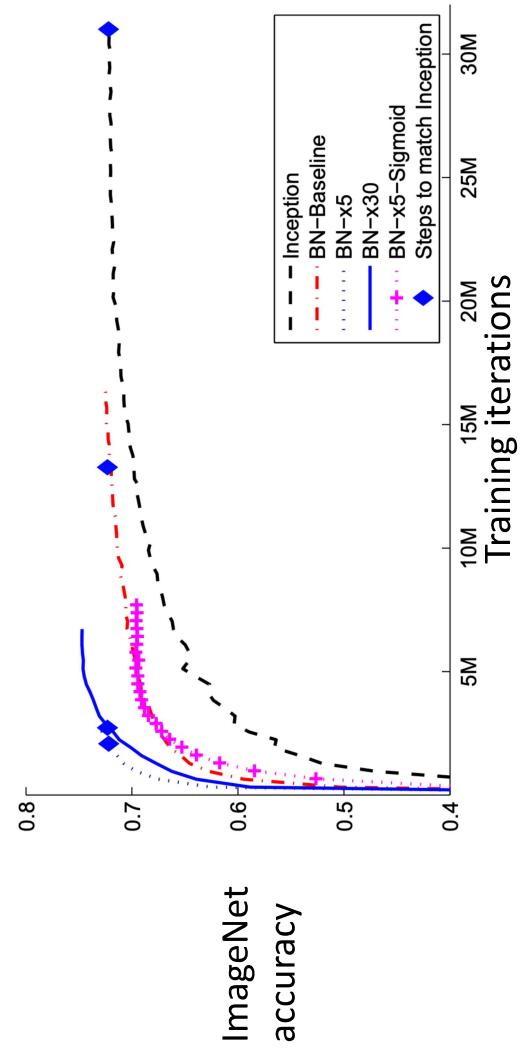


Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Batch Normalization

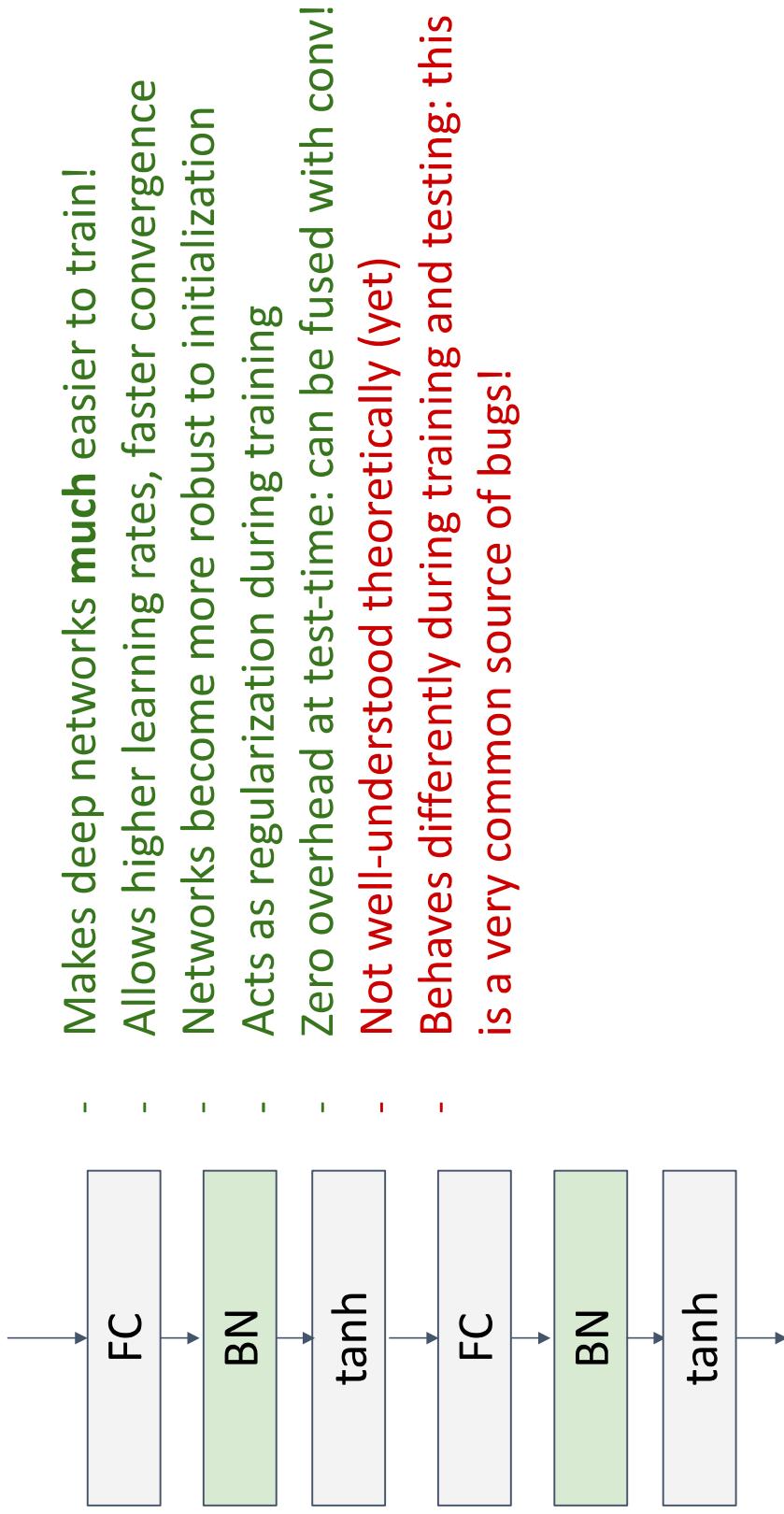


- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!



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

Batch Normalization



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

Layer Normalization

Batch Normalization for
fully-connected networks

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

Normalize

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

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

Layer Normalization for fully-connected networks

Same behavior at train and test!
Used in RNNs, Transformers

