



THE UNIVERSITY OF  
**MELBOURNE**

# **CNN Architectures**

**COMP90051 Statistical Machine Learning**

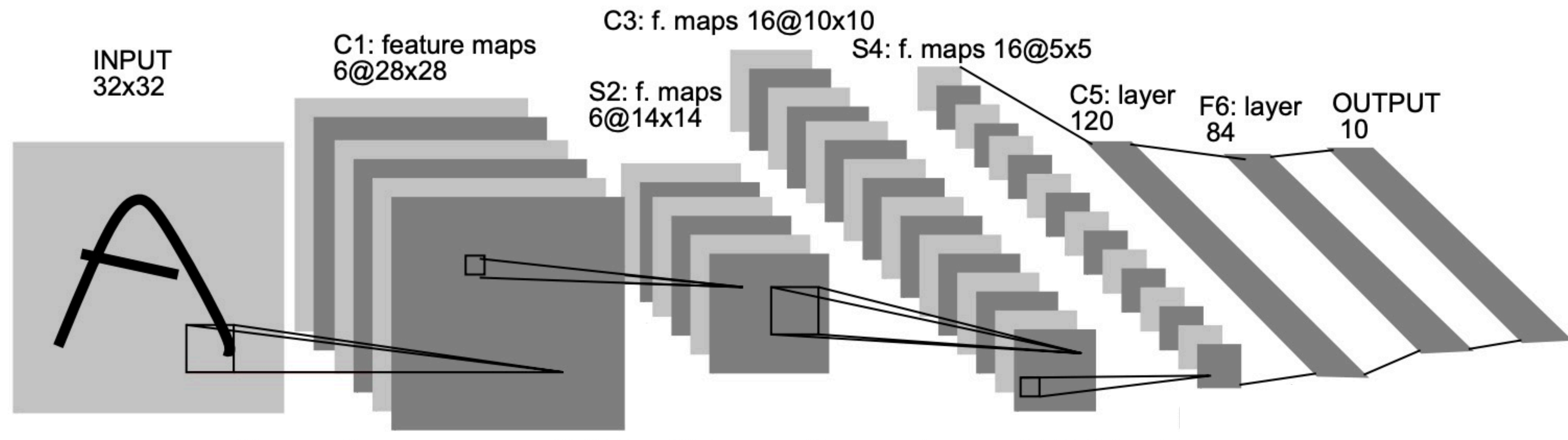
**QiuHong Ke**

Copyright: University of Melbourne

# Outline

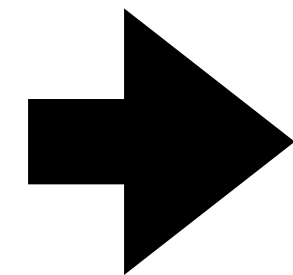
- LeNet5
- AlexNet
- VGG
- GoogleNet
- ResNet

# Architecture

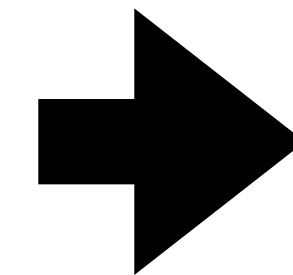


**1st layer**

32x32x1



**Convolutional layer**  
No.filters: 6  
Filter size: 5x5  
Padding: 0  
Stride:1



?

28x28x6

No padding:  $\text{output\_size} = \text{ceiling}((\text{input\_size} - \text{kernel\_size} + 1) / \text{stride})$

LeNet5

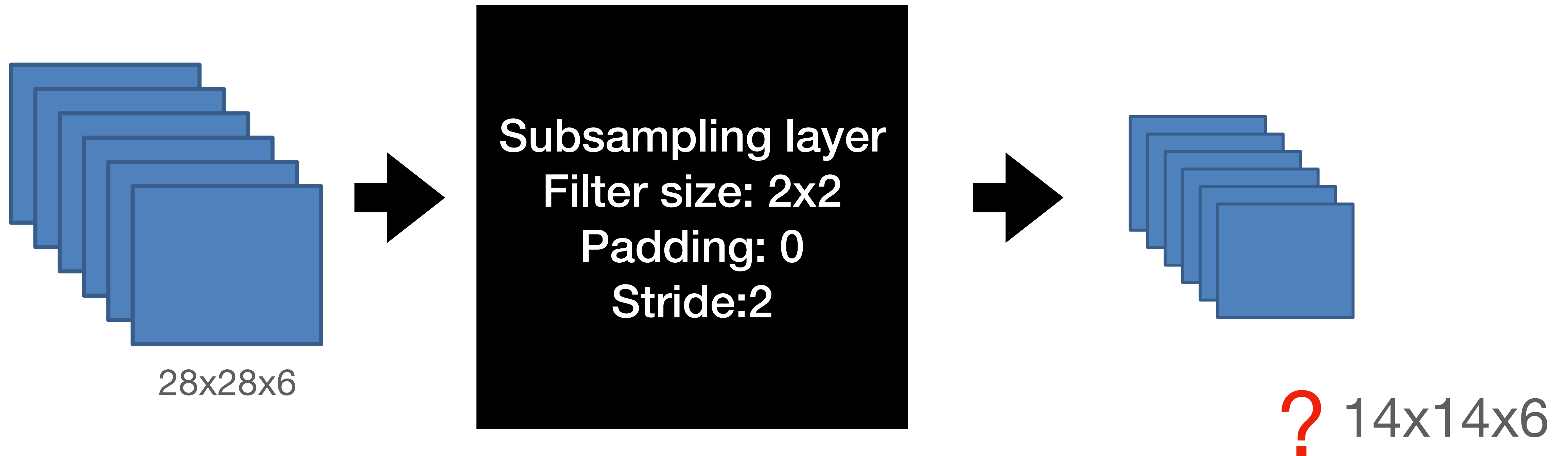
AlexNet

VGG

GoogleNet

ResNet

## 2nd layer



- Take the sum of all units in the  $2 \times 2$  window
- output of each patch =  $\text{sum} \times \text{coefficient (trainable)} + \text{bias (trainable)}$

No padding:  $\text{output\_size} = \text{ceiling} \left( \frac{\text{input\_size} - \text{kernel\_size} + 1}{\text{stride}} \right)$

LeNet5

AlexNet

VGG

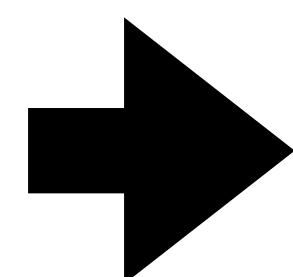
GoogleNet

ResNet

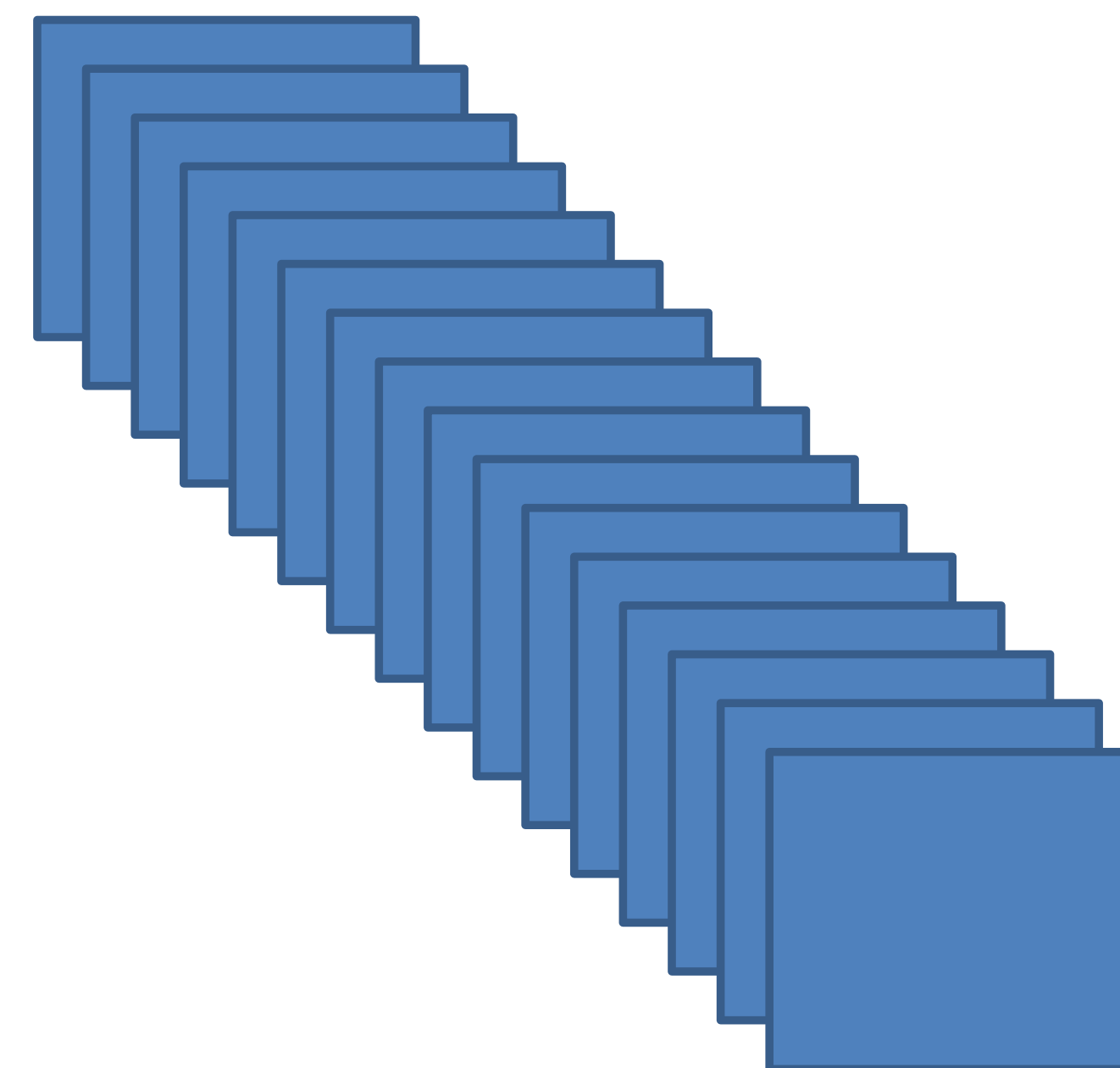
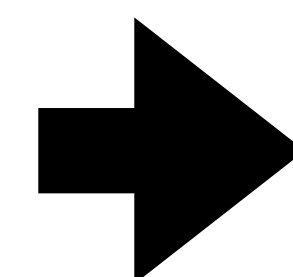
3rd layer



14x14x6



Convolutional layer  
No.filters: 16  
Filter size: 5x5  
Padding: 0  
Stride:1

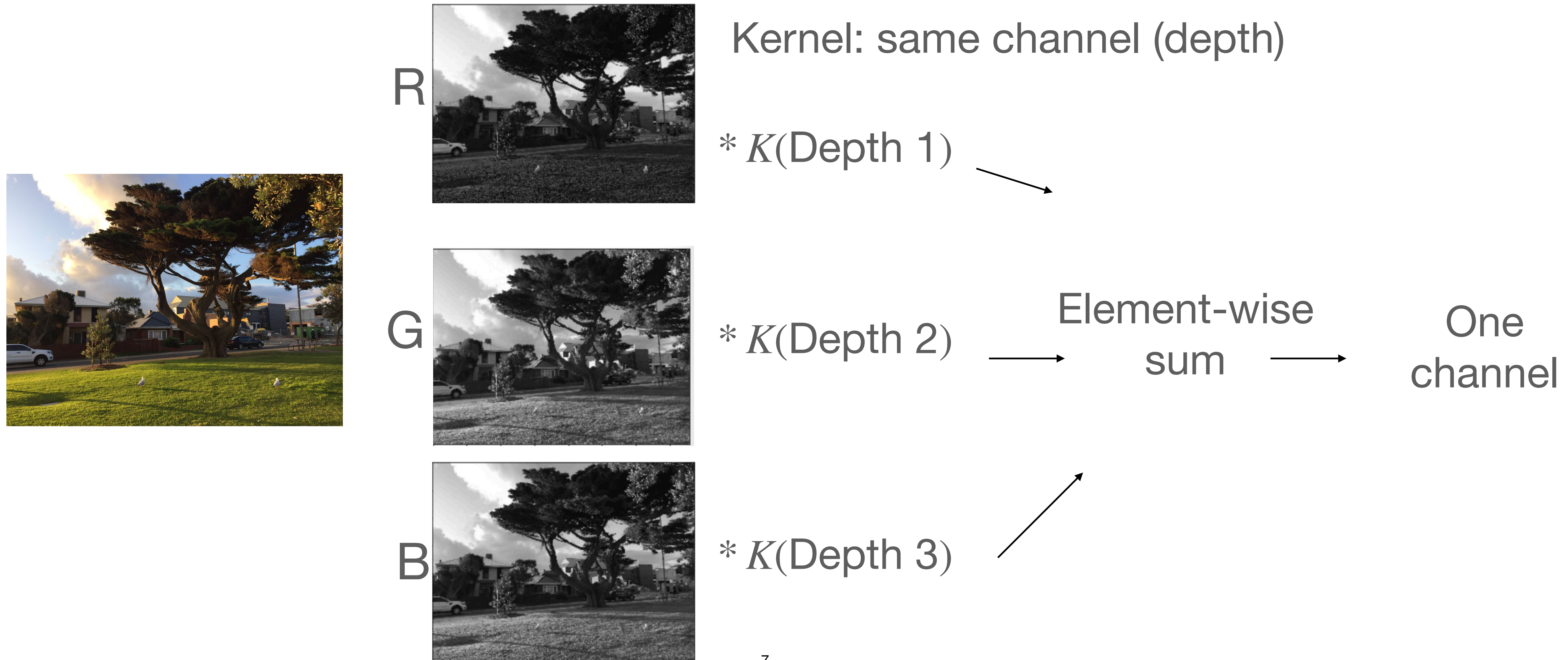


10x10x16 ?

No padding:  $\text{output\_size} = \text{ceiling} \left( \frac{\text{input\_size} - \text{kernel\_size} + 1}{\text{stride}} \right)$

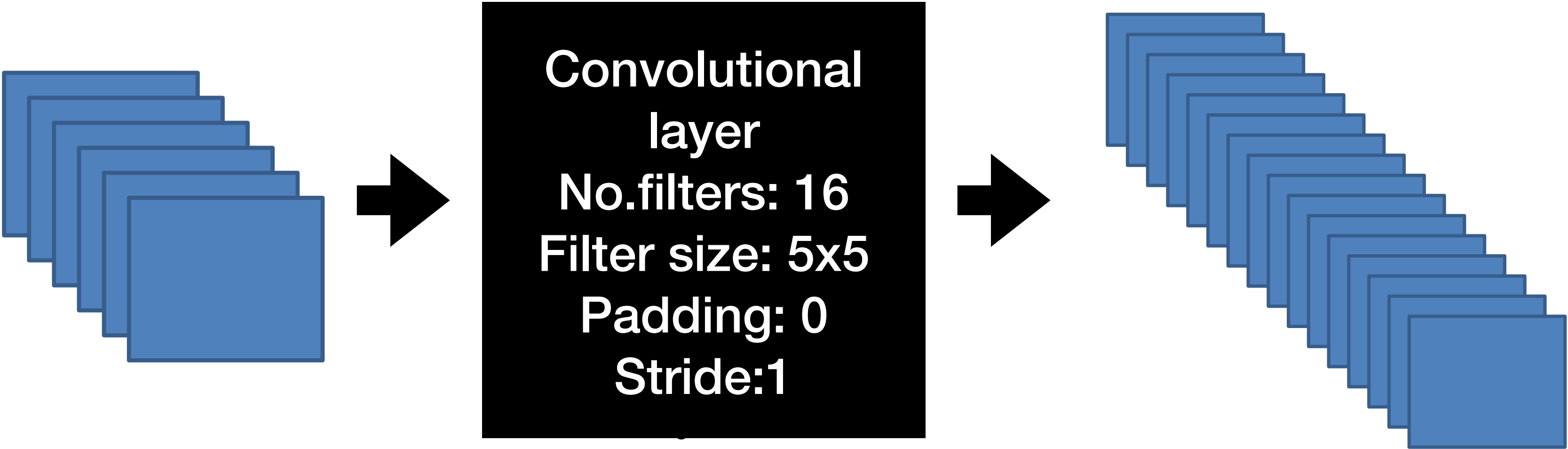


## Convolution on Multiple-channel input



3rd layer: Non-complete connection scheme

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	X				X	X	X			X	X	X	X		X	X
1	X	X				X	X	X			X	X	X	X		X
2	X	X	X				X	X	X			X		X	X	X
3		X	X	X			X	X	X	X			X		X	X
4			X	X	X			X	X	X	X		X	X		X
5				X	X	X			X	X	X	X		X	X	X





