

# Introduction to Deep Learning

Nandita Bhaskhar

Content adapted from CS231n and past CS229 teams  
April 29<sup>th</sup>, 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/>

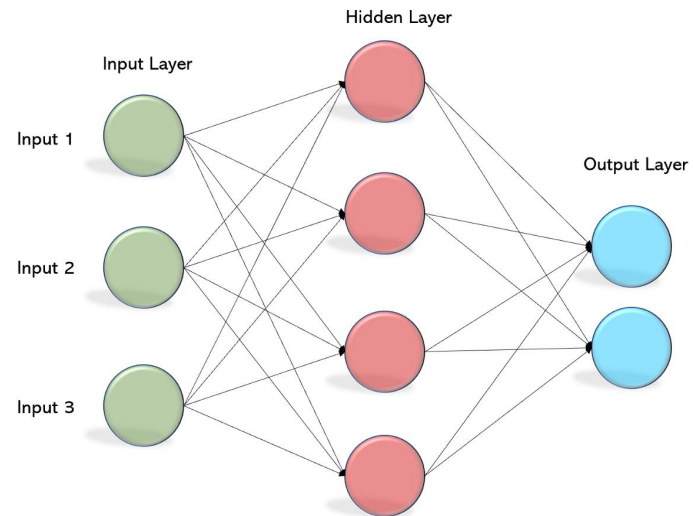


IN CS, IT CAN BE HARD TO EXPLAIN THE DIFFERENCE BETWEEN THE EASY AND THE VIRTUALLY IMPOSSIBLE.

## What about the MLPs we learnt in class?

Recall:

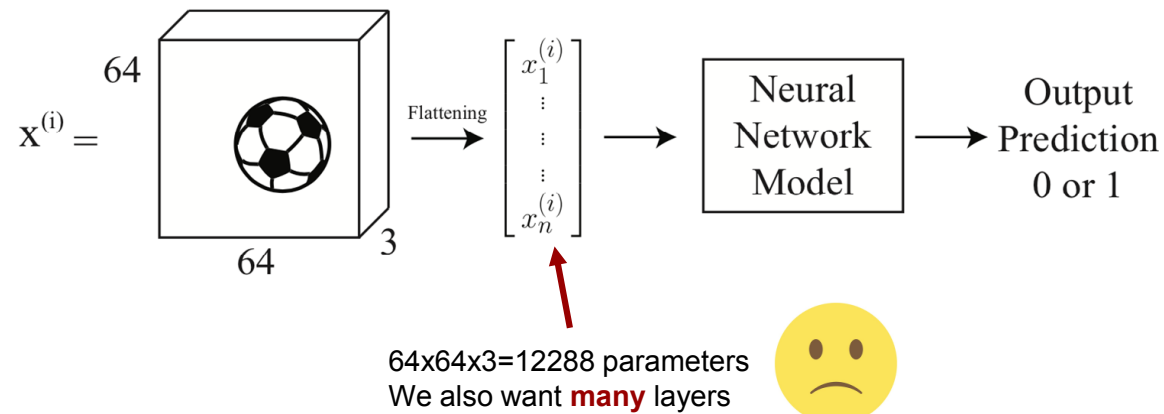
- Input Layer
- Hidden layer
- Activations
- Outputs



## 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!

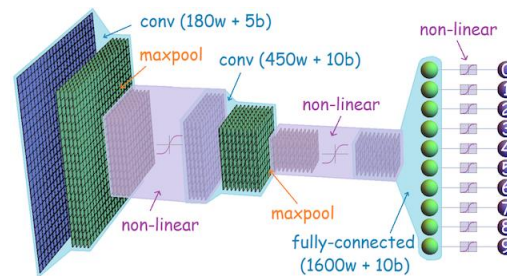


# Overview

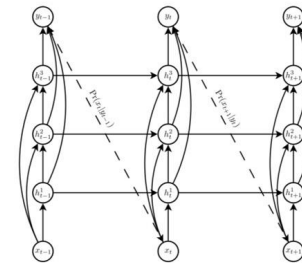
- Motivation for deep learning
- Areas in Deep Learning
- Convolutional neural networks
- Recurrent neural networks
- Deep learning tools

# What are different pillars of deep learning?

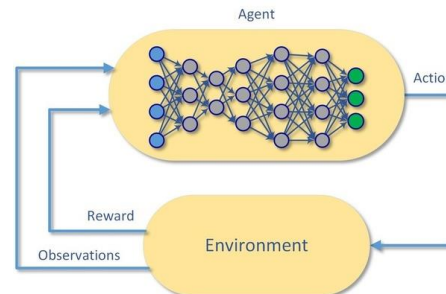
Convolutional NN  
Image



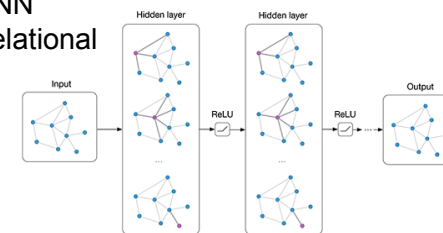
Recurrent NN  
Time Series



Deep RL  
Control System



Graph NN  
Networks/Relational



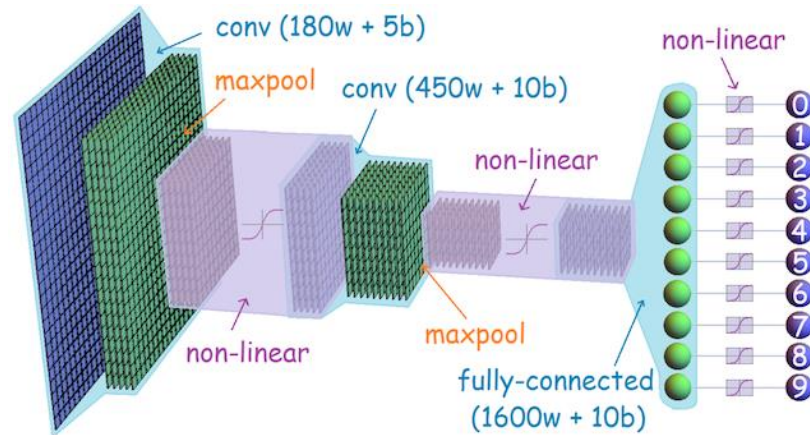
# Overview

- Motivation for deep learning
- Areas of Deep Learning
- Convolutional neural networks
- Recurrent neural networks
- Deep learning tools

