



CS3232 Fundamental of Deep learning

Credits: Parts of the contents are based on material by Dr. Andrew Ng and other resources on the Internet

Classic networks (architectures)

- LeNet-5 (60K parameters)
- AlexNet (60M parameters)
- VGG (VGG 16 has ~ 138 M parameters)

ResNet (152 layers) – some tricks on how to train them effectively

AlexNet: Introduced ReLU, dropout, and data augmentation; won ILSVRC 2012.

ZF-Net: Improved on AlexNet with better visualization techniques; won ILSVRC 2013.

VGGNet: Emphasized depth with small filters; influential design; top performance in ILSVRC 2014.

GoogLeNet: Introduced Inception modules; reduced parameters; won ILSVRC 2014.

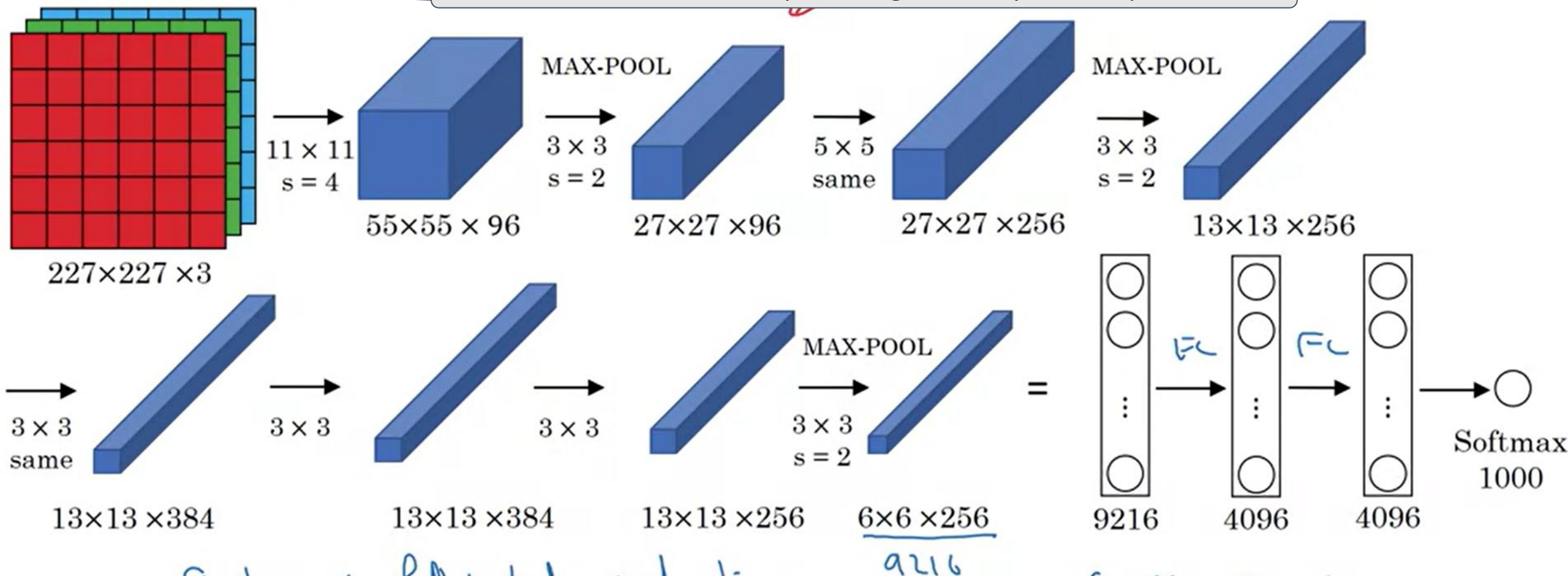
ResNet: Introduced residual connections; enabled very deep networks; won ILSVRC 2015.

Transfer Learning: Using pre-trained models to leverage learned features for new tasks.

Fine-Tuning: Adapting pre-trained models to specific tasks by continuing training on new data

AlexNet

Led to increased use of deep learning, even beyond computer vision



- Similar to LeNet, but much bigger.

- ReLU

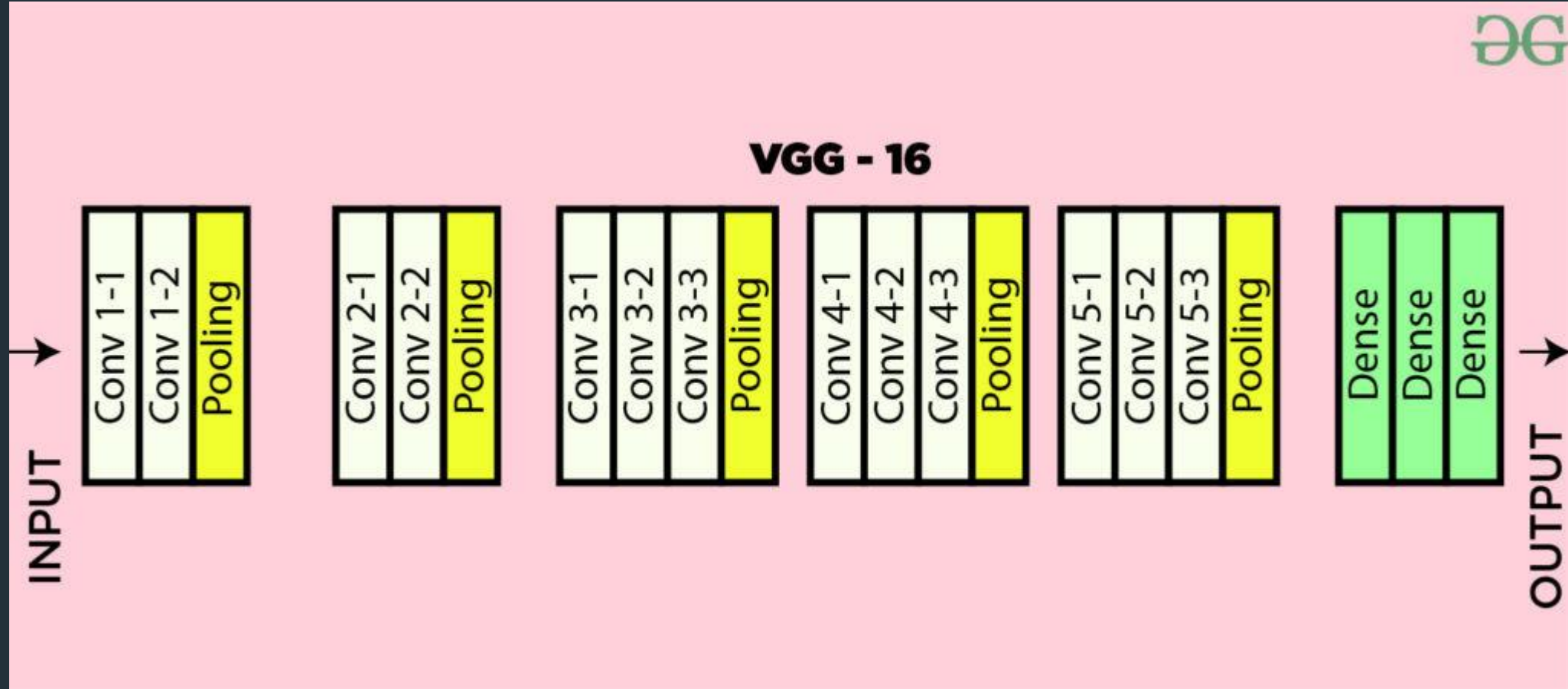
- Multiple GPUs.

- Local Response Normalization.

Not used any longer (as it does not help much)

~60M parameters

VGG Net



VGG-16

The VGG-16 architecture is a deep convolutional neural network (CNN) designed for image classification tasks. It was introduced by the Visual Geometry Group at the University of Oxford. VGG-16 is characterized by its simplicity and uniform architecture, making it easy to understand and implement.

The VGG-16 configuration typically consists of 16 layers, including 13 convolutional layers and 3 fully connected layers. These layers are organized into blocks, with each block containing multiple convolutional layers followed by a max-pooling layer for downsampling.