LeNet5

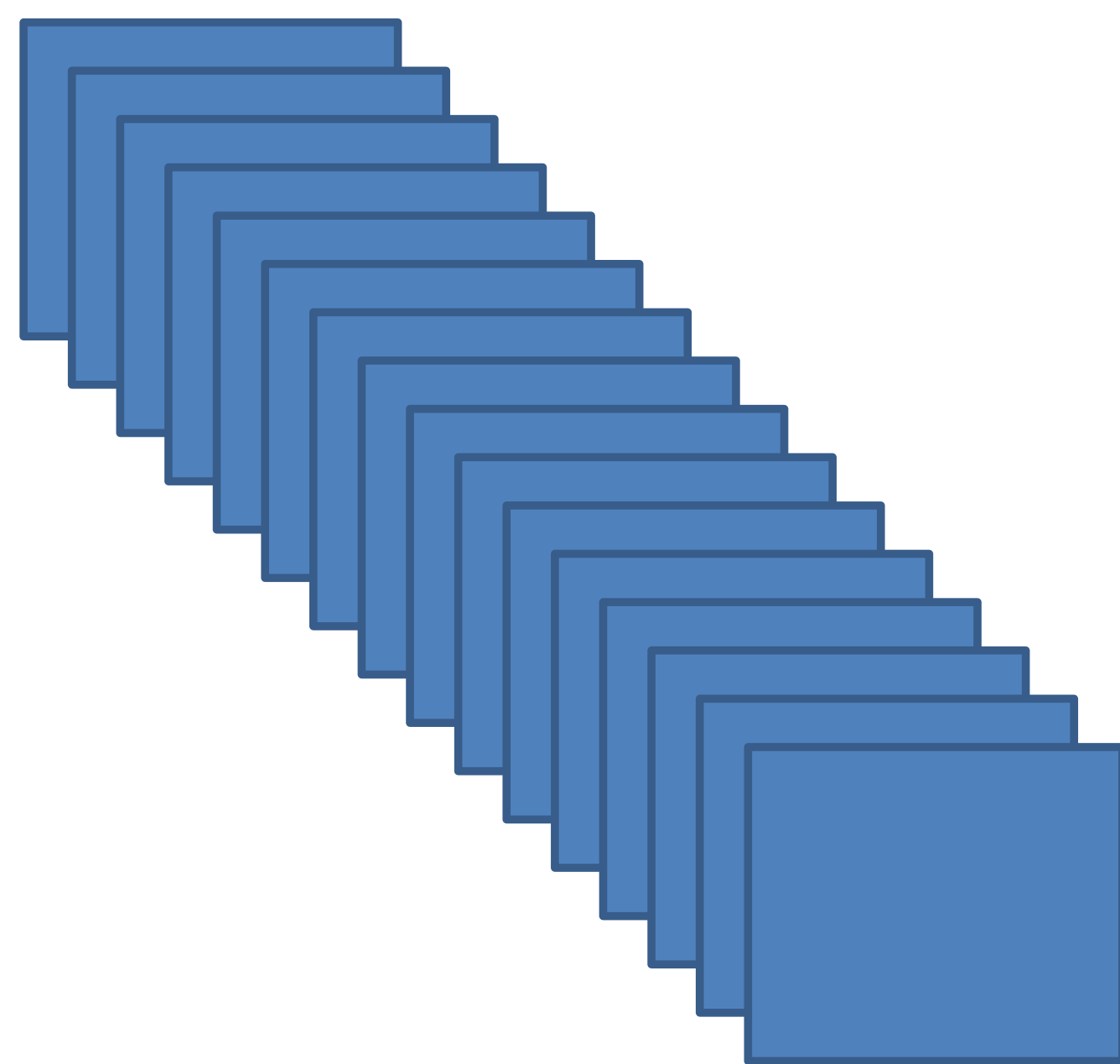
AlexNet

VGG

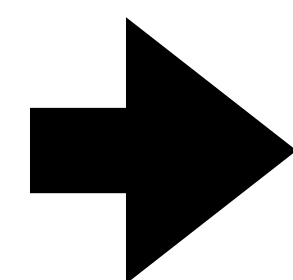
GoogleNet

ResNet

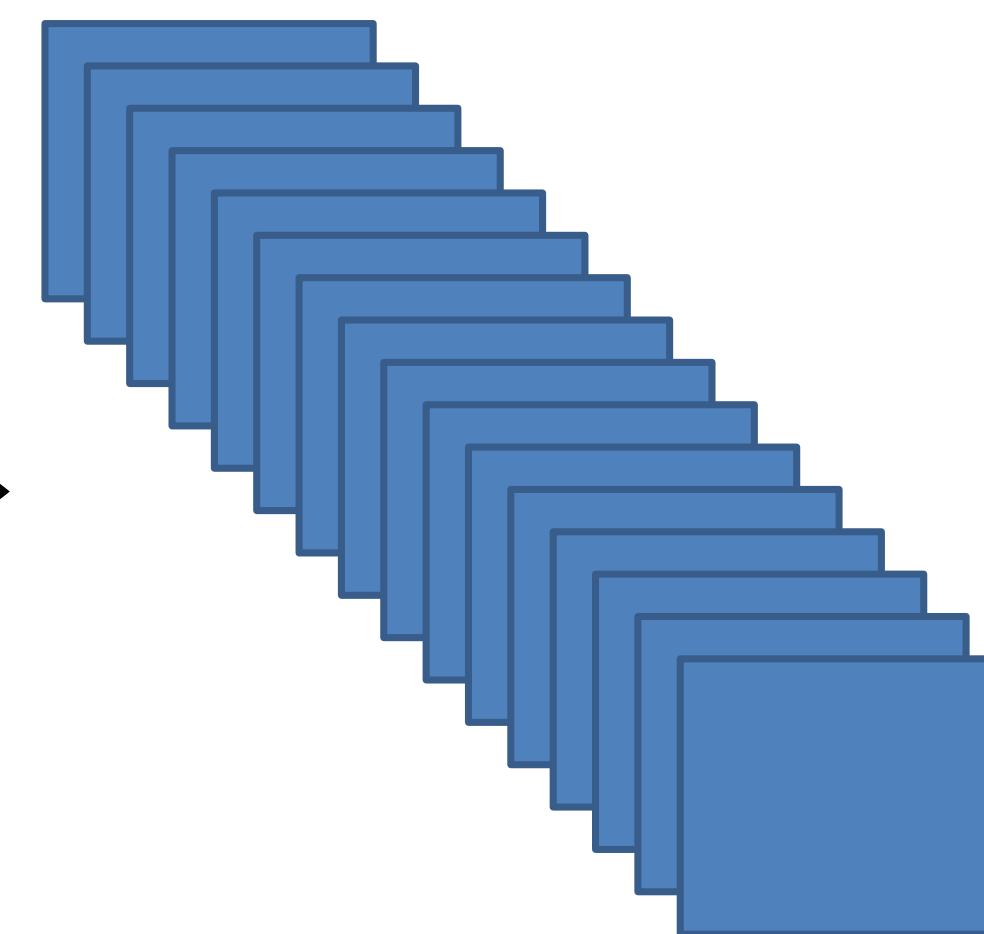
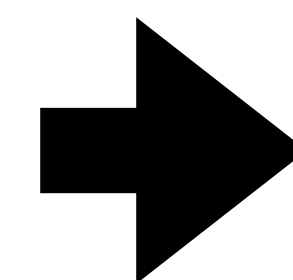
4th layer



10x10x16



Subsampling layer  
Filter size: 2x2  
Padding: 0  
Stride:2



5x5x16 ?

No padding:  $\text{output\_size} = \text{ceiling} \left( \frac{\text{input\_size} - \text{kernel\_size} + 1}{\text{stride}} \right)$

LeNet5

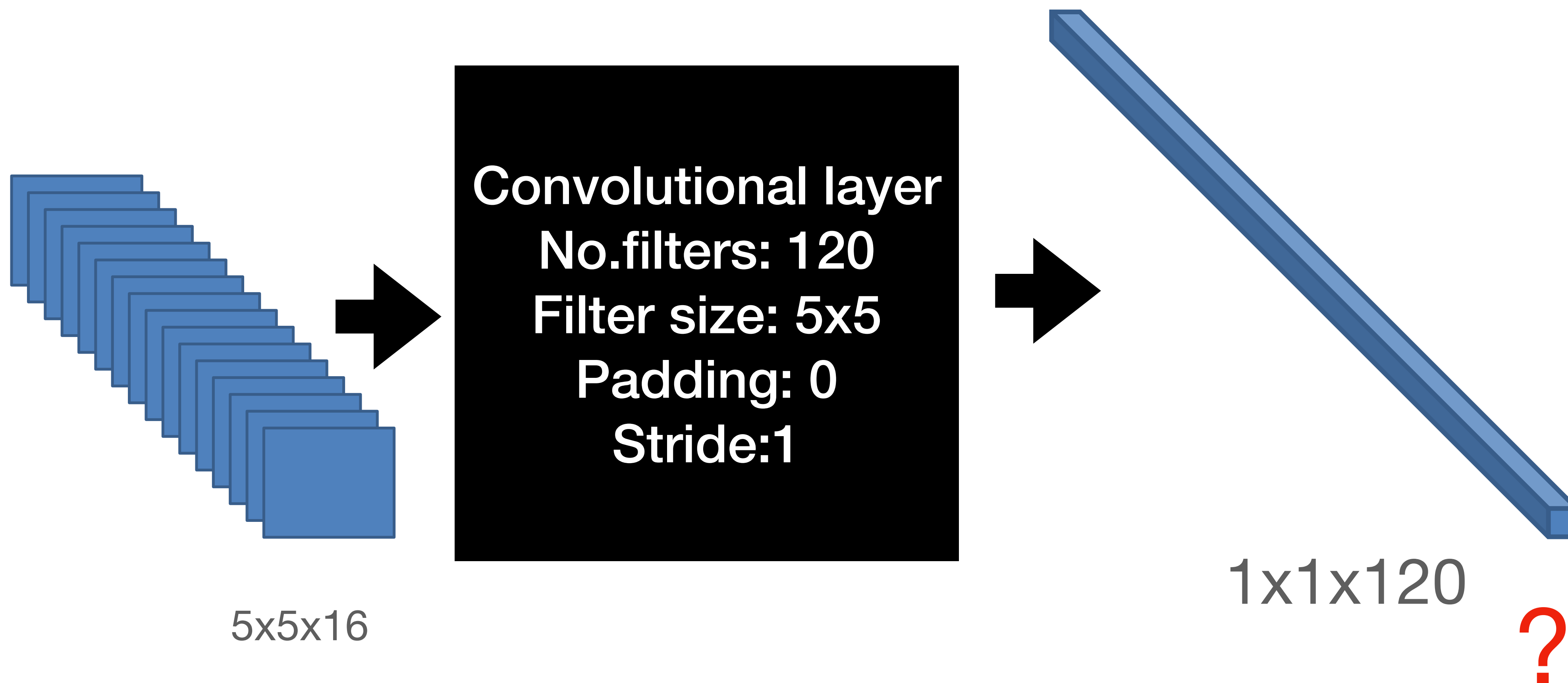
AlexNet

VGG

GoogleNet

ResNet

5th layer



No padding:  $\text{output\_size} = \text{ceiling} \left( \frac{\text{input\_size} - \text{kernel\_size} + 1}{\text{stride}} \right)$

LeNet5

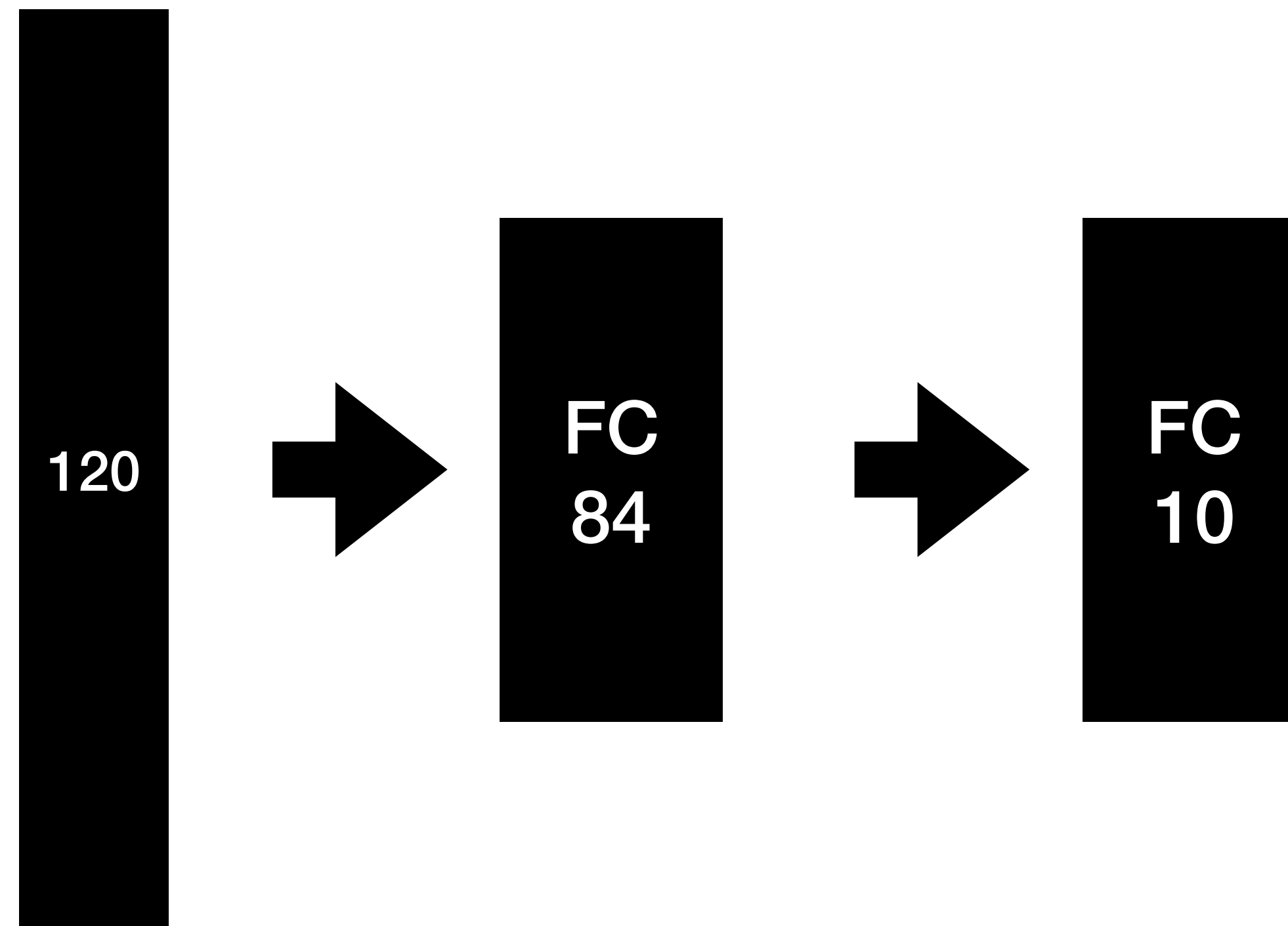
AlexNet

VGG

GoogleNet

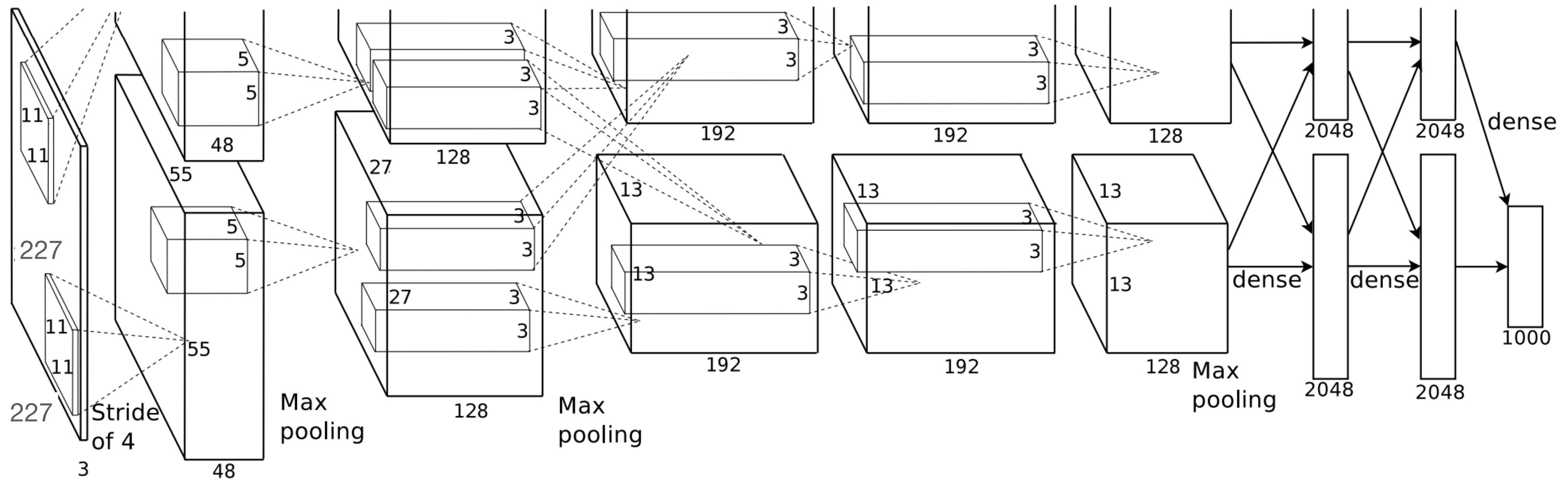
ResNet

**Following: Fully connected layers**



# Architecture

This training on 2 GPU where each stream represent a training stream on one GPU



LeNet5

AlexNet

VGG

GoogleNet

ResNet

1st layer

How many parameters?

$$11 \times 11 \times 3 \times 96 = 34,848$$

Filter size

Input channels

No. filters

Convolutional layer

No.filters: 96

Filter size: 11x11

Padding: 0

Stride: 4



227x227x3

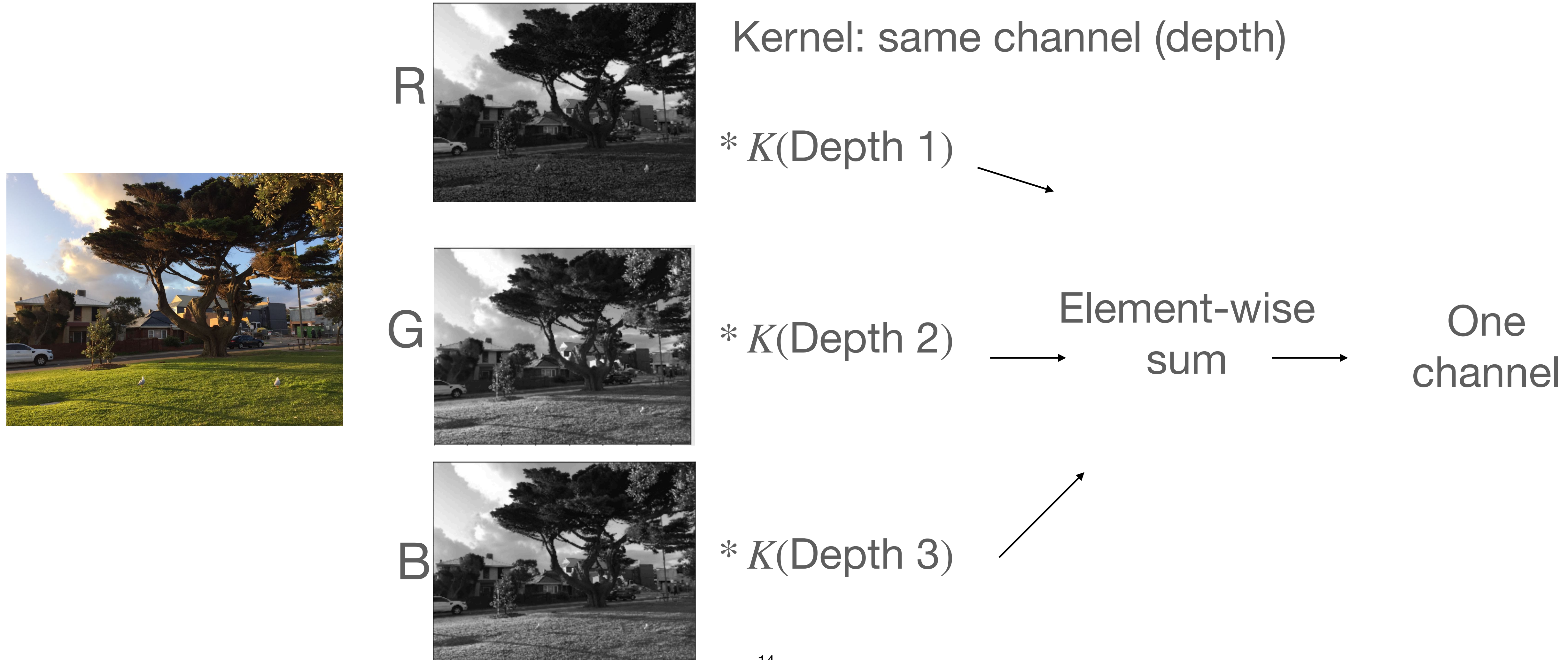


55x55x48

No padding:  $\text{output\_size} = \text{ceiling}((\text{input\_size} - \text{kernel\_size} + 1) / \text{stride})$