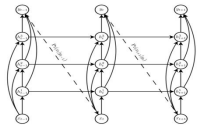


# Convolutional Neural Networks

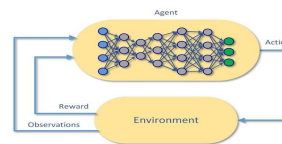
Convolutional Neural Network



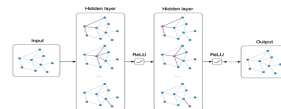
Recurrent NN



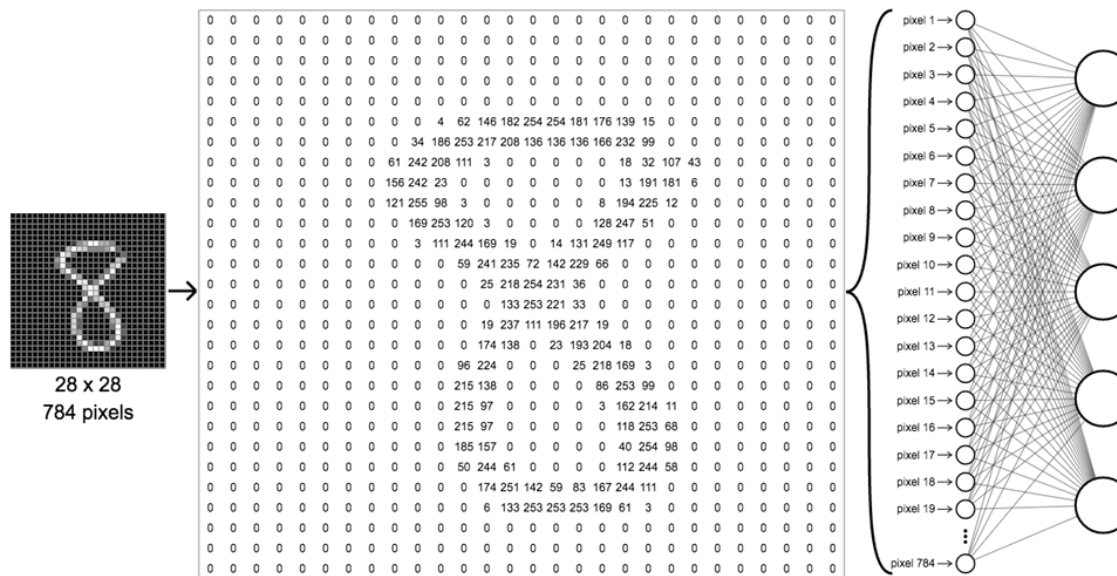
Deep RL



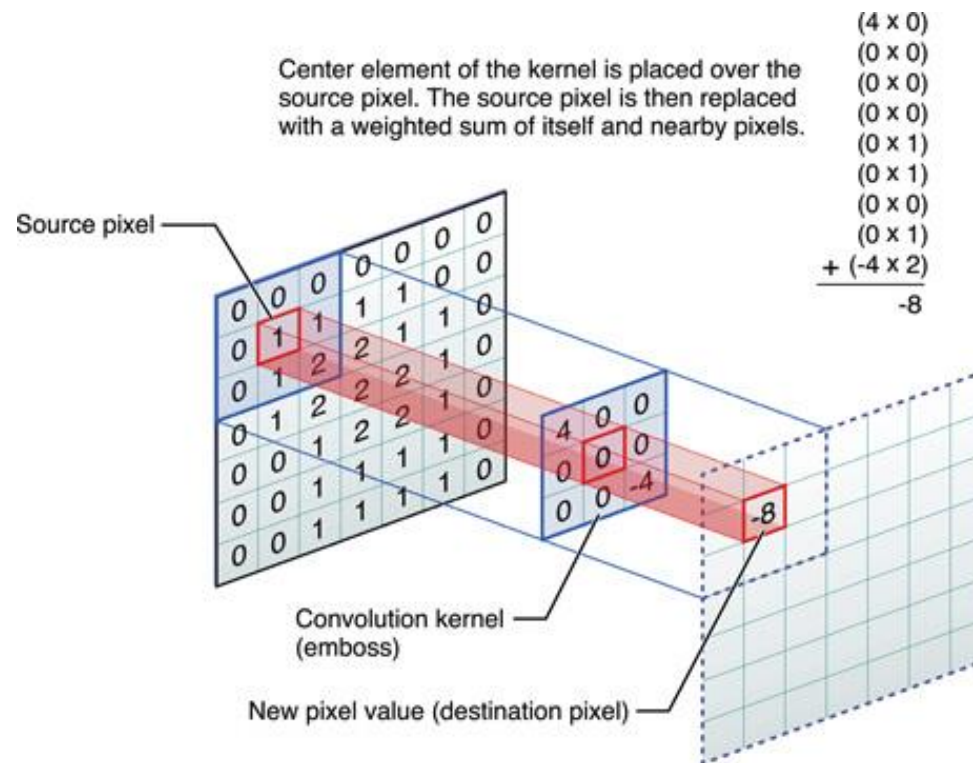
Graph NN



## Let us look at images in detail

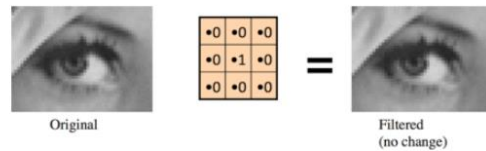


# 2D Convolution

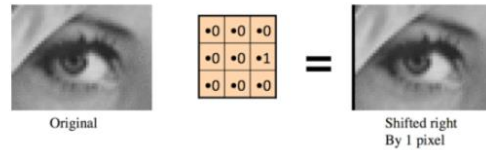


# Convolving Filters

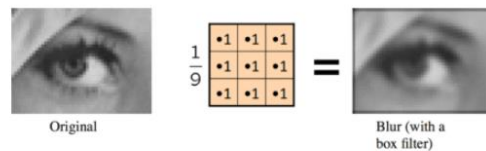
No change:



Shifted right by one pixel:



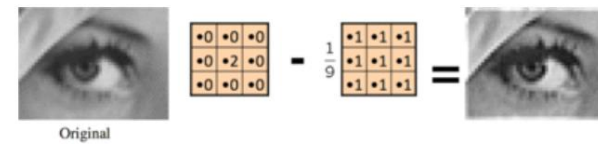
Blurred (you already saw this above):



Note the edge artifact.\*

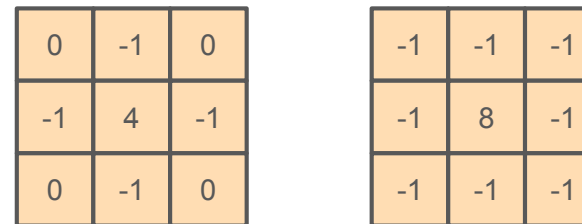
<https://ai.stanford.edu/~syYeung/cvweb/tutorials.html>

Sharpening



<https://ai.stanford.edu/~syYeung/cvweb/tutorials.html>

Edge Detection: Laplacian Filters



# Convolving Filters

- Why not extract features using filters?
- Better, why not let the data dictate what filters to use?
- Learnable filters!!



1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0	0
0 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	1	0
0 <sub>x1</sub>	0 <sub>x0</sub>	1 <sub>x1</sub>	1	1
0	0	1	1	0
0	1	1	0	0

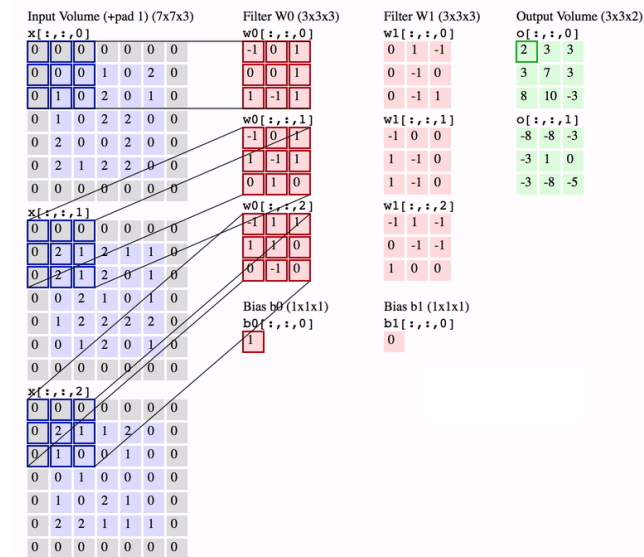
Image

4		

Convolved Feature

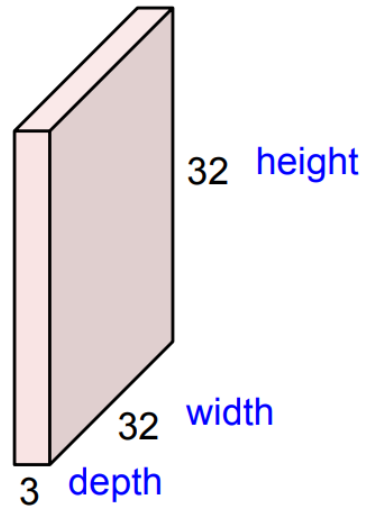
## Convolution on multiple channels

- 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.



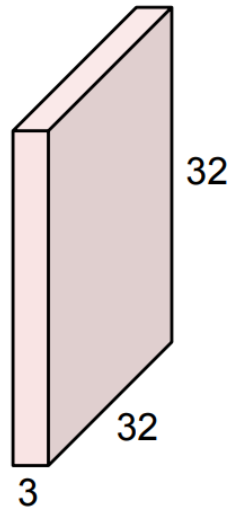
# Convolution Layer

32x32x3 image -> preserve spatial structure



# Convolution Layer

32x32x3 image



5x5x3 filter

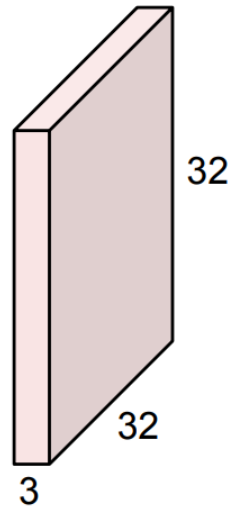


**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”



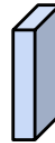
# Convolution Layer

32x32x3 image



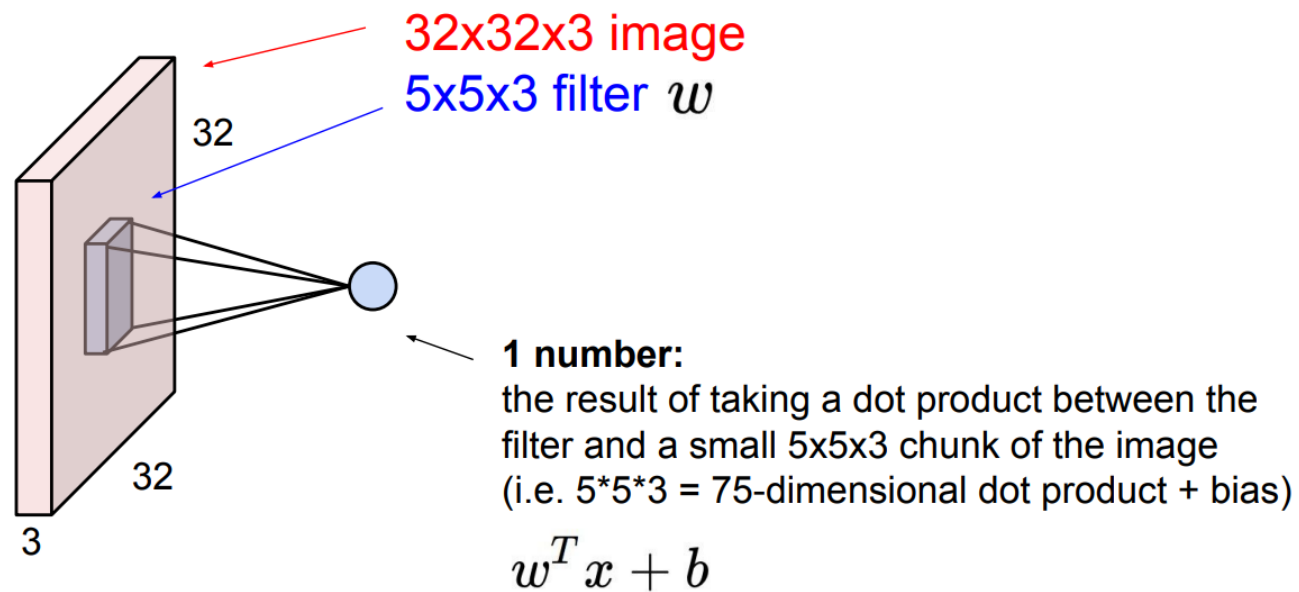
Filters always extend the full depth of the input volume

