# WEEK 2

In this week we will look at case studies to get a better understanding of CNNs.

Outline :

- Classic Networks:
  o LeNet - 5
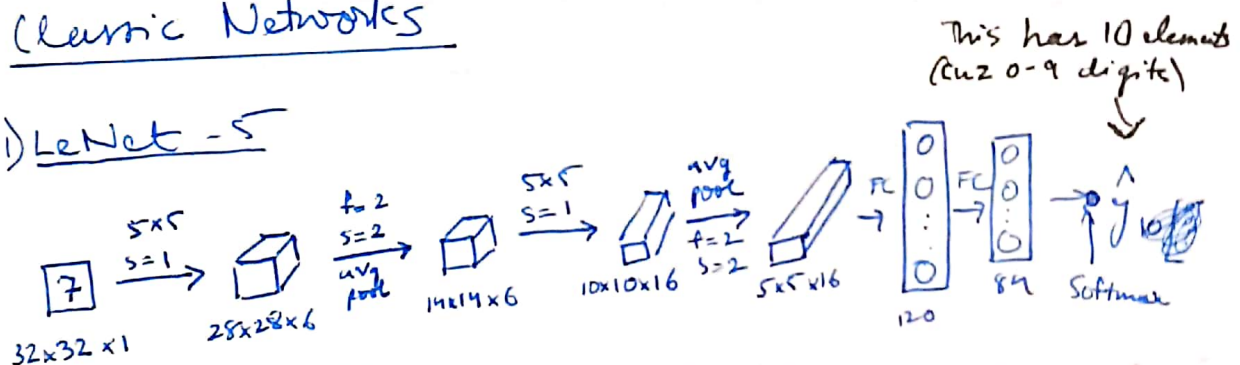  o AlexNet
  o VGG

- ResNet

- Inception

## Classic Networks

### 1) LeNet - 5

This has 10 elements (cuz 0-9 digits)



- 60 k parameters
- $n_H, n_w \downarrow$    $n_c \uparrow$
- conv$\to$pool$\to$ conv$\to$pool $\to$ fc $\to$fc $\to$ output
- Made for grey scale images
- It was written when computers were slow
- Didn't use ReLu

# 2) AlexNet



$227 \times 227 \times 3$ → $11 \times 11$, $s=4$ → $55 \times 55 \times 96$ → MAX-POOL $3 \times 3$, $s=2$ → $27 \times 27 \times 96$ → $5 \times 5$ same → $27 \times 27 \times 256$ → maxpool → $13 \times 13 \times 256$ → $3 \times 3$ same → $3 \times 3$ → $3 \times 3$ → maxpool $3 \times 3$, $s=2$ → = 9216 → FC 4096 → R 4096 → Softmax 1000

→ similar to leNet-5, but much bigger

→ ReLu

→ Multiple GPUs

→ ~~Lots~~ Used Local Response Normalisation (LRN) -
   this isn't used anymore since it sucks
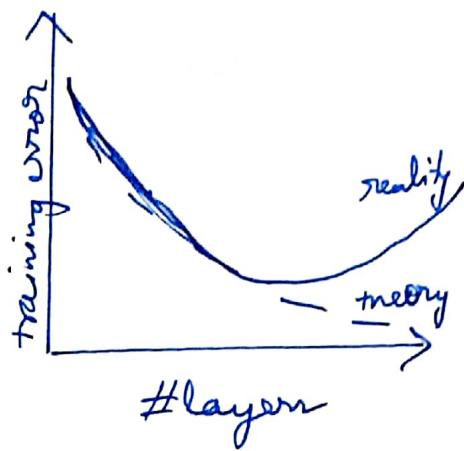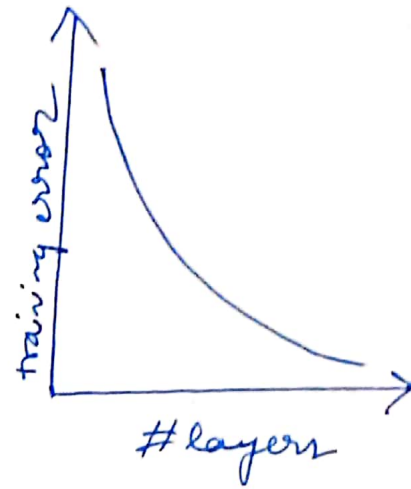
→ 60M parameters

