

## Week #2:- Deep Convolutional Models: Case Studies

Case Studies:- This week we will look at some case studies where CNNs worked better.

- Why Case Studies?
- Good way to see or get intuition about implementation of ConvNets is by looking at other people's work.

### Classic Networks:

- LeNet - 5
  - AlexNet
  - VGG
- } Revolutionized modern DL in CV.  
You'll see that, these ideas will probably be useful for your own application.

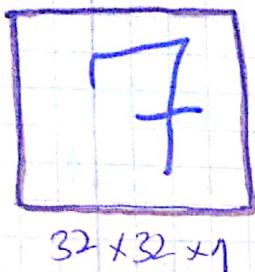
→ ResNet → Residual Network (152 layers).

→ Inception Neural Network.

## Classic Networks:

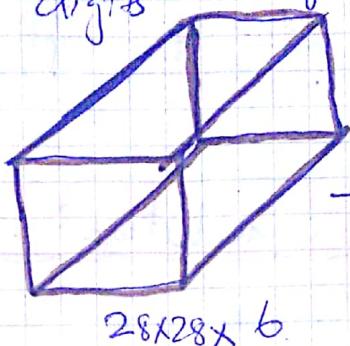
Let's go over some classic neural network architecture.

## LeNet - 5:



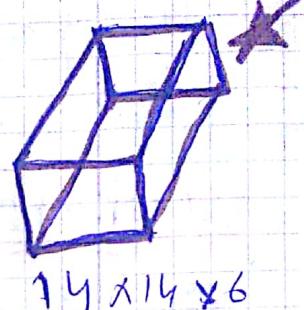
$32 \times 32 \times 1$

$f \in \mathbb{R}^{5 \times 5}$   
stride = 1  
6 filters  
padding = 0



$28 \times 28 \times 6$

avg pool  
 $f = 2 \times 2$   
 $S = 2$

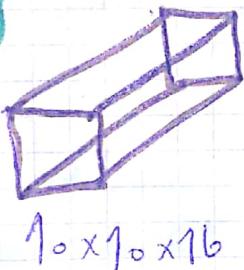


$14 \times 14 \times 6$

trained on grayscale

- \* back then when it was developed, developers used average pooling, now we will use max pooling.

~~\*~~  $f \in \mathbb{R}^{5 \times 5}$   
stride = 1  
16-filters



$10 \times 10 \times 16$

avg pool  
 $f = 2$   
 $S = 2$



$5 \times 5 \times 16$   
400

FC



$g$   
10 possible values

- When this paper was written, people didn't use padding that's why everytime you apply convolution the dimension shrinks.

- \* Modern version of this architecture would use Softmax layer as an output.

- As we go deep into the network, height and width tend to decrease while the number of channels increase

- Pattern in this architecture somewhat resembles the modern ones like

Conv  $\rightarrow$  Pool  $\xrightarrow{1 \text{ or more}}$  Conv  $\rightarrow$  Pool  $\rightarrow$  fc  $\rightarrow$  fc  $\rightarrow$  o/p

**Paper:-** Gradient-based learning applied to document recognition. [LeCun et al., 1998].

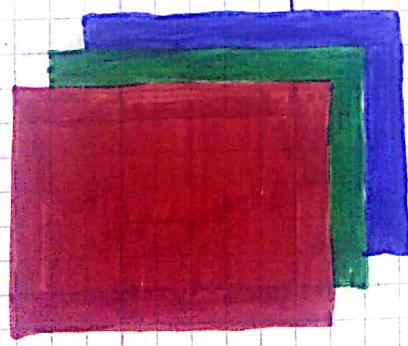
- Back then when this paper was written people were using Sigmoid and tanh non-linearities and weren't using ReLU non-linearities
- There were some funny wiring in the network if you look at the paper (at least by modern standards).
  - e.g. now we use  $(f \times f \times n_c)$  filter for  $(n_h \times n_w \times n_c)$  network where every filter looks at every one of these channels. Back then computers were much slower so, different filters were used assessing different channels of the input block.

### One more detail :-

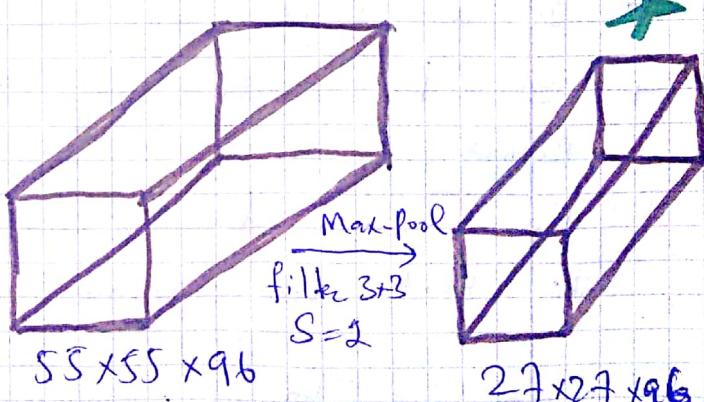
Original LeNet-5 had non-linearity after pooling.  
It uses Sigmoid non-linearity after the pooling layer.