## Convolution on Multiple-channel input



LeNet5

AlexNet

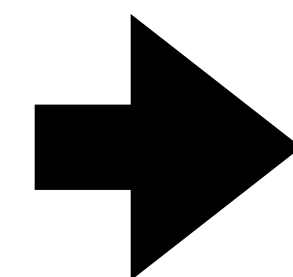
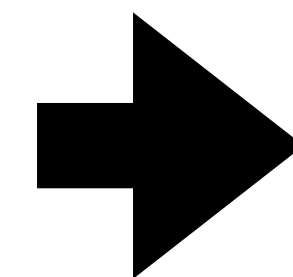
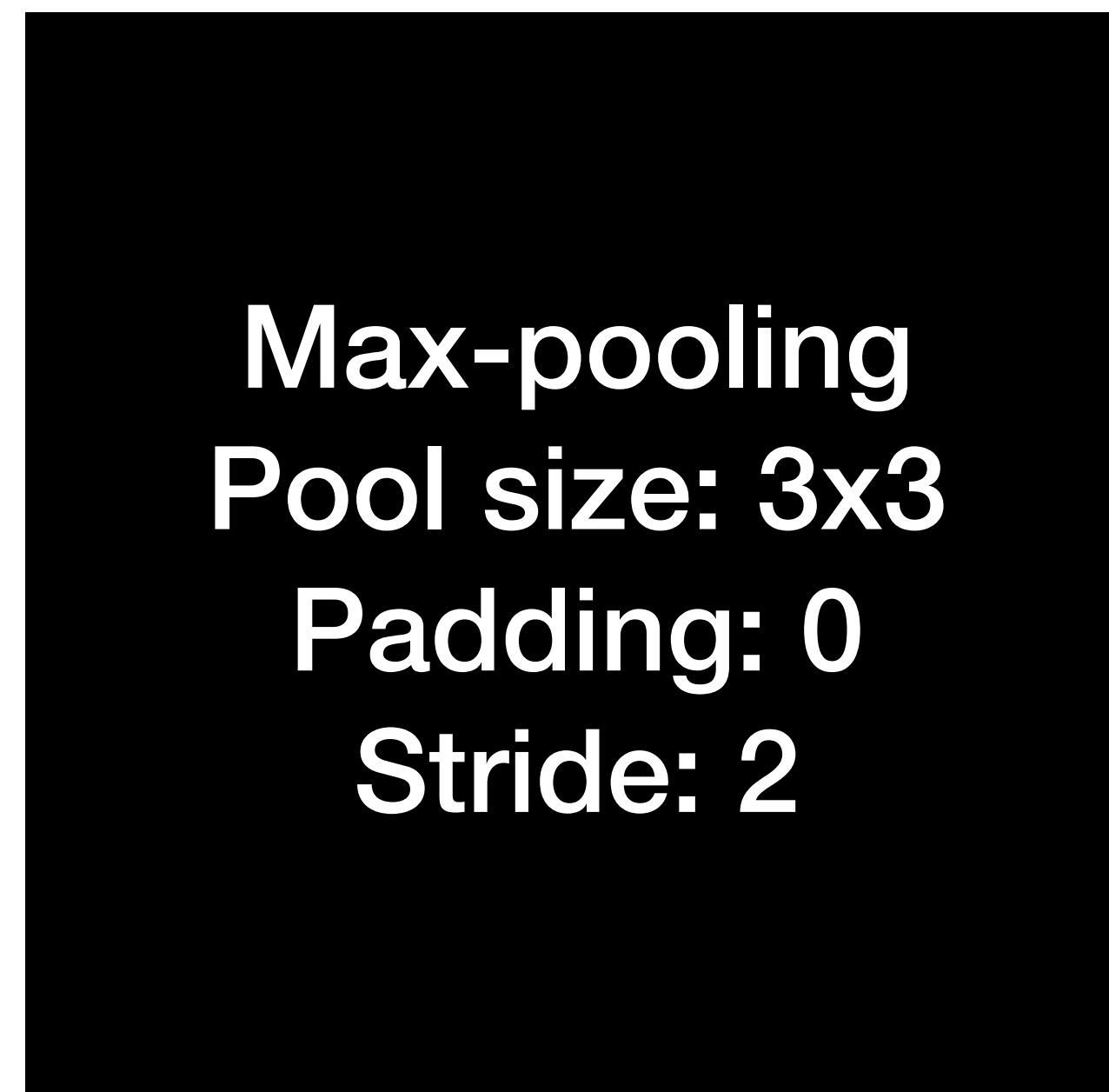
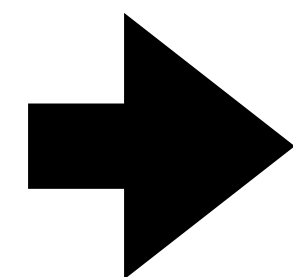
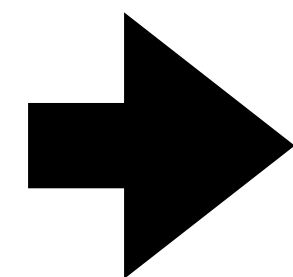
VGG

GoogleNet

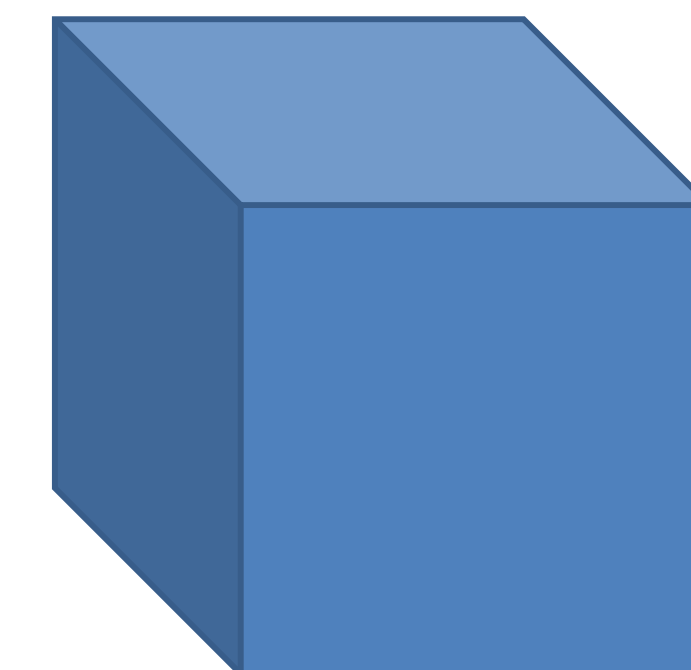
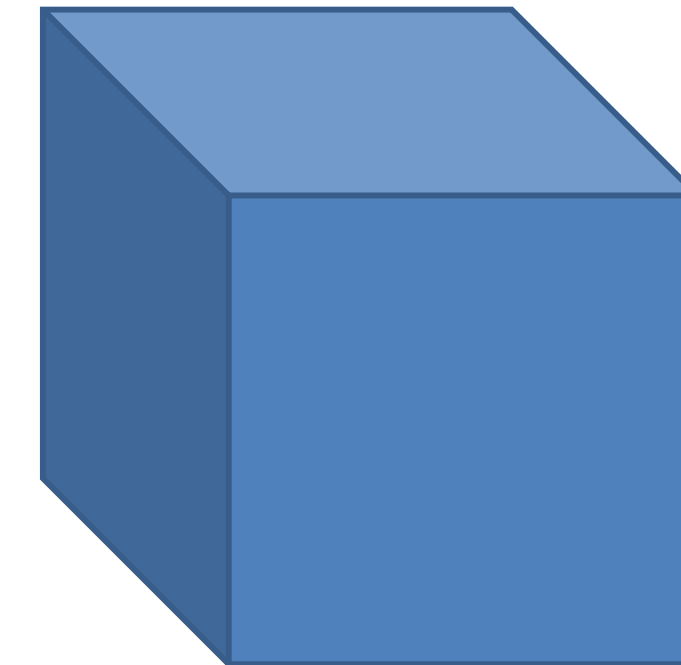
ResNet

2nd layer

How many parameters?



27x27x48



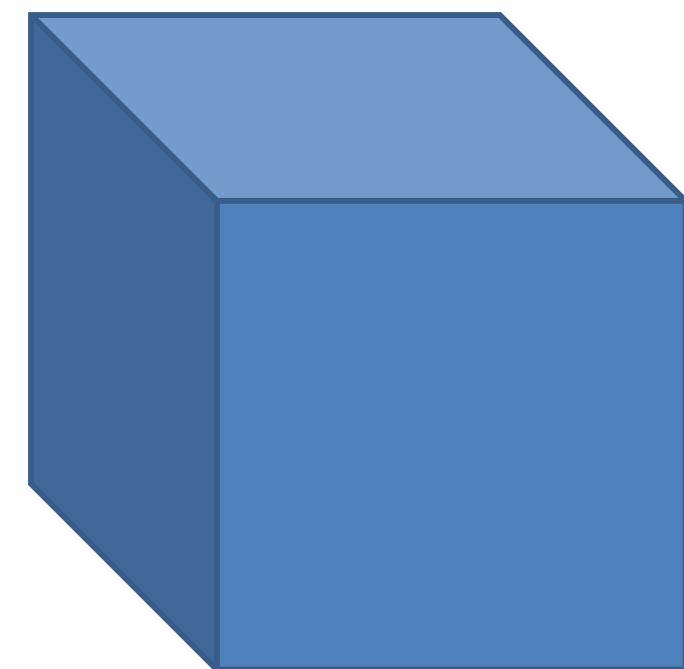
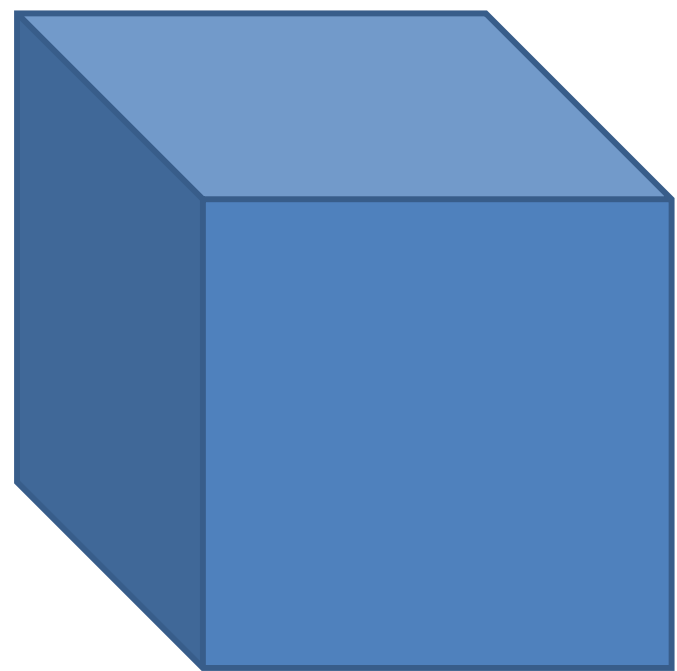
27x27x48 ?

55x55x48

No padding:  $\text{output\_size} = \text{ceiling} \left( \frac{\text{input\_size} - \text{kernel\_size} + 1}{\text{stride}} \right)$

