# Introduction to Deep Learning

## Nandita Bhaskhar

Content adapted from CS231n and past CS229 teams
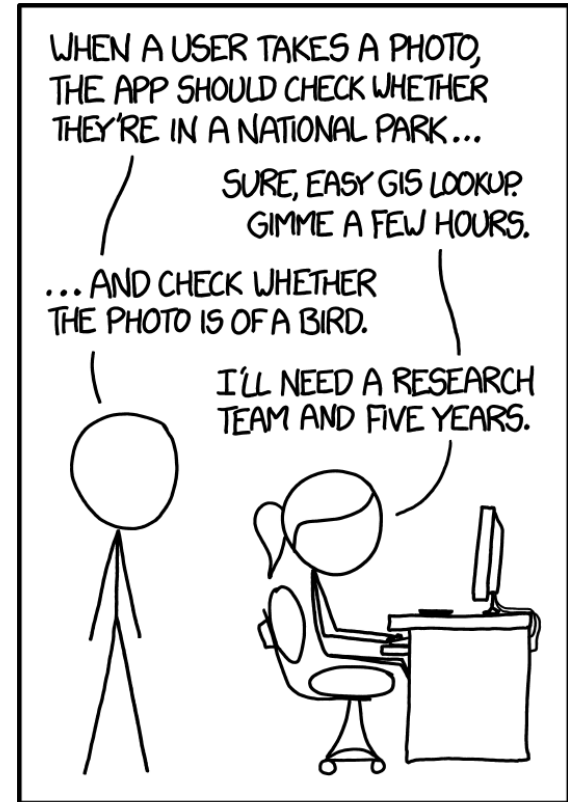April 29th, 2022

# Overview

- Motivation for deep learning

- Areas of Deep Learning

- Convolutional neural networks

- Recurrent neural networks

- Deep learning tools

# Classical Approaches Saturate!

- Computer vision is especially hard for conventional image processing techniques
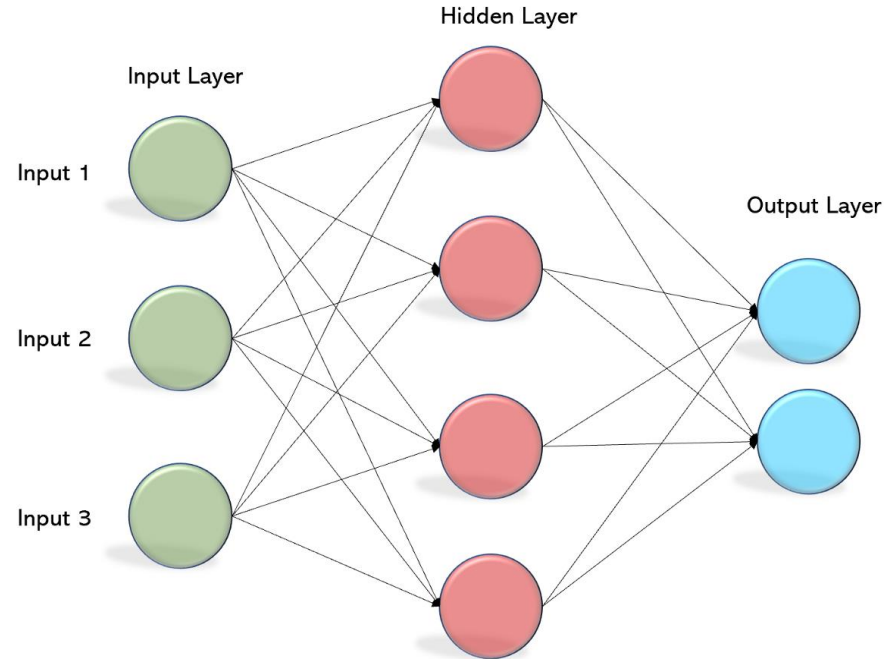- Humans are just intrinsically better at perceiving the world!

https://xkcd.com/1425/



3

# What about the MLPs we learnt in class?
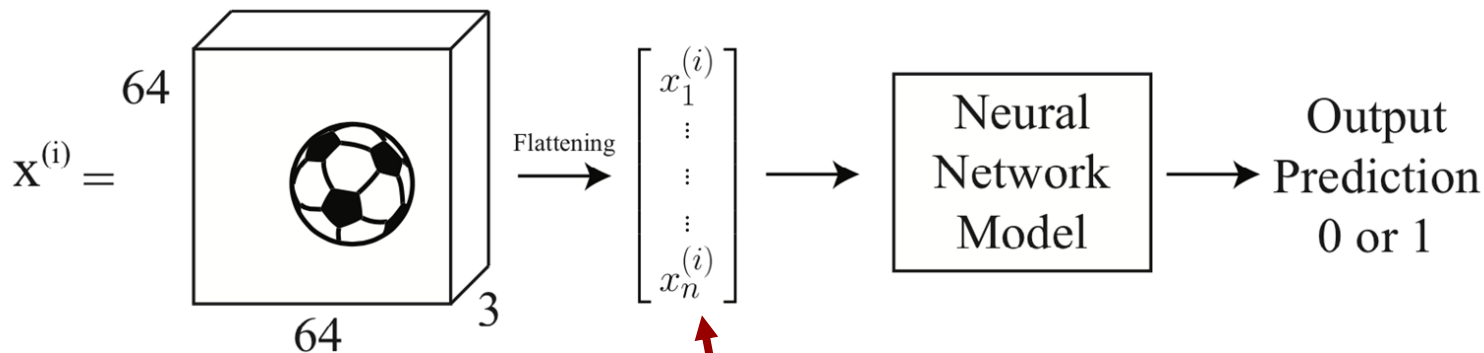
Recall:

- Input Layer
- Hidden layer
- Activations
- Outputs

# What about the MLPs we learnt in class?

Expensive to learn.  Will not generalize well

Does not exploit the order and local relations in the data!

$$\mathbf{x}^{(i)} = \boxed{\begin{array}{c} 64 \\ \\ 64 \end{array}} \xrightarrow{\text{Flattening}} \begin{bmatrix} x_1^{(i)} \\ \vdots \\ \vdots \\ \vdots \\ x_n^{(i)} \end{bmatrix} \longrightarrow \boxed{\begin{array}{c} \text{Neural} \\ \text{Network} \\ \text{Model} \end{array}} \longrightarrow \begin{array}{c} \text{Output} \\ \text{Prediction} \\ \text{0 or 1} \end{array}$$

64x64x3=12288 parameters
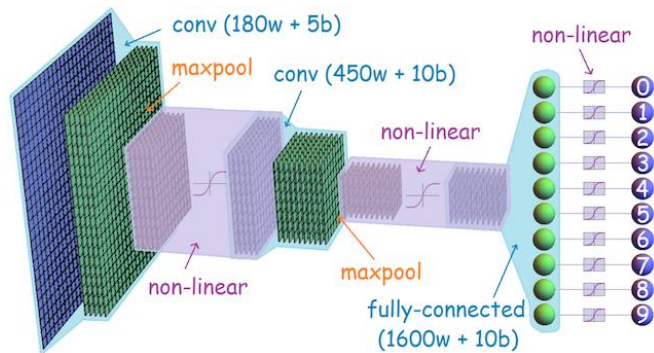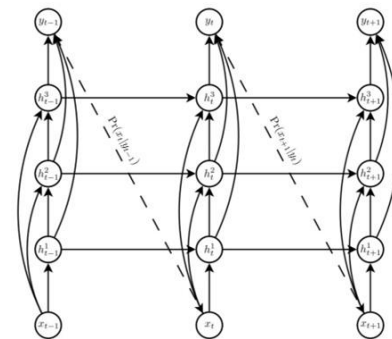We also want **many** layers

# Overview

- Motivation for deep learning

- **Areas in Deep Learning**

- Convolutional neural networks

- Recurrent neural networks

- Deep learning tools

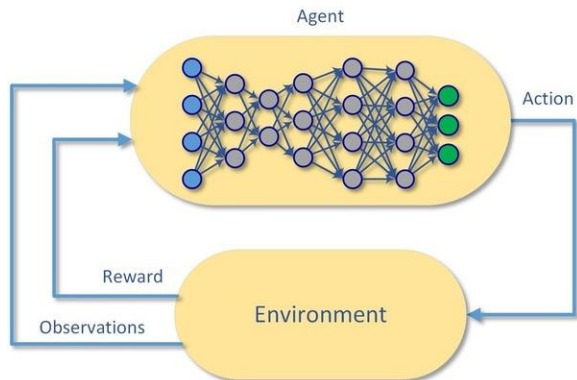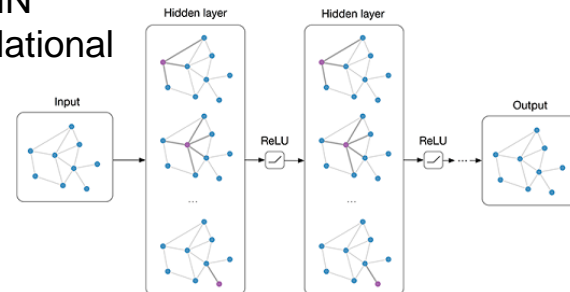# What are different pillars of deep learning?