**AlexNet:-** ImageNet classification with deep convolutional neural networks.  
Krizhevsky et al., 2012.

Start with input

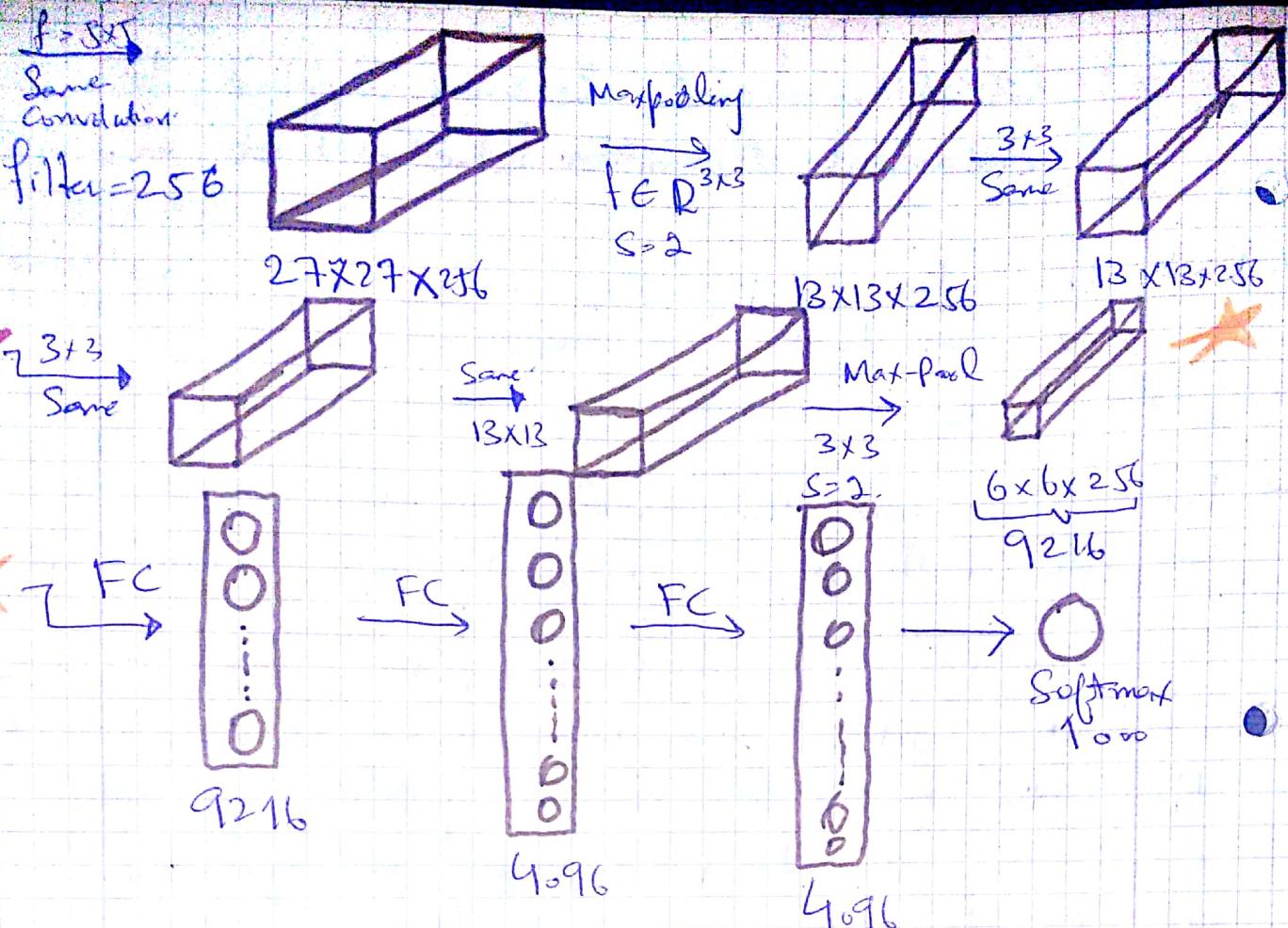


filter  $11 \times 11$   
stride = 4  
96-filters



People refers to

$224 \times 224 \times 3$  img. But 227 makes sense.



→ This NN had a lot of similarities with LeNet, but it was much bigger.

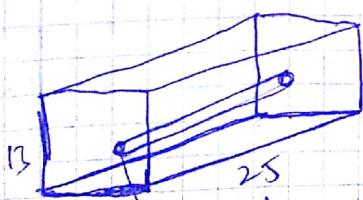
→ AlexNet had 60M parameters, this one had 60M parameters.

→ One aspect that made it much better to LeNet was using **ReLU non-linearity**.

→ It was trained on multiple GPUs.

→ It had another type of layer called **Local response normalization** → Not used much.

But the basic idea is that with anyone of the block take anyone like let's take  $13 \times 13 \times 256$ .



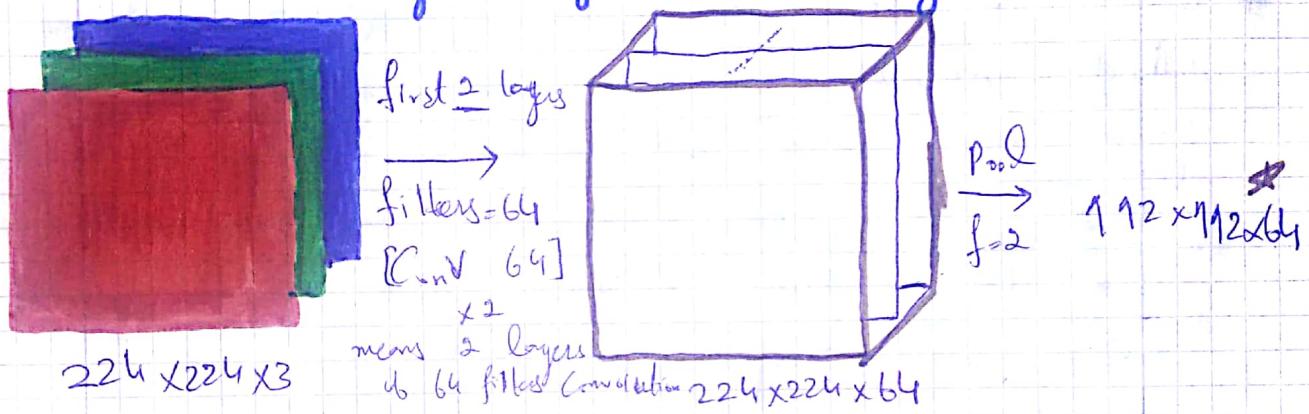
take any of the numbers in  $n_{14}$  and now look all the numbers in that row and normalize them.

→ Motivation for this is that, from each neuron may be you don't want high activation for some. But this doesn't help much.

