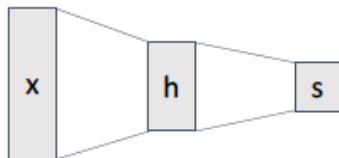


Lecture 7- Convolutional Networks

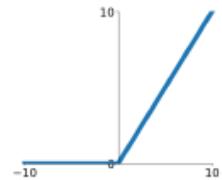
1. Components of a Network

- Components of a Fully-connected Network

Fully-Connected Layers



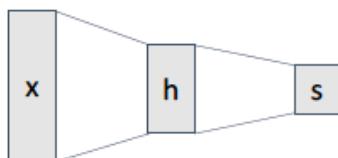
Activation Function



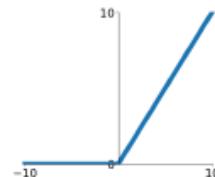
- Components of a Convolutional Network

--> Define new computational nodes that operate on images --> respect the spatial structure of images

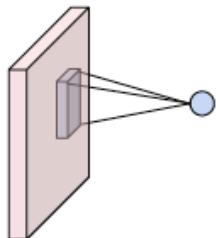
Fully-Connected Layers



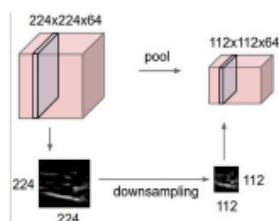
Activation Function



Convolution Layers



Pooling Layers



Normalization

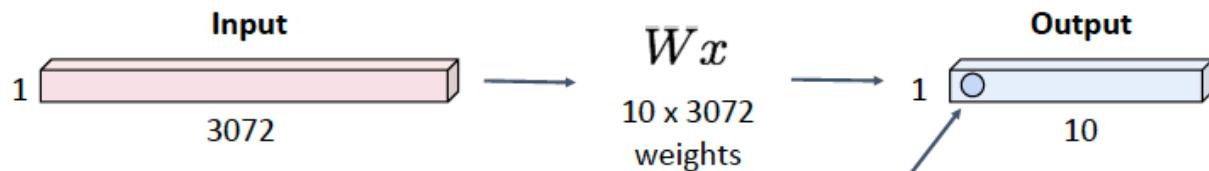
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

2. Components of a Convolutional Network

- **Fully-Connected Layer**

Input: Image height \times Image width \times Channel number

32x32x3 image -> stretch to 3072 x 1

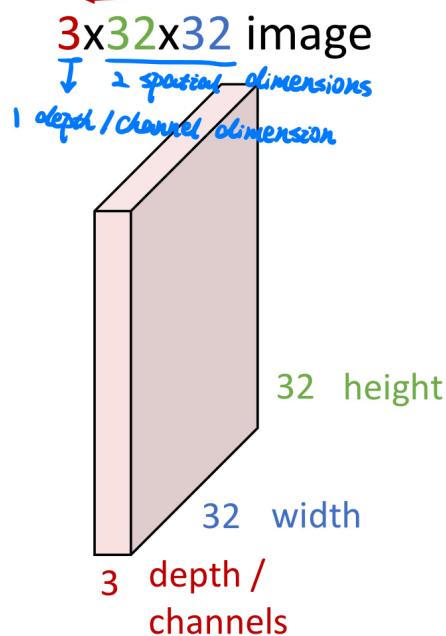


1 number:

the result of taking a dot product between a row of W and the input (a 3072-dimensional dot product)

- **Convolutional Layer**

Convolution Layer



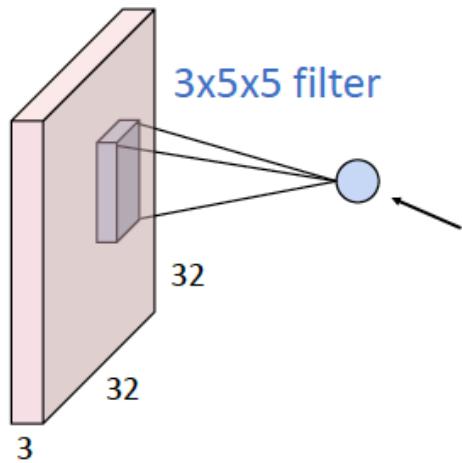
Filters always extend the full depth of the input volume

