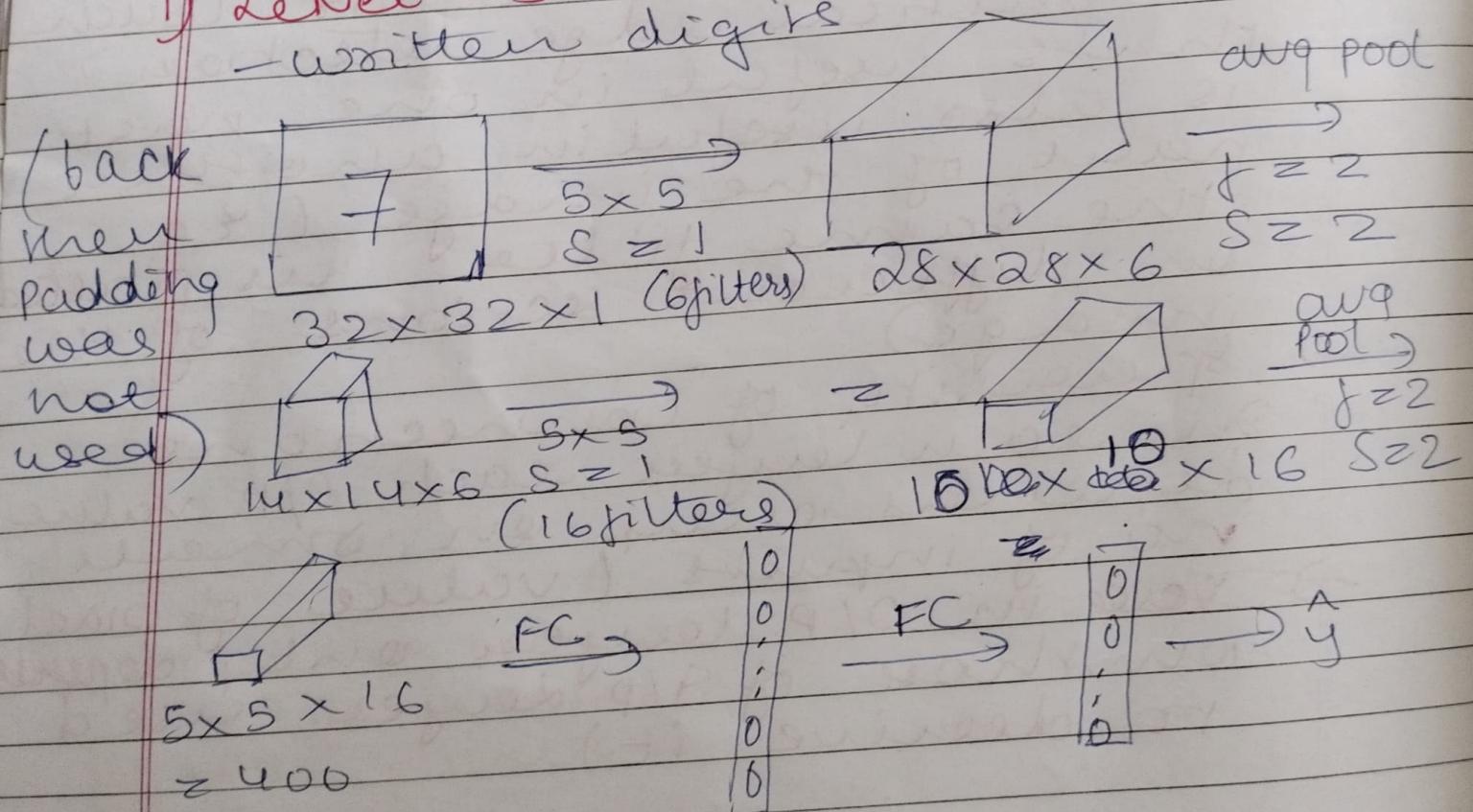


Case Studies

- Classic Networks:
 - LeNet - 5 → AlexNet → VGG
 - ResNet: (152 layer network)
 - Inception

Classic Networks

- 1) LeNet - 5 → Recognized hand written digits



- used sigmoid 120
- in a modern version of this layer we use softmax layer with 10 way classification O/P
- earlier it used a diff classifier O/P layer that's useless today
- about 60K parameters (quite small) $\rightarrow n_h, n_w \downarrow \rightarrow n_c \uparrow$

Then pool then come
con → fully con. - fully
output

PAGE No.	
DATE	/ /

Softmax function - converts

a vector of K real nos into
probability distribution of
 K possible outcomes.

→ It is a generalization of
logistic func (sigmoid) to multiple
dimensions used in multinomial
logistic regression

→ Often used as last activation
function of a neural network
to normalize o/p to a probability
distribution over predicted
o/p classes

→ Takes an input of K real nos
& normalizes it into a probability
distribution of K possibilities
resp. to the exponentials of
input nos

→ Each component after applying
softmax is in interval $(0, 1)$

→ All components add up
to 1 so they are interpreted
as probabilities

$\sigma: \mathbb{R}^K \rightarrow (0, 1)^K$ where $K \geq 1$
defined by formula

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i = 1, \dots, K$$

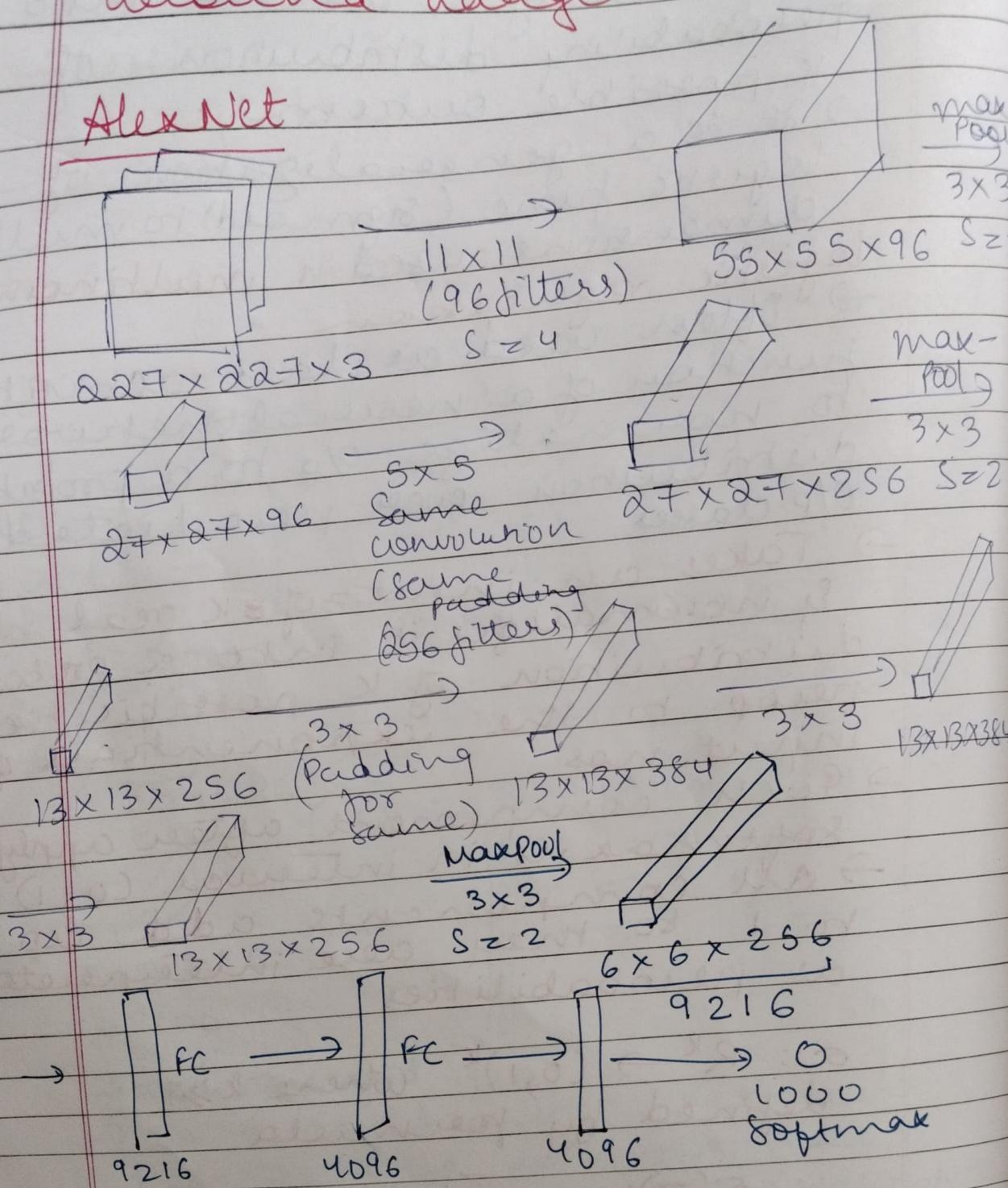
$$z_e = (z_1 \dots z_K) \in \mathbb{R}^K$$

