



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

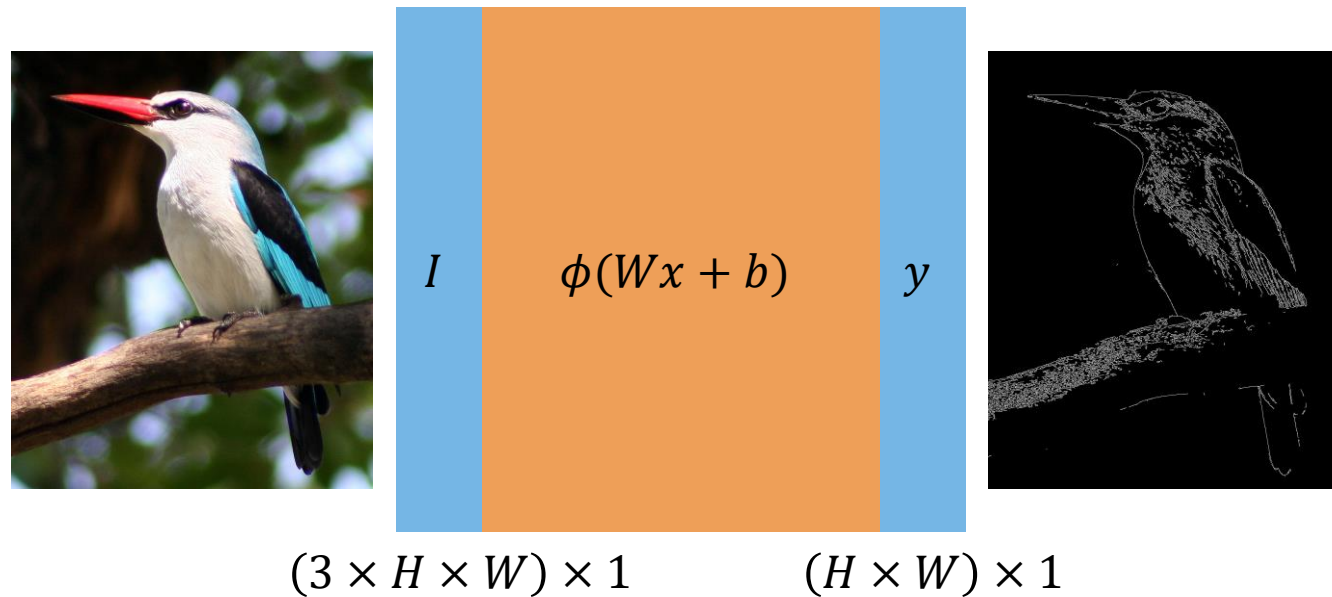
Deep Learning for Computer Vision

Luigi Di Stefano (luigi.distefano@unibo.it)

Department of Computer Science and Engineering (DISI)
University of Bologna, Italy

Limits of FC layers

Let's assume that we wish to rely on a FC layer to process an image, e.g. detect some kind of local features (e.g. edges, corners, ..)



$$H = W = 224 \rightarrow N = (3 \times H \times W) \times (H \times W) \approx 7.5 \times 10^9$$



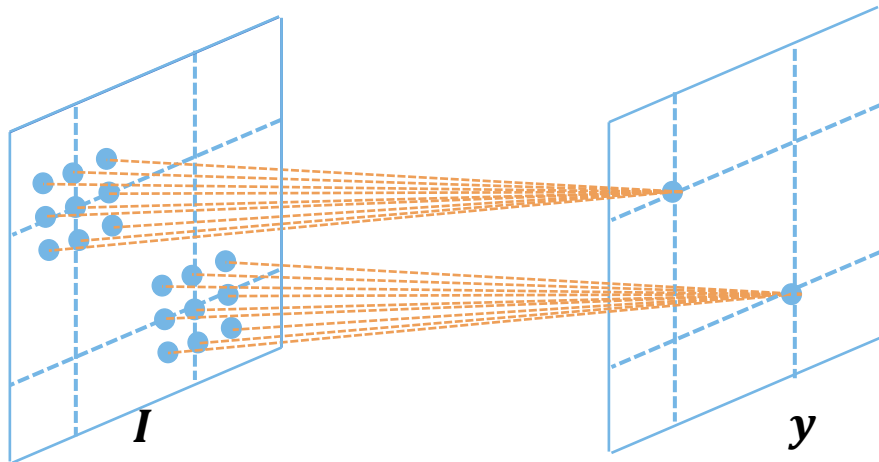
2 *Flops* (Multiply&Add) per param \approx 15 *Giga Flops*

$$H = W = 1024 \rightarrow N \approx 3.2 \times 10^{12} \rightarrow \approx 3.2/6.4 \text{ Tera Params/Flops}$$

FC layers require too many parameters (N) and *Flops* to process images (unless images are very small).

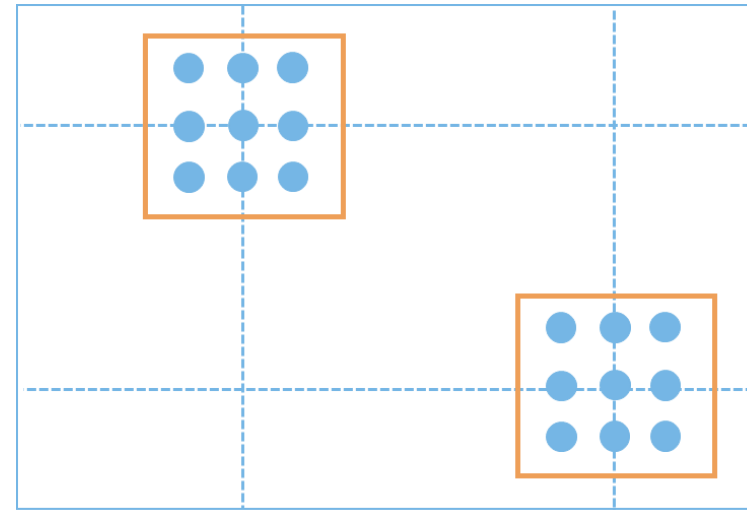
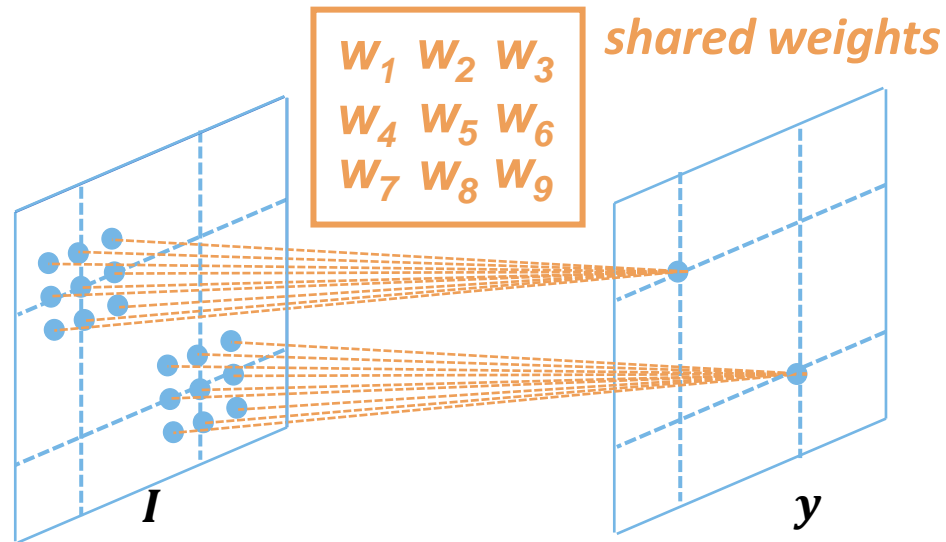
Convolution to the rescue

In image processing and classical computer vision we rely on convolution (correlation) to detect local features - as well as larger patterns- in images based on hand-crafted filters (kernels). Similarly, in deep learning we deploy **convolutional layers** to detect features and patterns based on **filters learned by minimizing a loss function**.



- In a **conv layer** the input and output are not flattened, i.e. they **preserve the spatial (2D) structure of images**.
- Unlike FC layers, in a conv layer each output unit is connected only to a –small- set of neighbouring input units. This realizes a so called **local receptive field**.
- Unlike FC layers, the weights associated with the connections between an output unit and its input neighbours are the same for all output units. Thus, **weights are** said to be **shared**.
- Conv layers embody **inductive biases** dealing with the structure of images: pixels exhibit **informative local patterns** that **may appear everywhere across the image**.

What a conv layer does actually compute ?



$$K = \begin{bmatrix} W_1 & W_2 & W_3 \\ W_4 & W_5 & W_6 \\ W_7 & W_8 & W_9 \end{bmatrix}$$

$$y(i, j) = K * I(i, j) = \sum_m \sum_l K(m, l) I(i + m, j + l)$$

➡ **Convolution (Correlation) !**

$$N = (3 \times 3)$$

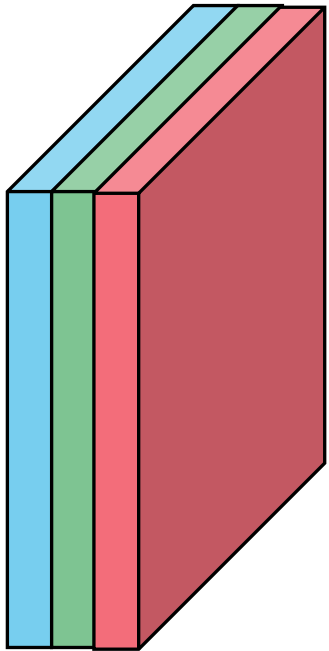
In a conv layer the number of parameters is small and independent of image size:

$$N = (H_k \times W_k)$$

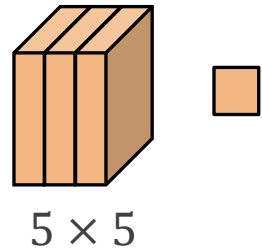
$$\text{Flops} = 2 \times (H_k \times W_k) \times (H \times W)$$

$$\text{e.g., } H = W = 1024, N = (3 \times 3) \rightarrow 9 \text{ Params, } \approx 18 \text{MFlops}$$

Multiple input channels



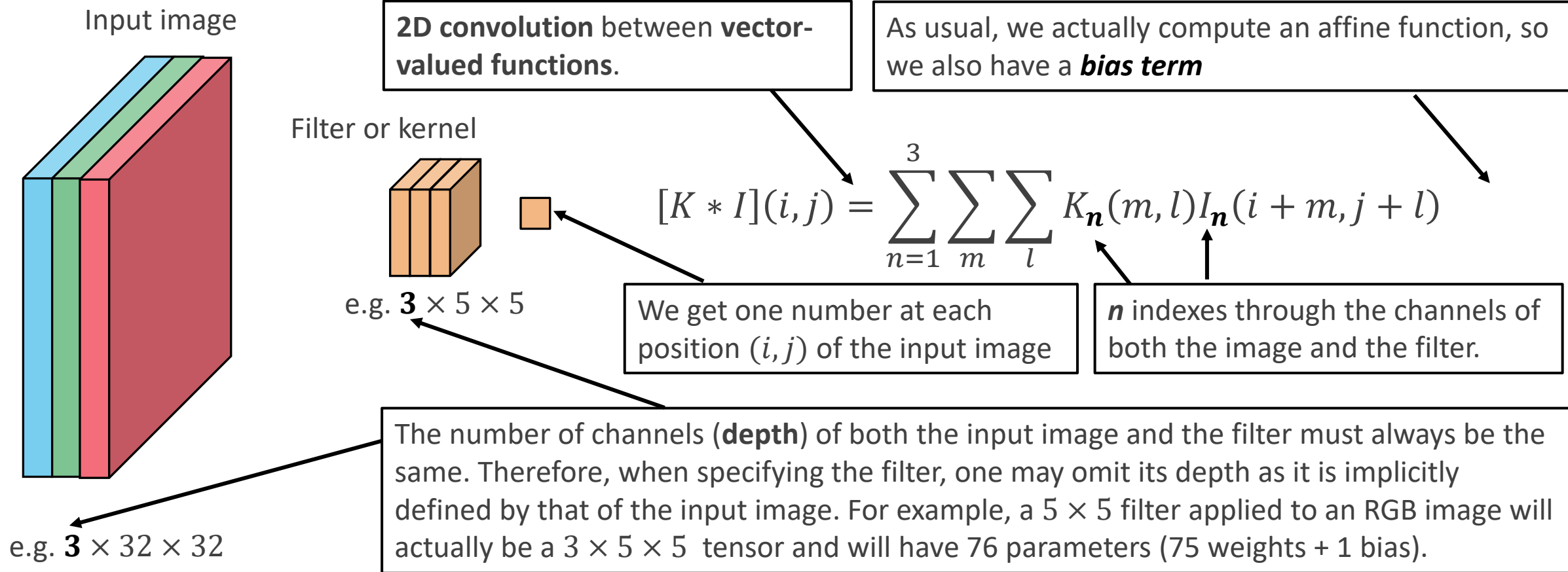
Filter or kernel



$$[K * I](i, j) = \sum_m \sum_l K(m, l) I(i + m, j + l)$$

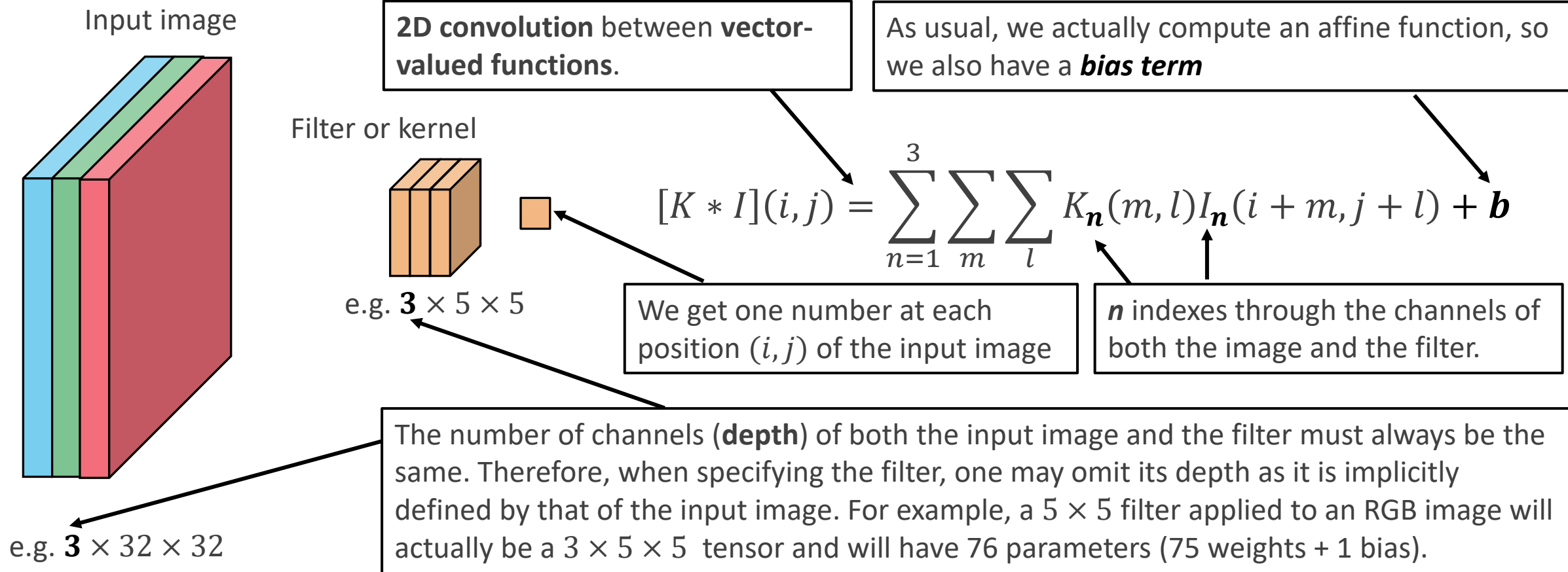
Multiple input channels

RGB images have 3 channels, so convolution kernels must be 3-dimensional tensors of size $3 \times H_K \times W_K$



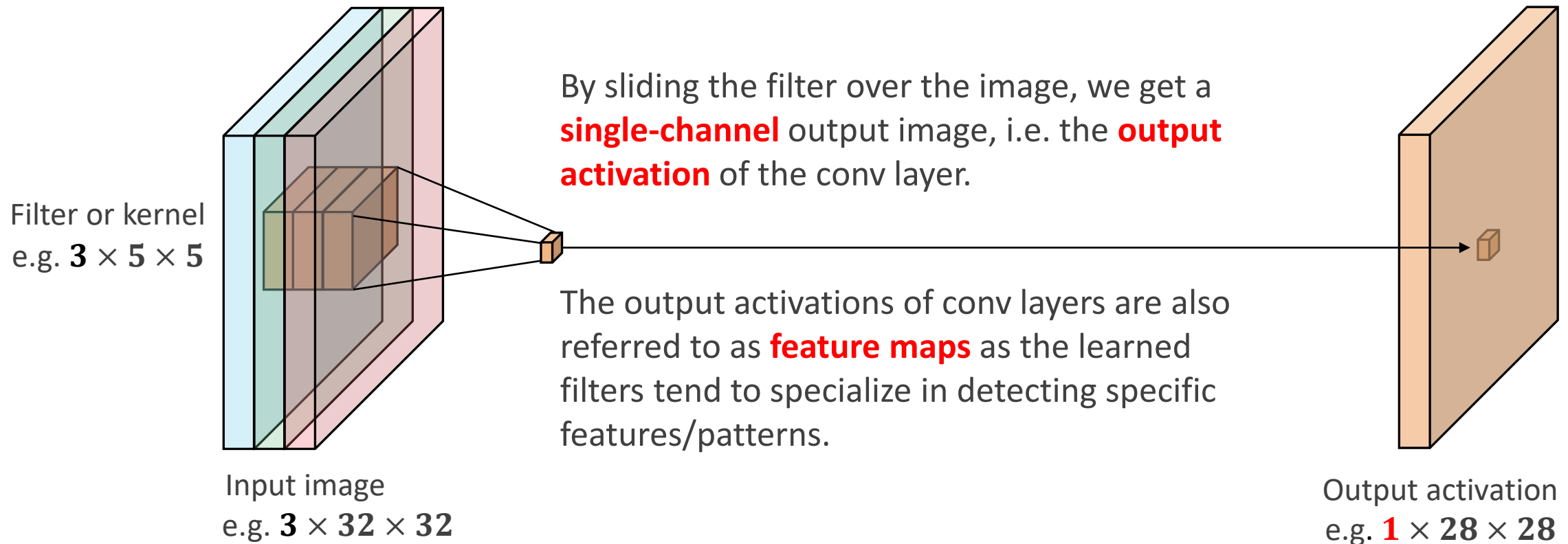
Multiple input channels

RGB images have 3 channels, so convolution kernels must be 3-dimensional tensors of size $3 \times H_K \times W_K$