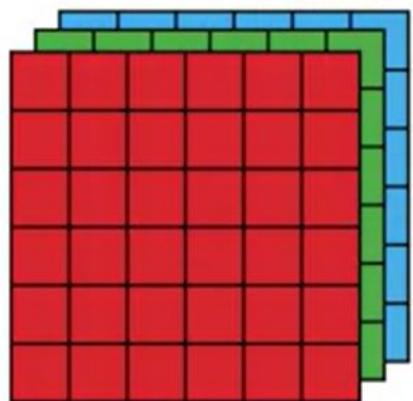
VGG - 16

16 layers with weights

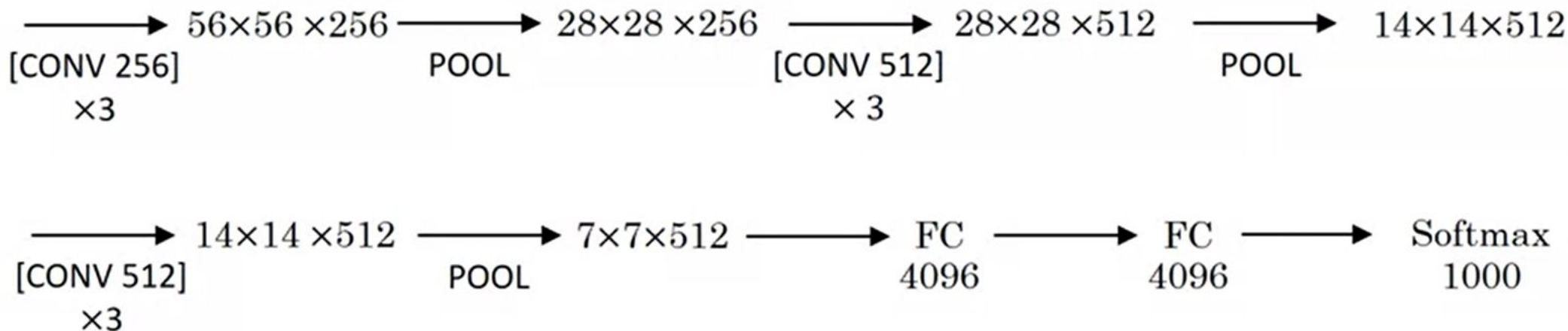
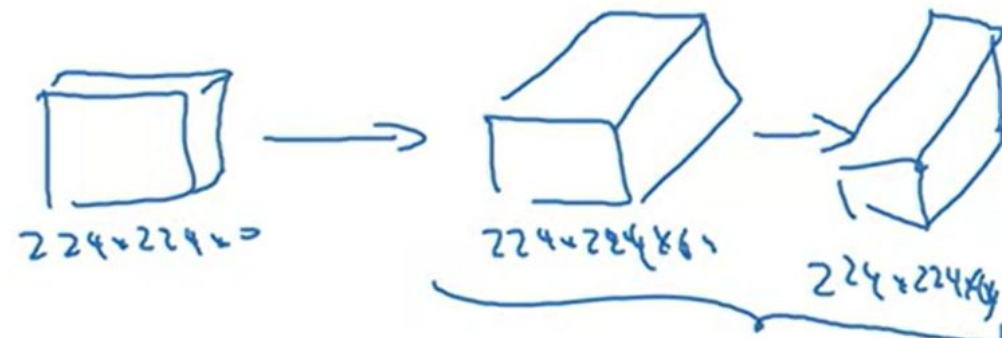
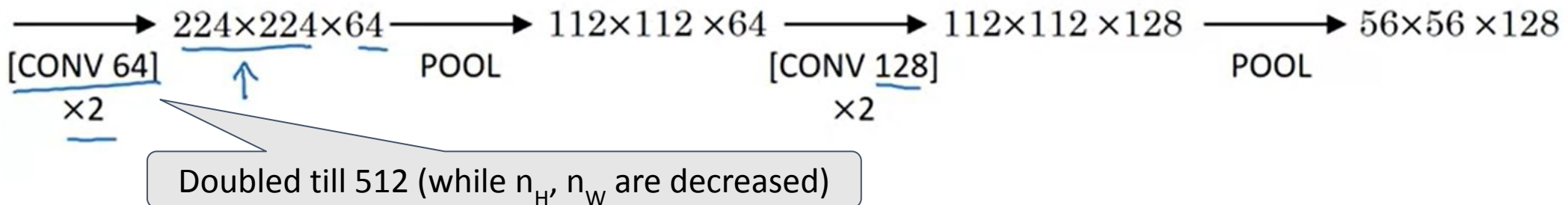
CONV = 3x3 filter, s = 1, same

MAX-POOL = 2x2, s = 2

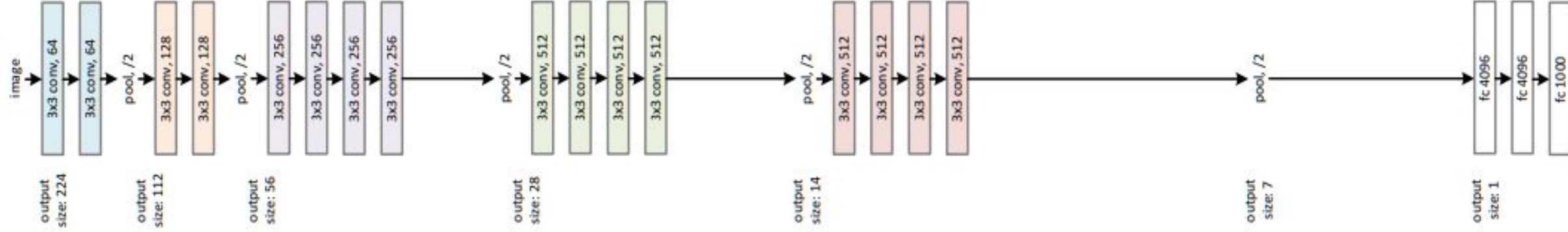
Simplified architecture (uniform)



224x224 x 3



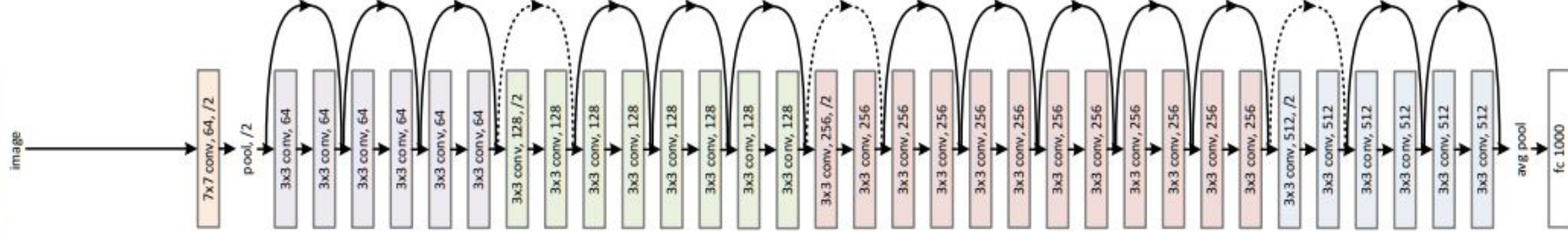
VGG-19



34-layer plain



34-layer residual



VGG-16 vs. VGG-19

- Since VGG-16 almost does as well as VGG-19, it is used by many

Model	Year	Architecture	Key Features	Impact
AlexNet	2012	8 layers (5 conv, 3 fully connected)	ReLU, dropout, data augmentation, won ILSVRC 2012	Popularized deep learning in computer vision
ZF-Net	2013	Similar to AlexNet with modified filter sizes	Deconvolutional layers for visualization, improved on AlexNet	Highlighted importance of feature visualization
VGGNet	2014	16-19 layers, 3x3 conv filters, 2x2 max pooling	Emphasized depth, consistent filter sizes, achieved top results	Influenced many future network designs with its simplicity
GoogLeNet (Inception)	2014	22 layers, Inception modules combining multiple filters	Reduced parameters significantly, auxiliary classifiers	Efficient network design, won ILSVRC 2014
ResNet	2015	50, 101, 152 layers with residual blocks	Residual connections, mitigated vanishing gradient problem	Enabled training of very deep networks, won ILSVRC 2015