Convolutional NN
Image

Recurrent NN
Time Series

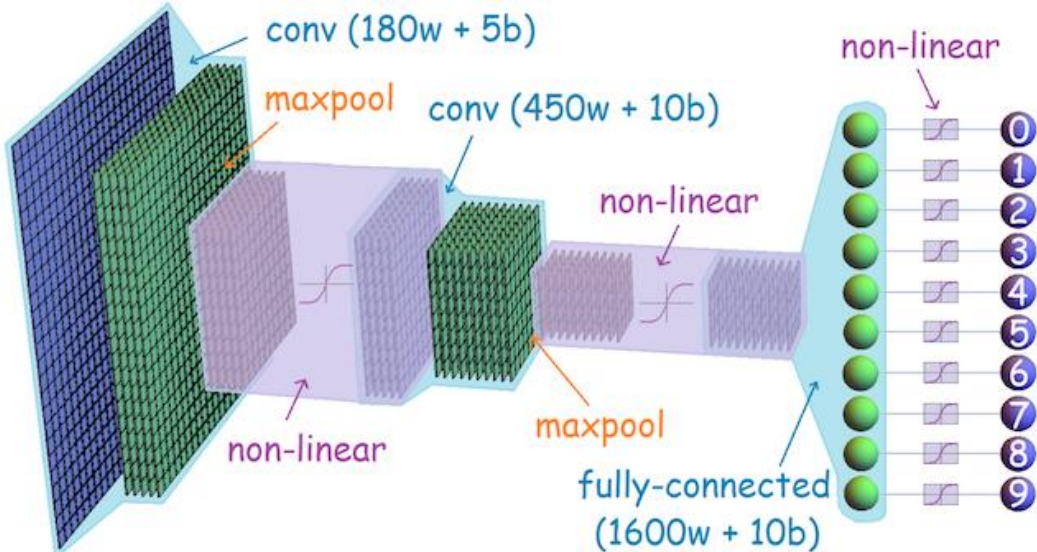Deep RL
Control System

Graph NN
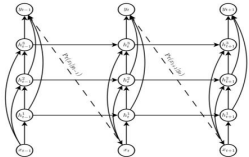Networks/Relational

# Overview

- Motivation for deep learning

- Areas of Deep Learning

- **Convolutional neural networks**

- Recurrent neural networks

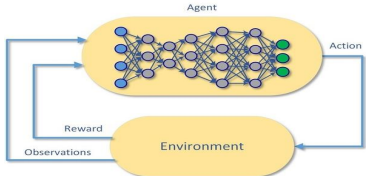- Deep learning tools

# Convolutional Neural Networks
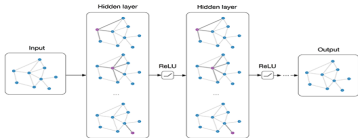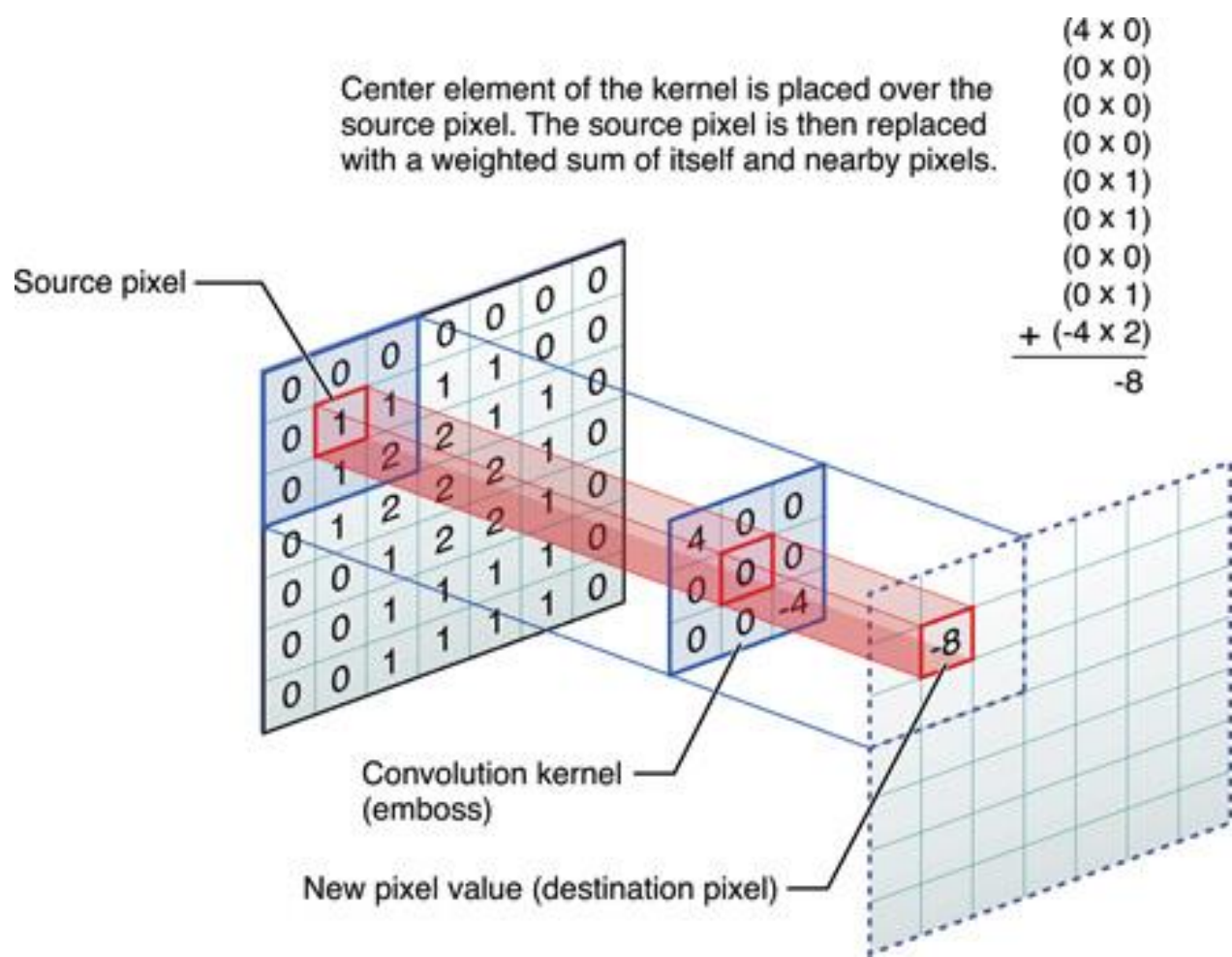


Convolutional Neural Network

conv (180w + 5b)
maxpool
conv (450w + 10b)
non-linear
non-linear
non-linear
maxpool
fully-connected (1600w + 10b)

Recurrent NN

Deep RL

Graph NN

# Let us look at images in detail



28 x 28
784 pixels

# 2D Convolution

Center element of the kernel is placed over the source pixel. The source pixel is then replaced with a weighted sum of itself and nearby pixels.

$$(4 \times 0)$$
$$(0 \times 0)$$
$$(0 \times 0)$$
$$(0 \times 0)$$
$$(0 \times 1)$$
$$(0 \times 1)$$
$$(0 \times 0)$$
$$(0 \times 1)$$
$$+ (-4 \times 2)$$
$$-8$$

Source pixel

Convolution kernel (emboss)

New pixel value (destination pixel)

# Convolving Filters

No change:



Original

Filtered
(no change)

Shifted right by one pixel:



Original

Shifted right
By 1 pixel

Blurred (you already saw this above):



$\frac{1}{9}$

Original

Blur (with a
box filter)

Note the edge artifact.*

Sharpening



Original

$-\frac{1}{9}$

$=$

Edge Detection: Laplacian Filters

| 0 | -1 | 0 |
|---|---|---|
| -1 | 4 | -1 |
| 0 | -1 | 0 |

| -1 | -1 | -1 |
|---|---|---|
| -1 | 8 | -1 |
| -1 | -1 | -1 |

# Convolving Filters

- Why not extract features using filters?

- Better, why not let the data dictate what filters to use?

- Learnable filters!!



Input



| $1_{\times1}$ | $1_{\times0}$ | $1_{\times1}$ | 0 | 0 |
| $0_{\times0}$ | $1_{\times1}$ | $1_{\times0}$ | 1 | 0 |
| $0_{\times1}$ | $0_{\times0}$ | $1_{\times1}$ | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Image

| 4 | | |
| | | |
| | | |

Convolved Feature

13

# Convolution on multiple channels

- Images are generally RGB !!

- How would a filter work on a image with RGB channels?

- The filter should also have 3 channels.

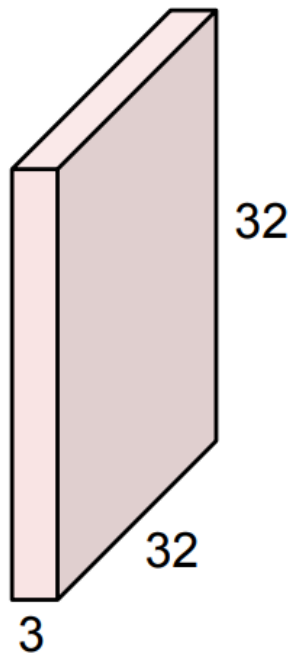- Now the output has a channel for every filter we have used.

# Convolution Layer
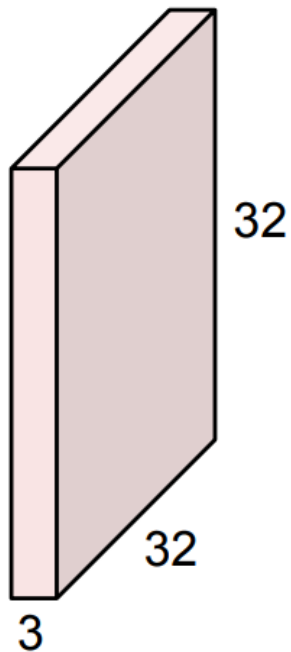
32x32x3 image -> preserve spatial structure



32 height

32 width

3 depth

# Convolution Layer

32x32x3 image

5x5x3 filter

32

32

3

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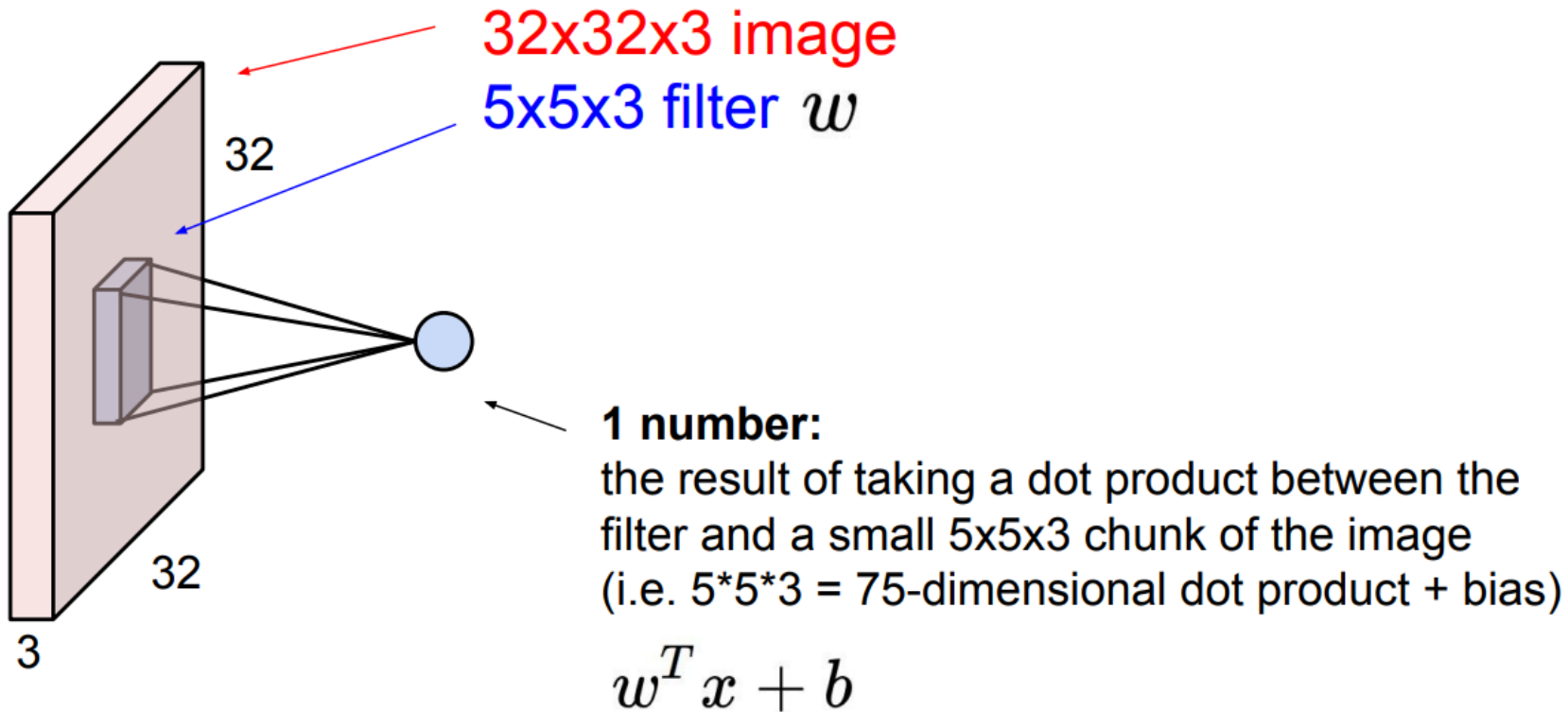computing dot products"