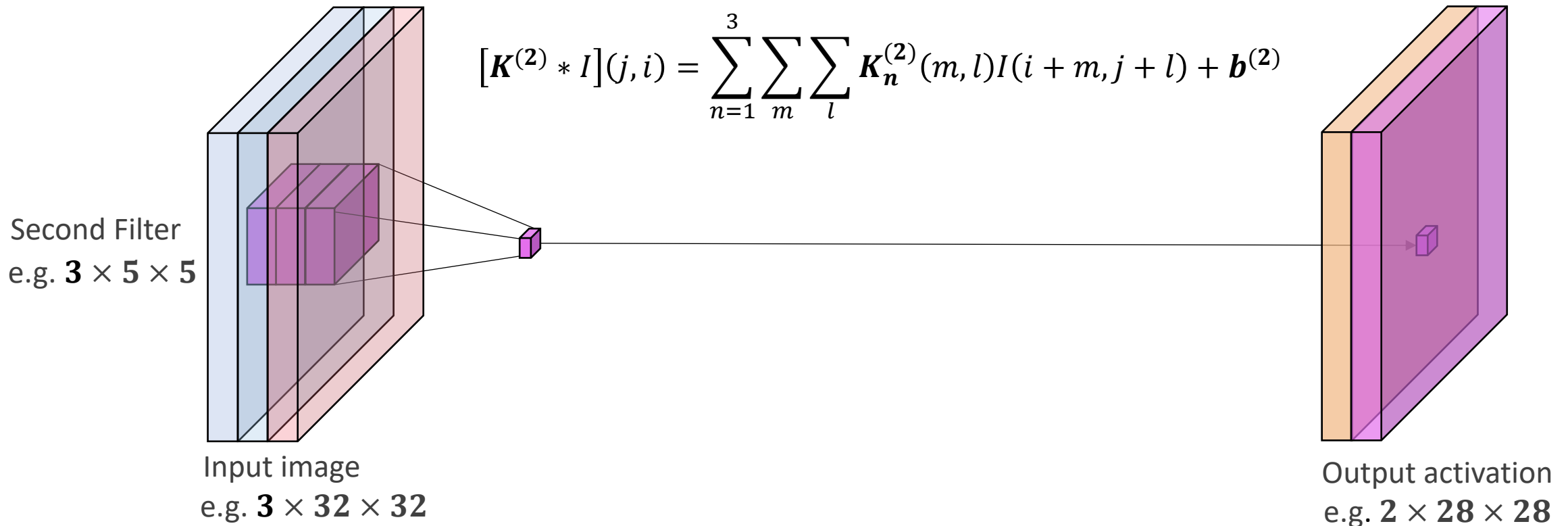
Output Activations – Feature Maps

$$[K * I](i, j) = \sum_{n=1}^3 \sum_m \sum_l K_n(m, l) I_n(i + m, j + l) + b$$



Multiple Output Channels (1)

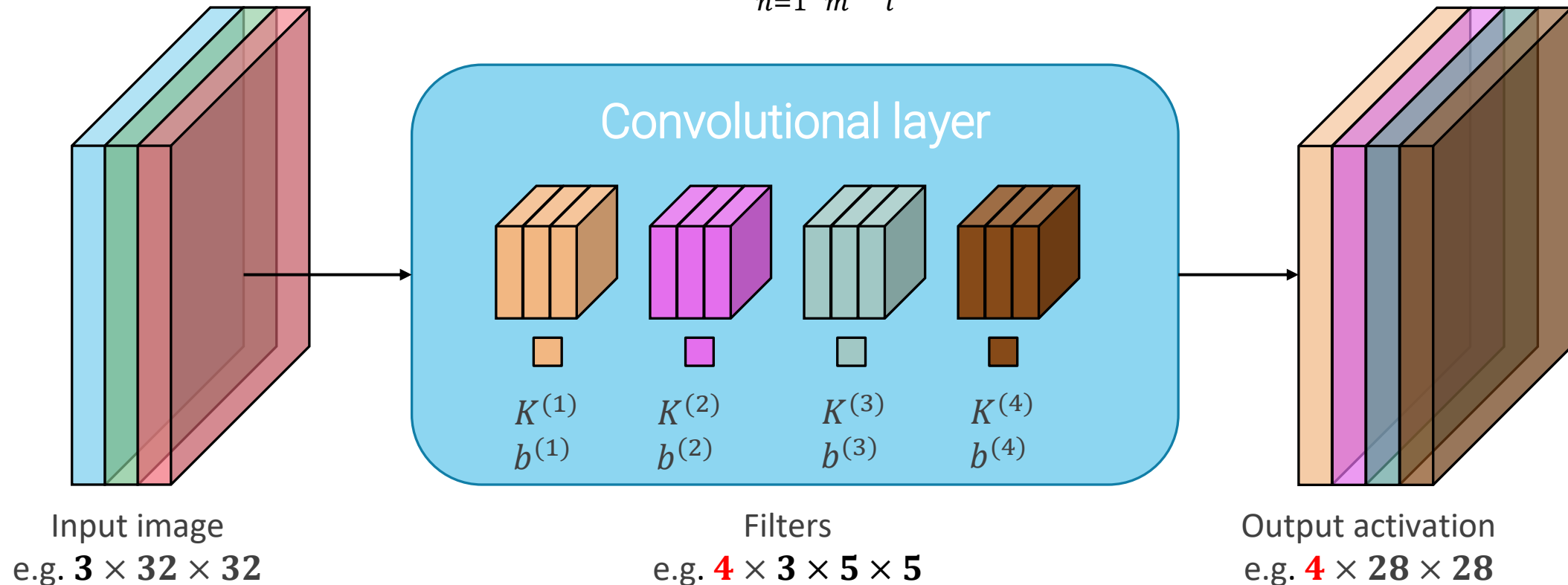
It may be useful to obtain a multi-channel activation by applying different filters with the same size and different weights within the same conv layer. For example, we may deploy two filters such that the conv layer would have the ability to detect two kinds of features, e.g. horizontal and vertical edges.



Multiple Output Channels (2)

If we want an even more powerful conv layer we may apply, e.g. , four filters. The whole operation realized by the layer can be expressed as:

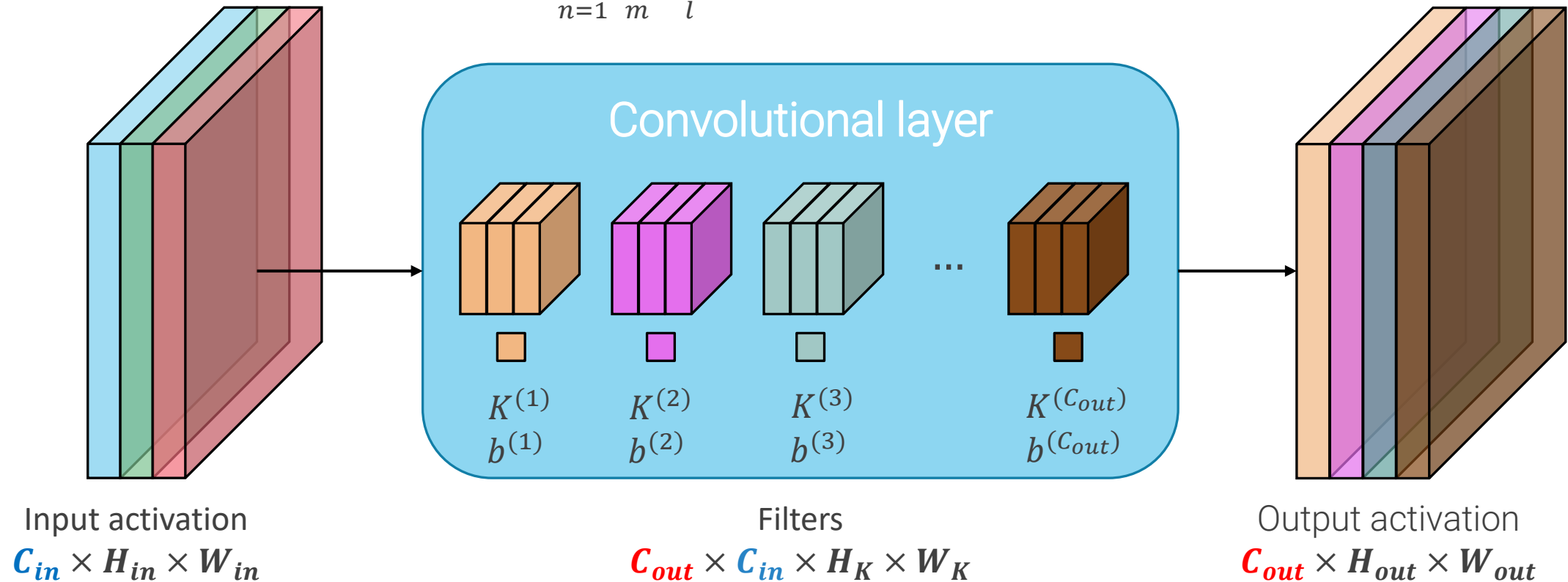
$$[K * I]_{\mathbf{k}}(j, i) = \sum_{n=1}^3 \sum_m \sum_l K_n^{(\mathbf{k})}(m, l) I(i + m, j + l) + b^{(\mathbf{k})} \quad \mathbf{k} = 1, \dots, 4$$



General structure of a convolutional layer

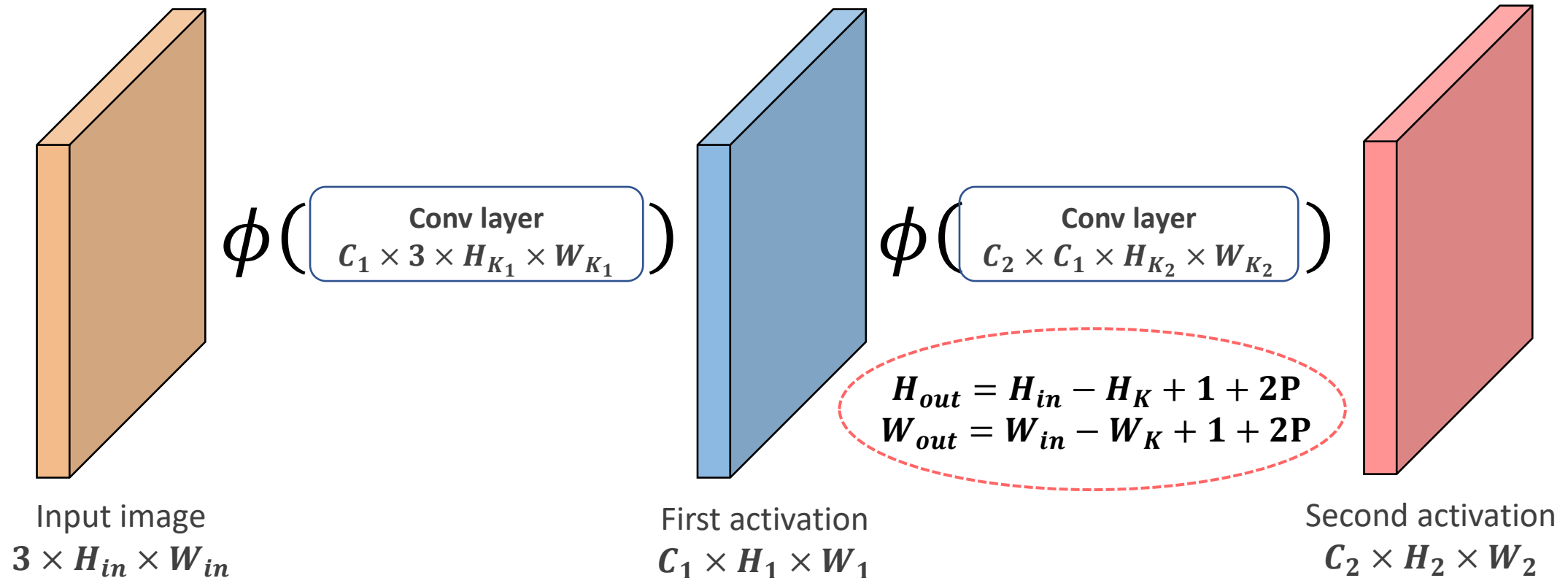
A conv layer receives a multi-channel (C_{in}) input activation and produces a multi-channel (C_{out}) output activation by applying as many filters as the output channels, with the depth of the filters given by the number of input channels.

$$[K * I]_k(j, i) = \sum_{n=1}^{C_{in}} \sum_m \sum_l K_n^{(k)}(m, l) I_n(i + m, j + l) + b^{(k)} \quad k = 1, \dots, C_{out}$$



Chaining convolutional layers

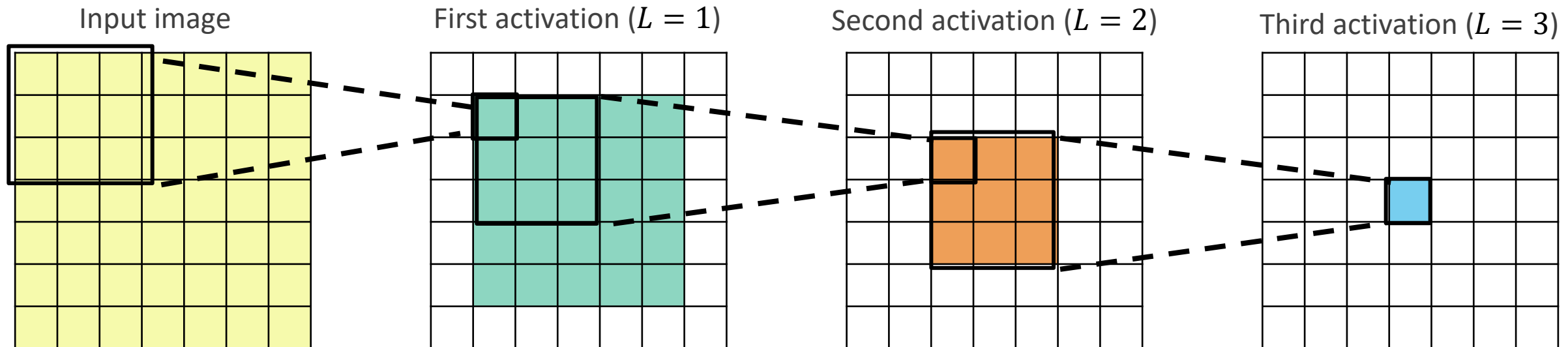
A convolutional layer is a special form of linear layer (indeed it can be expressed as a matrix multiplication). Thus, to take advantage of depth by chaining multiple layers we need to introduce non-linear activations (typically **ReLU**). Moreover, to avoid shrinking the activations along the chain we may **(zero)pad** the input to each layer.



Receptive Field

As already pointed out while studying classical computer vision algorithms, the set of **input pixels** affecting the value of a hidden unit is referred to as the **receptive field** of the unit. As we traverse a chain of conv layers the receptive field gets larger and larger, so as to compute features dealing with larger and larger image regions (from *local* to *global* features).