**3rd layer**

27x27x48

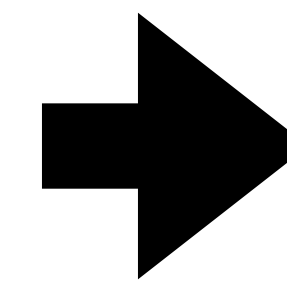


27x27x48

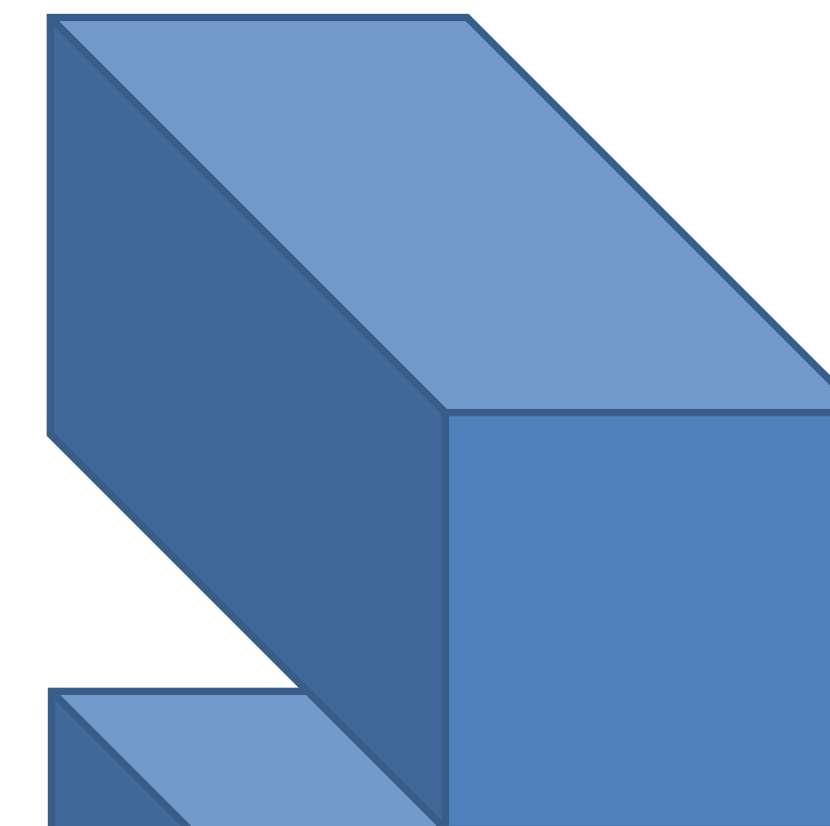
**How many parameters?**

$$5 \times 5 \times 48 \times 256 = 614,400$$

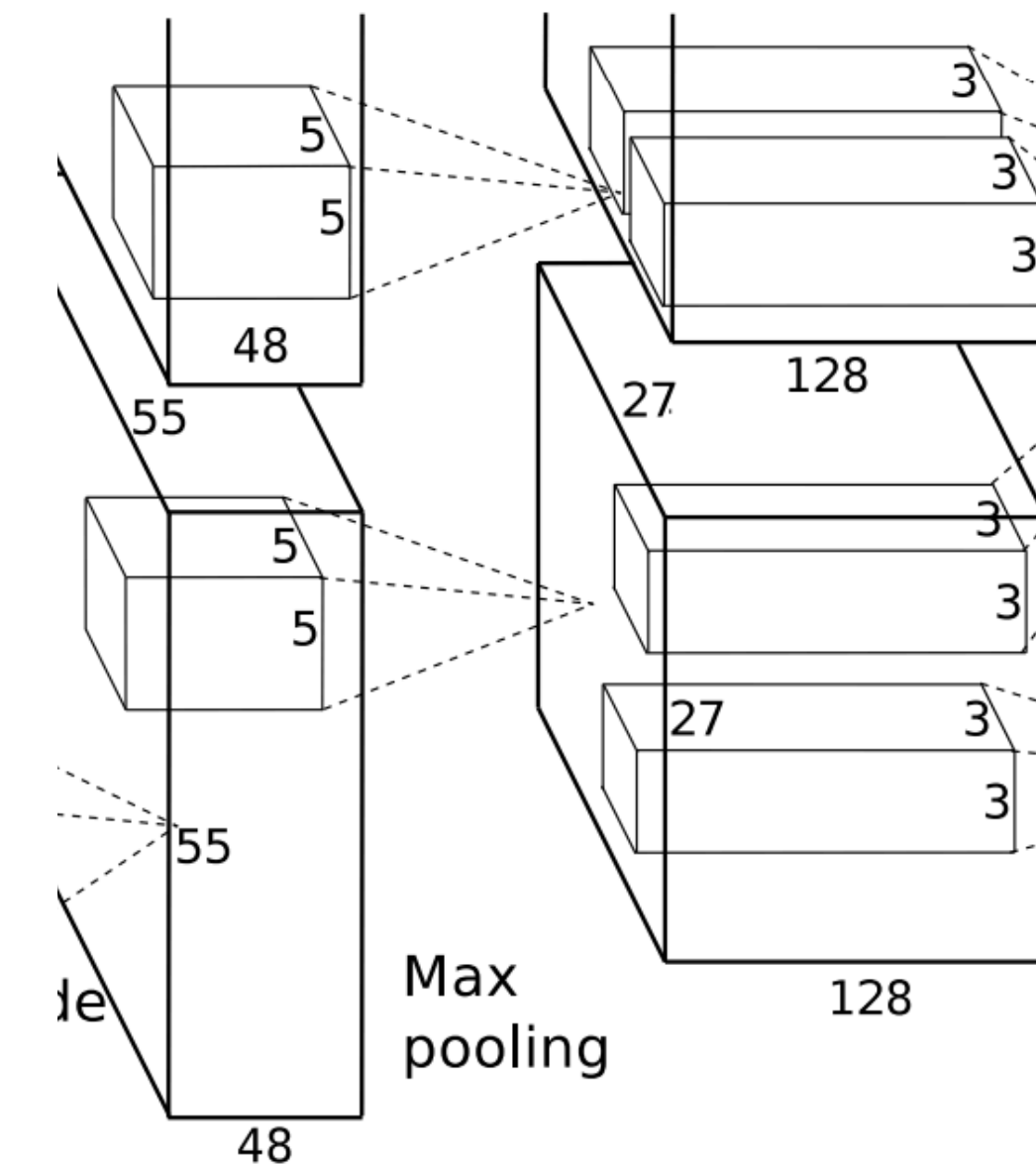
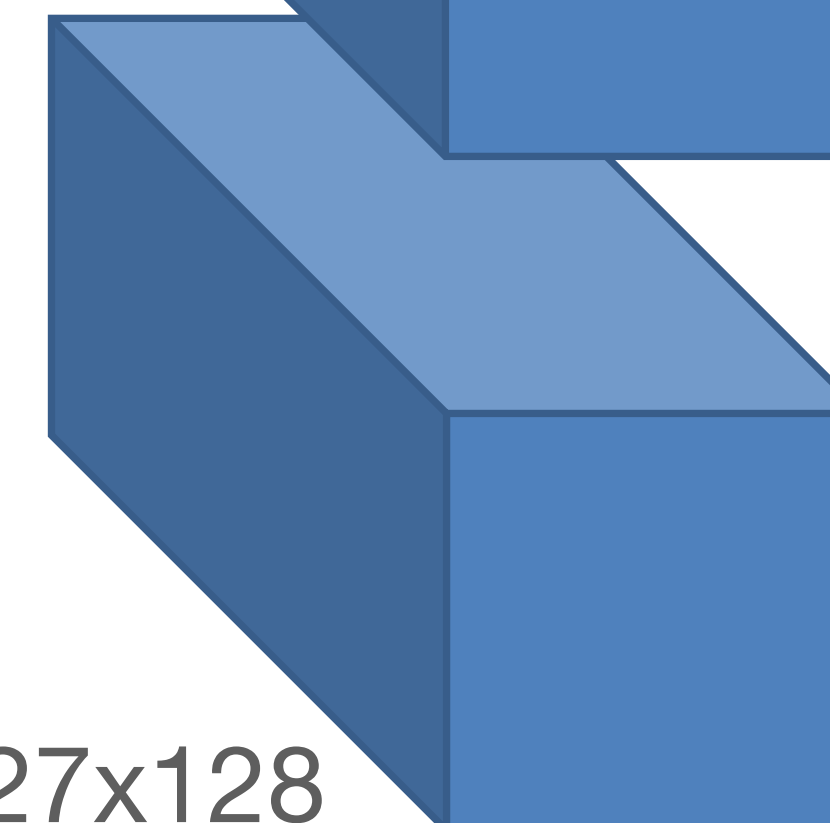
**Convolutional layer**  
 No.filters: 256  
 Filter size: 5x5  
 Padding: 2  
 Stride: 1



27x27x128



27x27x128

padding:  $\text{output\_size} = \text{ceiling}(\text{input\_size} / \text{stride})$ 

?



LeNet5

AlexNet

VGG

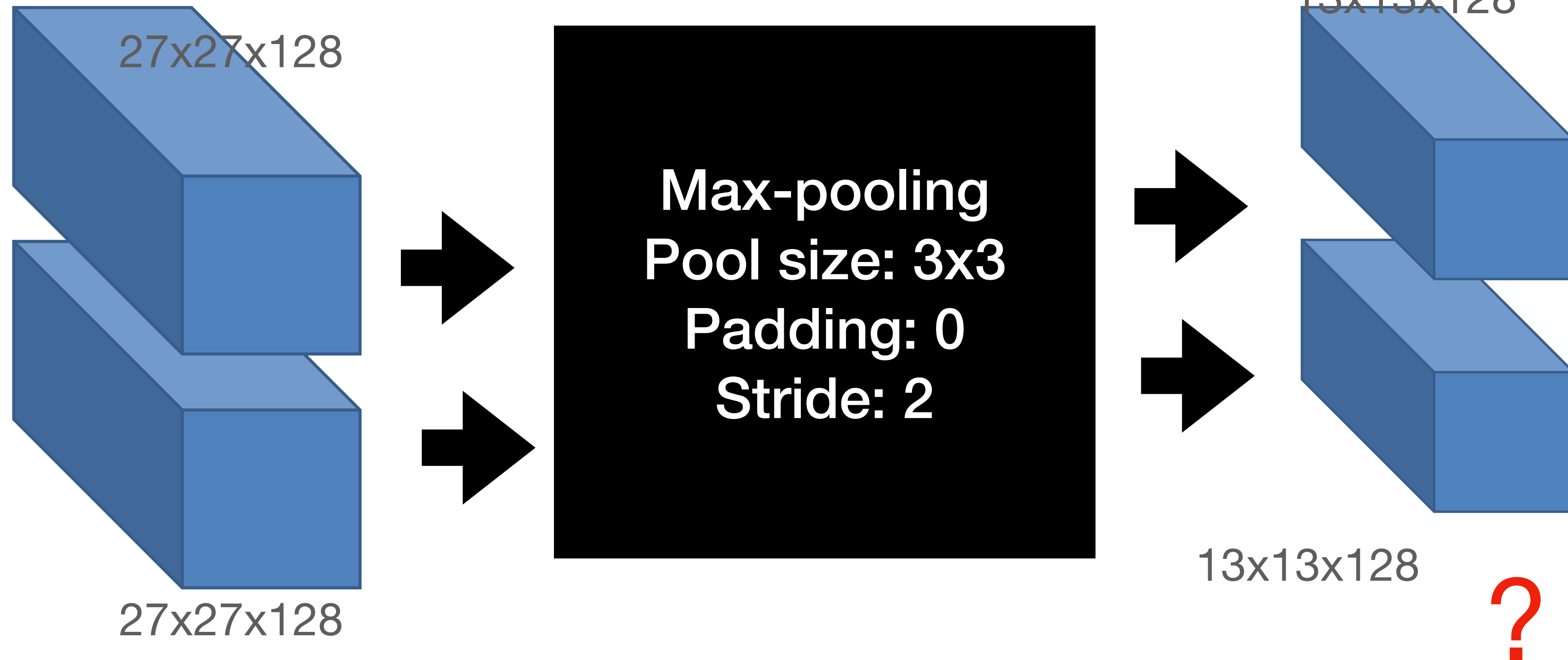
GoogleNet

ResNet

4th layer

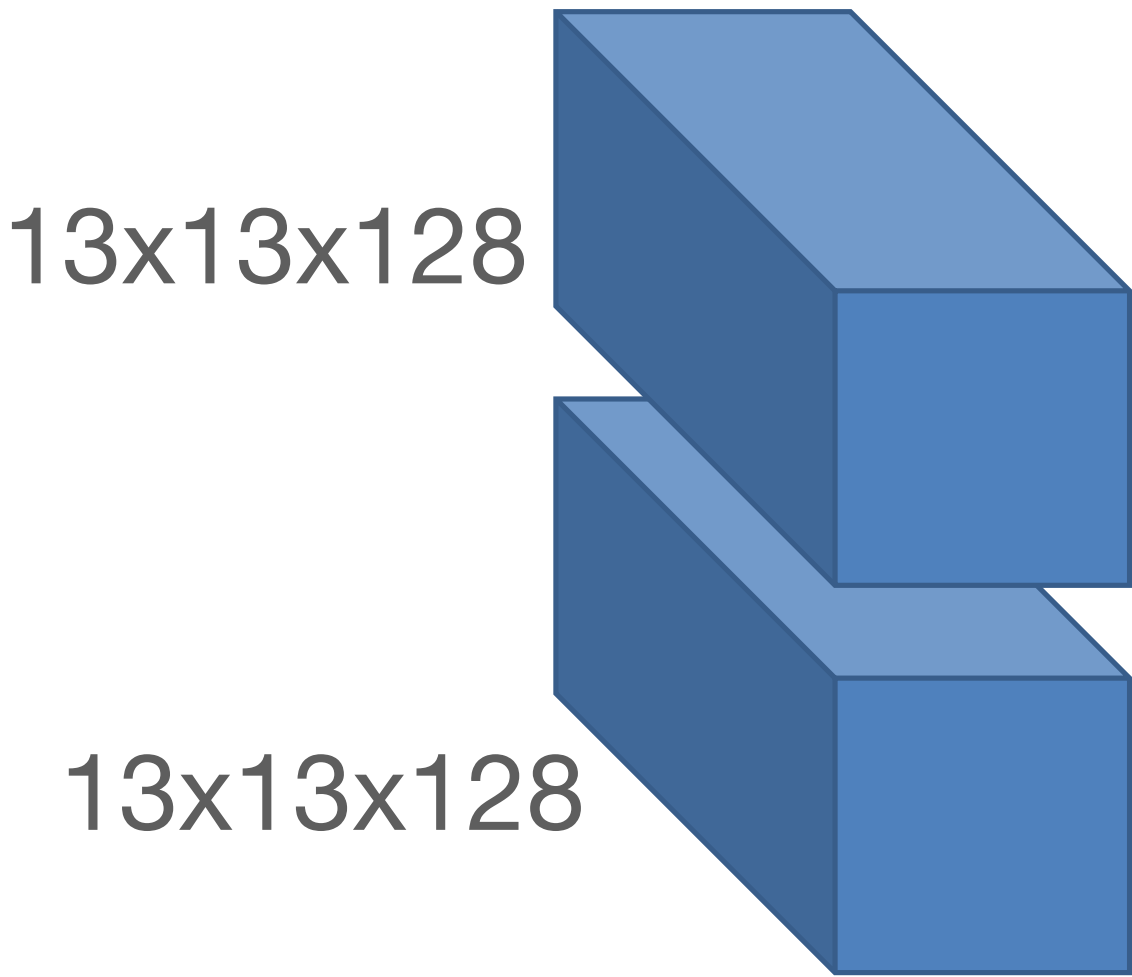
How many parameters?

0



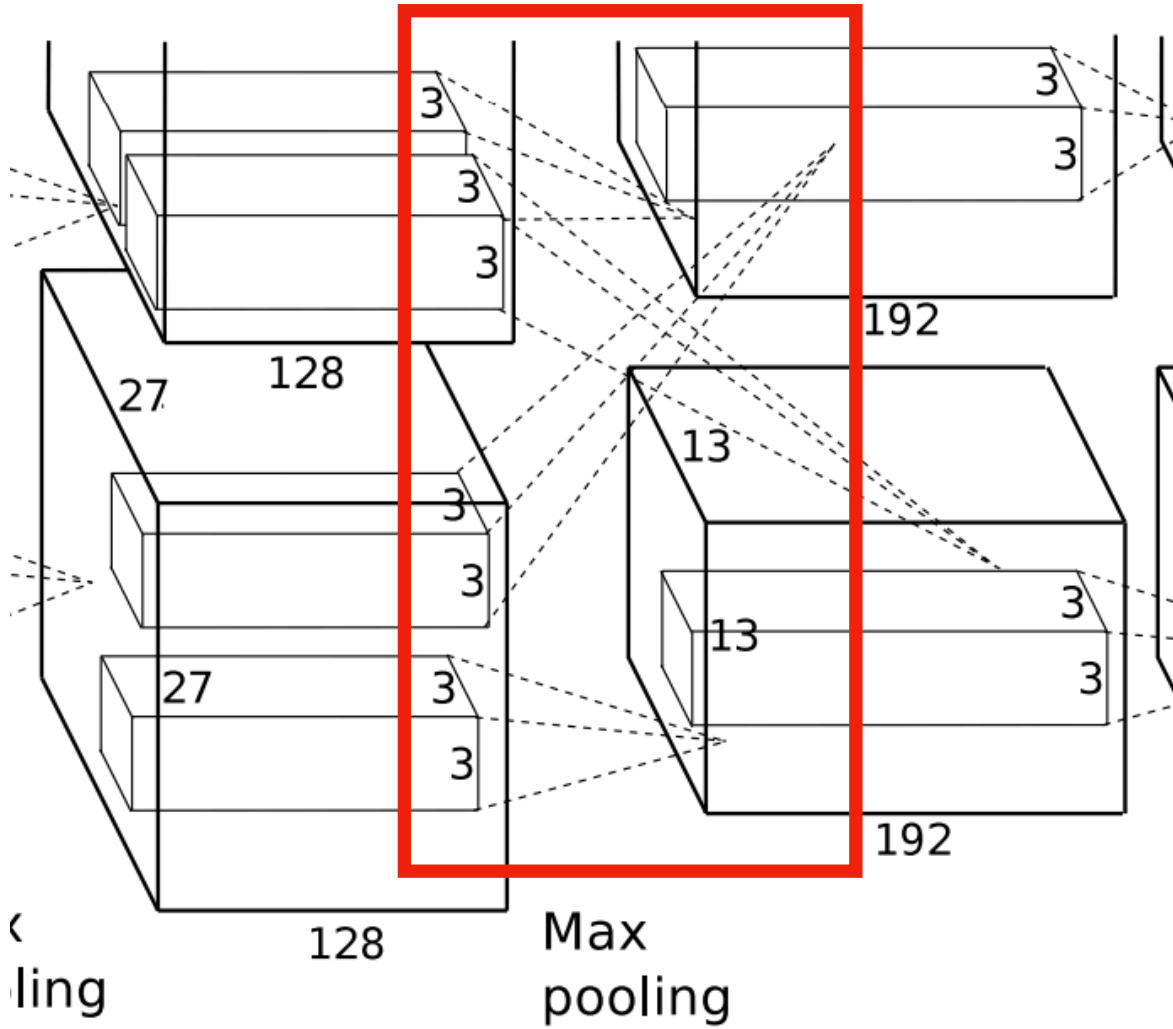
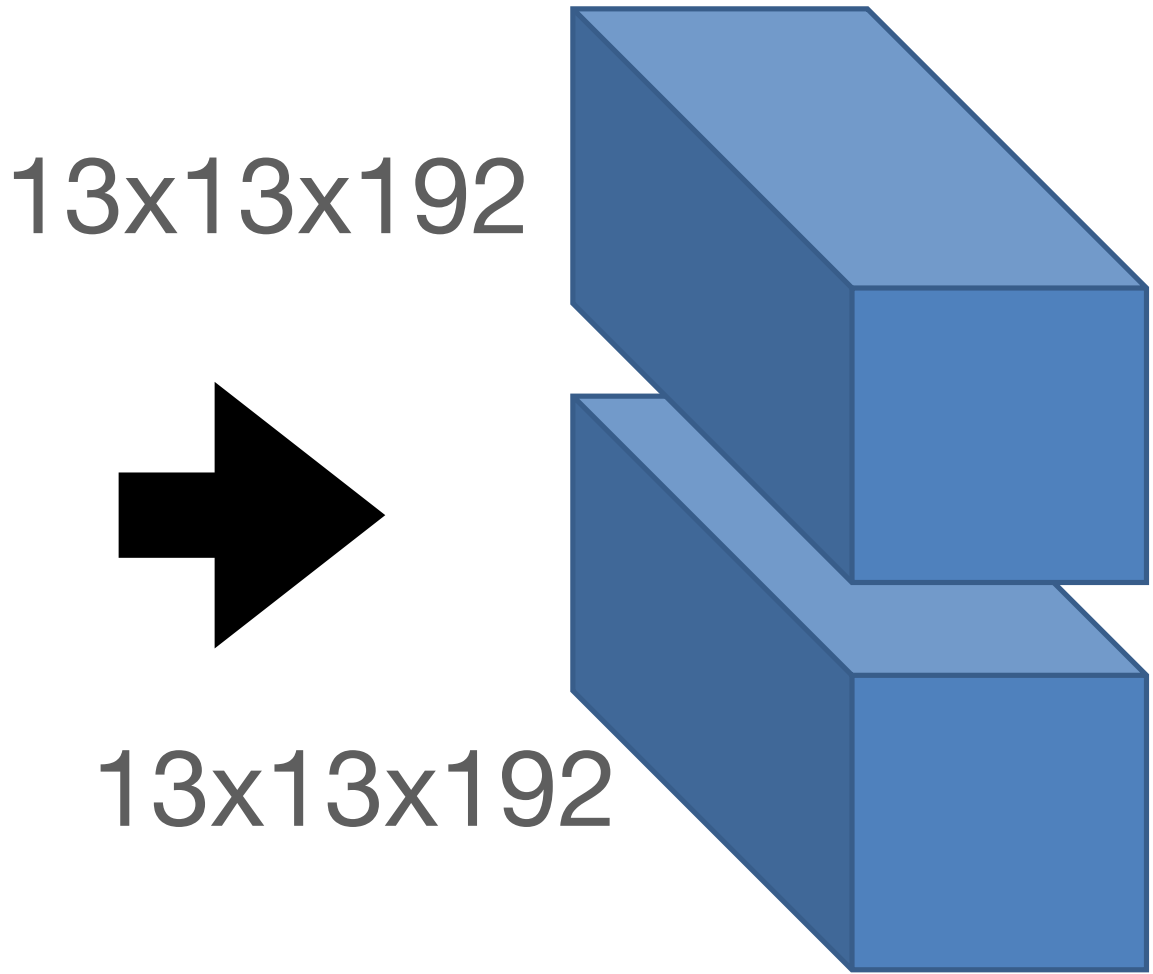
$$\text{output\_size} = \text{ceiling} \left( \frac{\text{input\_size} - \text{kernel\_size} + 1}{\text{stride}} \right)$$

5th layer



How many parameters?  
 $3 \times 3 \times 256 \times 384 = 884,736$

**Convolution layer**  
No.filters: 384  
Filter size:  $3 \times 3$   
Padding: 1  
Stride: 1



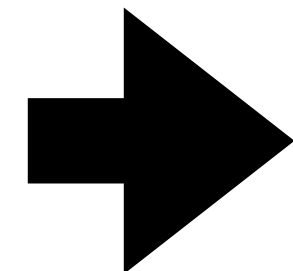


## Following convolutional layers

How many parameters?

