

Outline:

Classic Networks:

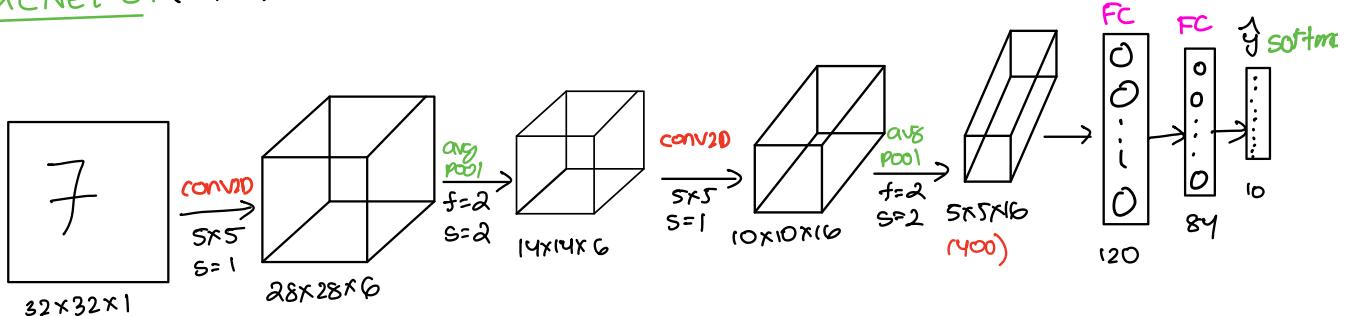
- LeNet-5
- AlexNet
- VGG

ResNet (152 layers)

Inception.

dec1: Classical Networks:

LeNet-5: (1998)



- parameters = 60k parameters

- $n_H, n_W \downarrow, n_C \uparrow$

- conv, pool, conv, pool, fc, fc output.

Advanced: used sigmoid / tanh function → **use ReLU**

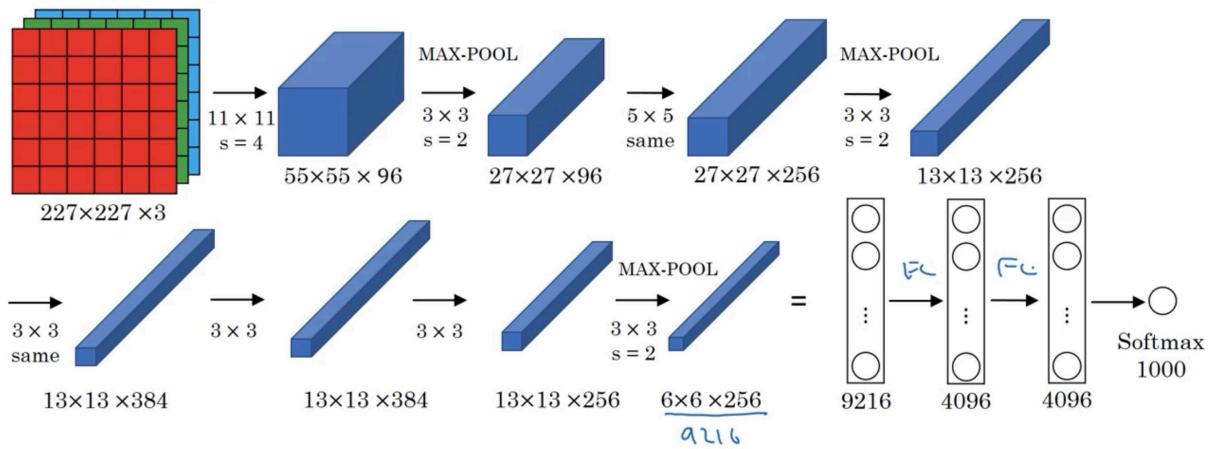
→ filter channel size is different than input channel size.

$$(n_H, n_W, n_C) \quad (f, f, n_C)$$

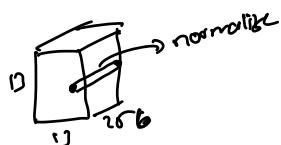
different (due to computation saving)

- use non-linearity after pooling (sigmoid).

AlexNet: (2012)



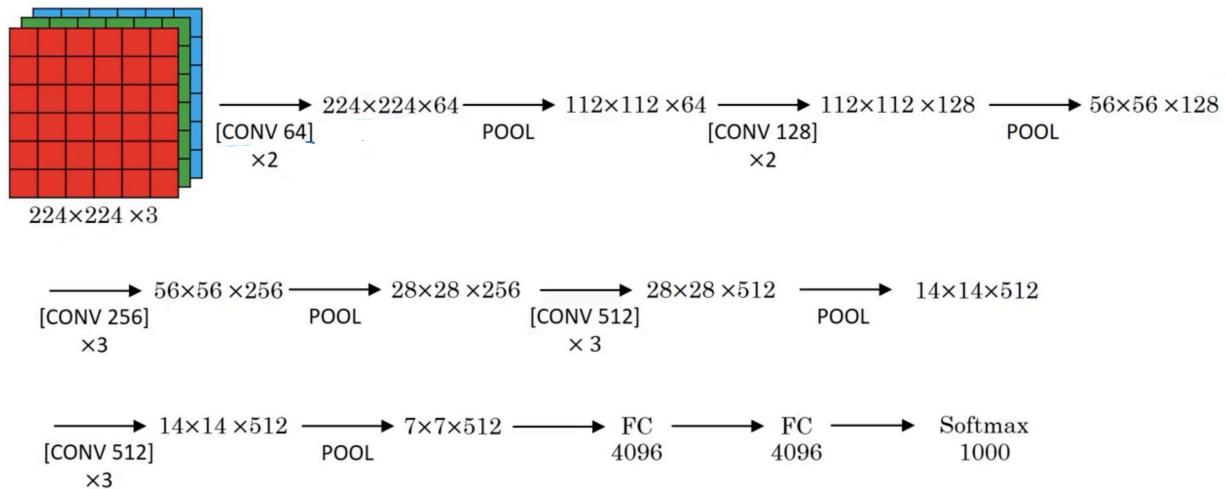
- Much bigger than deNet, has 60M parameters.
- use ReLU
- Multiple Gpu's training
- (less important) • local response Normalization