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

Normalize

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

$$\gamma, \beta : 1 \times D$$

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

Ba, Kiros, and Hinton, "Layer Normalization", arXiv 2016

Justin Johnson

Lecture 7 - 90

September 24, 2019

Instance Normalization

Batch Normalization for convolutional networks

Same behavior at train / test!

$$\begin{aligned}\mathbf{x} & : N \times C \times H \times W \\ & \text{Normalize} \\ \boldsymbol{\mu}, \boldsymbol{\sigma} & : 1 \times C \times 1 \times 1\end{aligned}$$

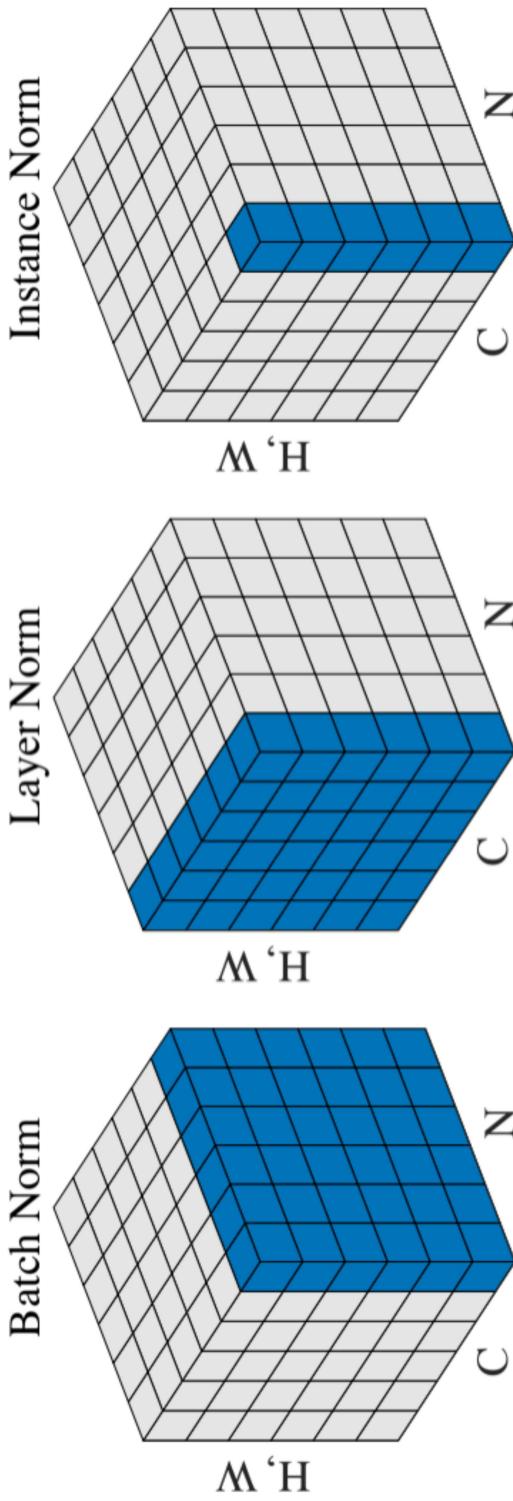
$$\begin{aligned}\gamma, \beta & : 1 \times C \times 1 \times 1 \\ y & = \gamma(x - \mu) / \sigma + \beta\end{aligned}$$