→ Instead of e diff base $0 < b < 1$

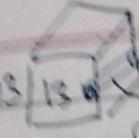
can be used which converges faster
o/p probabilities around smaller val

→ b>0 - Larger O/P probabilities around danger values

AlexNet



- similar to LeNet but much bigger
- Used ReLU activation
- Multiple GPU
- local response normalization layer

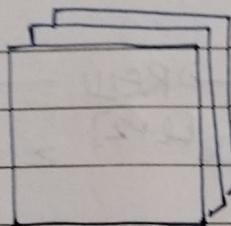

 look at 1 position height
 of width across all channels
 and normalize them
 ↳ Not really useful

VGG-16

$\text{CONV} = 3 \times 3 \text{ filters } S = 1$

always use padding

$\text{MAX-POOL} = 2 \times 2, S = 2$



$224 \times 224 \times 64$
 $(\text{Conv} 64)$
 $\times 2 \uparrow$

$224 \times 224 \times 3$

2 conv layers ( →  → )
 $224 \times 224 \times 64$
 $224 \times 224 \times 64$

POOL

$\rightarrow 112 \times 112 \times 64$ $\xrightarrow{(\text{Conv} 128)}$ $112 \times 112 \times 128$
 $\xrightarrow{\text{POOL}}$ $56 \times 56 \times 128$ $\xrightarrow{(\text{Conv} 256 \times 3)}$ $56 \times 56 \times 256$

$\xrightarrow{\text{POOL}}$ $28 \times 28 \times 256$ $\xrightarrow{(\text{Conv} 512 \times 3)}$ $28 \times 28 \times 512$

$\xrightarrow{\text{POOL}}$ $14 \times 14 \times 512$ $\xrightarrow{(\text{Conv} 512 \times 3)}$ $14 \times 14 \times 512$

$\xrightarrow{\text{POOL}}$ $7 \times 7 \times 512$ $\xrightarrow{FC 4096}$ $\xrightarrow{FC 4096}$ $\xrightarrow{\text{softmax}}$ 1000

Residual Networks (ResNets)

→ Built from a residual block

$$a^{[l]} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} a^{[l+1]} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow a^{[l+2]}$$

$$a^{[l]} \rightarrow \text{linear} \xrightarrow{\text{Relu}} z^{[l+1]} = W^{[l+1]} a^{[l]} + b^{[l+1]} \quad a^{[l+1]} = g(z^{[l+1]})$$

$$\rightarrow \text{linear} \xrightarrow{\text{ReLU}} z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]} \quad a^{[l+2]} = g(z^{[l+2]})$$

(ResNets)
 { a^[l] → a^[l+1] needs 4 steps
 "shortcut" / "skip connection" (2-lin, 2-relu)

$$\left\{ \begin{array}{l} a^{[l]} \xrightarrow{\text{linear}} \text{ReLU} \xrightarrow{\text{linear}} \text{ReLU} \xrightarrow{\oplus} a^{[l+2]} \\ \therefore a^{[l+2]} = g(z^{[l+2]} + a^{[l]}) \end{array} \right.$$

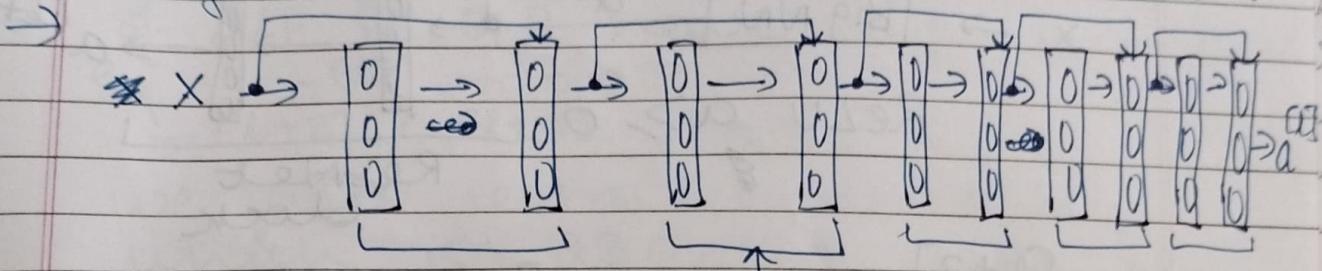
→ Rather than needing to follow the long path the info from a^[l] can follow shortcut to go much deeper into neural networks

ie instead of a^[l+2] = g(z^[l+2]) we have g(z^[l+2] + a^[l])

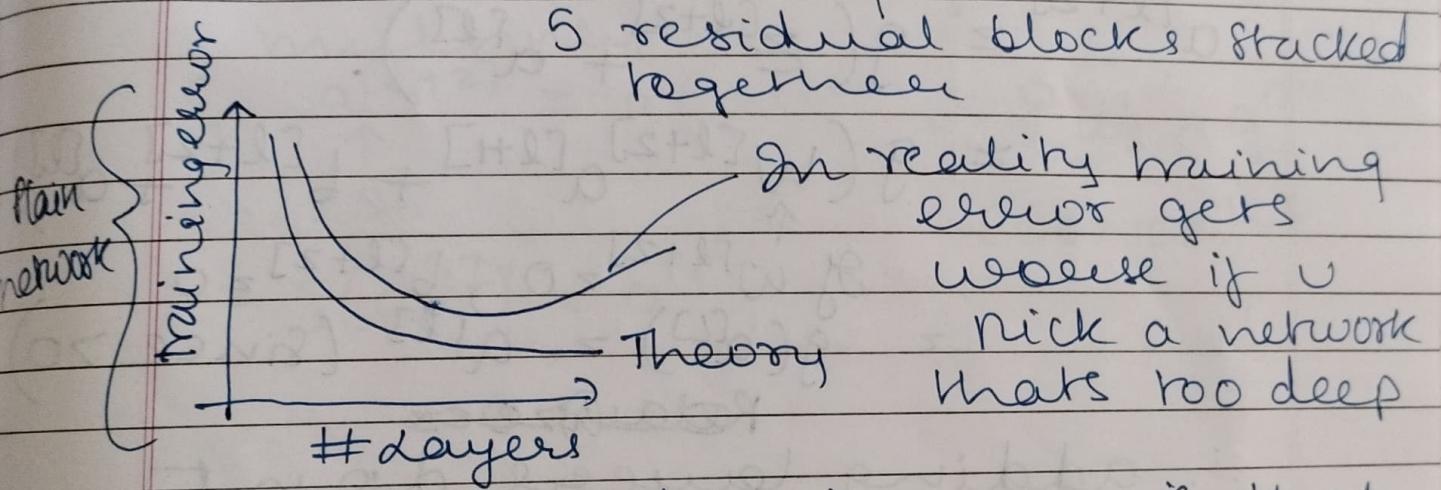
→ addition of a^[l] makes it a residual block