VGG-16 (2015) (VGG-19)

CONV = 3×3 filter, $s=1$, same

Max-POOL = 2×2 , $s=2$



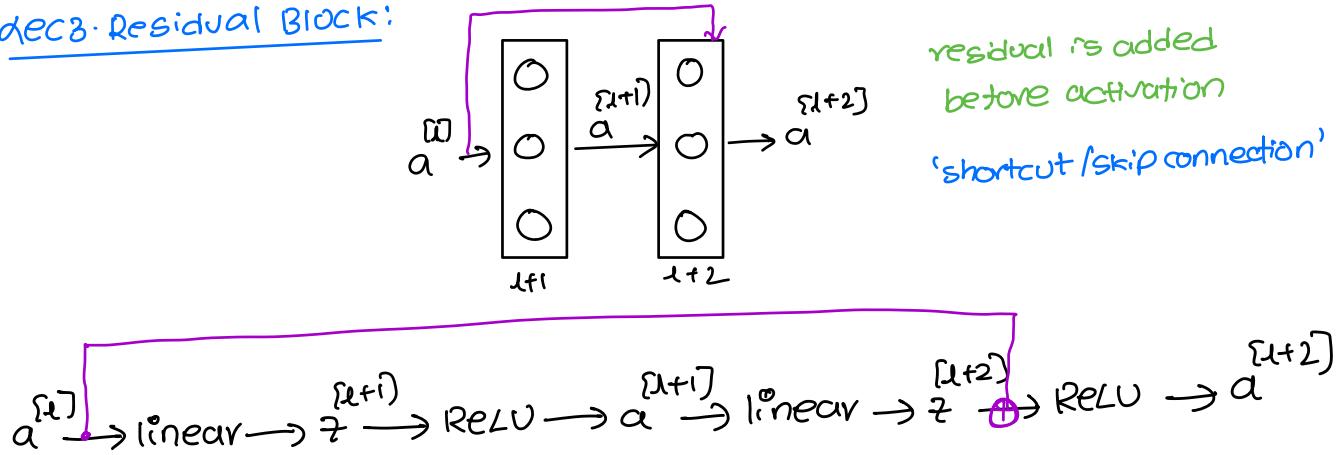
→ 16 layers which have weights.

→ simple hyperparameters to tune

→ filters (64, 128, 256, 512, 512)

→ large network to train. (138 M parameters)

dec3. Residual Block:

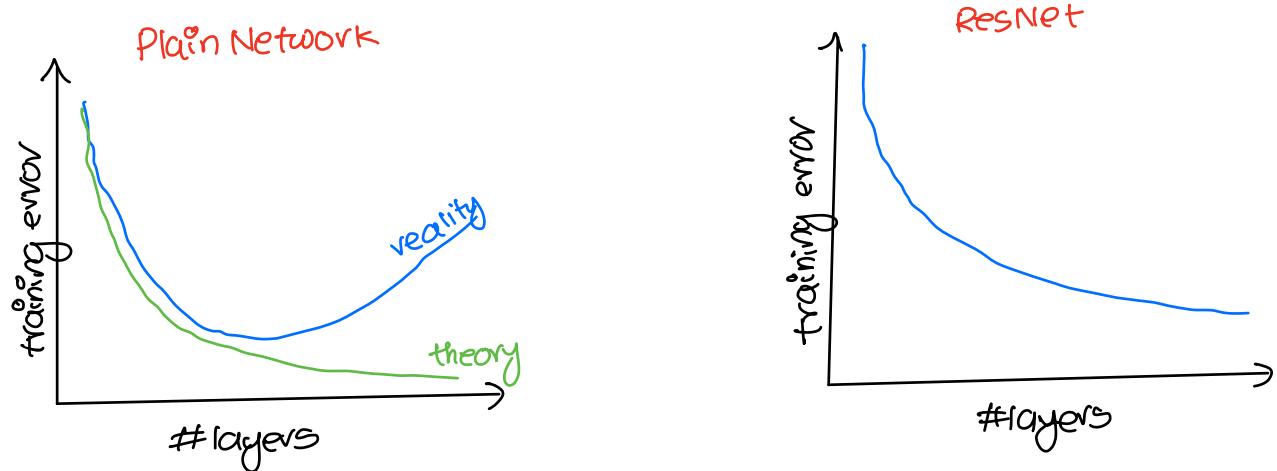
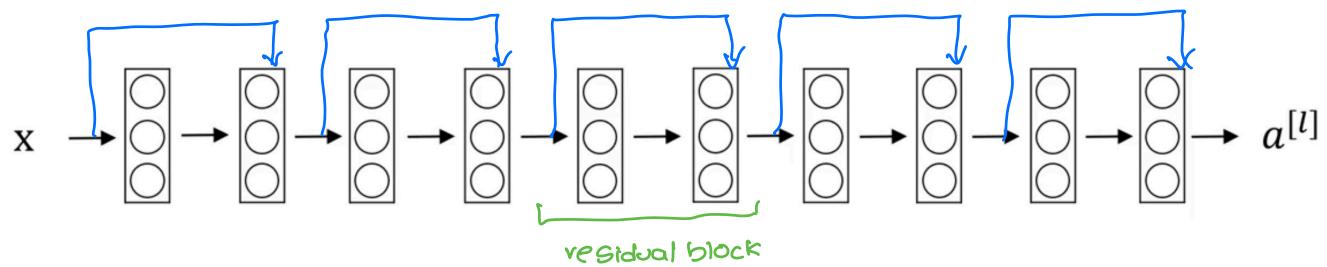


$$\begin{aligned}
 z^{[l+1]} &= W^{[l+1]} a^{[l]} + b^{[l+1]} \\
 a^{[l+1]} &= g(z^{[l+1]}) \\
 z^{[l+2]} &= W^{[l+2]} a^{[l+1]} + b^{[l+2]} \\
 a^{[l+2]} &= g(z^{[l+2]}) \quad \times \\
 a^{[l+2]} &= g(z^{[l+2]} + a^{[l]}) \quad \text{in red}
 \end{aligned}$$