Instance Normalization for convolutional networks

$$\begin{aligned}\mathbf{x} & : N \times C \times H \times W \\ & \text{Normalize} \\ \boldsymbol{\mu}, \boldsymbol{\sigma} & : N \times C \times 1 \times 1\end{aligned}$$

$$\begin{aligned}\gamma, \beta & : 1 \times C \times 1 \times 1 \\ y & = \gamma(x - \mu) / \sigma + \beta\end{aligned}$$

Comparison of Normalization Layers



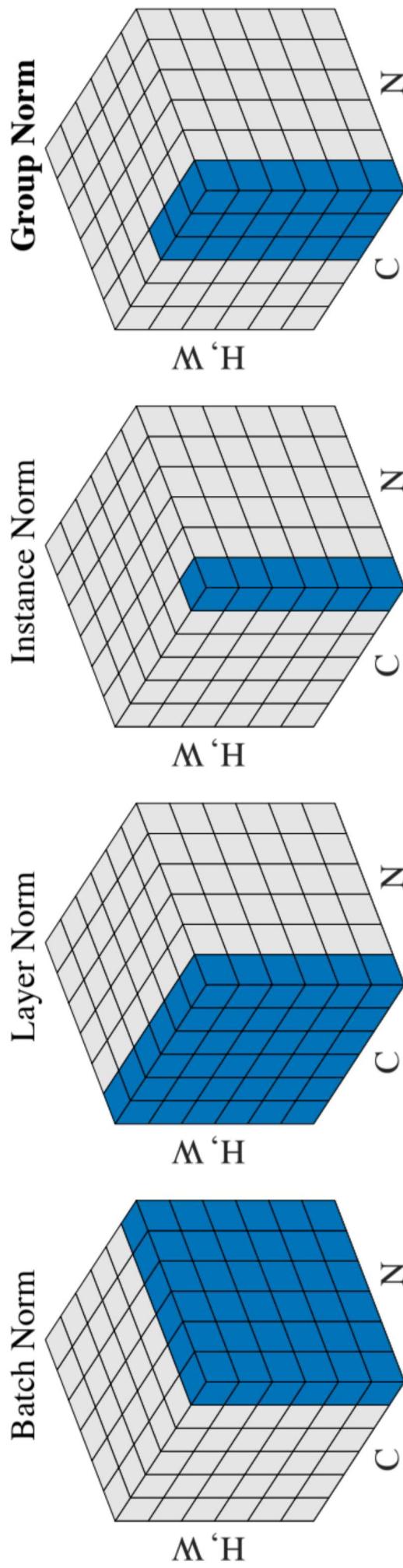
Wu and He, "Group Normalization", ECCV 2018