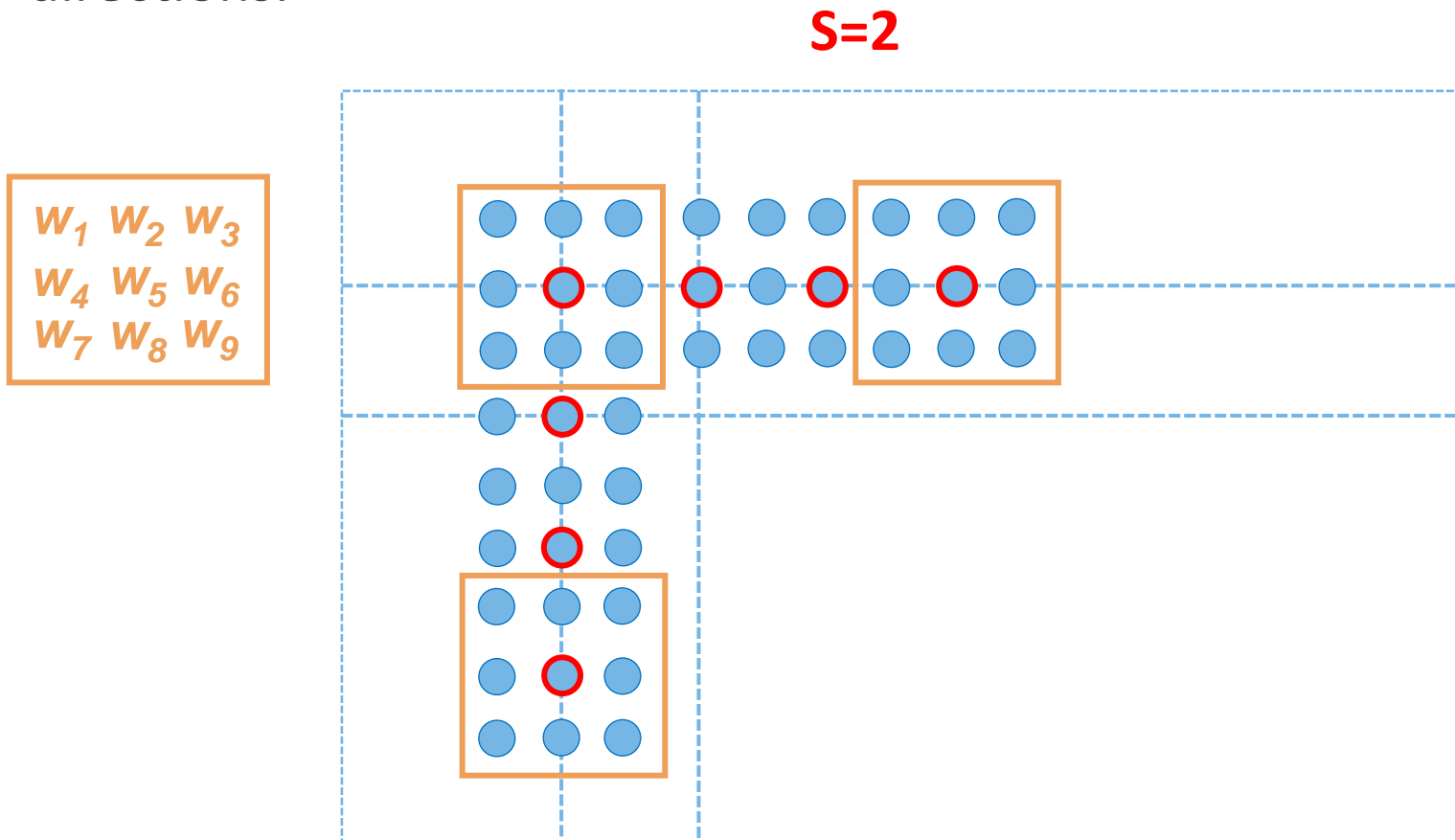
E.g., if the kernel size is $H_K \times W_K$, the size of the receptive field at the L -th activation is $[1 + L(H_K - 1)] \times [1 + L(W_K - 1)]$



Thus, both the height and width of the receptive field grow linearly with the number of layers. To obtain **larger receptive fields** with a limited number of layers we **down-sample the activations**.

Strided Convolution

Rather than densely, convolution may be computed every **S (stride)** positions in both directions.



If the input activation is zero-padded according to the size of the filter so to avoid shrinking the output, the actual size of the down-sampled activation computed by a strided convolution is given by:

$$H_{out} = \left\lfloor \frac{H_{in}-1}{S} \right\rfloor + 1$$

$$W_{out} = \left\lfloor \frac{W_{in}-1}{S} \right\rfloor + 1$$

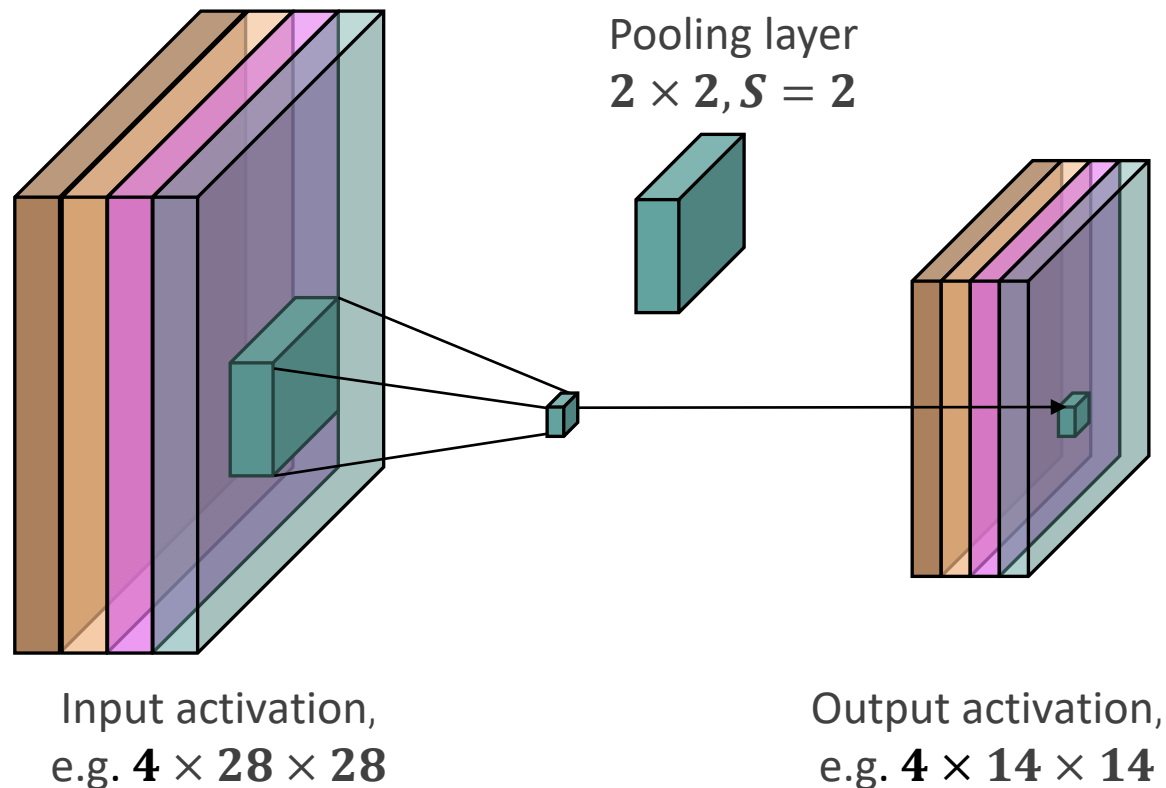
In general, with a kernel $H_k \times W_k$, padding P and stride S :

$$H_{out} = \left\lfloor \frac{H_{in}-H_k+2P}{S} \right\rfloor + 1$$

$$W_{out} = \left\lfloor \frac{W_{in}-W_k+2P}{S} \right\rfloor + 1$$

Pooling Layers

Aggregate neighbouring values into a single output by a specific hand-crafted function. The pooling kernel is applied **channel-wise** and with a **stride** ($s > 1$) to get a down-sampled output.



Hyper-parameters

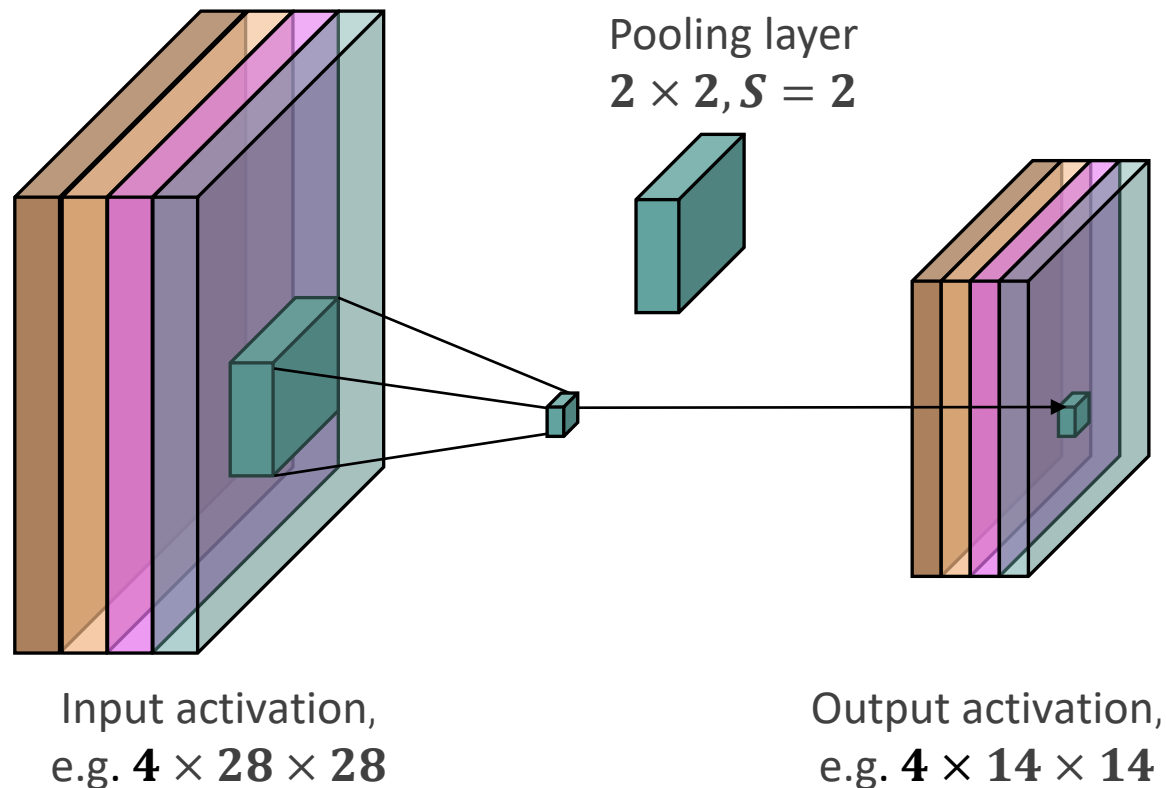
Kernel width and height: $W_K \times H_K$

Pooling function: max, avg, ...

Stride S (> 1)

Pooling Layers

Aggregate neighbouring values into a single output by a specific hand-crafted function. The pooling kernel is applied **channel-wise** and with a **stride** ($s > 1$) to get a down-sampled output.



Hyper-parameters

Kernel width and height: $W_K \times H_K$

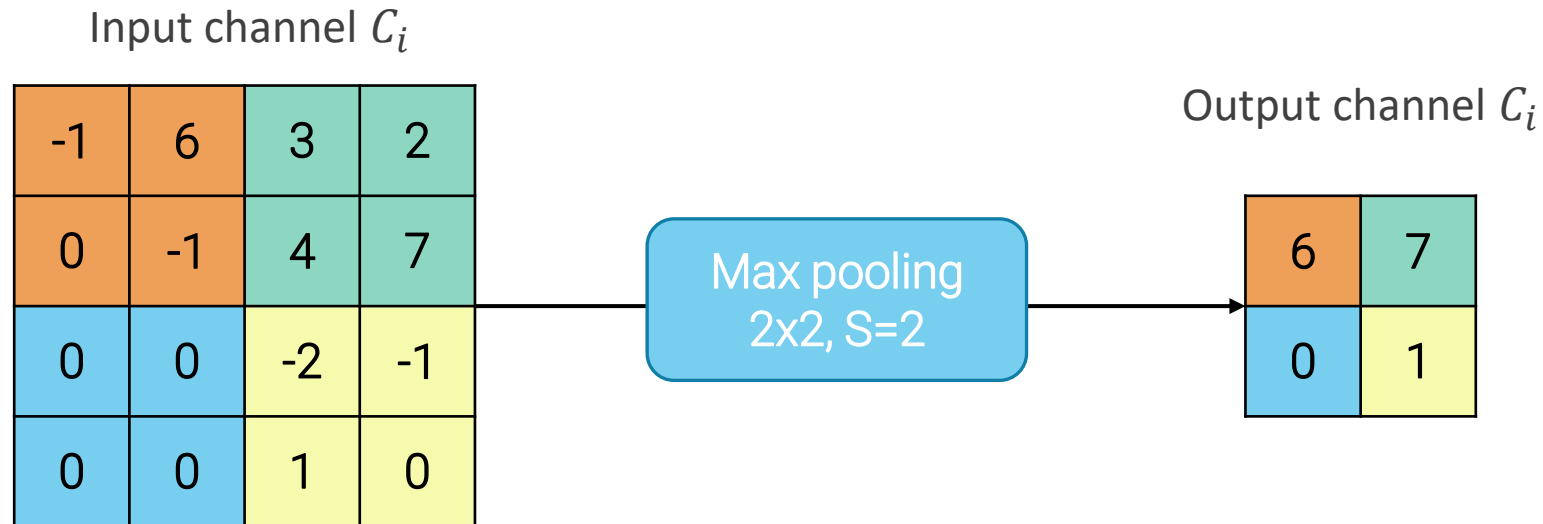
Pooling function: max, avg, ...

Stride S (> 1)

Most common choice:

$2 \times 2, S = 2, \text{max}$ (Max Pooling)

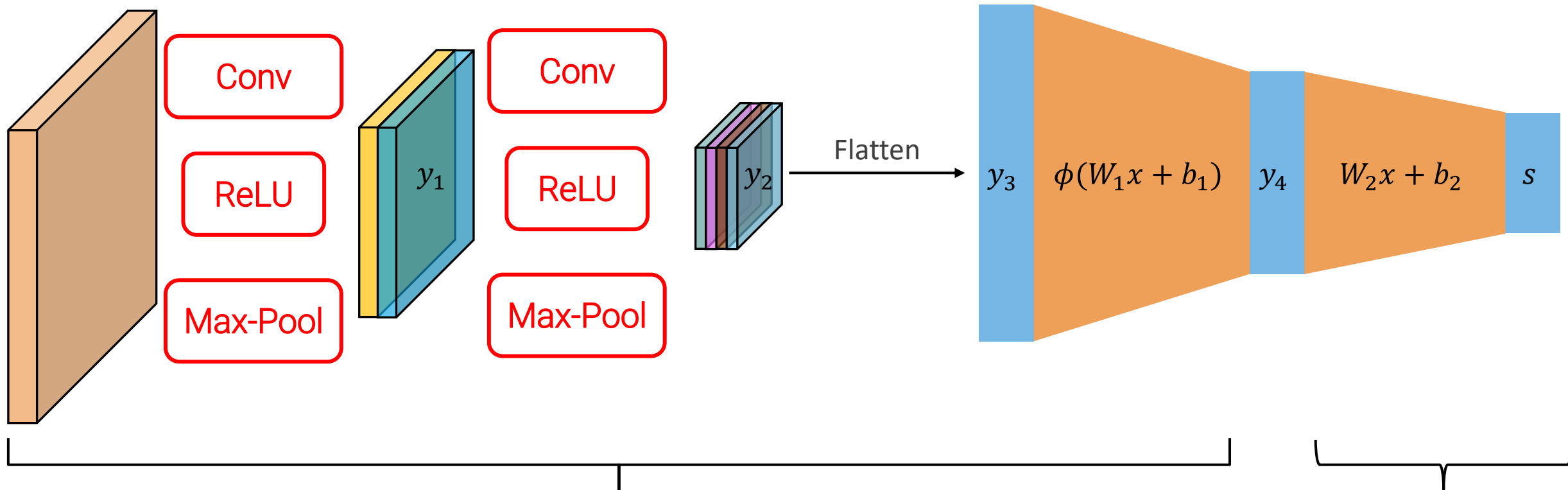
Max Pooling



Compared to strided convolutions, max pooling

- has no learnable parameters (pros and cons)
- provides **invariance** to small spatial shifts.

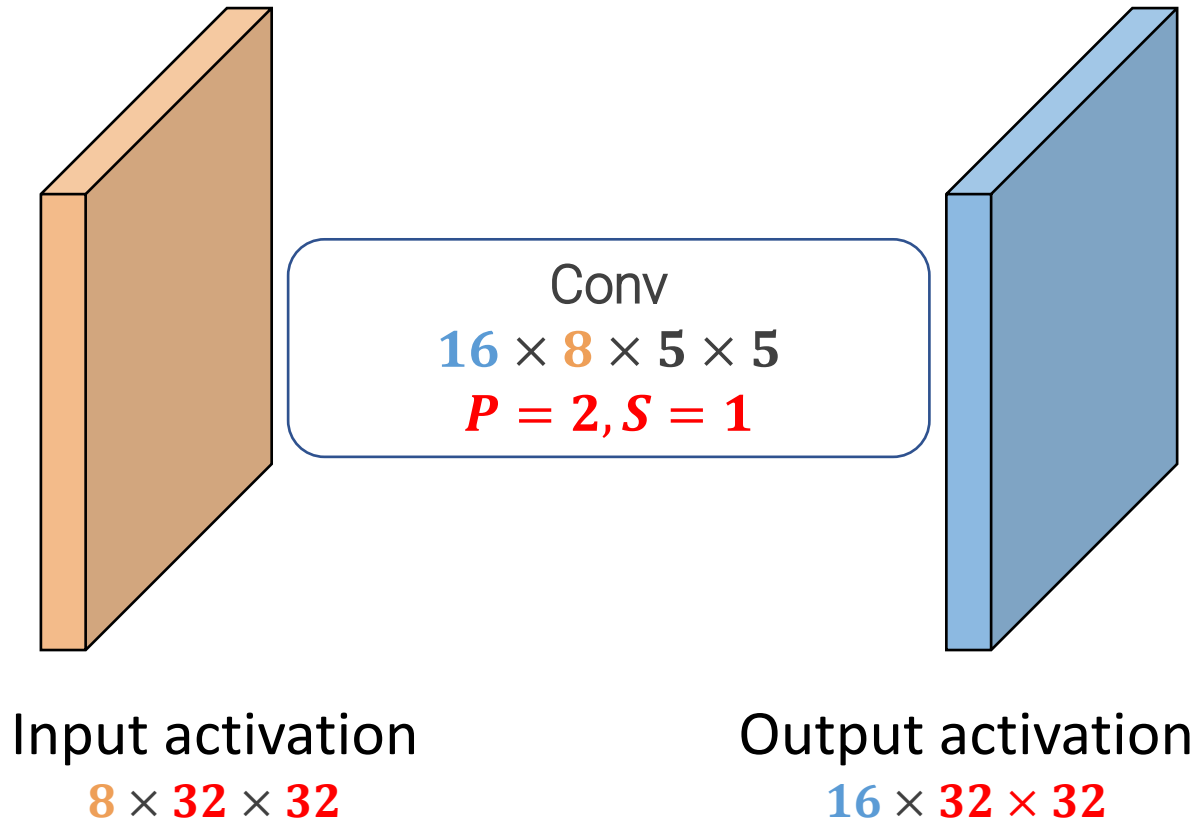
Convolutional Neural Networks (CNNs – convnets)



N **Conv+ReLU+Pooling** layers followed by M fully connected layers
This portion of the network is also called **feature extractor**.