# Convolution Layer

Filters always extend the full depth of the input volume

32x32x3 image

5x5x3 filter
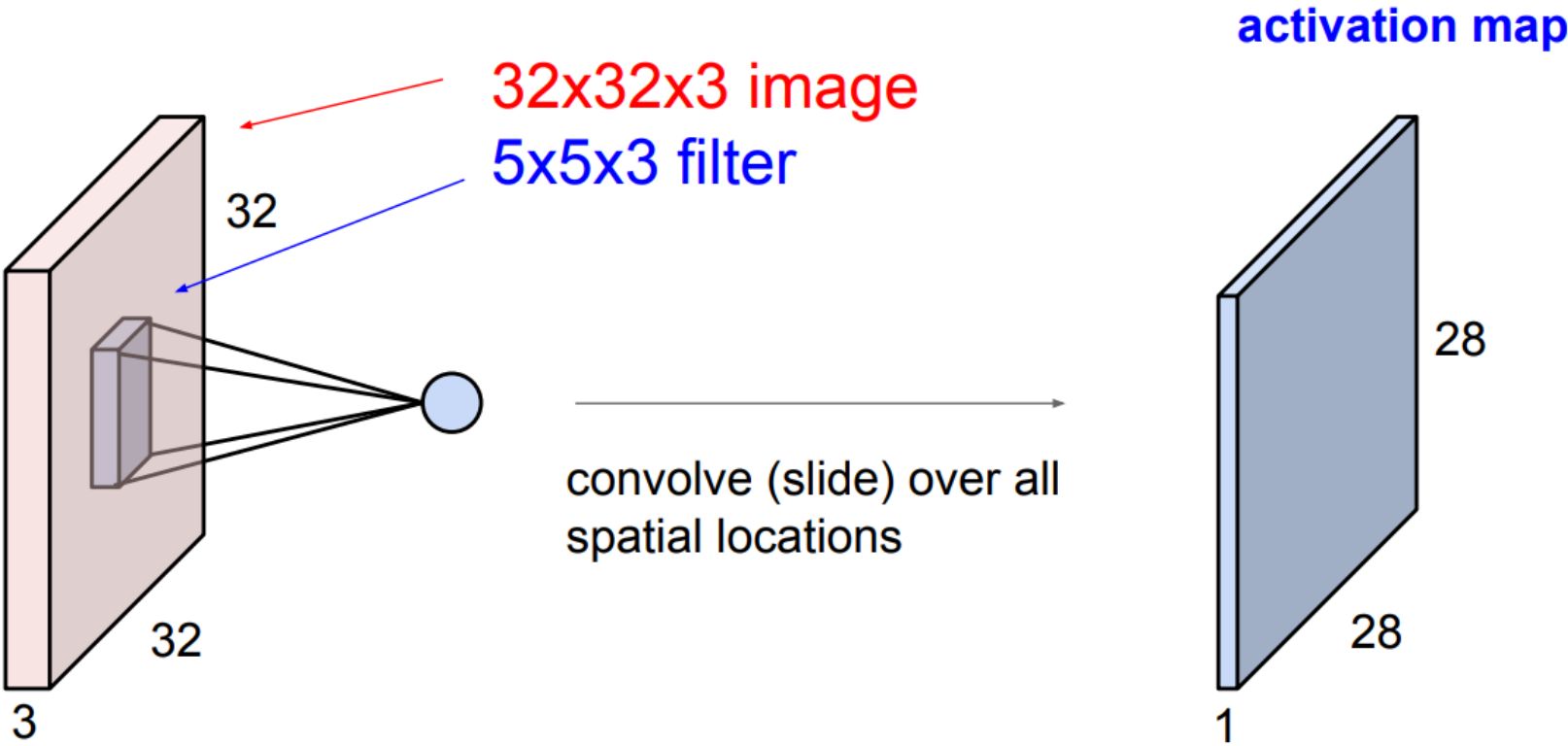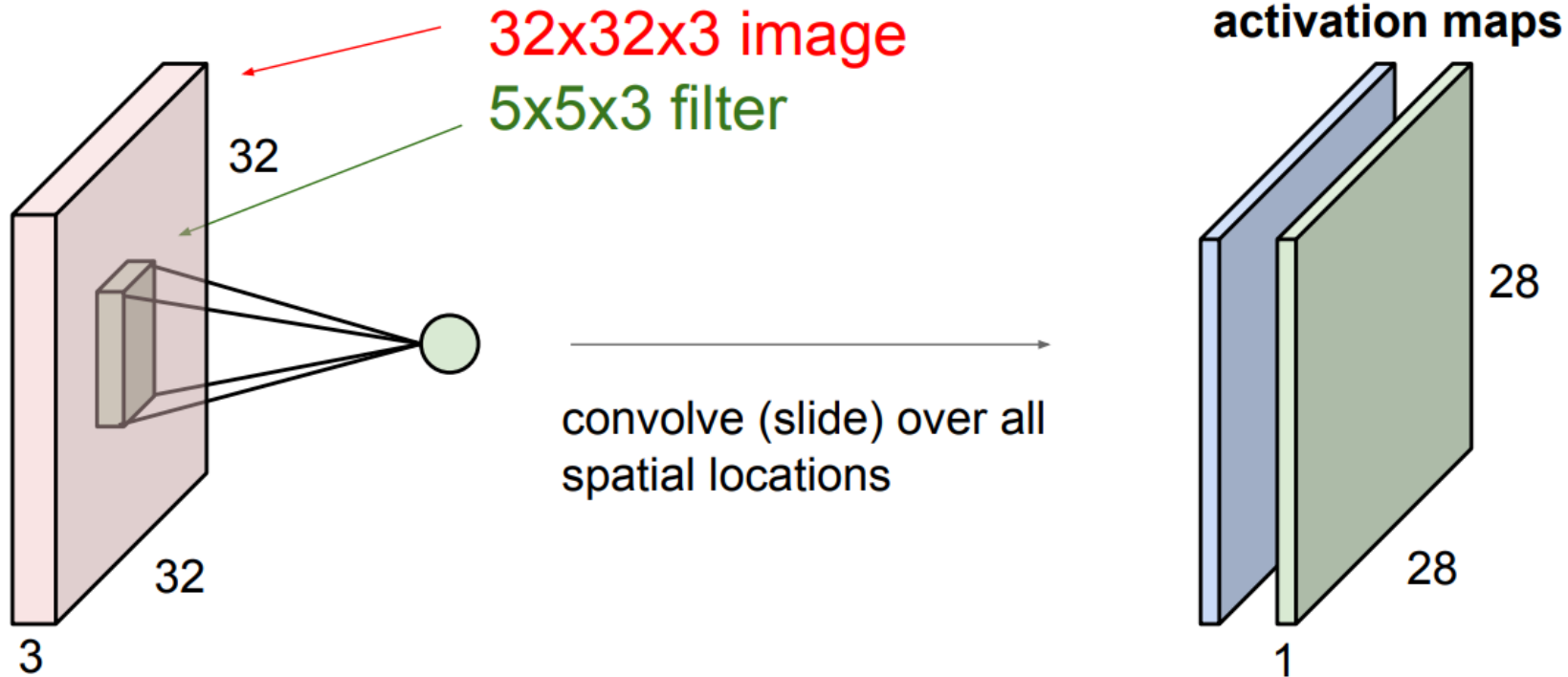
32

32

3

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"
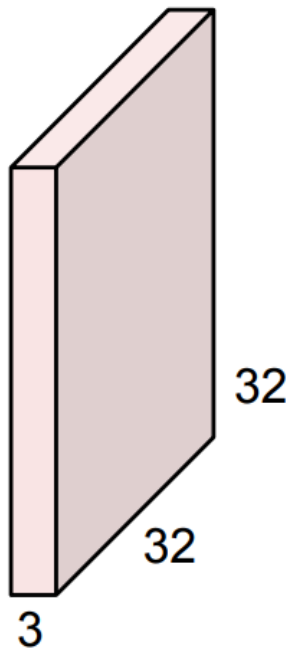
# Convolution Layer



32x32x3 image

5x5x3 filter $w$

**1 number:**
the result of taking a dot product between the filter and a small 5x5x3 chunk of the image (i.e. 5*5*3 = 75-dimensional dot product + bias)

$$w^T x + b$$

# Convolution Layer



32x32x3 image
5x5x3 filter

activation map

convolve (slide) over all spatial locations

# Convolution Layer

consider a second, green filter

**32x32x3 image**
**5x5x3 filter**

32

32

3

convolve (slide) over all spatial locations

**activation maps**

28

28

1

# Convolution Layer

**3x32x32 image**

Consider 6 filters,
each 3x5x5

6 activation maps,
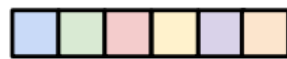each 1x28x28

32

32

3

Convolution
Layer

6x3x5x5
filters

Stack activations to get a
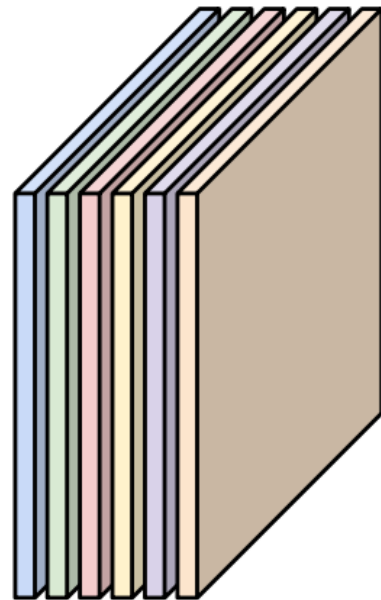6x28x28 output image!