Advantages:

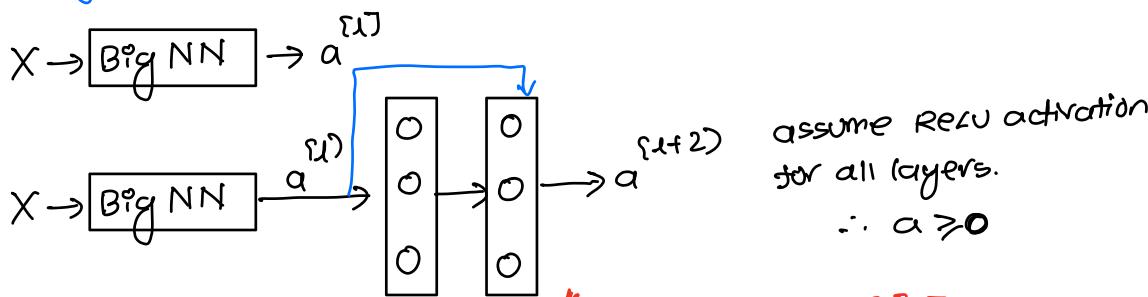
- using residual blocks allows you train much bigger deep learning networks.

Residual Network:



→ helps in vanishing/exploding gradient problem and
train much deeper networks.

dec 4: Why ResNets Work:



assume ReLU activation
for all layers.
 $\therefore a \geq 0$

$$\begin{aligned}
 a^{l+2} &= g(\tilde{z}^{l+2} + a^{l+1}) \\
 &= g(W^{l+2}a^{l+1} + b^{l+2} + a^{l+1}) \\
 \text{assume if } W^{l+2} &= 0 \quad b^{l+2} = 0 \\
 &= g(a^{l+1}) \\
 \therefore a^{l+2} &= a^{l+1} \quad \{\because a \geq 0\}
 \end{aligned}$$

[dim \neq a^{l+2}] *

if dim \neq a^{l+2}

$$a^{l+2} = g(\tilde{z}^{l+2} + a^{l+1})$$

256

odd W_s (256, 128)

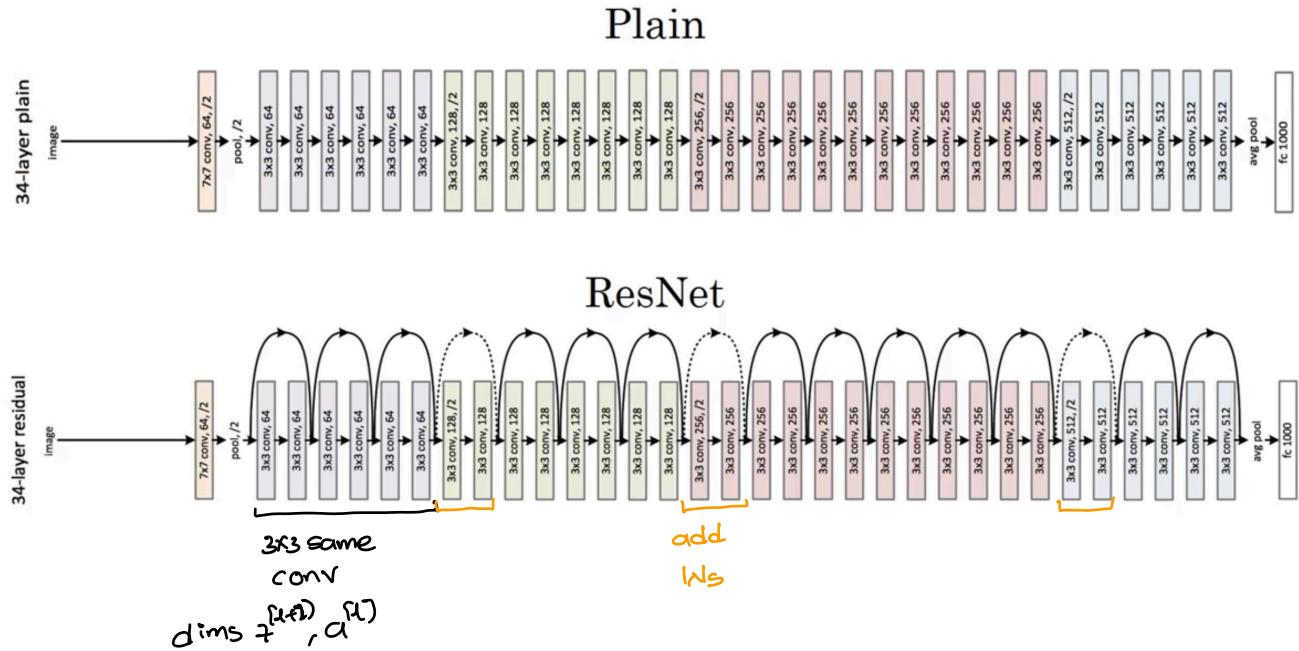
$$a^{l+2} = g(\tilde{z}^{l+2} + W_s a^{l+1})$$

W_s can be learned
 W_s can be fixed, and zero
pad a^{l+1}

Shows:

- Identity function is easy to learn for residual block, because of skip connection.
- adding extra two layers with skip connection, doesn't hurt your NN ability to do as well as simpler NN.
- Hence we can train bigger Networks, we can learn better than identity function.