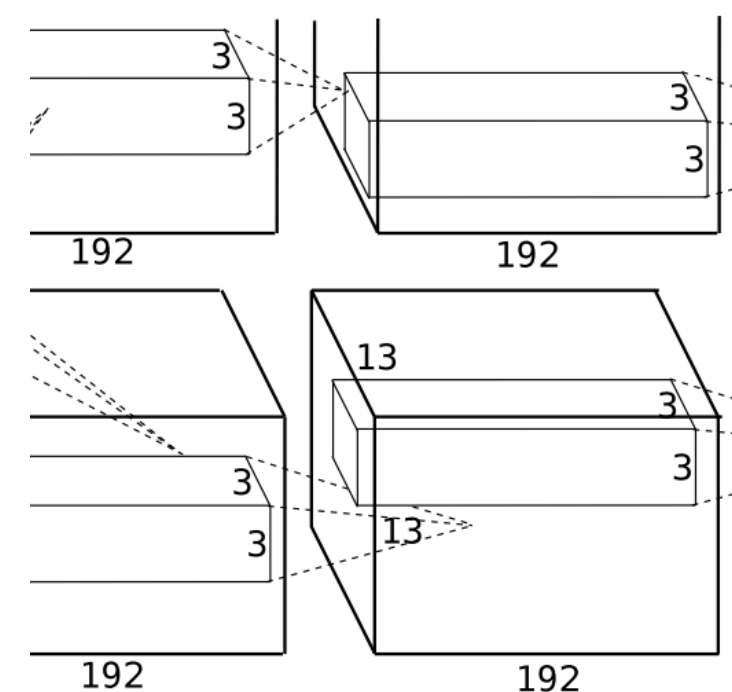
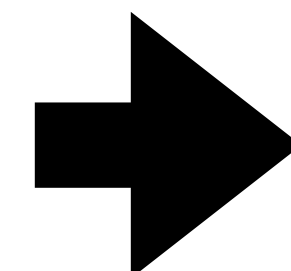
$$3 \times 3 \times 192 \times 384 = 663,552$$

13x13  
X192



**Convolution layer**  
No.filters: 384  
Filter size: 3x3  
Padding: 1  
Stride: 1

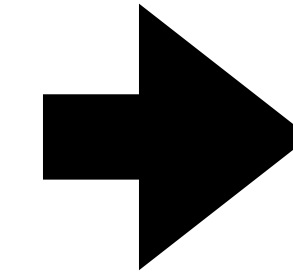
13x13  
X192



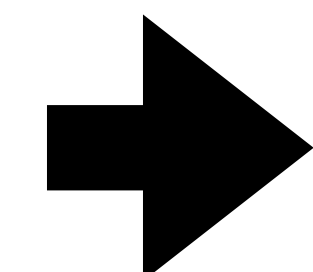
How many parameters?

$$3 \times 3 \times 192 \times 256 = 442,368$$

13x13  
X192

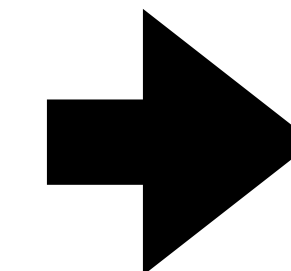


**Convolution layer**  
No.filters: 256  
Filter size: 3x3  
Padding: 1  
Stride: 1

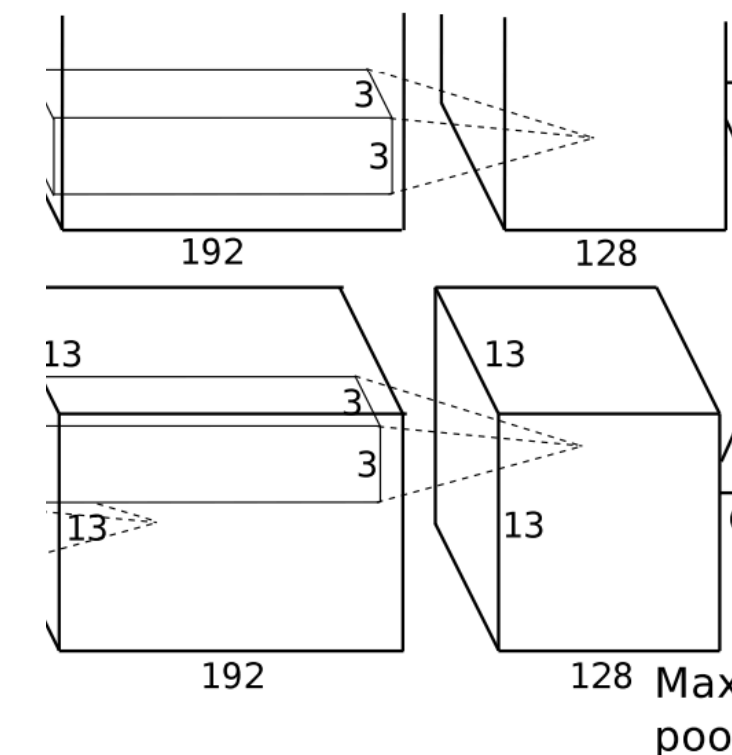
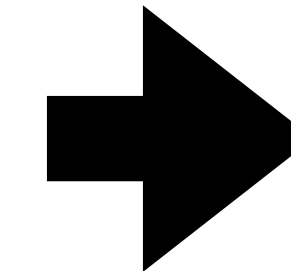


13x13  
X192

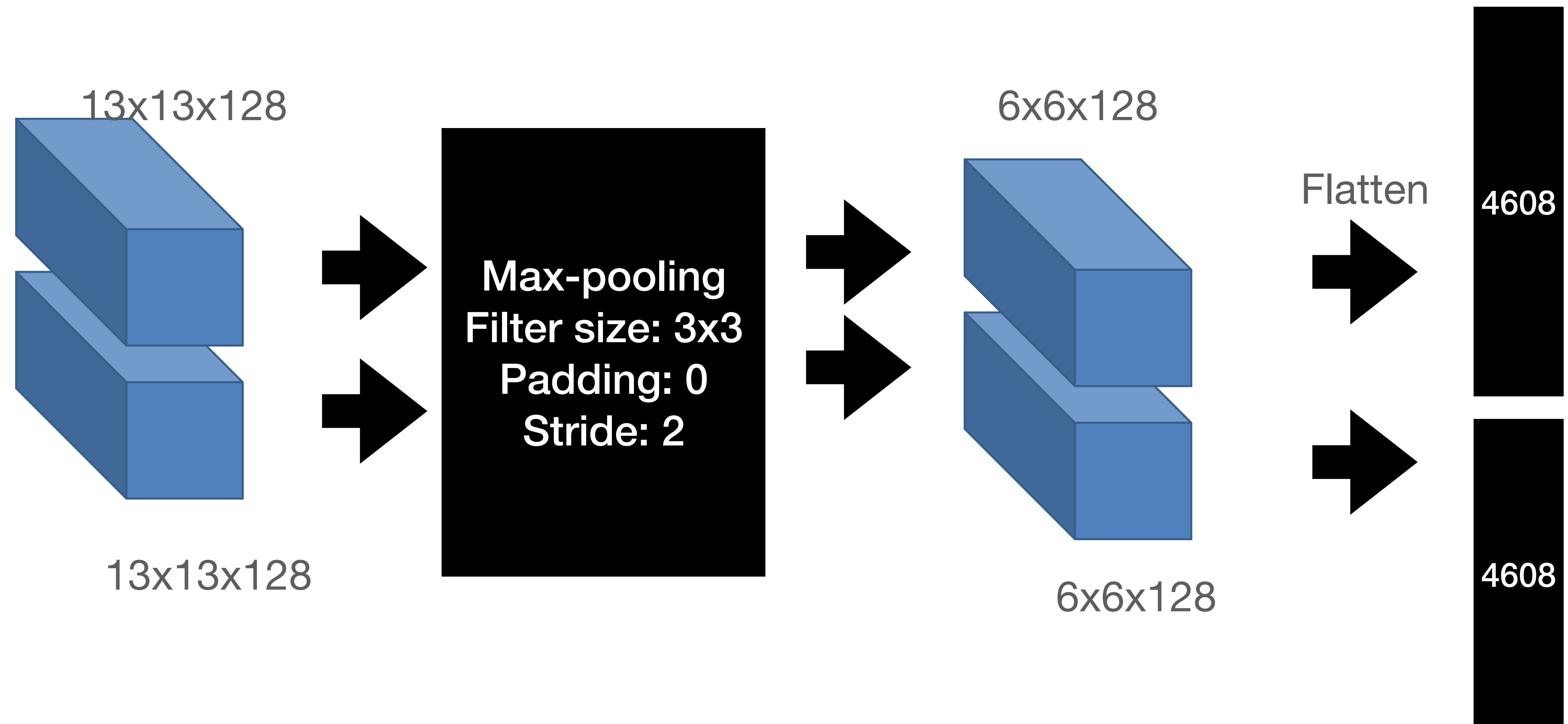
13x13  
X128



13x13  
X128

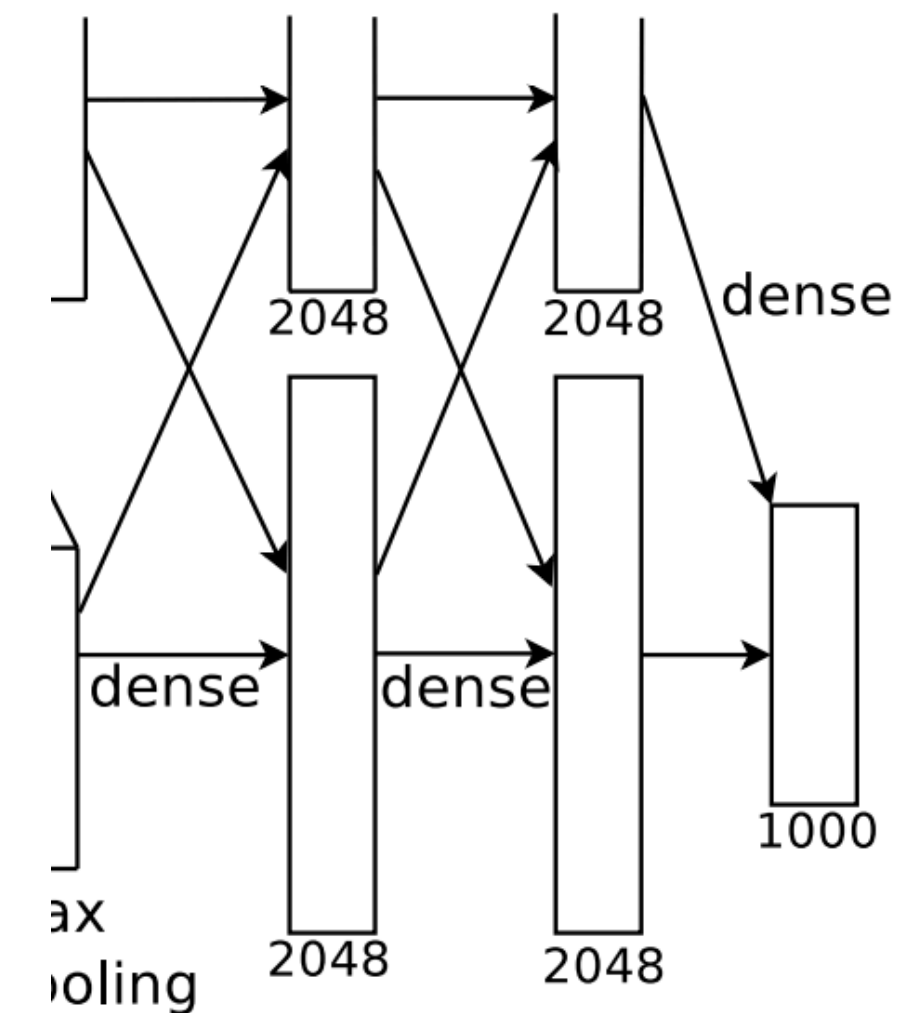
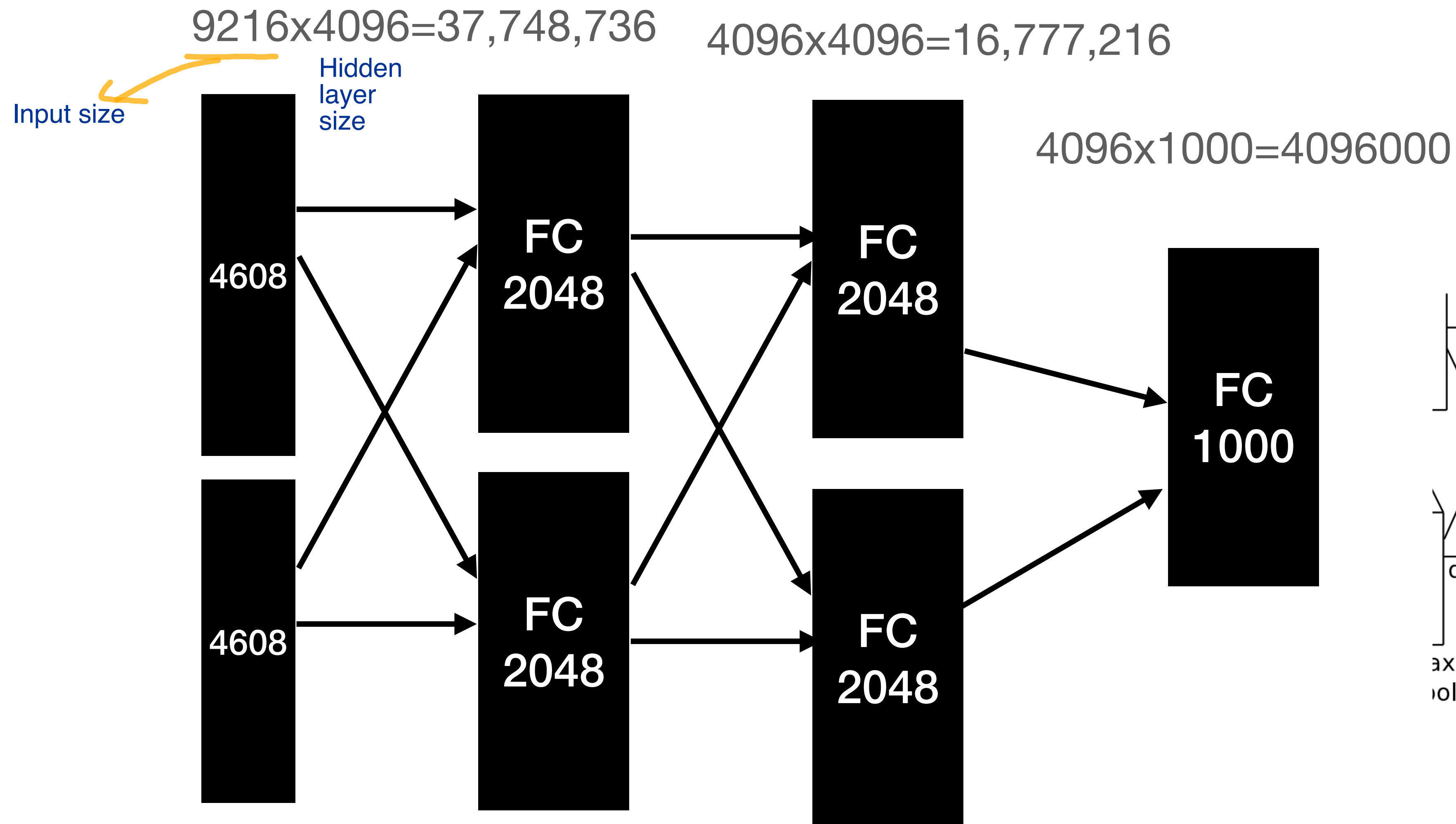


## Max-pooling and flatten



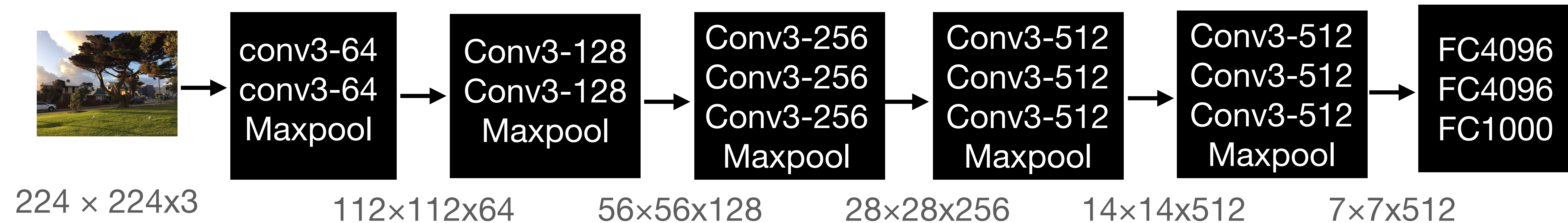
# Following fully connected layers

How many parameters?