# Convolution Layer

**3x32x32 image**

Also 6-dim bias vector:

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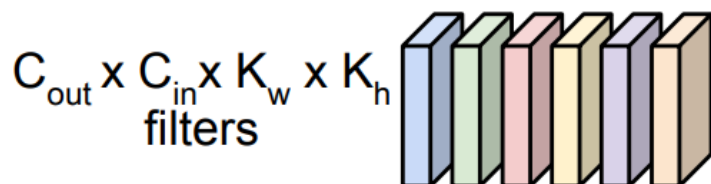
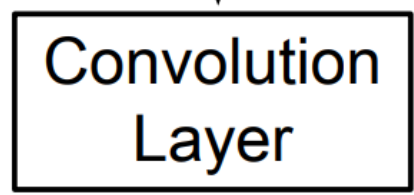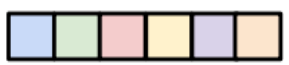Convolution Layer

32

32

3

6x3x5x5 filters

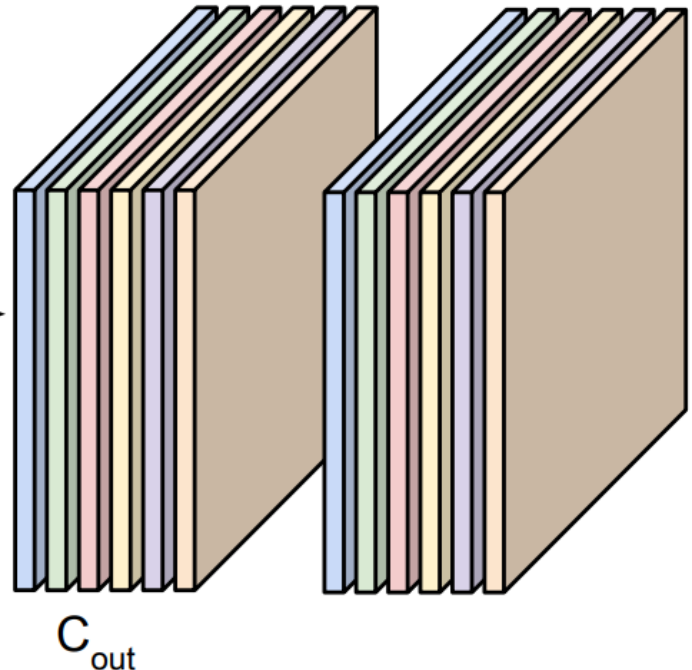Stack activations to get a 6x28x28 output image!

# Convolution Layer

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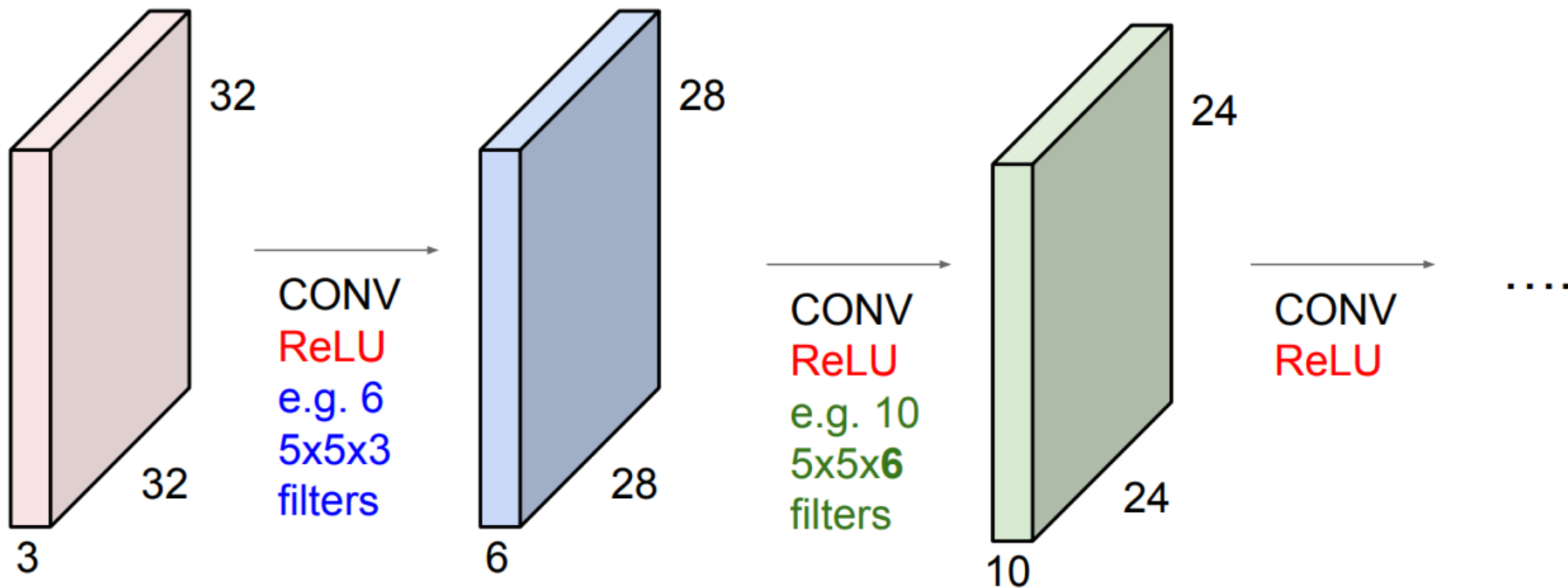Batch of images

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

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



Convolution Layer

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

H

W

$C_{in}$

$C_{out}$

Slide inspiration: Justin Johnson

**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions

# Convolution layer: summary

Let's assume input is $W_1$ x $H_1$ x C
Conv layer needs 4 hyperparameters:
- Number of filters **K**
- The filter size **F**
- The stride **S**
- The zero padding **P**

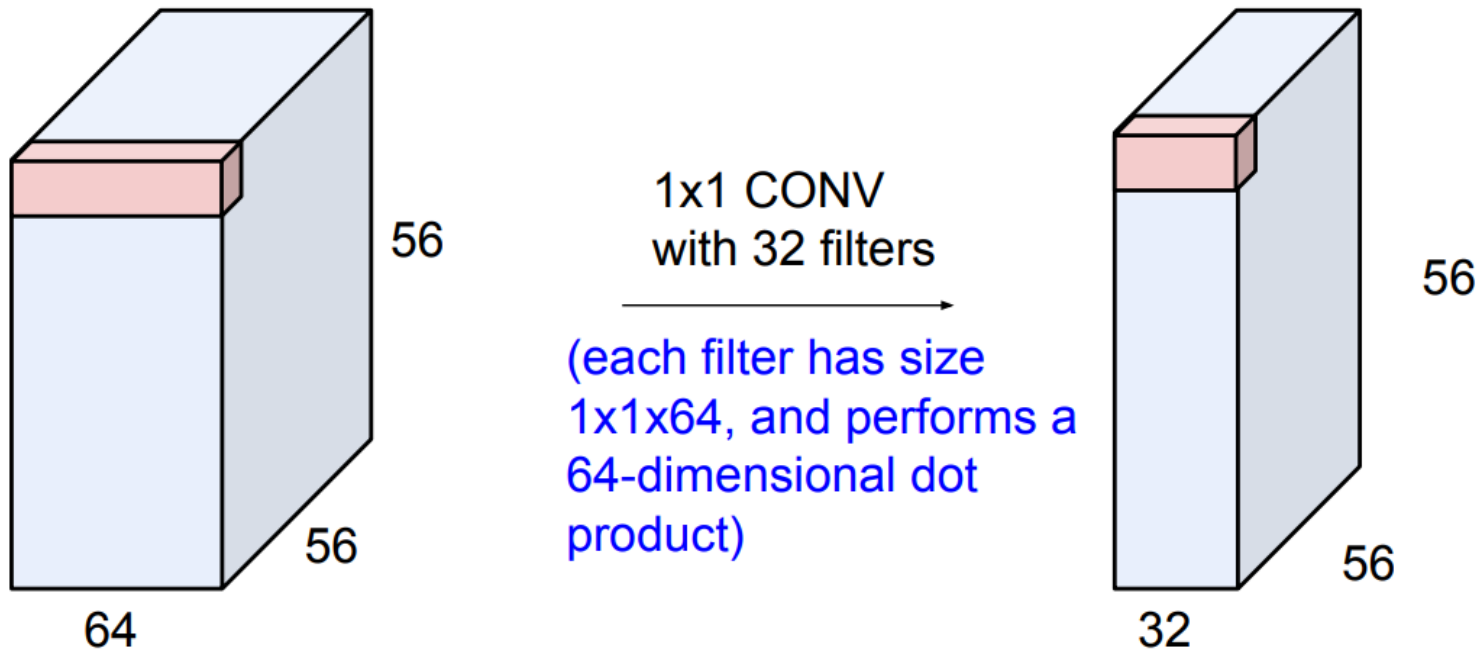This will produce an output of $W_2$ x $H_2$ x K
where:
- $W_2 = (W_1 - F + 2P)/S + 1$
- $H_2 = (H_1 - F + 2P)/S + 1$

Number of parameters: $F^2CK$ and K biases
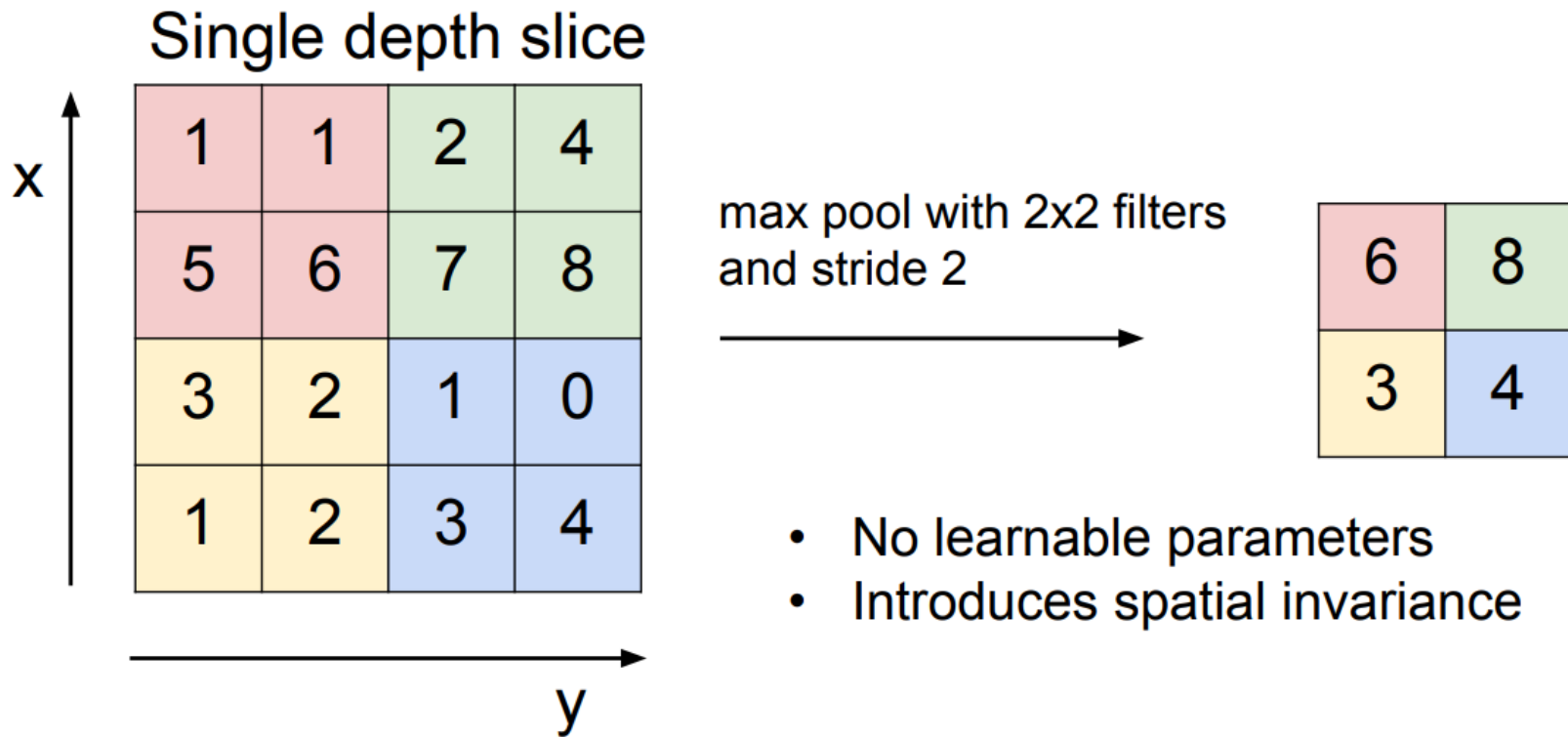
# (btw, 1x1 convolution layers make perfect sense)



