

# Why Case Studies? 見賢思齊

## Classic Networks:

- LeNet-5
- AlexNet
- VGG-16

If I want to read the above classic networks, Andrew suggests the following order (from easy to hard)

AlexNet → VGG-16 → LeNet-5

## ResNet

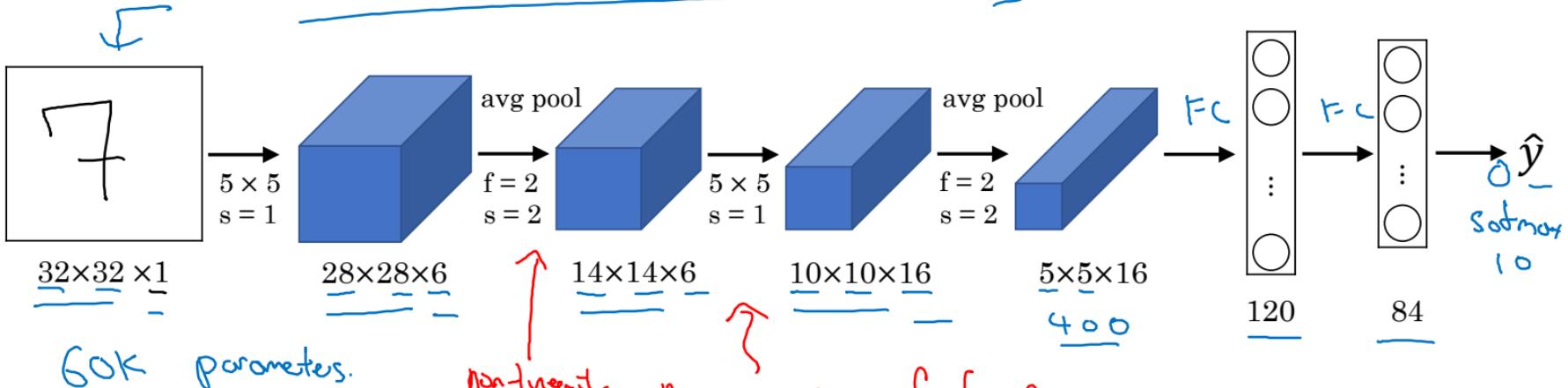
## Inception Network

MobileNet

EfficientNet

These two networks help to build neural networks in embedded systems.

# LeNet - 5 60k parameters



60K parameters.

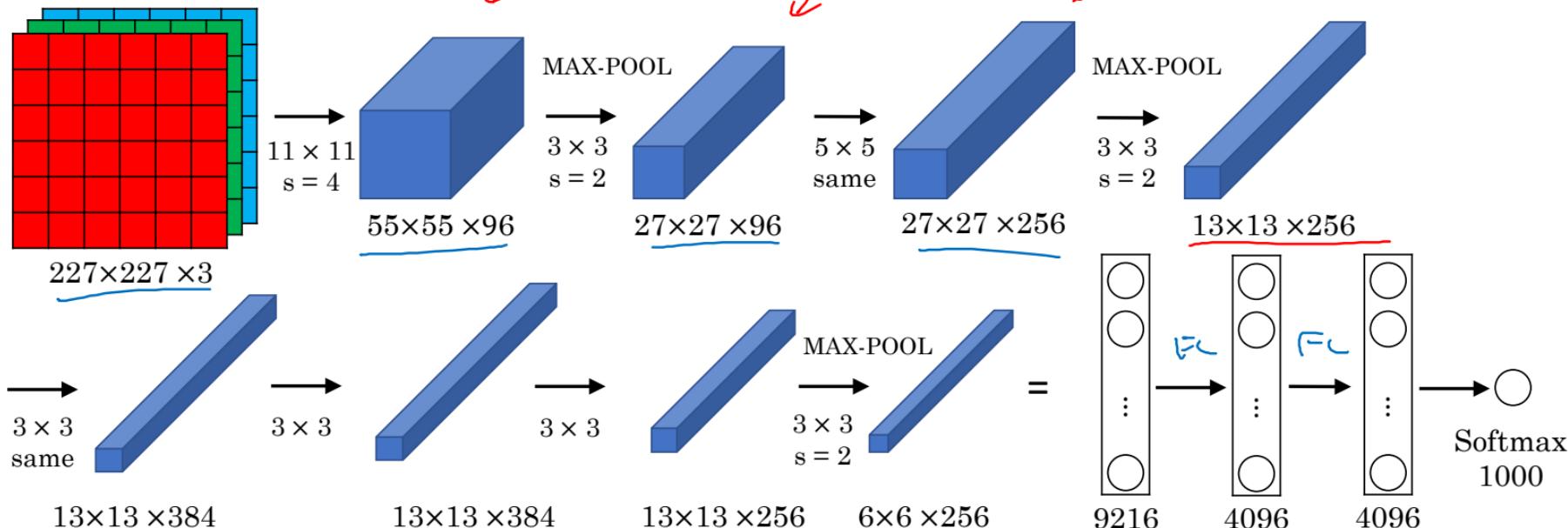
$n_H, n_w \downarrow$   $n_C \uparrow$

conv pool conv pool fc fc output

Activation: Sigmoid/tanh ReLU



# AlexNet 60m parameters



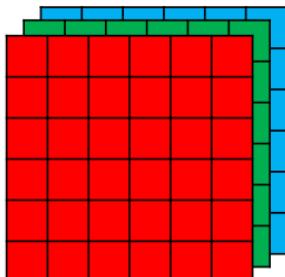
- Similar to LeNet, but much bigger.
  - ReLU
  - Multiple GPUs.
  - Local Response Normalization (LRN)
- $9216$
- $160M$  parameters
- $13 \times 13 \times 256$

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

Andrew Ng

# VGG - 16

CONV =  $3 \times 3$  filter, s = 1, same



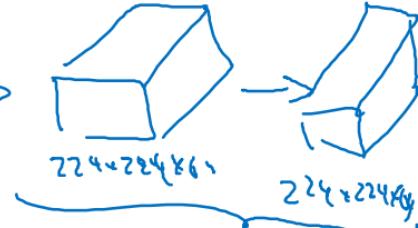
224x224x3

~~VGG-16~~ 138m parameters

MAX-POOL =  $2 \times 2$ , s = 2



$224 \times 224 \times 64$



$112 \times 112 \times 128$

$112 \times 112 \times 128$

[CONV 64]  
 $\times 2$

$224 \times 224 \times 64$   
POOL

[CONV 128]  
 $\times 2$

$112 \times 112 \times 128$   
POOL

224x224x3

[CONV 256]  
 $\times 3$

POOL

[CONV 512]  
 $\times 3$

POOL

$14 \times 14 \times 512$

$56 \times 56 \times 256$

POOL

$28 \times 28 \times 256$

POOL

$28 \times 28 \times 512$

POOL

$14 \times 14 \times 512$

[CONV 512]  
 $\times 3$

POOL

FC  
4096

FC  
4096

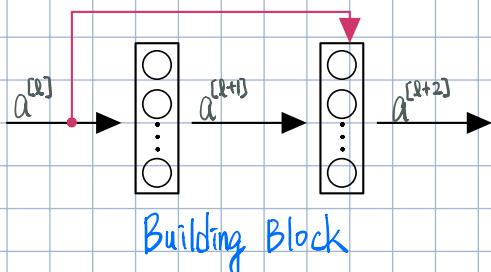
Softmax  
1000

$n_H, n_W \downarrow$

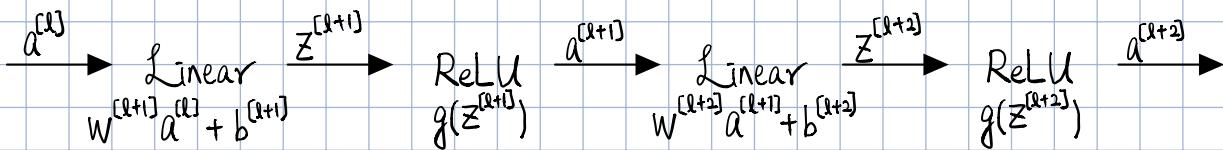
$n_C \uparrow$

$\sim 138M$

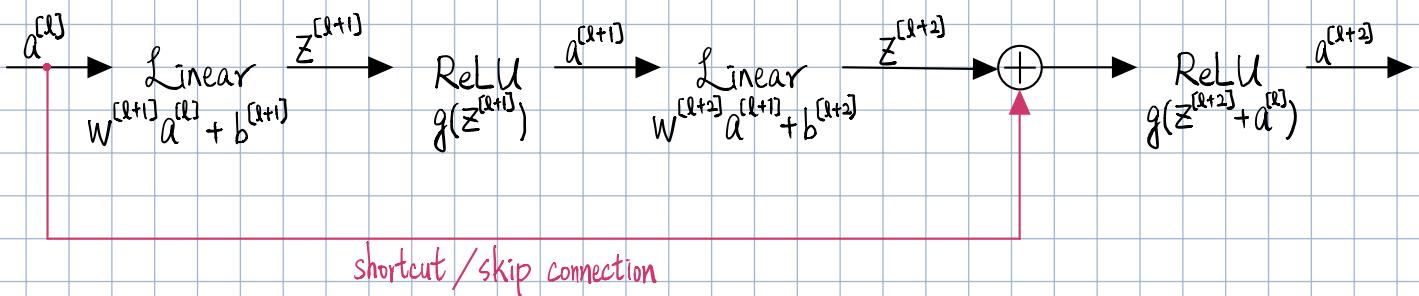
# Residual Neural Network (ResNet)



