

Purline

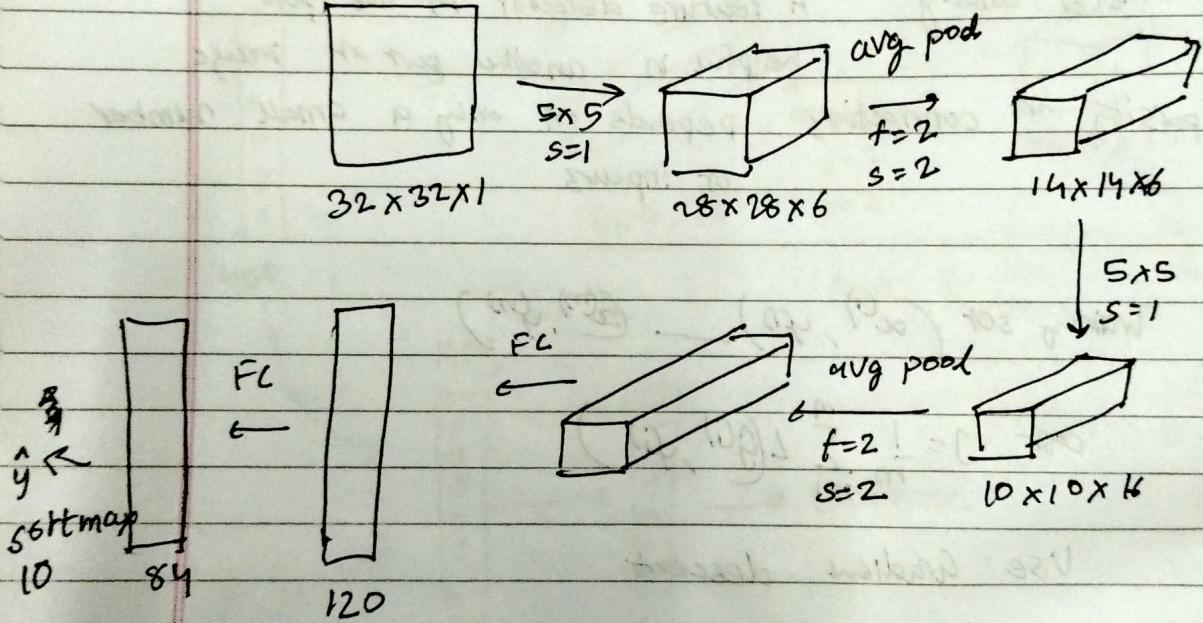
classic networks

- LeNet-5
- AlexNet
- VGG

ResNet (152)