1x1 CONV
with 32 filters

(each filter has size
1x1x64, and performs a
64-dimensional dot
product)

56

56

64

56

56

32

# Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently

# MAX POOLING

## Single depth slice

x →

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

y →

max pool with 2x2 filters
and stride 2

→

| 6 | 8 |
|---|---|
| 3 | 4 |

- No learnable parameters
- Introduces spatial invariance

# Pooling layer: summary

Let's assume input is $W_1$ x $H_1$ x C
Conv layer needs 2 hyperparameters:
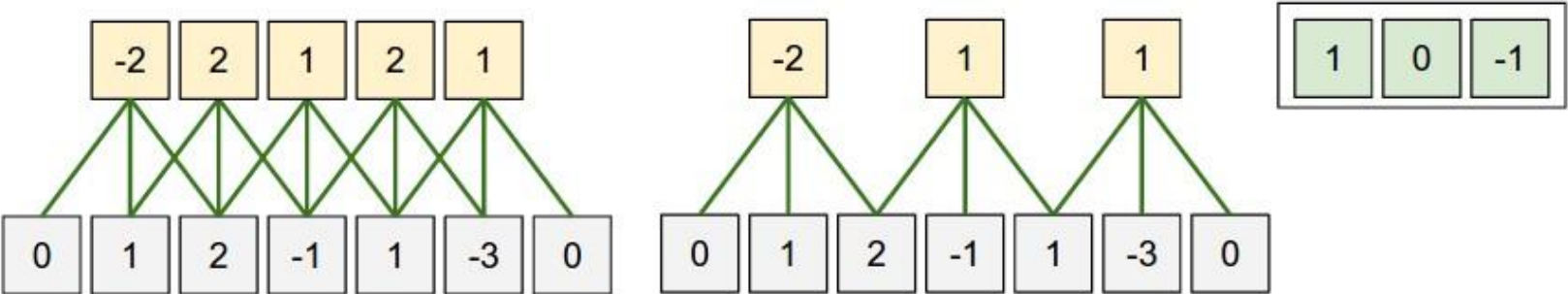- The spatial extent **F**
- The stride **S**

This will produce an output of $W_2$ x $H_2$ x C where:
- $W_2 = (W_1 - F)/S + 1$
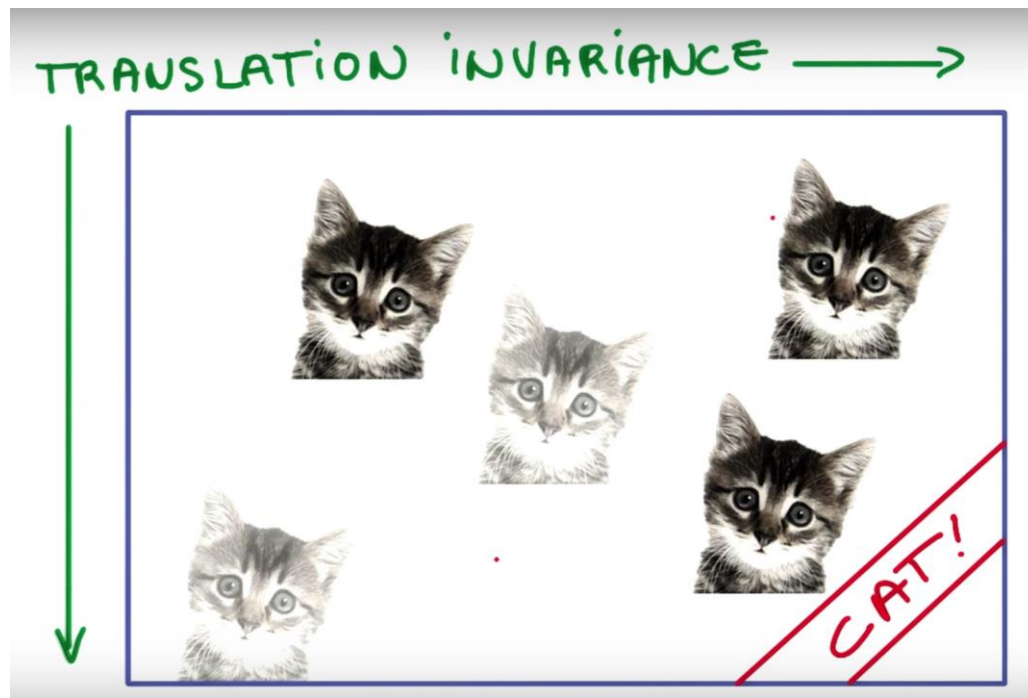- $H_2 = (H_1 - F)/S + 1$

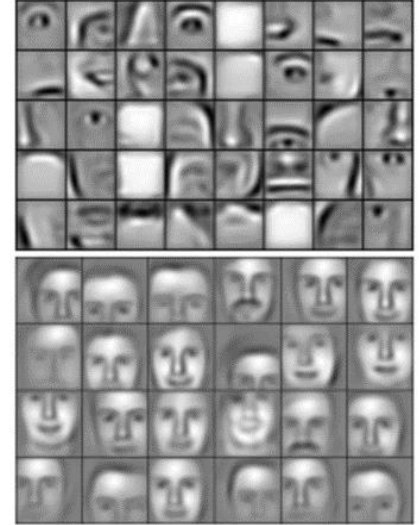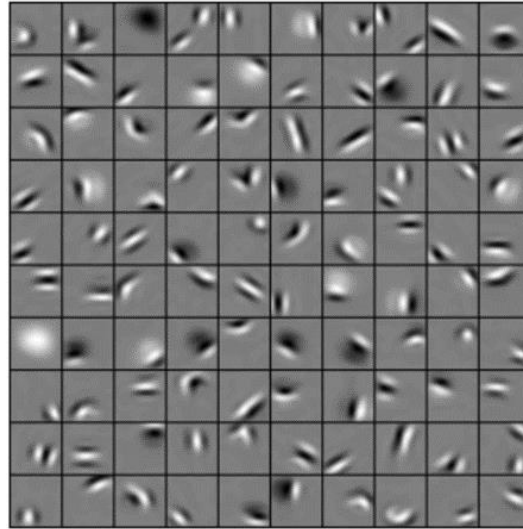Number of parameters: 0
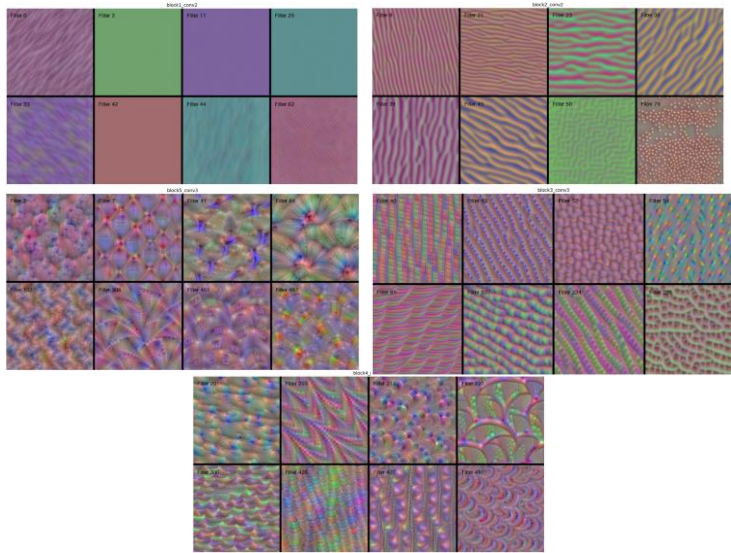
# Parameter Sharing



Lesser the parameters less computationally intensive the training. This is a win win as we are reusing parameters.