The final fully connected layer is also called the **classifier**.

Number of parameters (N) and Flops (1)



- The number of learnable parameters, N , for the Conv layer is
 $16 \times (8 \times 5 \times 5 + 1) = 16 \times 201 = 3,216$

- Hence, there are \rightarrow *bias*

$$16 \times 32 \times 32 = 16,384$$

values in the output activation.

- To compute each output value we need to compute a MAC (Multiply And Accumulate) per each parameter, which amounts to $2N$ Flops (or N MACs). In our example, the total number of Flops is therefore:

$$16,384 \times 8 \times 5 \times 5 \times 2 \cong 6.5\text{M flops}$$

Number of parameters (N) and Flops (2)

$$N = C_{out} \times (C_{in} \times H_K \times W_K + 1)$$

Output Channels

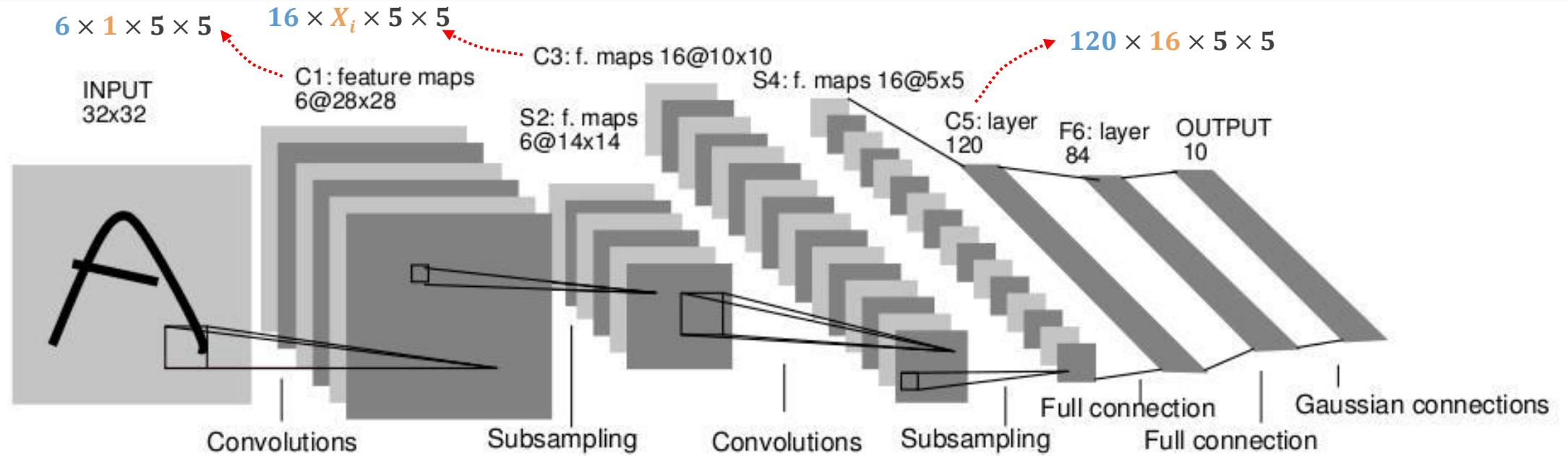
Input Channels

Kernel Size

$$\text{Flops} = 2 \times (C_{out} \times H \times W) \times (C_{in} \times H_K \times W_K)$$

Activation Size

LeNet-5



- Proposed to classify handwritten digits (MNIST dataset) and used in the US to read checks automatically.
- Alongside the layers, the spatial dimensions decrease and the number of channels increases. Normalization of inputs (zero mean and unit variance) to accelerate learning.
- 5x5 convolutional kernels, no padding, average pooling (with trainable scale and bias), tanh non-linearities.
- Sparse connection matrix in C3 (convs take input from a subset of input channels, as detailed in Tab. 1 of the paper).
- Two fully connected layers: F6, OUTPUT (10 RBF units: each unit compute the distance between its input vector and the corresponding parameter vector).

AlexNet (1)

- **Won the [ILSVRC 2012](#) bringing the Top-5 error from 25.8 to 16.4.**
- About 60M parameters, trained for 5-6 days on 2 GPUs.
- At training time, random-cropping of 224x224 patches (and their horizontal reflections) from the 256x256 RGB input images and *colour jittering* (**massive data augmentation**).
- At test time, averaging predictions (i.e. *softmax*) across 10 patches (central + 4 corner alongside their horizontal reflections).
- 8 layers with weights (5 Conv + 3 FC): most of the parameters are in the FC layers.
- All layers (Conv and FC) deploy ReLU non-linearities which yield faster training compared to saturating non-linearities (see Fig. 1 in the paper).
- First Conv layer has a stride of 4 ($S=4$): **stem layer** performing heavy reduction of the spatial size of activations, mainly to reduce memory and computation cost. In all other Conv layers $S=1$.
- Last FC layer has 1000 units (as many as the ILSVRC classes), the penultimate FC layer is the feature/representation layer and has a cardinality of 4096.

Layer	#Filters/ #Units	Filter Size	S	P	Activation Size
conv1	96	11x11	4	2	55x55
Pool1	1	3x3	2	0	27x27
conv2	256	5x5	1	2	27x27
pool2	1	3x3	2	0	13x13
conv3	384	3x3	1	1	13x13
conv4	384	3x3	1	1	13x13
conv5	256	3x3	1	1	13x13
pool3	1	3x3	2	0	6x6
flatten	0	0	0	0	1x1
fc6	4096	-	-	-	1x1
fc7	4096	-	-	-	1x1
fc8	1000	-	-	-	1x1

AlexNet (2)

Layer	#Filters/ #Units	Filter Size	S	P	Activation Size
conv1	96	11x11	4	2	55x55
Pool1	1	3x3	2	0	27x27
conv2	256	5x5	1	2	27x27
pool2	1	3x3	2	0	13x13
conv3	384	3x3	1	1	13x13
conv4	384	3x3	1	1	13x13
conv5	256	3x3	1	1	13x13
pool3	1	3x3	2	0	6x6
flatten	0	0	0	0	1x1
fc6	4096	-	-	-	1x1
fc7	4096	-	-	-	1x1
fc8	1000	-	-	-	1x1

Totals: #params \cong 61M, GFlops \cong 2.3

overlapping (max) pooling

$$\text{---} \rightarrow 96 \times 3 \times 11 \times 11 \rightarrow \text{\#params} \cong 35\text{K}, \text{MFlops} \cong 211$$

$$\text{---} \rightarrow 256 \times 96 \times 5 \times 5 \rightarrow \text{\#params} \cong 615\text{K}, \text{MFlops} \cong 896$$

$$\text{---} \rightarrow 384 \times 256 \times 3 \times 3 \rightarrow \text{\#params} \cong 885\text{K}, \text{MFlops} \cong 299$$

$$\text{---} \rightarrow 384 \times 384 \times 3 \times 3 \rightarrow \text{\#params} \cong 1.2\text{M}, \text{MFlops} \cong 448$$

$$\text{---} \rightarrow 256 \times 384 \times 3 \times 3 \rightarrow \text{\#params} \cong 885\text{K}, \text{MFlops} \cong 299$$

$$\text{---} \rightarrow 4096 \times (6 \times 6 \times 256) \rightarrow \text{\#params} \cong 37.5\text{M}, \text{MFlops} \cong 75$$

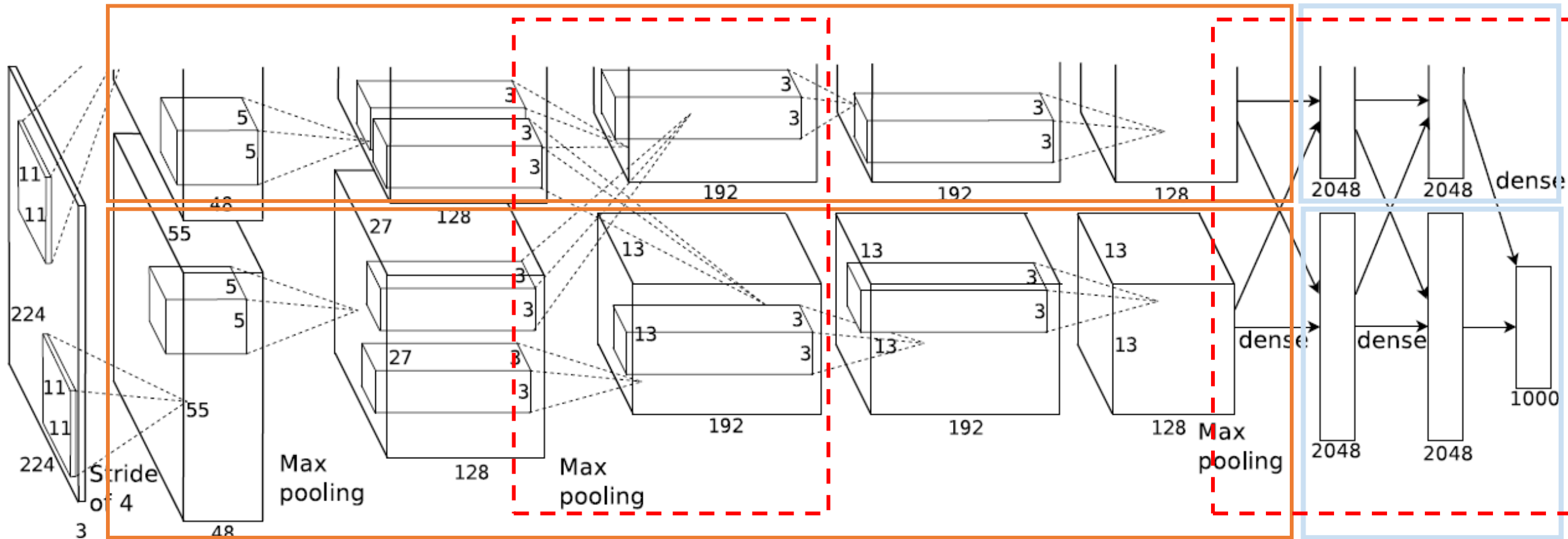
$$\text{---} \rightarrow 4096 \times (4096) \rightarrow \text{\#params} \cong 16.7\text{M}, \text{MFlops} \cong 33$$

$$\text{---} \rightarrow 1000 \times (4096) \rightarrow \text{\#params} \cong 4\text{M}, \text{MFlops} \cong 8$$

- **Local Response Normalization** (after conv1 and conv2): activations are normalized by the sum of those at the same spatial position in a few ($n=5$) *adjacent* channels (mimics [lateral inhibition](#) in real neurons).
- **Dropout** (fc6,fc7): at training time the output of each unit is set to zero with probability 0.5. This forces units to learn more robust features since none of them can rely on the presence of particular other ones (**Regularization**)

AlexNet (3)

In the original implementation the computational load was split between two GPUs



The red boxes with dashed lines highlight layers that take input from both GPUs

AlexNet – Training Recipe

Hyperparameter	Values
Optimizer	SGD with B=128
Epochs	90
Learning Rate	0.01, divided 3 times by 10 when the validation error stopped to improve
Weight Decay	0.0005
Momentum	0.9
Data Augmentation	Colour Jittering, Random Crop, Horizontal Flip
Initialization	Weights $\sim N(0, 0.01)$, Biases: 1 (conv2,conv4,conv5, fc6,fc7,fc8) or 0 (conv1,conv3)
Normalization	Centering (subtraction of the mean RGB colour in the training set)

- **Second place in ILSVRC 2014 (Top-5 error: 7.5 %)**
- Explores the benefits of deep and regular architectures based on a few simple design choices:
 - 3x3 conv layers with S=1, P=1
 - 2x2 max-pooling, S=2, P=0
 - #Filters (#channels) double after every pool
- The architecture is designed as a repetition of **stages**: a chain of layers that *process activations at the same spatial resolution (conv-conv-pool, conv-conv-conv-pool and conv-conv-conv-conv-pool)*.
 - A **stage** has the **same receptive field** as a single larger convolution but, given the same number of input/output channels, introduces **more non-linearities** and requires **less parameters** and **less computation**. A stage requires **more memory** to store the activations, though.

VGG-16 VGG-19

conv3-64	conv3-64
conv3-64	conv3-64
maxpool	maxpool
conv3-128	conv3-128
conv3-128	conv3-128
maxpool	maxpool
conv3-256	conv3-256
conv3-256	conv3-256
conv3-256	conv3-256
maxpool	maxpool
conv3-512	conv3-512
conv3-512	conv3-512
conv3-512	conv3-512
maxpool	maxpool
conv3-512	conv3-512
conv3-512	conv3-512
conv3-512	conv3-512
maxpool	maxpool
FC-4096	FC-4096
FC-4096	FC-4096
FC-1000	FC-1000

- **Second place in ILSVRC 2014 (Top-5 error: 7.5 %)**
- Explores the benefits of deep and regular architectures based on a few simple design choices:
 - 3x3 conv layers with $S=1$, $P=1$
 - 2x2 max-pooling, $S=2$, $P=0$
 - #Filters (#channels) double after every pool
- The architecture is designed as a repetition of **stages**: a chain of layers that *process activations at the same spatial resolution (conv-conv-pool, conv-conv-conv-pool and conv-conv-conv-conv-pool)*.
 - A **stage** has the **same receptive field** as a single larger convolution but, given the same number of input/output channels, introduces **more non-linearities** and requires **less parameters** and **less computation**. A stage requires **more memory** to store the activations, though.
 - For example, a single $C \times C \times 5 \times 5$ conv layer:

$$\text{\#params} = C \times (C \times 5 \times 5 + 1) = 25 \times C^2 + C$$

$$\text{\#Flops} = (C \times W \times H) \times C \times 5 \times 5 \times 2 = 50 \times C^2 \times W \times H$$

$$\text{\#activations} = C \times W \times H$$
 - while a stage consisting of 2 staked $C \times C \times 3 \times 3$ conv layers (same receptive field):

$$\text{\#params} = 2 \times C \times (C \times 3 \times 3 + 1) = 18 \times C^2 + 2C$$

$$\text{\#Flops} = 2 \times (C \times W \times H) \times C \times 3 \times 3 \times 2 = 36 \times C^2 \times W \times H$$

$$\text{\#activations} = 2 \times C \times W \times H$$

VGG-16 VGG-19

conv3-64	conv3-64
conv3-64	conv3-64
maxpool	maxpool
conv3-128	conv3-128
conv3-128	conv3-128
maxpool	maxpool
conv3-256	conv3-256
conv3-256	conv3-256
conv3-256	conv3-256
conv3-256	conv3-256
maxpool	maxpool
conv3-512	conv3-512
conv3-512	conv3-512
conv3-512	conv3-512
conv3-512	conv3-512
maxpool	maxpool
conv3-512	conv3-512
conv3-512	conv3-512
conv3-512	conv3-512
conv3-512	conv3-512
maxpool	maxpool
FC-4096	FC-4096
FC-4096	FC-4096
FC-1000	FC-1000

VGG-16

							Totals: #params \cong 138M, GFlops \cong 39	
	Layer	#Filters/ #Units	Filter Size	S	P	Activation Size		
S1	Conv1	64	3x3	1	1	224x224	\rightarrow	$64 \times 3 \times 3 \times 3 \rightarrow$ #params \cong 1.8K, MFlops \cong 173
	Conv2	64	3x3	1	1	224x224	\rightarrow	$64 \times 64 \times 3 \times 3 \rightarrow$ #params \cong 37K, GFlops \cong 3.7
	Pool1	1	2x2	2	0	112x112		
S2	Conv3	128	3x3	1	1	112x112	\rightarrow	$128 \times 64 \times 3 \times 3 \rightarrow$ #params \cong 74K, GFlops \cong 1.85
	Conv4	128	3x3	1	1	112x112	\rightarrow	$128 \times 128 \times 3 \times 3 \rightarrow$ #params \cong 147.6K, GFlops \cong 3.7
	Pool2	1	2x2	2	0	56x56		
S3	Conv5	256	3x3	1	1	56x56	\rightarrow	$256 \times 128 \times 3 \times 3 \rightarrow$ #params \cong 295K, GFlops \cong 1.85
	Conv6	256	3x3	1	1	56x56	\rightarrow	$2 \times 256 \times 256 \times 3 \times 3 \rightarrow$ #params \cong 1.19M, GFlops \cong 7.4
	Conv7	256	3x3	1	1	56x56		
	Pool3	1	2x2	2	0	28x28		
S4	Conv8	512	3x3	1	1	28x28	\rightarrow	$512 \times 256 \times 3 \times 3 \rightarrow$ #params \cong 1.18M, GFlops \cong 1.85
	Conv9	512	3x3	1	1	28x28	\rightarrow	$2 \times 512 \times 512 \times 3 \times 3 \rightarrow$ #params \cong 4.7M, GFlops \cong 7.4
	Conv10	512	3x3	1	1	28x28		
	Pool4	1	2x2	2	0	14x14		
S5	Conv11	512	3x3	1	1	14x14	\rightarrow	$3 \times 512 \times 512 \times 3 \times 3 \rightarrow$ #params \cong 7M, GFlops \cong 11.1
	Conv12	512	3x3	1	1	14x14		
	Conv13	512	3x3	1	1	14x14		
	Pool5	1	2x2	2	0	7x7		
Dropout	Flatten	0	0	0	0	1x1		
	Fc14	4096	-	-	-	1x1	\rightarrow	$4096 \times (7 \times 7 \times 2508) \rightarrow$ #params \cong 102.7M, MFlops \cong 205
	Fc15	4096	-	-	-	1x1	\rightarrow	$4096 \times (4096) \rightarrow$ #params \cong 16.7M, MFlops \cong 33
	fc16	1000	-	-	-	1x1	\rightarrow	$1000 \times (4096) \rightarrow$ #params \cong 4M, MFlops \cong 8