3x5x5 filter



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

3x32x32 image



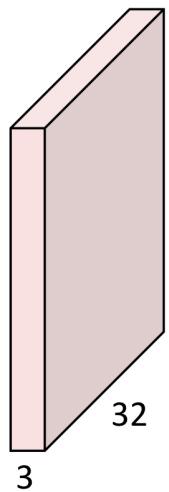
1 number:

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

$$w^T x + b$$

Convolution Layer

3x32x32 image



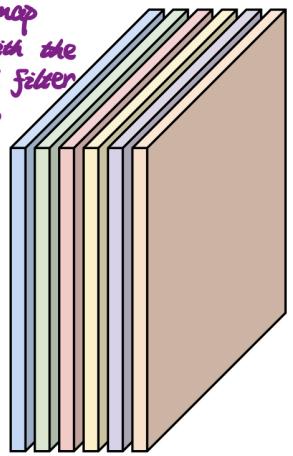
Also 6-dim bias vector:

$$\begin{array}{c} \text{blue} \\ \text{green} \\ \text{red} \\ \text{yellow} \\ \text{purple} \\ \text{brown} \end{array} b$$

⑤

6 activation maps,
each 1x28x28

→ each map
corresponds with the
response of each filter
→ each feature



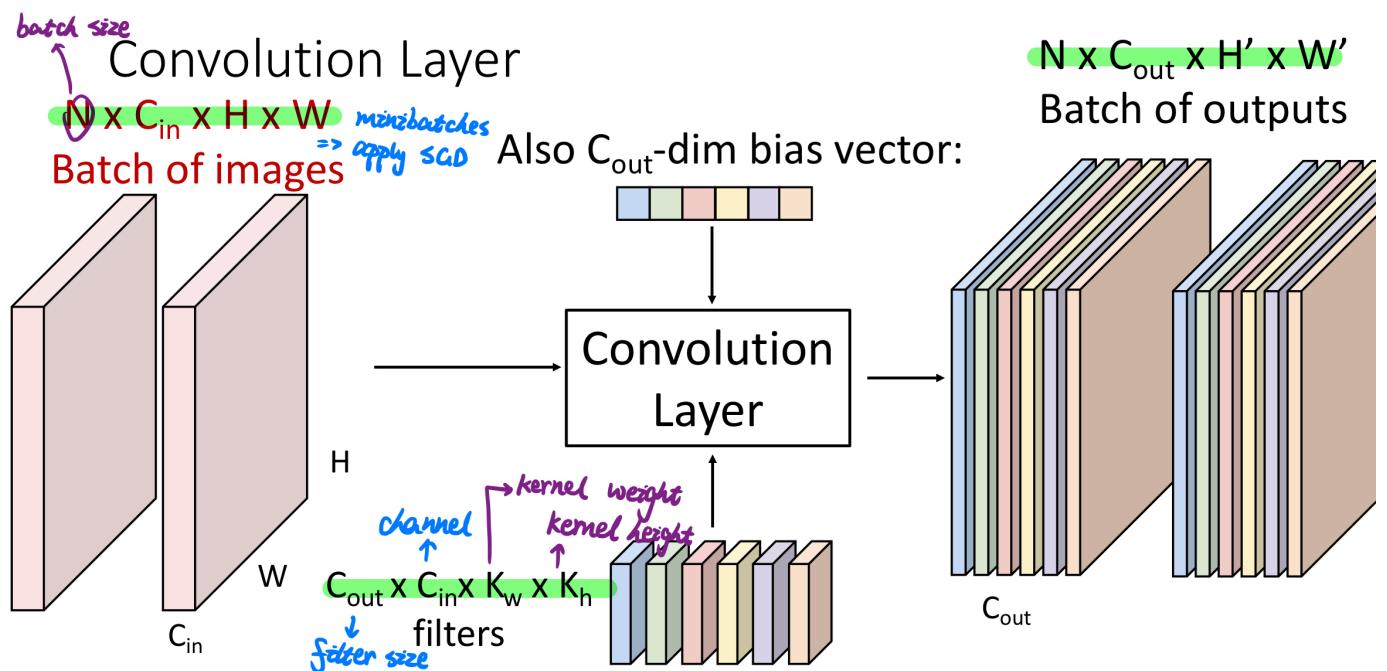
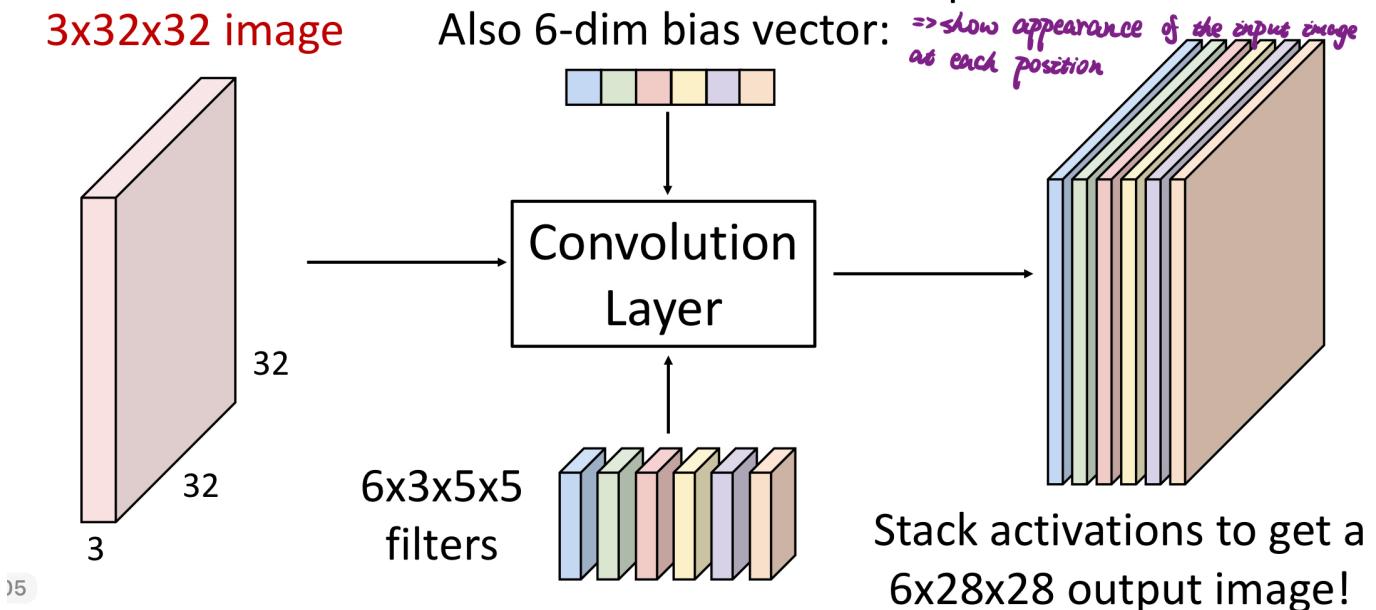
6 filters
 $\rightarrow w^T x$