5x5x3 filter

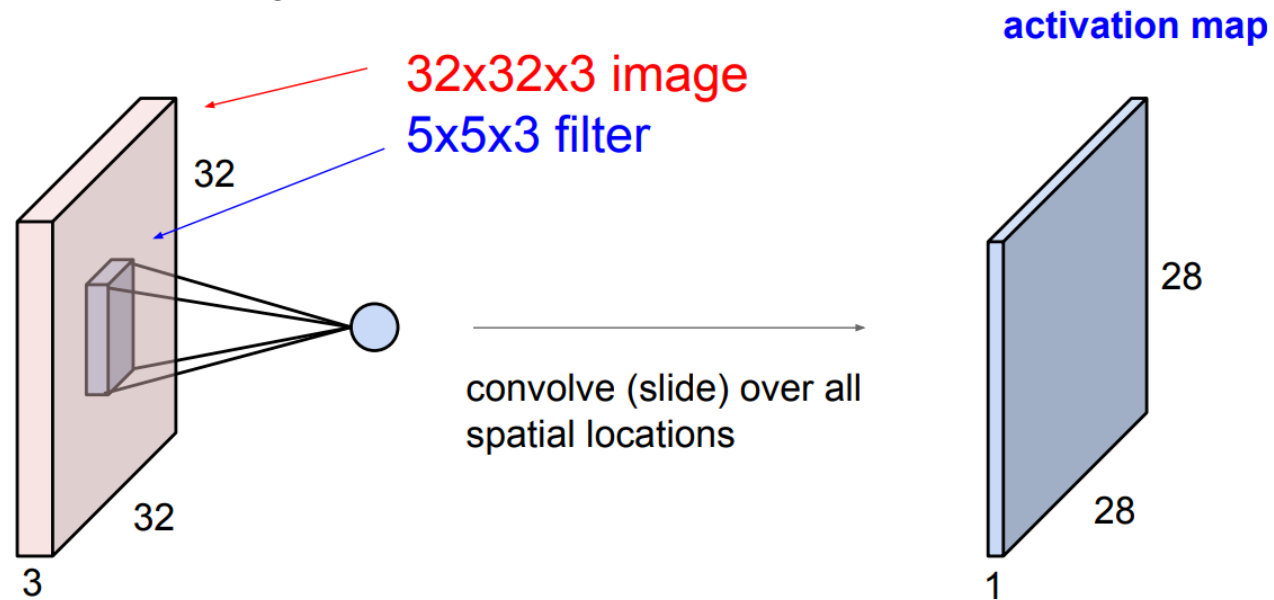


**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolution Layer

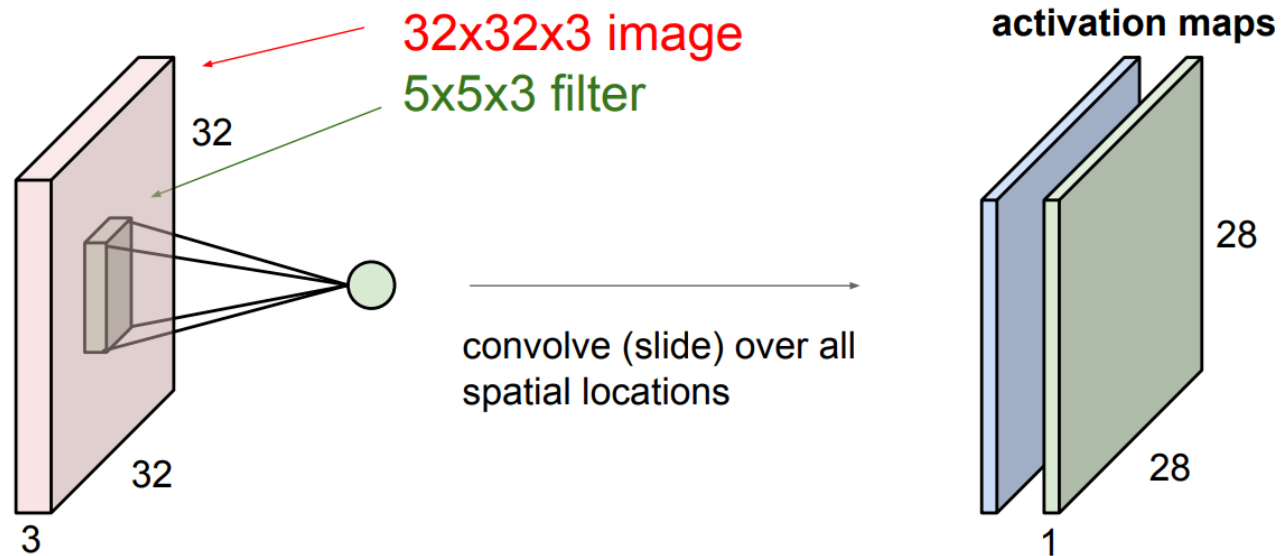


# Convolution Layer

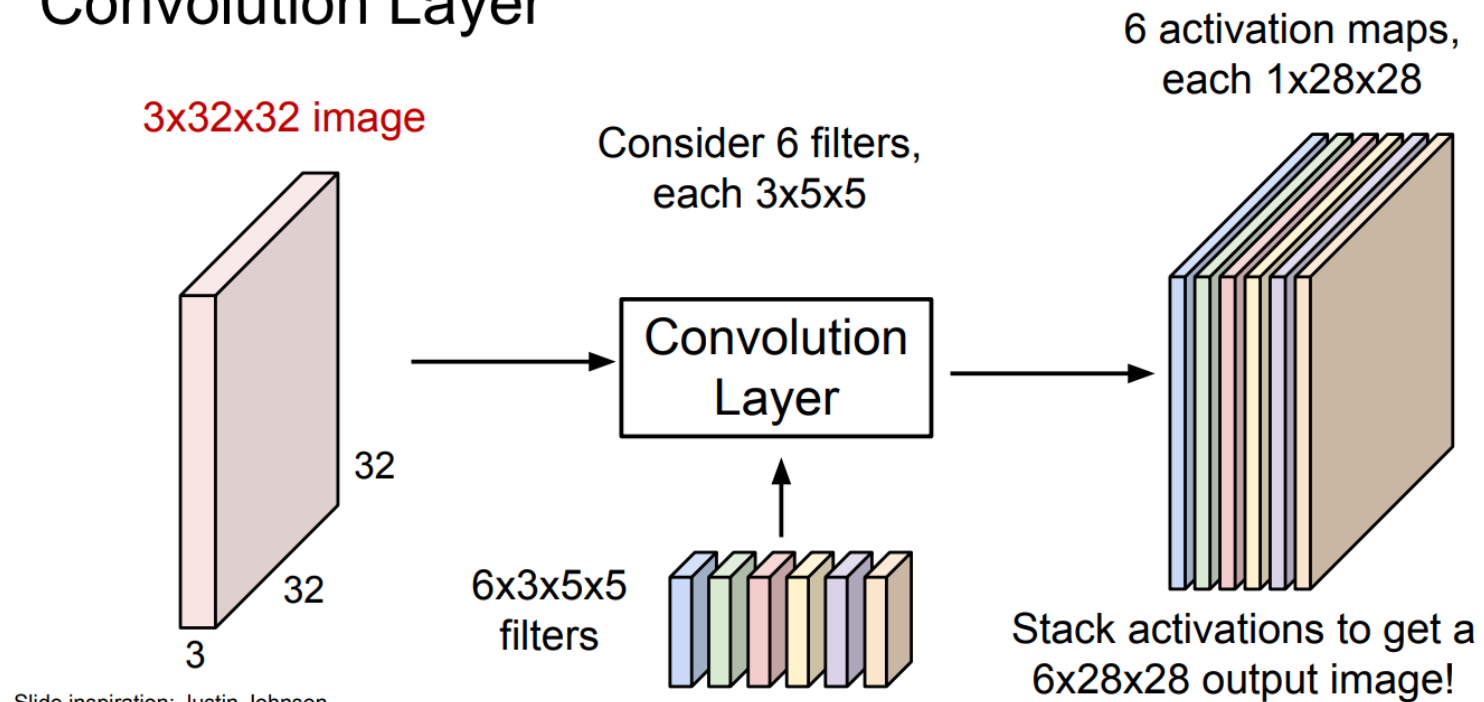


# Convolution Layer

consider a second, **green** filter



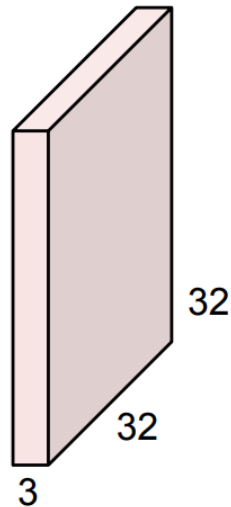
# Convolution Layer



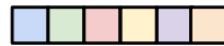
Slide inspiration: Justin Johnson  
Slide Credit: CS231n