ResNet: (2015)



dec 5: 1×1 convolutions!

1	2	3	6	5	8
3	5	5	1	3	4
2	1	3	4	9	3
4	7	8	5	7	9
1	5	3	7	4	8
5	4	9	8	3	5

$6 \times 6 \times 1$

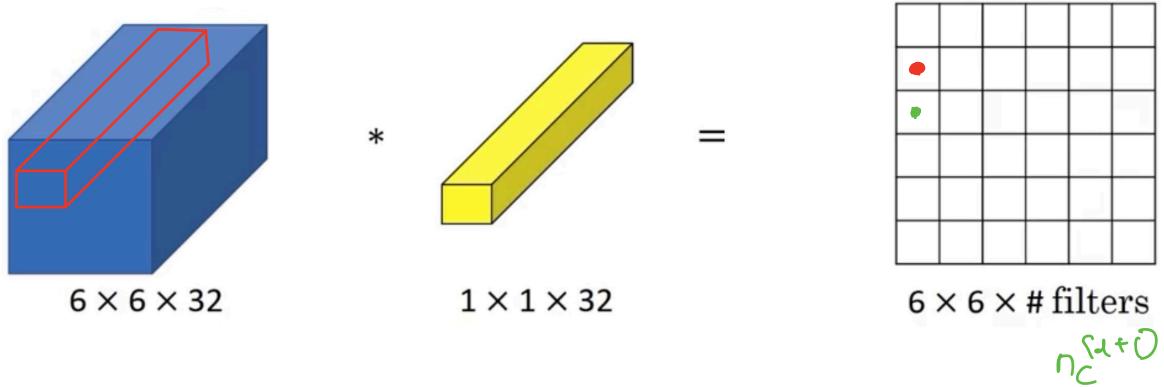
multiply each element by 2.

2

=

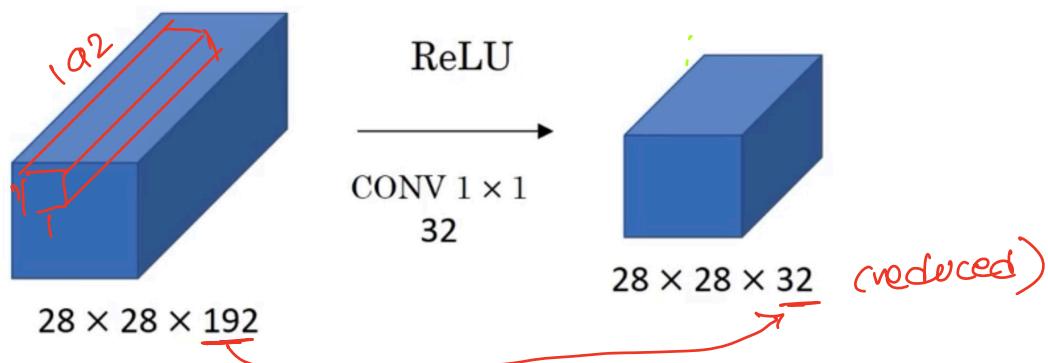
2	4	6	12	10	16
6					
4					
8					
2					
0					

Multiple channel: (Network in Network)



→ element wise product between 32 numbers
and apply ReLU, gives single scalar.

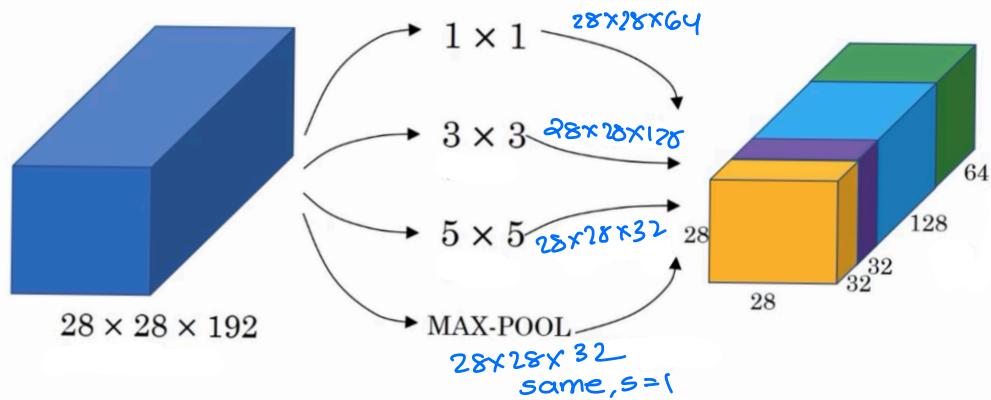
Using 1×1 convolutions:



Filter dim: $1 \times 1 \times 192$ (reduces the no. of channels)

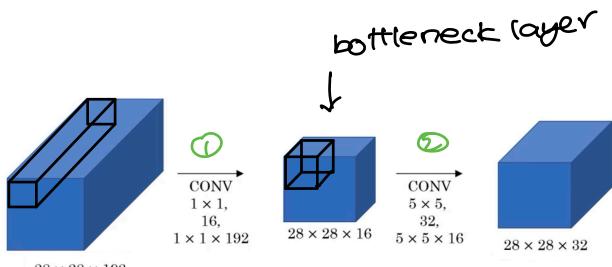
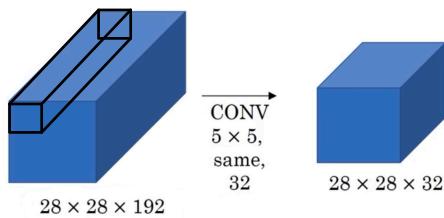
no. of filters: 32

dec6: Inception Network :



- Instead of choosing what filter size you want, conv, or pooling
use them all.