# 3) VGG-16



$224 \times 224 \times 3$ [CONV 64] ×2 → $224 \times 224 \times 64$ → POOL → $112 \times 112 \times 64$ → [CONV 128] ×2

$112 \times 112 \times 128$ → POOL → $56 \times 56 \times 128$ → [CONV 256] ×3 → $56 \times 56 \times 256$

→ POOL → $28 \times 28 \times 256$ → [CONV 512] ×3 → $28 \times 28 \times 512$ → POOL → $14 \times 14 \times 512$

→ [CONV 512] ×3 → $14 \times 14 \times 512$ → POOL → $7 \times 7 \times 512$ → FC 4096 → FC 4096 → Softmax 1000

→ Here all filters, CONV = $3 \times 3$ filters, $s=1$, same

→ MAX-POOL : $2 \times 2$, $s=2$

→ ~138M parameters

→ $n_H$, $n_W$ ↓ by 2 while $n_c$ ↑ - very uniform

# Residual Network

We can train very deep networks without a drop in performance using this.



Plain

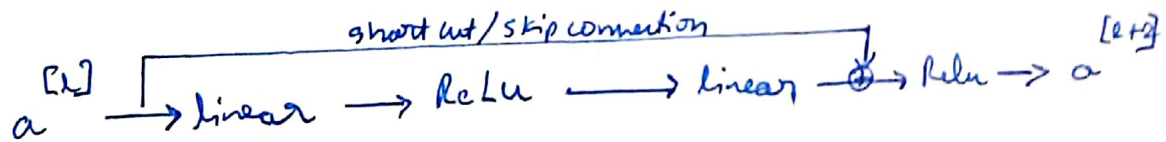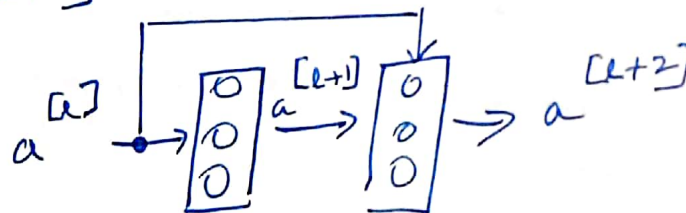ResNet
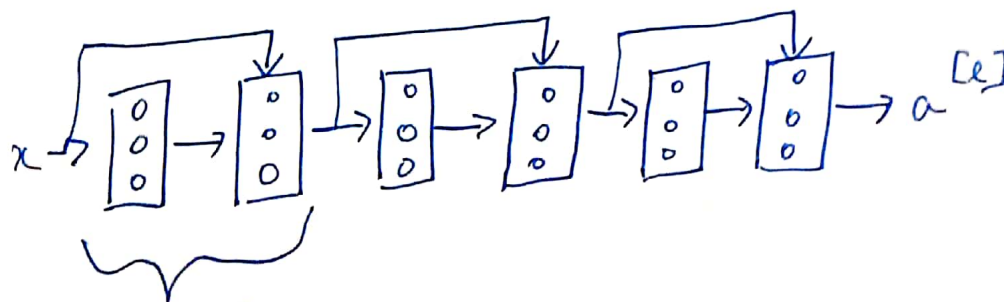
We do this by adding an older $a^{[\ell]}$ to $z^{[\ell+2]}$



short cut/skip connection

$$a^{[\ell]} \longrightarrow \text{linear} \longrightarrow \text{ReLu} \longrightarrow \text{linear} \xrightarrow{\oplus} \text{Relu} \longrightarrow a^{[\ell+2]}$$