→ Using residual block you can train much deeper neural networks

→ The way to build ResNet is by taking many of these blocks like those and stacking them together to form a deep network

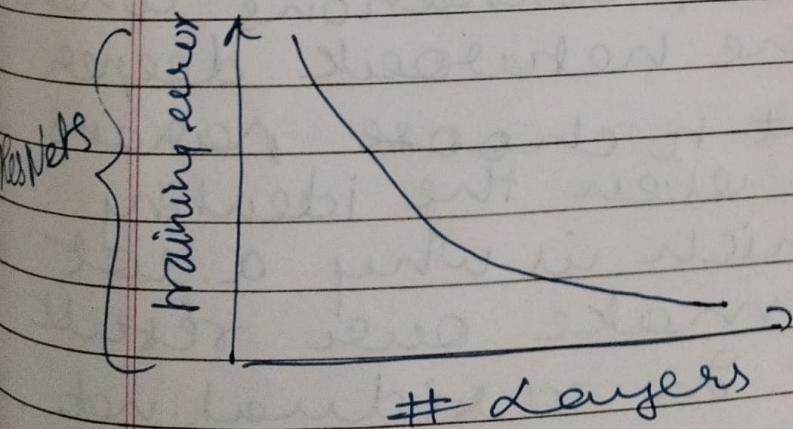


5 residual blocks stacked together

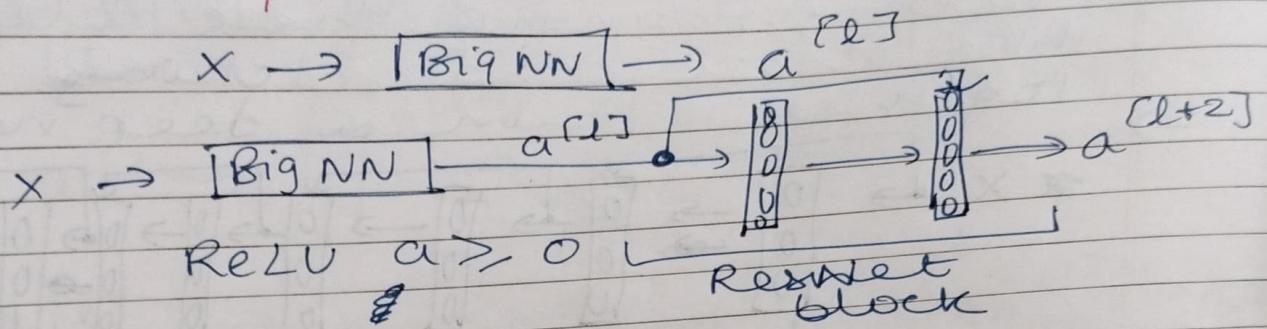


In reality training error gets worse if u pick a network that's too deep

→ In ResNet what happens is that even as no of layers get deeper u can have reperformance of training error keep going down



Why Do residual Networks work?



$$\begin{aligned}
 a^{[l+2]} &= g(z^{[l+2]} + a^{[l]}) \\
 &= g(w^{[l+2]} a^{[l+1]} + b^{[l+2]} + a^{[l]}) \\
 &\quad \text{if } w^{[l+2]} = 0 \Rightarrow b^{[l+2]} = 0 \\
 &\Rightarrow g(a^{[l]}) = a^{[l]} \quad (\text{since } > 0) \\
 &\quad \text{ReLU}
 \end{aligned}$$

\therefore adding layers don't cause error increase

\rightarrow what goes wrong in deep learning plain net is without residual or skip connections when we make the network deeper its difficult to choose params that learn even the identity function which is why a lot of layers make our result worse \rightarrow In residual net. its easy for extra layers to learn identity func so it doesn't hurt performance

→ The alpha dimensions of $z^{(l+2)}$ & $a^{(l)}$ must be same so we see lot of same order convolutions.