VGG – Training Recipe



Hyperparameter	Values
Optimizer, Learning Rate, Weight Decay, Normalization, Epochs	Same as AlexNet, but for B=256 and Epochs=74.
Initialization*	Deep nets are hard to train with randomly initialized weights due to <u>instability</u> of gradients. They train a VGG-11 with Weights $\sim N(0, 0.01)$, Biases=0. Then train VGG-16 and VGG-19 by initializing the first 4 conv layers and the last 3 fc layers with the pre-trained weights of the corresponding layers of VGG-11.
Data Augmentation	Same as AlexNet plus Scale Jittering (randomly rescale the input image to $S \times S$, with S in $[256, 512]$)

* After the paper submission the authors found that they could train the deeper VGGs without pre-training by the so called Glorot (Xavier) initialization procedure:

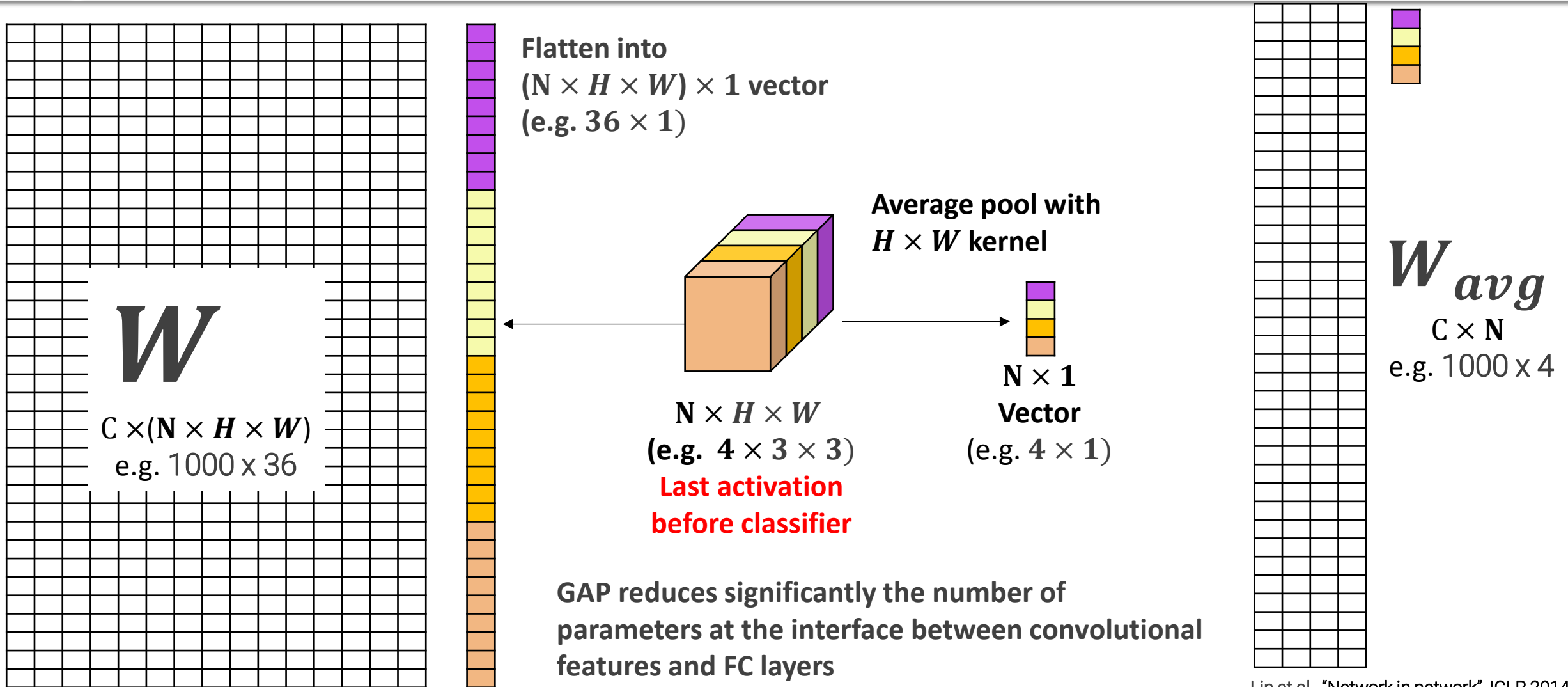
Glorot, X. and Bengio, Y. “Understanding the difficulty of training deep feedforward neural networks”. In Proc. AISTATS, 2010.

which was later extended by He to work better with ReLu (He initialization)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun “Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification”, ICCV 2015

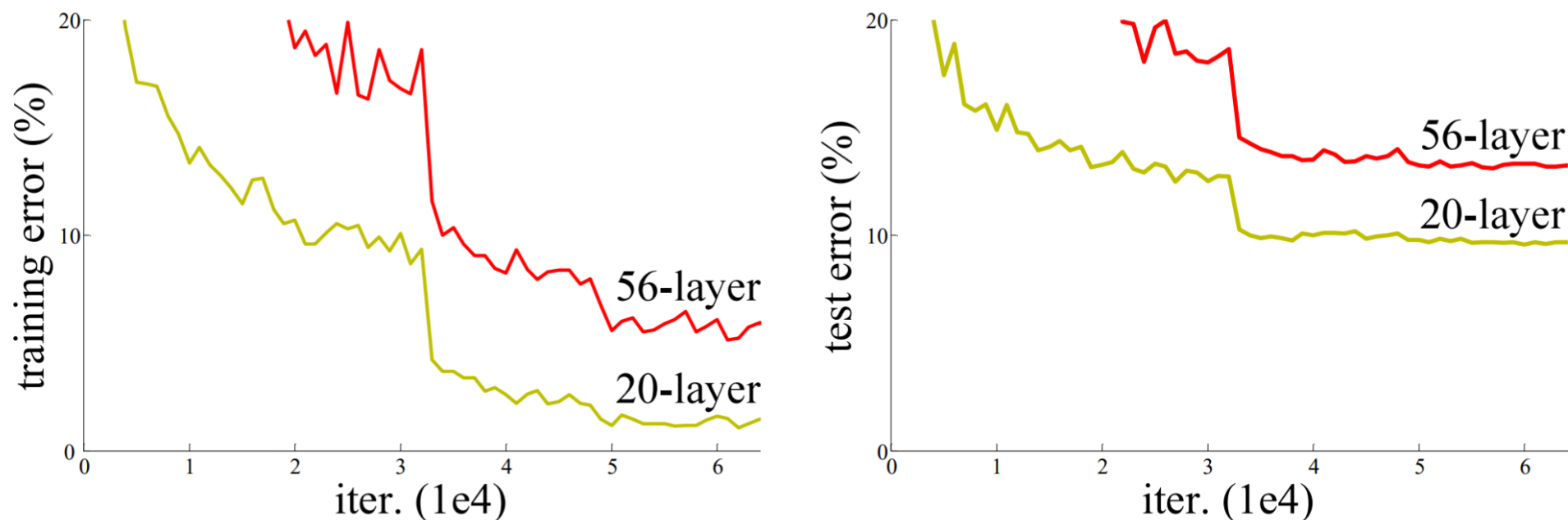
[Weight Initialization Schemes - Xavier \(Glorot\) and He | Mustafa Murat ARAT \(mmuratarat.github.io\)](https://github.com/muratarat/vgg16)

Global Average Pooling (GAP)



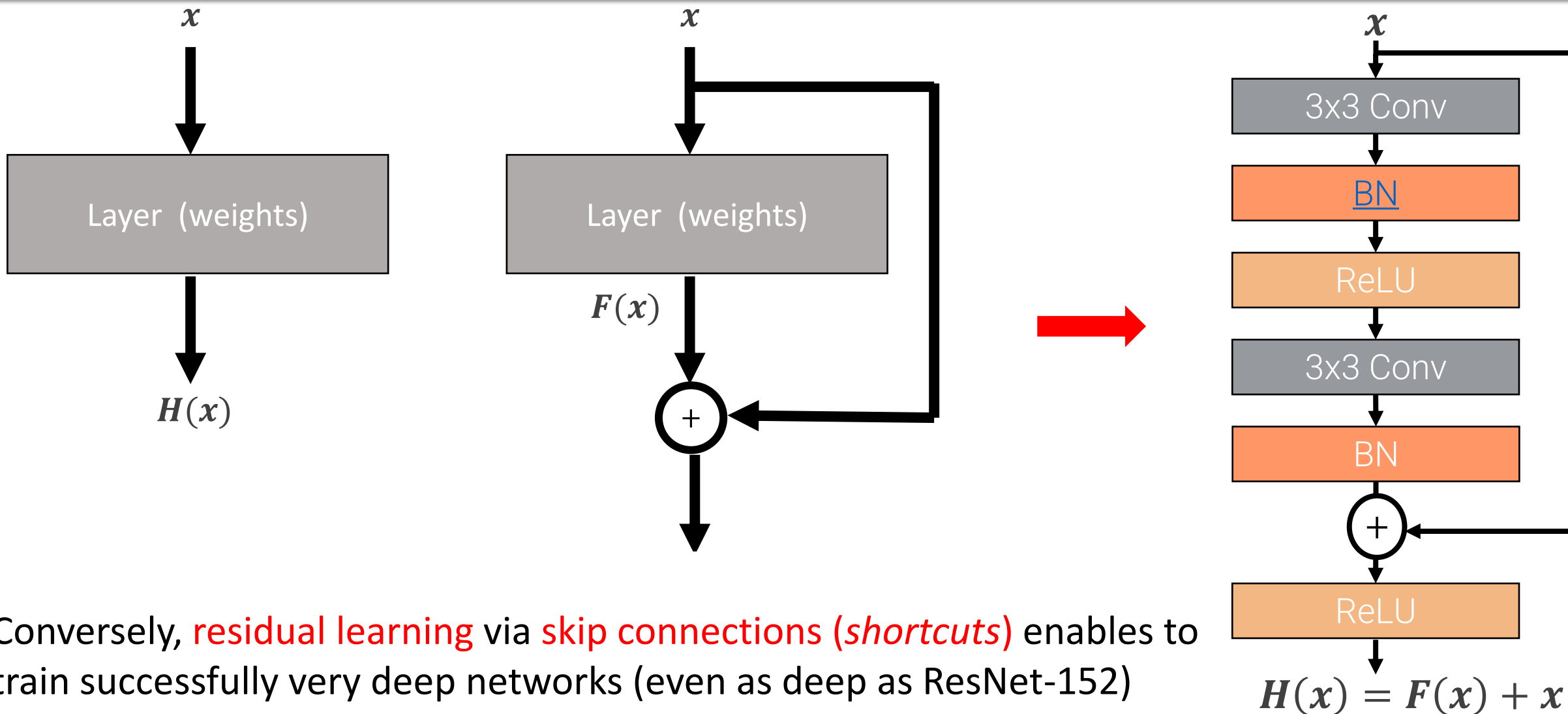
Residual Networks - Motivation

The success of the VGG design would suggest to increase depth to improve performance, but..

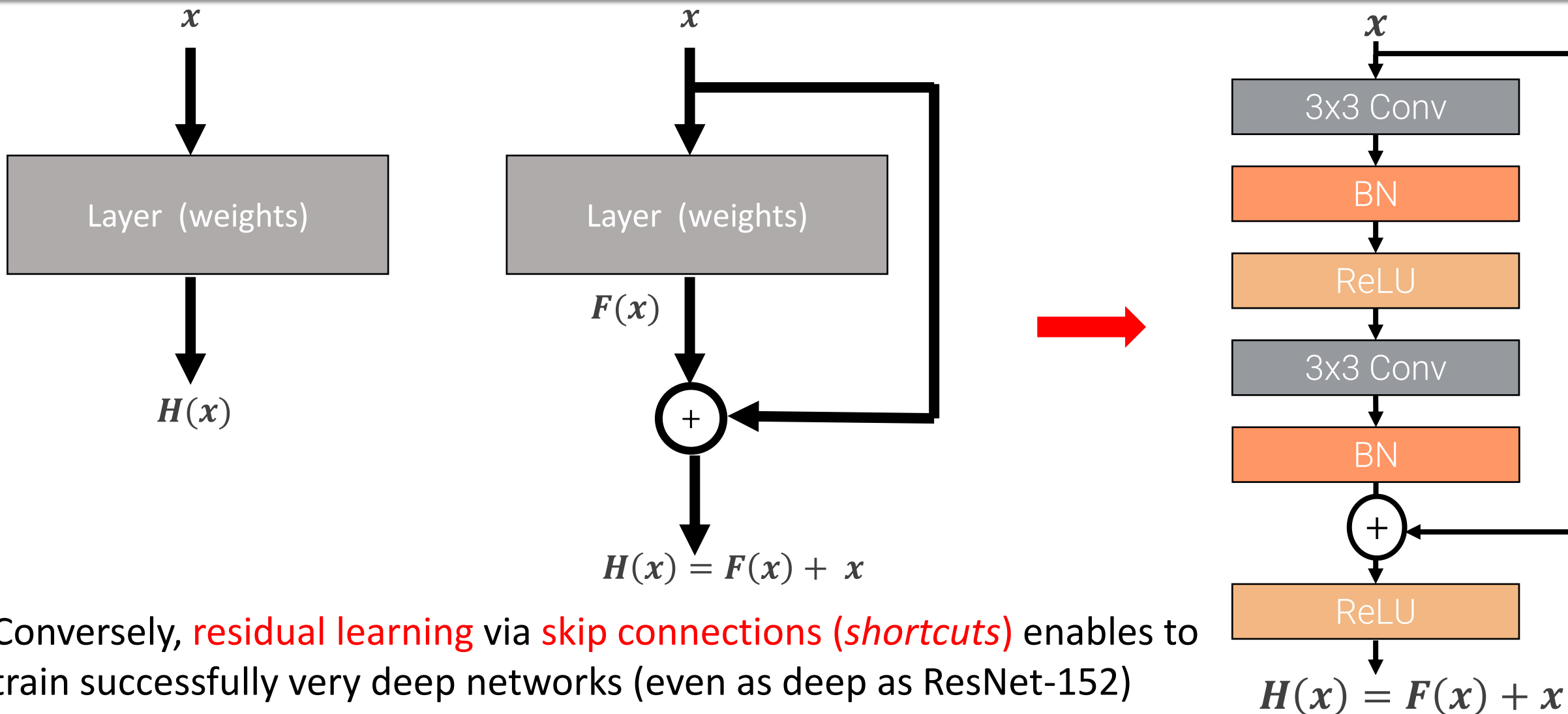


The problem is not (only) overfitting, the training error is larger for the deeper network !
Training very deep networks turns out to be inherently hard !