Problem: computational cost:



cost:

32 filters

$$\text{filter size} = 5 \times 5 \times 192$$

$$\text{output size} = 28 \times 28 \times 32$$

$$\text{no.of multi} = 28 \times 28 \times 32 \times 5 \times 5 \times 192$$

$$= 120M$$

cost:

① 16 filters

$$\text{filter size} = 1 \times 1 \times 192$$

$$\text{outsize} = 28 \times 28 \times 16$$

$$\text{no.of multi} = 28 \times 28 \times 16 \times 1 \times 1 \times 192 \\ = 2.4M$$

② 32 filters

$$\text{filter size} = 5 \times 5 \times 16$$

$$\text{outsize} = 28 \times 28 \times 32$$

$$\text{no.of multi} = 28 \times 28 \times 32 \times 5 \times 5 \times 16 \\ = 10.0M$$

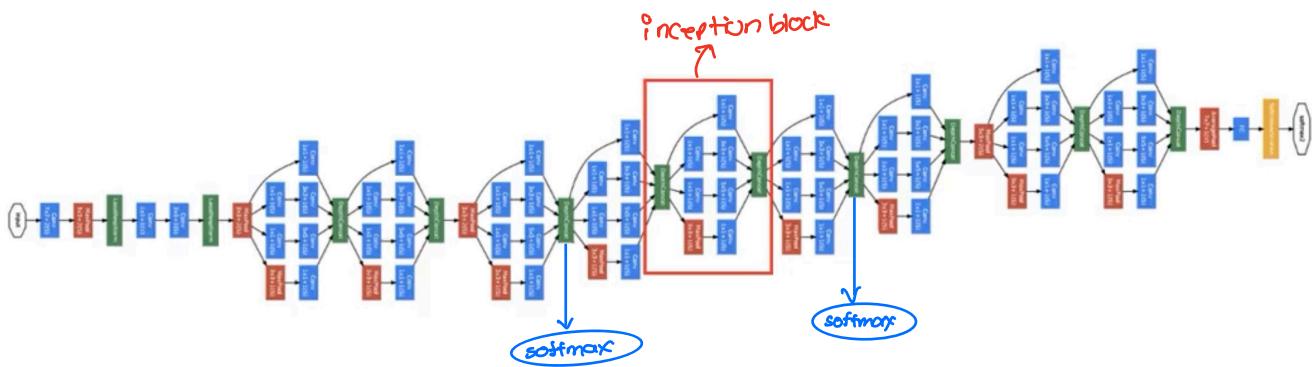
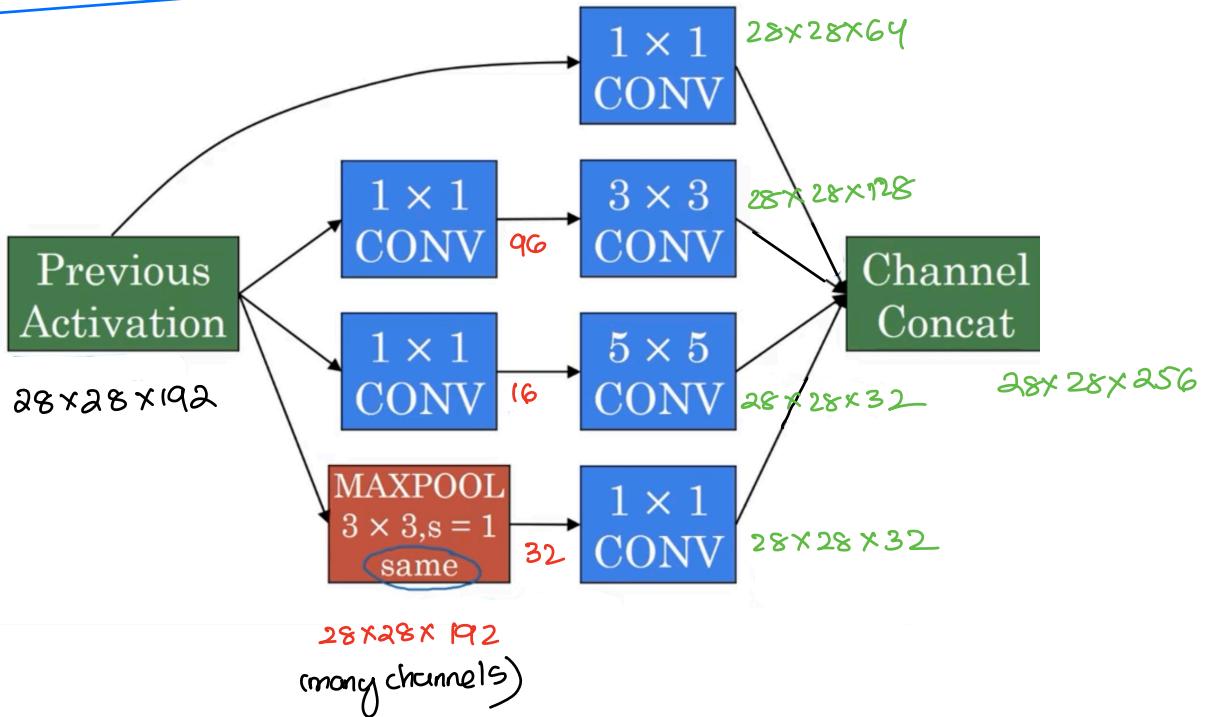
$$\textcircled{1} = 120M$$

Reduce the cost

$$\textcircled{1} + \textcircled{2} = 12.4M$$

dec7: Inception Network: (Google Net)

Inception module:



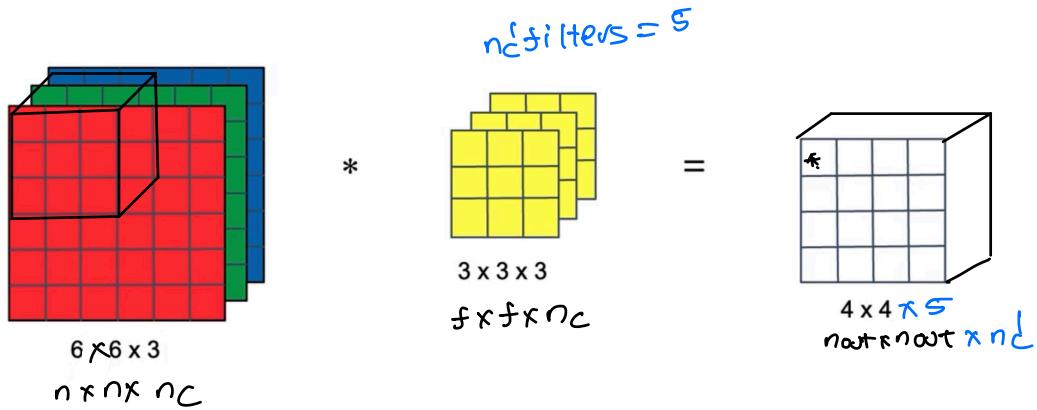
- multiple inception blocks
- side branches, for better regularization effect.



Inception Network → Inception Movie.

AEC8: MobileNet:

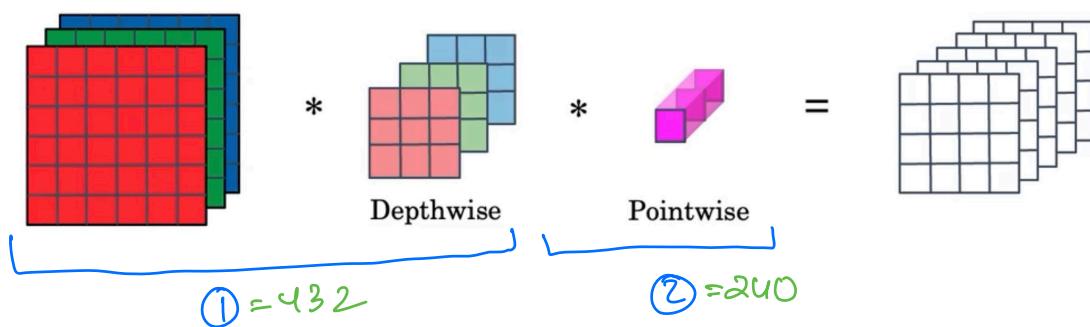
key idea: Normal vs depthwise separable convolutions



computation cost:

$$\begin{aligned}
 &= \text{filter rows} \times \text{filter positions} \times \text{no. of filters} \\
 &= 3 \times 3 \times 3 \times 4 \times 4 \times 5 = 2160
 \end{aligned}$$

Depthwise Separable Convolution



compute cost: 672 multiplications

$$\frac{672}{2160} = 0.31 \quad \text{save } 3 \times \text{compute cost}$$

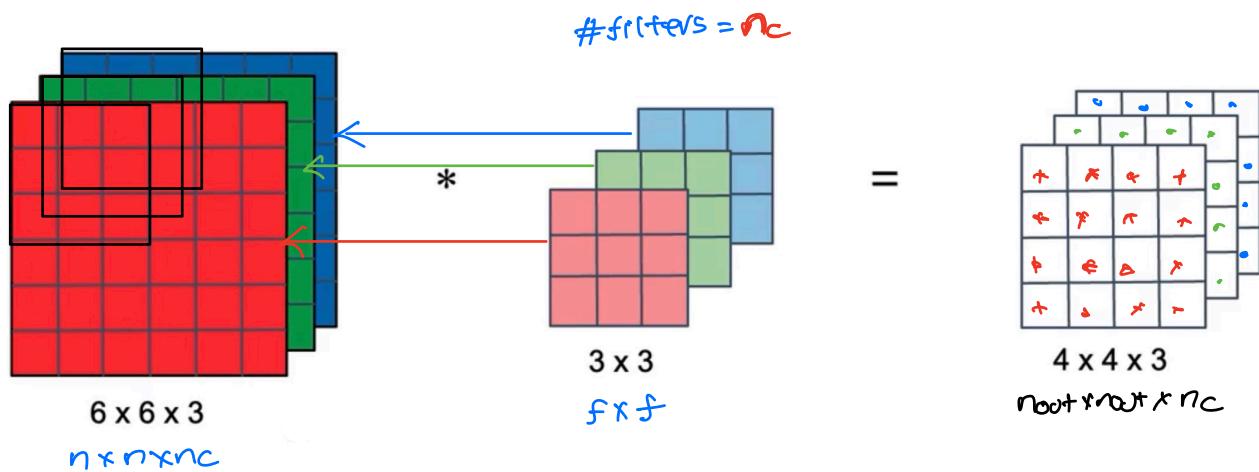
General cost ratio = $\frac{\text{depth}}{\text{normal}}$

$$\Rightarrow \frac{1}{n_C'} + \frac{1}{f^2} = \frac{1}{5} + \frac{1}{32} \approx 0.31$$

If $n_C' = 5 \times 2$, $f = 3$

cost ratio = 0.11 $\Rightarrow 10 \times$ time cheaper

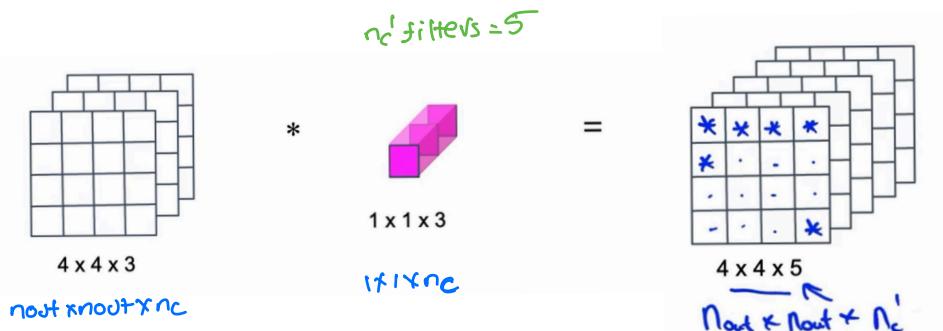
Depthwise Convolution:



compute cost:

$$\begin{aligned}
 &= \text{filter params} \times \text{filter position} \times \text{no. of filters} \\
 &= 3 \times 3 \times 4 \times 4 \times 3 \\
 &= 432
 \end{aligned}$$