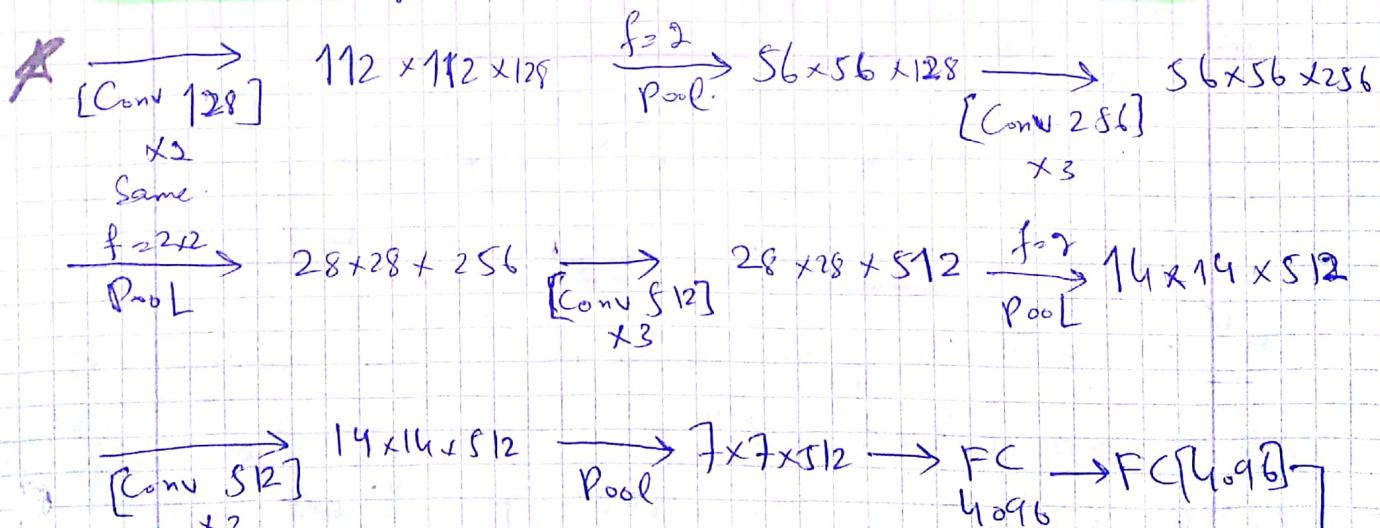
\* This was the paper which motivated DL community to take this area seriously.

→ Complex Network, lots and lots of hyperparameters.

VGG-16: Very deep CNNs for large-scale image recognition. [Simonyan & Zisserman] 15



Remarkable thing :- Instead of having a lot of hyperparameters, let's use a simpler network with Conv  $f = 3 \times 3$  and padding always the 'Same', i.e.  $g = 1$  and all Max pooling layers  $\Rightarrow f \in \mathbb{R}^{2 \times 2}$  with stride of 2.



VGG-16 → refers that it has 16 layers that has weights.

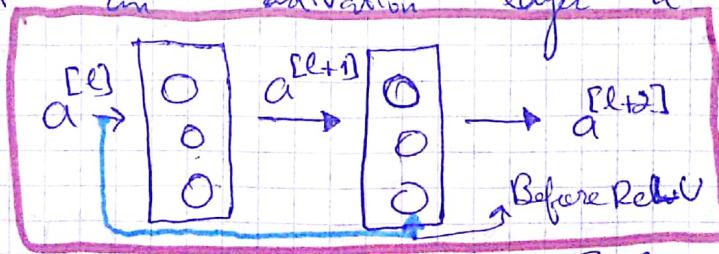
- 1.38 M parameters
- Relative uniformity makes it quite appealing.
- Main downside is that, it has quite large number of parameters.

## ResNets [Residual Networks].

- Very deep neural networks are difficult to train due to vanishing and exploding gradients.
- Let's go through the process of **skip connections** which allows you to take the activation from one layer and suddenly feed it to another layer even much deeper in the neural network and using that to build ResNets which enable you to train very deep neural networks, sometimes even networks of 100 layers.
- Resnets are built over residual block.

## Residual Block:-

Here are 2 layers of N.N where you start off with an activation layer  $a^{[l]}$ .



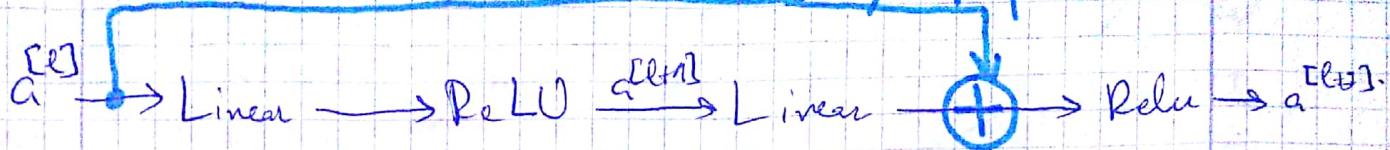
$$a^{[l]} \xrightarrow{\text{linear operator}} a^{[l+1]} \xrightarrow{\text{ReLU}} a^{[l+2]}$$

$$w^{[l+1]} a^{[l]} + b^{[l+1]} \quad a^{[l+1]} = g(z^{[l+1]}) \quad a^{[l+2]} = g(z^{[l+2]})$$

In order to go from  $a^{[l]}$  to  $a^{[l+2]}$ , you'll have to go through all of aforementioned steps.

→ In residual network, we are gonna make changes to this.

## "Short cut" / Skip connection



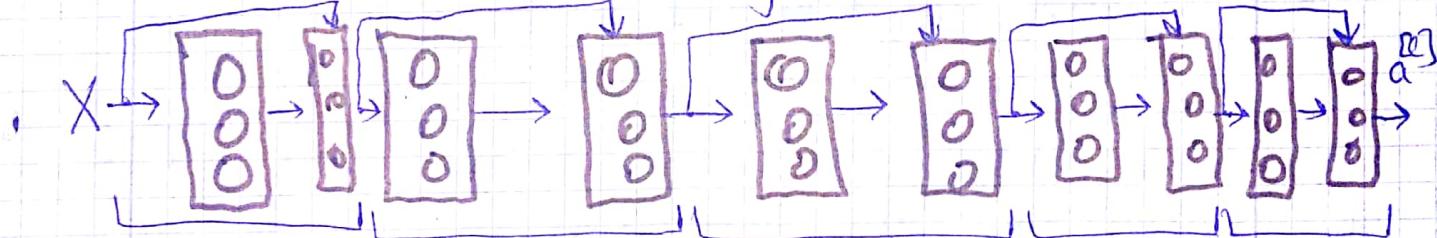
So that means for the last layer instead of using  $a^{[l+2]} = g(z^{[l+2]})$  we are gonna use:

$$a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

↓  
residual block.

- What developers of ResNets found is that, Using residuals allow you to train very deep neural networks
- The way you build ResNets is that, you take ~~these~~ many of these residual blocks and stack them together to form a deep network.