# Convolution Layer

3x32x32 image



Also 6-dim bias vector:

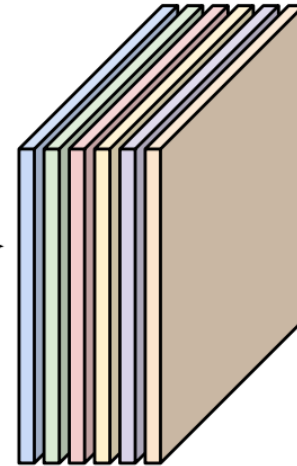


Convolution Layer

6x3x5x5 filters



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



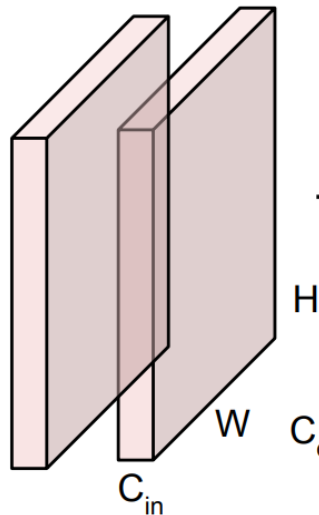
Stack activations to get a 6x28x28 output image!

Slide inspiration: Justin Johnson

Slide Credit: CS231n

# Convolution Layer

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



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

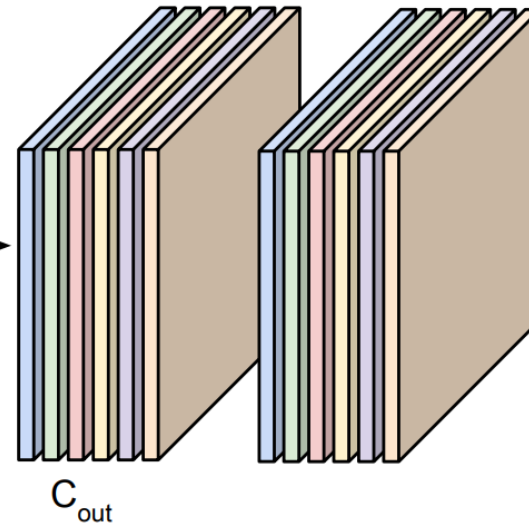


Convolution  
Layer

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



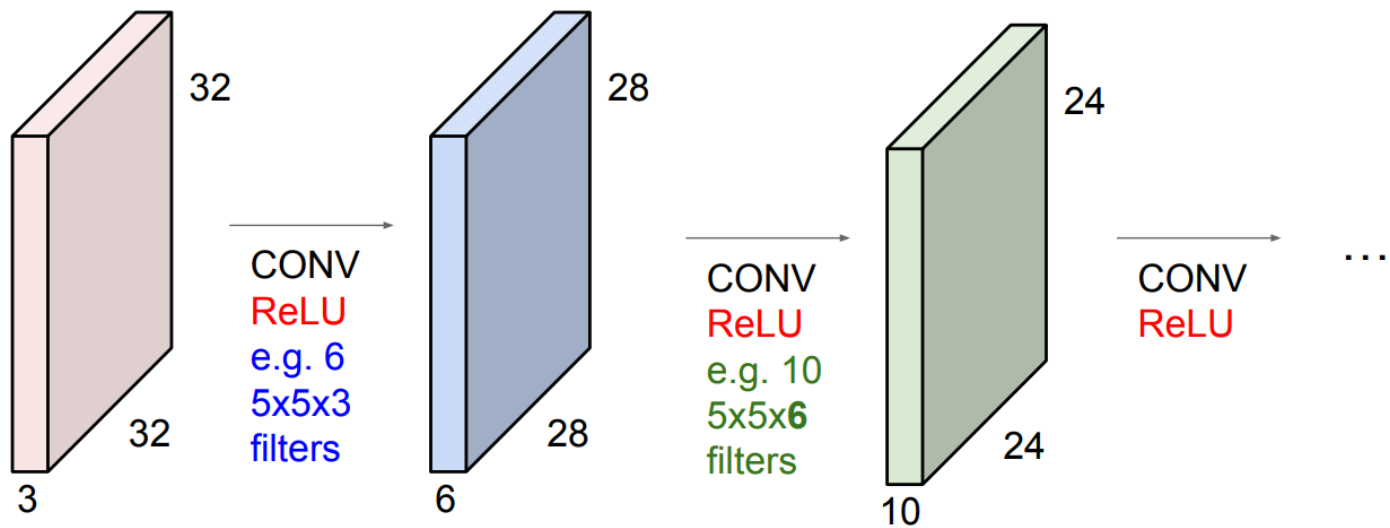
$N \times C_{out} \times H' \times W'$   
Batch of outputs



Slide inspiration: Justin Johnson

Slide Credit: CS231n

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





## Convolution layer: summary

Common settings:

Let's assume input is  $W_1 \times H_1 \times C$

Conv layer needs 4 hyperparameters:

- Number of filters **K**
- The filter size **F**
- The stride **S**
- The zero padding **P**

**K** = (powers of 2, e.g. 32, 64, 128, 512)

- **F** = 3, **S** = 1, **P** = 1
- **F** = 5, **S** = 1, **P** = 2
- **F** = 5, **S** = 2, **P** = ? (whatever fits)
- **F** = 1, **S** = 1, **P** = 0

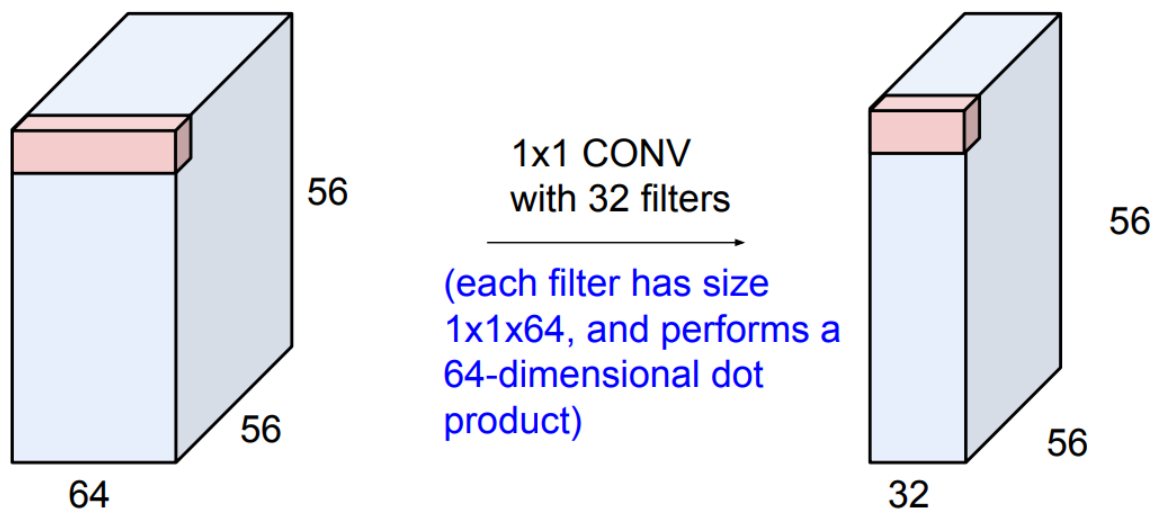
This will produce an output of  $W_2 \times H_2 \times K$

where:

- $W_2 = (W_1 - F + 2P)/S + 1$
- $H_2 = (H_1 - F + 2P)/S + 1$

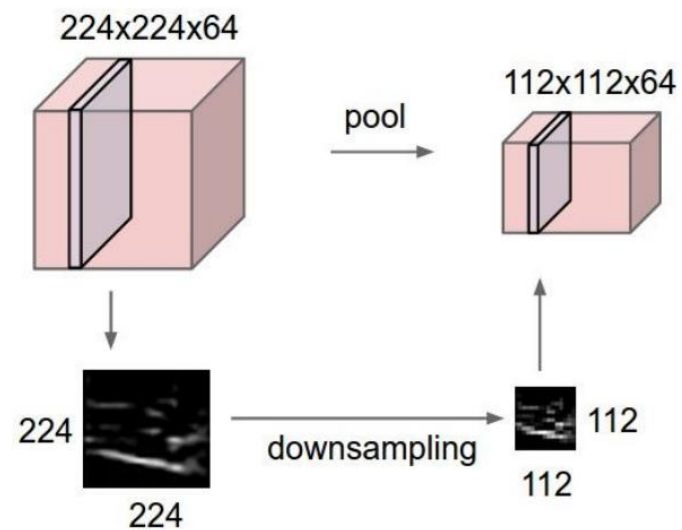
Number of parameters:  $F^2CK$  and  $K$  biases