Inception

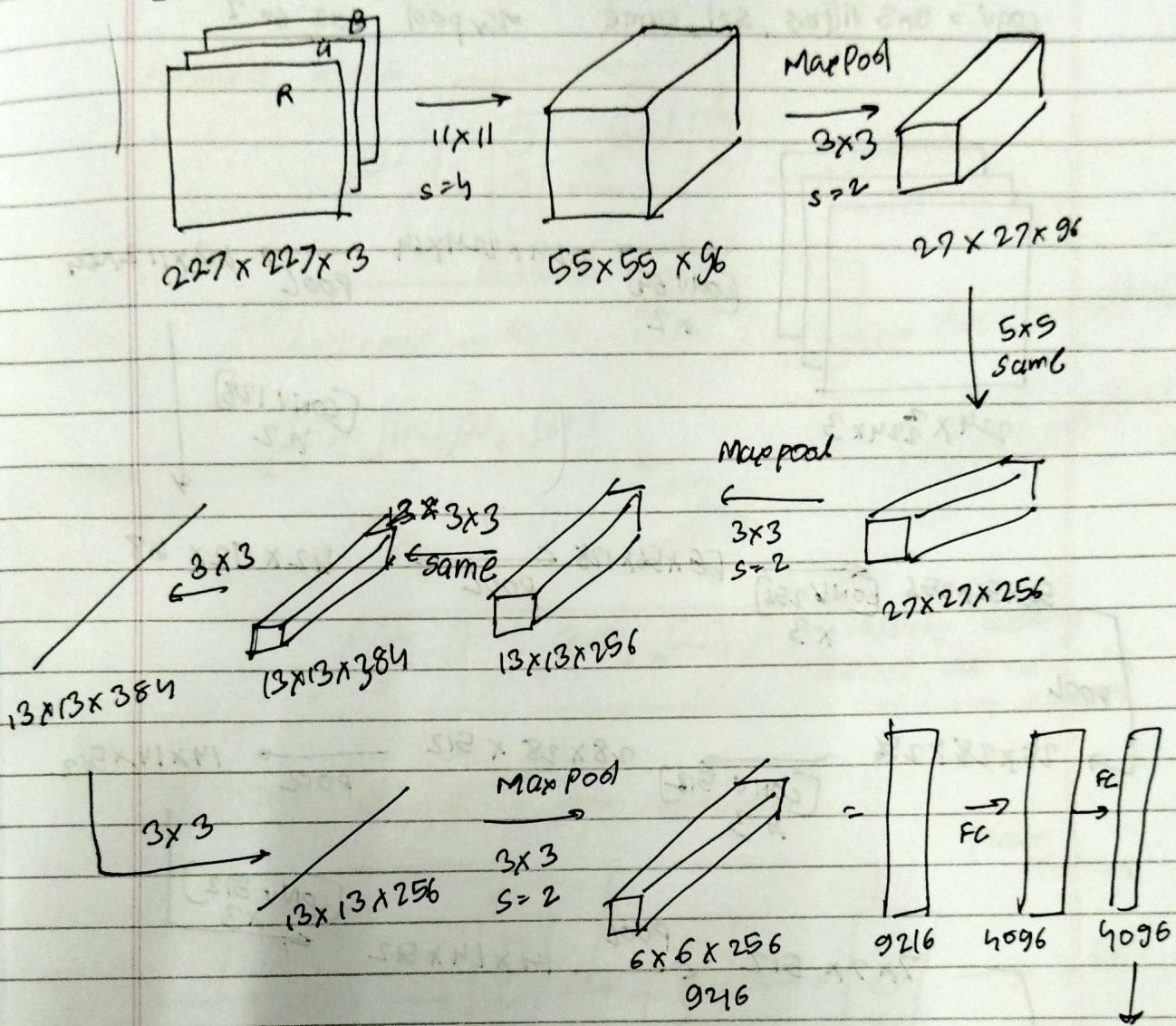
CLASSIC NETWORKS

LeNet-5 $n_h, n_w \downarrow n_c$

conv pool conv pool fc fc output

 $n_h \times n_w \times n_c \quad f \times f \times n_c$

AlexNet



60M parameters

softmax
1000

→ similar to LeNet but much bigger
 $60,000 \rightarrow 60M$ per epoch

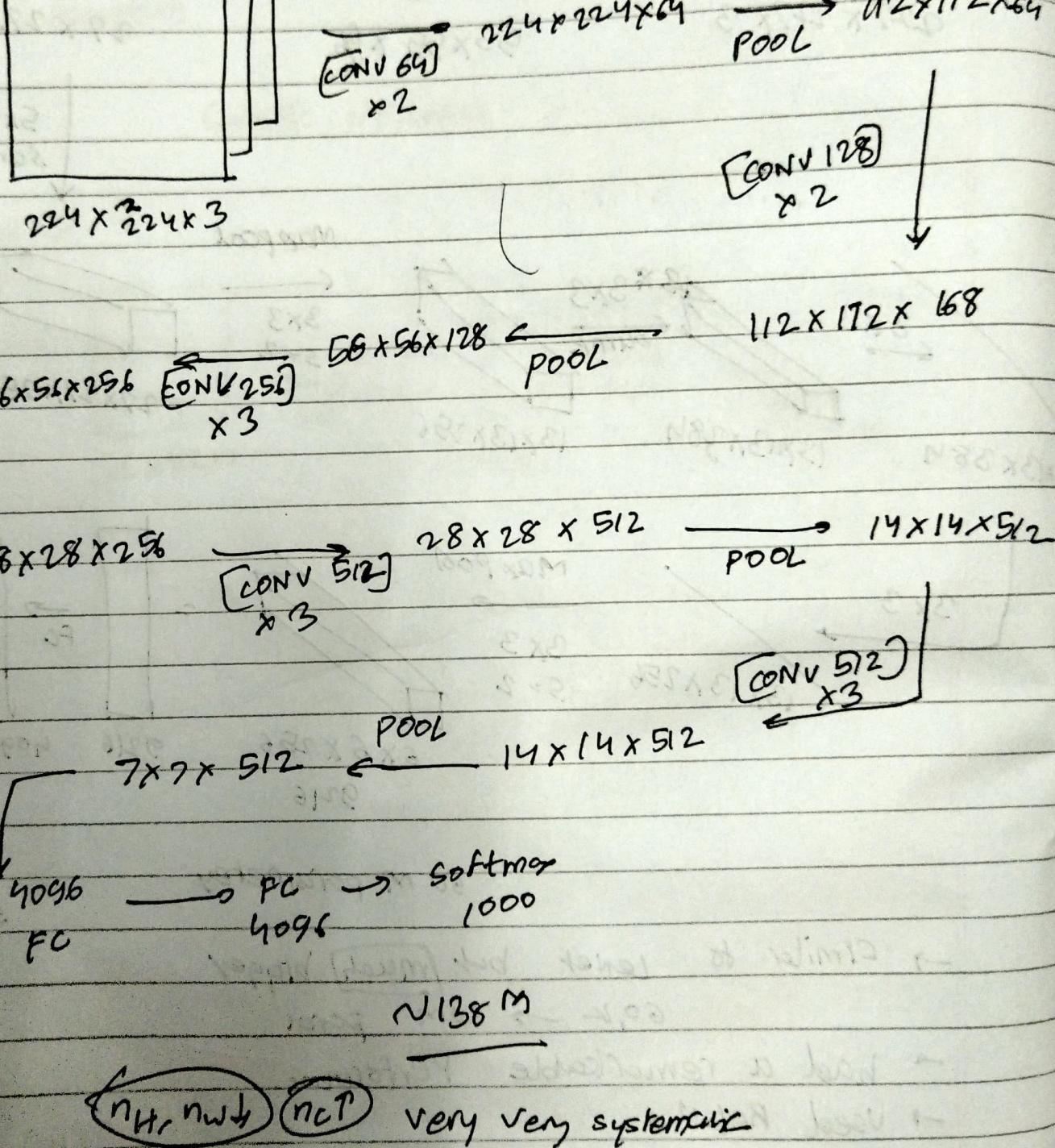