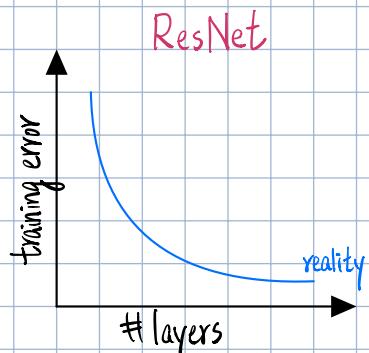
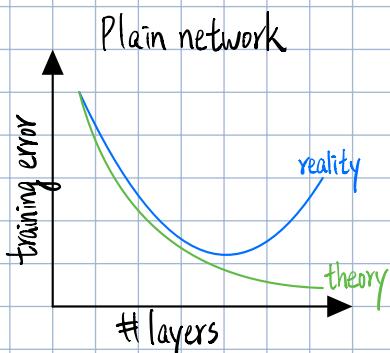
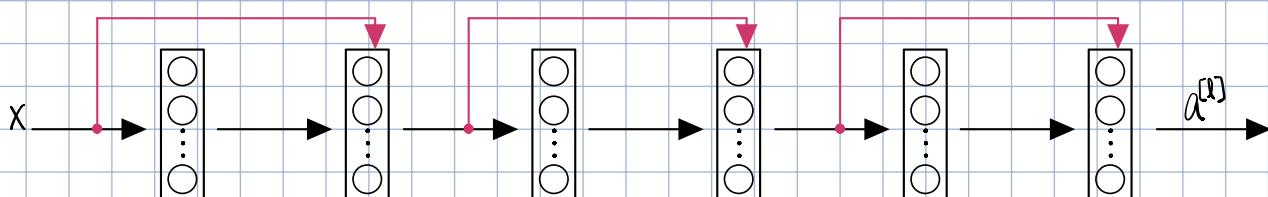
The path of neural network without shortcut:



The path of neural network without shortcut:

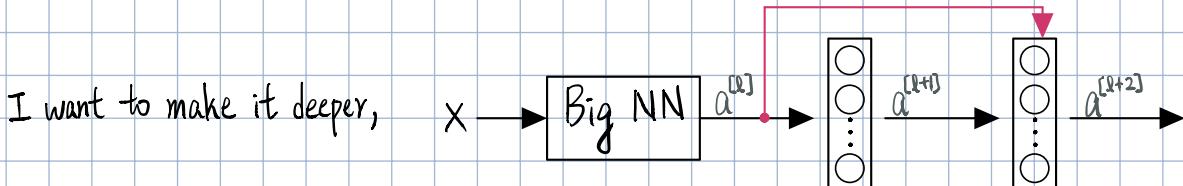


Turn a Plain network into a ResNet.



Why ResNet works? To gain intuition in ResNet.

Given a big neural network,  $X \rightarrow \text{Big NN} \rightarrow a^{[0]}$



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

If  $W^{[l+2]} = 0$  and  $b^{[l+2]} = 0$ , then  $a^{[l+2]} = a^{[l]}$

It shows that the identity function is easy for residual block to learn  $a^{[l+2]} = a^{[l]}$ . What it means is that adding the addition residual block in my neural network, it doesn't hurt my neural network's ability to do as well as the simpler network without these two extra layers.

Note: To make  $g(z^{[l+2]} + a^{[l]})$  work,  $z^{[l+2]}$  and  $a^{[l]}$  need to have the same dimension.

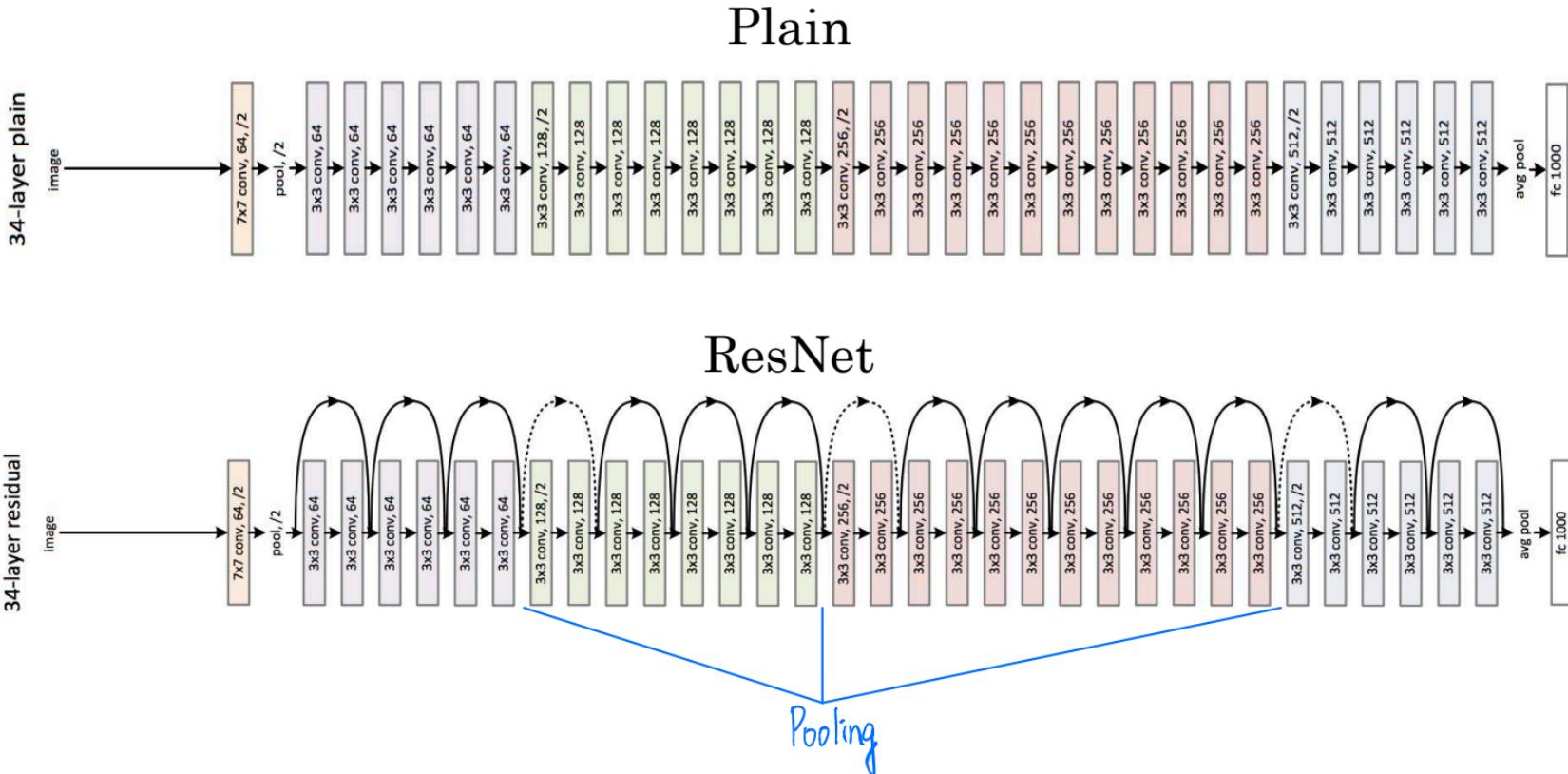
In case of  $z^{[l+2]}$  and  $a^{[l]}$  having different dimensions, e.g.,  $a^{[l]}$  is 128 and  $a^{[l+2]}$  is 256, we can add an extra matrix before  $a^{[l]}$ .

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

$R^{256 \times 128}$   
↑  
 $W_s$

or just use zero-padding to  $a^{[l]}$ .

## ResNet



[He et al., 2015. Deep residual networks for image recognition]

Andrew Ng

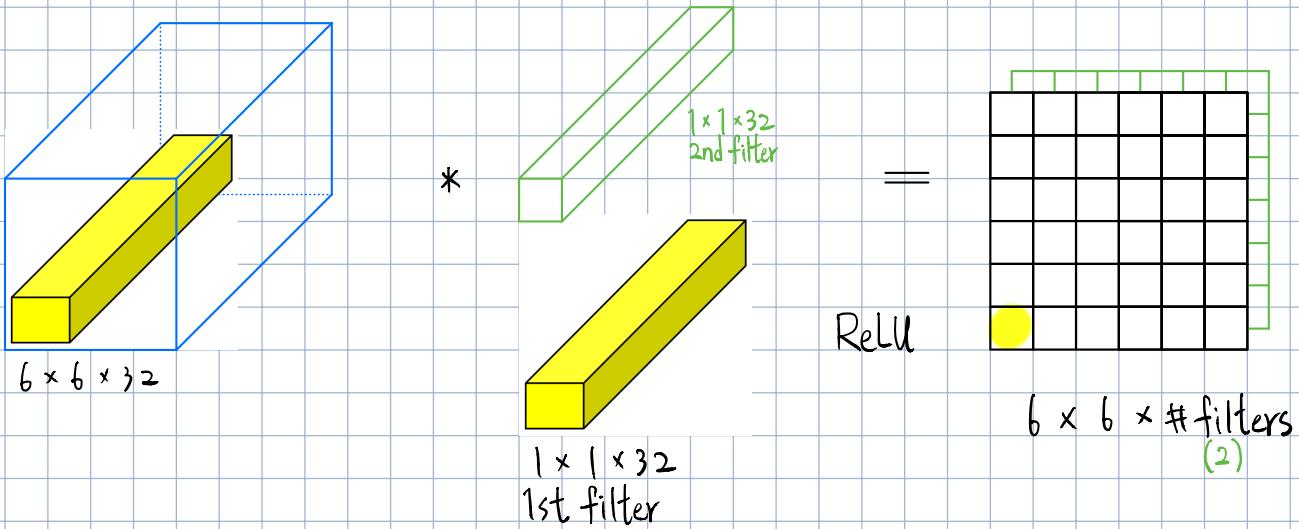
## Networks in Networks and $1 \times 1$ Convolutions