## Architecture

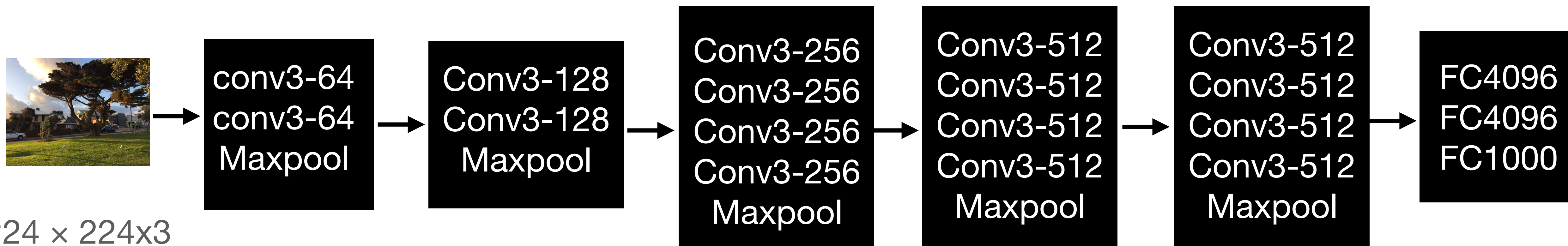
- VGG16: 16 weight layers



Conv layer: kernel size 3x3, pad 1, stride 1  
 Maxpooling layer: 2x2, stride 2

## Architecture

- VGG19: 19 weight layers



Conv layer: kernel size 3x3, pad 1, stride 1  
Maxpooling layer: 2x2, stride 2

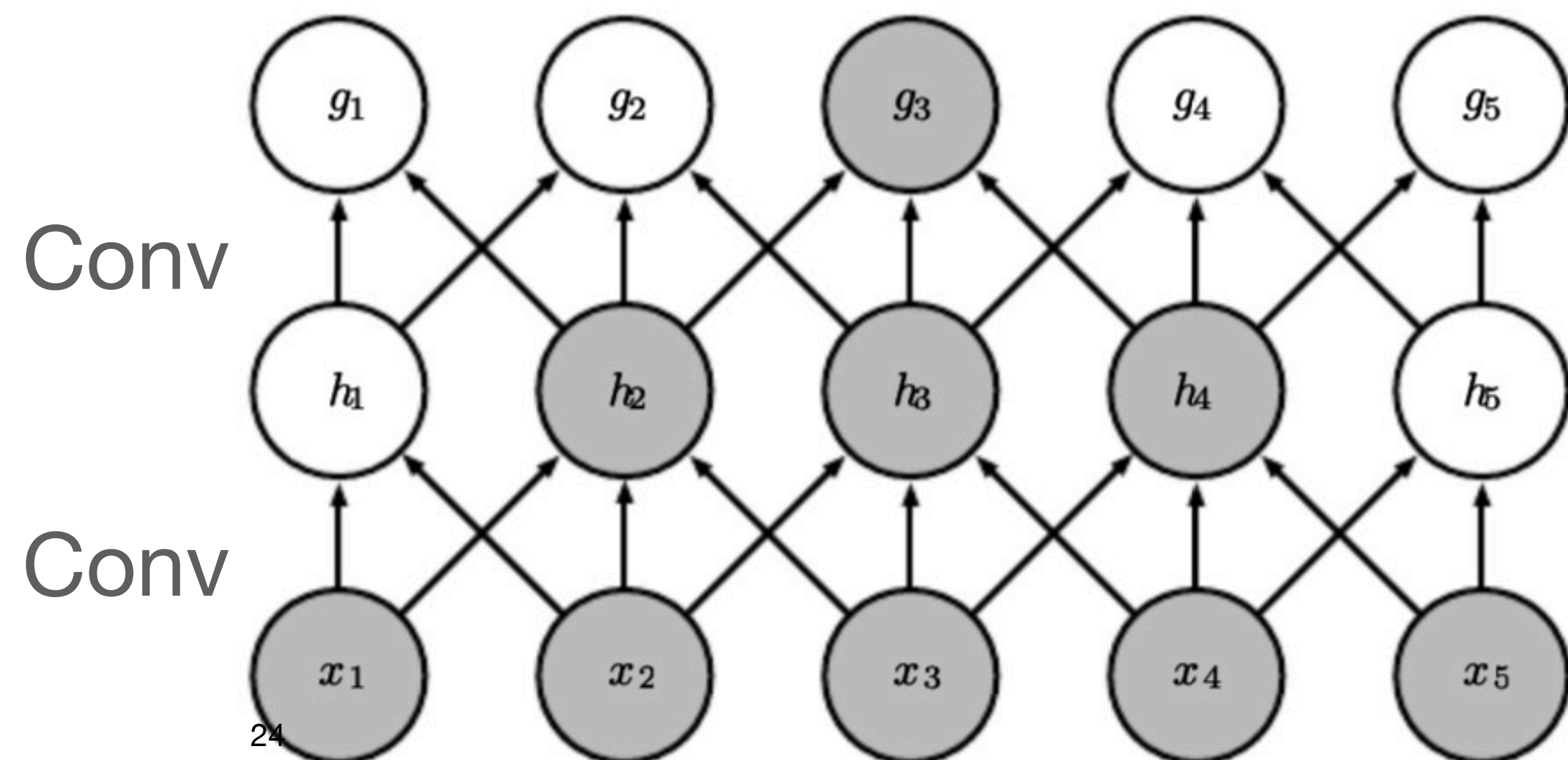


## Stacking multiple 3x3 conv layers

- a stack of two 3x3 conv. layers has an effective receptive field of 5x5
- a stack of 3 3x3 conv. layers has an effective receptive field of 7x7

If you add an additional convolutional layer with kernel size K, the receptive field is increased by (K-1)

More layers: larger size of receptive field  
(larger window of the input is seen)



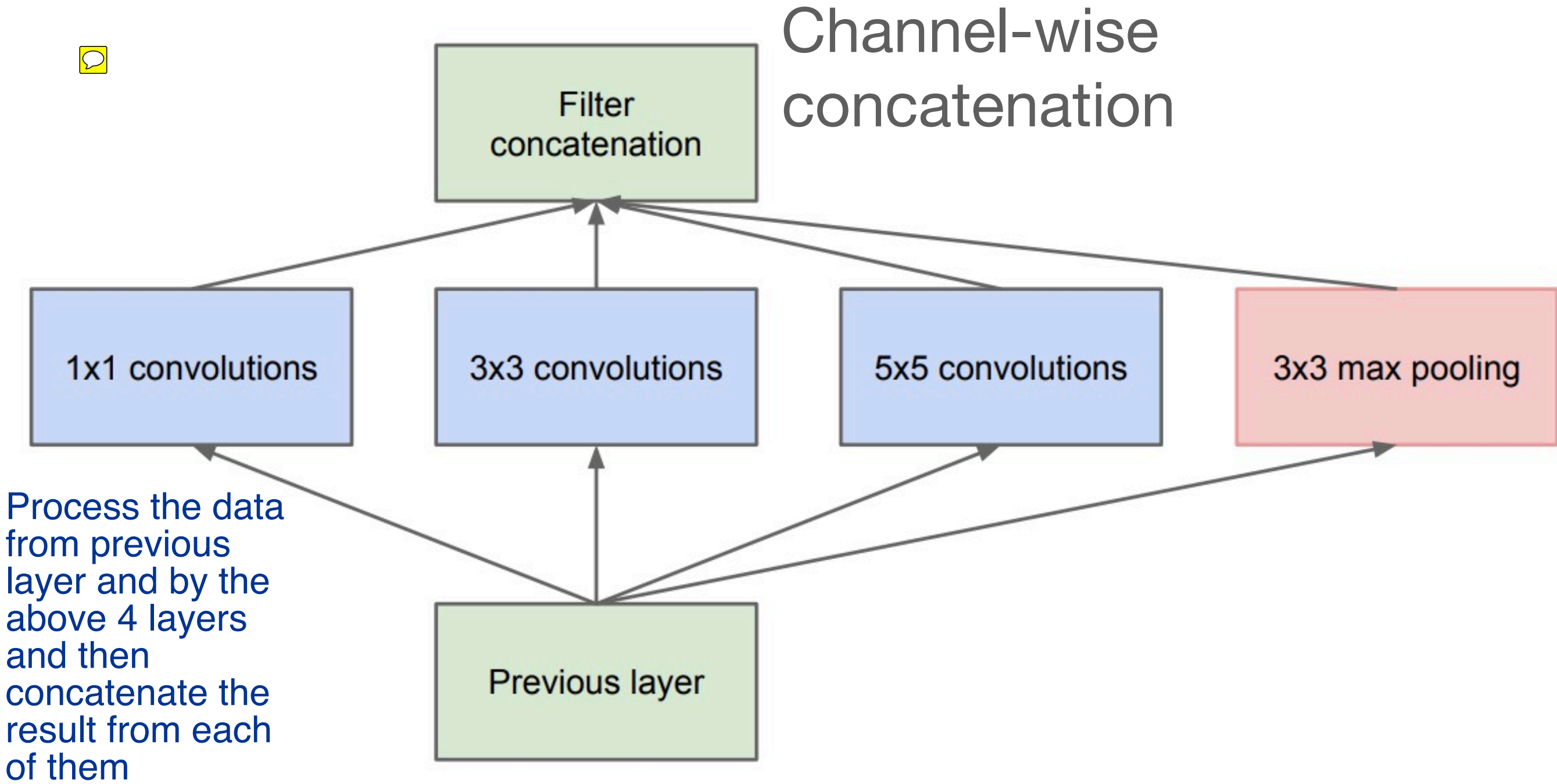
## Why Stacking multiple 3x3 conv layers instead of large filter size?

- Reduce parameter:
  - $5 \times 5 = 25$ ,  $3 \times 3 \times 2 = 18$
  - $7 \times 7 = 49$ ,  $3 \times 3 \times 3 = 27$

We could achieve a similar receptive field by stacking 3 3\*3 conv layer but with less parameters but more powerful
- Each conv use ReLU as the activation function. More layers, more non-linear rectification layers ➡ More powerful network

Architecture

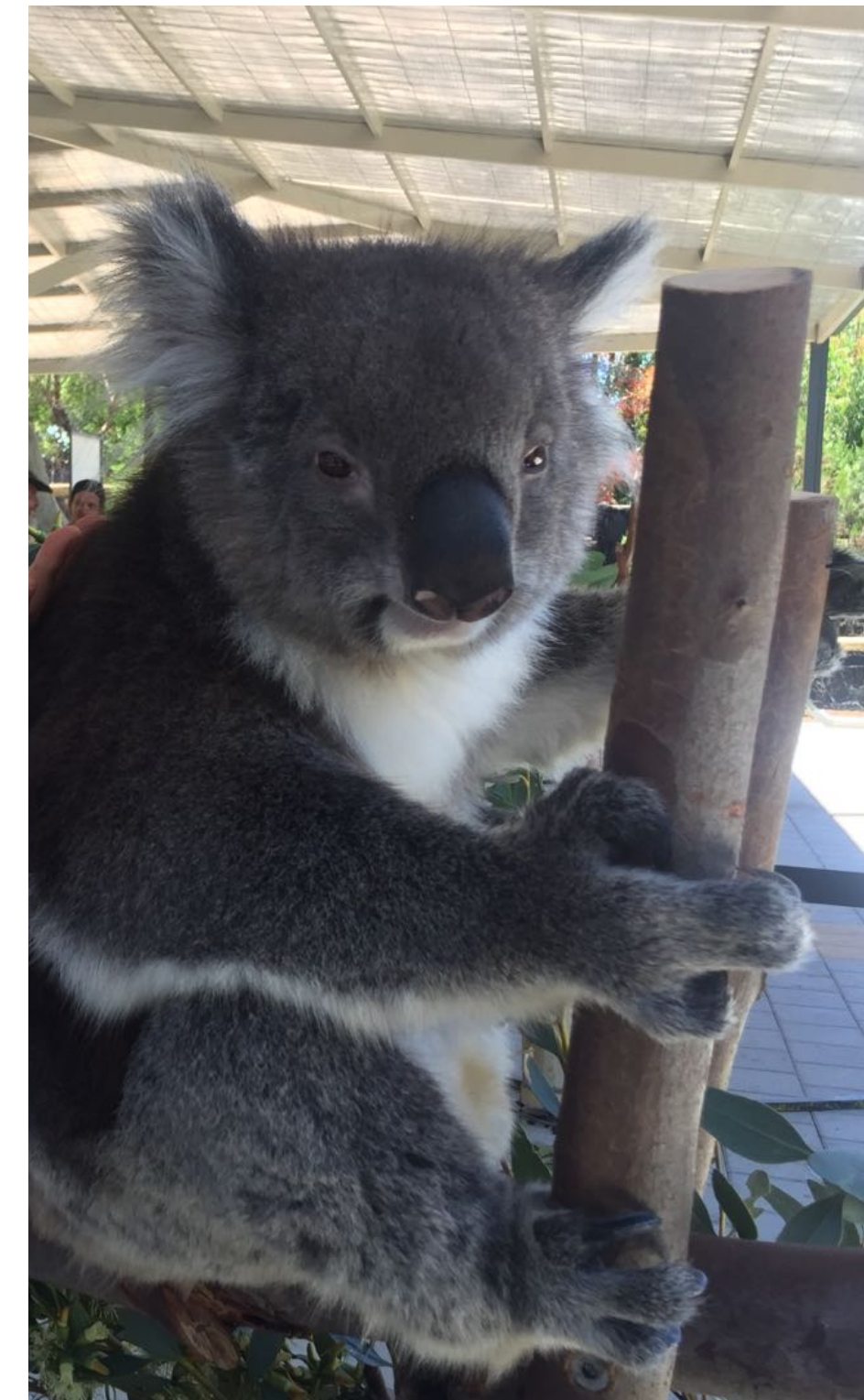
Inception module (Naive version):



type	patch size/ stride	output size	depth
convolution	7×7/2	112×112×64	1
max pool	3×3/2	56×56×64	0
convolution	3×3/1	56×56×192	2
max pool	3×3/2	28×28×192	0
inception (3a)		28×28×256	2
inception (3b)		28×28×480	2
max pool	3×3/2	14×14×480	0
inception (4a)		14×14×512	2
inception (4b)		14×14×512	2
inception (4c)		14×14×512	2
inception (4d)		14×14×528	2
inception (4e)		14×14×832	2
max pool	3×3/2	7×7×832	0
inception (5a)		7×7×832	2
inception (5b)		7×7×1024	2
avg pool	7×7/1	1×1×1024	0
dropout (40%)		1×1×1024	0
linear		1×1×1000	1
softmax		1×1×1000	0



## Different scales of data require different convolutional filter sizes



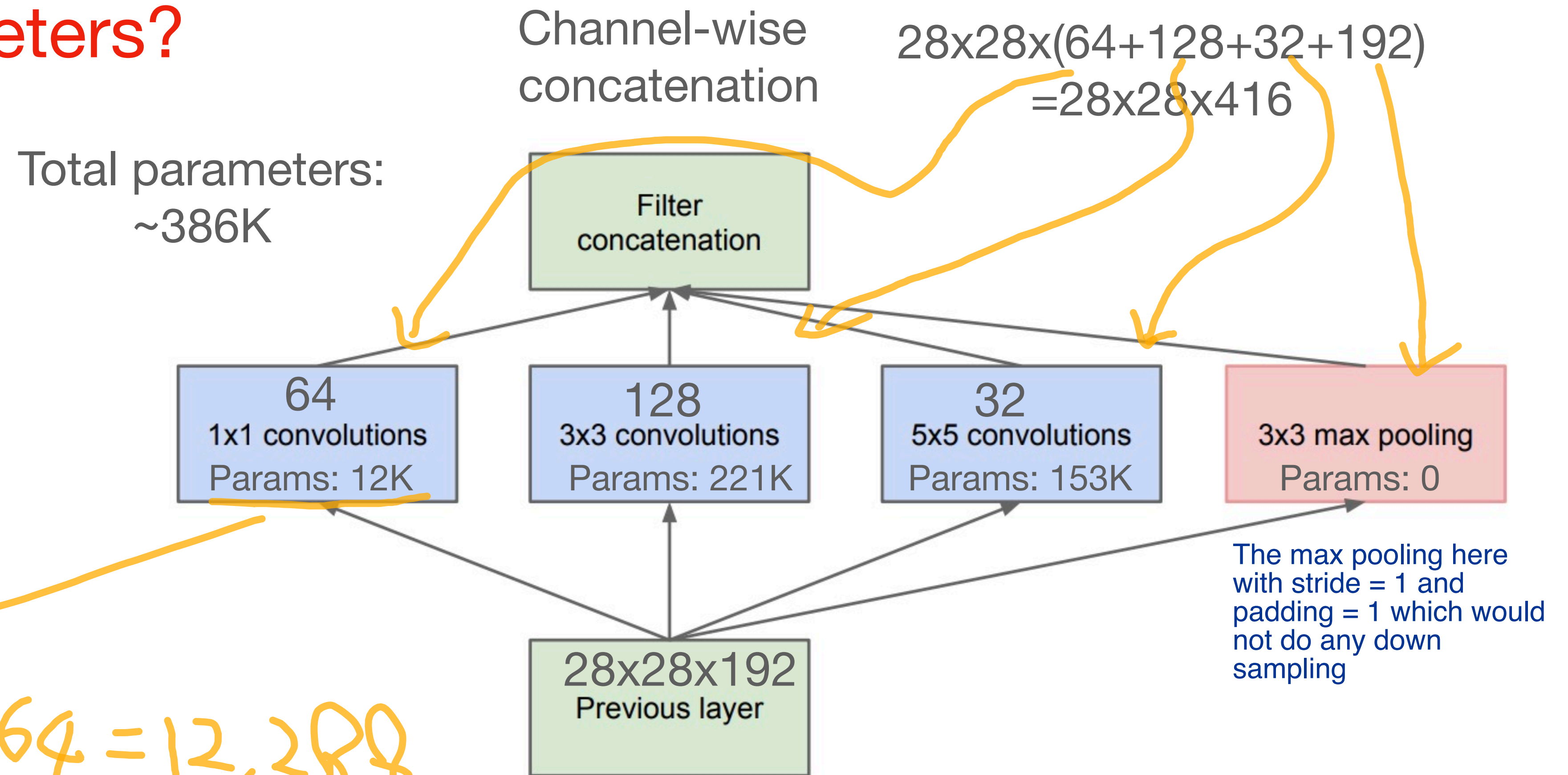


## Inception module (Naive version)

### How many parameters?

Input\_channel x  
Kernel\_size x  
NO.filters  
(output\_channel)

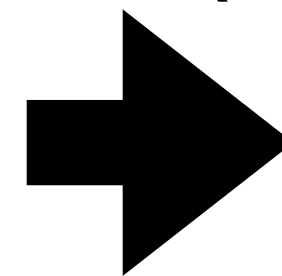
$$192 \times 1 \times 1 \times 64 = 12,288$$