PointWise Convolution:



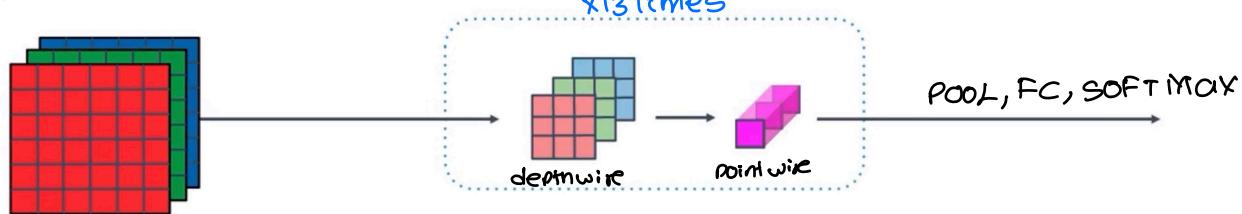
compute cost:

$$\begin{aligned}
 &= \text{filter params} \times \text{filter positions} \times \text{no. of filters} \\
 &= 1 \times 1 \times 3 \times 4 \times 4 \times 5 \\
 &= 240
 \end{aligned}$$

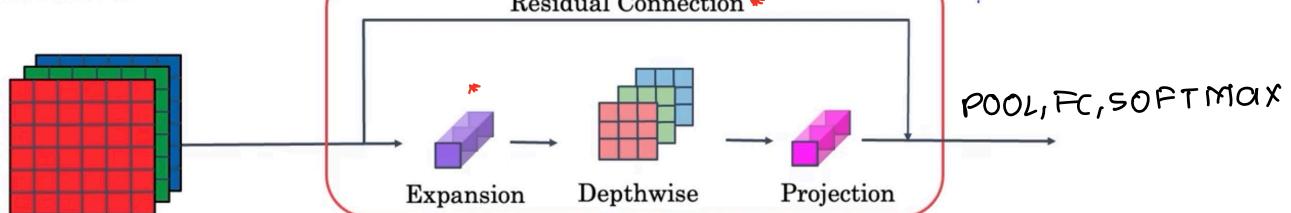
240

decq: MobileNet: (2019)

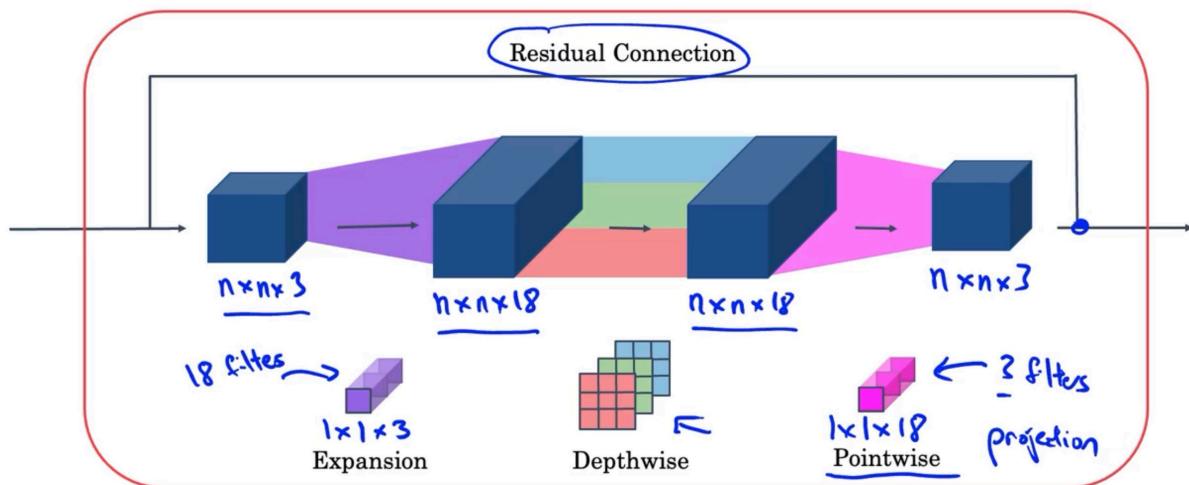
v1



MobileNet v2



v2: Inverted residuals and linear bottlenecks.



Expansion factor: 6

$3 \rightarrow 18$

$1 \times 1 \times 3 \times 18$

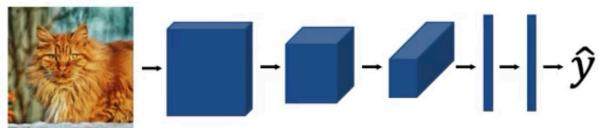
$3 \times 3 \times 18$

$1 \times 1 \times 18 \times 18$

dec10: EfficientNet (2019)

- Model scaling for CNN.

Baseline

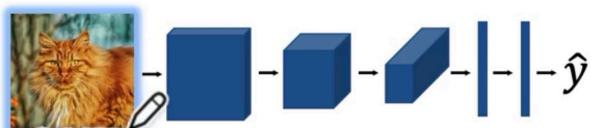


r - resolution

d - depth of NN

w - width of layers

Higher Resolution

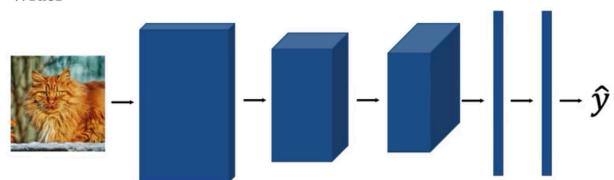


Question?