→ If I/P & O/P have diff dimensions

$$\text{Eq } a^{(l)} = 128$$

$$z^{(l+2)} / a^{(l+2)} = 256$$

we add extra matrix W_s

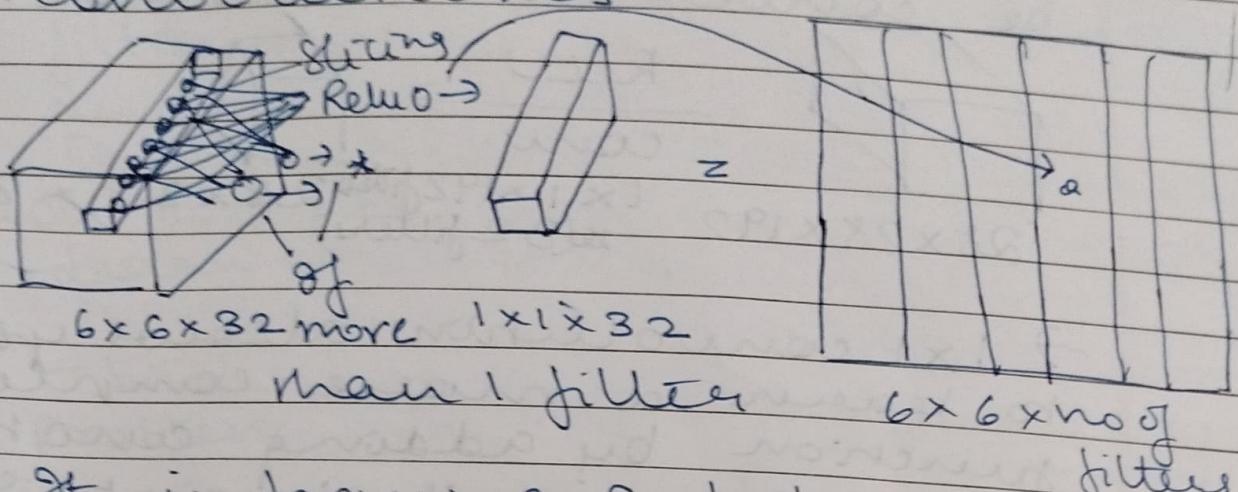
which is 256×128 since if

multiply to $a^{(l)}$

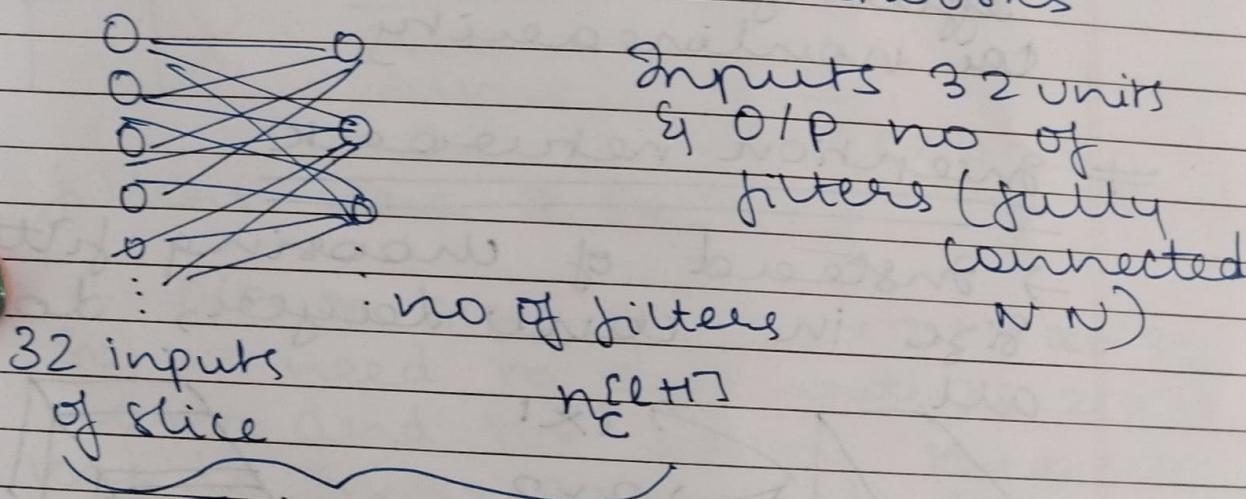
$$\text{Eq } \begin{matrix} a^{(l+2)} \\ \uparrow \\ 256 \end{matrix} = g \left(\begin{matrix} z^{(l+2)} \\ \uparrow \\ 256 \end{matrix} + \begin{matrix} W_s \\ \uparrow \\ 256 \times 128 \end{matrix} \times \begin{matrix} a^{(l)} \\ \uparrow \\ 128 \end{matrix} \right)$$

$$= g(a^{(l)}) = a^{(l)}$$

Network in network 1×1 convolutions



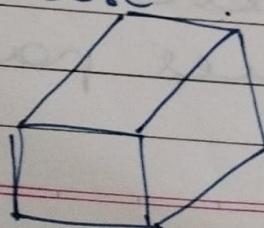
→ It is having a fully connected neural network i.e. for each input slice ie 36 networks



This occurs 6×6 times to give total of $6 \times 6 \times n_c^{l+1}$ outputs

→ This is called a network in network

Eg of use:

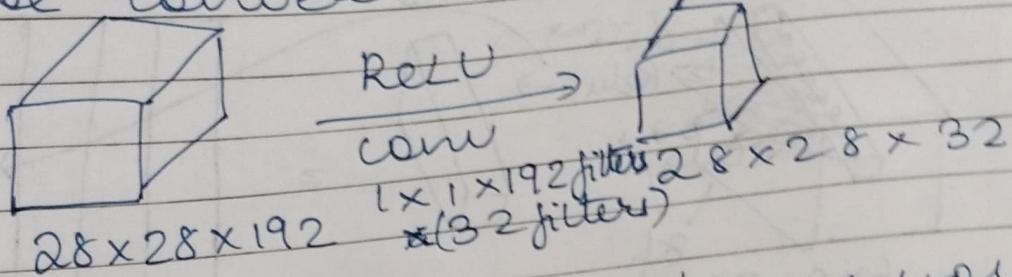


To shrink it we can use Pooling but if no of channels

Same convolution

- O/P size = Input size by padding

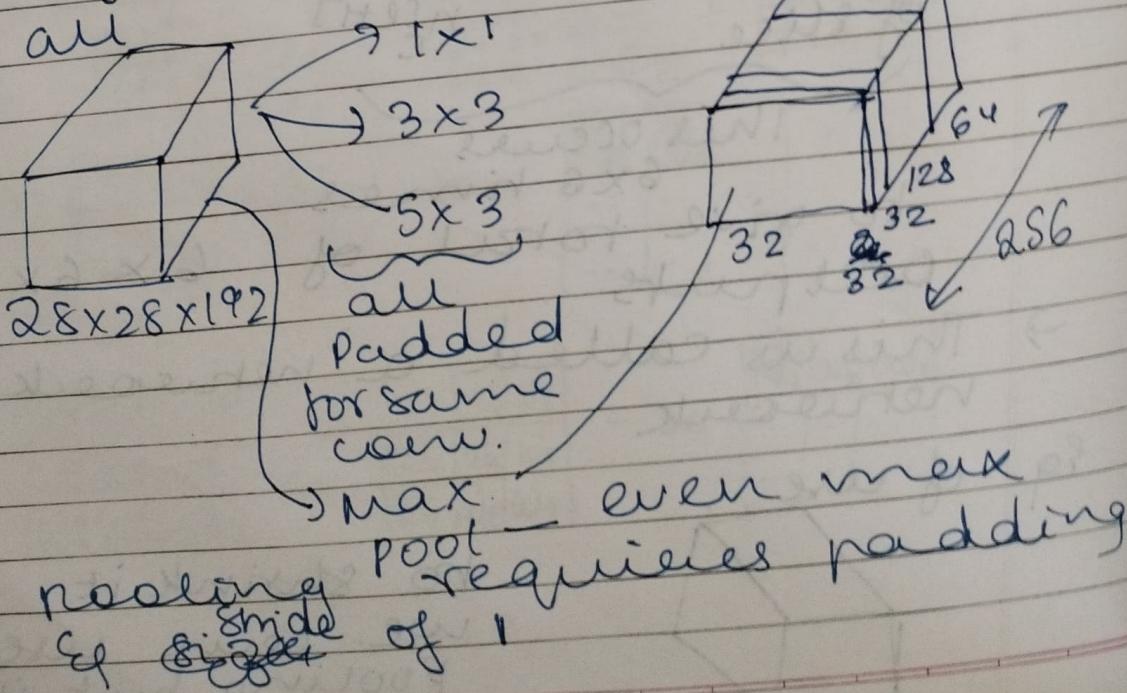
is too big ie it has to be converted to $28 \times 28 \times 32$



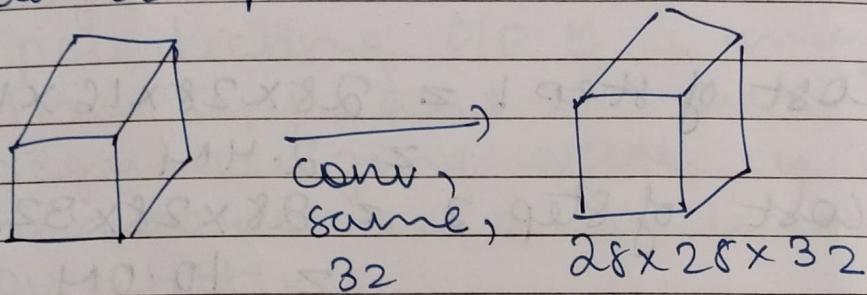
→ 1×1 convolution allows you to learn a more complex function by adding another layer (even if input & output dim. are same). The effect of 1×1 conv is it just has one nonlinearity

