Classic neural network architecture:-

- LeNet-5
- AlexNet
- VGG
- ResNet (155 layers)

## LetNet-5

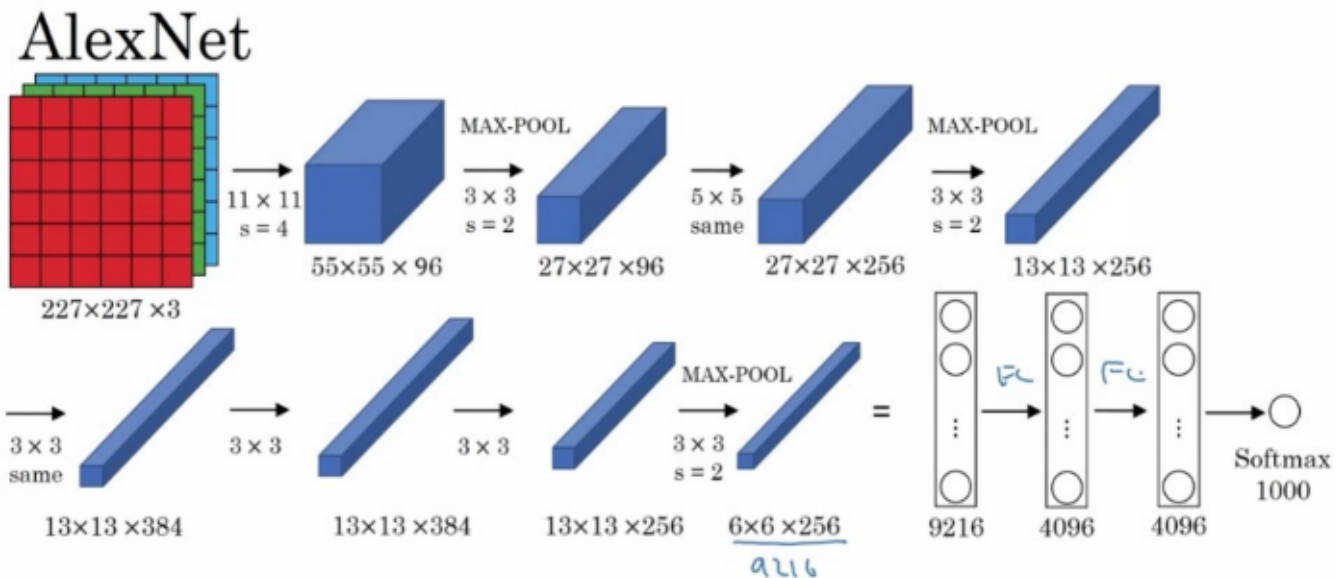*average pooling was more popular back then.*



avg pool

$f = 5 \times 5$
$s = 1$

32×32×1    28×28×6

$f = 2$
$s = 2$

14×14×6    $5 \times 5$ $s = 1$    10×10×16

avg pool
$f = 2$
$s = 2$

6×5×16    FC

120

non liniearity was done after pooling.
sigmoid/tanh was used.

60K parameters

No padding was used.

$\hat{y}$

softmax wasn't used.

84

## AlexNet :-



AlexNet

227×227×3
$11 \times 11$ $s = 4$    55×55×96
MAX-POOL $3 \times 3$ $s = 2$    27×27×96
$5 \times 5$ same    27×27×256
MAX-POOL $3 \times 3$ $s = 2$    13×13×256

$3 \times 3$ same    13×13×384
$3 \times 3$    13×13×384
$3 \times 3$    13×13×256
MAX-POOL $3 \times 3$ $s = 2$    6×6×256