$$z^{[\ell+1]} = W^{[\ell+1]} a^{[\ell]} + b^{[\ell+1]}, \quad a^{[\ell+1]} = g\left(z^{[\ell+1]}\right)$$

$$z^{[\ell+2]} = W^{[\ell+2]} a^{[\ell+1]} + b^{[\ell+2]}, \quad a^{[\ell+2]} = g\left(z^{[\ell+2]} + a^{[\ell]}\right)$$
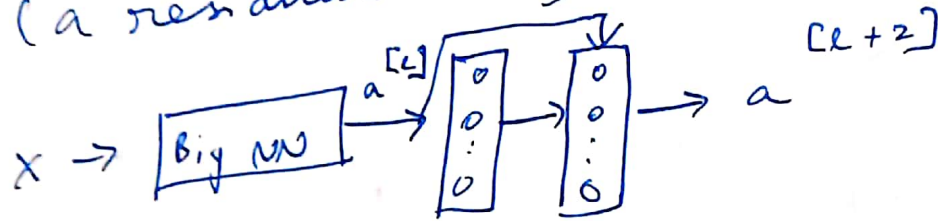


residual block

Why do ResNets work?

Consider a network,

$$x \Longrightarrow \boxed{Big\ NN} \longrightarrow a^{[\ell]}$$

Now suppose we add 2 more layers
(a residual block)

$$x \rightarrow \boxed{Big\ NN} \xrightarrow{a^{[\ell]}} \boxed{} \rightarrow \boxed{} \longrightarrow a^{[\ell+2]}$$

$$a^{[\ell+2]} = g\left(z^{[\ell+2]} + a^{[\ell]}\right)$$

$$= g\left(w^{[\ell+2]} a^{[\ell+1]} + b^{[\ell+2]} + a^{[\ell]}\right)$$

In deeper NN, the layers find it hard
to detect parameters, so suppose

$$w^{[\ell+2]} = 0, \quad b^{[\ell+2]} = 0$$

Then $\quad a^{[\ell+2]} = a^{[\ell]}$

So if no parameters are found, the
performance won't be affected.

# 1x1 convolution

$$\begin{array}{|c|c|c|}\hline 1 & 2 & 3 \\\hline 4 & 5 & 6 \\\hline 7 & 8 & 9 \\\hline\end{array} \quad * \quad \boxed{2} \quad = \quad \begin{array}{|c|c|c|}\hline 2 & 4 & 6 \\\hline 8 & 10 & 12 \\\hline 14 & 16 & 18 \\\hline\end{array}$$

i.e. we multiply all the elements by the 1x1 filter

This seems useless, but when we consider a 3D network, it's different

**Filter 1**



$6 \times 6 \times 32$  $*$  $1 \times 1 \times 32$  $=$  $6 \times 6 \times 1$

Here  $*$ $\Box$ $=$ $\Box$

So we multiply element wise and add it all up $\Box$ $*$ $\Box$ $=$ $\Box$

This is similar to  relu → $\Box$

**Filter 2**



$6 \times 6 \times 32$  $*$  $1 \times 1 \times 32$  $=$  $6 \times 6 \times 1$

stack these up to get $\Box$ $6 \times 6 \times \# filters$

This is useful to reduce no. of filters:



$27 \times 28 \times 192$  $\xrightarrow[\text{CONV } 1 \times 1]{\text{ReLu}}$  $32$  $28 \times 28 \times 32$

192 reduced to 32

# Motivation for inception network

Rather than us choosing the filter size, we try all:



28x28x192          s=1 ↑
                   we need
              to add padding so
              it matches 28x28

1x1  64 filters
same  3x3  128 filters
5x5  32 filters
same
MAX POOL  32 filters
same

28x28x256

The problem of computation cost:



28x28x192   CONV 5x5 same, 32   28x28x32
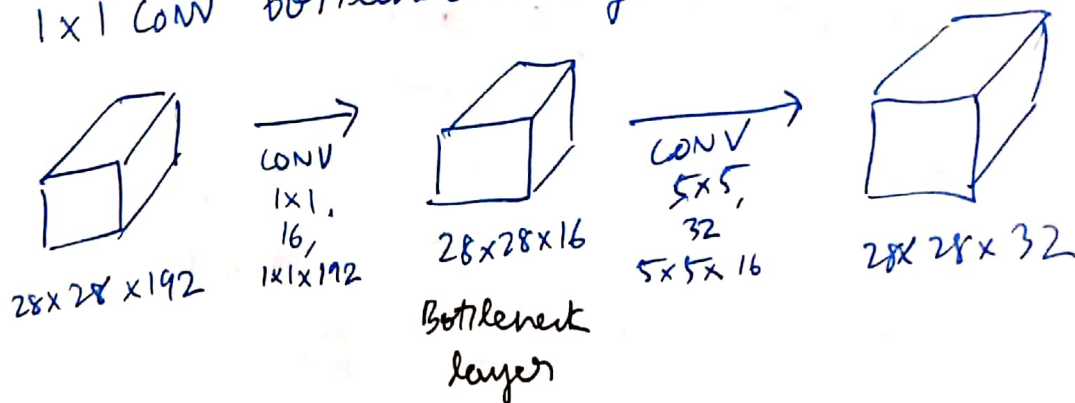
Consider the 32 5x5x192 filters above

The computation will be $28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120M$

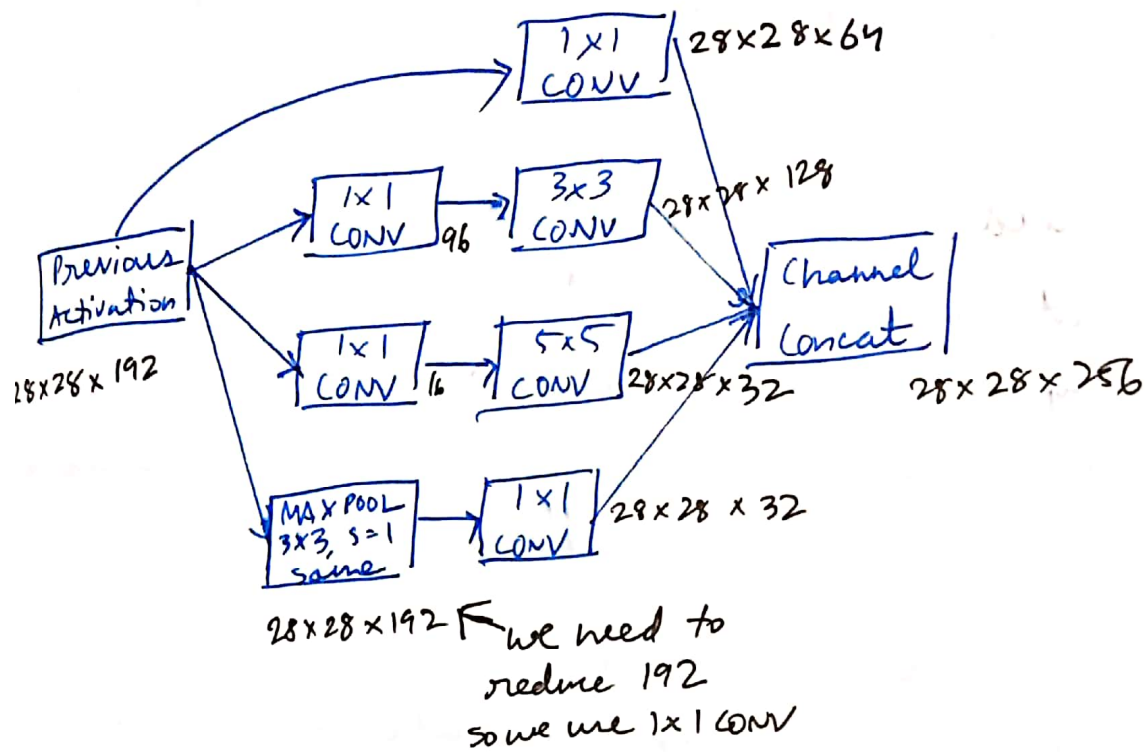However we can reduce this by using a 1x1 CONV bottleneck layer:



28x28x192   CONV 1x1, 16, 1x1x192   28x28x16   CONV 5x5, 32, 5x5x16   28x28x32

Bottleneck layer

$28 \times 28 \times 16 \times 192 = 2.4M$  +  $28 \times 28 \times 32 \times 5 \times 5 \times 16 = 10M$

$= 12.4M$

Hence its reduced.

# Inception module / GoogleNet



We need to reduce 192
so we use 1x1 conv

→ The inception network contains many inception modules connected to each other in series

→ It also contains additional side-branches that each make a prediction, and all these predictions are later connected to a softmax output

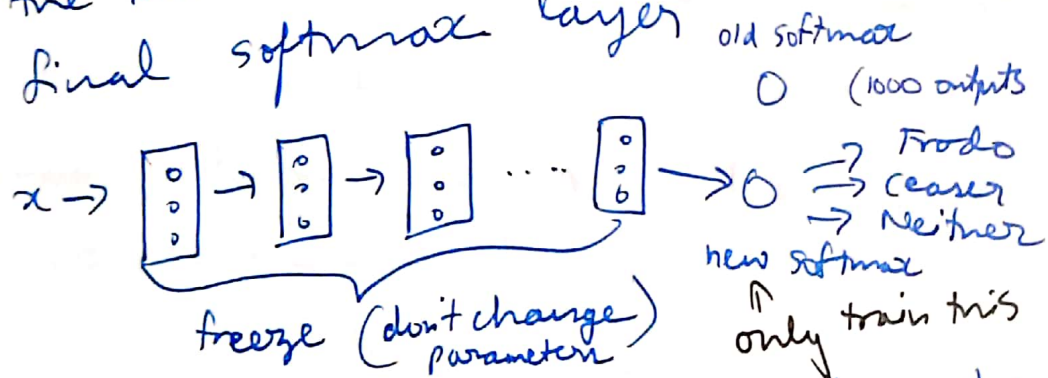→ The names comes from the movie Inception. "We need to go deeper"

# Practical Advice for using ConvNets

1) Using Open-Source implementation:
Rather than studying a paper and building the network from scratch, check out github to find the open source implementation. In most cases, the authors of the paper would have made their code open-source.

2) Transfer Learning:
Let's say you need to make a dog detector (Frodo, Ceaser, Neither). You can download an existing pre-trained model from the internet, and only replace the final softmax layer.

old softmax

$x \rightarrow$ [ : ] $\rightarrow$ [ : ] $\rightarrow$ [ : ] $\cdots$ [ : ] $\rightarrow$ O

O  (1000 outputs

$\rightarrow$ Frodo
$\rightarrow$ Ceaser
$\rightarrow$ Neither

new softmax

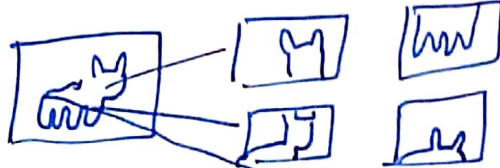freeze (don't change parameter)    only train this

If we have more data, we can reduce the no of layers we freeze

### 3) Data Augmentation:

We can use this to increase our data.