→ had a remarkable Performance

→ used ReLU

→ had a complicated way of training on two GPU's

→ Local Response Normalization

↳ This doesn't help that much

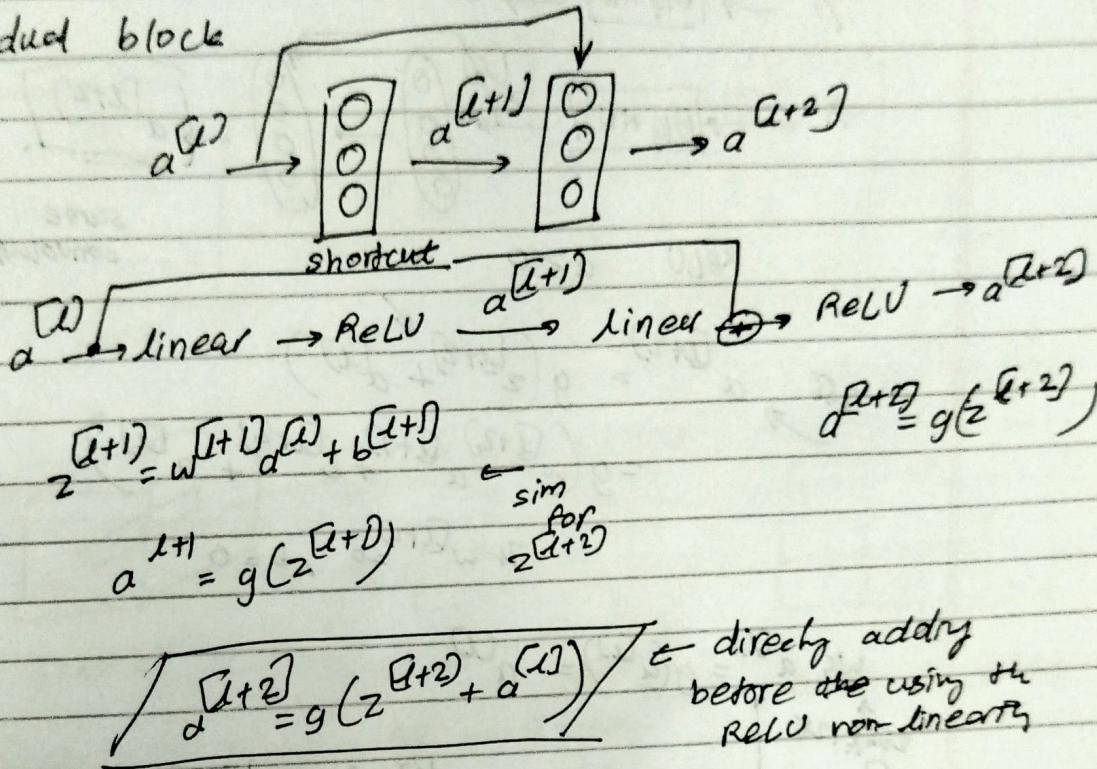


Read AlexNet & VGG 16 Papers

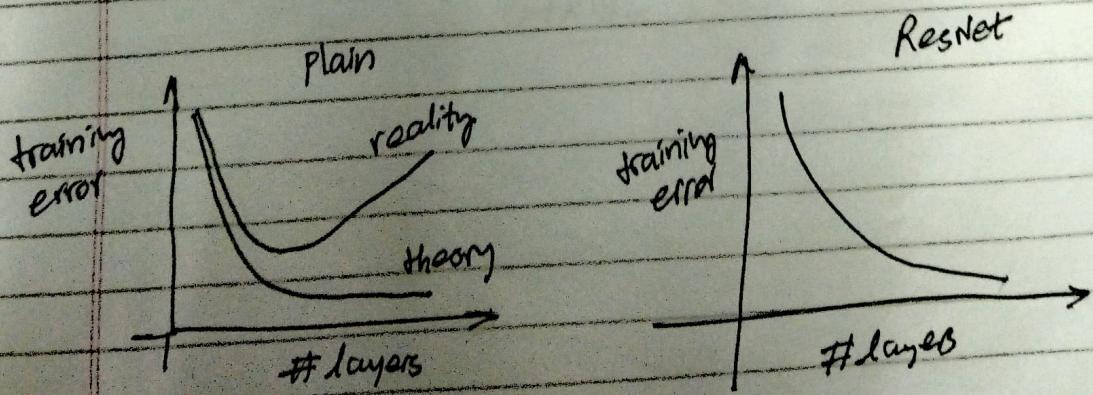
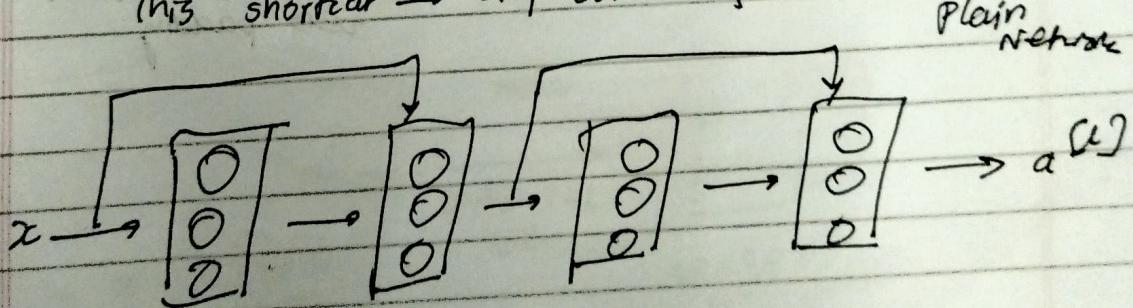
Residual Networks