Inception network

→ Instead of choosing filter size in conv layers, do them all



- basic idea is instead of you needing to pick one of these filter sizes or pooling you want and committing to that you can do them all, ~~concatenate~~ concatenate outputs and let the network learn whatever parameters it wants to use, whatever combinations of filter sizes it wants
- problem with inception layer is computational cost

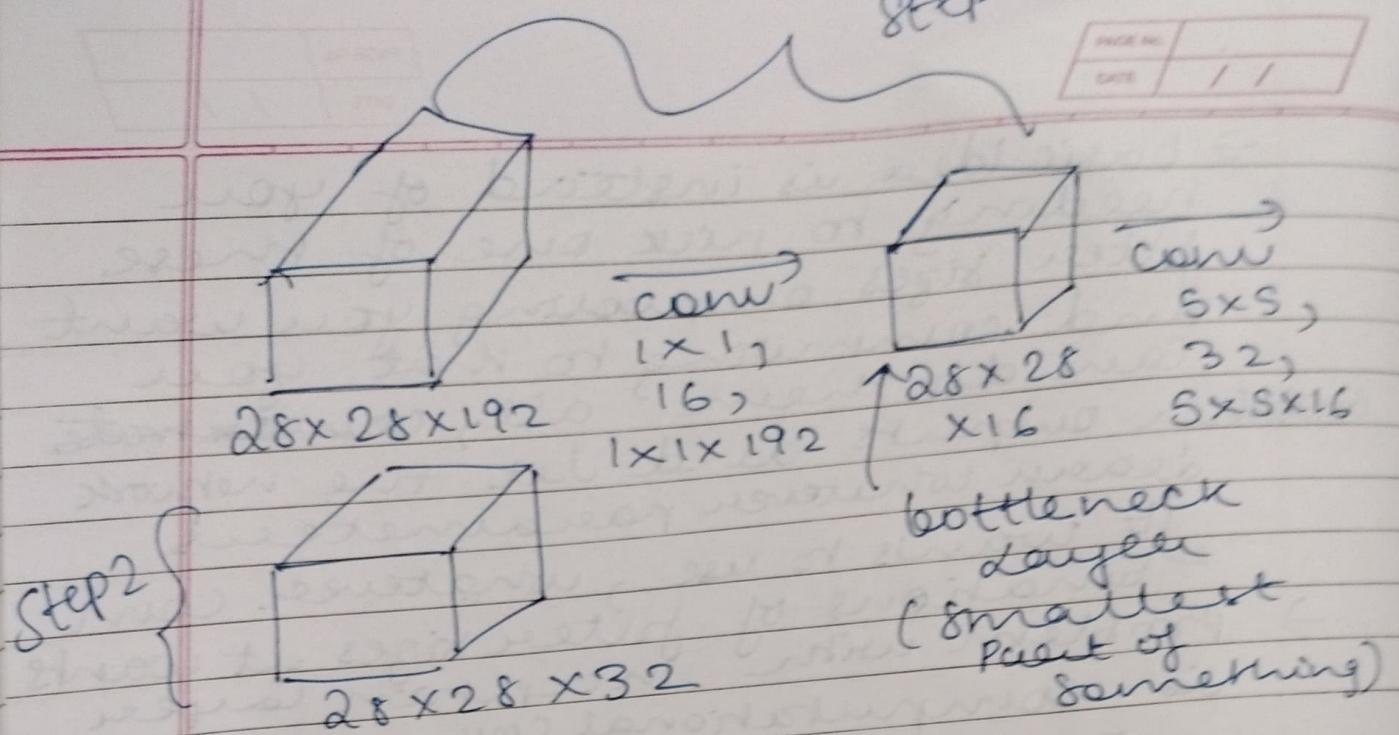


filters are $5 \times 5 \times 192$

\therefore O/P size is $28 \times 28 \times 32$
 so we need to compute $28 \times 28 \times 32$
 nos, and for each nos. we
 have to do $5 \times 5 \times 192$ multiplications
 \therefore Total multiplication = $28 \times 28 \times 32 \times 5 \times 5 \times 192$
 ≈ 120 mill. (This is expensive)

Total multiplies = no of multiplies
 per each O/P. \times no of O/P

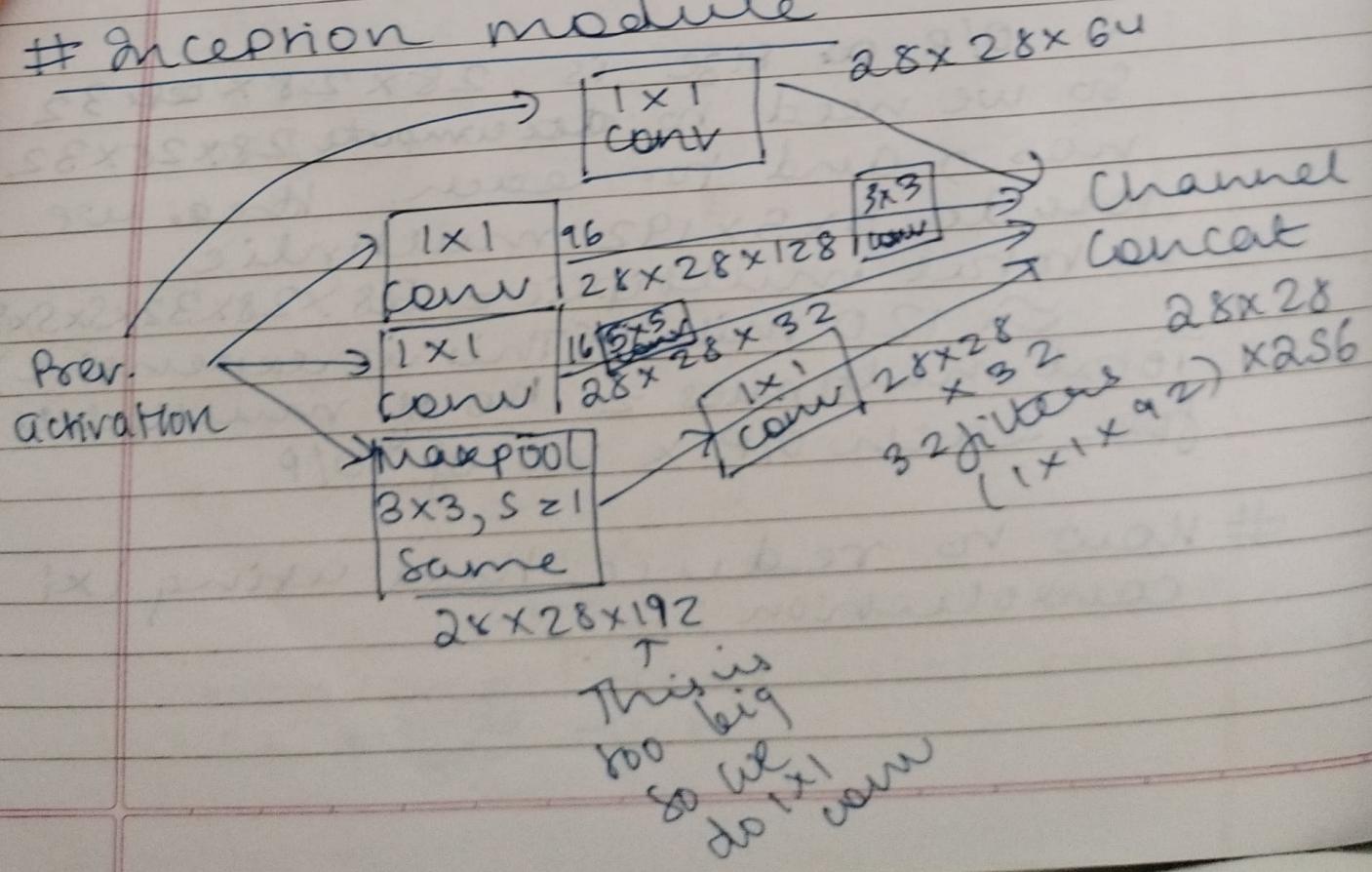
How to reduce this using 1×1 convolution :-



$$\text{cost of step 1} = 28 \times 28 \times 16 \times 1 \times 1 \times 192 \\ = 2.4M$$

$$\text{cost of step 2} = 28 \times 28 \times 32 \times 5 \times 5 \times 16 \\ = 10.0M$$