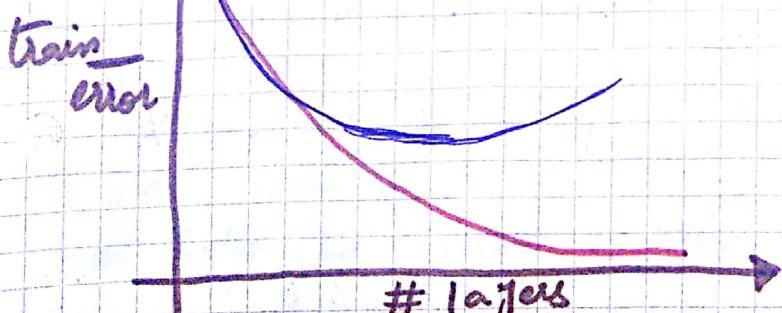
Let's look at the following (not residual but "Plain network")



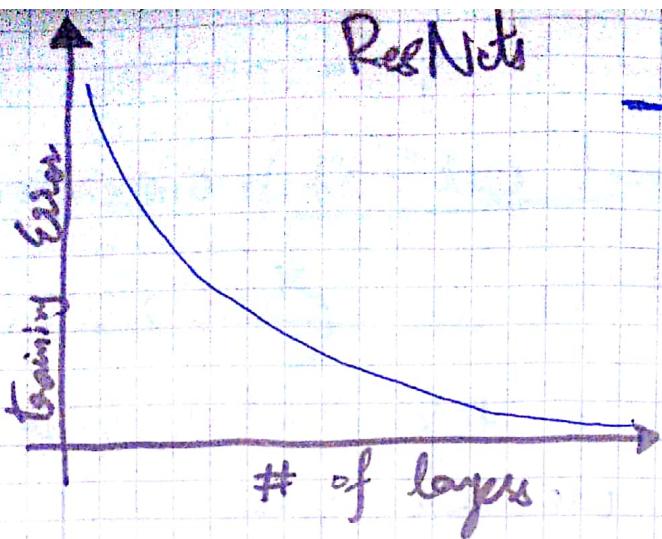
- In order to turn this into residual networks you add all those skip connections:

Plain

— original —  
— theoretical —



Optimization algorithm will have much harder time reducing the cost for very deep network.



- Reality

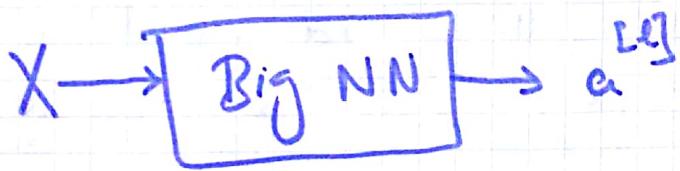
# ResNets' skip connections really helps with vanishing and exploding gradients and allow you to train very deep networks.

### Why ResNets Work?

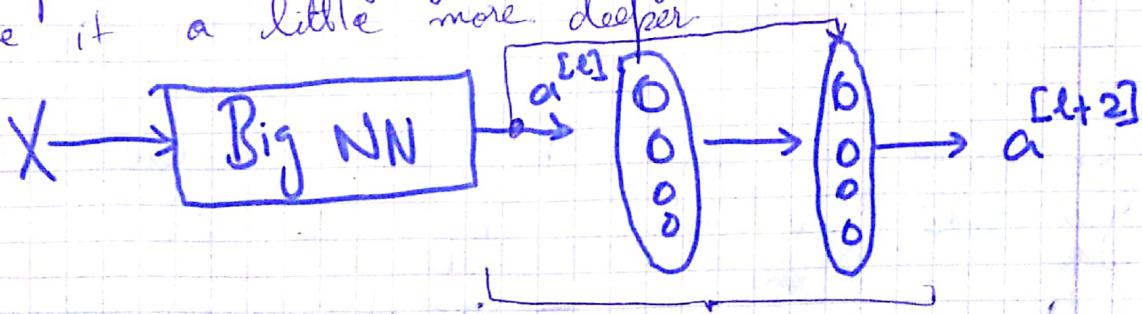
Let's see an example:-

- As we have seen that, if you make a network deeper, it can hurt your ability to train the network to do well on the training set.
- But, this is not the case when you are training a ResNet.

e.g:-



For example, let's say you modify the normal network to make it a little more deeper



Residual Block

all the activations are  $\geq 0$  cause we are using ReLU, with the possible exceptions of the input  $X$ .

Let's look at  $a^{[l+2]}$ :

$$a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

$$a^{[l+2]} = g(w^{[l+2]} a^{[l+1]} + b^{[l+2]} + a^{[l]})$$

- If you are using L2 regularization or weight decay that would tend to shrink the value of  $w^{[l+2]}$ .
- If you apply weight decay to  $b$ , that will also shrink them.

But  $W$  is the term to pay attention here.

Suppose  $w^{[l+2]} \approx 0$  and  $b \approx 0$  so

$$a^{[l+2]} = g(a^{[l]}) = a^{[l]}$$

because we are using ReLU activation function here.

→ This means that, identity function for residual block is easy to learn, and it is easy to get  $a^{[l+2]} = a^{[l]}$  because of this skip connection.

→ So what this means is that, adding block to your neural network doesn't really hurt your neural network's ability as well as this simpler network without these two extra layers, because it is easy for it to learn the identity function.

→ So, adding this block in the middle of your big network doesn't really hurt your performance.

→ But our goal is not just to "not hurt performance" but help these units to learn something better rather than identity.

→ What goes wrong in these very big nets (Plain) is that, it is actually very difficult for the network to choose parameters that learn even the identity function which is why a lot of layers end up making your result worse rather than better.

→ I think, the main reason ResNets work is that, it can learn identity function which doesn't hurt performance and a lot of times you get lucky and it even helps performance.

In resNets our baseline is considered as "not hurting performance" and optimizer can only improve from there.

$$\rightarrow \text{In } a^{[l+2]} = g(\underline{z}^{[l+2]} + \underline{a}^{[l]})$$

We are assuming that,  $\underline{z}^{[l+2]}$  and  $a^{[l]}$  have the same dimensions, that's why there is a dot of usage of "some" convolution.

→ Suppose  $\underline{z}^{[l+2]}$  and  $a^{[l]}$  have dimension clash then you can add another matrix 'W' as a multiple with  $a^{[l]}$

$$\text{Suppose: } \underline{z}^{[l+2]} \in \mathbb{R}^{J \times K} \quad a^{[l]} \in \mathbb{R}^X$$

then  $W \in \mathbb{R}^{J \times X}$

$$\text{and } a^{[l+2]} = g(\underline{z}^{[l+2]} + W a^{[l]})$$

$W$  can be learned or remain there as a fixed matrix which adds zero padding in  $a^{[l]}$ .