skip connections

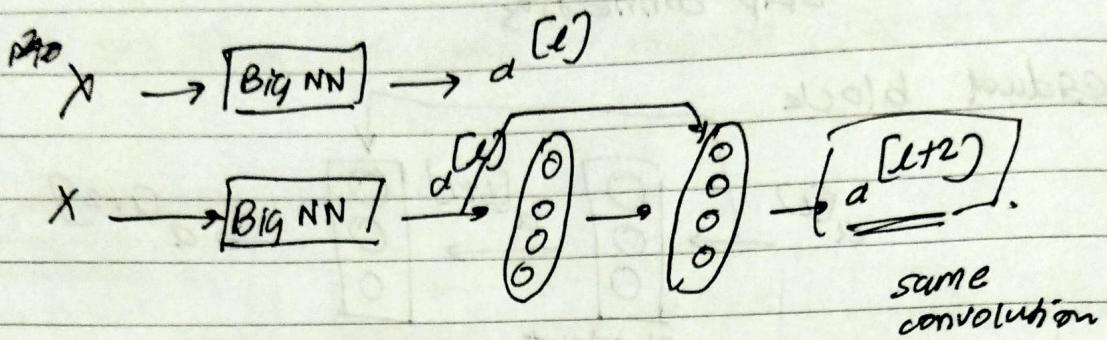
Residual block



This shortcut \rightarrow skip connections



why Resnets work so well?



$$\begin{aligned}
 & \xrightarrow{256} a^{(l+2)} = g(z^{(l+2)} + a^{(l)}) \\
 &= g(w^{(l+2)} a^{(l+1)} + b^{(l+2)} + a^{(l)}) \\
 &\quad \text{if } w^{(l+2)} = 0, b = 0
 \end{aligned}$$

$$\begin{matrix}
 & \uparrow \\
 & 256 \times 128 \\
 \text{R} & | \\
 & 128
 \end{matrix}
 \quad \text{ws } a^{(l)} = g(a^{(l)}) = a^{(l)}$$

ws \rightarrow matrix or parameter
 or \rightarrow that just be a matrix with zero padding

$\times 1$ convolution in NN

what does $1 \times N$ convolutions do?

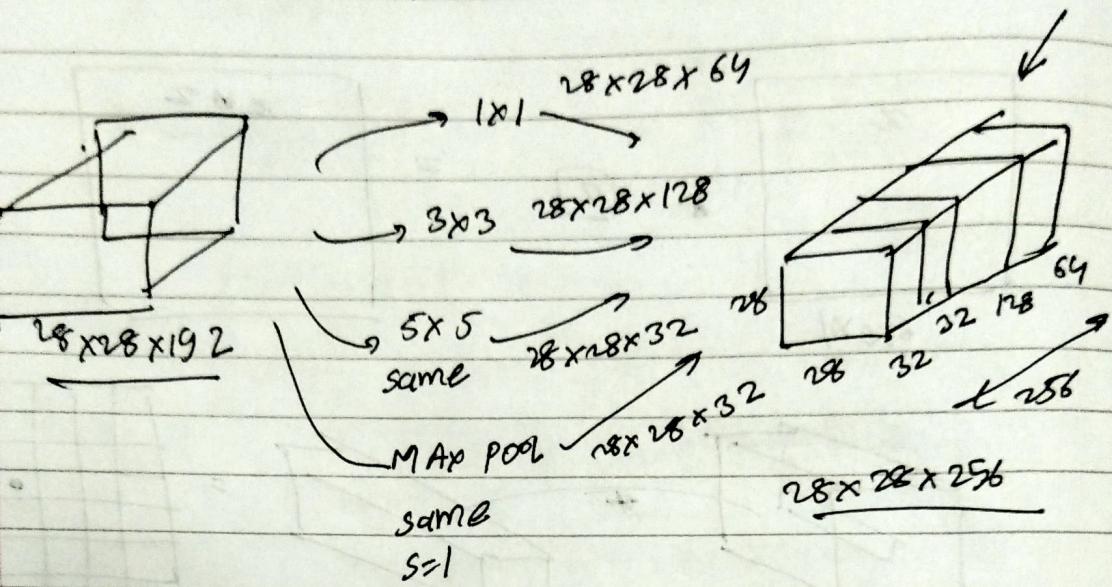
$$\begin{array}{c} x \\ \boxed{6 \times 6 \times 1} \\ \times \boxed{2} = \boxed{x \times 2} \end{array}$$

$$\begin{array}{c} \text{Input: } \boxed{6 \times 6 \times 32} \\ \text{Filter: } \boxed{1 \times 1 \times 32} \\ \text{Output: } \boxed{6 \times 6 \times \# \text{ filters}} \\ \text{Note: } 32 \rightarrow \# \text{ filters} \\ n_c(1+1) \end{array}$$

$$\begin{array}{c} \text{Input: } \boxed{28 \times 28 \times 192} \\ \text{Operation: } \xrightarrow{\text{RELU}} \text{CONV } 1 \times 1 \\ \text{Output: } \boxed{28 \times 28 \times 32} \\ \text{Note: } 192 \rightarrow 32 \end{array}$$