$$6 \times 6 \times 1$$

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

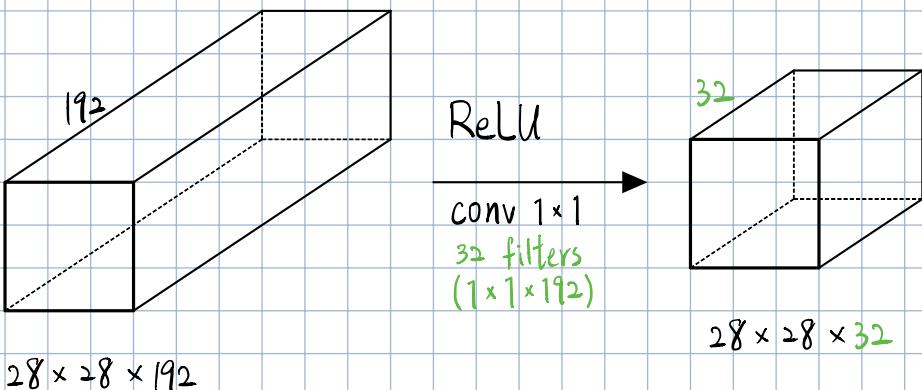
\*      [2]      =

2	4	6	...



Calculate the dot product of two  $1 \times 1 \times 32$  vectors, then apply ReLU on it.

Using  $1 \times 1$  convolutions in an example.

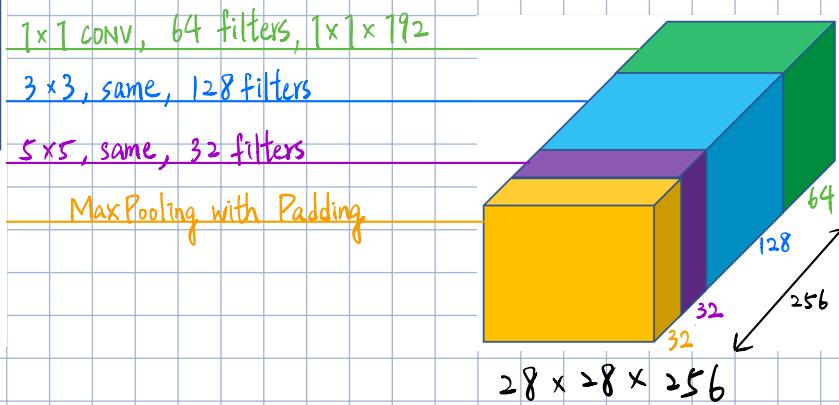
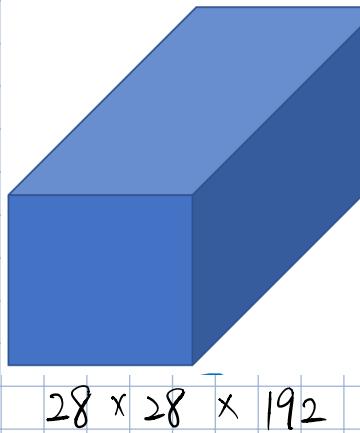


I know that the height and width can be shrunk by using pooling layer, what about # channels?

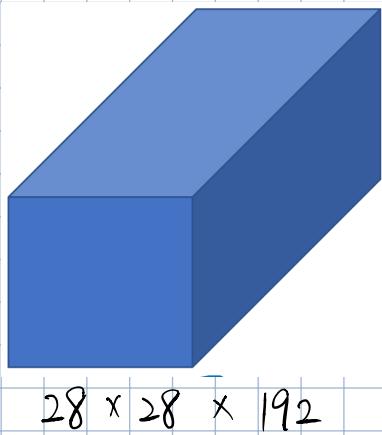
$1 \times 1$  convolutions can shrink the number of channels.

# Inception Network

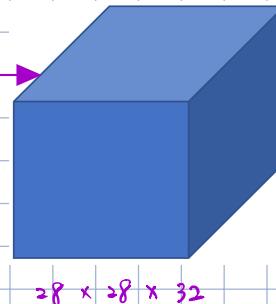
Motivation for Inception Network



The problem of computation cost



CONV  
 $5 \times 5$   
same  
32 filters

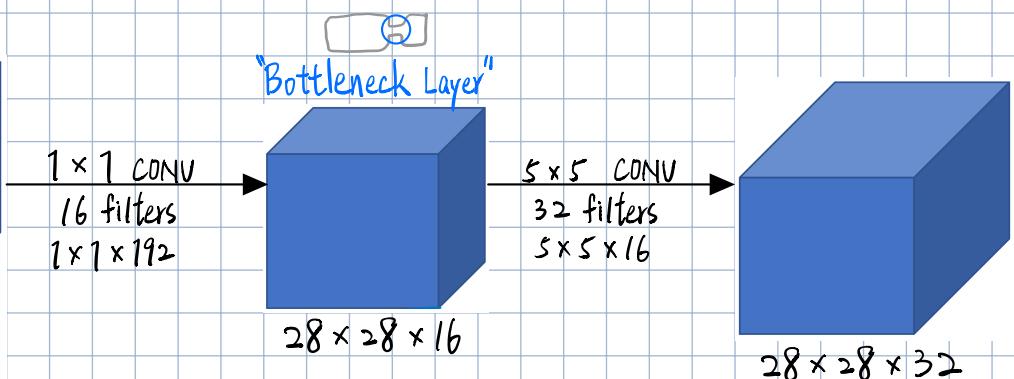
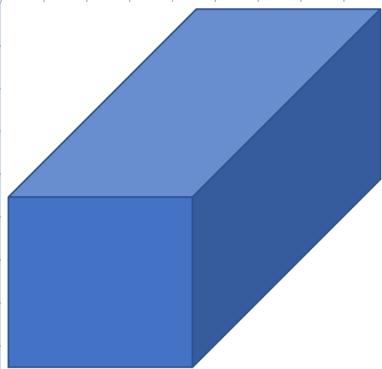


There are 32 filters. And each filter is  $5 \times 5 \times 192$ .

The computation cost is  $(28 \times 28 \times 32) \cdot (5 \times 5 \times 192) = 120 \cdot 10^6$

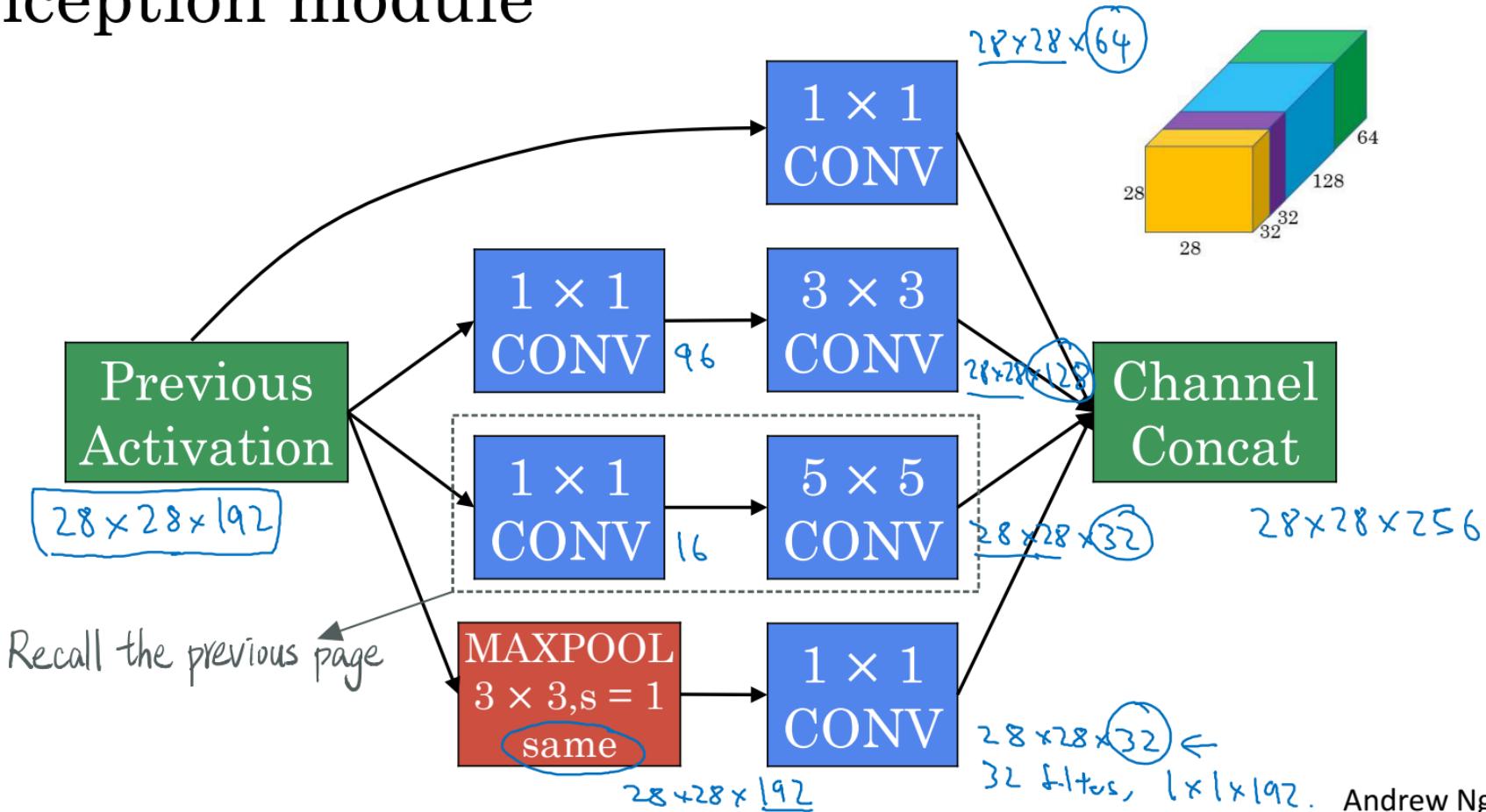
$120 \cdot 10^6$  is too large. We need to reduce it by using  $1 \times 1$  convolution below.