# Paper: Deep residual networks for image recognition. [He et al]

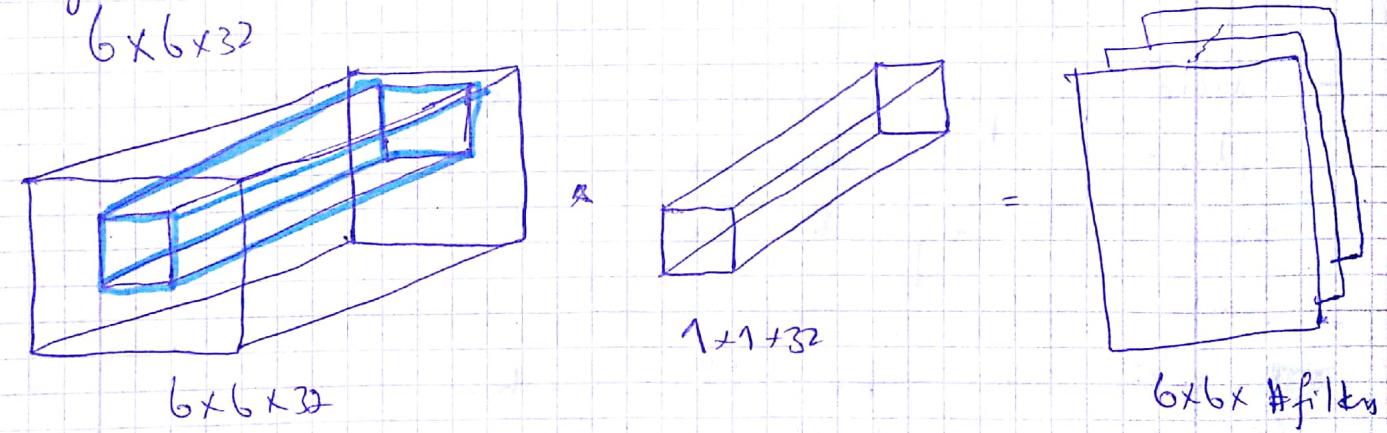
- Add residual layers to plain network.  
[Look up this paper online].
- whenever there is padding dimensions are gonna change  
then add 'Ws' there as a multiple of 'a' with  $W_{\text{res}}$

## Network in Network and $1 \times 1$ Convolutions:-

- In terms of designing ConvNet one of the ideas that might help is  $1 \times 1$  convolution.

### What does a $1 \times 1$ Convolution do?

$6 \times 6$  image convolve with  $1 \times 1$  filters is just like multiplying  $6 \times 6$  image with a number. But imagine if we had an ip block of



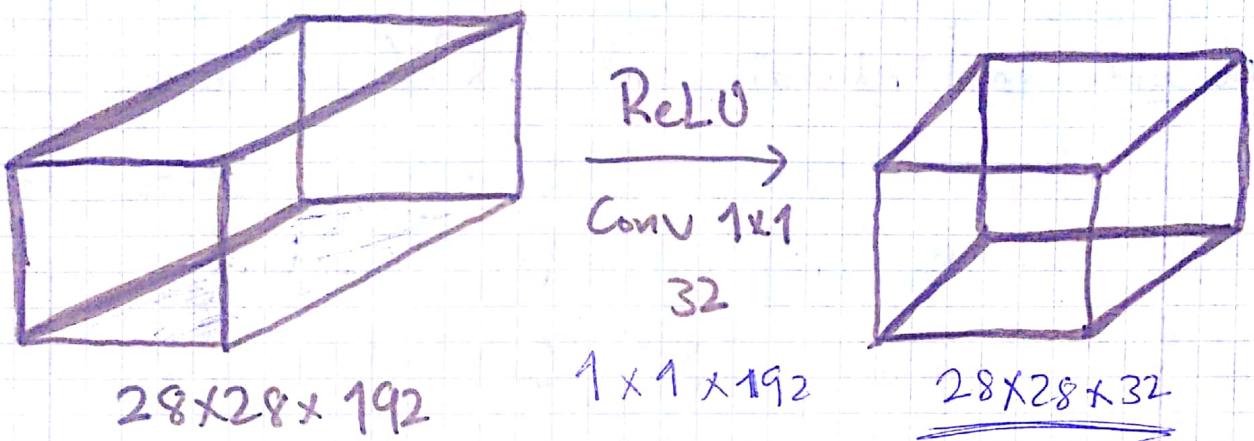
So, in this case multiply element wise and add them to get a sum number.

- It can be thought of as one neuron taking 32 numbers processing it applying ReLU and output non-linearity.

This idea is sometimes called Network in Network.

Here is an example where  $1 \times 1$  Convolution is useful.

If your height and width of image get too big you can use pooling what if channels increase, then use  $1 \times 1$  convolution.



### Inception Network motivation :-

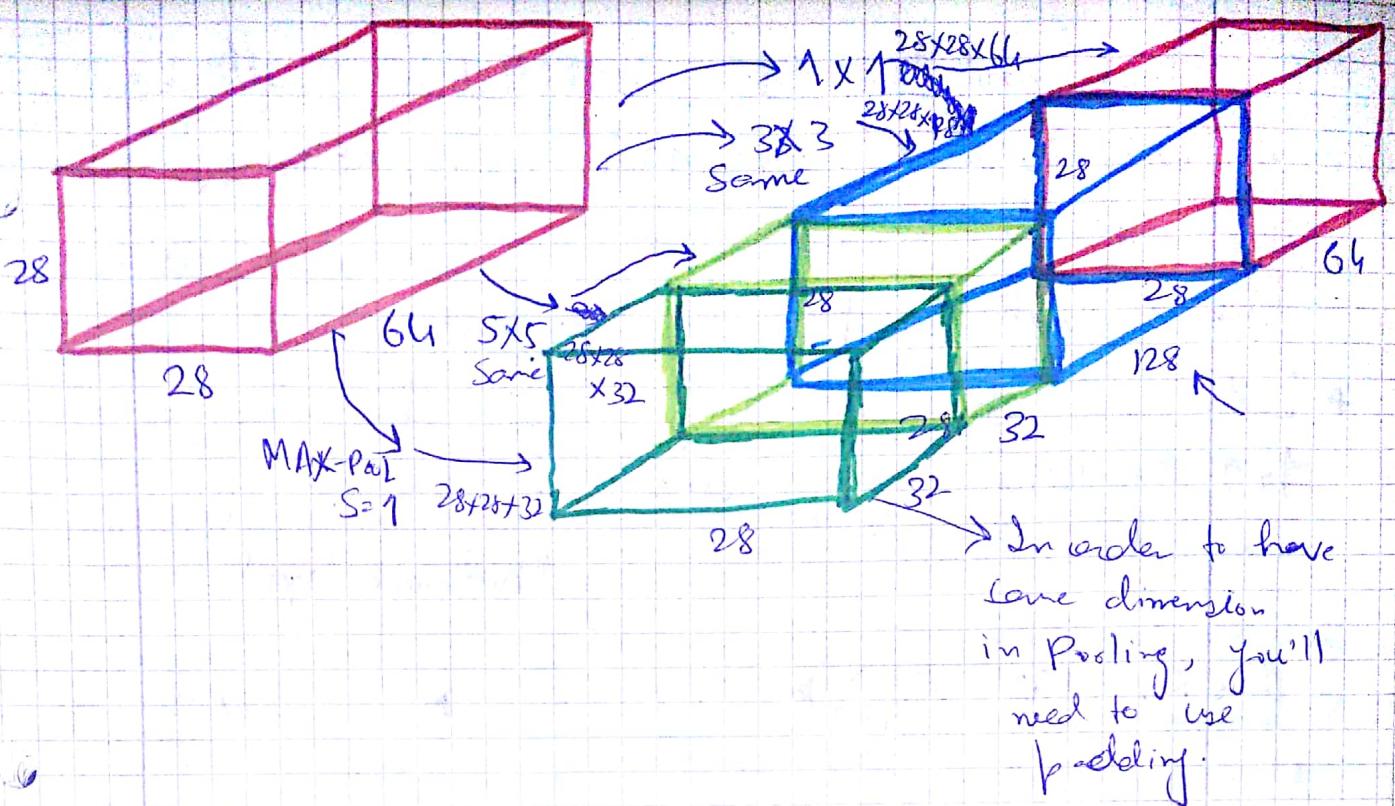
When designing a ConvNet you have to pick which filter you want,  $3 \times 3$ ,  $5 \times 5$  --- etc