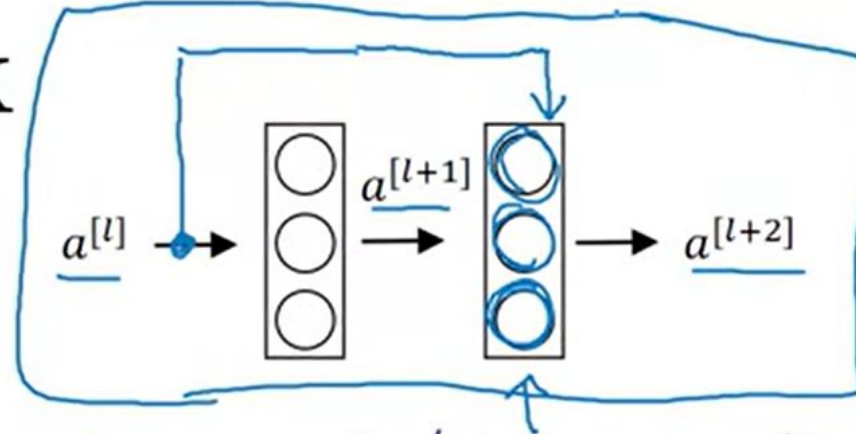
Modern networks

- ResNet
- Inception network (GoogLeNet)
- MobileNet

ResNet

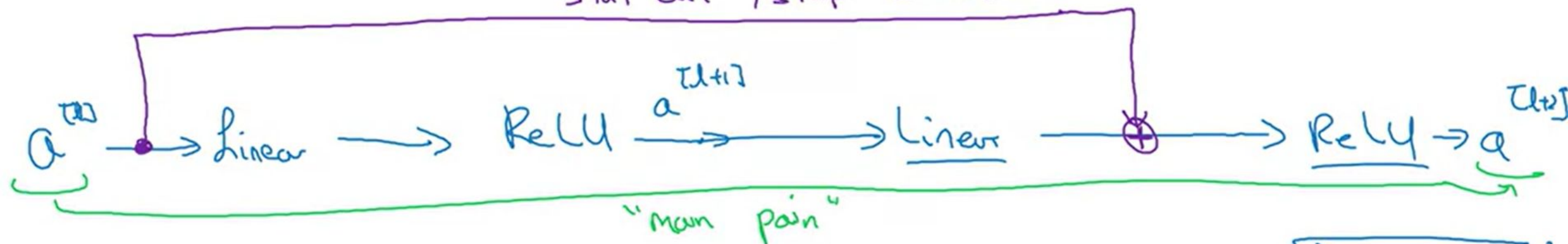
- Very deep networks (> 100 layers) are hard to train because of problems
 - Vanishing, exploding gradients etc.
- Skip connections
 - Take activation of one layer and feed it in another layer much deeper in the network
- Use of residual blocks help with training very deep networks

Residual block



Stacked together to form ResNet

"short cut" / skip connection



$$\underline{z^{[l+1]}} = \underline{W^{[l+1]}} \underline{a^{[l]}} + \underline{b^{[l+1]}}$$

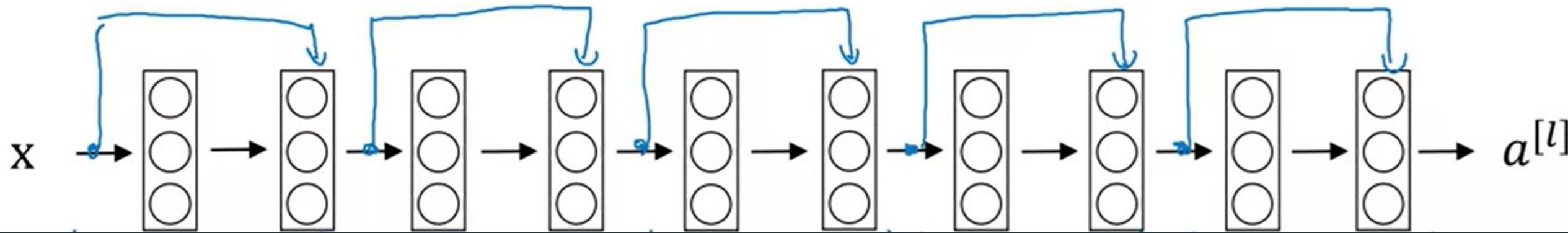
$$\underline{a^{[l+1]}} = g(\underline{z^{[l+1]}})$$

$$\underline{z^{[l+2]}} = \underline{W^{[l+2]}} \underline{a^{[l+1]}} + \underline{b^{[l+2]}}$$

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

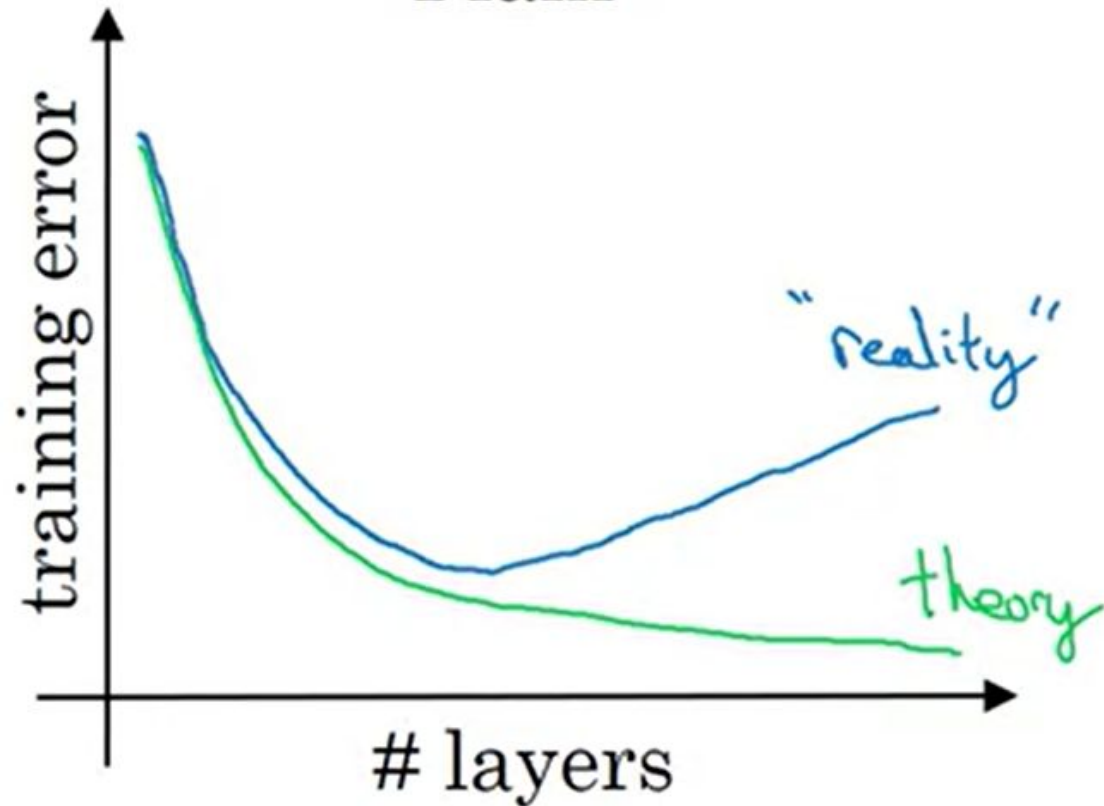
$$a^{[l+2]} = g(z^{[l+2]} + \underbrace{a^{[l]}})$$