Using  $1 \times 1$  convolution

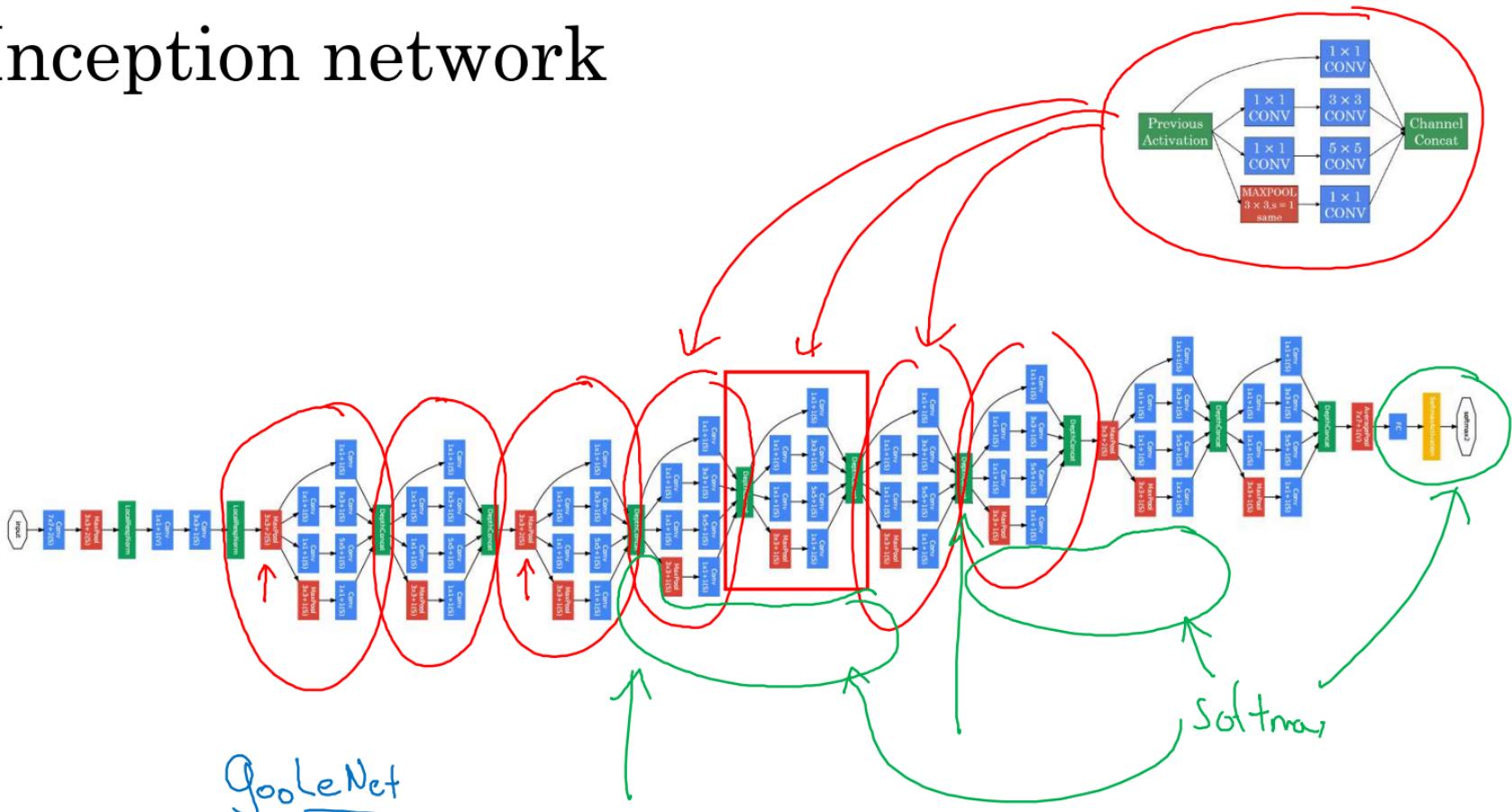


$$\begin{aligned} \text{Computation Cost: } & 28 \times 28 \times 16 \times 192 \\ & + 28 \times 28 \times 32 \cdot 5 \times 5 \times 16 \\ & = 2.4M + 10M = 12.4M \ll 120 \times 10^6 \end{aligned}$$

# Inception module



# Inception network



[Szegedy et al., 2014, Going Deeper with Convolutions]

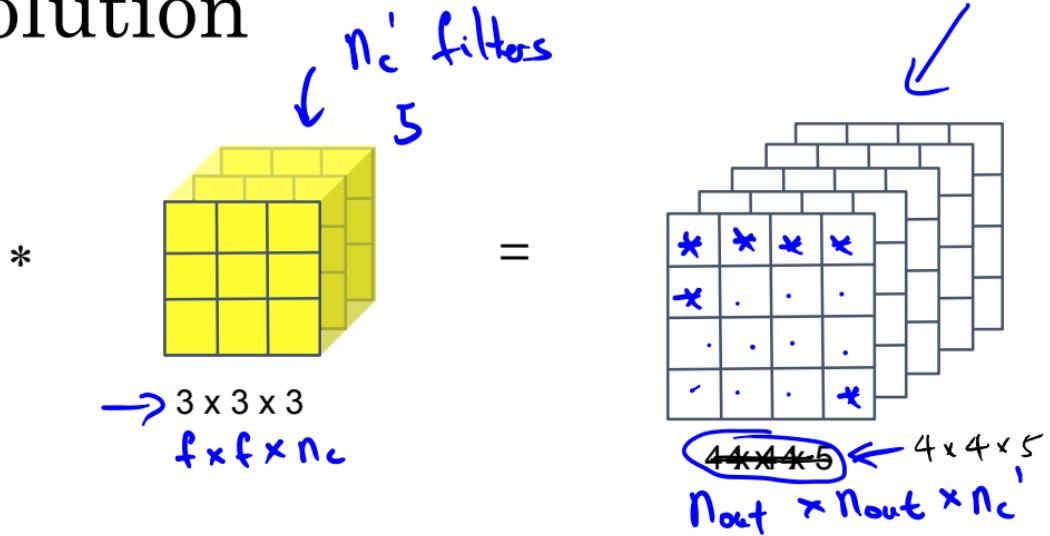
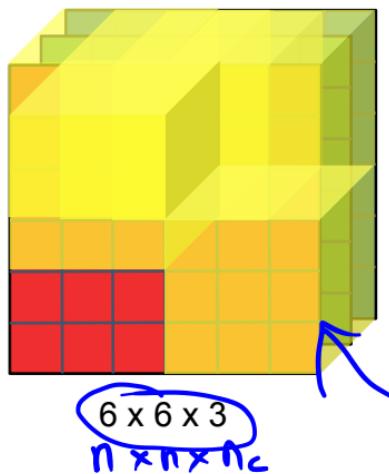
Andrew Ng

# Motivation for MobileNets

- Low computational cost at deployment
- Useful for mobile and embedded vision applications
- Key idea: Normal vs. depthwise-separable convolutions



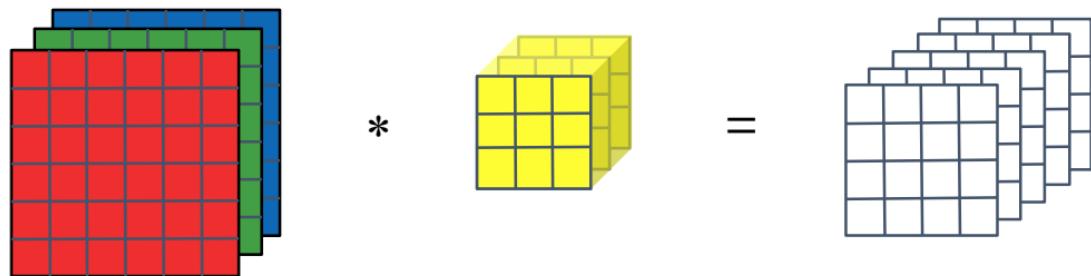
# Normal Convolution



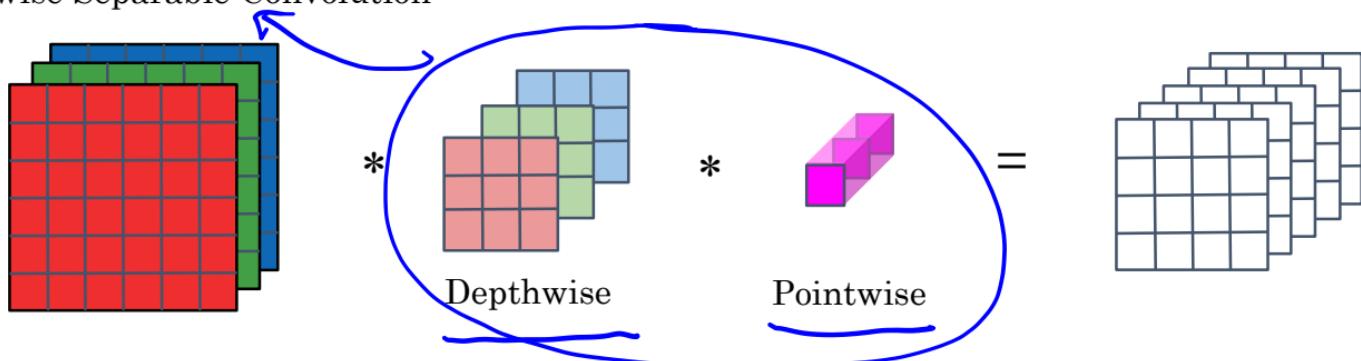
$$\begin{aligned}
 \text{Computational cost} &= \frac{\text{\#filter params}}{3 \times 3 \times 3} \times \frac{\text{\# filter positions}}{4 \times 4} \times \text{\# of filters } (n_c') \\
 \rightarrow \underline{2160} &= \underline{\uparrow} \quad \underline{\uparrow} \quad \underline{\uparrow} \\
 \end{aligned}$$