→ What inception network suggests is that, why don't you try them all, this makes network more complicated but it works remarkably well.

→ Let's see how this works.

(Paper: Szegedy et al. 2014. Going deeper with convolutions)

Suppose, we have  $28 \times 28 \times 192$ . Now what filter size to use, whether to do convolution or pooling, inception network: Well let's do them all and see which one works best.



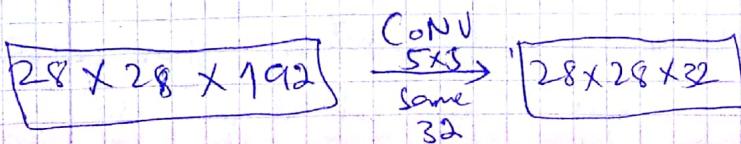
→ You had  $28 \times 28 \times 64$  i/p and you o/p  $28 \times 28 \times 256$  by trying different things. **This is the heart of Inception Networks.**

"Try them all and stack up the results.

Then see and let the network learn which combination it want to go with.

- There is a problem with inception networks and that is, **Computational Cost.**

From above example let's figure out what problem the  $5 \times 5$  convolution could pose in terms of computational cost.

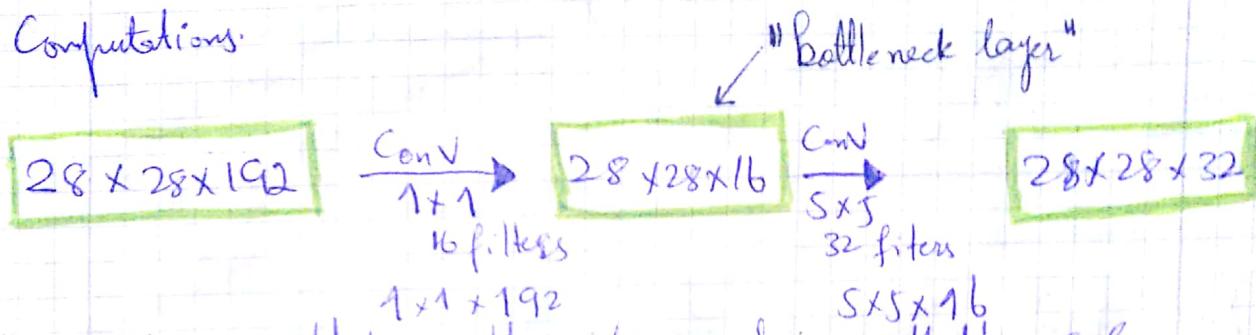


→ You have 32 filters - each filter is  $5 \times 5 \times 192$ . You need to compute  $28 \times 28 \times 32$  numbers. For each of them you need to do  $5 \times 5 \times 192$  multiplications.

So, total multiplications will be  
 $(28 \times 28 \times 32) \times (5 \times 5 \times 192) = 120 \text{ Million.}$

Seems like a normal task for today's computer.  
This is still a pretty big operation.

→ Let's see how can we utilize the idea of 1x1 convolution and reduce the amount of computations.



Observe that the i/p and output are still the same as previous example.

Let's look at the computational cost involved.

For computing bottleneck layer

$$(28 \times 28 \times 192) \times (1 \times 1 \times 192) \text{ multiplications} = 2.4 \text{ M.}$$

How about the second layer

$$(28 \times 28 \times 32) \times (5 \times 5 \times 16) \text{ multiplications} = 10.0 \text{ M.}$$

So Total Multiplications are 12.4 M.

Wooh! (P), Right ;)