Residual Network

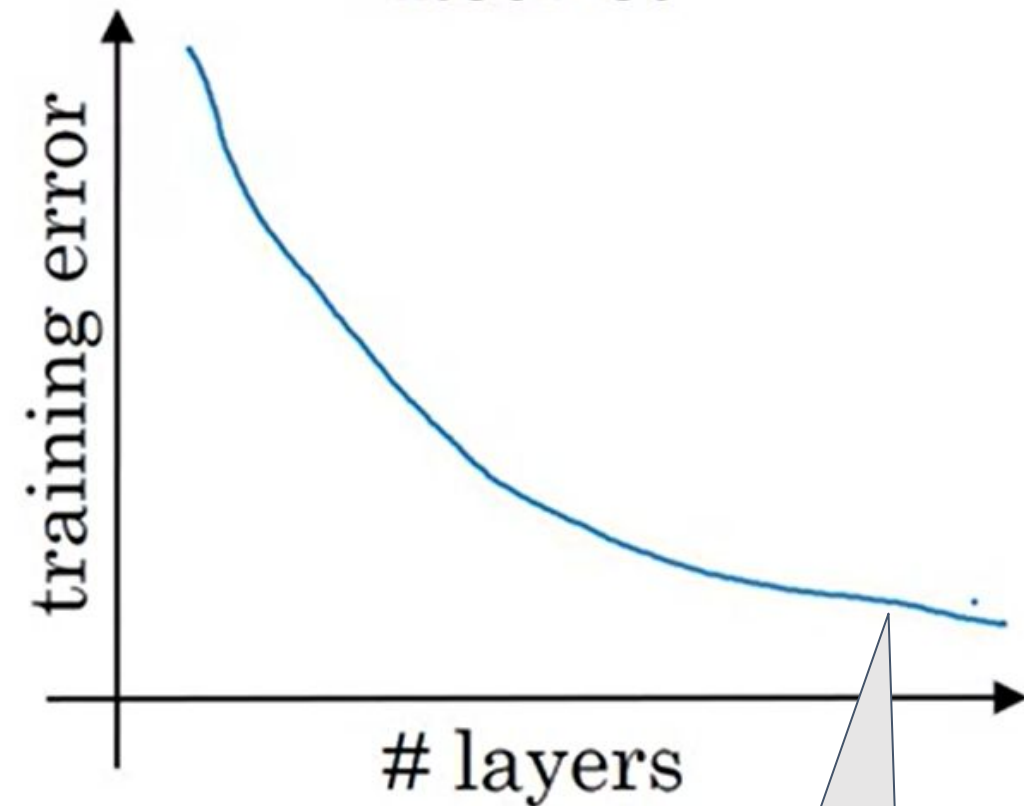


ResNet

Plain

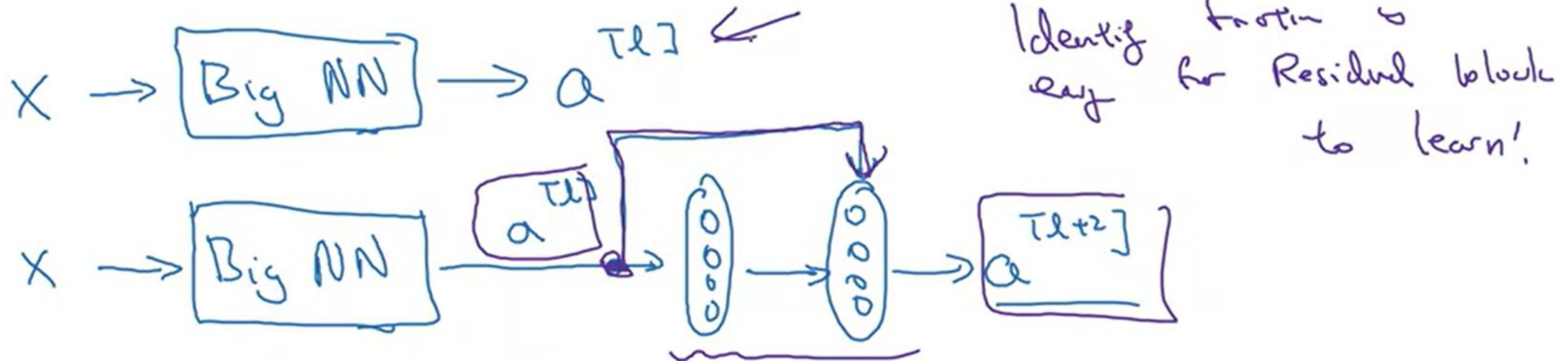


ResNet



Plateau

Why do residual networks work?