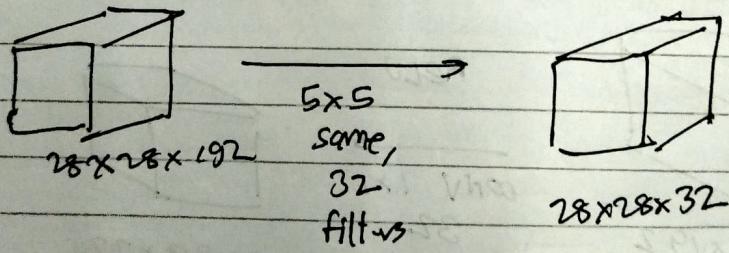
INCEPTION NETWORK

Motivation for inception network



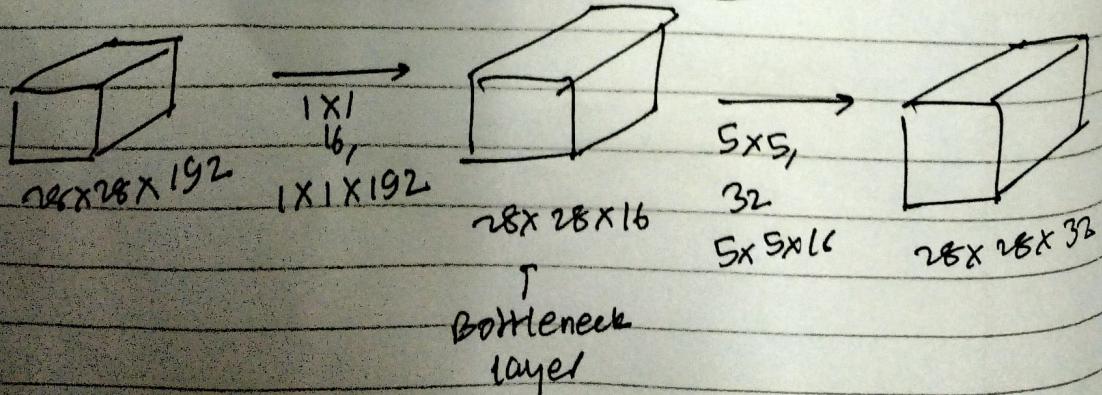
computational cost

5×5 part



$$28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120 \text{ M}$$

using 1×1 conv



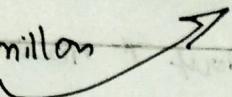
$$28 \times 28 \times 16 \times 192 = 2.4 \text{ million} \quad - \textcircled{1}$$

2nd

$$28 \times 28 \times 32 \times 55 \times 5 \times 16 = 10.0 \text{ million} \quad - \textcircled{2}$$

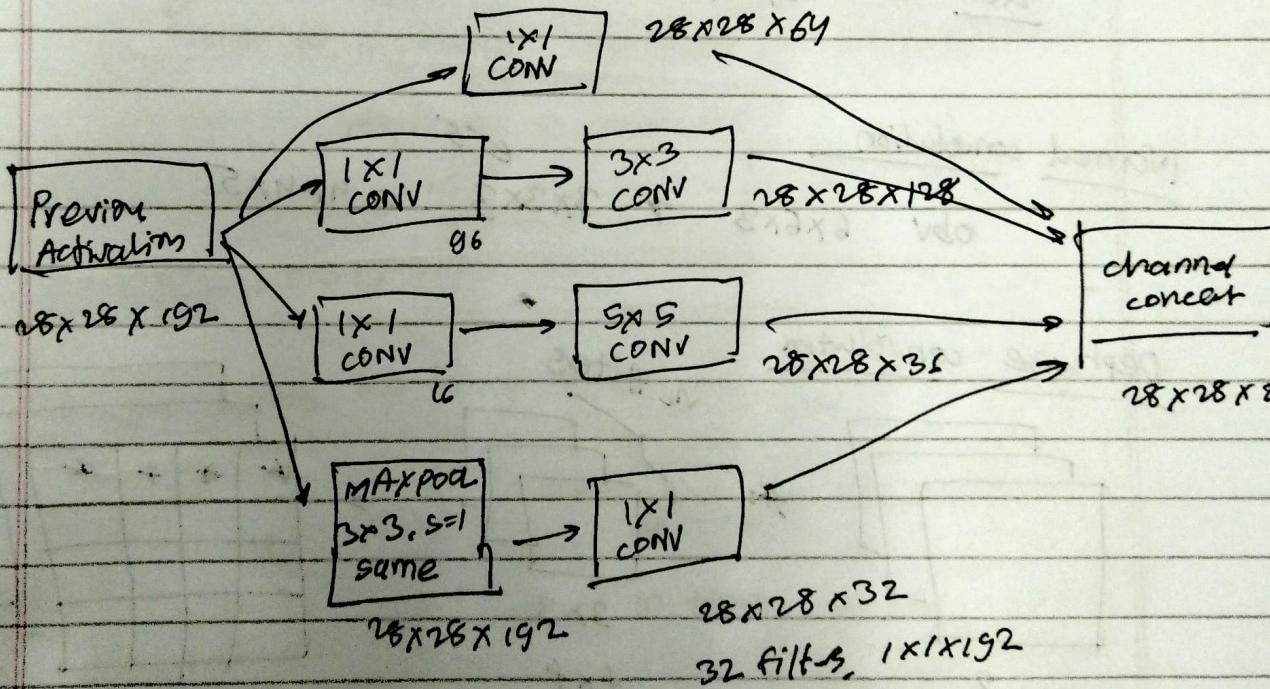
① + ②

12.4 million

120 million 

Does not affect much the performance/accuracy
(Bottlenecks)

Inception Network



Puts a lots of these module together

Read Inception Paper

MobileNet

for computing is low computing cost

→ useful for mobile & embedded system application

$$\begin{array}{c}
 \text{6x6x3} * \underset{\substack{\text{# filter params} \\ \times \# \\ \text{filter positions}}}{\cancel{3 \times 3 \times 3}} = \underset{\substack{\text{no. of filters} \\ \times 5}}{\cancel{4 \times 4 \times 5}} \\
 \text{n} \times \text{n} \times \text{n}_c \qquad \qquad \qquad \qquad \qquad \qquad \text{n}_{out} * \text{n}_{out} \times \text{n}'_c
 \end{array}$$

computational

$$\text{cost} = \# \text{ filter params} \times \# \text{ filter positions} \times \# \text{ of filters}$$

$$\underline{2160} = 3 \times 3 \times 3 \times 4 \times 4 \times 5$$

normal convolution

$$\text{obv } \underset{\substack{\text{6x6x3} \\ \times 5}}{6 \times 6 \times 3} * \underset{\substack{\text{3x3x3} \\ \times 5}}{3 \times 3 \times 3} = \underset{\substack{\text{4x4x5}}}{}{4 \times 4 \times 5}$$