# Depthwise Separable Convolution

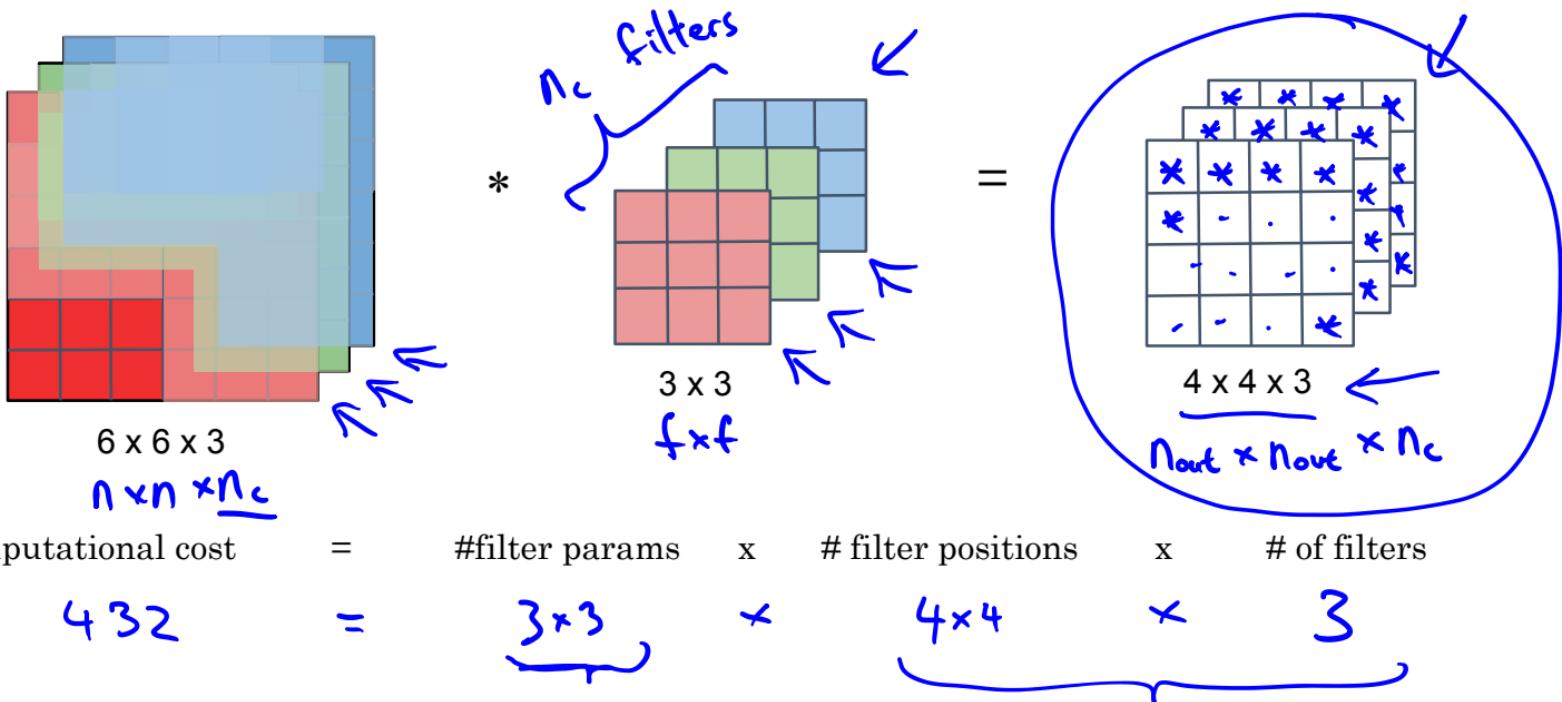
Normal Convolution



Depthwise Separable Convolution

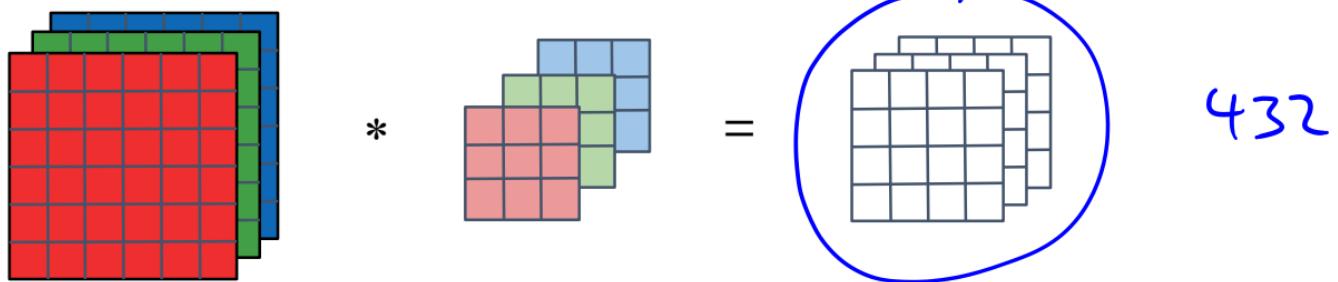


# Depthwise Convolution

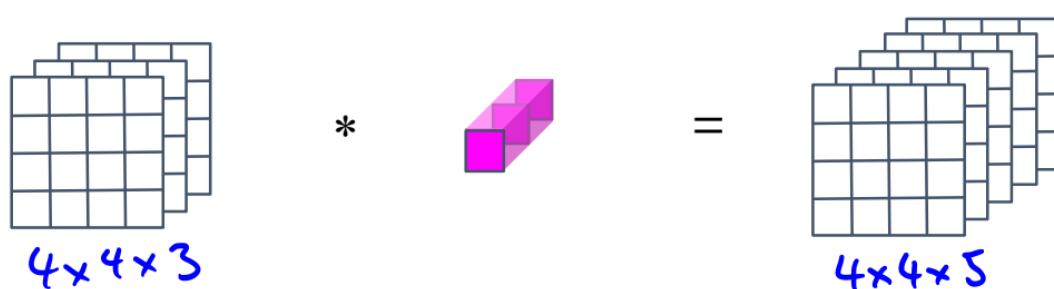


# Depthwise Separable Convolution

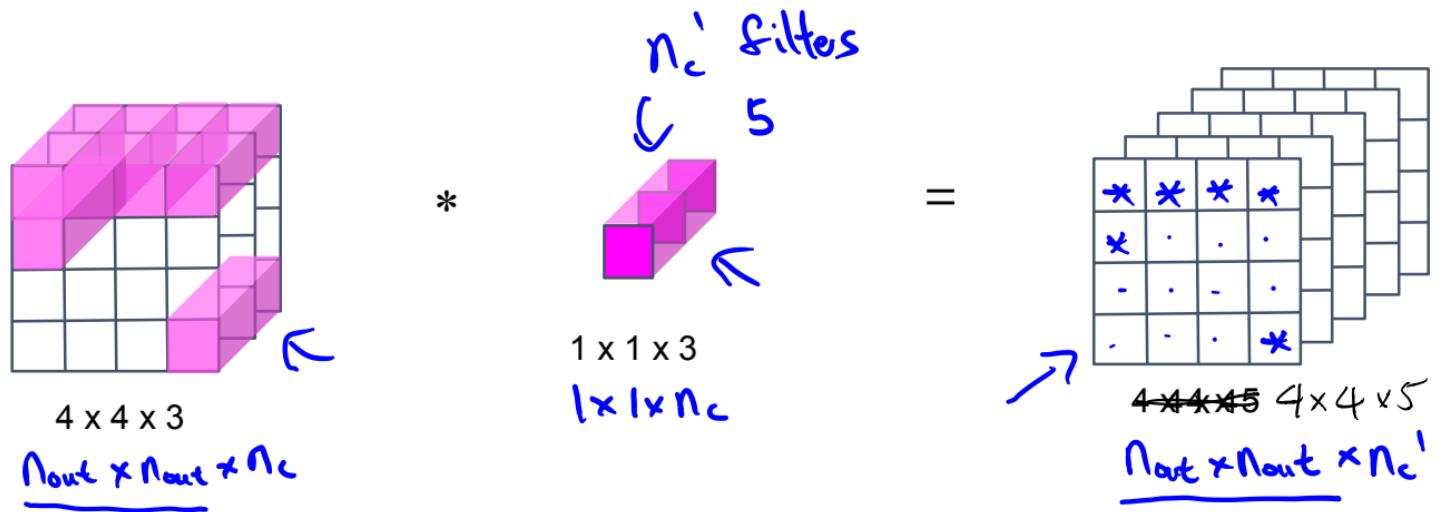
Depthwise Convolution



Pointwise Convolution



# Pointwise Convolution

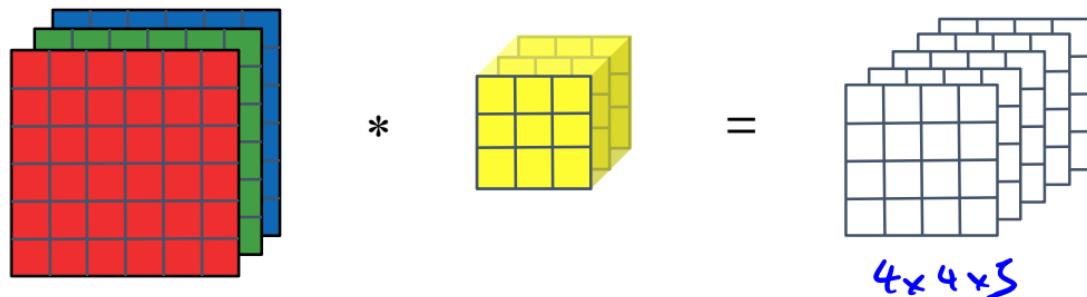


Computational cost = #filter params  $\times$  # filter positions  $\times$  # of filters