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

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

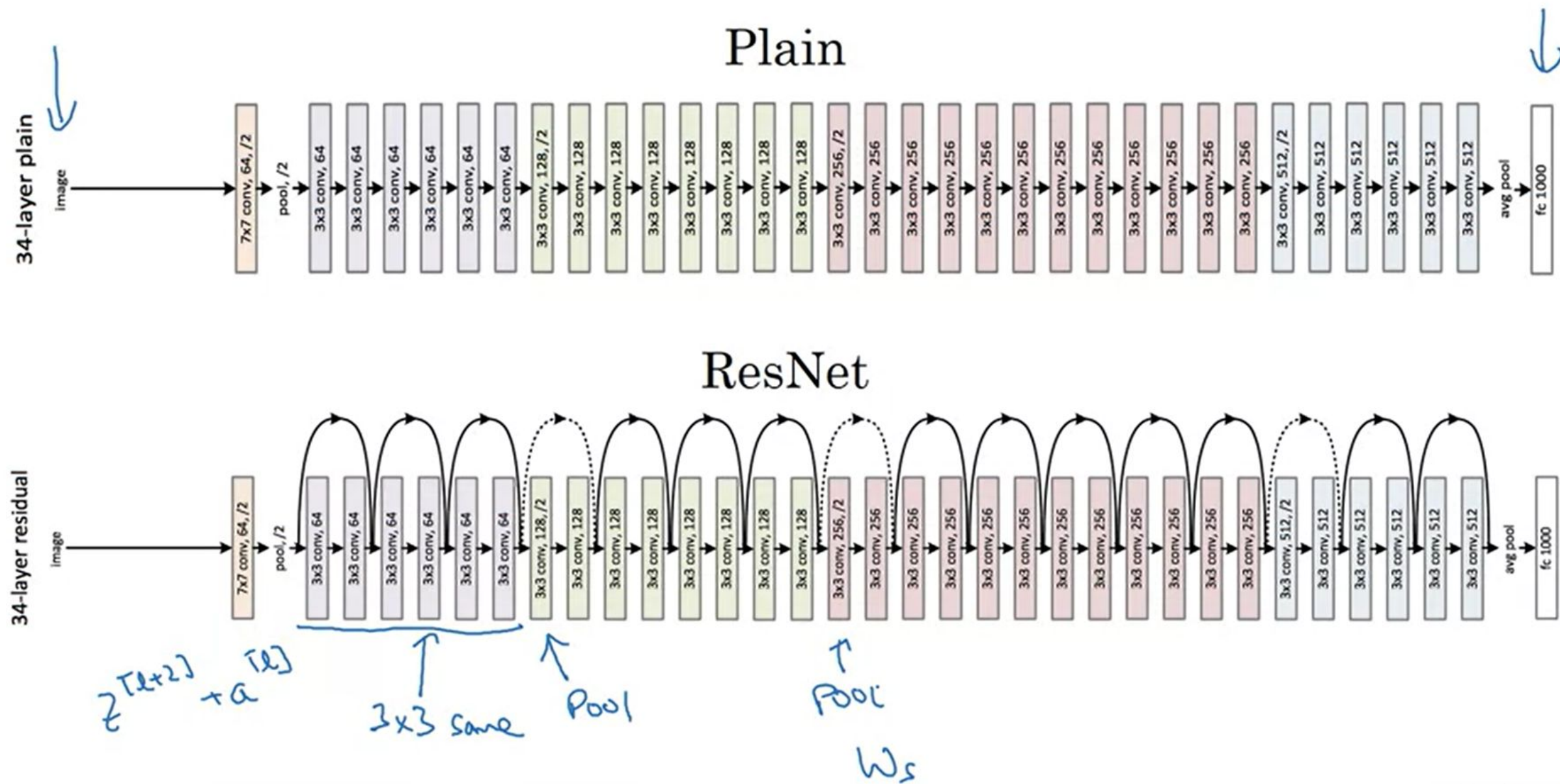
If $w^{[l+2]} = 0, b^{[l+2]} = 0$

Does not hurt performance because $a^{[l+2]} = a^{[l]}$ but may help too

To have same dimensions, "same" convolution is used

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

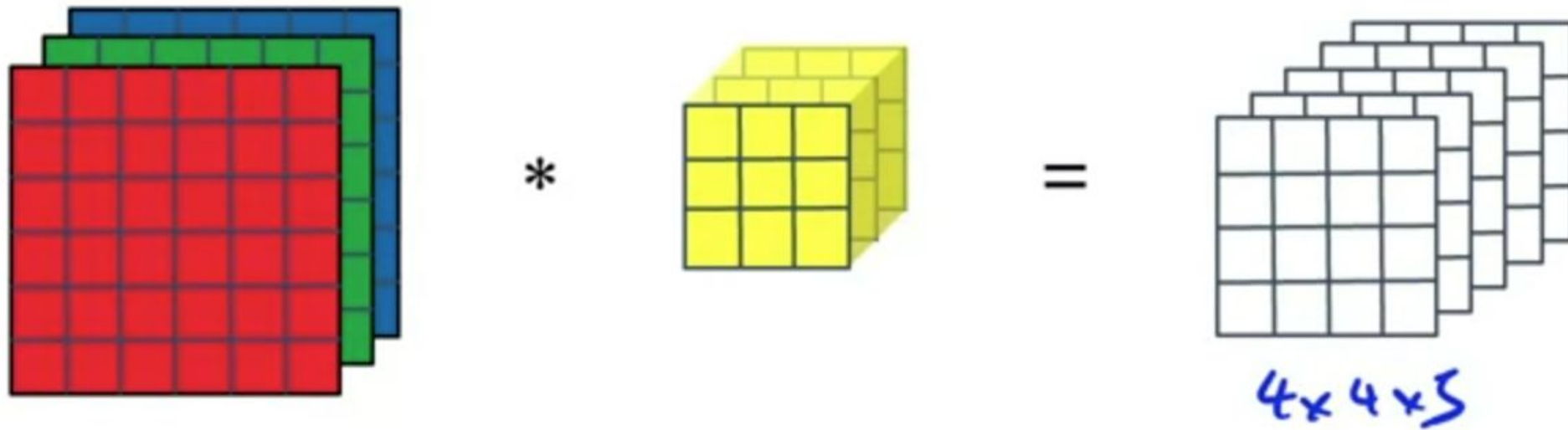
ResNet



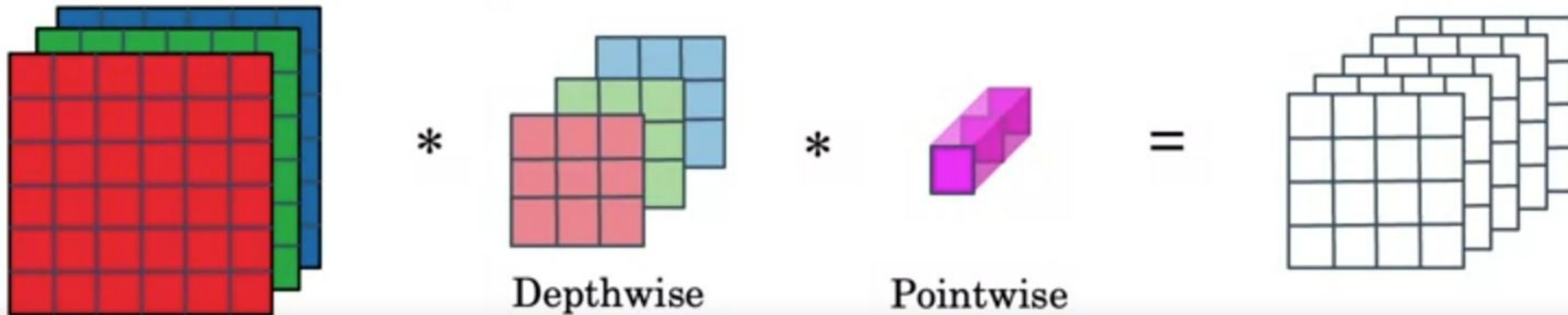
MobileNet

- Can be built and deployed even in low compute environments
 - Mobile phones, embedded vision applications
- Key idea
 - Normal vs. depth-wise separable convolutions