Convolution
Layer

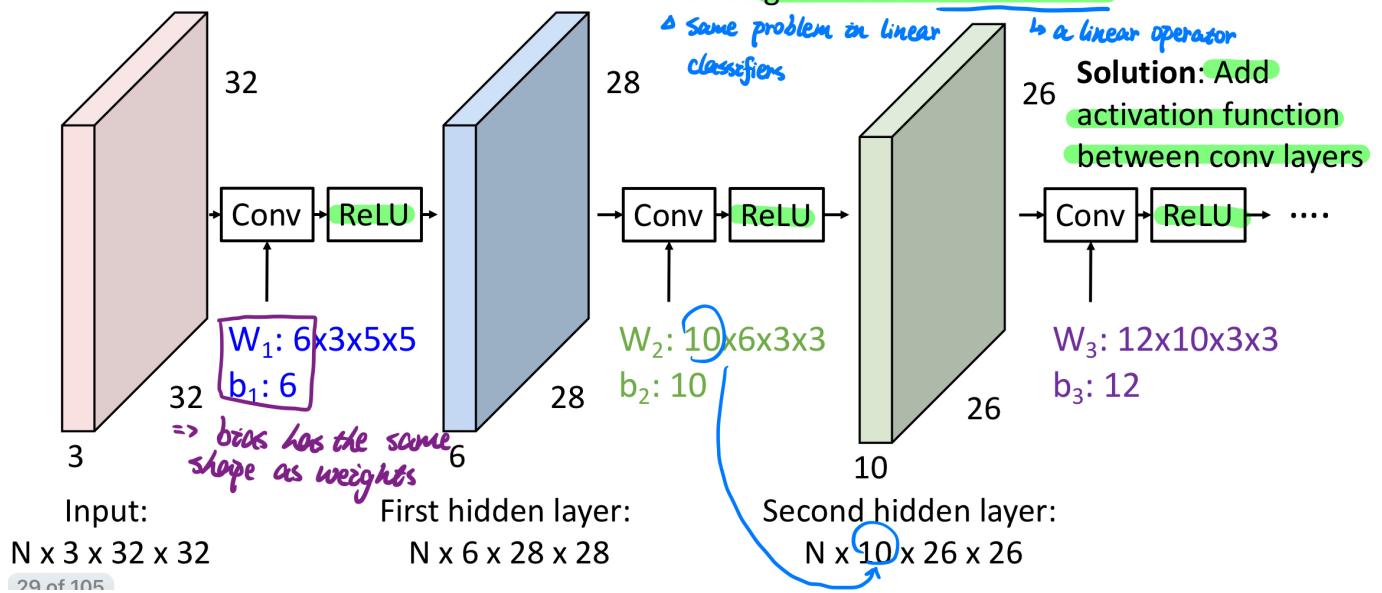
Stack activations to get a
6x28x28 output image!