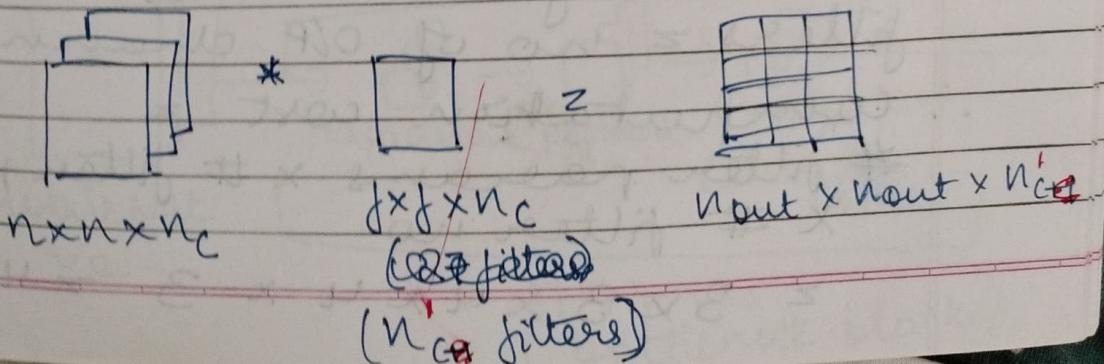
Inception module



- Inception mode network
- Inception mode network combines several such modules together
 - There are sidebranches that try to evaluate ~~the next~~ results using softmax.
 - This ensures that features computed even in the hidden units even at intermediate layers are not too bad for predicting O/P of a image. This appears to have a regularizing effect on inception network and prevents overfitting.
 - Dev. by google (googlenet)

MOBILENET

- lot of N.N are quite computation ally expensive can't run on phones / embedded vision
- Key idea : Normal vs depthwise separable conv.
- Normal conv

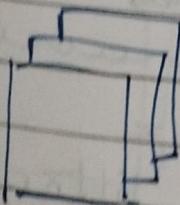


- Inception makes network combines several modules together.
- There are side branches that try to evaluate branches results using softmax. This ensures that features computed even in the deeper units even at it becomes difficult predicting O/P of a image. This appears to have a regularizing effect on inception network and prevents overfitting.
- Dev. by google (Googlenet)

MOBILENET

- lot of N.N are quite computation ally expensive can't run on phones / embedded vision
- **Key idea**: Normal vs depthwise separable conv.

Normal conv



$n \times n \times n_c$

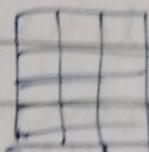
*



$f \times f \times n_c$

(~~one filter~~)

(n_{out} filters)



$n_{out} \times n_{out} \times n_{c1}$

$f \times f \times n_c$

computation cost
 $= \text{filter params} \times \text{filter positions} \times \text{no of filters}$

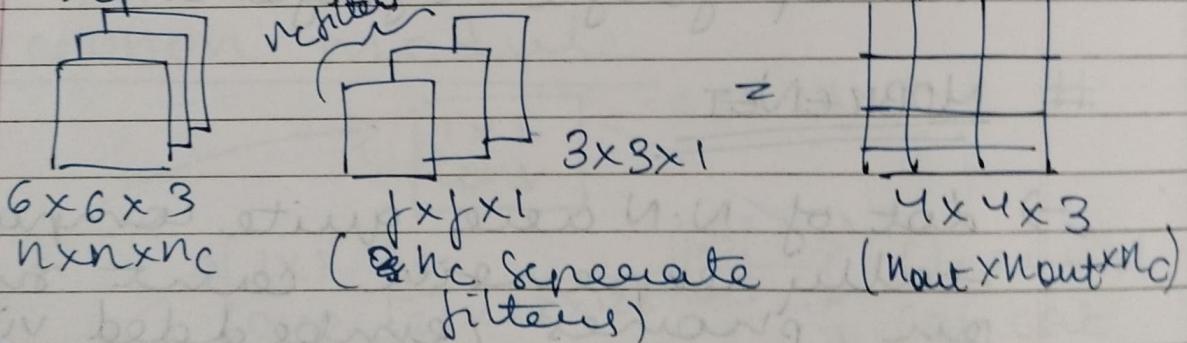
↑
 how many times filter is placed on input img
 ie $n_{out} \times n_{out}$

$$= f \times f \times n_c \times n_{out} \times n_{out} \times n_d$$

Depthwise Separable Convolution

→ Depthwise conv. followed by a pointwise conv.

1) Depthwise conv



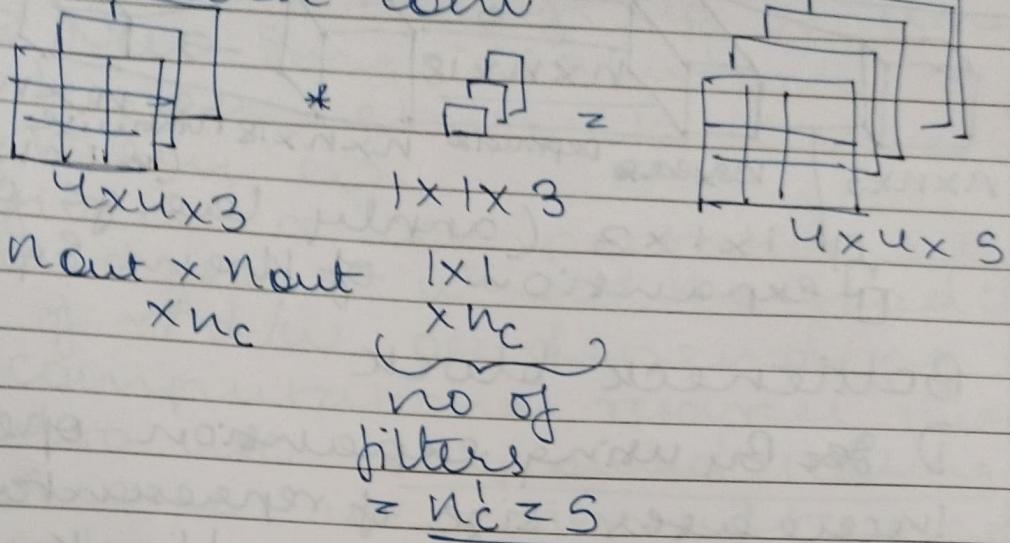
→ 1st filter placed on 1st channel of I/P to give 1st channel of O/P and so on
 \therefore no of I/P channels = no of filters = no of O/P channels