Normal Convolution

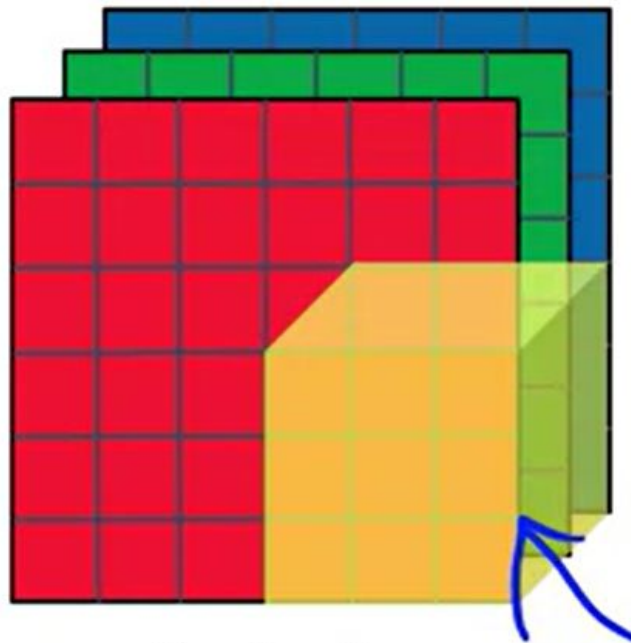


Depthwise Separable Convolution



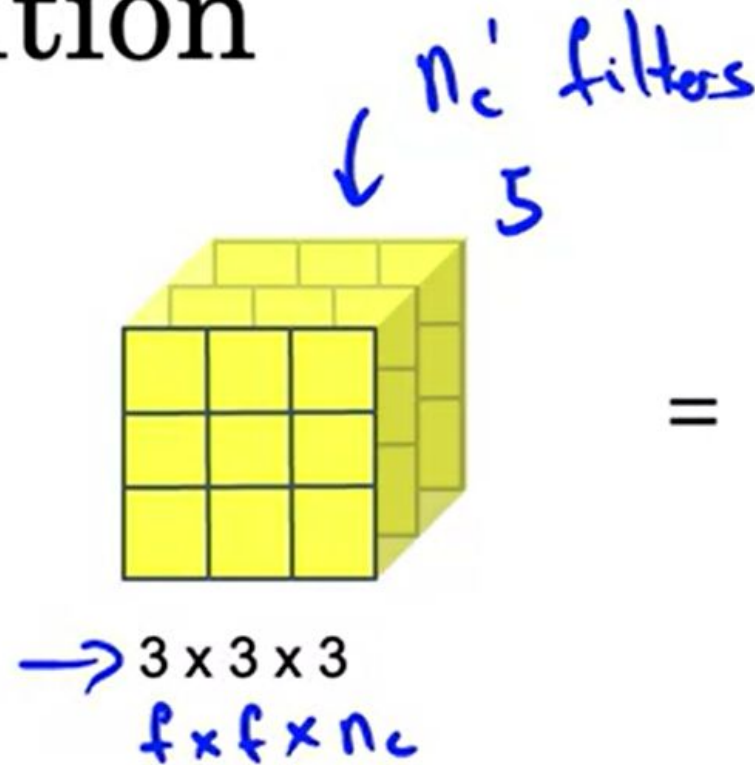
Computations

Normal Convolution

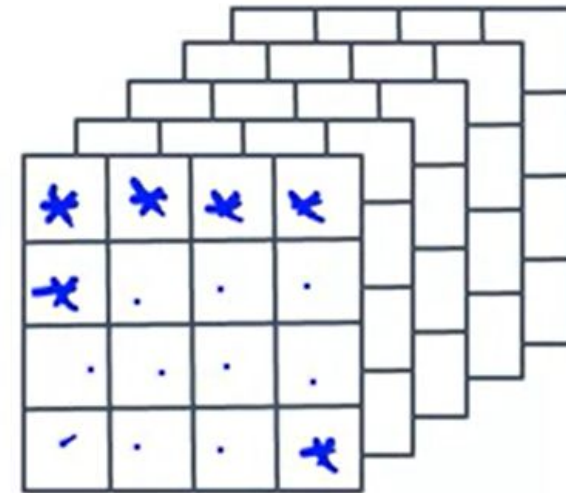


$6 \times 6 \times 3$
 $n \times n \times n_c$

*



=

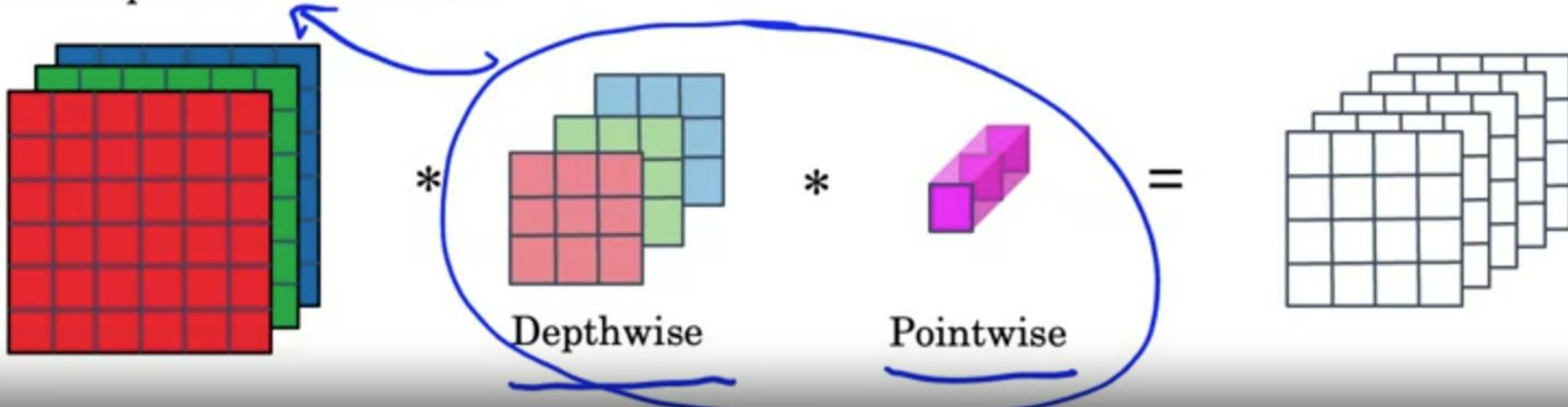


$4 \times 4 \times 5$
 $n_{out} \times n_{out} \times n_c'$

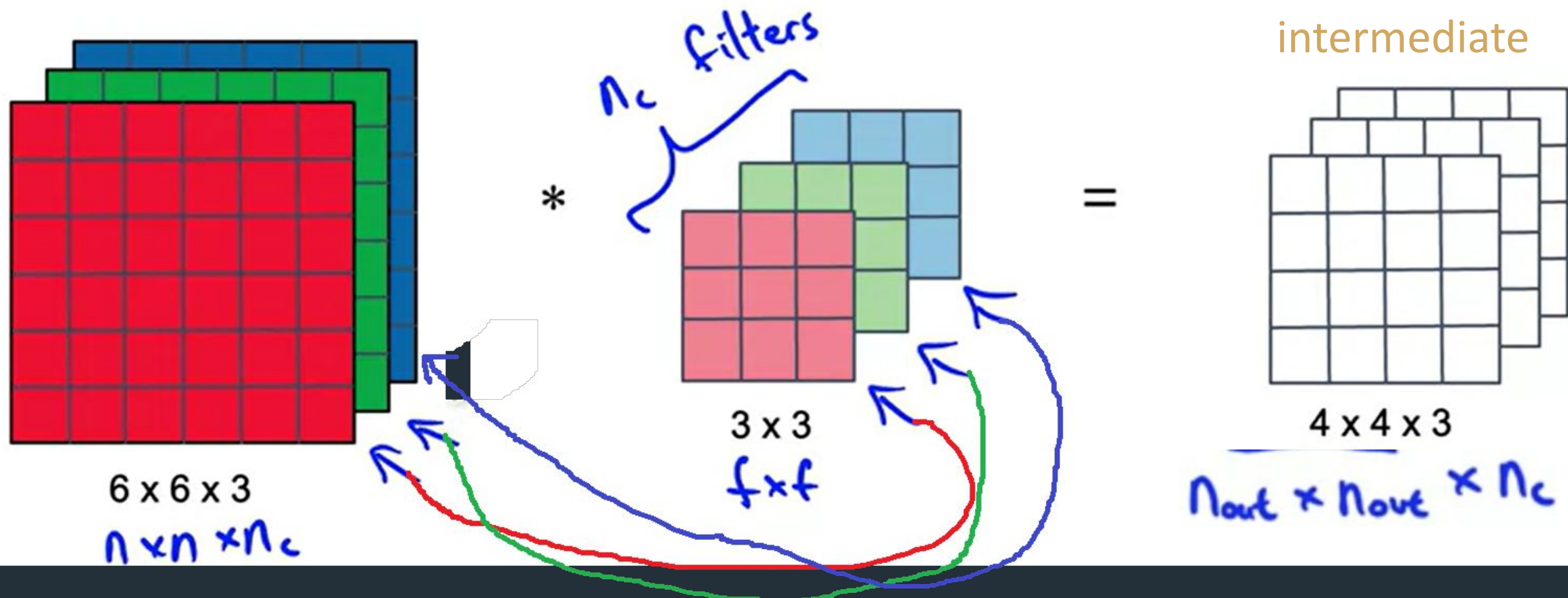
Computational cost = $\frac{\text{\#filter params}}{3 \times 3 \times 3} \times \frac{\text{\# filter positions}}{4 \times 4} \times \text{\# of filters}$

$2160 = 3 \times 3 \times 3 \times 4 \times 4 \times 5$

Depthwise Separable Convolution

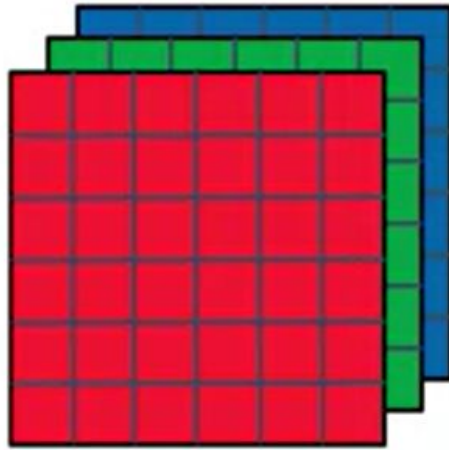


Depthwise Convolution

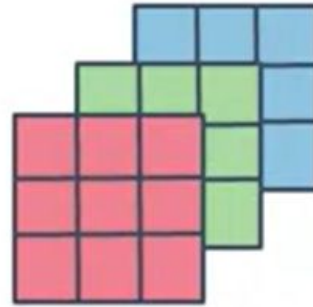


Computational cost	=	#filter params	x	# filter positions	x	# of filters
432	=	3×3	x	4×4	x	3