9216

= 9216    FC    4096    FC    4096    Softmax 1000

[Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks]
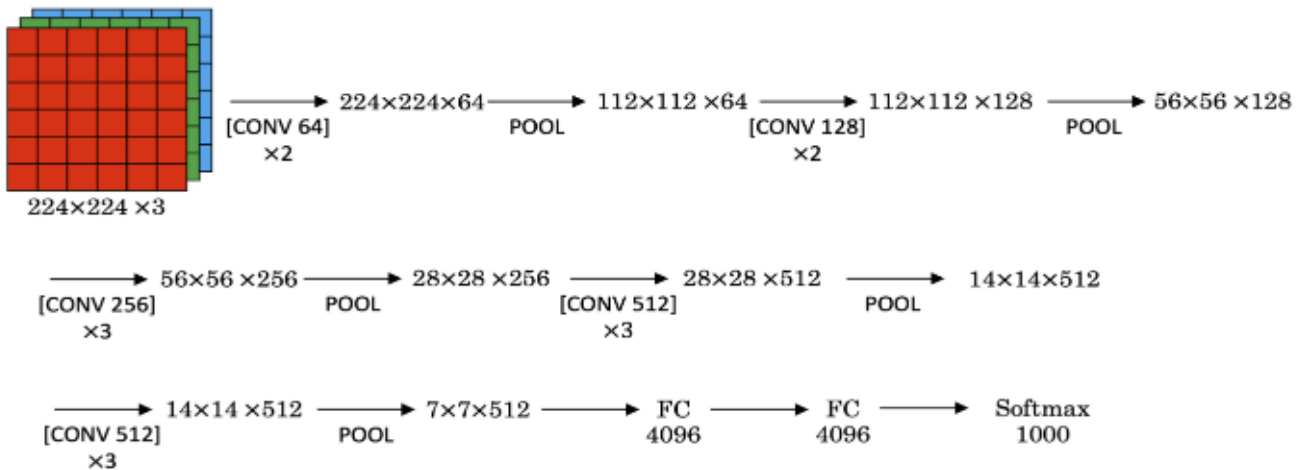
Andrew Ng

-60M parameter.

- ReLu
- Multiple GcPU.
- Local respose normalization (LRN) - it normalizes along the channel.

## VGG-16



**VGG - 16**

CONV = 3×3 filter, s = 1, same        MAX-POOL = 2×2 , s = 2

224×224 ×3

→ 224×224×64 → 112×112 ×64 → 112×112 ×128 → 56×56 ×128
[CONV 64]      POOL      [CONV 128]      POOL
×2                        ×2

→ 56×56 ×256 → 28×28 ×256 → 28×28 ×512 → 14×14×512
[CONV 256]      POOL      [CONV 512]      POOL
×3                        ×3

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

[Simonyan & Zisserman 2015. Very deep convolutional networks for large-scale image recognition]                    Andrew Ng
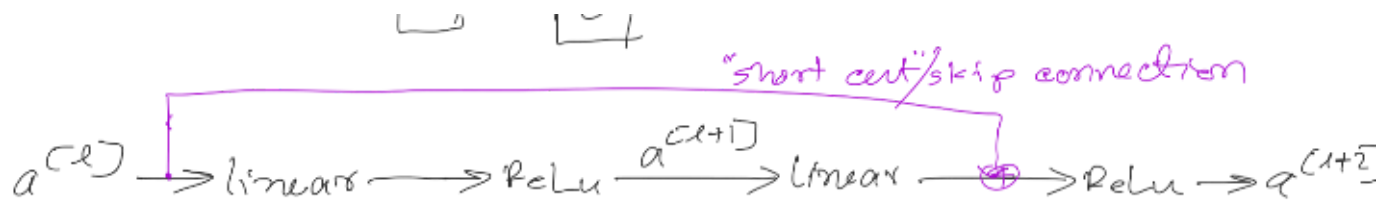
160M parameters.
it's simplicity makes it populas (like same filter size)