(btw, 1x1 convolution layers make perfect sense)

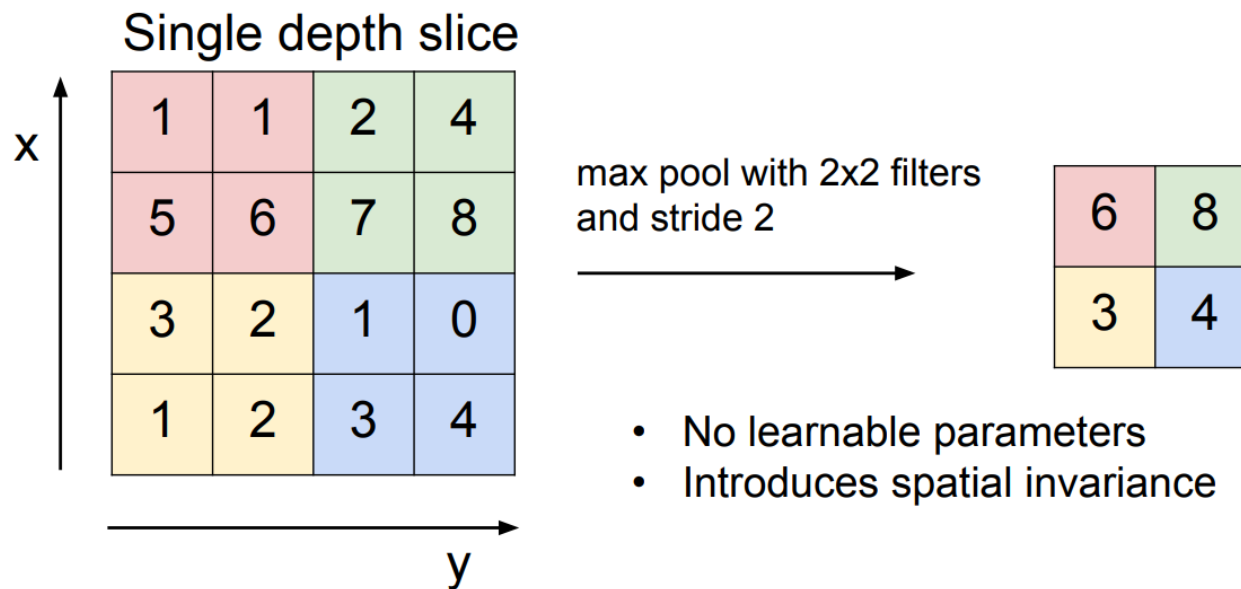


## Pooling layer

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



# MAX POOLING



## Pooling layer: summary

Let's assume input is  $W_1 \times H_1 \times C$

Conv layer needs 2 hyperparameters:

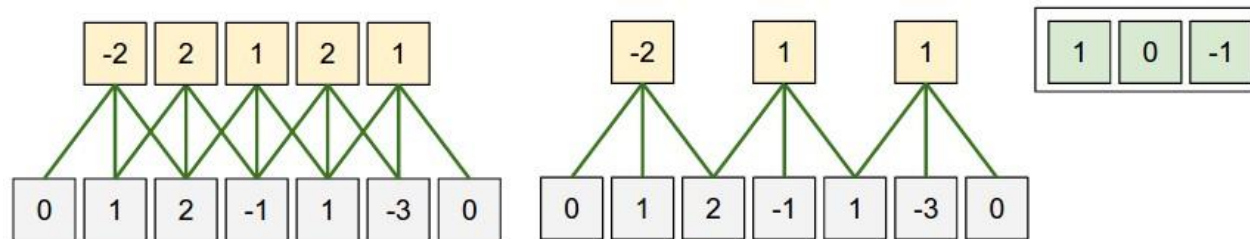
- The spatial extent **F**
- The stride **S**

This will produce an output of  $W_2 \times H_2 \times C$  where:

- $W_2 = (W_1 - F) / S + 1$
- $H_2 = (H_1 - F) / S + 1$

Number of parameters: 0

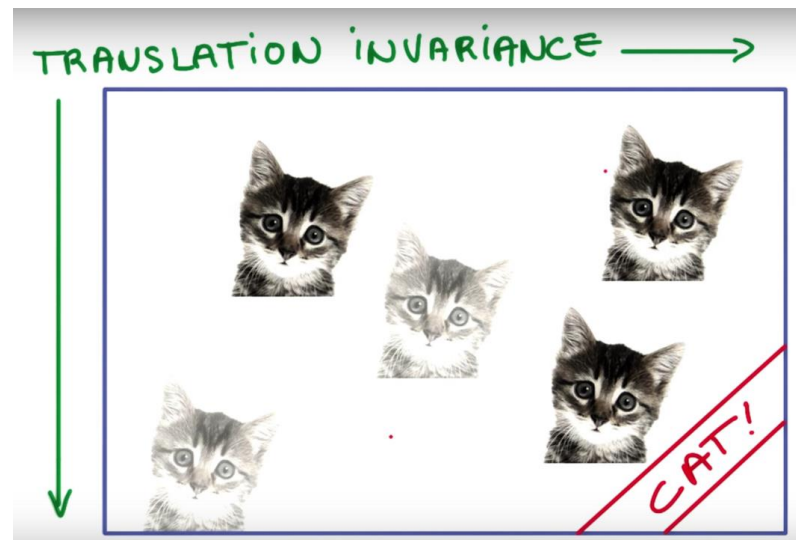
## Parameter Sharing



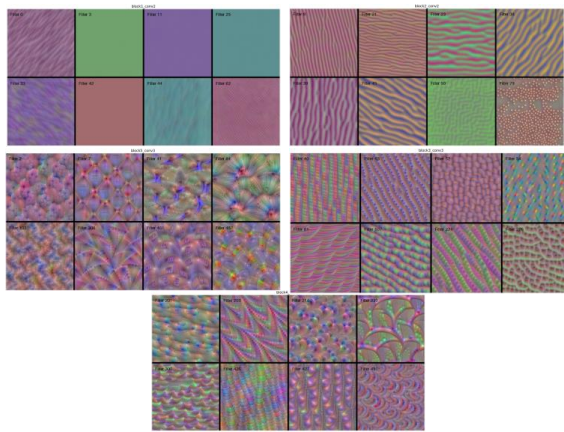
Lesser the parameters less computationally intensive the training. This is a win win as we are reusing parameters.

## Translational invariance

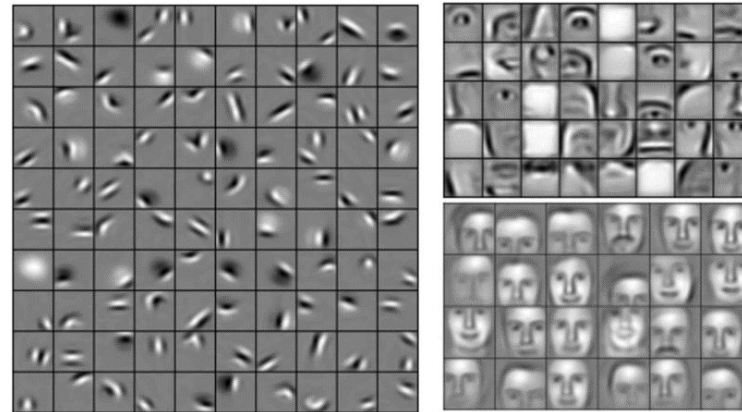
Since we are training filters to detect cats and the moving these filters over the data, a differently positioned cat will also get detected by the same set of filters.