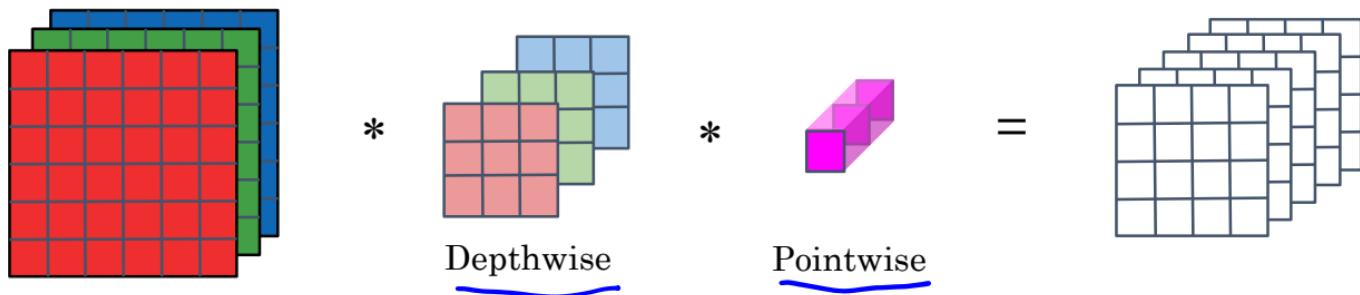
$$240 = 1 \times 1 \times 3 \times 4 \times 4 \times 5$$

# Depthwise Separable Convolution

Normal Convolution



Depthwise Separable Convolution



# Cost Summary

Cost of normal convolution

$$\leftarrow 2160$$

Cost of depthwise separable convolution

$$\begin{array}{l} \text{depthwise} + \text{pointwise} \\ 432 + 240 = 672 \end{array}$$

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

$$= \frac{1}{n_c} + \frac{1}{f^2}$$

$$\boxed{\frac{1}{5} + \frac{1}{9}}$$

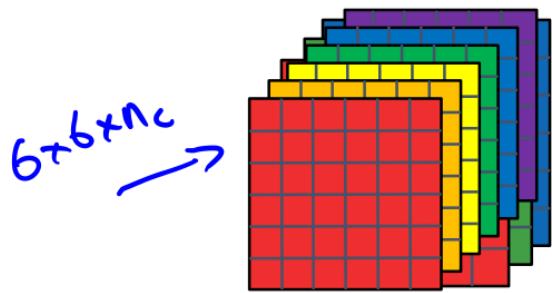
In our case,  
31% cheaper

$$= \frac{1}{512} + \frac{1}{3^2}$$

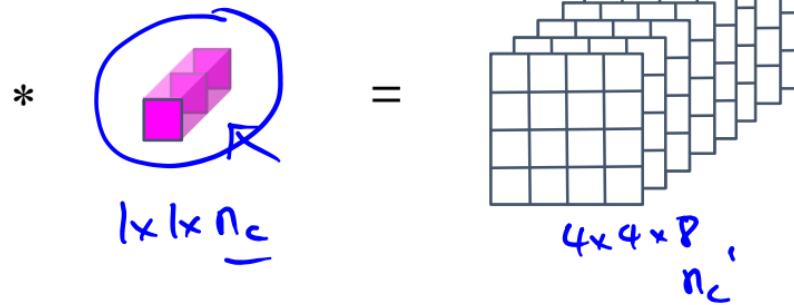
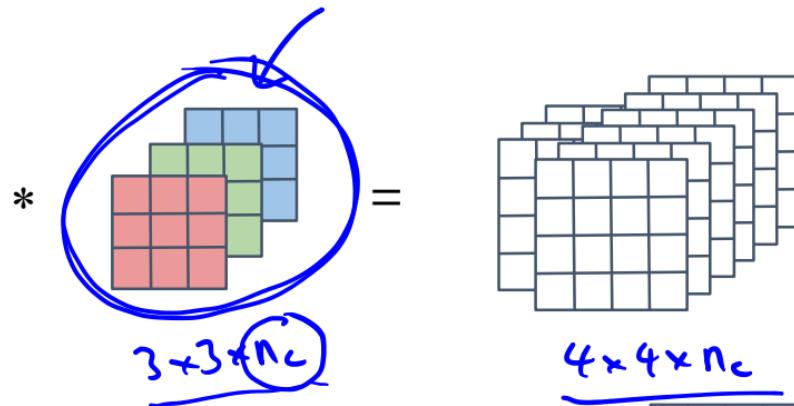
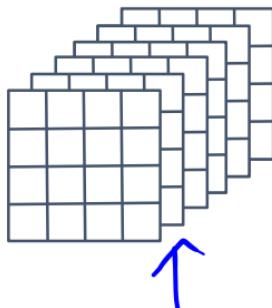
$\approx 10$  times cheaper

# Depthwise Separable Convolution

Depthwise Convolution



Pointwise Convolution





deeplearning.ai

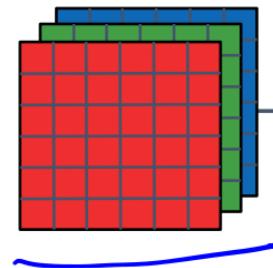
# Convolutional Neural Networks

---

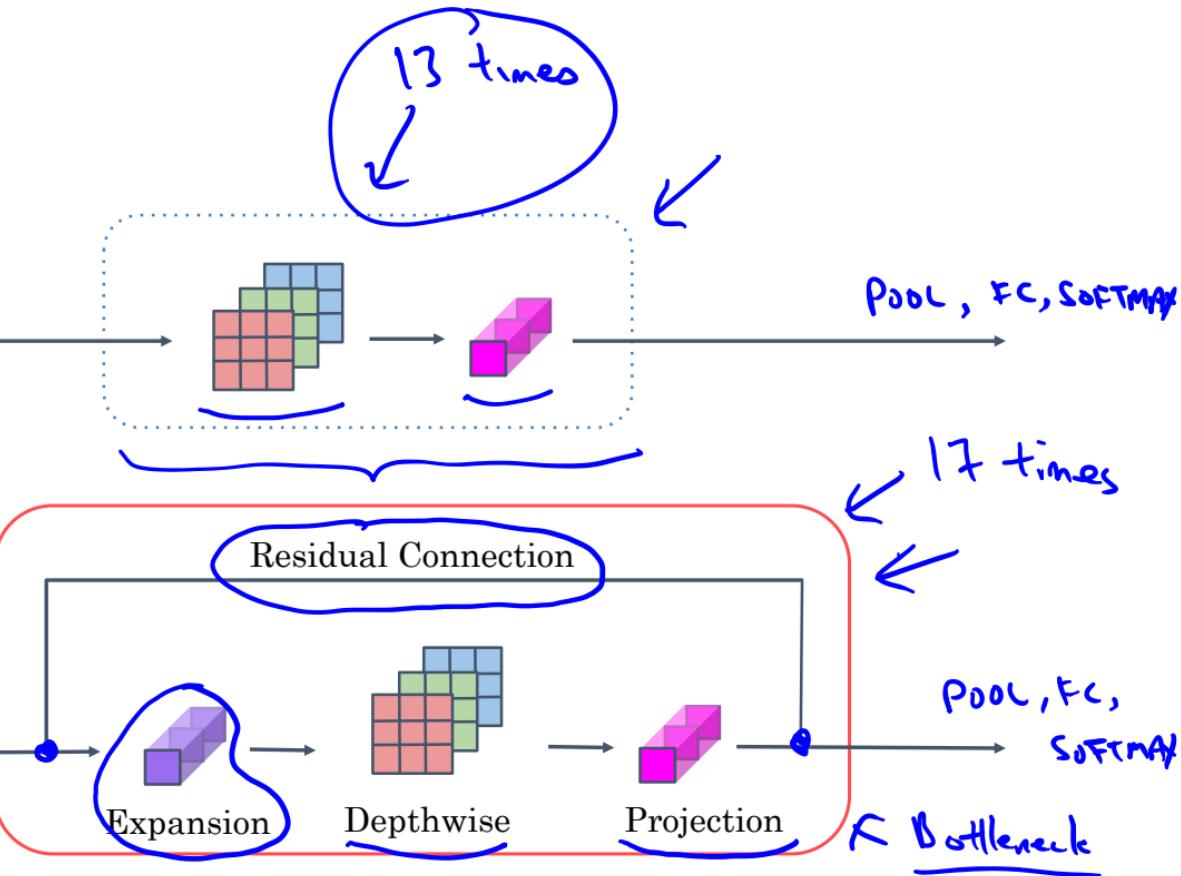
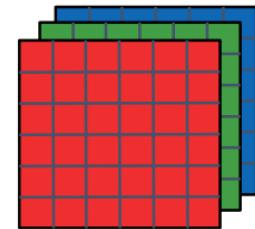
## MobileNet Architecture