From 120 M  $\rightarrow$  12.4 M round about 90% decrease

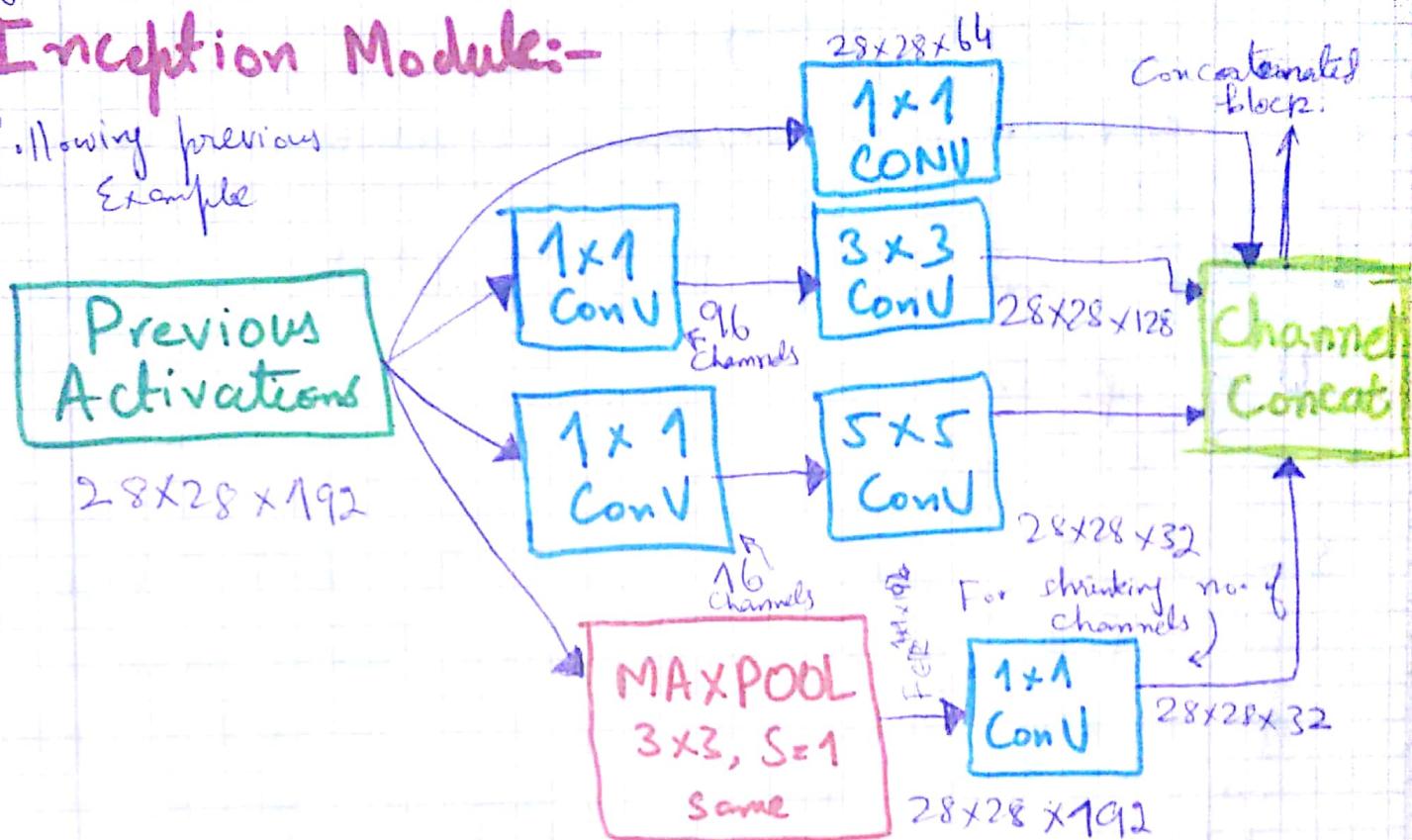
→ You must be thinking that, whether your performance remains the same or not by shrinking the size so dramatically. It turns out that, so long as you implement this strategy reasonably it doesn't hurt your performance.

## Inception Network :-

We have seen the building block for Inception network, let's see how can we utilize all of them to build a full network!

### Inception Module:-

following previous example



- This is one inception module.
- What Inception network does is that, take a lot of these modules together.

In full big Inception Network, along the network to the output, in between we try to compute/predict the output through the layers in between.

Just another detail of Inception network. What it does is that, it helps ensure that the features computed even in the intermediate layers, that they are not too bad for predicting the output. Cost of the image.

This appears to have regularization effect which helps prevent this layer from overfitting.

## Where does the word Inception come from?

Have you seen "Inception" ;). I think you can understand lol :)

<http://knowyourmeme.com/memes/we-need-to-go-deeper>

## Practical advices for Using ConvNets:-

### Using open-source Implementations:-

git clone " \_\_\_\_\_ "

Some of the open source material's authors have used high power GPUs and massive datasets for training their models. So, you might not have access to that kind of computational power or data, so replicating the network might do the trick for you (Transfer Learning).

## Transfer Learning:-

If you are building a CV application, it is often a good way to start by looking someone else's implementation and use those parameters/ hyperparameters in order to sort your work and start learning from there.

Let's take a deeper look and start with an example:-

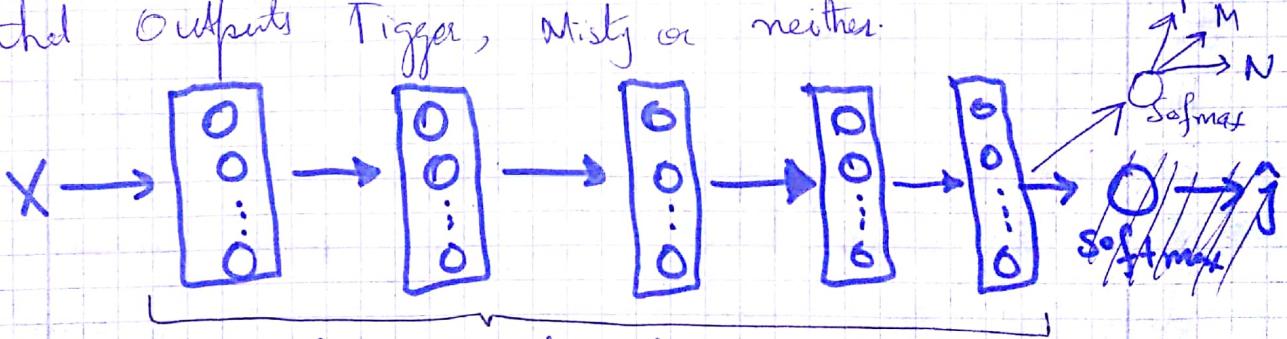
Let's say you are building a Cat classifier.

According to internet Tigger and Misty are two common cat names.

- You have a classification problem with 3 classes, i.e. is it a Tigger, Misty or None in the image?
- Now probably you don't have a lot of pictures of Tigger or Misty, so your training data would be small, now what can you do?
- Go online and download some open source implementation of Neural Networks, not just the code but also the weights.

There are lot of online NN architecture which you can download e.g. image net data has thousand different classes So network will have 1000 softmax outputs.

- What you can do is that, get rid of that layer and create your own softmax layer that outputs Tigger, Misty or neither.



Consider all these layers as frozen

- only just train parameters associated with softmax layer.