Residual Learning and Residual Blocks



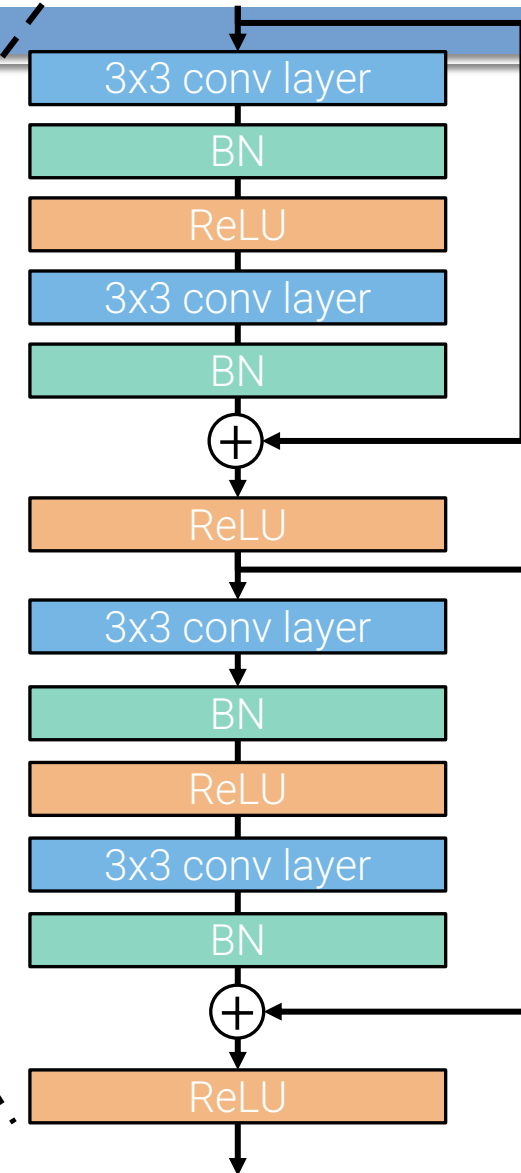
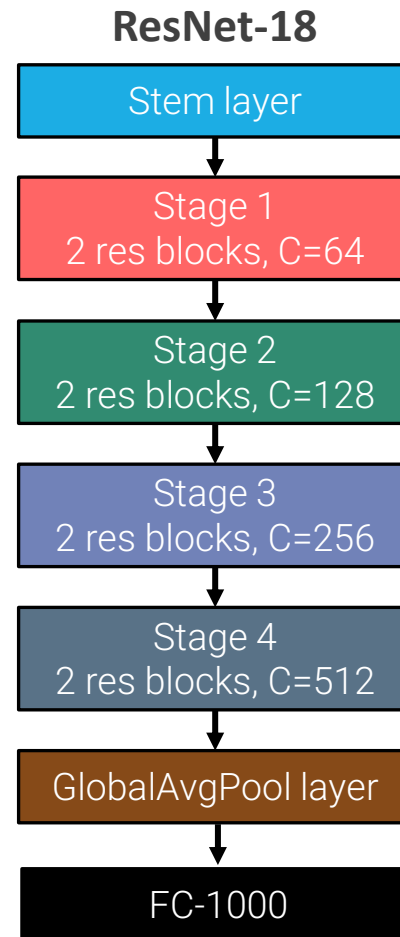
Residual Learning and Residual Blocks



Conversely, **residual learning** via **skip connections (shortcuts)** enables to train successfully very deep networks (even as deep as ResNet-152)

Residual Networks – Architecture

- Inspired by VGG: a few simple design choices, repetition of stages.
- Stages are stacks of Residual Blocks (RB). Each RB includes two 3x3 conv layers.
- The first RB in *most* stages halves the spatial resolution ($S=2$) and doubles the number of channels.
- Initial Stem layer ($S=2$ conv followed by $S=2$ max pool) to quickly down-sample the input image and Global Average Pooling to efficiently interface the final 1000-ways FC layer.
- Naming notation similar to VGG: ResNet-X denotes a ResNet architecture having X layers with learnable parameters.



ResNet18

						Totals: #params \cong 11.5M, GFlops \cong 3.6	
Layer						#Filters/ #Units	Filter Size
S						P	Activation Size
S1	Conv1	64	7x7	2	3	112x112	
	Pool1	1	3x3	2	1	56x56	
	Conv2_1	64	3x3	1	1	56x56	
	Conv2_2	64	3x3	1	1	56x56	
	Conv2_3	64	3x3	1	1	56x56	
	Conv2_4	64	3x3	1	1	56x56	
S2	Conv3_1	128	3x3	2	1	28x28	
	Conv3_2	128	3x3	1	1	28x28	
	Conv3_3	128	3x3	1	1	28x28	
	Conv3_4	128	3x3	1	1	28x28	
S3	Conv4_1	256	3x3	2	1	14x14	
	Conv4_2	256	3x3	1	1	14x14	
	Conv4_3	256	3x3	1	1	14x14	
	Conv4_4	256	3x3	1	1	14x14	
S4	Conv5_1	512	3x3	2	1	7x7	
	Conv5_2	512	3x3	1	1	7x7	
	Conv5_3	512	3x3	1	1	7x7	
	Conv5_4	512	3x3	1	1	7x7	
AvgPool		1	7x7	1	0	1x1	
fc1		1000	-	-	-	1x1	

$64 \times 3 \times 7 \times 7 \rightarrow \text{\#params} \cong 9.5\text{K}, \text{MFlops} \cong 236$

$4 \times 64 \times 64 \times 3 \times 3 \rightarrow \text{\#params} \cong 148 \text{ K}, \text{GFlops} \cong 0.9$

$128 \times 64 \times 3 \times 3 \rightarrow \text{\#params} \cong 74 \text{ K}, \text{MFlops} \cong 115$

$3 \times 128 \times 128 \times 3 \times 3 \rightarrow \text{\#params} \cong 443 \text{ K}, \text{GFlops} \cong 0.7$

$256 \times 128 \times 3 \times 3 \rightarrow \text{\#params} \cong 295 \text{ K}, \text{MFlops} \cong 115$

$3 \times 256 \times 256 \times 3 \times 3 \rightarrow \text{\#params} \cong 1.8 \text{ M}, \text{GFlops} \cong 0.7$

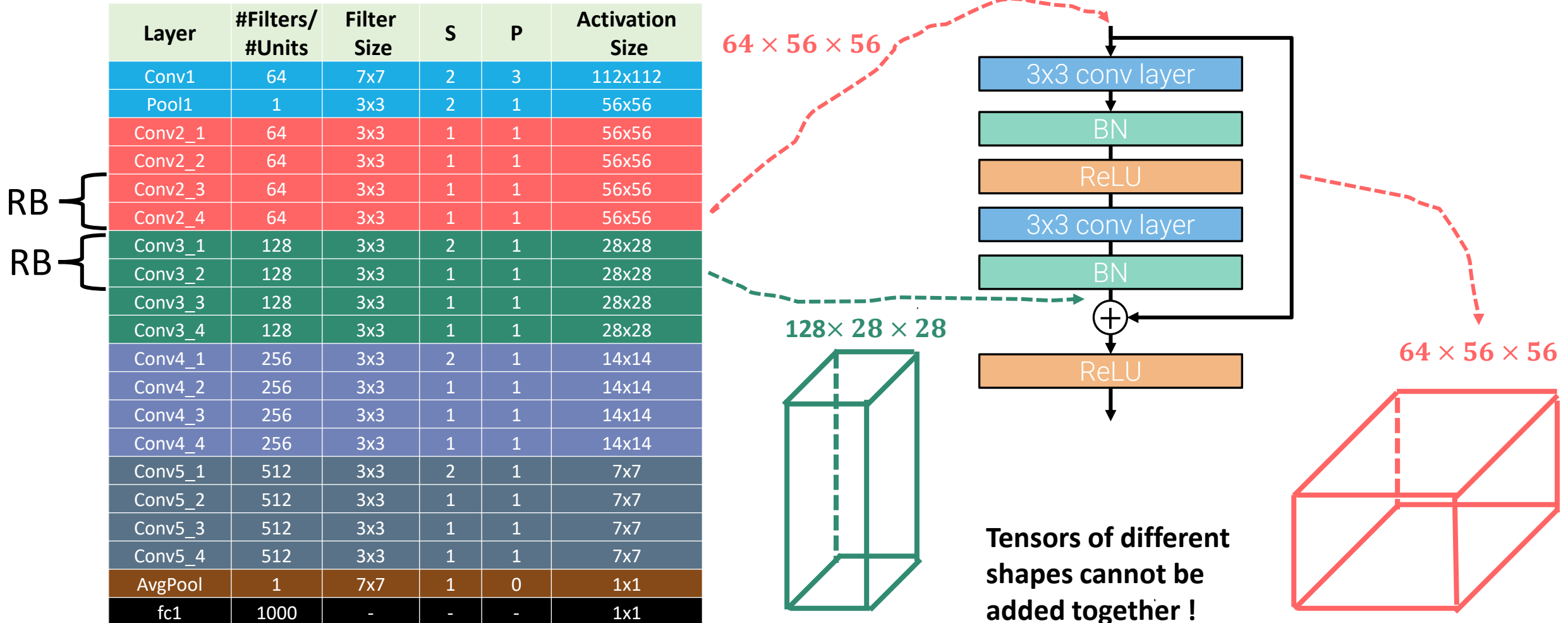
$512 \times 256 \times 3 \times 3 \rightarrow \text{\#params} \cong 1.18 \text{ M}, \text{MFlops} \cong 115$

$3 \times 512 \times 512 \times 3 \times 3 \rightarrow \text{\#params} \cong 7 \text{ M}, \text{GFlops} \cong 0.7$

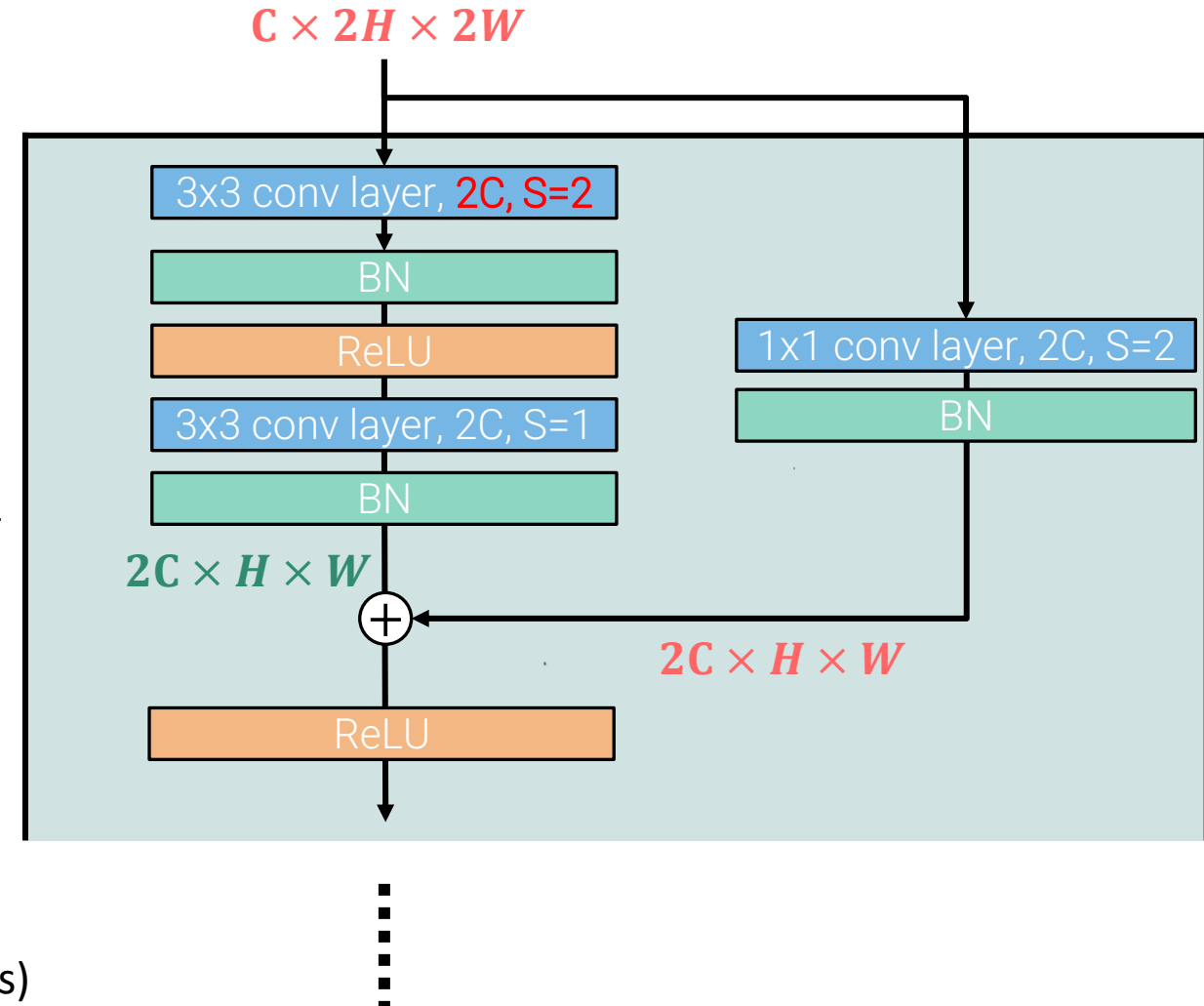
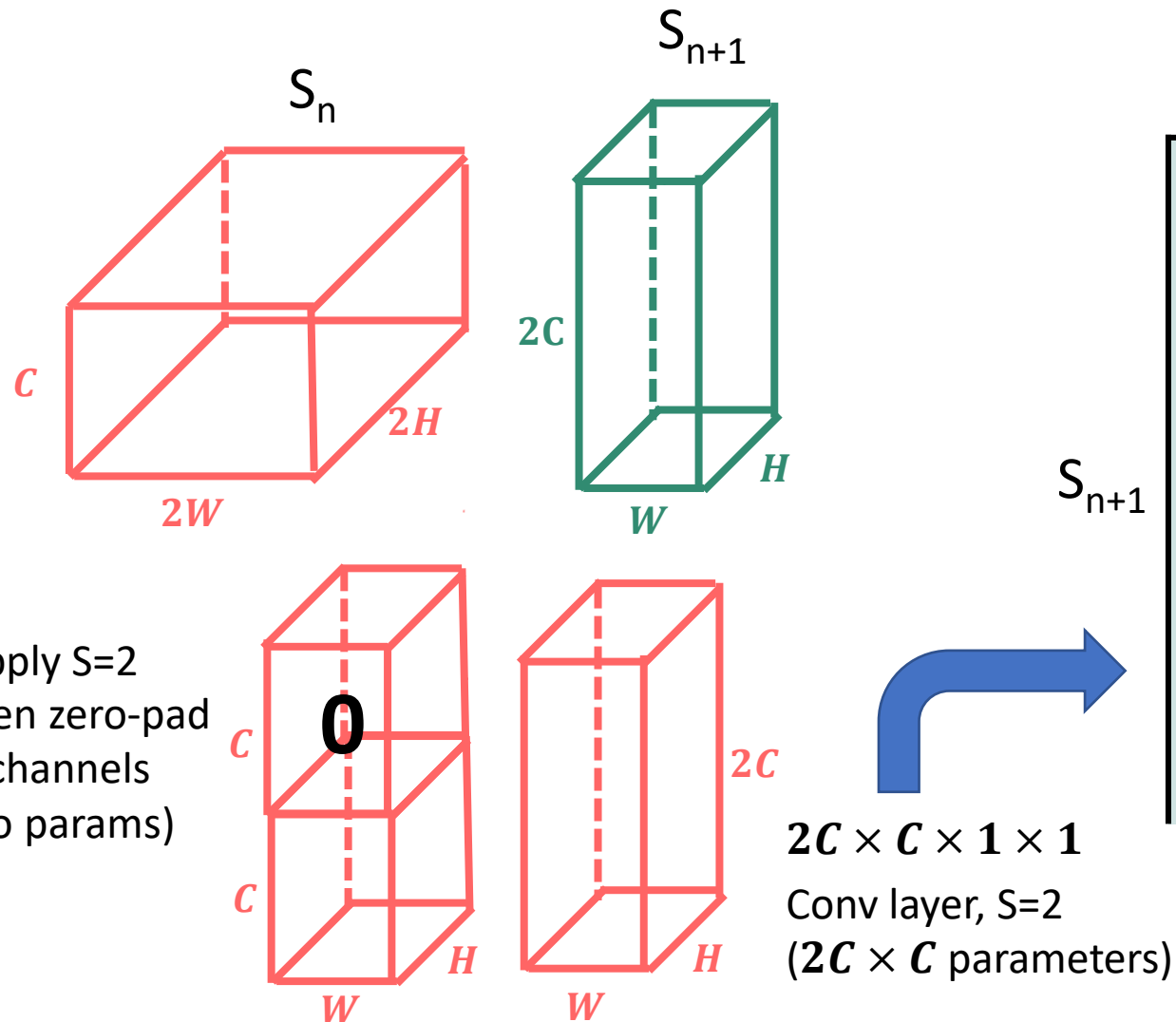
$1000 \times 512 \rightarrow \text{\#params} \cong 0.5 \text{ M}, \text{MFlops} \cong 1$

About one order of magnitude less parameters and Flops than VGG networks of comparable depths

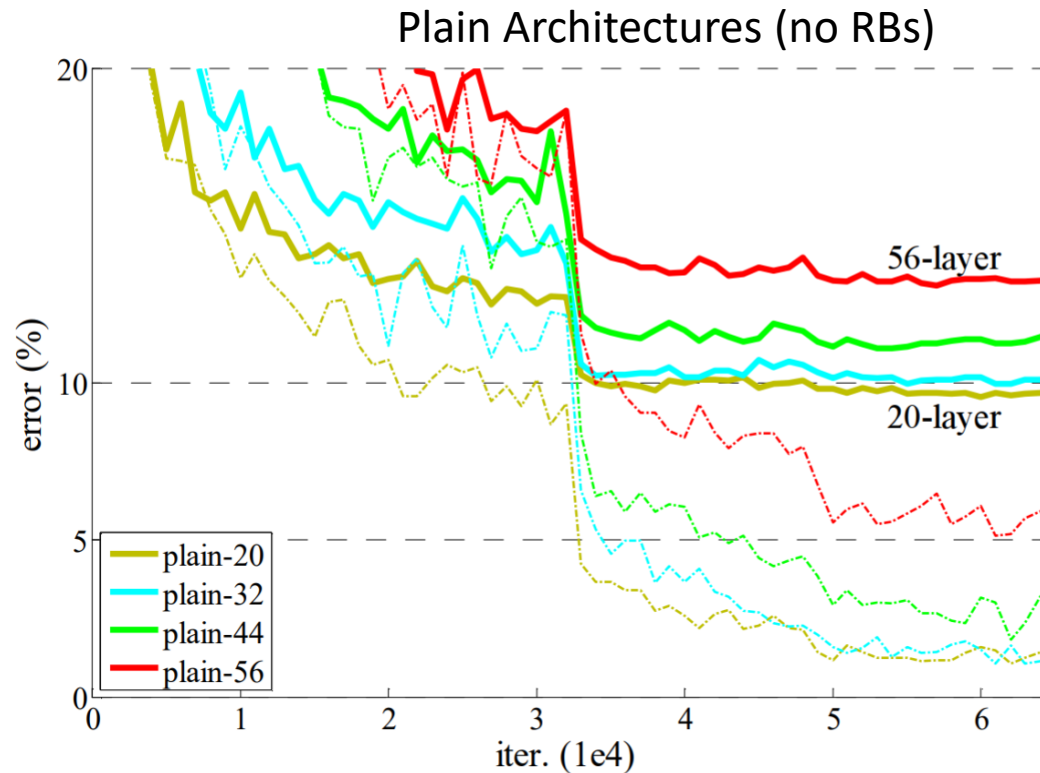
First RB in a stage & shape of skip connections