# Translational invariance

Since we are training filters to detect cats and the moving these filters over the data, a differently positioned cat will also get detected by the same set of filters.
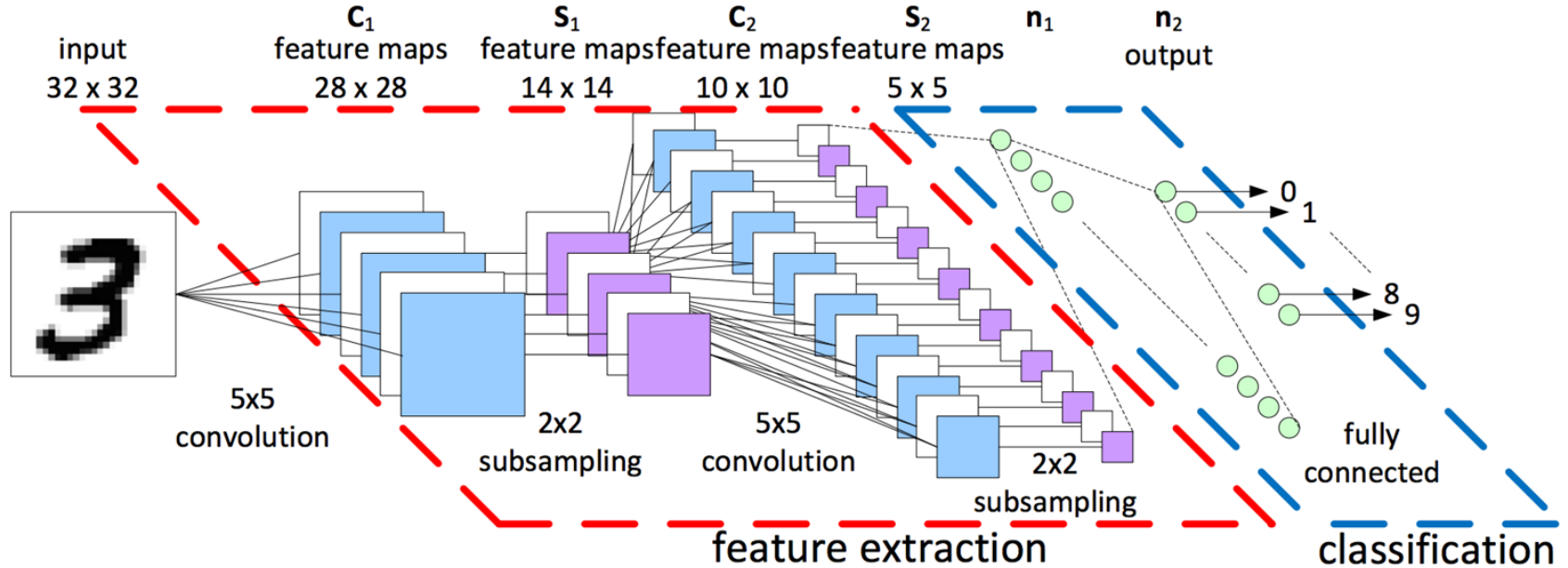
# Filteres? Layers of filters?



Images that maximize filter outputs at certain layers. We observe that the images get more complex as filters are situated deeper

How deeper layers can learn deeper embeddings. How an eye is made up of multiple curves and a face is made up of two eyes.

# How do we use convolutions?



Image credit: LeCun et al. (1998)

Let convolutions extract features!

# Fun Fact: Convolution really is just a linear operation

- In fact convolution is a giant matrix multiplication.

- We can expand the 2 dimensional image into a vector and the conv operation into a matrix.

$$\begin{pmatrix} x1 & x2 & x3 \\ x4 & x5 & x6 \\ x7 & x8 & x9 \end{pmatrix} * \begin{pmatrix} k1 & k2 \\ k3 & k4 \end{pmatrix}$$

$$\begin{pmatrix} k1 & k2 & 0 & k3 & k4 & 0 & 0 & 0 & 0 \\ 0 & k1 & k2 & 0 & k3 & k4 & 0 & 0 & 0 \\ 0 & 0 & 0 & k1 & k2 & 0 & k3 & k4 & 0 \\ 0 & 0 & 0 & 0 & k1 & k2 & 0 & k3 & k4 \end{pmatrix} \cdot \begin{pmatrix} x1 \\ x2 \\ x3 \\ x4 \\ x5 \\ x6 \\ x7 \\ x8 \\ x9 \end{pmatrix}$$

$$\begin{pmatrix} k1\,x1 + k2\,x2 + k3\,x4 + k4\,x5 \\ k1\,x2 + k2\,x3 + k3\,x5 + k4\,x6 \\ k1\,x4 + k2\,x5 + k3\,x7 + k4\,x8 \\ k1\,x5 + k2\,x6 + k3\,x8 + k4\,x9 \end{pmatrix}$$

# How do we learn?

We now have a network with:
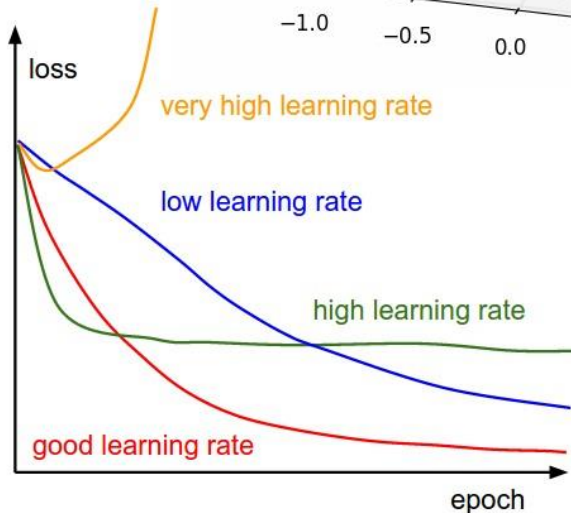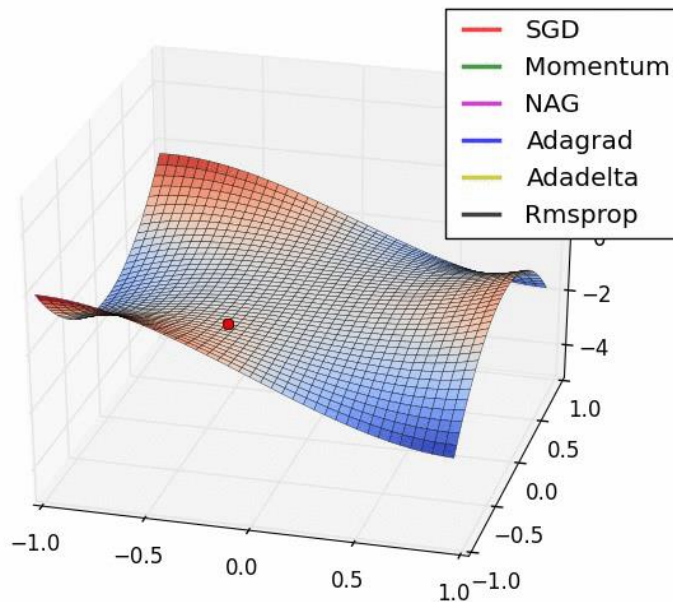- a bunch of weights
- a loss function

To learn:
- Just do gradient descent and backpropagate the error derivates

# How do we learn?

Instead of $\theta := \theta + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x^{(i)}$

There are "optimizers"

- Momentum: Gradient + Momentum
- Nestrov: Momentum + Gradients
- Adagrad: Normalize with sum of sq
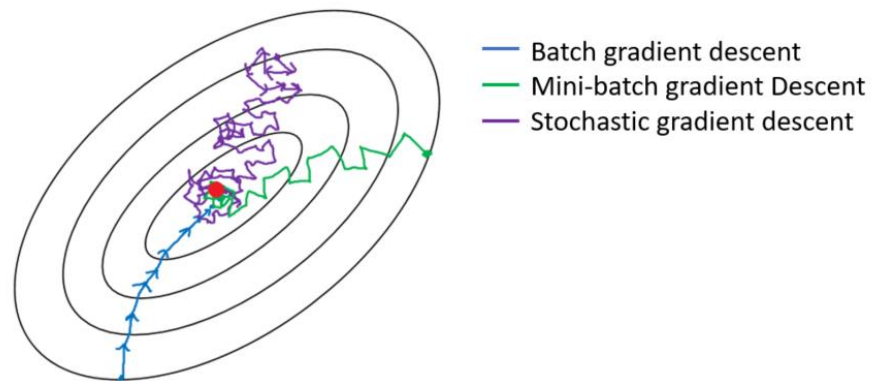- RMSprop: Normalize with moving avg of sum of squares
- ADAM: RMsprop + momentum

# **Mini**-batch Gradient Descent

Expensive to compute gradient for large dataset

      Memory size

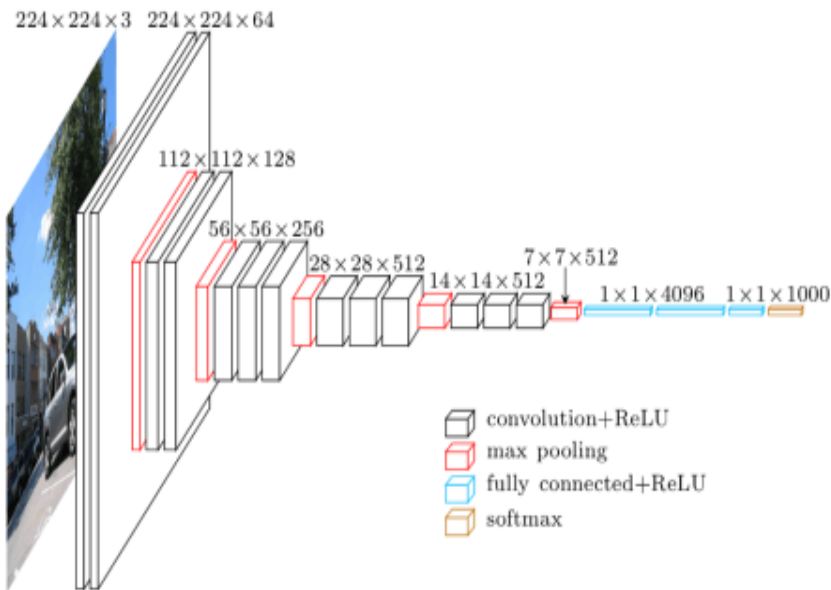      Compute time

Mini-batch: takes a sample of training data

      How to we sample intelligently?



— Batch gradient descent
— Mini-batch gradient Descent
— Stochastic gradient descent

# Is deeper better?

Deeper networks seem to be more powerful but harder to train.

- Loss of information during forward propagation
- Loss of gradient info during back propagation

There are many ways to "keep the gradient going"

# Solution

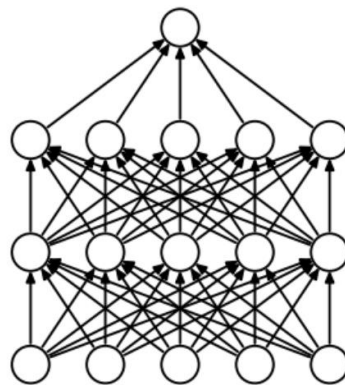Connect the layers, create a gradient highway or information highway.



ResNet (2015)

Image credit: He et al. (2015)

# Initialization

- Can we initialize all neurons to zero?

- If all the weights are same we will not be able to break symmetry of the network and all filters will end up learning the same thing.

- Large numbers, might knock relu units out.

- Relu units once knocked out and their output is zero, their gradient flow also becomes zero.

- We need small random numbers at initialization.

- Variance : 1/sqrt(n)
- Mean: 0

Popular initialization setups

(Xavier, He) (Uniform, Normal)

# Dropout

- What does cutting off some network connections do?

- Trains multiple smaller networks in an ensemble.