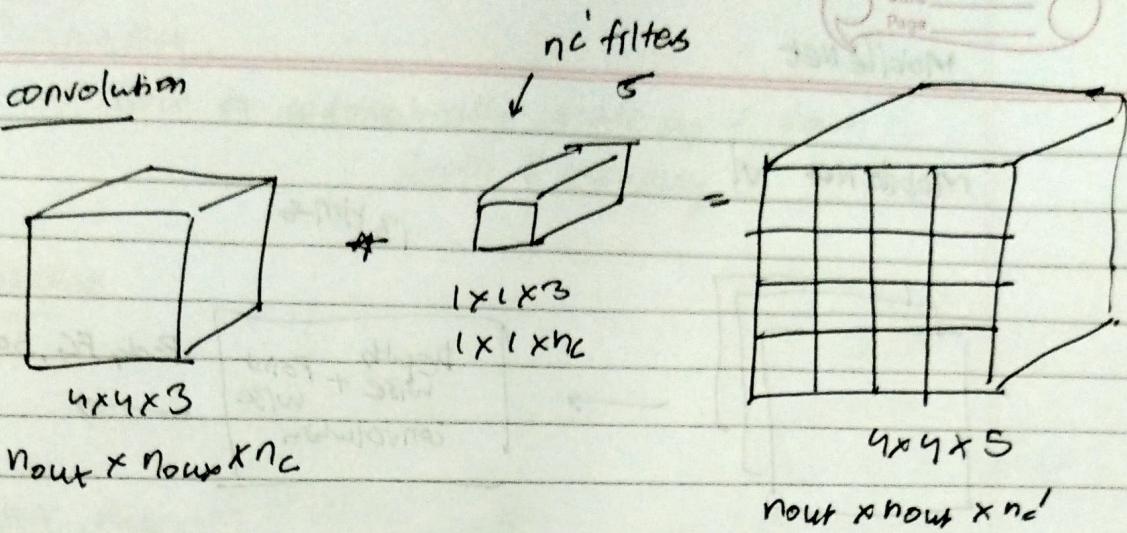
depthwise convolution

$$\begin{array}{c}
 \text{6x6x3} * \underset{\substack{\text{no. filters} \\ \times 3 \times 3 \\ \times f \times f}}{\cancel{\text{3x3x3}}} = \underset{\substack{\text{4x4x3}}}{}{4 \times 4 \times 3} \\
 \text{n} \times \text{n} \times \text{n}_c \qquad \qquad \qquad \qquad \qquad \qquad \text{n}_{out} \times \text{n}_{out} \times \text{n}'_c
 \end{array}$$

$$\text{computational cost} = \# \text{ filter params} \times \# \text{ filter positions} \times \# \text{ filters}$$

$$432 = 3 \times 3 \times 3 \times 4 \times 4 \times 3$$

Pointwise convolution



$$CC = \# \text{filterparams} \times \# \text{filter ps}$$

$$240 = 1 \times 1 \times 3 \times 4 \times 4 \times 5$$

cost of normal convolution

(2160)

cost of depthwise + pointwise convolution

depthwise + pointwise

$$432 + 240 = (672)$$

$$\frac{672}{2160} = 0.31$$

→ This comes out to be $= \left[\frac{1}{nc} + \frac{1}{f^2} \right]$

$$= \frac{1}{5} + \frac{1}{9} = 0.31$$

$\frac{1}{512} + \frac{1}{3^2} \approx 10 \text{ times cheaper}$

⇒ Depthwise (dimensions)

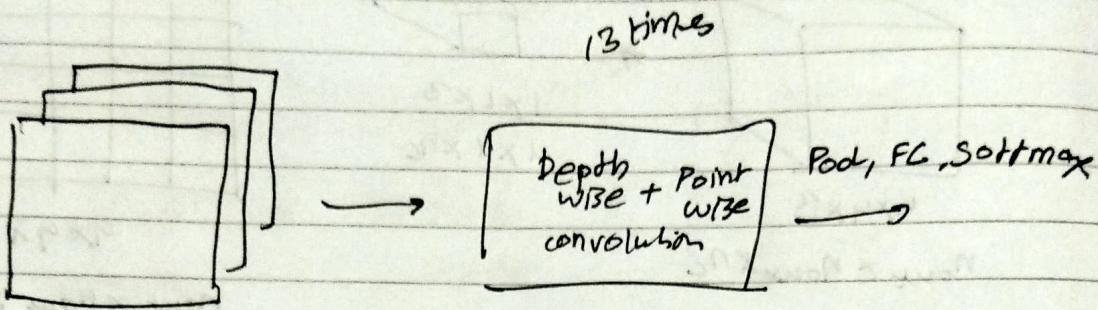
↳ $6 \times 6 \times nc \star 3 \times 3 \times nc = 4 \times 4 \times nc$