## Inception module (Naive version)

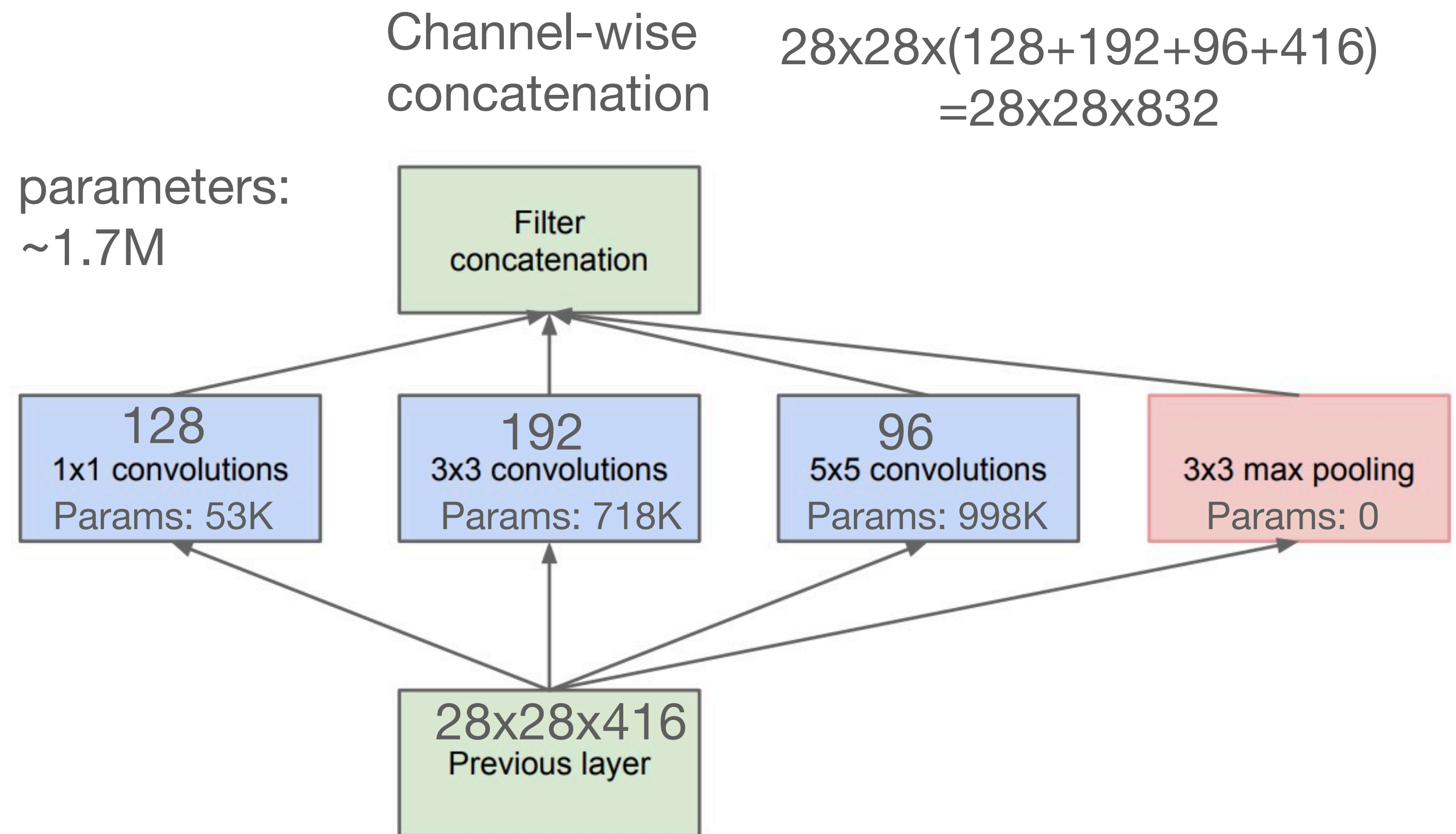
More modules



More parameters

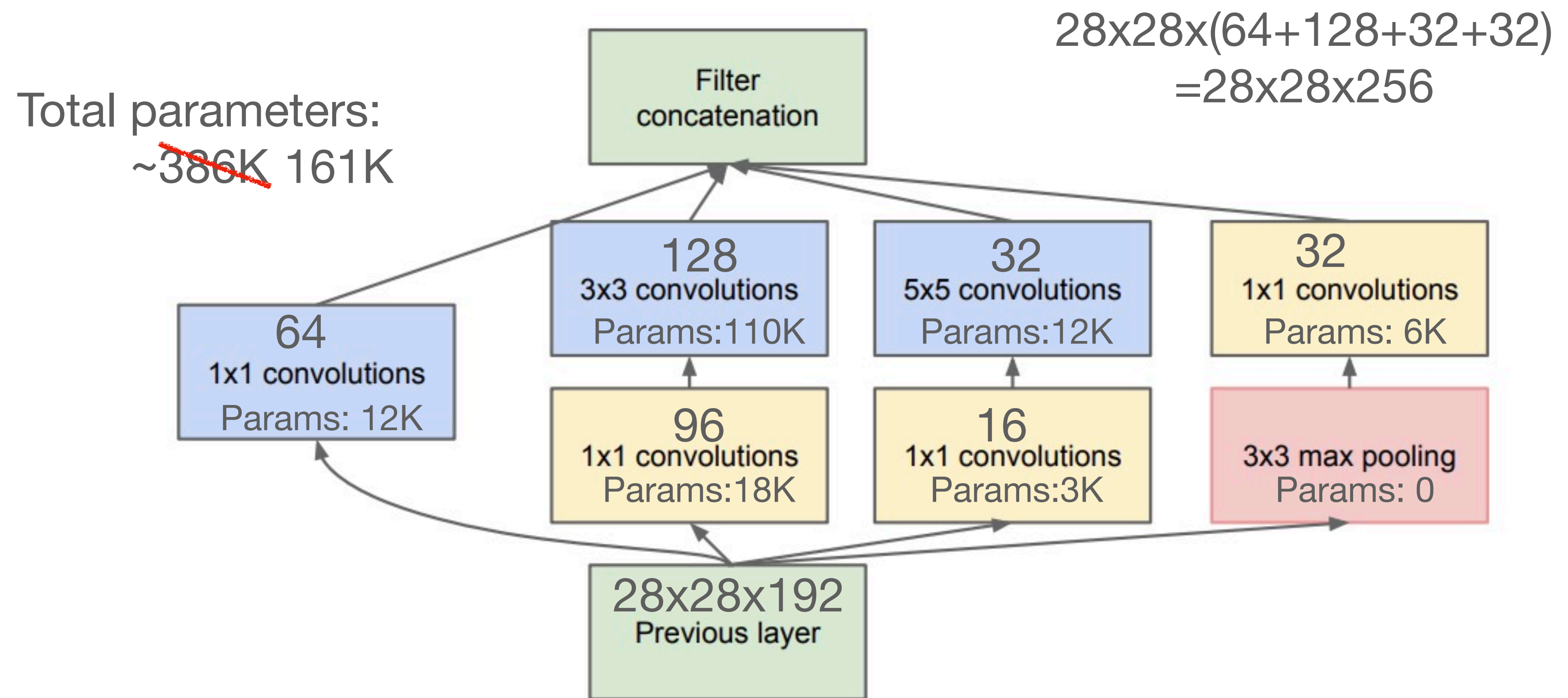
More computations

Total parameters:  
~1.7M



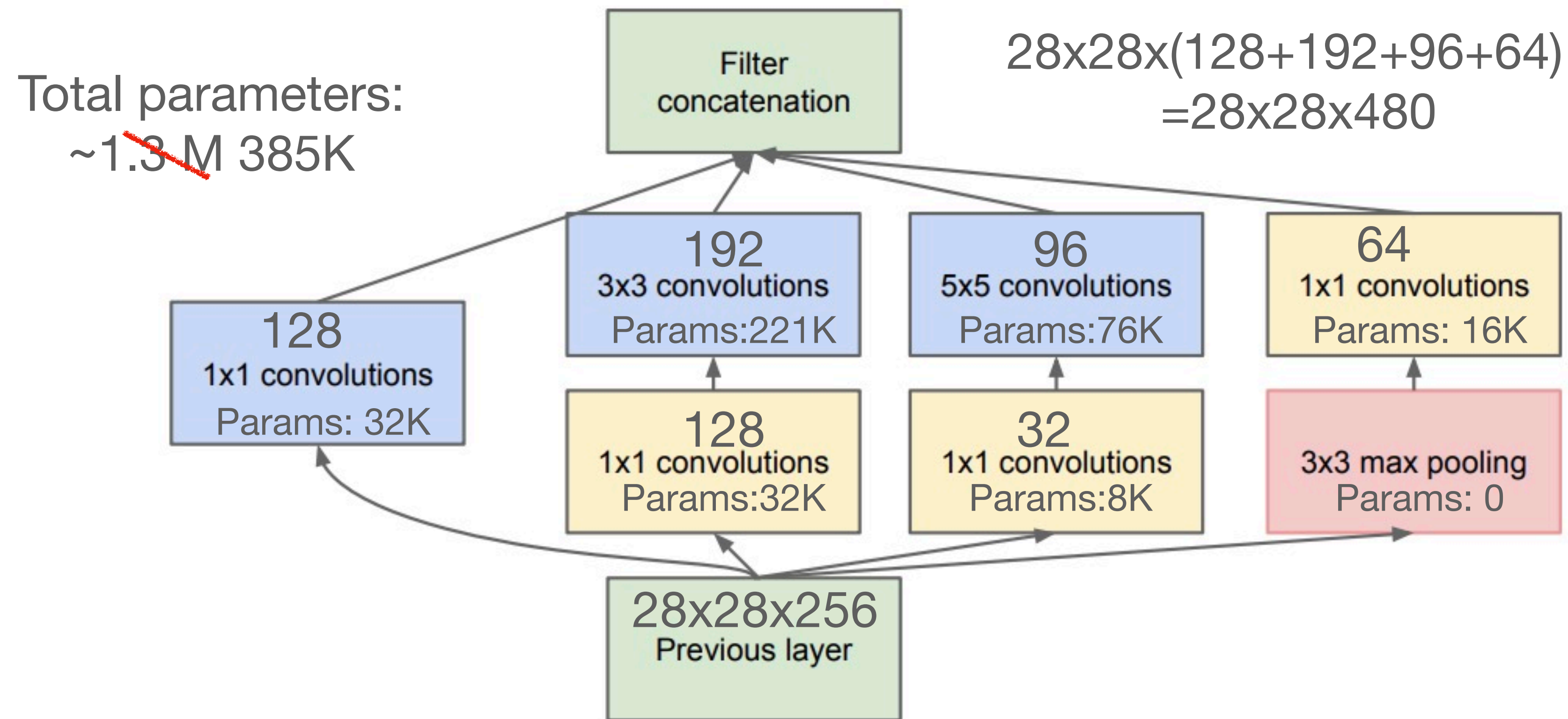
## Inception module with dimensionality reduction

Use 1x1 convolution to reduce channels



## Inception module with dimensionality reduction

Use 1x1 convolution to reduce channels





# Large Scale Visual Recognition Challenge

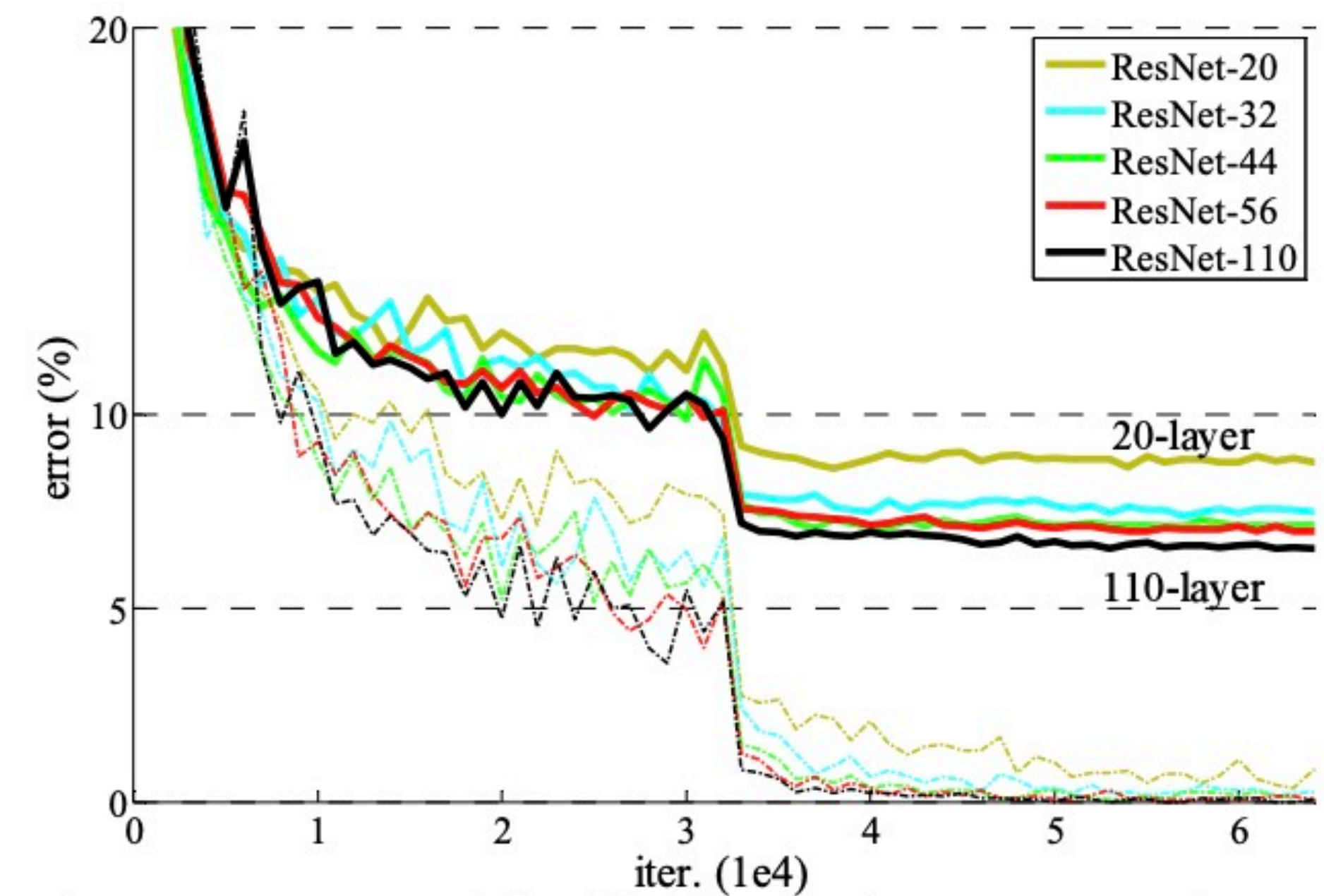
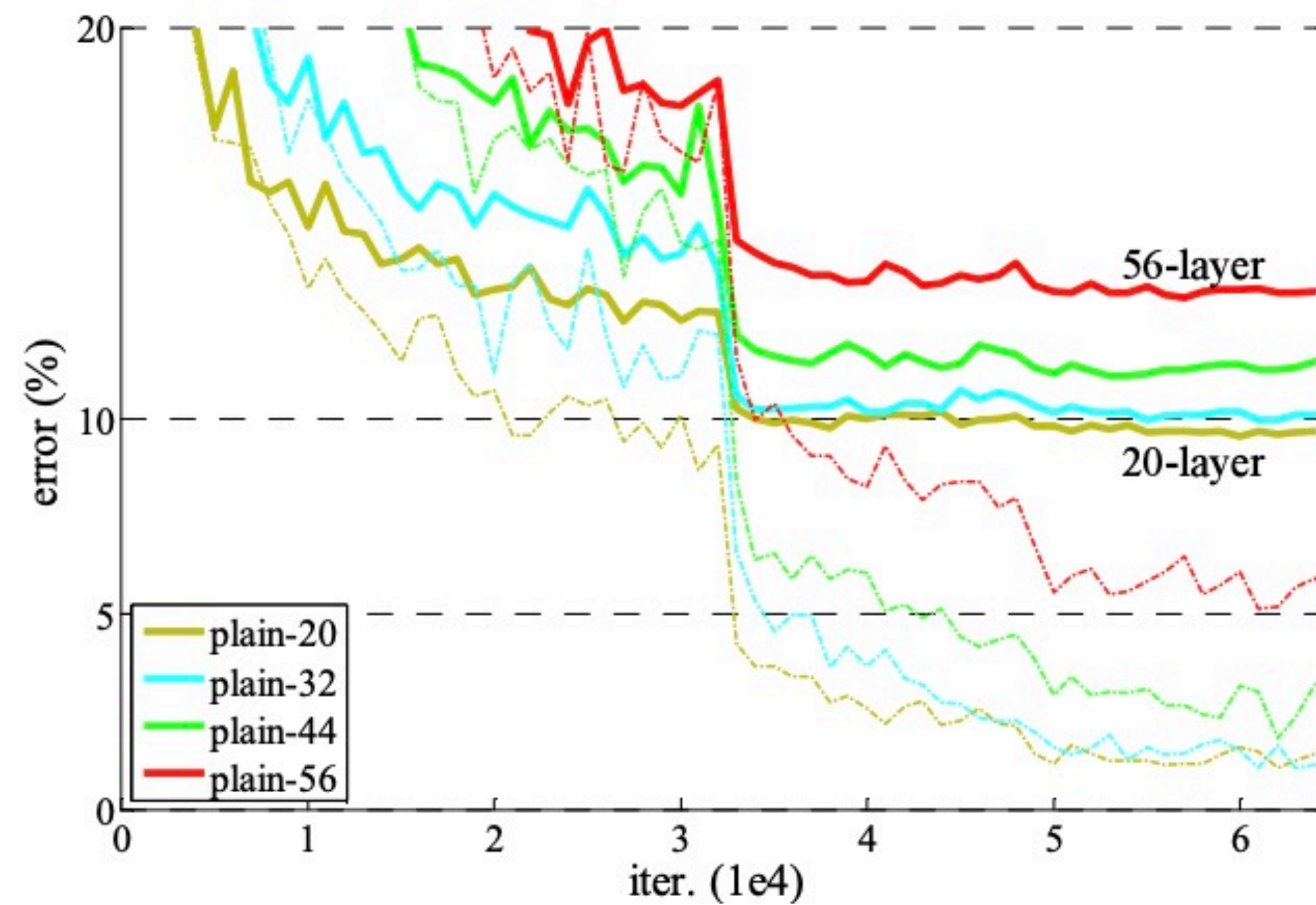
## Comparison of different architectures

- Top-5 error: the proportion of images that the ground-truth category is outside the top-5 predicted categories of the model.