Given a computational budget,
what is good choice of
(r, d, w)?

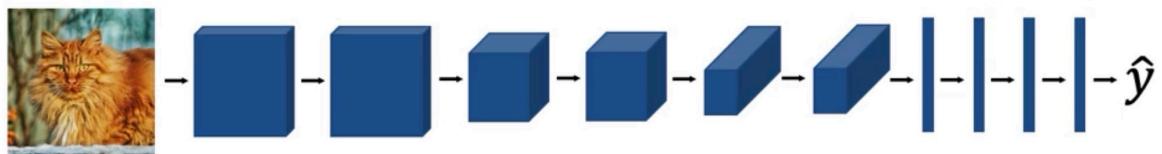
$$\text{compound} = \alpha_1 r + \alpha_2 d + \alpha_3 w$$

Wider

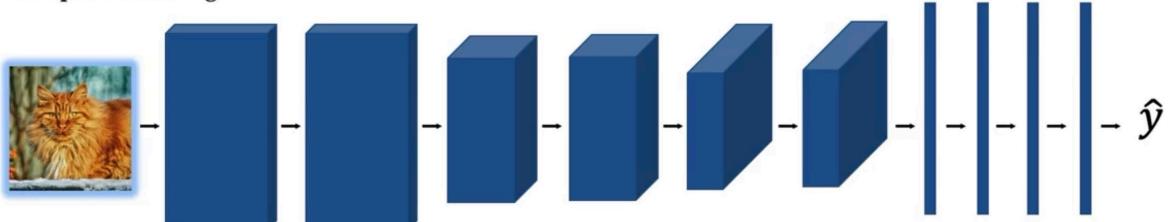


- look at efficient Net code.

Deeper



Compound scaling



Advice:

- use open source code from github.
- use pretrained weights, do transfer learning. (from ImageNet)
 - freeze = 1
 - trainable parameter = 0
 - train softmax layer weights.
 - precompute activation for all training examples (before softmax).
 - if m > large, then train bit more layers.

if you have a lot of data

- train the whole network with your training data.

Data Augmentation!

- more data, better NN training for CV model
- mirroring vertical
- random cropping (reasonable large subset).
- color shifting (+20, -20, +20), (-20, +20, +20) ...

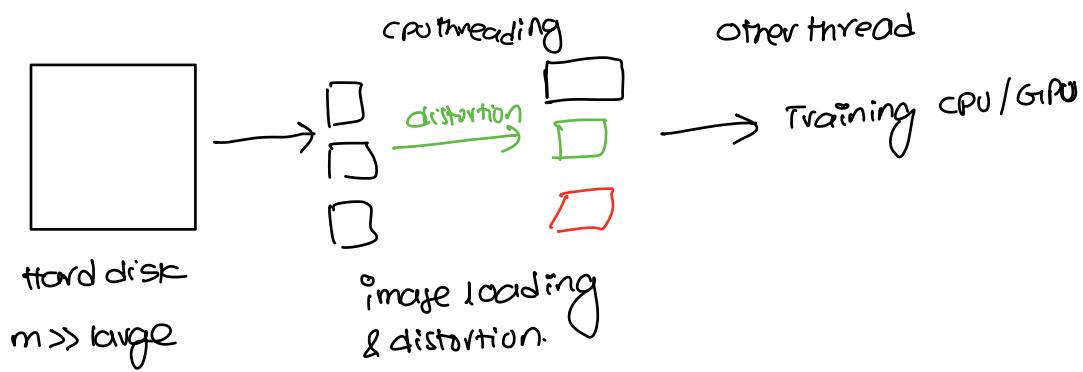
sampling R, G, B

- PCA color augmentation (AlexNet paper)

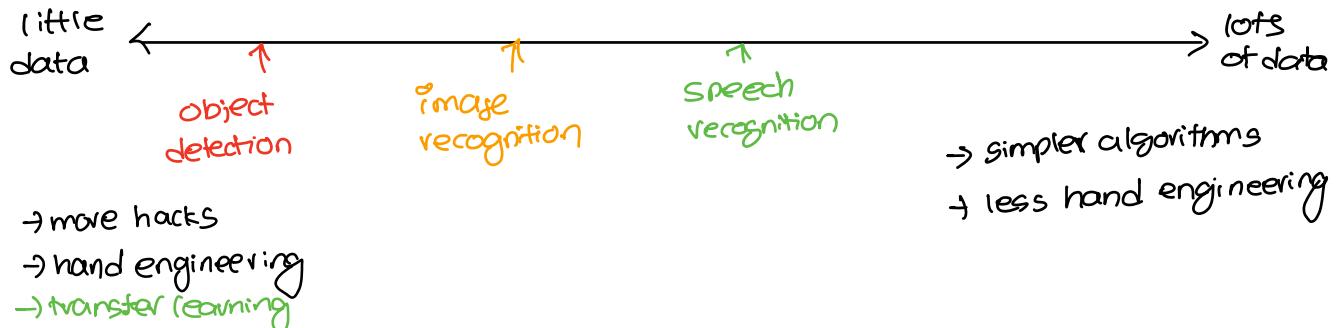
if image has more R, B and less G

then PCA add, sub a lot in R, B, and less G

⇒ overall color remains same.



State of Computer Vision:



Source of knowledge:

① labeled data (x, y)

② hand engineered features / NN architecture / other components ..

Doing well on benchmarks / winning competitions:

① ensembling (2 to 15 NN)

- train several NN independently and average their output. (\bar{y})

② Multicrop at test time

- run classifier on multiple version of test image and average results. (multicrop)

③ use open source code.