# MobileNet

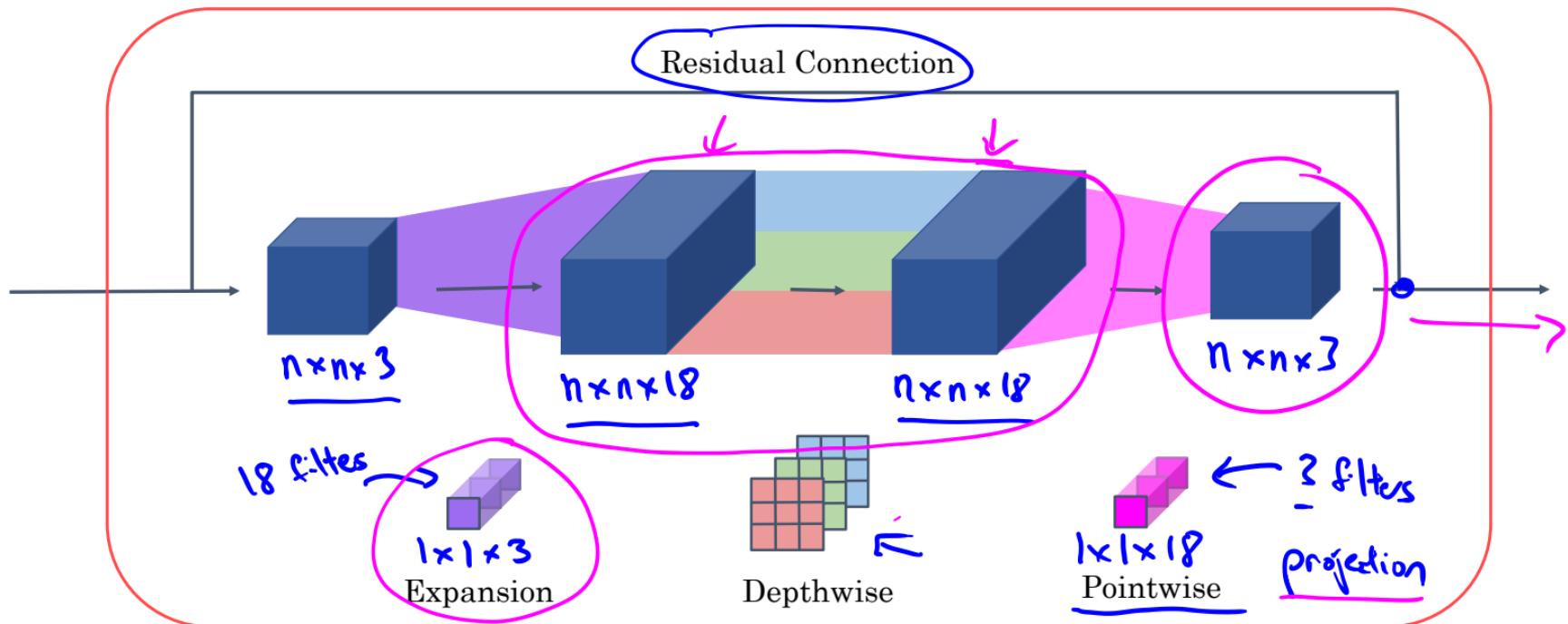
MobileNet v1



MobileNet v2

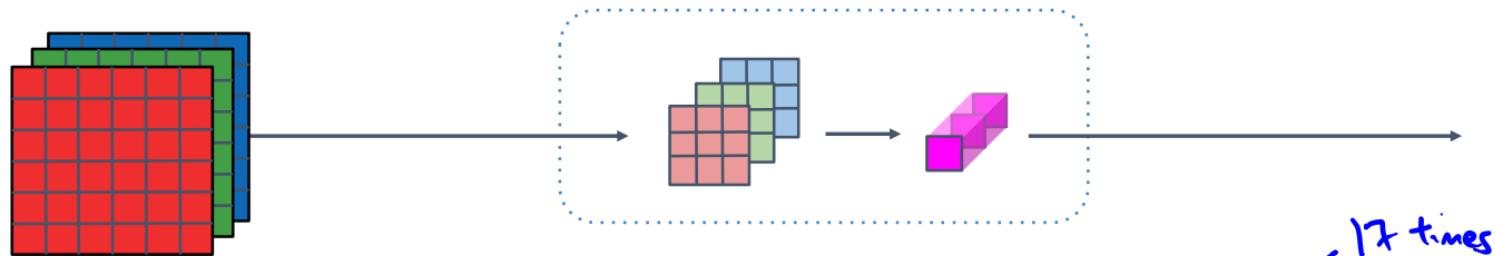


# MobileNet v2 Bottleneck

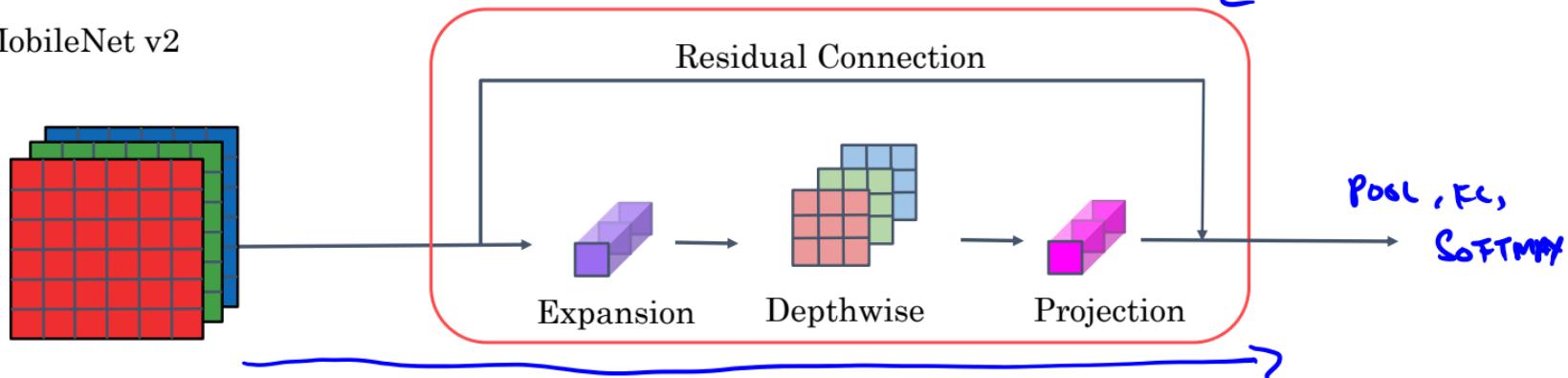


# MobileNet

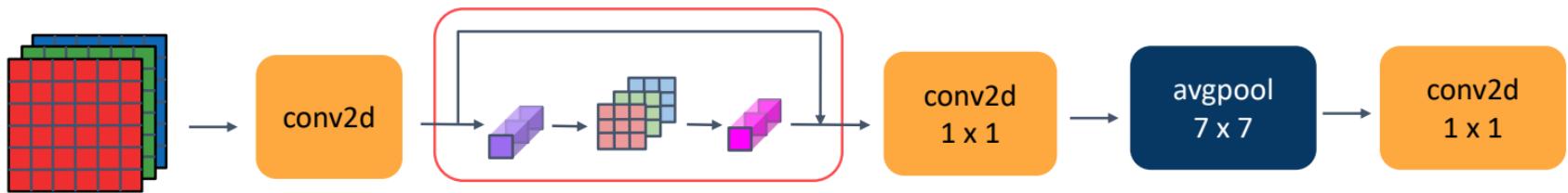
MobileNet v1



MobileNet v2



# MobileNet v2 Full Architecture





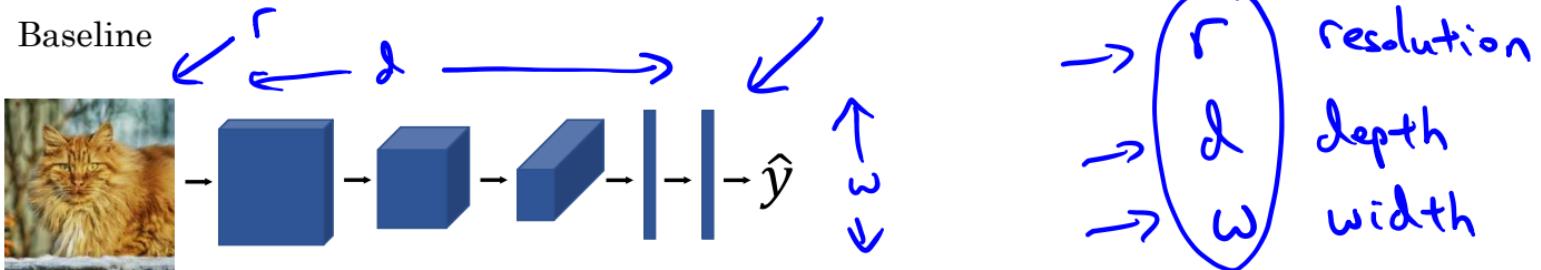
deeplearning.ai

# Convolutional Neural Networks

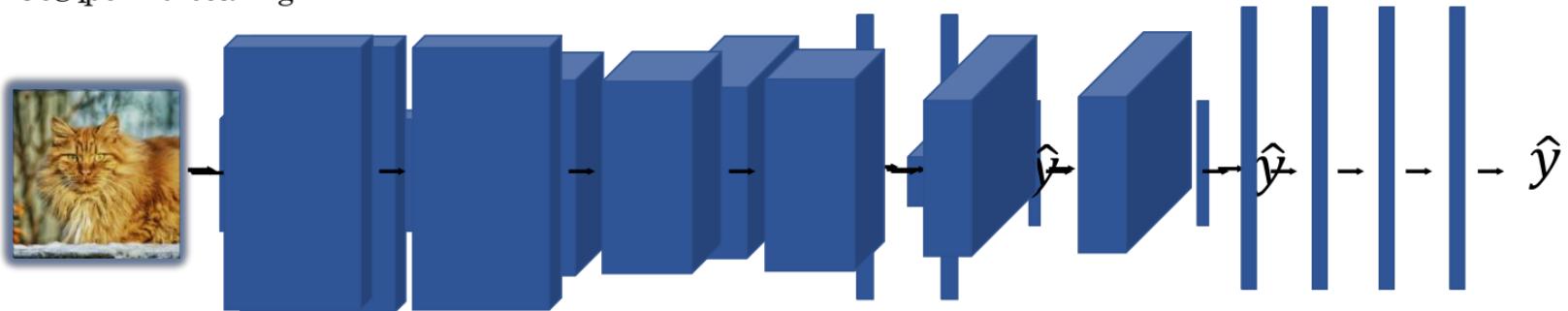
---

## EfficientNet

# EfficientNet



## Width Resolution Compound Scaling

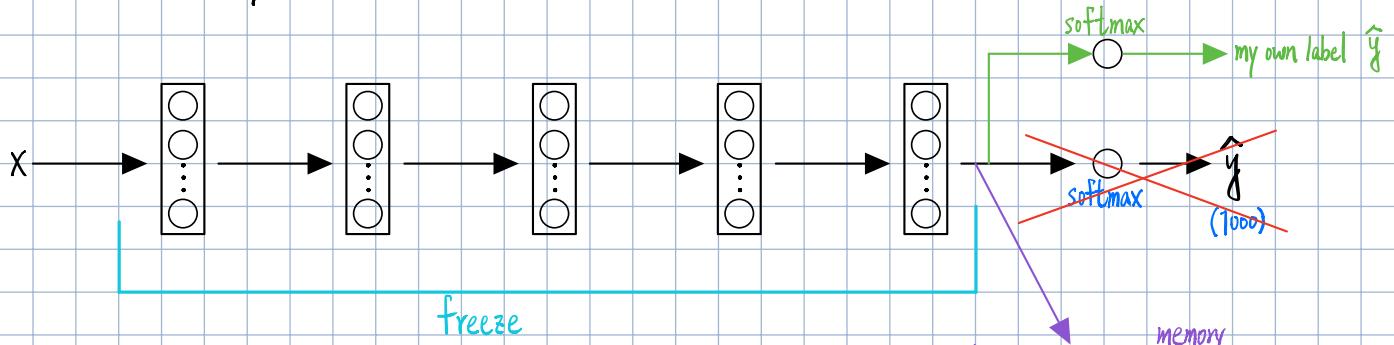


# Transfer Learning

*Andrew suggests strongly considering Transfer Learning for Computer Vision tasks.*

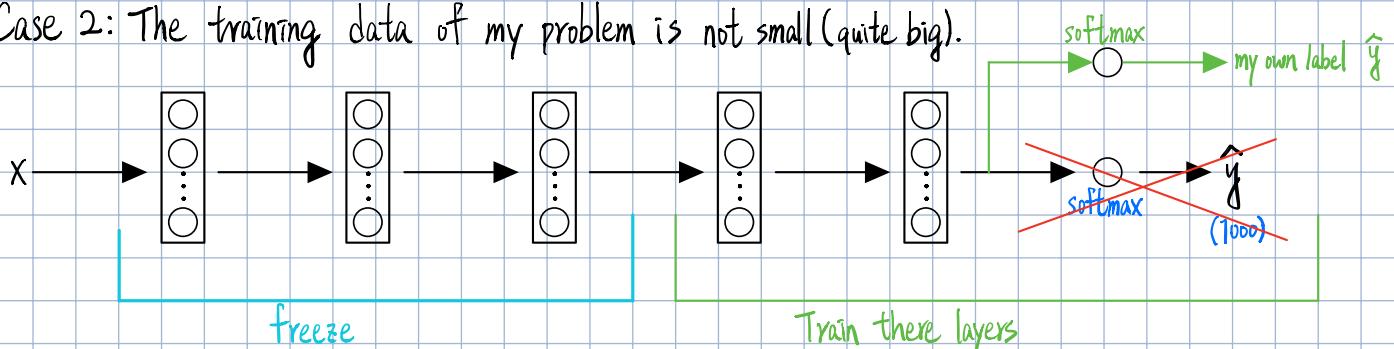
There are various ways to do transfer learning. Andrew categorizes different cases as follow:

Case 1: The training data of my problem is small. Freeze all front layers. Only train the softmax layer.



One trick to save computation time is to precompute the output values of the last frozen layer and save them to disk. During training, I just use the saved precomputed value as input to train the softmax layer.

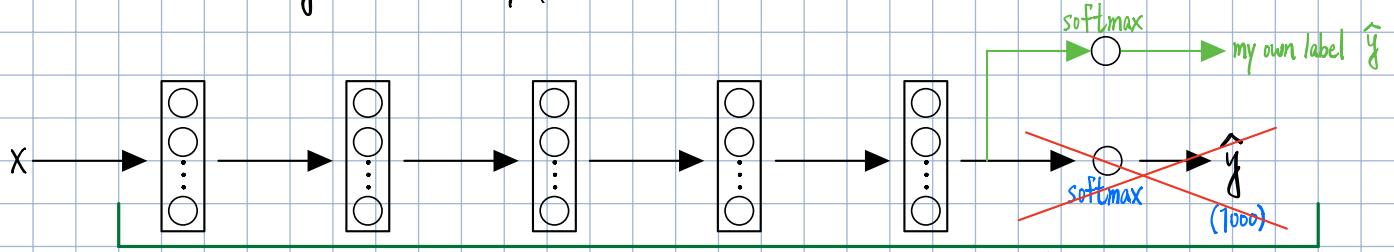
Case 2: The training data of my problem is not small (quite big).



Option 1: Use the pretrained weights as initialization

Option 2: Scratch the pretrained layers or weights.  
Start from scratch. Design new layers.  
Use random initialization.

Case 3: The training data of my problem is A LOT.



Use the whole pretrained weights as initialization and train the whole network.



deeplearning.ai

Practical advice for  
using ConvNets

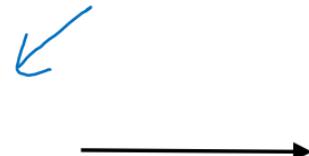
---

Data augmentation

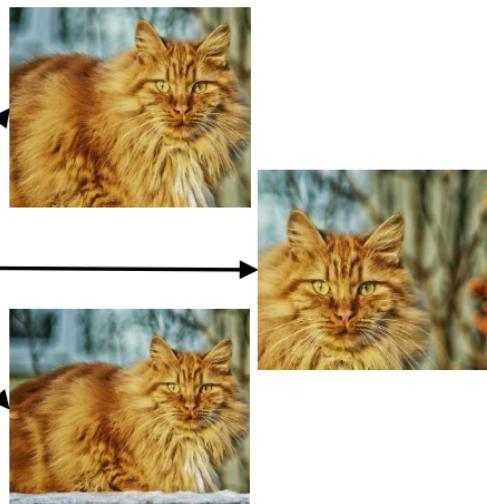
# Common augmentation method

Remember to use multithreading to implement data augmentation during Training.

Mirroring



Random Cropping



{ Rotation  
Shearing  
Local warping  
...



# Color shifting



y

R G B  
↓ ↓ ↓  
+20,-20,+20



-20,+20,+20



+5,0,+50



Advanced:

PCA

ml-class.org