## Filters? Layers of filters?



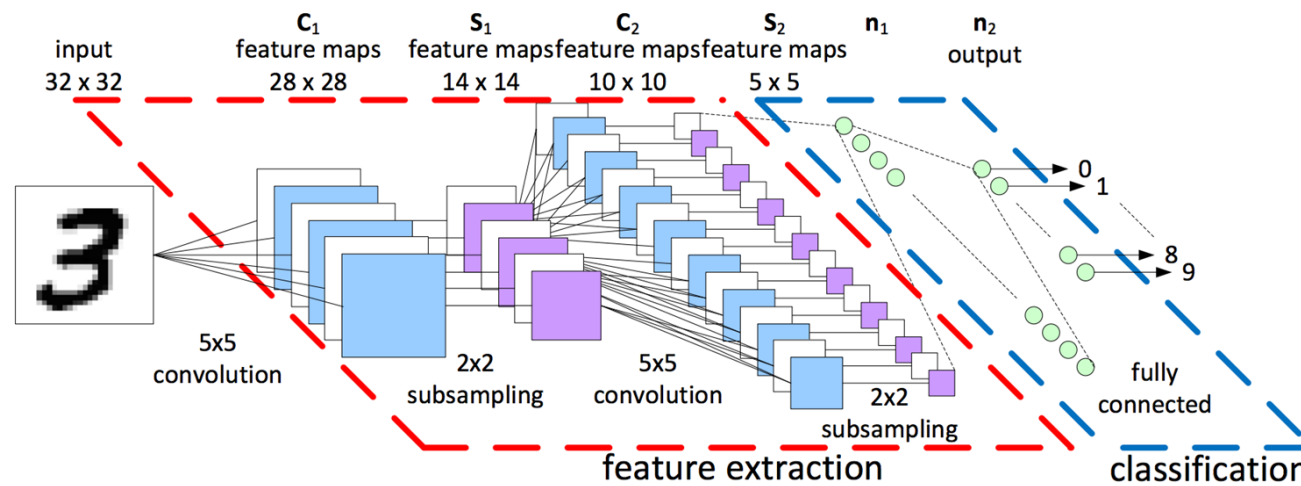
Images that maximize filter outputs at certain layers. We observe that the images get more complex as filters are situated deeper



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?



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:

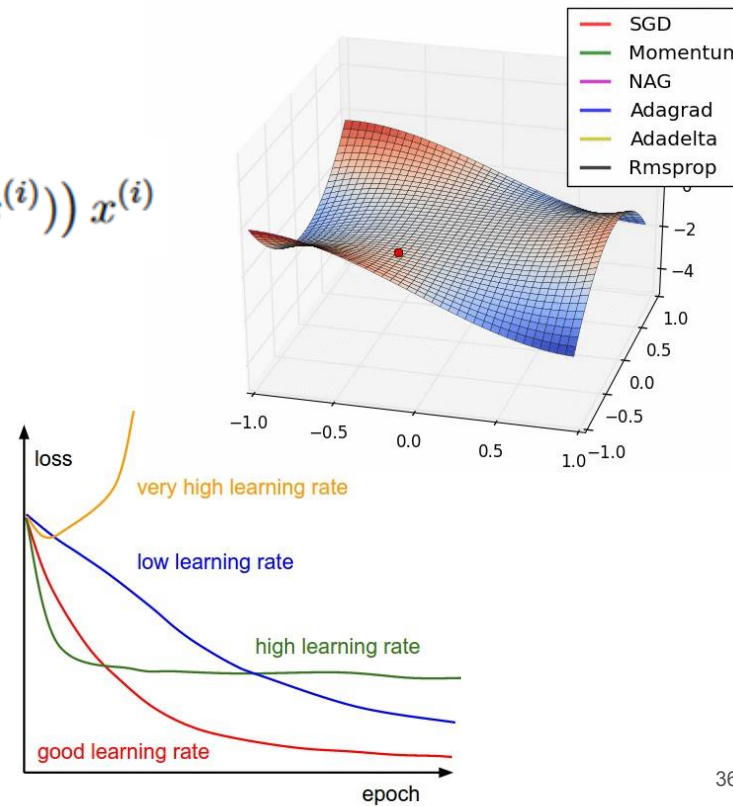
- Just do gradient descent and backpropagate the error derivatives

## How do we learn?

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

There are “optimizers”

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



# Mini-batch Gradient Descent

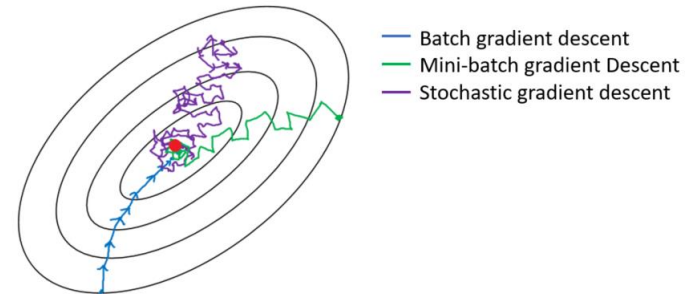
Expensive to compute gradient for large dataset

Memory size

Compute time

Mini-batch: takes a sample of training data

How to we sample intelligently?

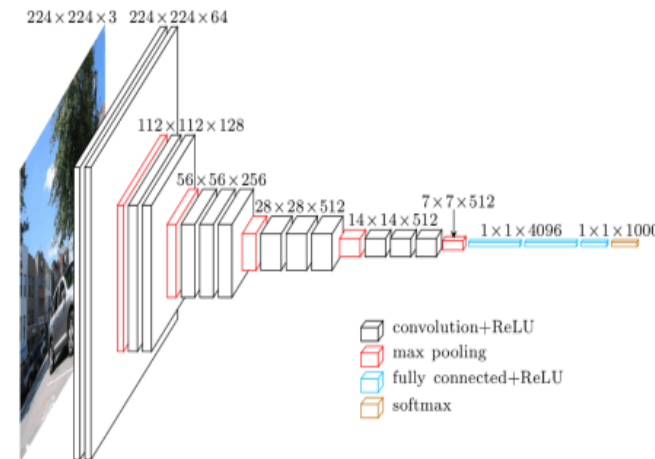


## Is deeper better?

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

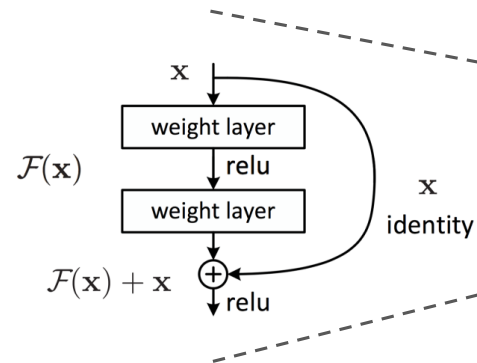
- Loss of information during forward propagation
- Loss of gradient info during back propagation

There are many ways to “keep the gradient going”



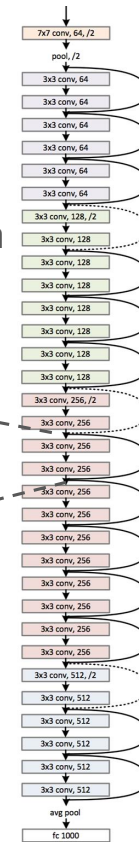
# Solution

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

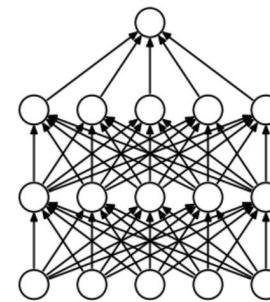
Popular initialization setups

(Xavier, He) (Uniform, Normal)

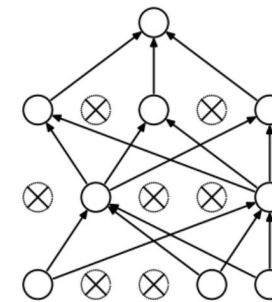


# Dropout

- 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



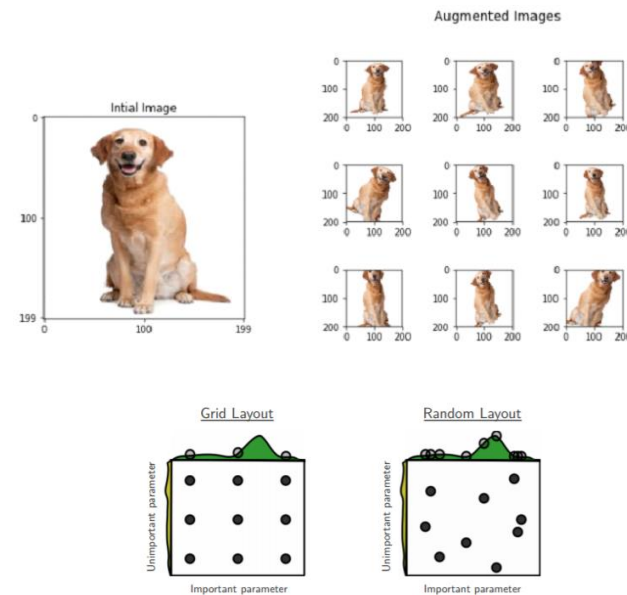
(a) Standard Neural Net



(b) After applying dropout.

## Tricks for training

- 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?