## Residual Network (ResNet):-

Very, very deep neural networs are harder to train because of vanishing & exploding gradient.
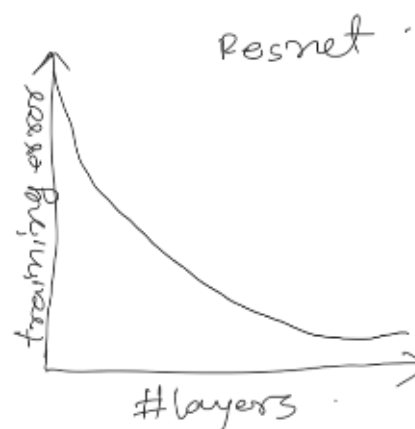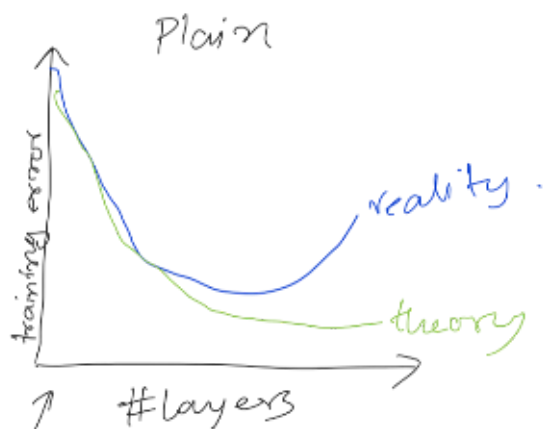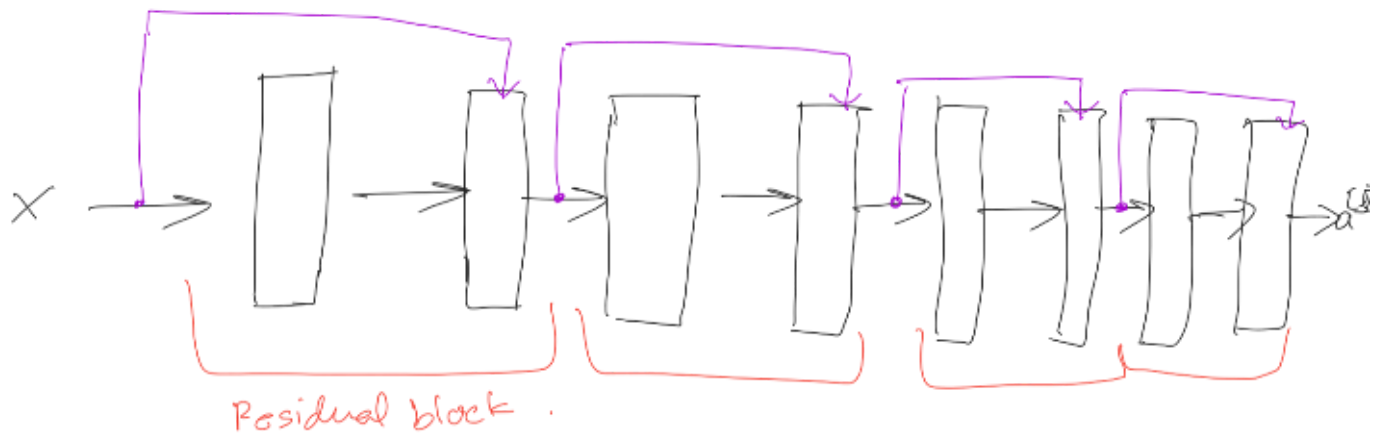
Residual block :-

$$a^{[\ell]} \rightarrow \boxed{\begin{matrix}o\\o\\o\end{matrix}} a^{[\ell+1]} \rightarrow \boxed{\begin{matrix}o\\o\\o\end{matrix}} \rightarrow a^{[\ell+2]}$$

"short cut"/skip connection

$$a^{[l]} \longrightarrow \text{linear} \longrightarrow \text{Relu} \xrightarrow{a^{[l+1]}} \text{linear} \longrightarrow \oplus \longrightarrow \text{Relu} \longrightarrow a^{[l+2]}$$

$$z^{[l+1]} = w^{[l+1]} a^{[l]} + b^{[l+1]}$$
$$z^{[l+2]} = w^{[l+2]} a^{[l+1]} + b^{[l+2]}$$