[ AlexNet paper

[ "PCA color  
augmentation."

R B

G

# The State of Computer Vision

- Use open source implementations of networks architectures published in the literature.
- Use pre-trained models and fine-tune on your dataset.

Little Data

More hand-crafted features

Utilize Transfer Learning  
(freeze front layers and train rear layers)

Object  
Detection

Image  
Recognition

Speech  
Recognition

Lots of Data

Simpler Algorithms

Less hand-crafted features  
Giant neural networks

Machine Learning algorithms can be seen as having two sources of knowledge:

1. Labeled data  $(x, y)$  pair.
2. Hand engineered features/network architecture/other components

Andrew says hand engineering data is getting less important if there are lots of data.

I agree: Lots of data can improve accuracy a lot. It's definitely good.

I do NOT agree: Hand engineering insights can help to optimize a system's performance (memory usage, CPU usage, accuracy) with the same amount of data. It doesn't matter there is a lot or a few

data. As long as hand engineering insights are used in a system, I am sure that my system has reached its theoretically best performance.  
(car)

Eg, in object detection, an ML algorithm with dynamic system model is definitely better than a "pure" ML algorithm without the knowledge of the dynamic system model.

Tips for doing well on benchmarks / winning competition. Andrew also says these two tips will probably not be helpful for deploying an actual application. So, these two tips are probably irrelevant to me. Another implication is that be suspicious about the results in papers. Papers will just use whatever means they can to show they are better when in actuality they are not. Thus, I must learn how to recognize such papers and throw them away. Anyway, here are the two tips:

- Ensembling: Train several networks independently and average their outputs  $\hat{y}$ .
- Multi-crop at test time (10-crop): Run classifier on multiple versions of test images and average results.