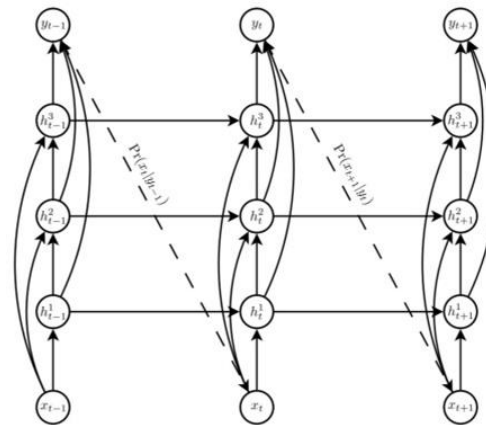


# Overview

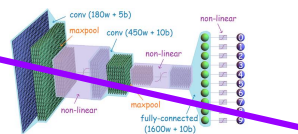
- Motivation for deep learning
- Areas of Deep Learning
- Convolutional neural networks
- Recurrent neural networks
- Deep learning tools

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

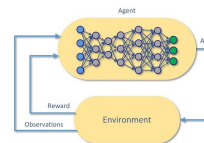
Recurrent NN  
Time Series



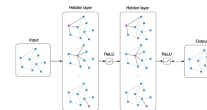
Convolutional NN



Deep RL

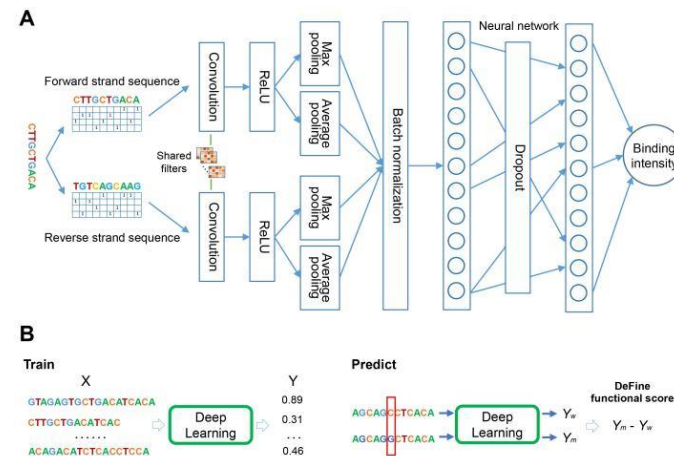


Graph NN



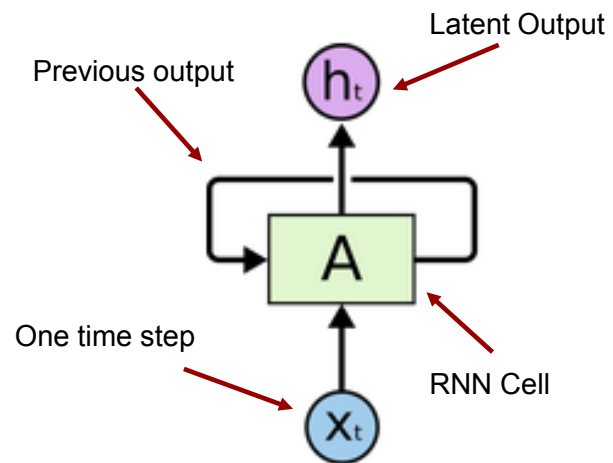
## We can also have 1D architectures (remember this)

- 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**?

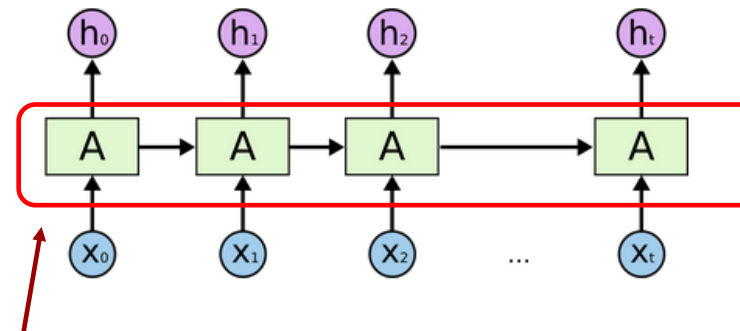


## To address sequential dependency?

Use recurrent neural network (RNN)