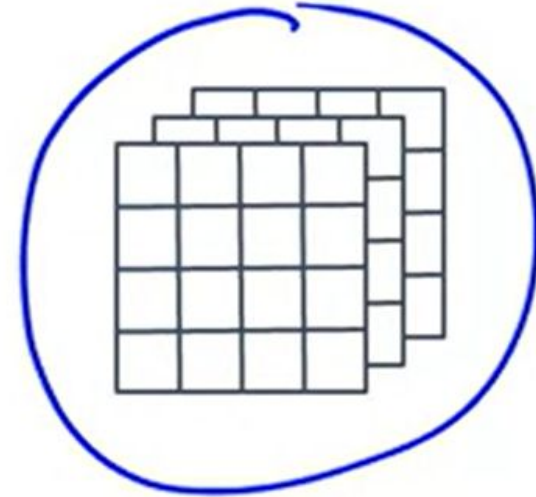
Depthwise Convolution



*

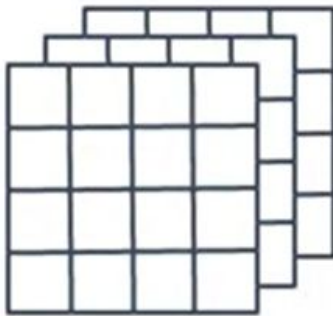


=



432

Pointwise Convolution

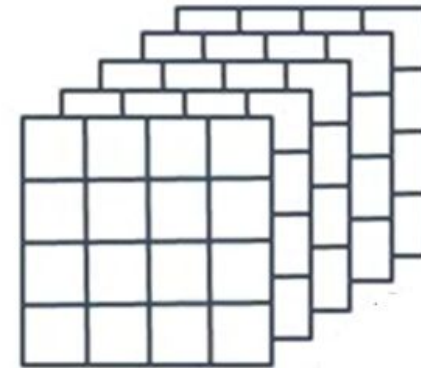


4x4x3

*

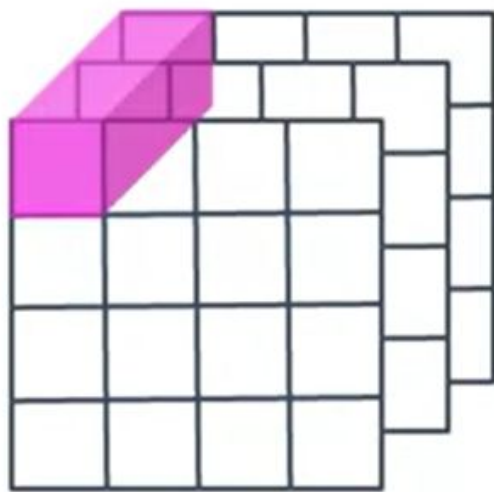


=



4x4x5

Pointwise Convolution



$4 \times 4 \times 3$

$n_{out} \times n_{out} \times n_c$

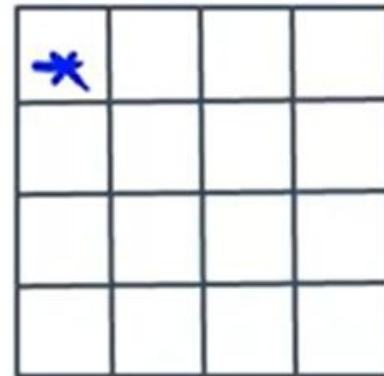
*



$1 \times 1 \times 3$

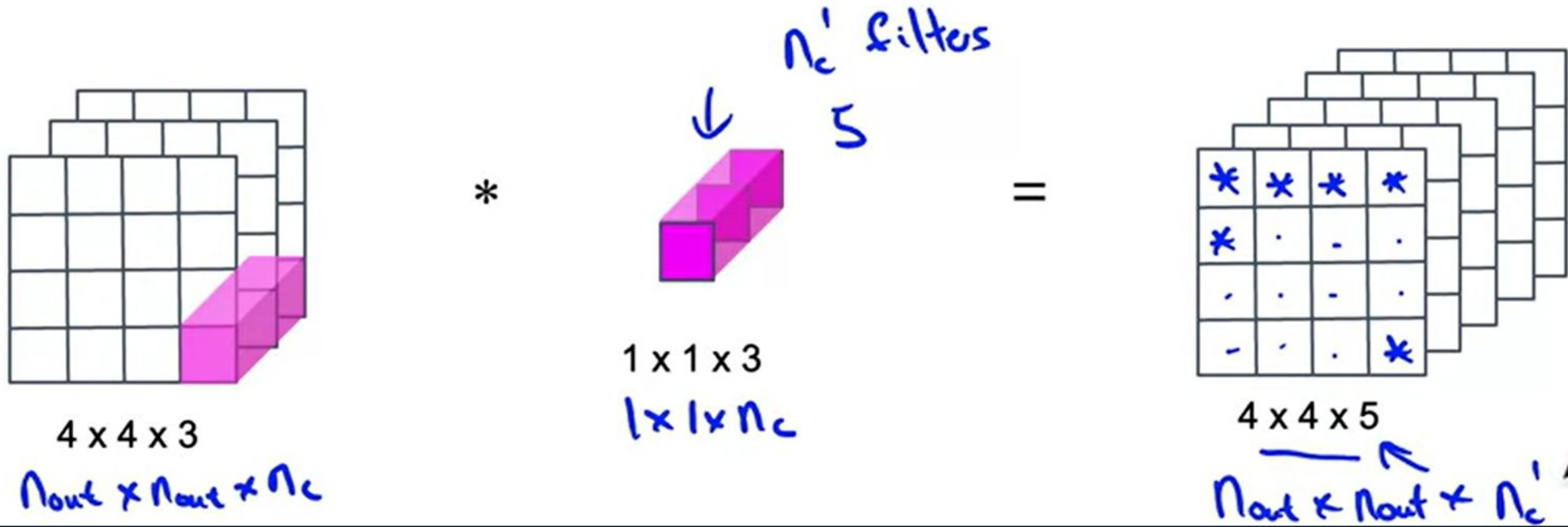
$1 \times 1 \times n_c$

=



4×4

Pointwise Convolution



Computational cost	=	#filter params	x	# filter positions	x	# of filters
240	=	$1 \times 1 \times 3$	x	4×4	x	5

Cost

Cost of normal convolution 2160

Cost of depthwise separable convolution

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

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

Cost

Cost of normal convolution 2160

Cost of depthwise separable convolution

depthwise + pointwise
432 + 240 = 672

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

In general, ratio:

$$\frac{1}{n_c^1} + \frac{1}{f^2}$$
$$\frac{1}{5} + \frac{1}{9}$$

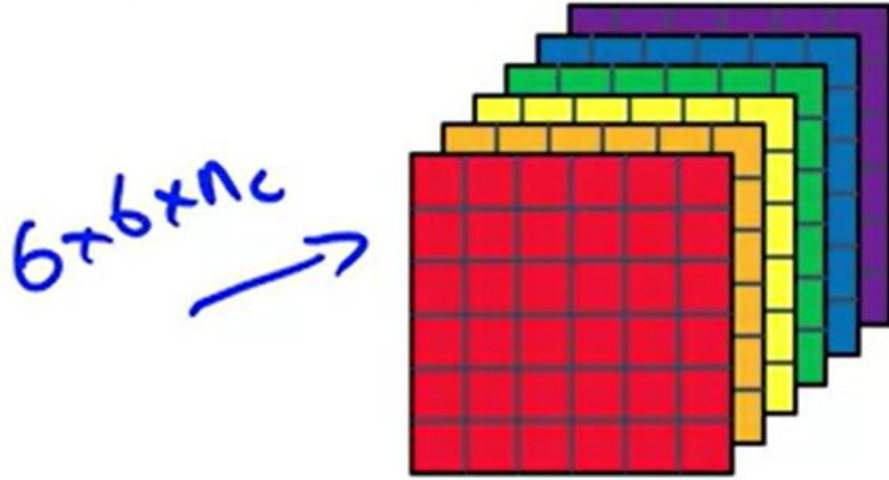
More generally:

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

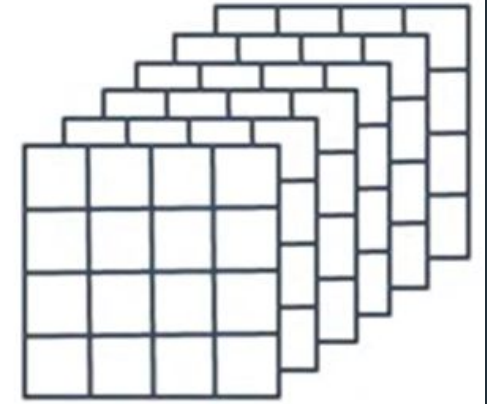
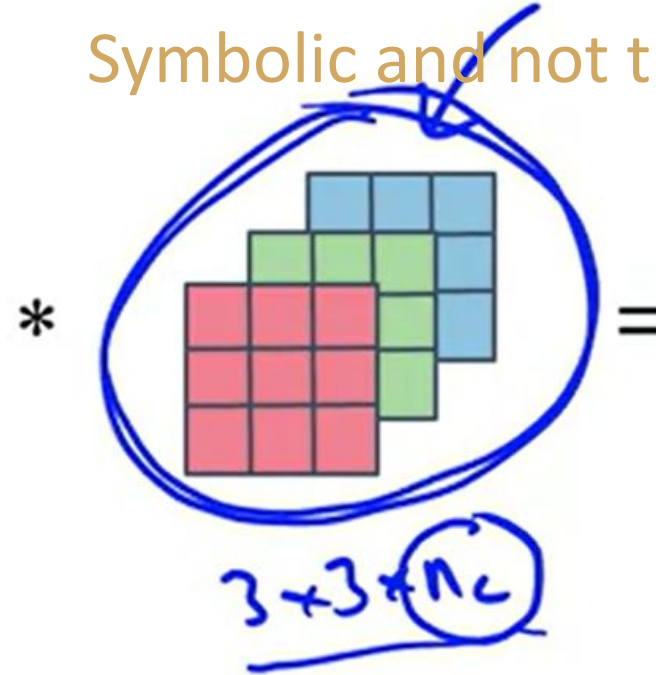
~10 times cheaper