- Can drop entire layer too!
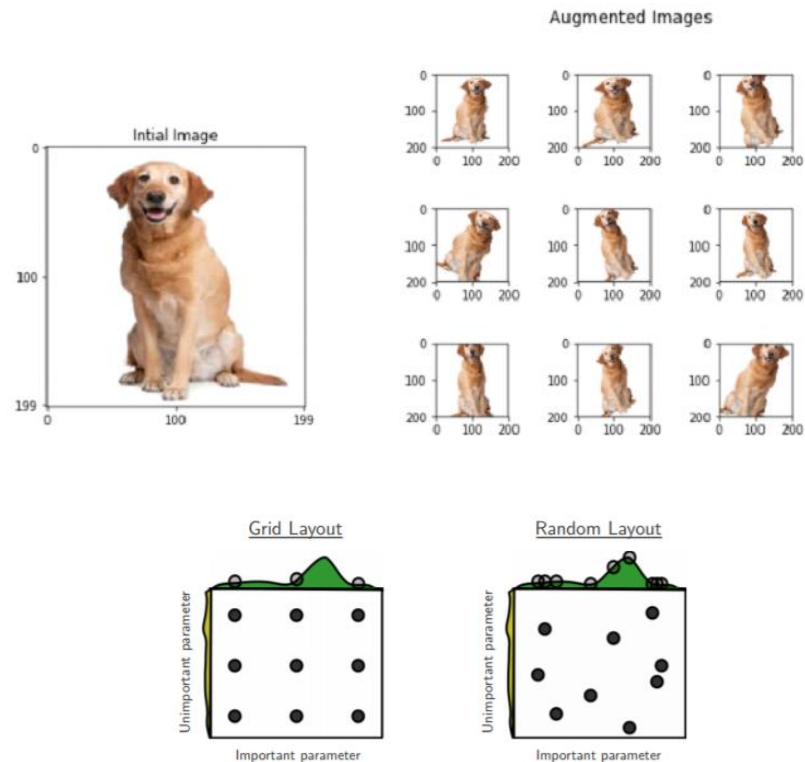
- Acts like a really good regularizer



(a) Standard Neural Net

(b) After applying dropout.

# Tricks for training


Augmented Images

- Data augmentation if your data set is smaller. This helps the network generalize more.

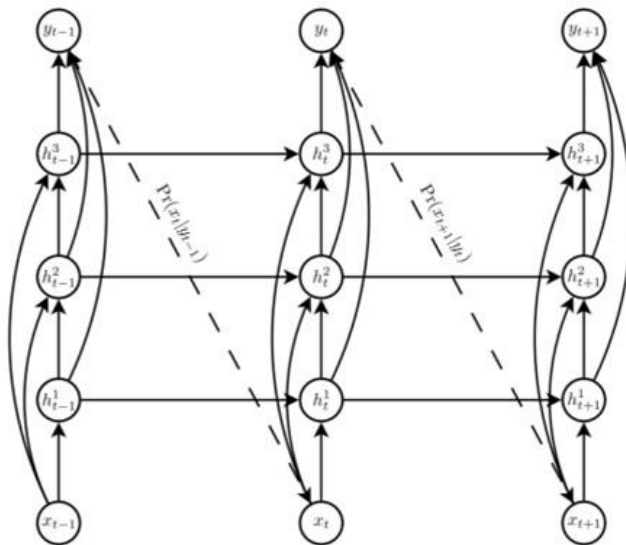- Early stopping if training loss goes above validation loss.

- Random hyperparameter search or grid search?


Grid Layout     Random Layout

# Overview

- Motivation for deep learning

- Areas of Deep Learning

- Convolutional neural networks

- **Recurrent neural networks**
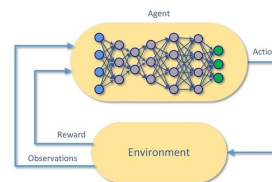
- Deep learning tools

# CNN sounds like fun!
# What are some deep learning pillars?

Recurrent NN
Time Series

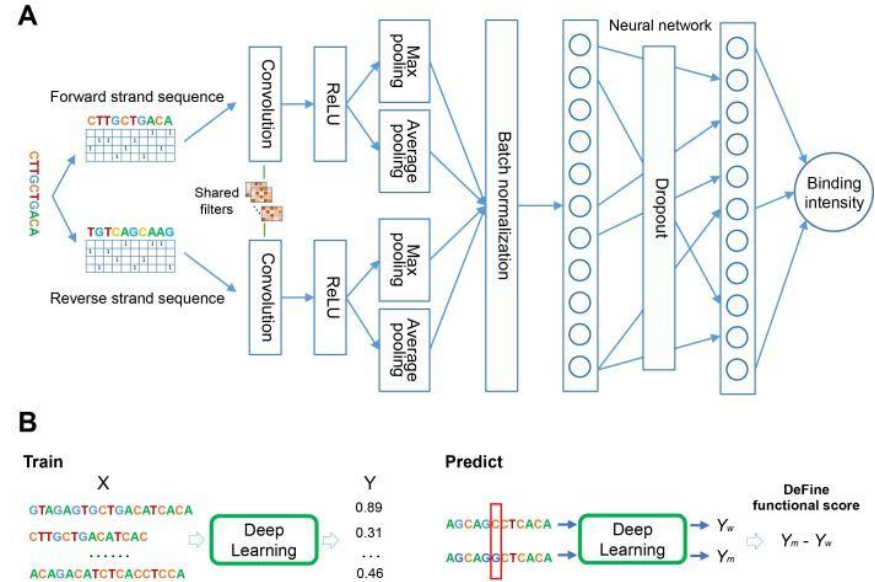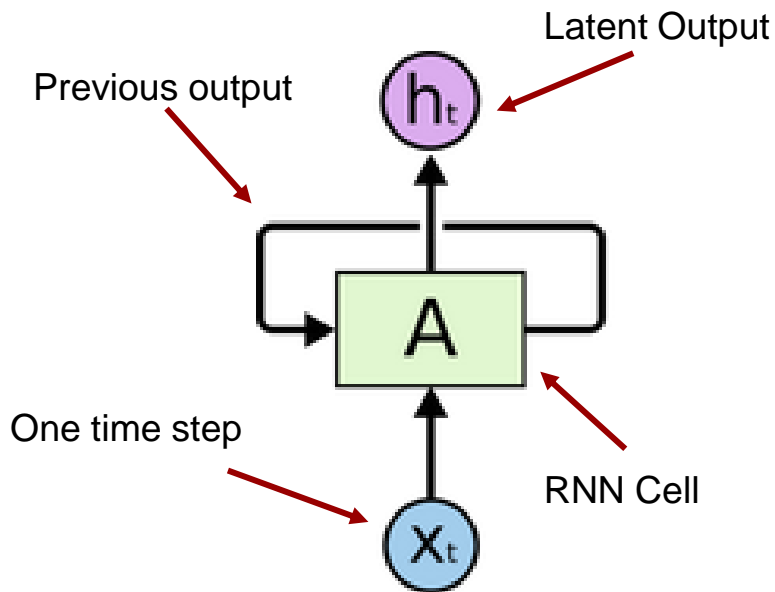$\Pr(x_t|y_{t-1})$

$\Pr(x_{t+1}|y_t)$

Convolutional NN

Deep RL

Graph NN

# We can also have 1D architectures (remember this)

- CNN works on any data where there is a local pattern

- We use 1D convolutions on DNA sequences, text sequences and music notes

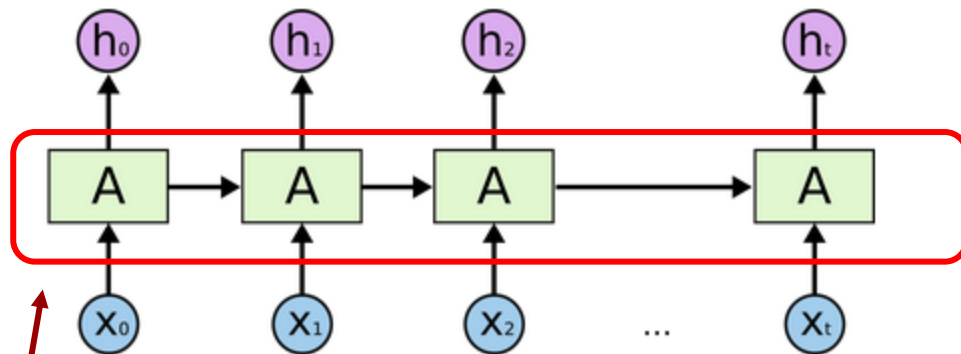- But what if time series has **causal dependency** or any kind of **sequential dependency**?

# To address sequential dependency?
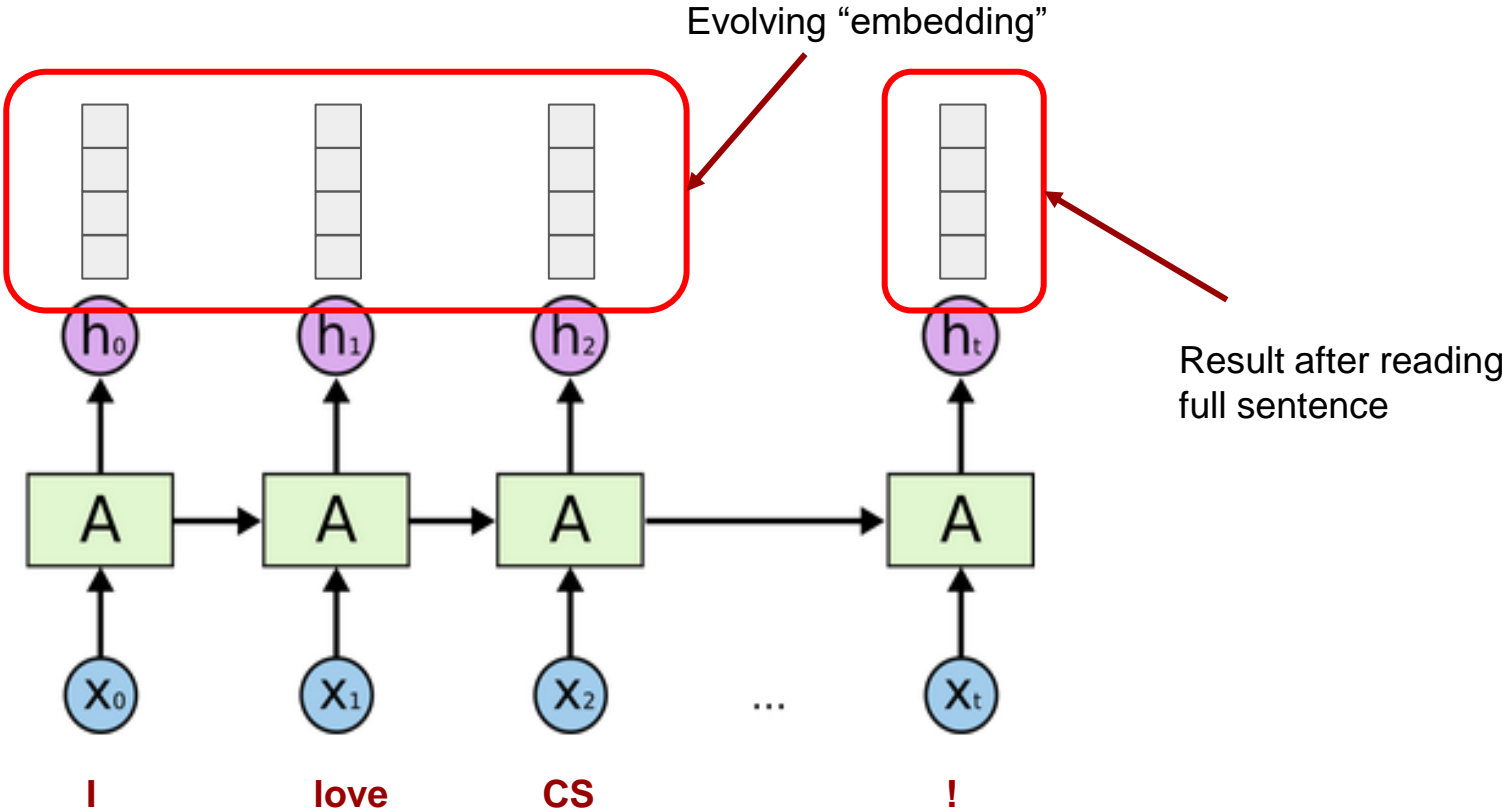
Use recurrent neural network (RNN)

Latent Output

Previous output

Unrolling an RNN

One time step

RNN Cell

They are really the same cell,
NOT many different cells like kernels of CNN

46

# How does RNN produce result?



Evolving "embedding"

Result after reading
full sentence

I          love          CS                    !