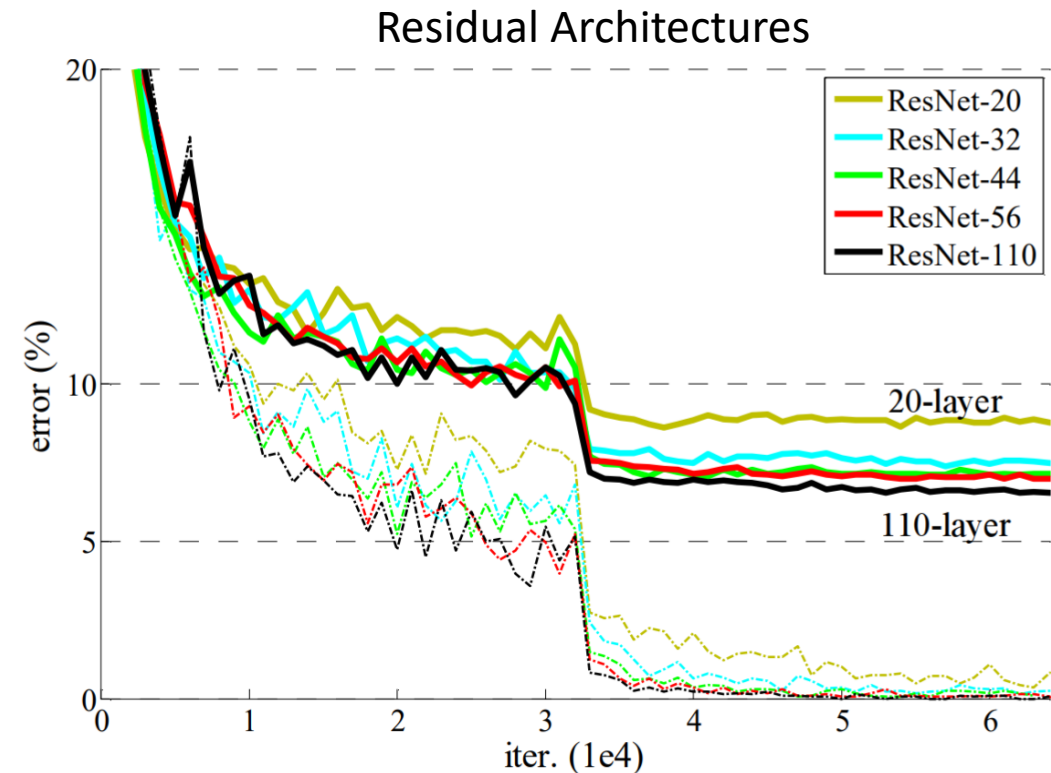
Modified first RB in a stage



Residual vs Plain Architectures



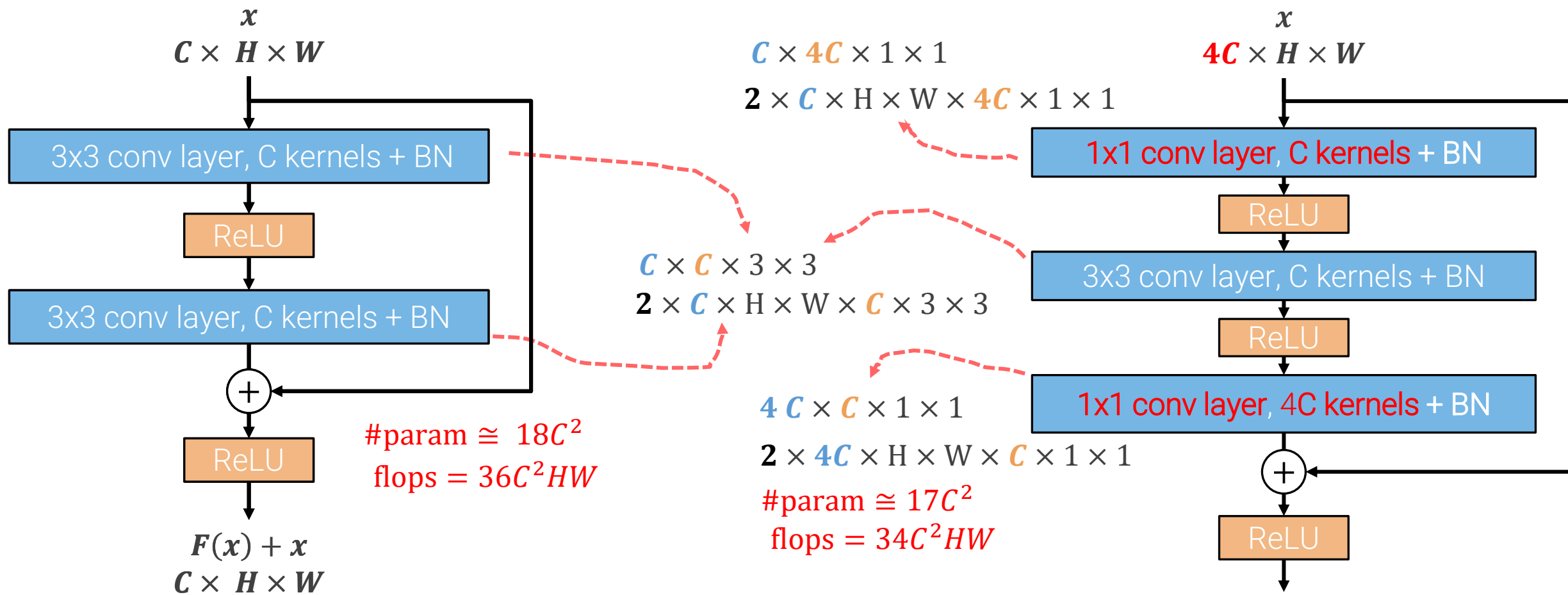
The shallower the better.....



....the deeper the better !

Residual Blocks allow for training deep CNNs. When properly trained, deep architectures outperform shallower ones. In 2015, ResNets won all the main computer vision competitions by large margins. In ILSVRC, an ensemble of ResNets (including two models with 152 layers) brought the Top-5 error from 6.7 to 3.6.

Bottleneck Residual Blocks for deeper ResNets

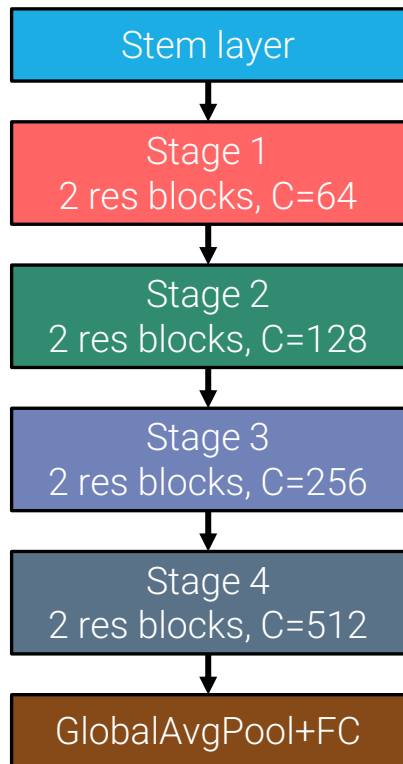


Bottleneck Residual Blocks (BRB) realize a cheaper design favoured by the authors when training deeper residual networks (i.e. with 50, 101 and 152 trainable layers).

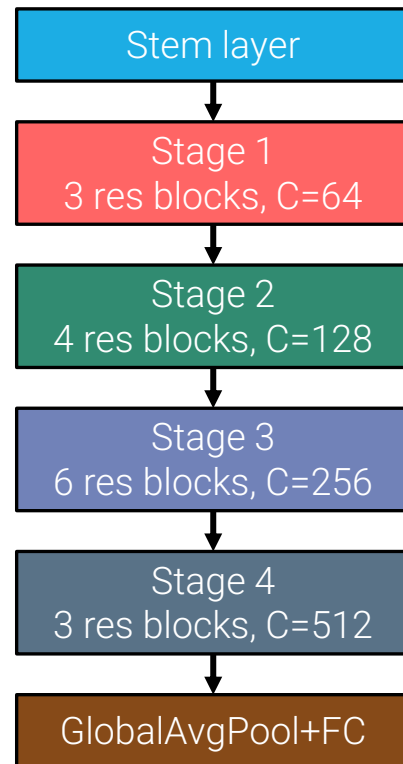
Main ResNet Architectures

Common choice in many settings

ResNet-18

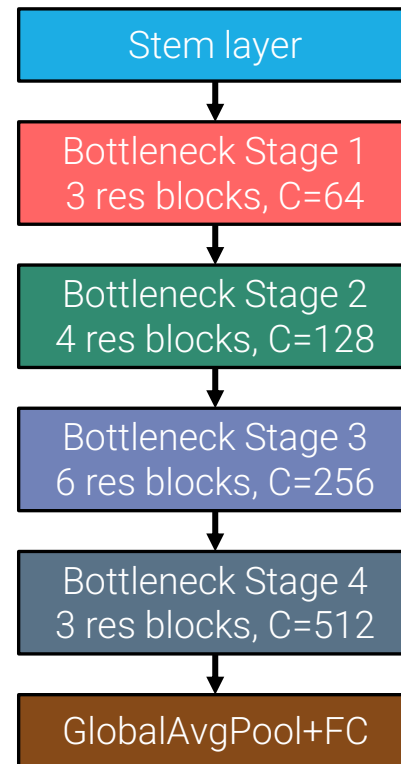


ResNet-34



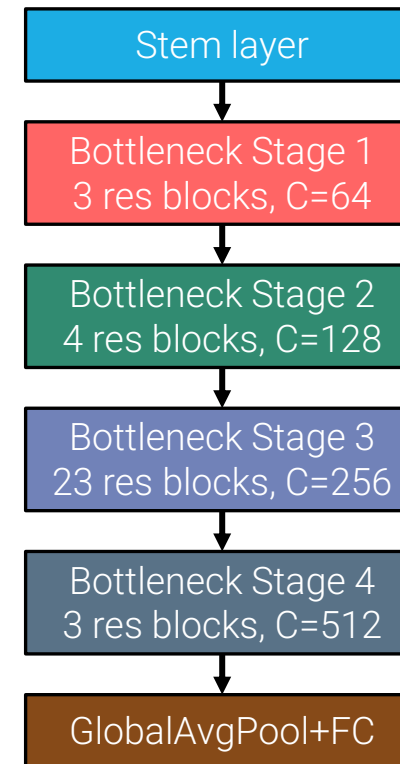
~7.3 Gflops
~22M params
7.46 top-5 error

ResNet-50



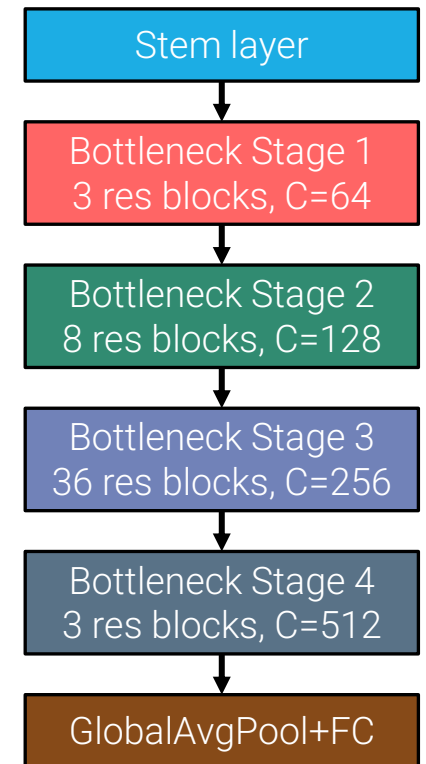
~7 Gflops
~23M params
6.71 top-5 error

ResNet-101



~14.5 Gflops
~42M params
6.05 top-5 error

ResNet-152



~22.7 Gflops
~60M params
5.71 top-5 error

Errors are on the validation set with 10 crops

ResNets (from the paper)

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112	7×7, 64, stride 2				
conv2_x	56×56	3×3 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	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1	average pool, 1000-d fc, softmax				
	FLOPs	1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9

MACs

Down-sampling is performed by conv3_1, conv4_1 and conv5_1.

The first BRB in a stage performs either no compression (conv2) or compression by 2 (conv3, conv4, conv5). All other BRB compress and expand by 4.

ResNet – Training Recipe

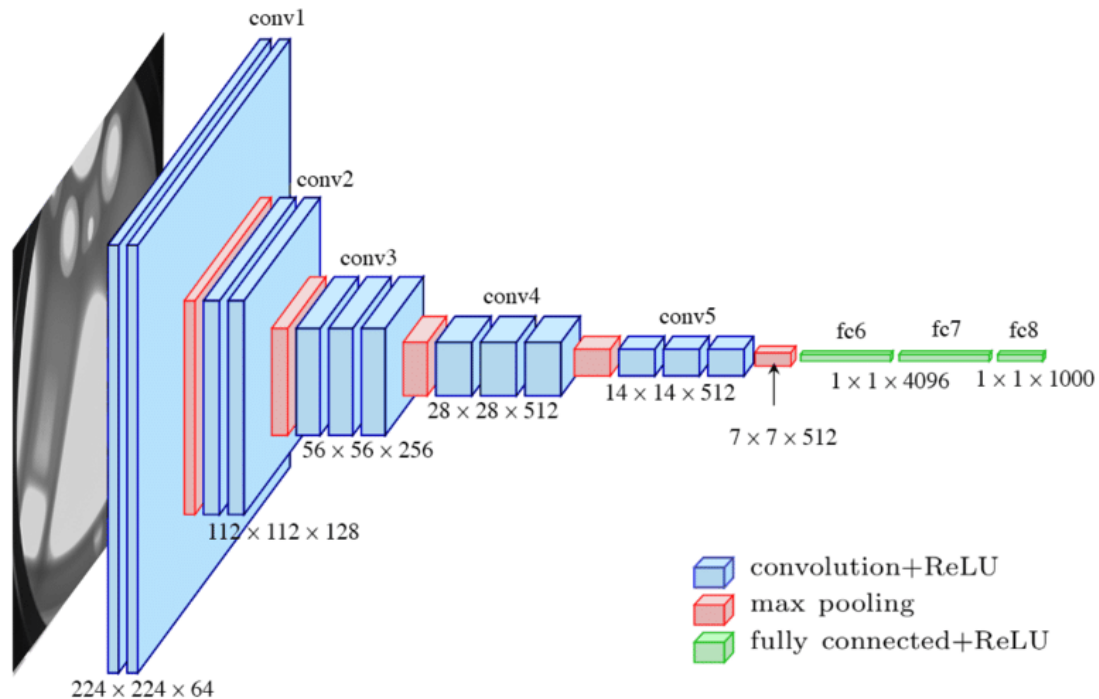
Hyperparameter	Values
Optimizer	SGD with B=256
Training Iterations	60×10^4
Learning Rate	0.1, divided by 10 when the validation error plateaus
Weight Decay	0.0001
Momentum	0.9
Data Augmentation, Normalization	Same as VGG
Initialization	He initialization

The use of a higher learning rate is motivated by the adoption of BN (Batch Normalization). Similarly, dropout is not deployed because experimental evidence suggests BN to act as a regularizer. This may also motivate the choice of a smaller weight decay parameter.

S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In ICML, 2015.