⇒ Pointwise

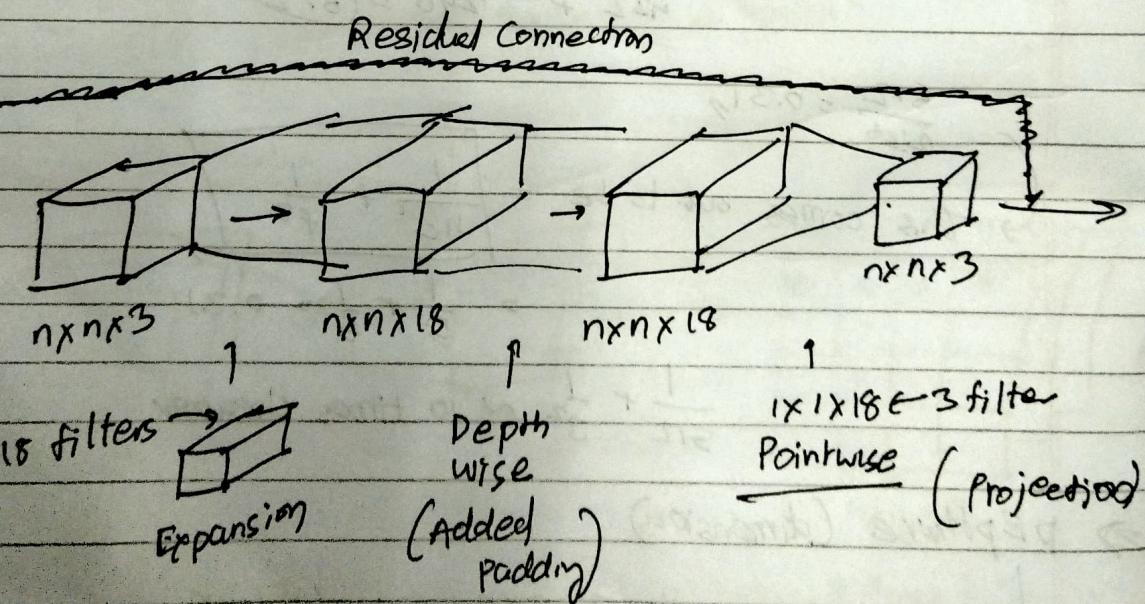
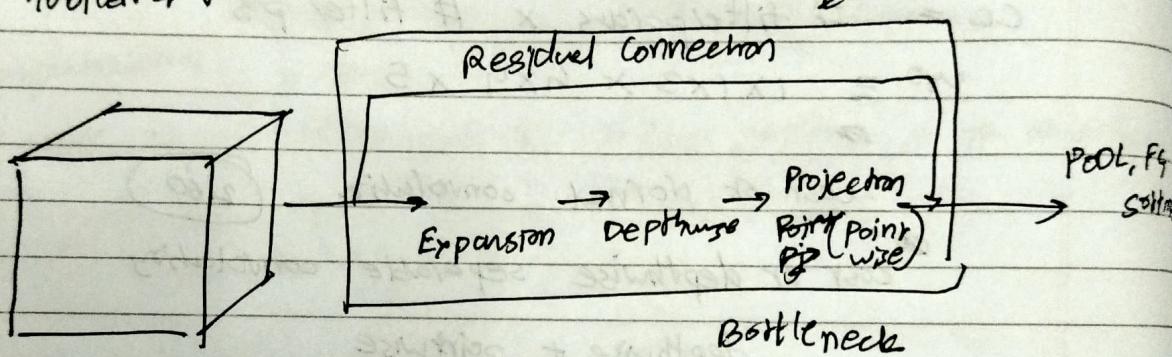
↳ $n \times n \times nc \star 1 \times 1 \times nc = 4 \times 4 \times nc'$

MobileNet

MobileNet v1



MobileNet v2



Why Bottlenecks?

→ Between two uses heavy computation
243 243

→ Keeping the amount of memory and computation as small as possible

EfficientNet

→ How to automatically scale up & down
depth & accuracy?

Baseline

→ (r, d, w)
resource
depth
wider

Higher Resolution / Deeper / wider

what's good choice
according to the
resources we have?

what to increase & how much?

}] How to adapt for a particular device?

look at one of the open source implementation of
efficientNet

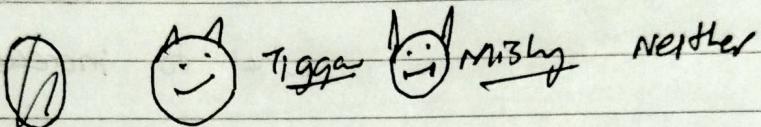
Practical Advices

- Using open-source implementations

Github

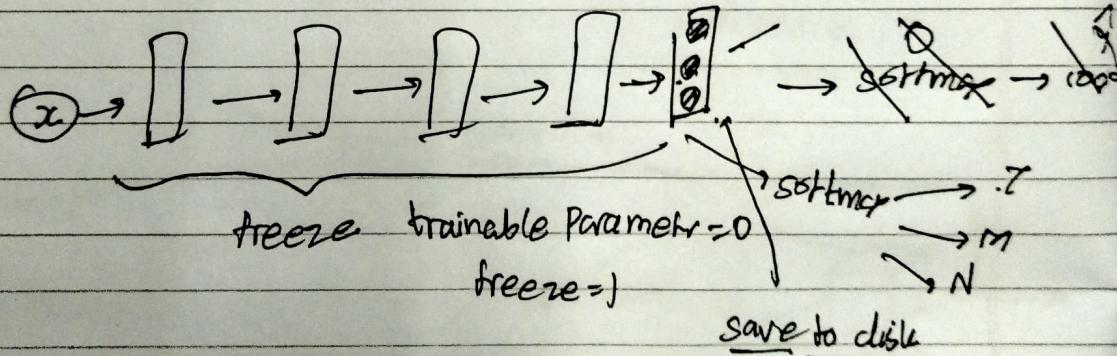
git clone url
to download on local harddisk