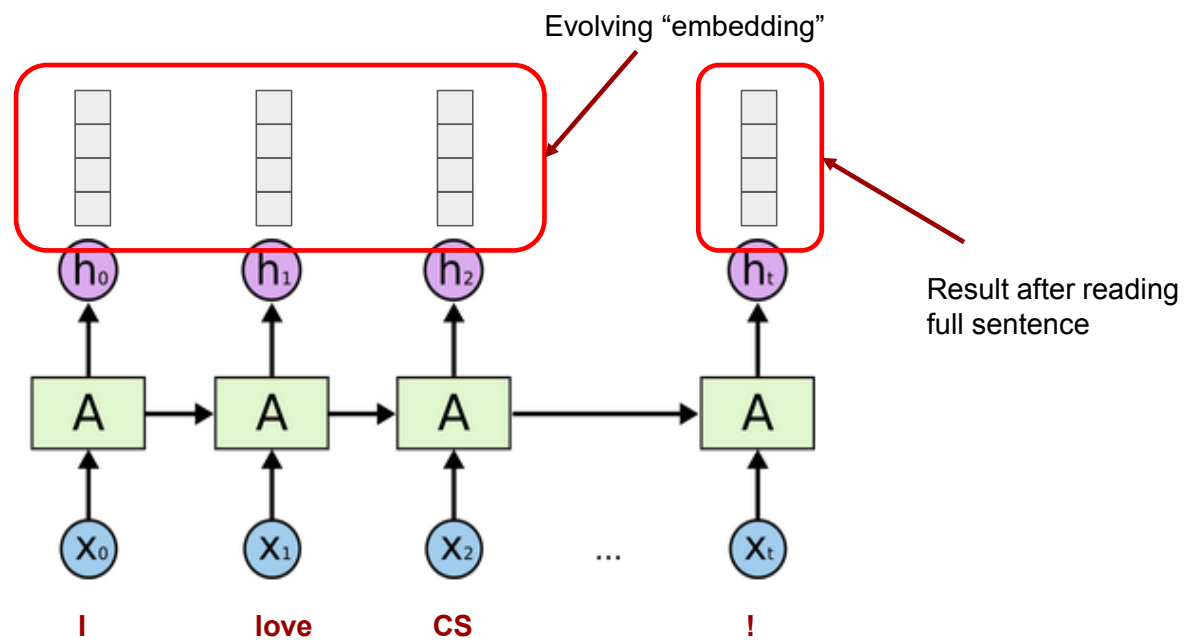


Unrolling an RNN

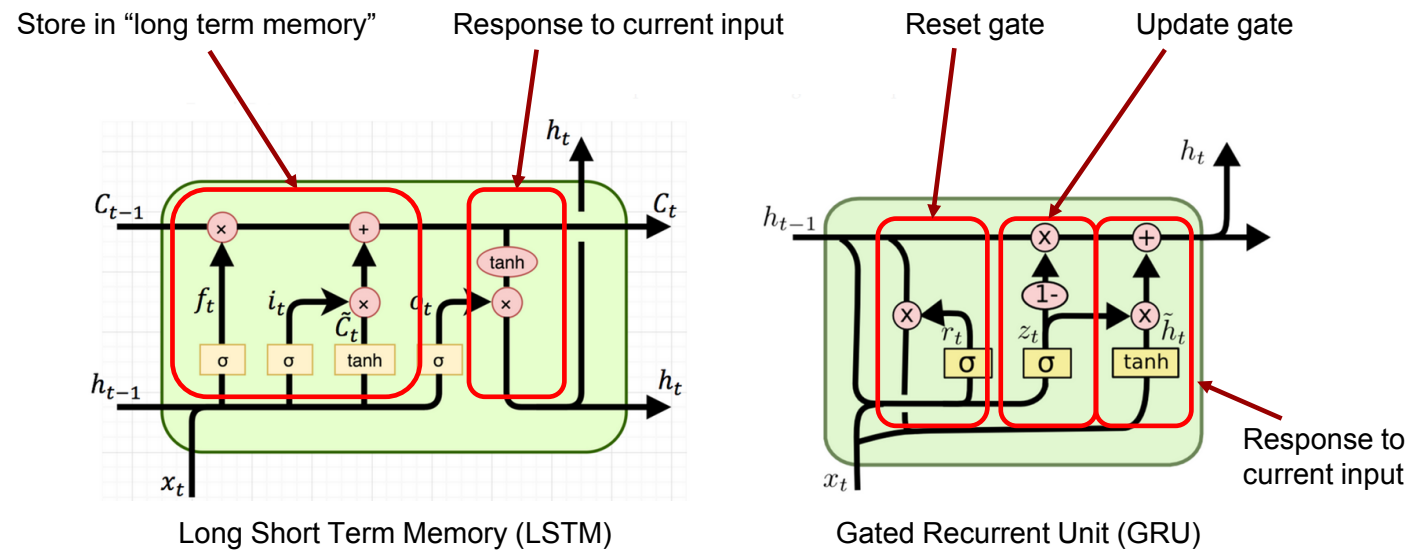


They are really the same cell,  
NOT many different cells like kernels of CNN

## How does RNN produce result?

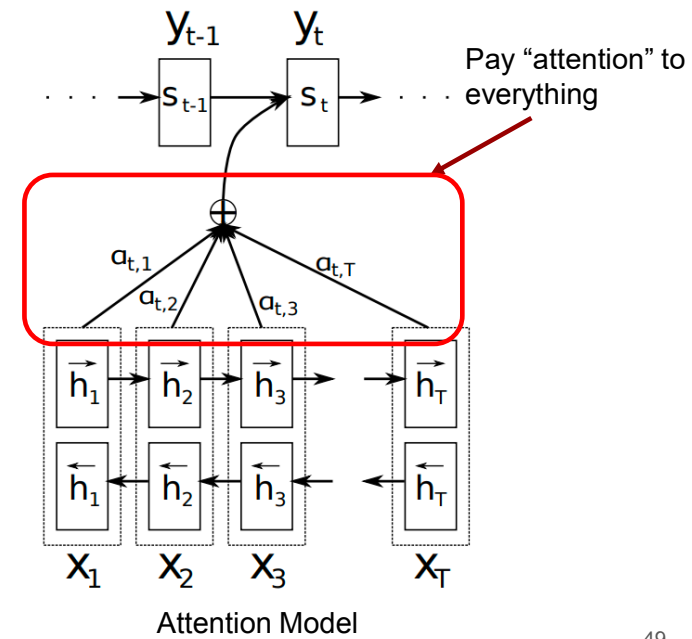
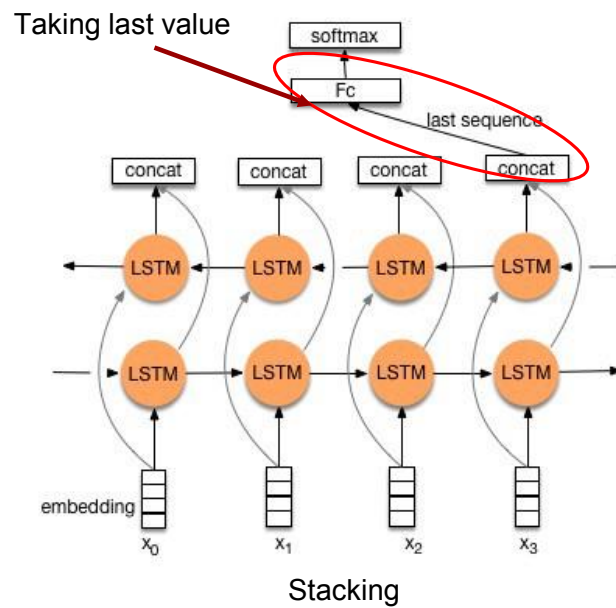


## There are 2 types of RNN cells





# Recurrent AND deep?



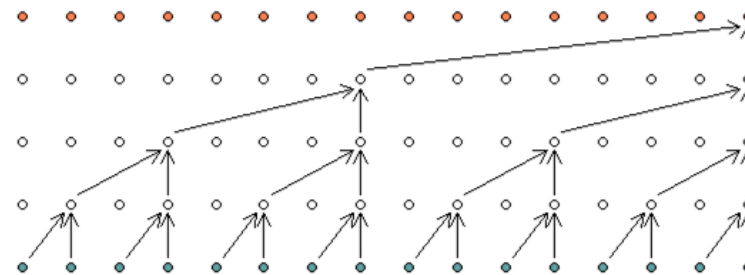
# “Recurrent” AND convolutional?

## Temporal convolutional network

Temporal dependency achieved through  
“one-sided” convolution

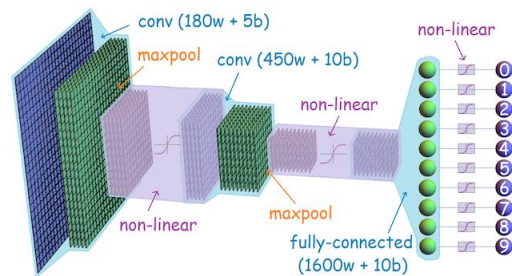
More efficient because deep learning  
packages are optimized for matrix  
multiplication = convolution

No hard dependency

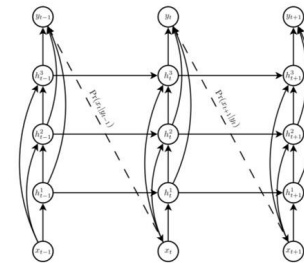


More? Take CS230, CS236, CS231N, CS224N

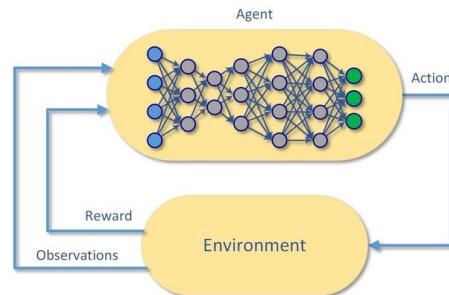
Convolutional NN  
Image



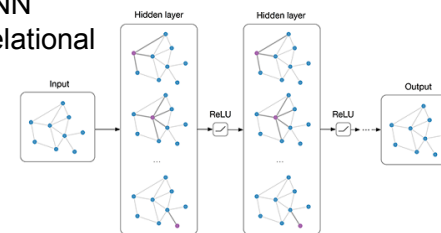
Recurrent NN  
Time Series



Deep RL  
Control System

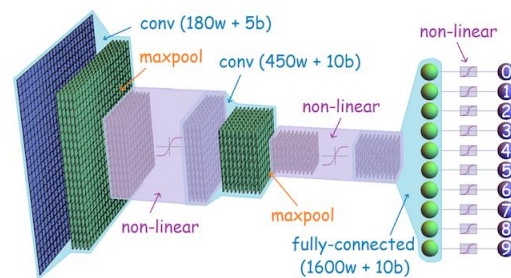


Graph NN  
Networks/Relational

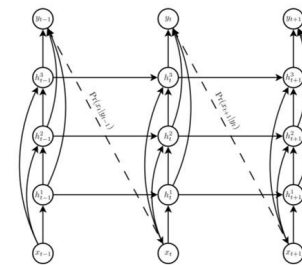


Not today, but take CS234 and CS224W

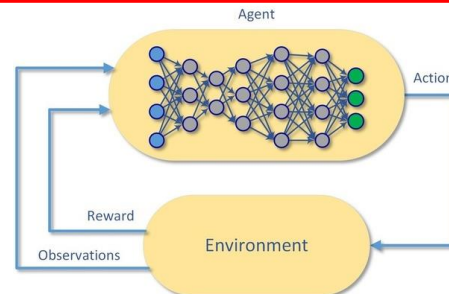
Convolutional NN  
Image



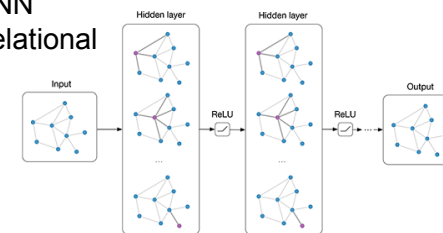
Recurrent NN  
Time Series



Deep RL  
Control System



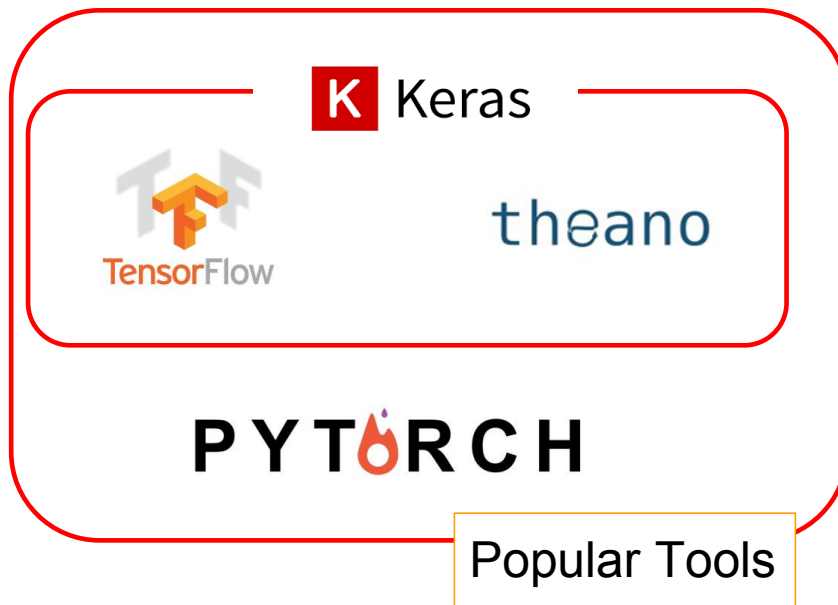
Graph NN  
Networks/Relational



# Overview

- Motivation for deep learning
- Areas of Deep Learning
- Convolutional neural networks
- Recurrent neural networks
- Deep learning tools

## Tools for deep learning



## Specialized Groups



# Where can I get free stuff?

Google Colab

Free (limited-ish) GPU access

Works nicely with Tensorflow

Links to Google Drive



Azure Notebook

Kaggle kernel???

Amazon SageMaker?

Register a new Google Cloud account

=> Instant \$300??

=> AWS free tier (limited compute)

=> Azure education account, \$200?

To **SAVE** money

**CLOSE** your GPU instance

**~\$1** an hour

Good luck!  
Well, have fun too :D