Depthwise Separable Convolution

Depthwise Convolution

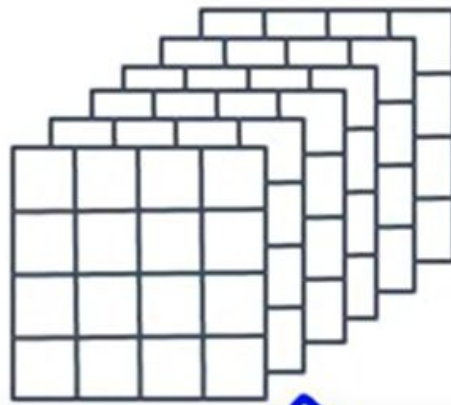


Symbolic and not the actual n_c

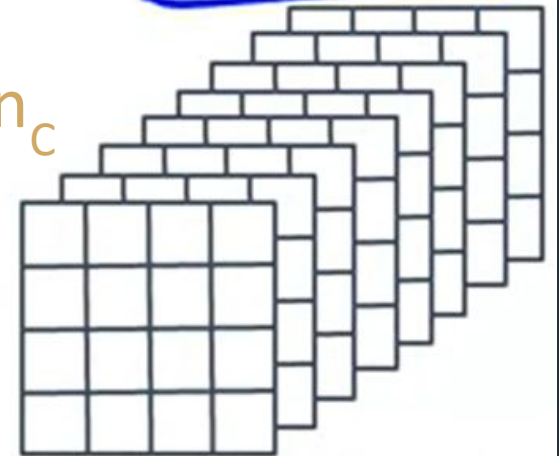
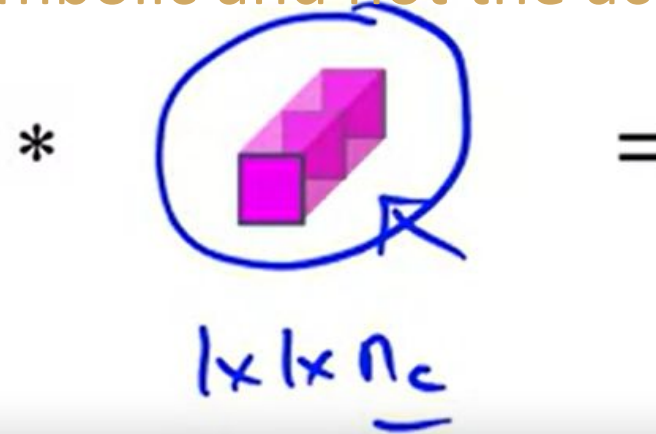


$4 \times 4 \times n_c$

Pointwise Convolution



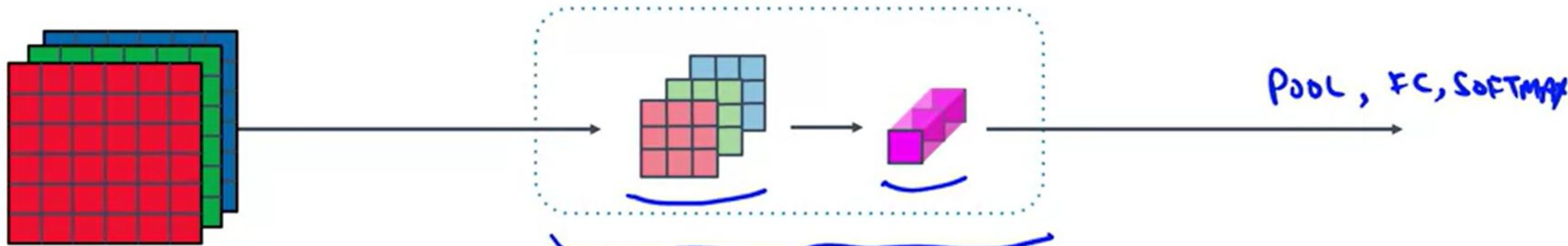
Symbolic and not the actual n_c



$4 \times 4 \times 8$
 n_c

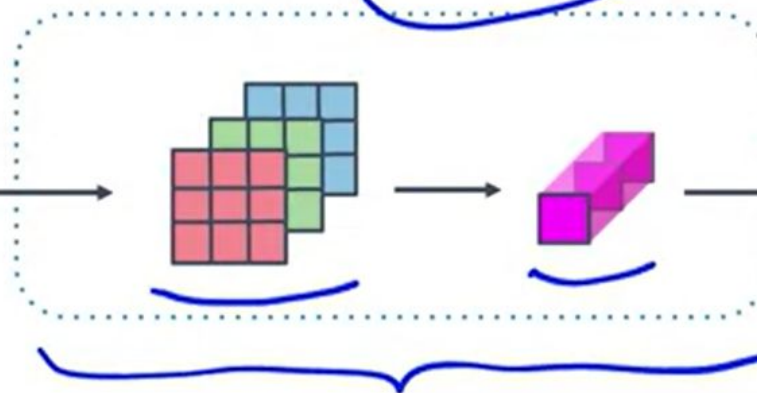
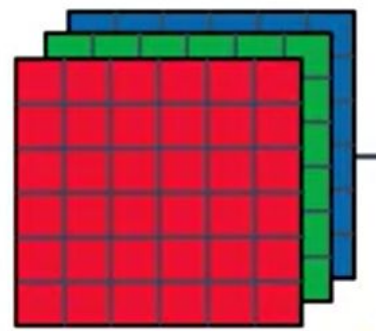
MobileNet

MobileNet v1



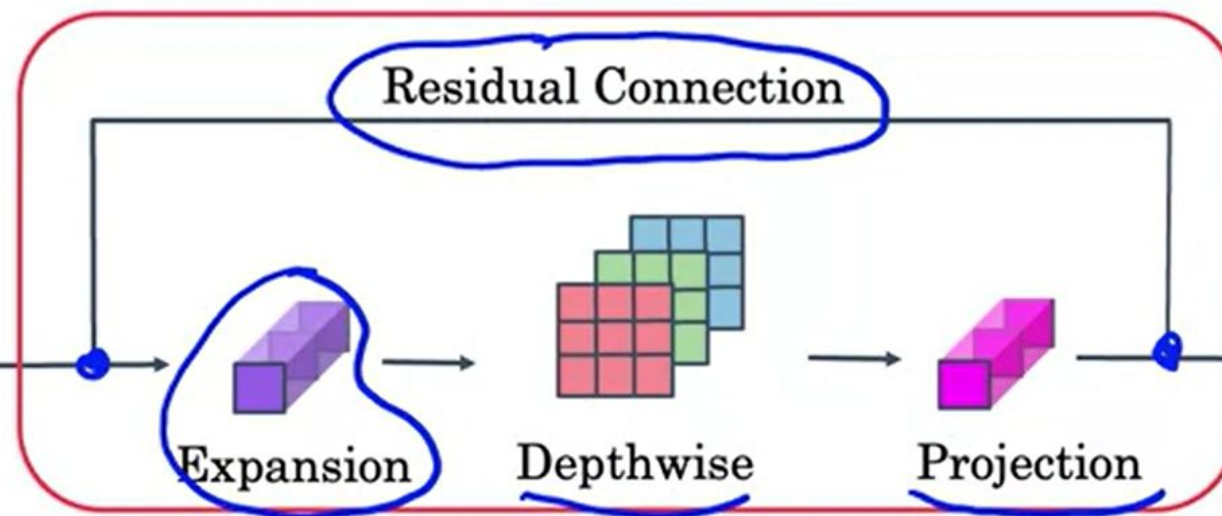
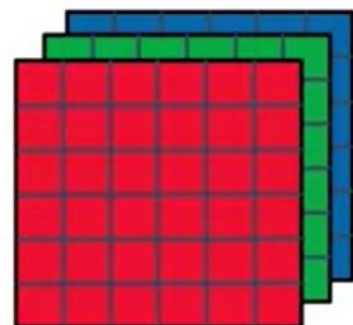
MobileNet

MobileNet v1



Pool, FC, Softmax

MobileNet v2



Pool, FC, Softmax

[Sandler et al. 2019, MobileNetV2: Inverted Residuals and Linear Bottlenecks]

Bottleneck