\therefore computation cost =

filter params \times # filter positions
 \times # filter no

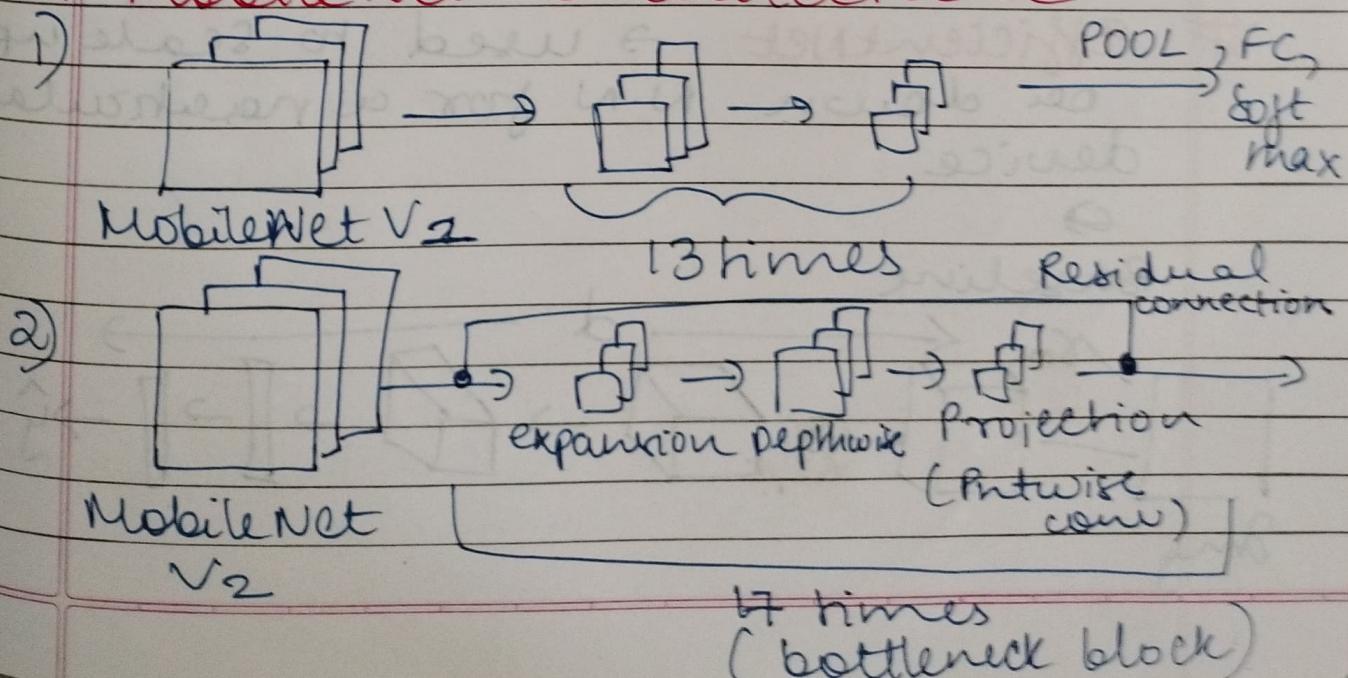
$$\approx 3 \times 3 \times 4 \times 4 \times 3 = 432$$

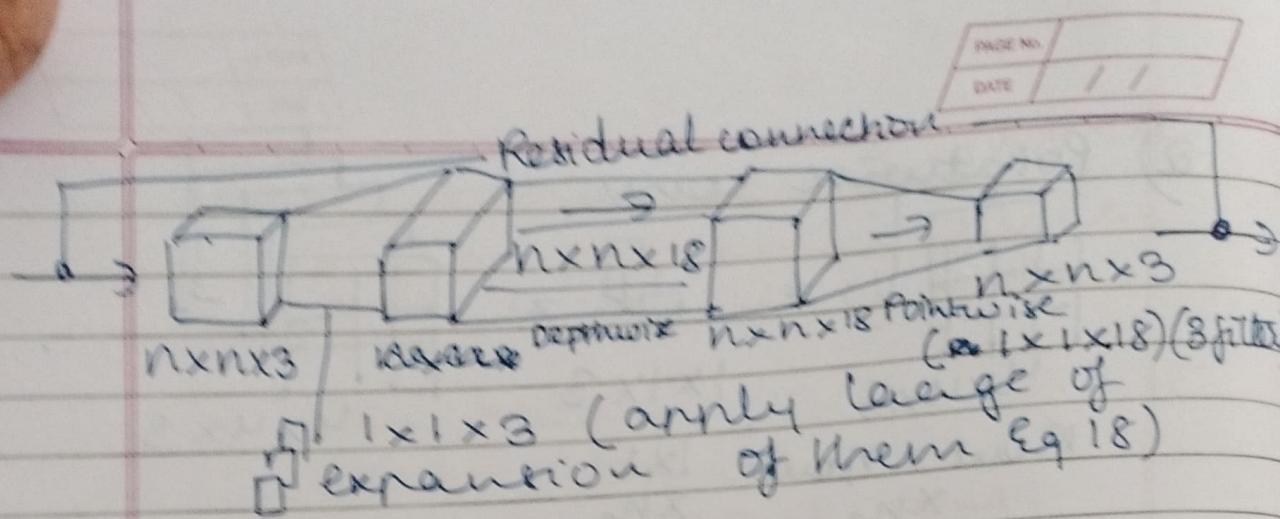
2) Pointwise conv



- cost of normal conv = 2160
- cost of depthwise = depthwise + Pointwise = $432 + 240 = 672$
- cost of depthwise+pointwise = $\frac{1}{n_c} + \frac{1}{f^2}$
 cost of normal conv $= \frac{1}{5} + \frac{1}{9}$
 $= \frac{672}{2160}$

MobileNet Architecture





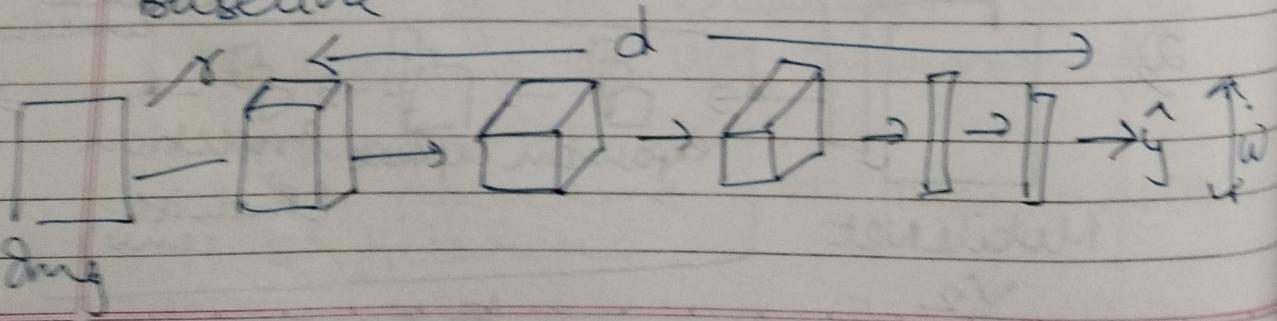
Bottleneck block

- 1) So by using expansion operation increases size of representation within bottleneck block this allows N.N to learn a richer function

2) But in mobile there is memory constraint. Bottleneck uses pointwise convolution projection to project it back to smaller set of values so when we pass this to next block amt of memory to store these values is reduced back down