--> One filter is a template for extracting one feature

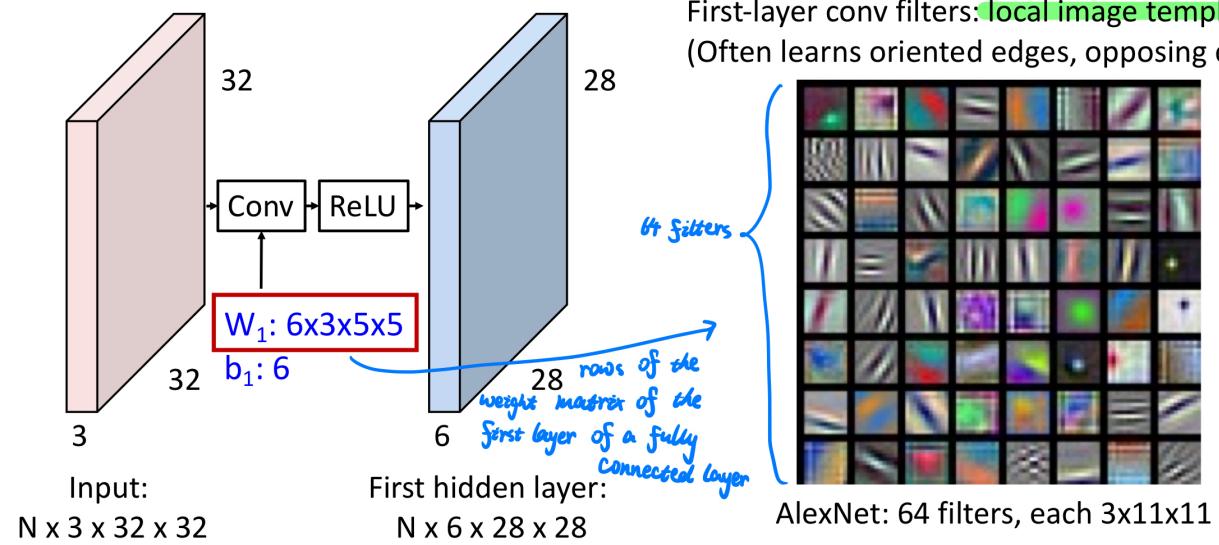
Convolution Layer



Stacking Convolutions



- Convolutional Filters



- 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

(W, K, P are all sizes)

- o Receptive Fields

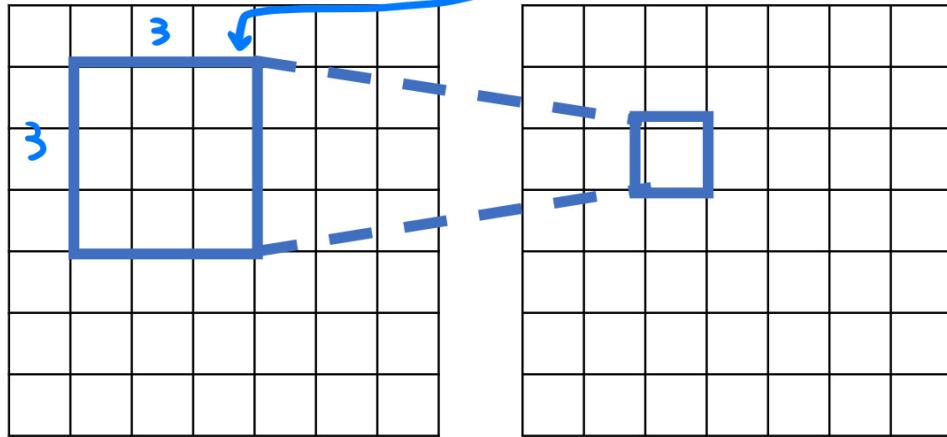
- Definition:

When dealing with high-dimensional inputs such as images, it is impractical to connect neurons to all neurons in the previous volume. Instead, we will **connect each neuron to only a local region of the input volume.**

The spatial extent of this connectivity is a hyperparameter called the **receptive field** of the neuron (equivalently this is the **filter size**).

- In the input:

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



Input

Output

- In the previous layers:

Input: 7x7

Filter: 3x3

Output: 5x5

In general:

Input: W

Filter: K

Padding: P

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

Very common:

[Same padding] Set $P = (K - 1) / 2$ to

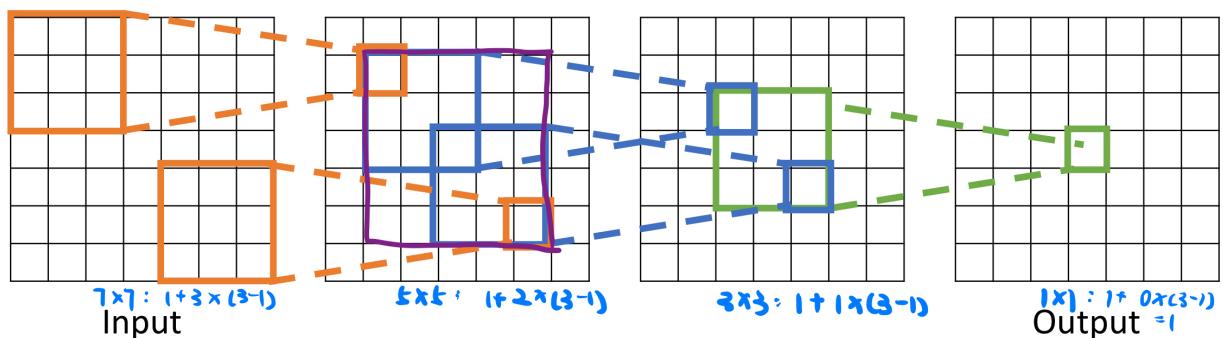
make output have same size as input!

Receptive Fields

$\rightarrow k: \text{filter size, in this case } k=3$

Each successive convolution adds $K-1$ to the receptive field size

With L layers the receptive field size is $1 + L * (K-1)$

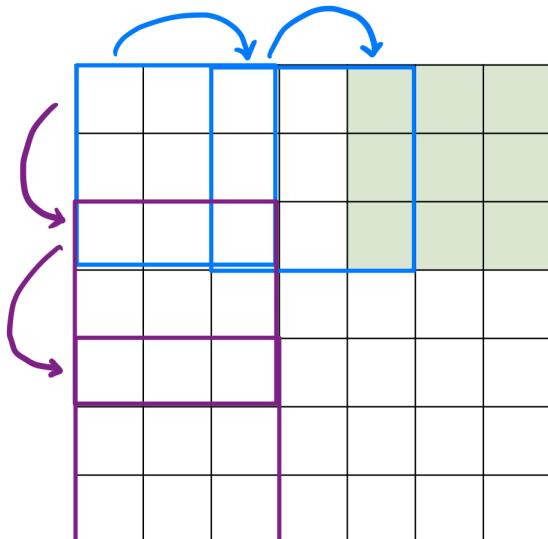


--> Given kernel width k and stride s , for m adjacent pixels in the activation output, **cummulative receptive field n** with respect to layer input is $n = k + s(m - 1)$

- Problem

For large images, many layer for each output to "see" the whole image image --> result in computational inefficiency --> Solution: **Downsample** inside the network --> *Strided Convolution*

- Strided Convolution



Input: 7x7

Filter: 3x3

Output: 3x3

Stride: 2

In general:

Input: W

Filter: K

Padding: P

Stride: S

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

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

10 5x5 filters with stride 1, pad 2

$$c(W - k + 2P) / S + 1$$

Output volume size:

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

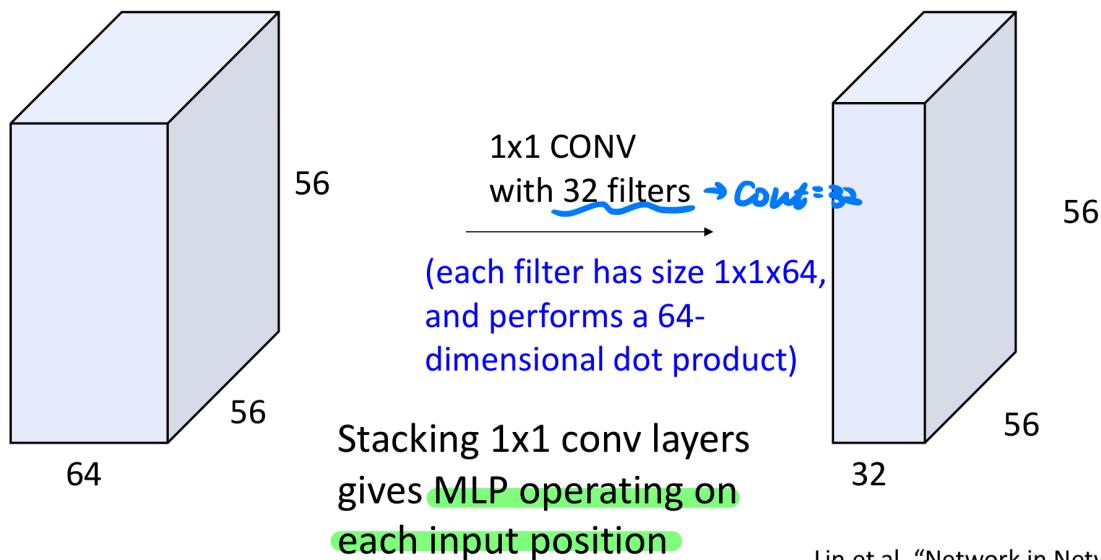
10 x 32 x 32

Number of learnable parameters: /bu
 $\{$ Parameters per filter: $3 \times 5 \times 5 + 1$ (for bias) = 76
 10 filters, so total is $10 \times 76 = 760$

Number of multiply-add operations: 768,000 convolution
 $\Delta 10 \times 32 \times 32 = 10,240$ outputs; each output is the inner product of two $3 \times 5 \times 5$ tensors (75 elems); total = $75 \times 10240 = 768K$

- 1x1 convolution

--> implement 3-dimensional dot products



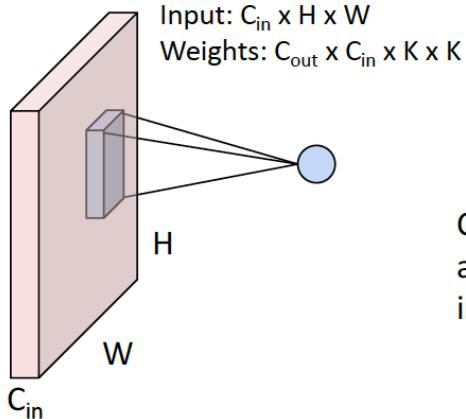
- 1D / 3D convolution

1D Convolution

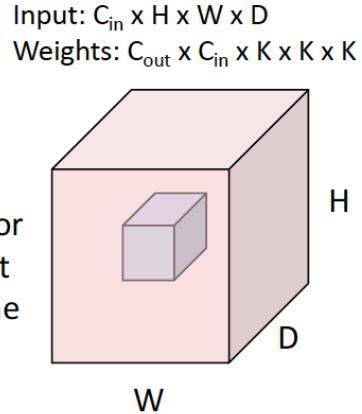
Input: $C_{in} \times W$
 Weights: $C_{out} \times C_{in} \times K$



So far: 2D Convolution



3D 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:

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

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

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

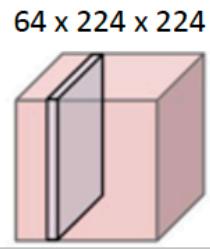
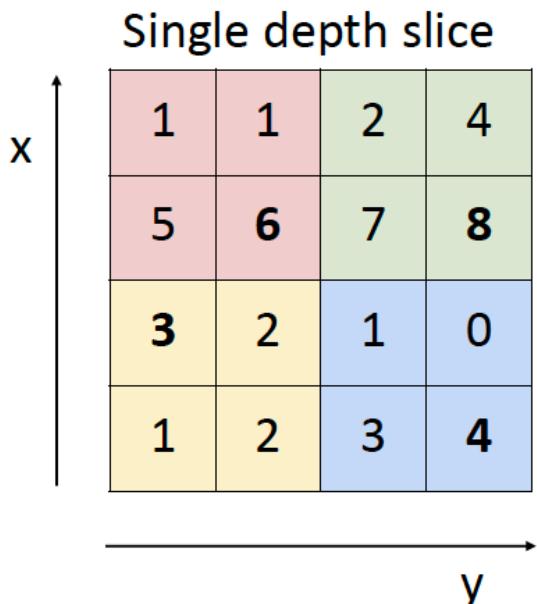
• **Pooling Layers: Another way to downsample**

- Properties:

- No learnable parameters
- Make the image smaller --> reduce computational complexity
- Hyperparameters: **Kernel size, Stride, Pooling function**
- Downampling methods:
 - Pooling: cheap functions, easy to implement

- Strided convolution: learnable downsample

Max Pooling



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!

- Convolutional Networks

↳ **Downsample**:

- pooling: cheap functions, easy to implement
- strided convolution: learnable downsample

✳ Common settings:

max, $K = 2, S = 2$ ($k=5$)

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

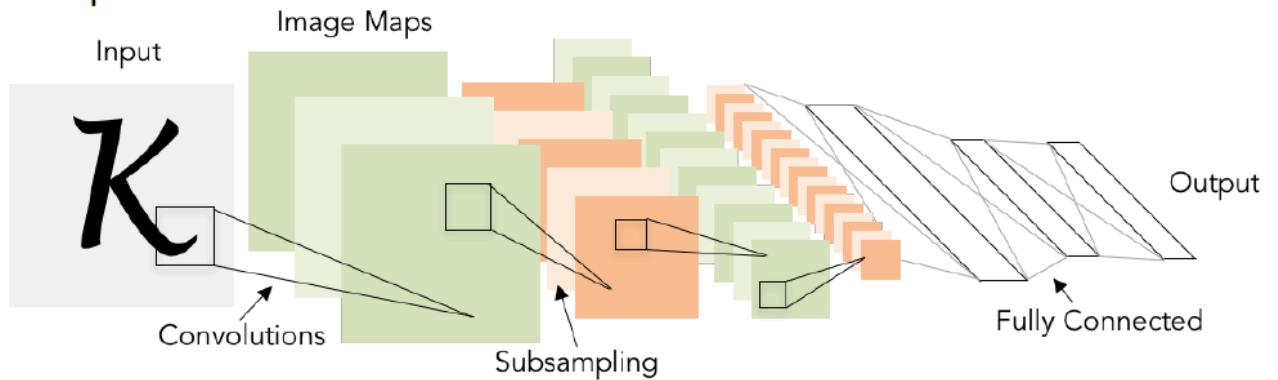
↳ **average pooling** (linear classifiers)

→ a certain convolutional layer

=> can implement average pooling via convolution

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

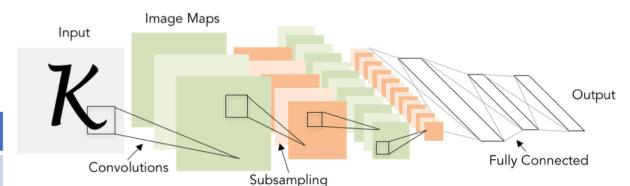
Example: LeNet-5



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 \rightarrow 500)	500	2450×500
ReLU	500	
Linear (500 \rightarrow 10)	10	500×10

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



As we go through the network:

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

 Number of channels **increases**
 (total "volume" is preserved!)
 => more learnable parameters
 => learn more complex patterns
 Some modern architectures
 break this trend -- stay tuned!

- **Batch Normalization**

- Problem: deep networks very hard to train! --> as input matrix can have large values as well as negligible values --> want to process more smoothly

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

\Rightarrow add an extra normalization layer

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!

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

Batch Normalization
At training time

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

Learnable scale and shift parameters:

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

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

Problem: Estimates depend on

minibatch; can't do this at test-time!