Justin Johnson

Lecture 7 - 92

September 24, 2019

Group Normalization



Wu and He, "Group Normalization", ECCV 2018

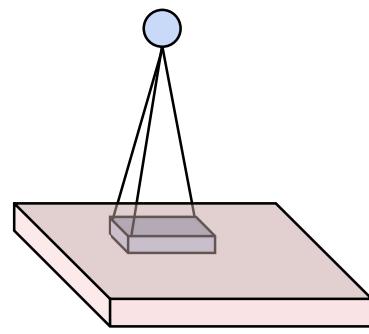
Justin Johnson

Lecture 7 - 93

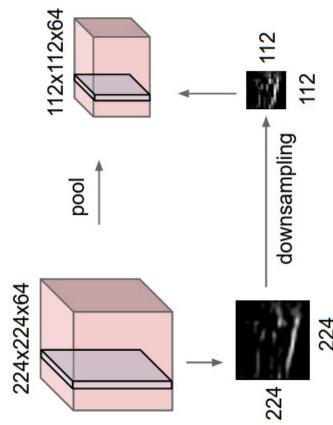
September 24, 2019

Components of a Convolutional Network

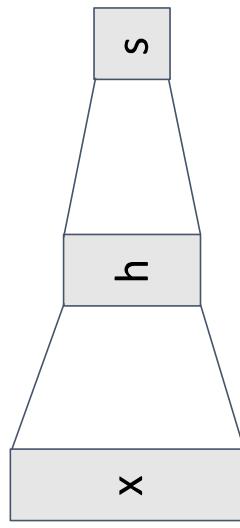
Convolution Layers



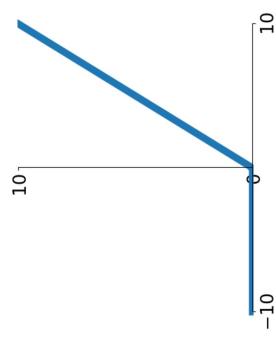
Pooling Layers



Fully-Connected Layers



Activation Function

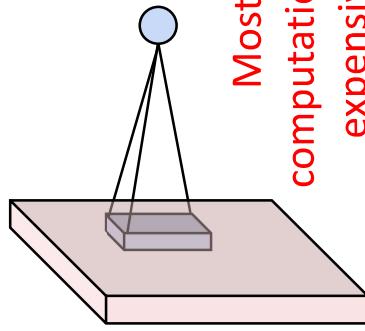


Normalization

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

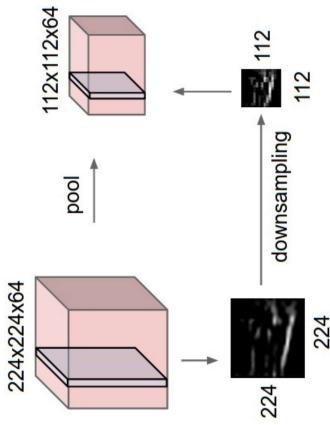
Components of a Convolutional Network

Convolution Layers

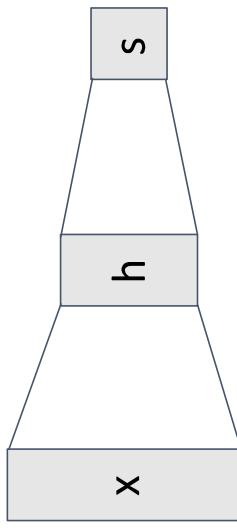


Most computationally expensive!