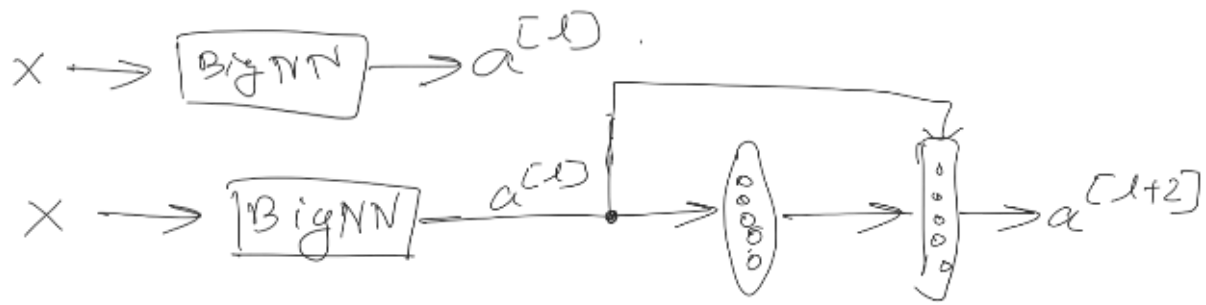$$a^{[l+1]} = g(z^{[l+1]})$$
$$\cancel{a^{[l+2]} = g(z^{[l+2]})}$$

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

## Residual network :-



Residual block .



Plain

reality .

theory

#layers

↳ this might happen for vanishing exploding gradients.

Resnet .

#layers .

## why resNetwork :-

$$X \longrightarrow \boxed{BigNN} \longrightarrow a^{[l]}$$

$$X \longrightarrow \boxed{BigNN} \longleftarrow a^{[l]} \longrightarrow a^{[l+2]}$$

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

$$\underset{256}{\uparrow}$$

$\longrightarrow 0$ if this term becomes 0 (vanishing gradient/ regularizati)

$$\longrightarrow g( w^{[l+2]} a^{[l+1]} + b^{[l+2]} + \underset{256 \times 128}{\underset{\downarrow}{W_s}} \underset{128}{a^{[l]}} )$$

$$\longrightarrow g(a^{[l]})$$

Identity function is easy to earn for residual block.

if dimension of $a^{[l+2]}$ & $a^{[l]}$ is different then we have to we have to multiply the term $a^{[l]}$ with $W_s$. $W_s$ might be fixed value matrix. or it may 0 pads the $a^{[l]}$ to match the dimensions of $a^{[l+2]}$
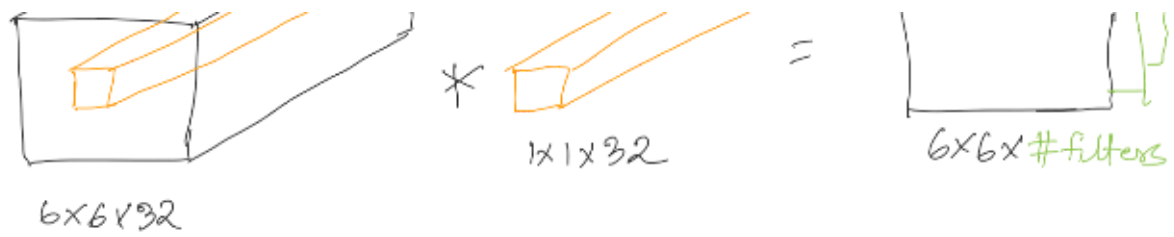
what does 1x1 filter do?

| 1 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|
| 2 | 3 | 6 |   |   |   |

$* \boxed{2} =$

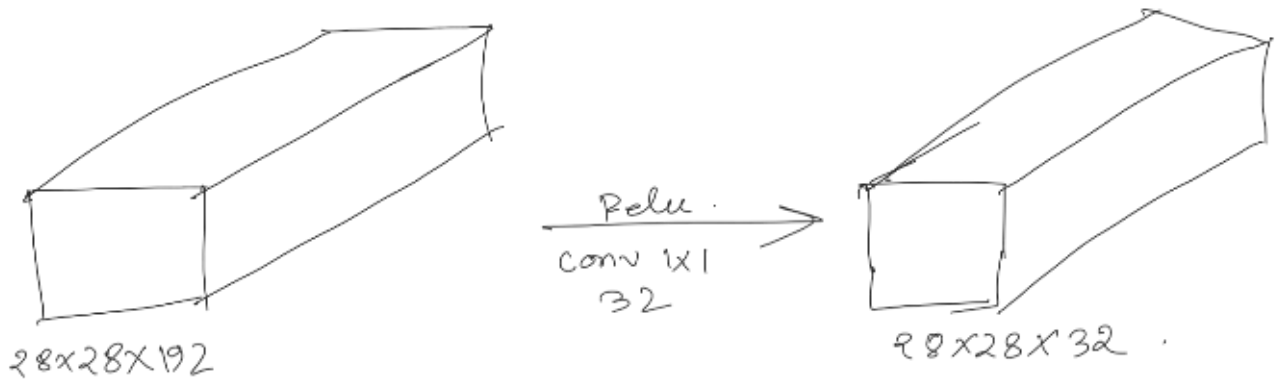| 2 | 6 | 8 | 10 | 12 | 14 |
|---|---|---|----|----|----|
| 4 | 6 | 12 |   |    |    |