~~=> can't calculate the parameters~~

~~may mix information across mini-batch~~

Per-channel mean, shape is D

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

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

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

Normalized x,
Shape is N x D

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

Output,
Shape is N x D

$$\mu_j^{test} = 0$$

For each training iteration:

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

$$\mu_j^{test} = 0.99 \mu_j^{test} + 0.01 \mu_j$$

(Similar for σ)

μ_j : (Running) average of values seen during training

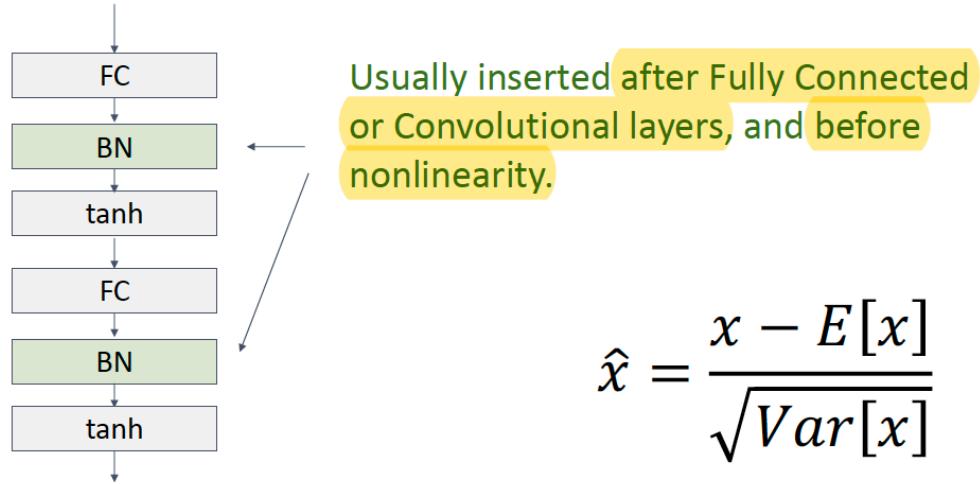
σ_j^2 : (Running) average of values seen during training

■ Batch Normalization during test-time

- Batch normalization during test time

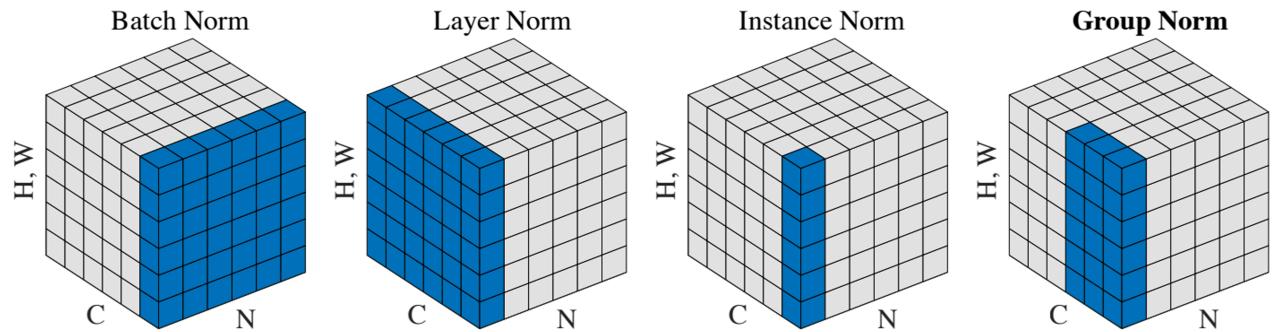
--> become a **linear operator** and can be fused with the previous
fully-connected or conv layer --> no runtime during test time

- o Properties



- 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!

- o Comparison of Normalization Layers



Layer Normalization

Batch Normalization for **fully-connected** networks *no subtraction along all elements of the mini-batch*

$$x : N \times D$$

Normalize

$$\mu, \sigma : 1 \times D$$

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

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

Batch Normalization for convolutional networks

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

Normalize

$$\mu, \sigma : 1 \times C \times 1 \times 1$$

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

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

Layer Normalization for fully-connected networks

Same behavior at train and test!

Used in RNNs, Transformers

cannot fine-tune FC/Conv

$$x : N \times D$$

Normalize

$$\mu, \sigma : N \times 1$$

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

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

Instance Normalization for convolutional networks

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

Normalize

$$\mu, \sigma : N \times C \times 1 \times 1$$

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

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