- By using someone else's work, you might be able to get pretty good performance even with small dataset.
- Deep learning frameworks support this sort of things like by setting trainable parameters = 0, freeze = 1. You specify whether or not to train the weights associated with a particular layer.
- One neat trick that might help with some implementations is that, because all of these early layers are frozen, there is some fixed function that doesn't change, because you are not changing it, so, it can be seen as you input any image and maps it to output Softmax layer.
- One thing that could speed up your training is that, Just compute the values for the layer before Softmax for all the training examples you have then start training it.  
This way you won't need to recompute those activations on every optimization iteration.
- On the other hand:- If you have massive data set then what you can do is freeze some layers and train the rest or keep the frozen layers and replace the rest and attach your own setting of layers there.  
Either ways are fine.

Bottom-line:- If you have big enough dataset, no of layers frozen could be smaller and unfrozen greater and the opposite is true if you have small dataset.

→ For big enough dataset you can also train the whole thing.

## Data Augmentations:-

- Most Computer Vision tasks could use more data.
- Data augmentation is one of the techniques to improve the performance of Computer Vision System.  
This is true ~~also~~ regardless of the approach you are following i.e. ↑ Transfer learning or building everything from scratch. either

→ Most common method is mirroring the image.

- Another is Random cropping

↳ Won't work in some cases, you have to make sure that cropping settings are legible.

→ Rotation.

→ Shearing.

• → Local warping.

} In practice they are used less, perhaps due to the complexity.

## Color Shifting:

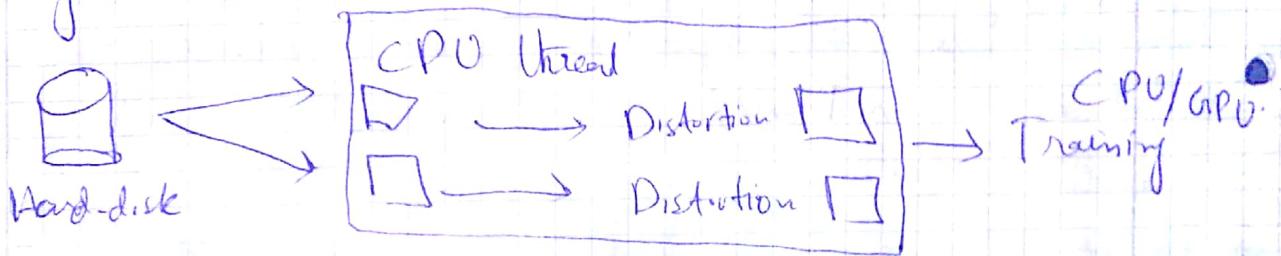
Let's say you have an image and you add some distortion to every channel like  $[(R+20), (G-20), (B+20)]$ . You can set other settings like  $[(R+5), G, (B+5)]$ .

→ In practice these values are drawn from some probability distribution.

There are different ways to sample R, G, and B.

One of the algorithm is called PCA  $\Rightarrow$  Principle Component Analysis.

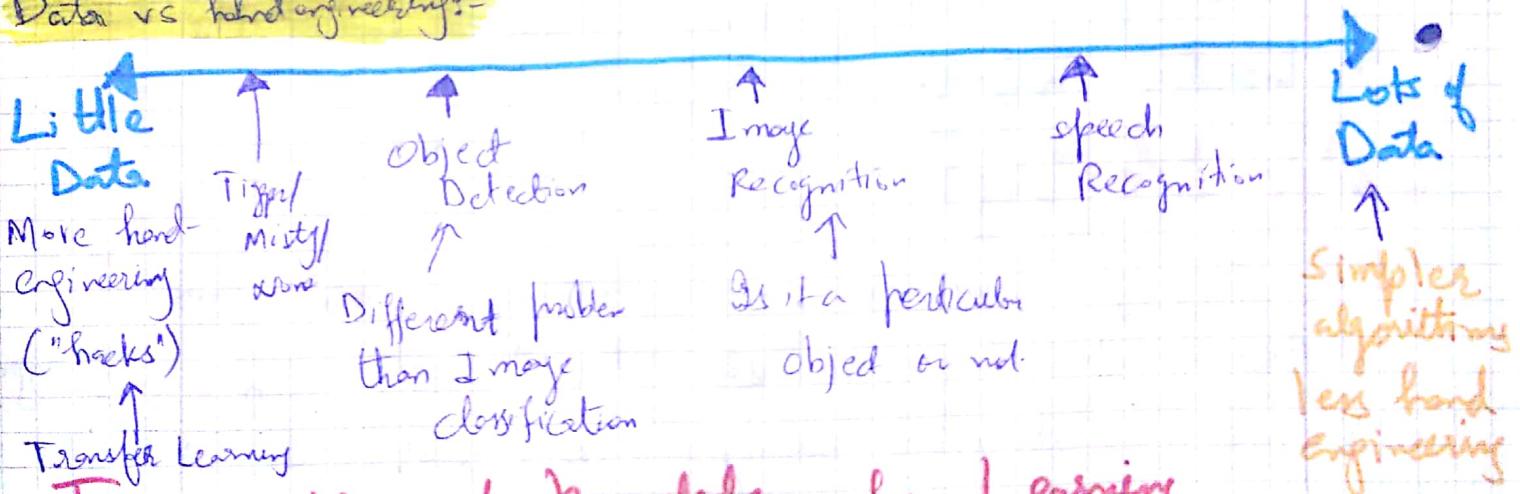
PCA Color augmentation  $\Rightarrow$  Rough idea is that, if your image is mainly purple ie ~~mainly~~ red or blue pixels then PCA will add or subtract a lot of Red or blue images in order to give overall color the same.



### State of Computer Vision:-

Deep learning has a lot of applications, but there are few things about deep learning that are unique just for Deep learning, let's see.

### Data vs hand engineering:-



### Two sources of knowledge for Learning

#### Algorithms:-

- Labeled Data ( $n, j$ )
- (Hand engineering / network architecture / other components).

- Since we don't have much data for CV applications  
that's why we use more complex networks.

## Tips for doing well on benchmarks / winning Competitions

- Ensembling:

- ↳ Train several networks independently and average their outputs. (**not weights**) augments 1-2% performance.



- Multi-Crop at test time:

- ↳ Run classifiers on multiple versions of test images and average results.

- It is a form of providing data augmentation at test time to your test images.

like **10-crop** Suppose you have an image and its mirror reflection. Take 1 each image of a cropped part from both. Now take 4 images each from both of them with a fixed setting from 4 sides. That makes up your 10 images.



Run them through your classifier then average the results.

- These tips are used for benchmarking not in production.

### Use open Source Code:-

In CV application, transfer learning works.

- Use architectures of networks published in the literature
- Use open source implementations if possible.
- Use pretrained models and fine-tune on your dataset.
  - ↳ Because lot of little details like learning rate, decay rate are tuned