it seems 1x1 filter just multiply the orignal data. which is not usefull.
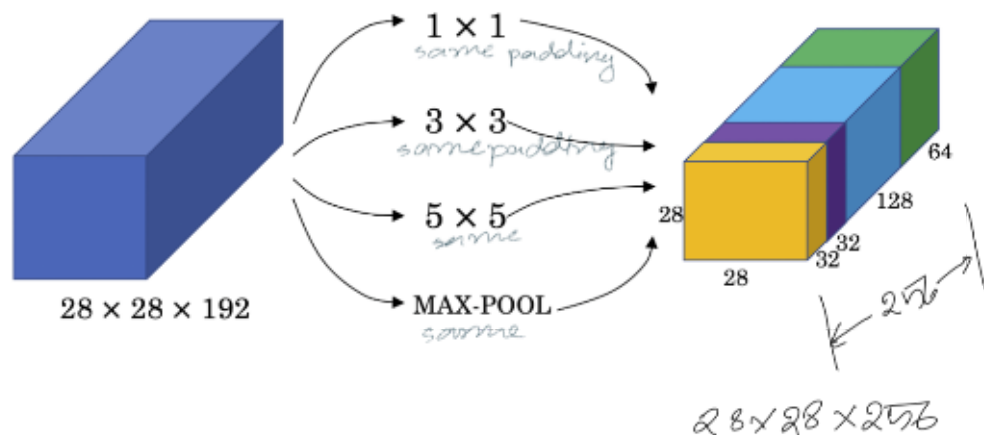But it is useful in multidimensional image

multidimension :-

$6 \times 6 \times 32$ $\quad * \quad$ $1 \times 1 \times 32$ $\quad = \quad$ $6 \times 6 \times$ #filters

it is also called networ in network architecture.



$28 \times 28 \times 192$ $\quad\xrightarrow[\substack{conv \ 1 \times 1 \\ 32}]{Relu}\quad$ $28 \times 28 \times 32$ .

if we have a large channel and want to shrink the channel then we can use $1 \times 1$ filter to shrink the channel pooling layer is used to shrink the height & width.
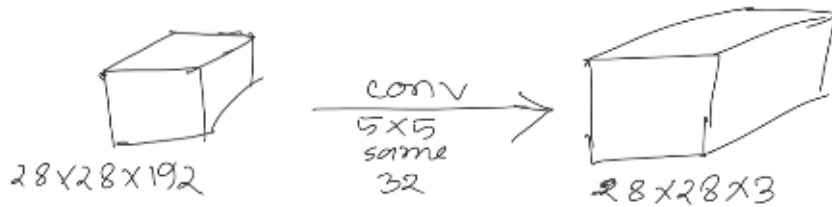
## Inception network :- (Google Net)

Instead of choosing the filter size / max pooling / avg pooling. inception network does all of them at the same time. And it let the parameter learn whatever it wants to learn



$28 \times 28 \times 256$

# Computational cost is higher for inception :-

Let's focus on the computational cost of 5x5 filter



28×28×192 → conv 5×5 same 32 → 28×28×3

each of the 32 filters are 5×5×192

so total computation. $28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120 M.$

We can use 1×1 convolution to reduce the computation cost :-



← bottleneck layer.

28×28×192 → conv 1×1. 16, 1×1×192 → 28×28×16 → CONV 5×5, 32 5×5×16 → 28×28×32.

Cost :-

$28 \times 28 \times 16 \times 1 \times 1 \times 192 = 2.4 M.$     $28 \times 28 \times 32 \times 5 \times 5 \times 16 = 10 M.$

$+$

$\boxed{12.4 M}$