- Mirroring  → 

- Random cropping 

However this may not be good since what if it crops to something like  But usually it works fine

- Rotation
- Shearing }  Used less
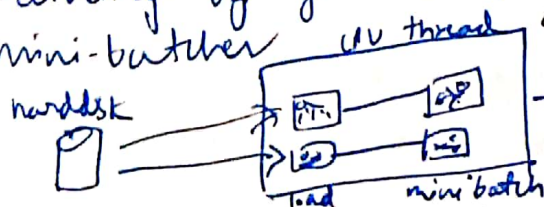- Local Warping

- Color shifting — we add/minus RGB values by small values

This makes your model more robust to colour changes (sunlight, etc)
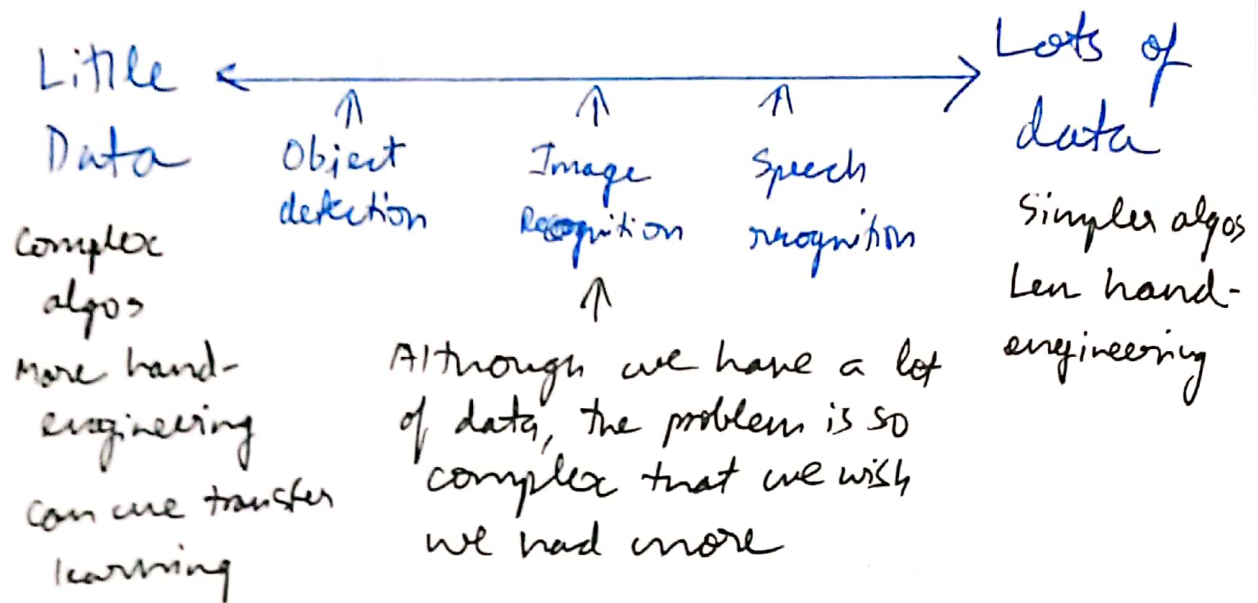
   — PCA color augmentation: If the image has more R & B values, it'll subtract a lot to R & B than G. It's used in Alexnet paper.

We can implement distortions during training by generating a stream of distorted mini-batches  These 2 happen parallely

harddisk → Training

# 4) State of Computer Vision

Little ⟵————————————————⟶ Lots of
Data    ↑         ↑         ↑      data

Complex   Object   Image    Speech    Simpler algos
algos    detection Recognition recognition   Less hand-
more hand-                          engineering
engineering         ↑
              Although we have a lot
can use transfer   of data, the problem is so
learning           complex that we wish
                  we had more

Two sources of knowledge:

→ Labeled data $(x,y)$

→ Hand engineering features, network
   architecture, other components

Tips for doing well on benchmarks/
winning competitions

→ Ensembling: Train several networks
   independently (simultaneously) and
   average their outputs

→ Multi-crop at test time: Run classifier on
   multiple versions of test images and
   average results (10·crop method)

Use open source code:
   → Use architecture of networks published in the
     literature
   → Use open source implementations if possible
   → Use pretrained models and fine-tune on your
     dataset