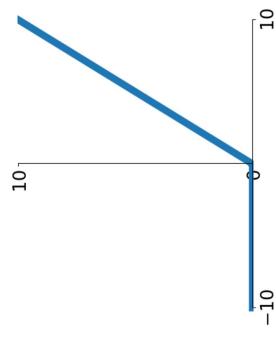
Pooling Layers



Fully-Connected Layers



Activation Function

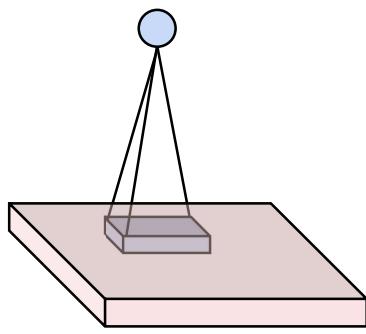


Normalization

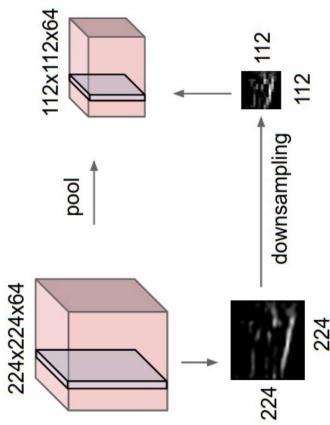
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Summary: Components of a Convolutional Network

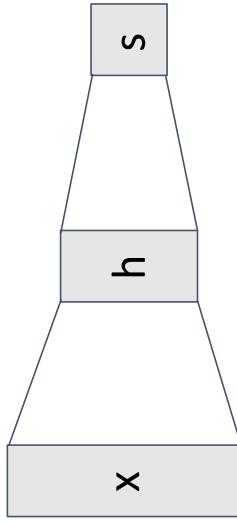
Convolution Layers



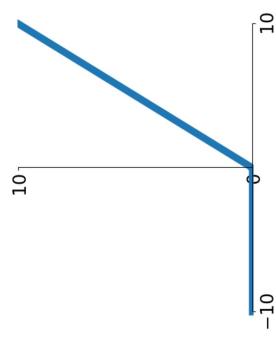
Pooling Layers



Fully-Connected Layers



Activation Function

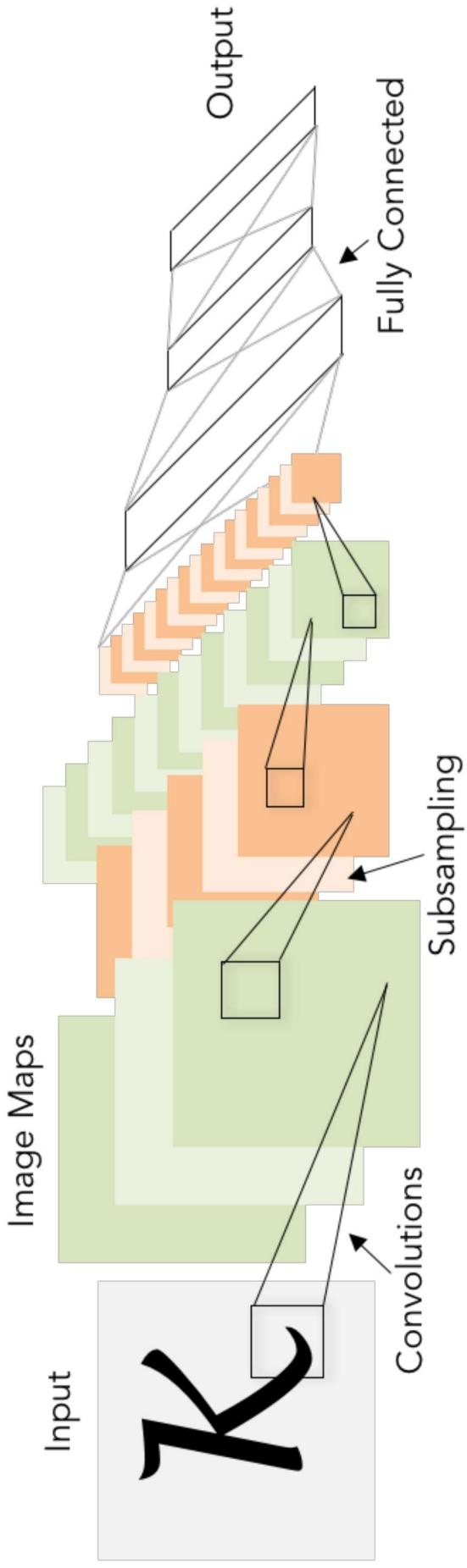


Normalization

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

Summary: Components of a Convolutional Network

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



Next time:
CNN Architectures