Transfer learning



Download open source implementation

↳ code + weights



Precompute till last layer
maps its to last layer using your training set

If you have larger dataset then
freeze only some layers (first few)
and use large few layers to train

for lots of data

You can train the whole unit data
so just use weights as initialisation

data augmentation

some common techniques

Mirror

cropping (Be aware)

Rotation

shearing

local warping

^{PS} more complicated don't worry

color shifting

→ Make training data robust

→ By making slight changes in the colors of images

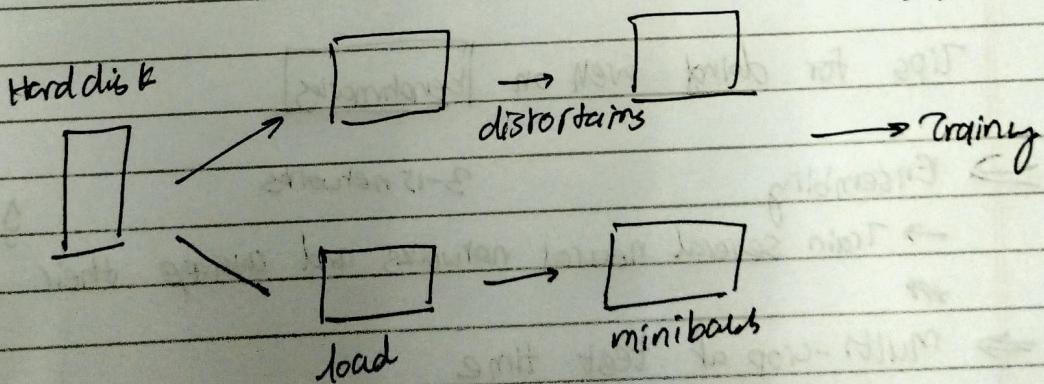
R G B

Distort the channel a little bit and model should work good even then

or try to make the RGB gradient a bit centered before implementing the training

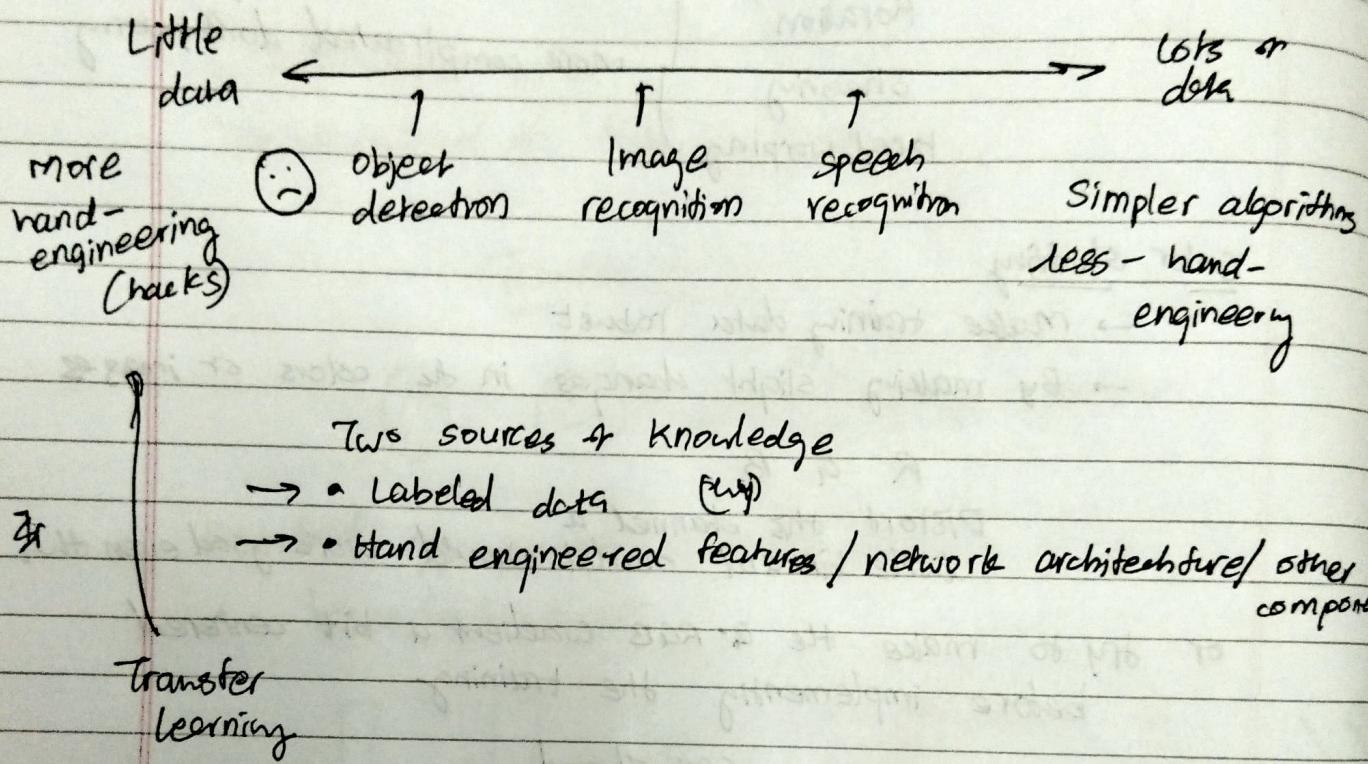
CPU thread

CPU/CPU



State of computer vision

Data vs hand engineering



Tips for doing well on benchmarks

⇒ Ensembling

3-15 networks

1, 2% better

→ Train several neural networks and average their outputs

#

⇒ Multi-crop at test time

→ Run classifier on multiple versions of test images and average outputs