CNN Architecture	Layers	Top-5 error
AlexNet	8	16.4%
VGG-19	19	7.3%
GoogleNet	22	6.7%

# More layers?

## Residual network: more layers & better performance

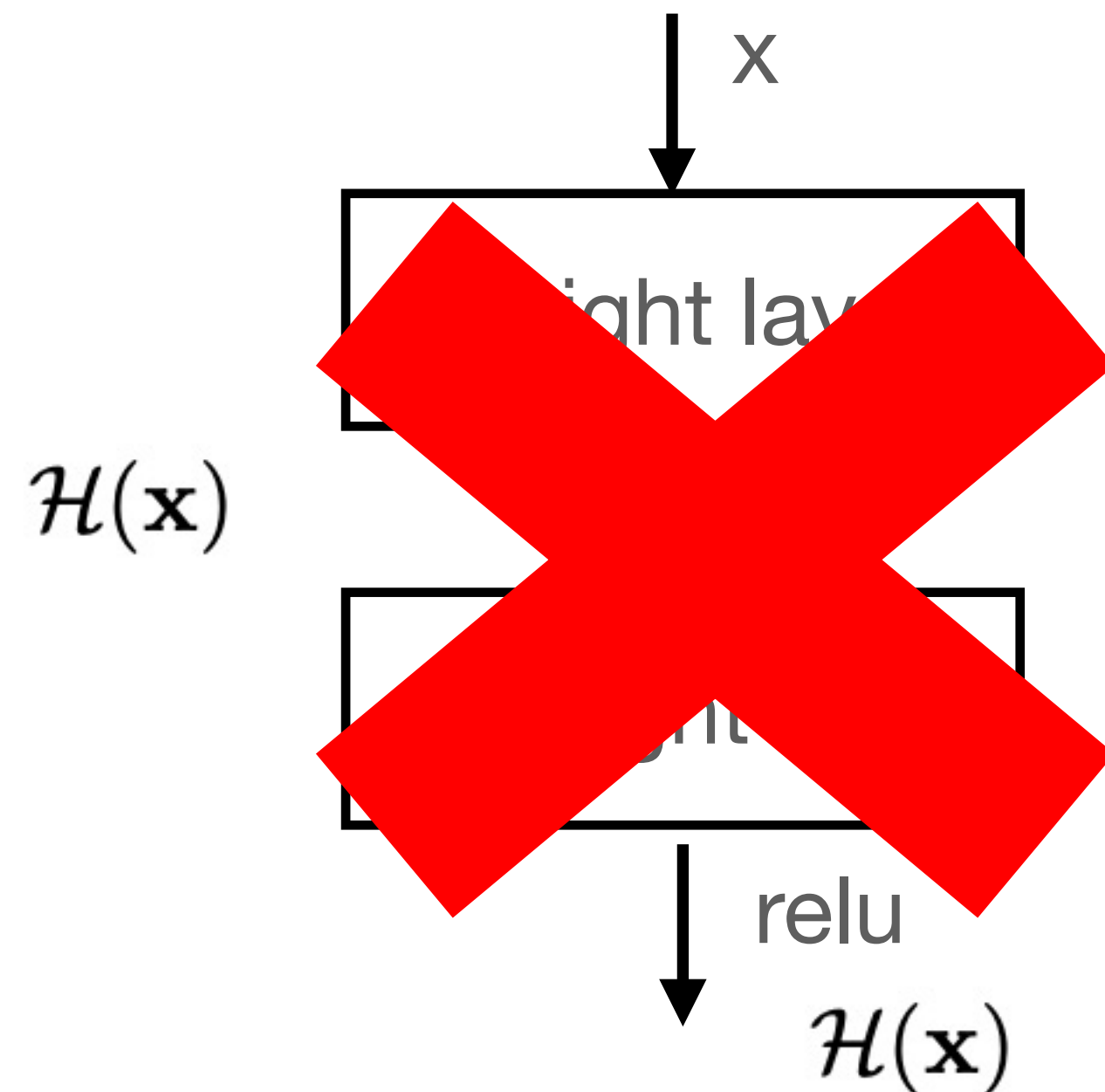


Training on CIFAR-10. Dashed lines denote training error, and bold lines denote testing error

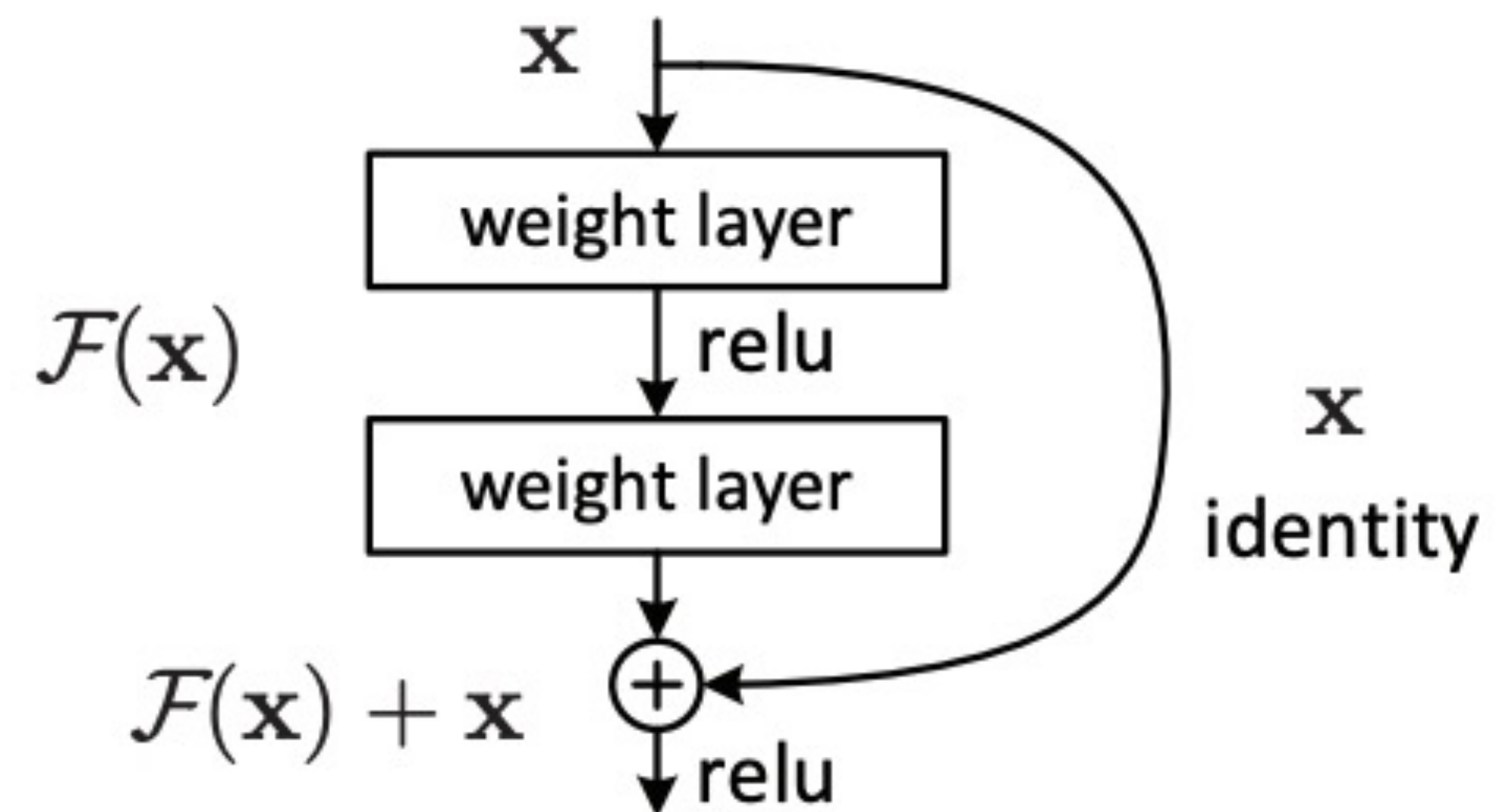


## Residual network

Hypothesis: residual mapping  $\mathcal{F}(\mathbf{x}) := \mathcal{H}(\mathbf{x}) - \mathbf{x}$  is easier to optimise

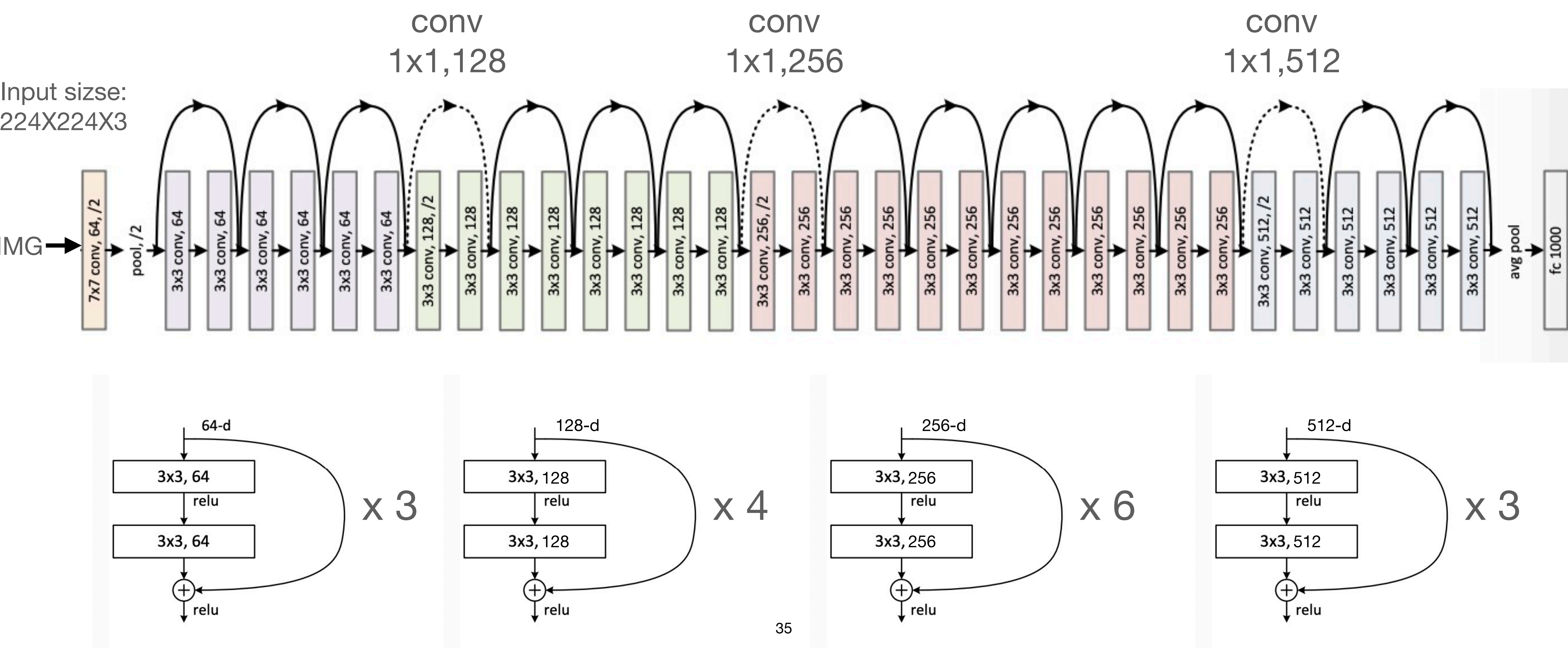


Unreferenced mapping:  
directly fit the desired underlying mapping



Residual learning:  
let these layers fit a residual mapping

# Residual network (34 layers)



# Residual network

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	$112 \times 112$	$7 \times 7, 64, \text{stride } 2$				
conv2_x	$56 \times 56$	$3 \times 3 \text{ max pool, stride } 2$				
		$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	$28 \times 28$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	$14 \times 14$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	$7 \times 7$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	$1 \times 1$	average pool, 1000-d fc, softmax				

Downsampling is performed by conv3 1, conv4 1, and conv5 1 with a stride of 2



# Large Scale Visual Recognition Challenge

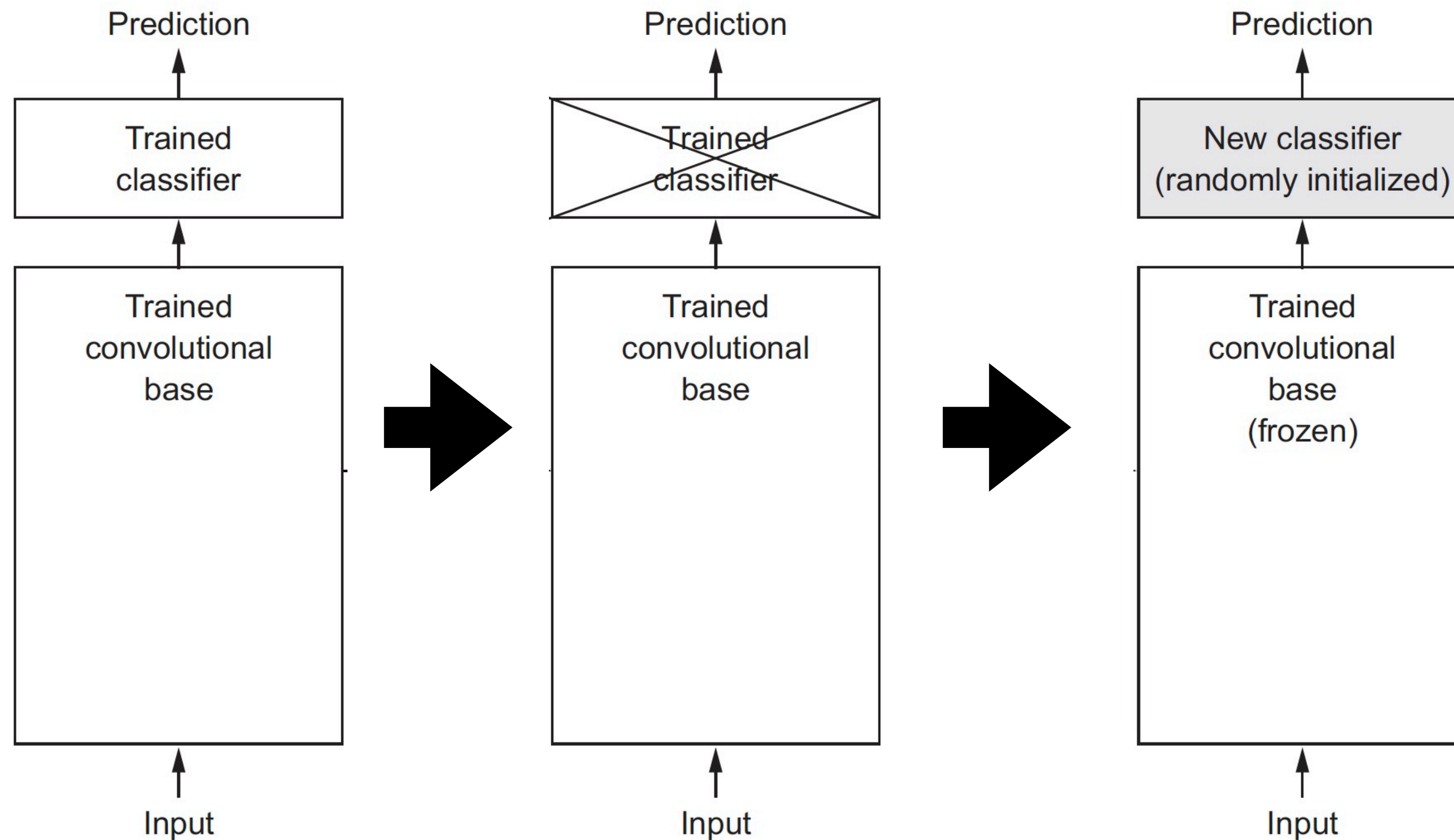
## Comparison of different architectures

- Top-5 error: the proportion of images that the ground-truth category is outside the top-5 predicted categories of the model.

CNN Architecture	Layers	Top-5 error
AlexNet	8	16.4%
VGG-19	19	7.3%
GoogLeNet	22	6.7%
ResNet	152	3.57%

# Use the pretrained CNN model as feature extractor

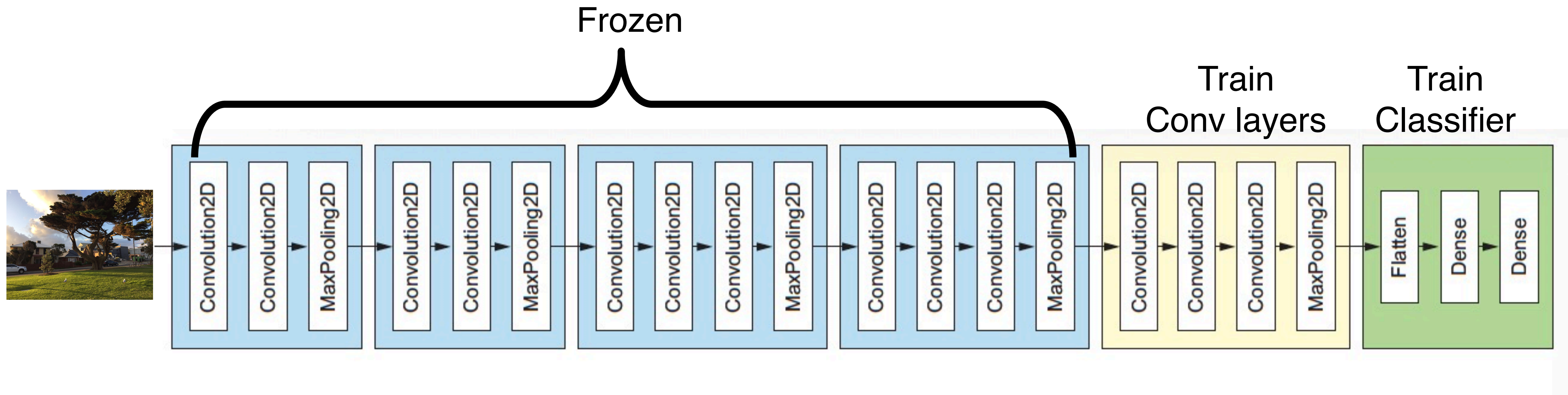
Train a new classifier for output





# If you have quite a lot of data: fine-tuning

Slightly train a few more top layers



# Summary

- How to calculate the NO. parameters & size of output feature map ?
- Difference of the architectures
- Key idea of VGG: how to increase receptive field?
- Key idea of GoogleNet: how to reduce parameters?
- Key idea of ResNet: how to increase layers with better performance?