# There are 2 types of RNN cells

Store in "long term memory"          Response to current input

Reset gate          Update gate
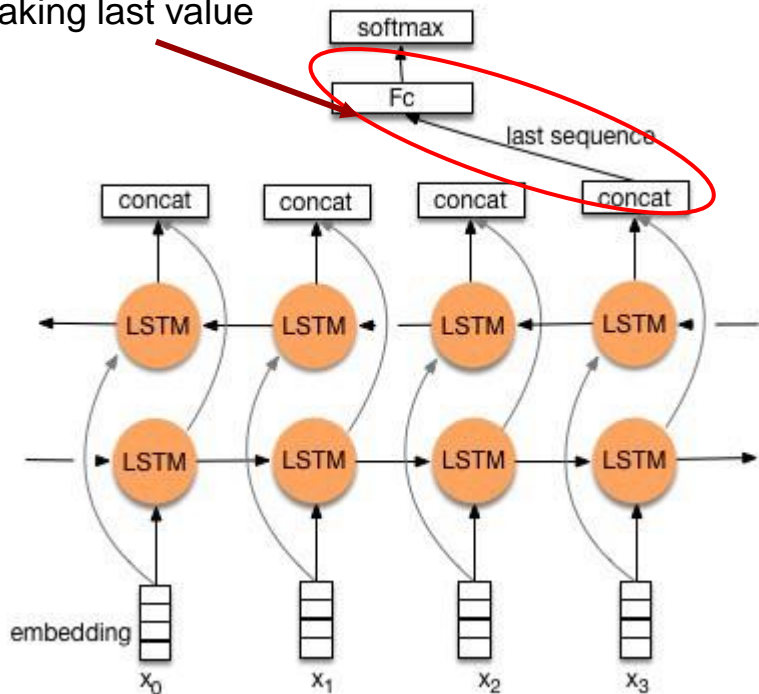


Long Short Term Memory (LSTM)          Gated Recurrent Unit (GRU)

Response to current input

# Recurrent AND deep?

Taking last value
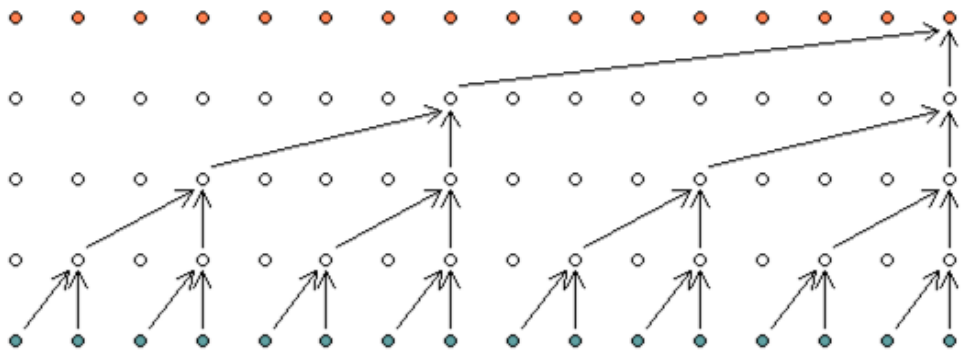
Pay "attention" to everything

Stacking

Attention Model

# "Recurrent" AND convolutional?

Temporal convolutional network

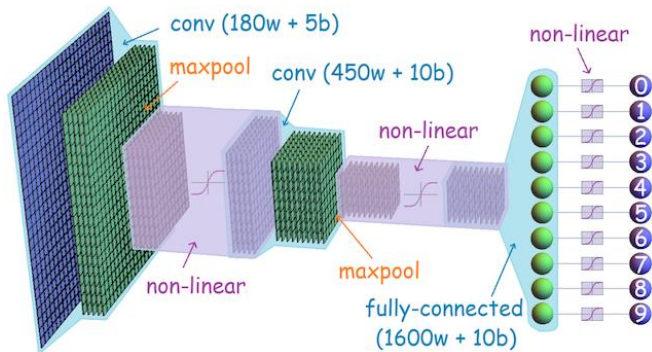Temporal dependency achieved through
"one-sided" convolution

More efficient because deep learning
packages are optimized for matrix
multiplication = convolution
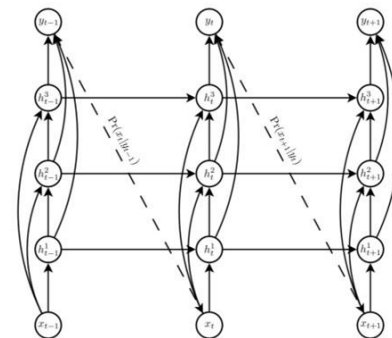
No hard dependency

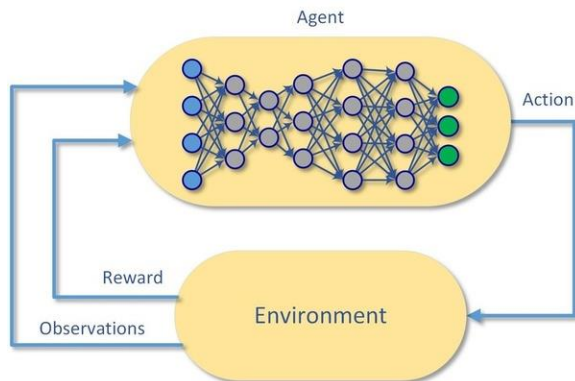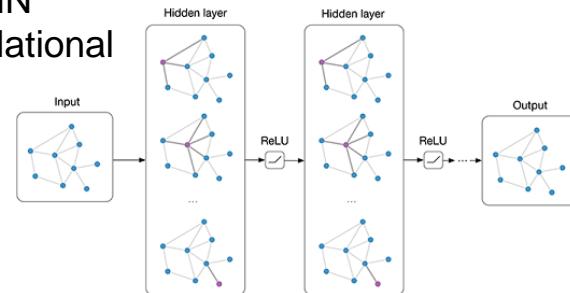# More? Take CS230, CS236, CS231N, CS224N



Convolutional NN
Image

Recurrent NN
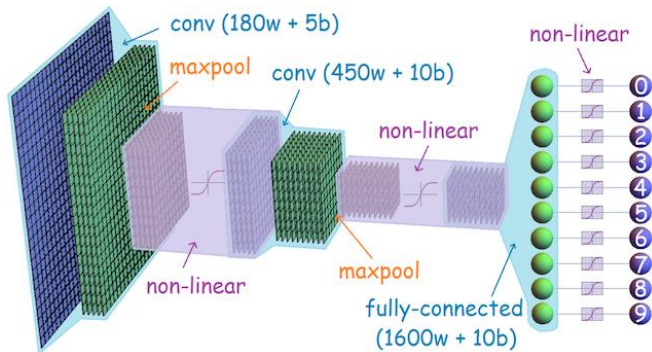Time Series

Deep RL
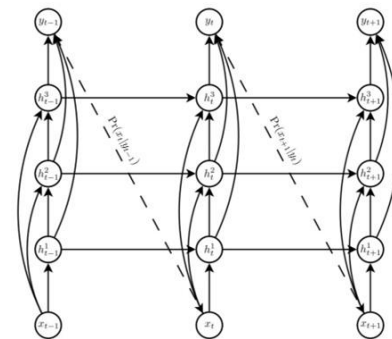Control System

Graph NN
Networks/Relational

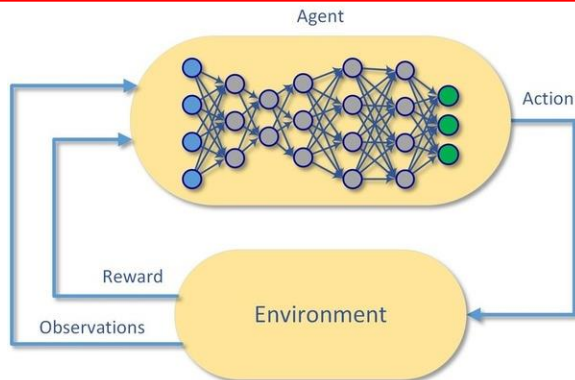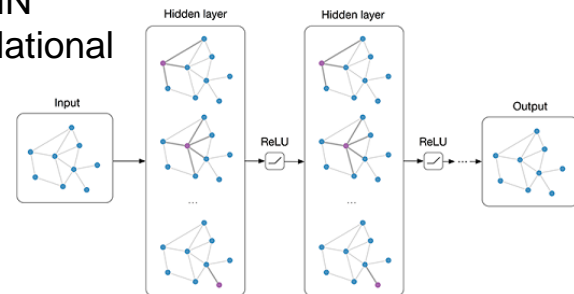# Not today, but take CS234 and CS224W



Convolutional NN
Image

Recurrent NN
Time Series

Deep RL
Control System

Graph NN
Networks/Relational

# Overview

- Motivation for deep learning

- Areas of Deep Learning

- Convolutional neural networks

- Recurrent neural networks

- Deep learning tools

# Tools for deep learning



Keras

TensorFlow

theano

PYTORCH

Popular Tools
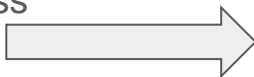
Specialized Groups

Caffe2

mxnet

Microsoft CNTK

# Where can I get free stuff?

Google Colab

      Free (limited-ish) GPU access

      Works nicely with Tensorflow

      Links to Google Drive

Register a new Google Cloud account

      => Instant $300??

      => AWS free tier (limited compute)

      => Azure education account, $200?

Azure Notebook

Kaggle kernel???

Amazon SageMaker?

To **SAVE** money

**CLOSE** your GPU instance

**~$1** an hour

Good luck!
Well, have fun too :D