EfficientNet → used to scale up or down N.N for a particular device

Baseline

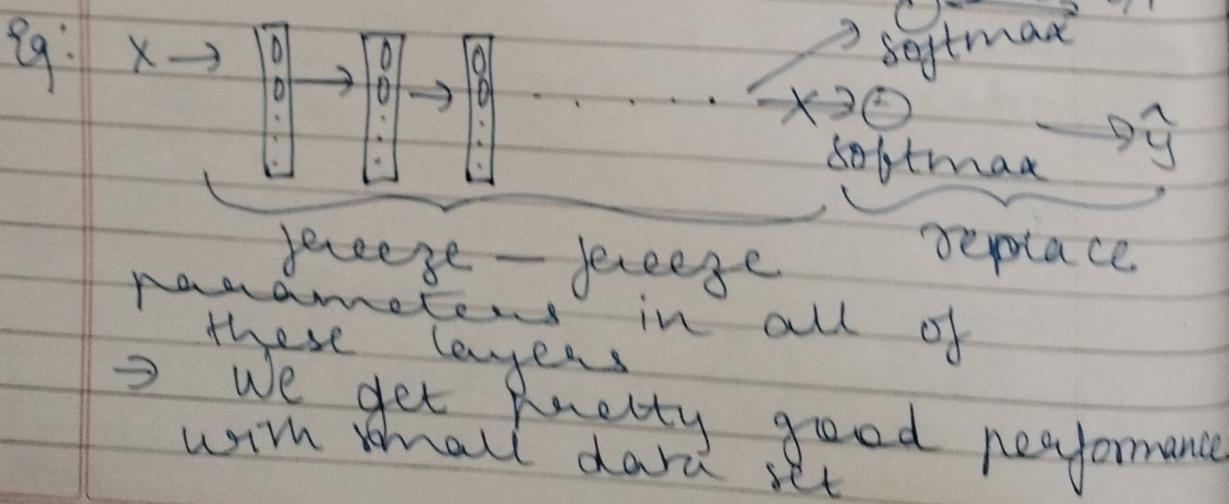


- to scale things up or down
- New img resolution & /
deeper network ie vary d/
vary w
- Given a computational
budget — what's the good choice
of r/d/w or depending on
computational resources , we can
also use compound scaling
where we simultaneously
scale up or down the
resolution of the img &
depth / width of NN
- Rate at which we should
scale up r, d/w — look at
one of the open src. implementations
of EfficientNet which will help
choose a good trade off btw r,d,w

Practical Advice for Using ConvNets

Transfer Learning

- For building Comp. vision application rather than training the way from scratch from random initialization its faster to download way that someone has already trained and use as pre training and transfer that to new task
- Types of datasets - ImageNet
- MS COCO or Pascal
- If we don't have lot of training data then we download the open src implementation of a NN (code + weights)
- Then we can get rid of softmax layer and create our own softmax layer ^{custom} O/P



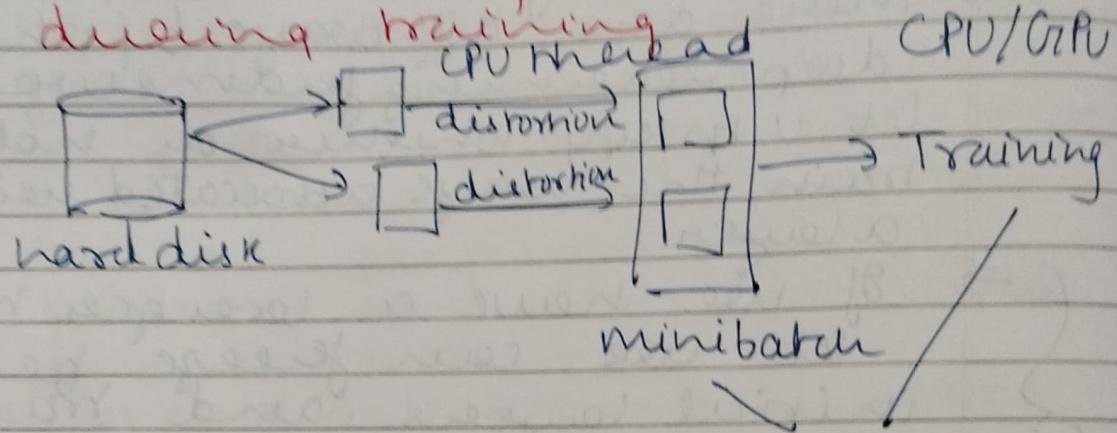
- Depending on framework it might have things like trainable parameters ≈ 0 we might set that for some early layers
- There are diff ways by which deep learning frameworks let us specify whether or not to train the weights associated with a layer
- If we have a larger network we can freeze the initial layers and train the rest of the layers along with own O/P unit
- Couple of ways to do this:
 - Take last few layers and use that as initialization and do gradient descent \rightarrow blow away the last few layers and use ~~gradient descent~~ own hidden units with own softmax/o/p
 - If we have lot of data we can use the network as initialization and train whole network

Data Augmentation

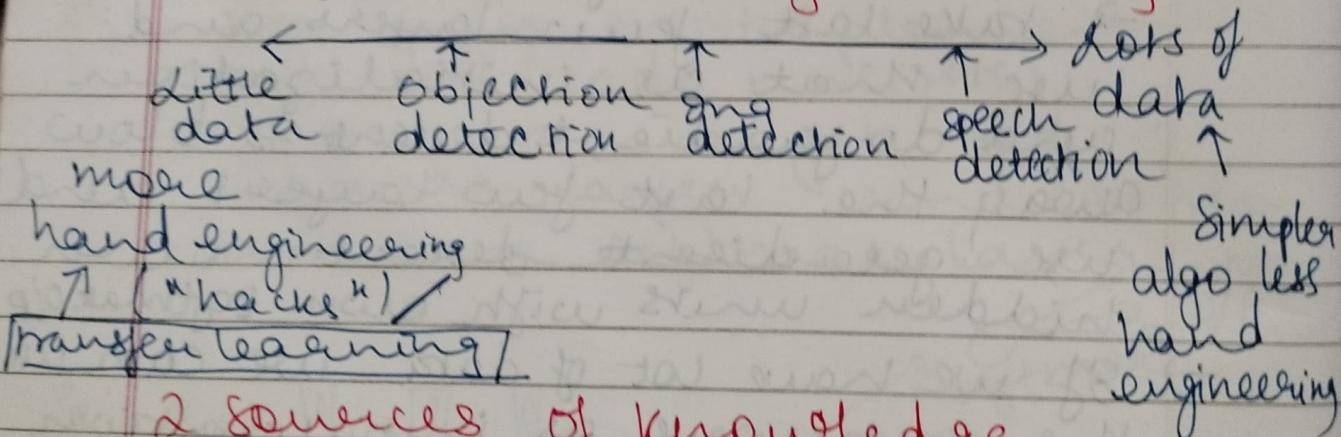
- The simplest form is mirroring on vertical axis
- Another is random cropping
- Rotation, shearing, local warping

→ And type commonly used
is colour shifting → add
distortions to R, G and B channels

Implementing distortion during training



Data vs hand engineering



2 sources of knowledge

- Labelled data (x, y)
- Hand engineering features/
network architecture /
other components

→ Computer vision relies more on hand engineering & has developed a more complex network architectures.

- * Ensembling - Train several networks independently & avg their O/P (\hat{Y})
- * Multicrop @ test time - Run classifier on multiple versions of test images & average results