CNNs learn to detect hierarchies of features

The typical architecture of a CNN consists in using **small** (e.g. 3x3), **fixed-size filters** and **increasing the number of channels while down-sampling the activations**. In other words, as we move from shallower to deeper layers, activations get “*smaller*” and “*taller*”. **Why is that ?**

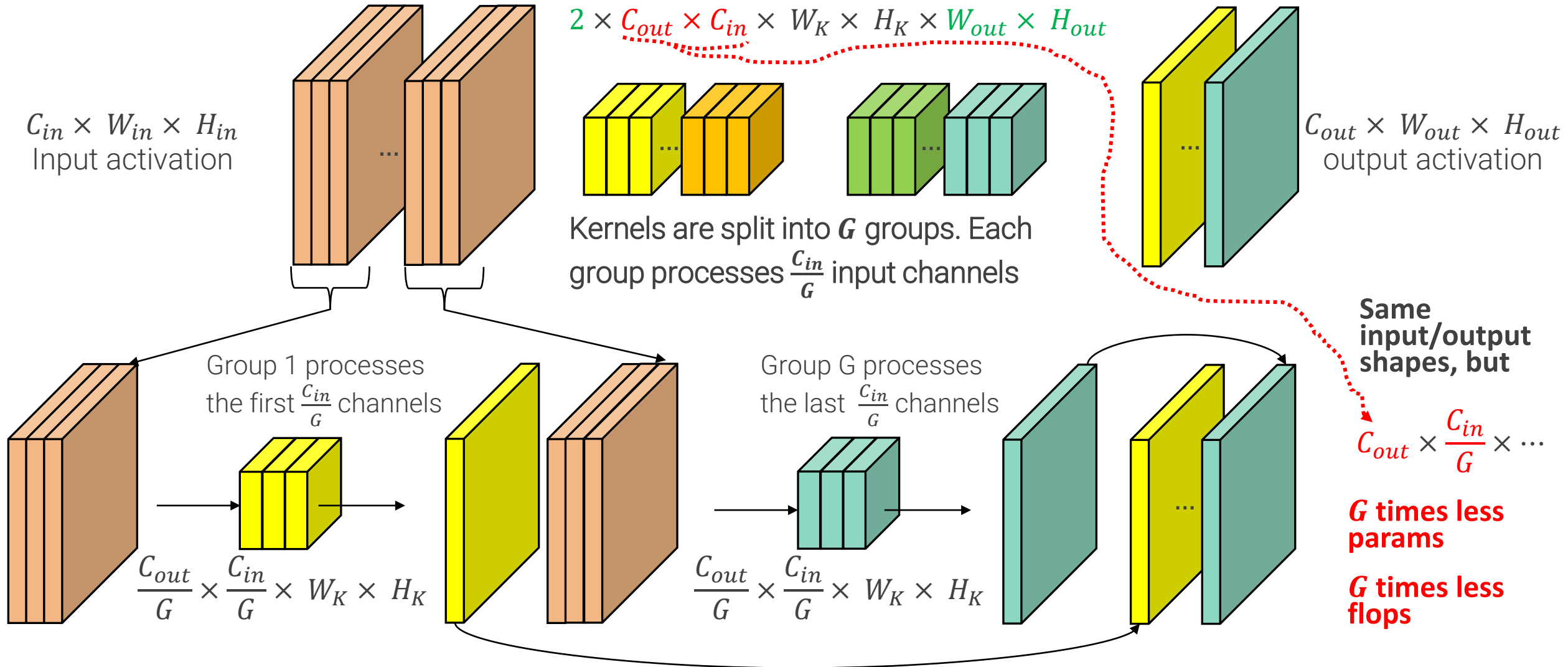


According to this design, the filters in the deeper layers have **larger receptive fields** than those in the shallower ones, that is they analyze larger regions of the input image to seek to find interesting patterns.

Hence, the filters in the initial layers can detect “small” patterns, i.e. **local features** such as, e.g. edges, corners, color blobs. As we move towards deeper layers, filters gain the ability to detect larger patterns, i.e. more **global features** (e.g. eyes, lips, mouth and then, further on, whole faces).

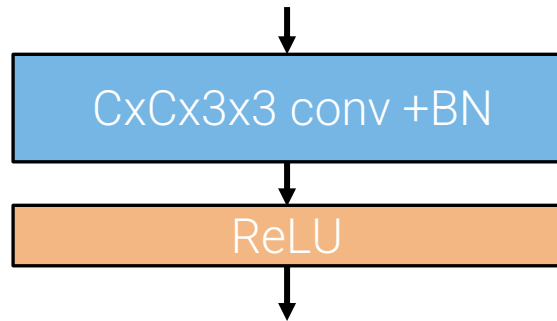
As the number of interesting patterns that may be found grows inherently with the receptive field, it makes a lot of sense to increase the number of detectors after down-sampling the activations.

Grouped Convolutions



Depthwise Separable Convolutions

Standard Convolution



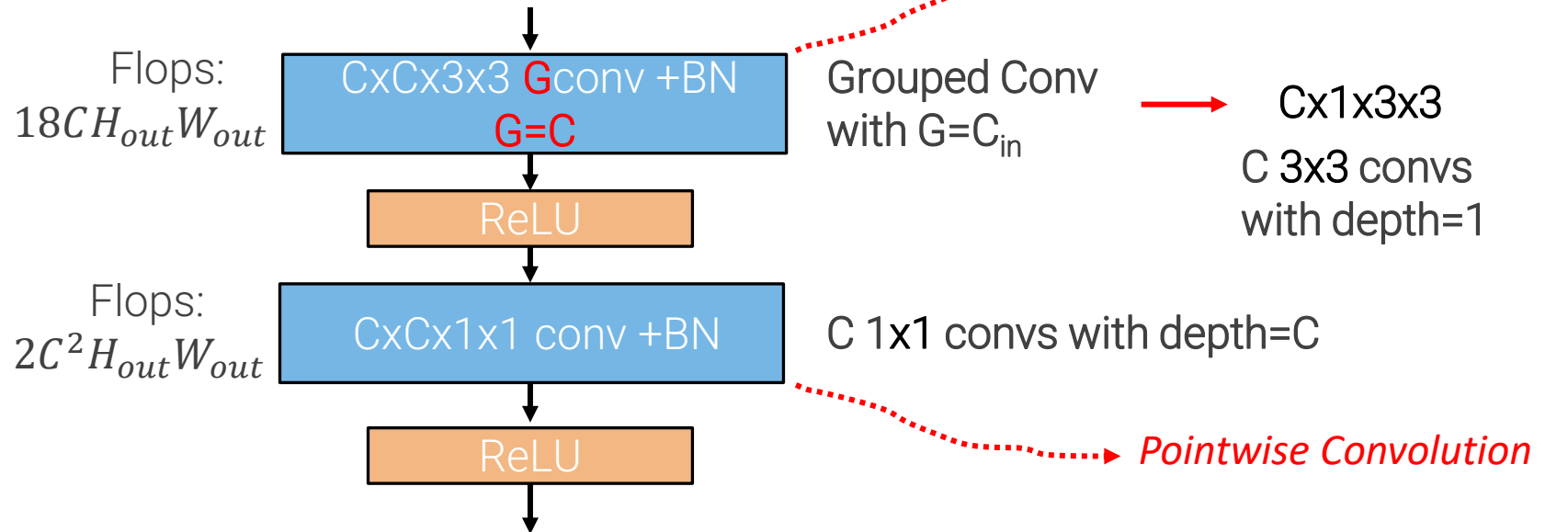
Flops: $18C^2 H_{out} W_{out}$

Computational Savigs:

$$\frac{18C^2}{18C + 2C^2} = \frac{18C}{18 + 2C} \cong [8, 8.85]$$

for C in $[64, 512]$

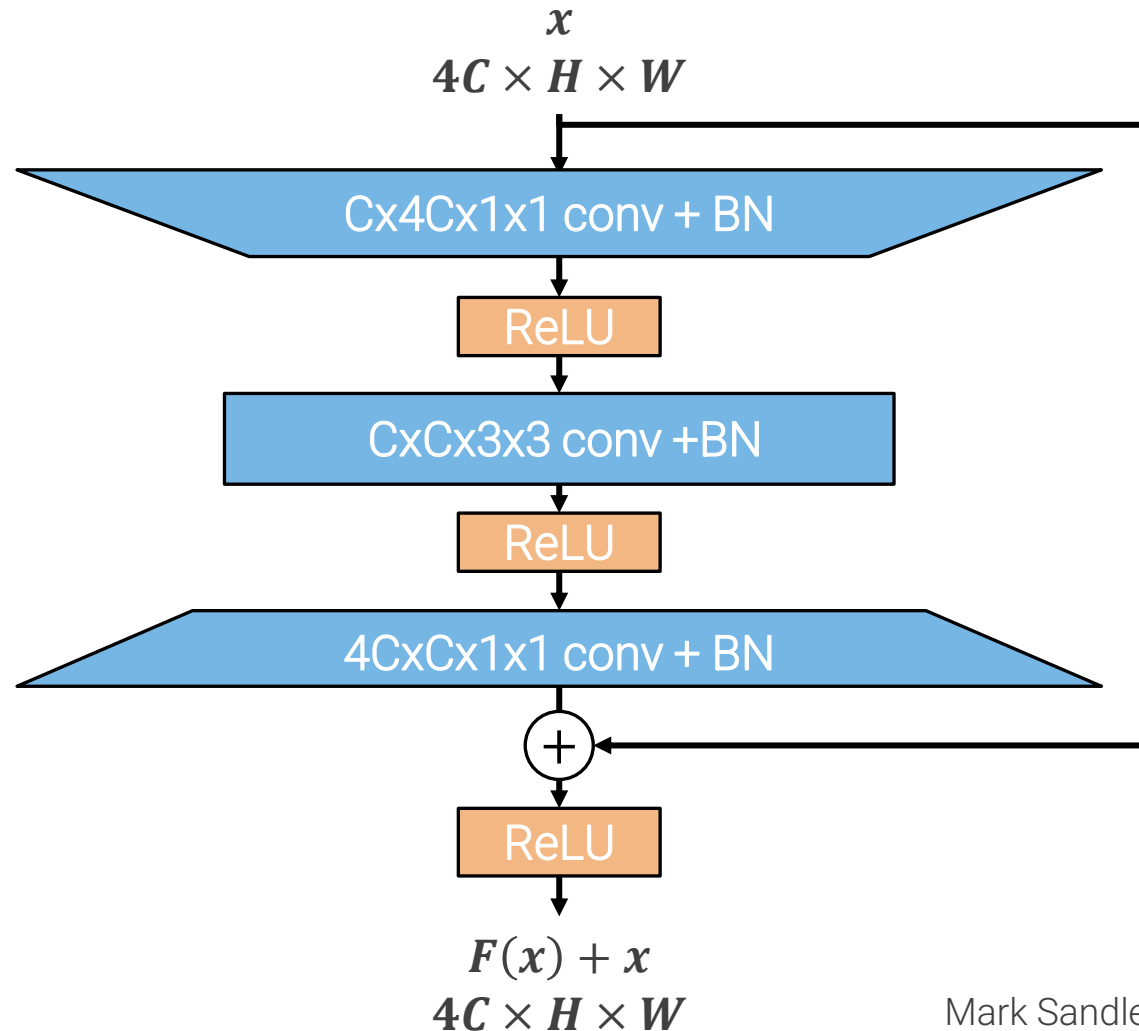
Depthwise Separable Convolution



Standard convolutions used in CNNs *filter* features *spatially* while *combining* them to produce new representations.

With **Depthwise Separable Convolutions** the two steps are split and carried out sequentially to gain substantial computational savings.

A closer look at Bottleneck Residual Blocks

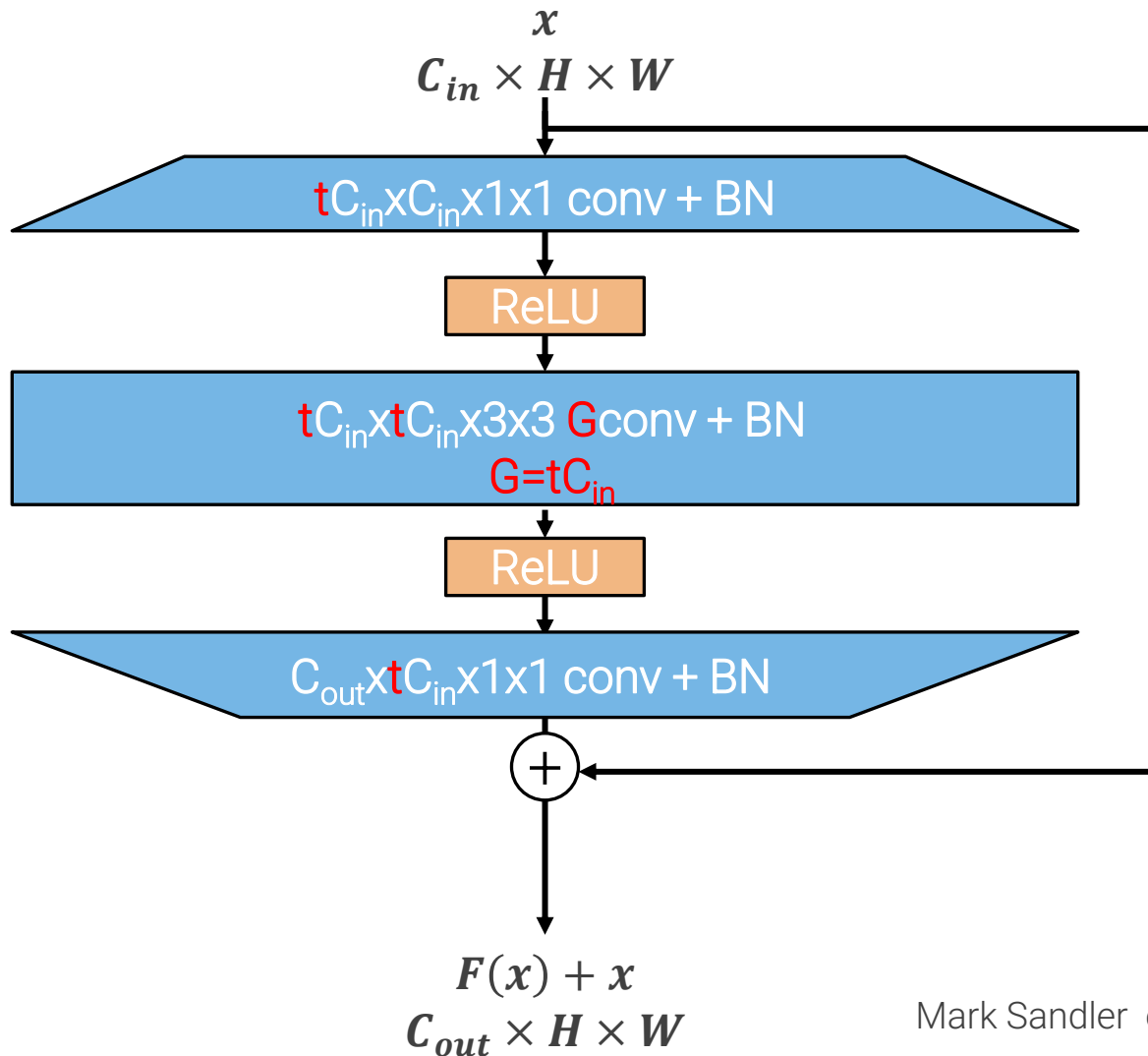


The bottleneck residual block was introduced to scale up the depth of ResNets by increasing significantly the number of blocks per stage without growing too much the computation and number of parameters.

Purposely, it uses a pair of 1x1 convs, where **the first compresses** the number of channel and **the second expands** them.

Hence, the 3x3 convolution that processes spatial information, i.e. the core representation learning function performed by the block, is carried out in a **compressed** domain. This may result in **information loss**.

Inverted Residual Block (IRB)



To avoid the potential information loss of standard bottleneck blocks, **MobileNet-v2** proposes to use **Inverted Residual Blocks**.

In such blocks, the first 1×1 conv expands the channels according to a chosen **expansion factor t** , while the second compresses them back (to the same or a different number of channels, i.e. C_{out} may be different than C_{in}).

To limit the increase in computation, the inner 3×3 convolution is realized as a **depthwise convolution**.

The last difference wrt standard bottleneck blocks is the **removal of non-linearities between residual blocks**: this is motivated by a theoretical study and experimentally verified.

Compared to the standard bottleneck block, the inverted design is **considerably more memory efficient** at inference time.

MobileNet-v2



- **MobileNet-v2** is a deep architecture specifically tailored for **mobile and resource constrained platforms**. It is based on a stack of **Inverted Residual Blocks** and features 54 layers with parameters. It deploys only 1x1 and 3x3 convs.
- The **number of channels** grows slowly compared to previous architectures to keep the complexity low. A low number of channels does not require an heavy size reduction in the **stem layer** ($s=2$). Yet, the representation must be expanded by a **pointwise convolution** before feeding it to final **k-way classifier** via **global average pooling**.
- As for the stack of **Inverted Residual Blocks**, each line in the table can be seen as a *stage*. When a stage down-samples the activation, it does so by applying $s=2$ in the inner 3x3 conv of the first Inverted Residual Block.
- Whenever spatial dimensions or number of channels do not match between input and output of a block, there are no skip connections.

Data taken
from the
paper

Top-1	#params	MACs	CPU
72.0	3.4M	300M	75 ms

Input	Operator	t	c	n	s
$224^2 \times 3$	conv2d	-	32	1	2
$112^2 \times 32$	bottleneck	1	16	1	1
$112^2 \times 16$	bottleneck	6	24	2	2
$56^2 \times 24$	bottleneck	6	32	3	2
$28^2 \times 32$	bottleneck	6	64	4	2
$14^2 \times 64$	bottleneck	6	96	3	1
$14^2 \times 96$	bottleneck	6	160	3	2
$7^2 \times 160$	bottleneck	6	320	1	1
$7^2 \times 320$	conv2d 1x1	-	1280	1	1
$7^2 \times 1280$	avgpool 7x7	-	-	1	-
$1 \times 1 \times 1280$	conv2d 1x1	-	k	-	-

Google Pixel 1

e.g. $k=1000$ (ILSVRC)

Number of IRBs

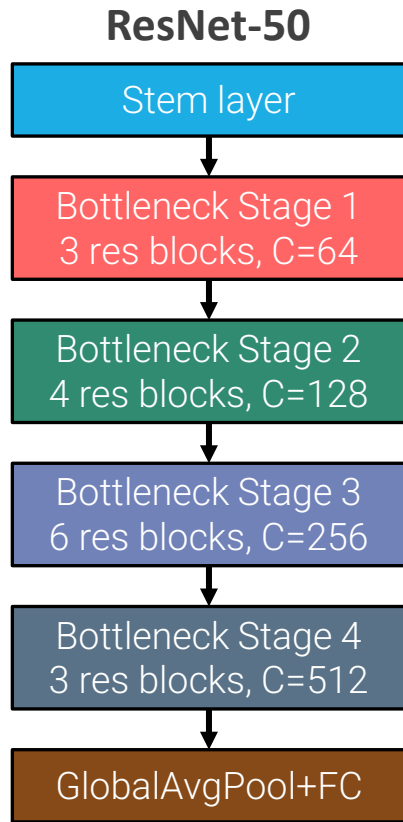
Transfer Learning (1)



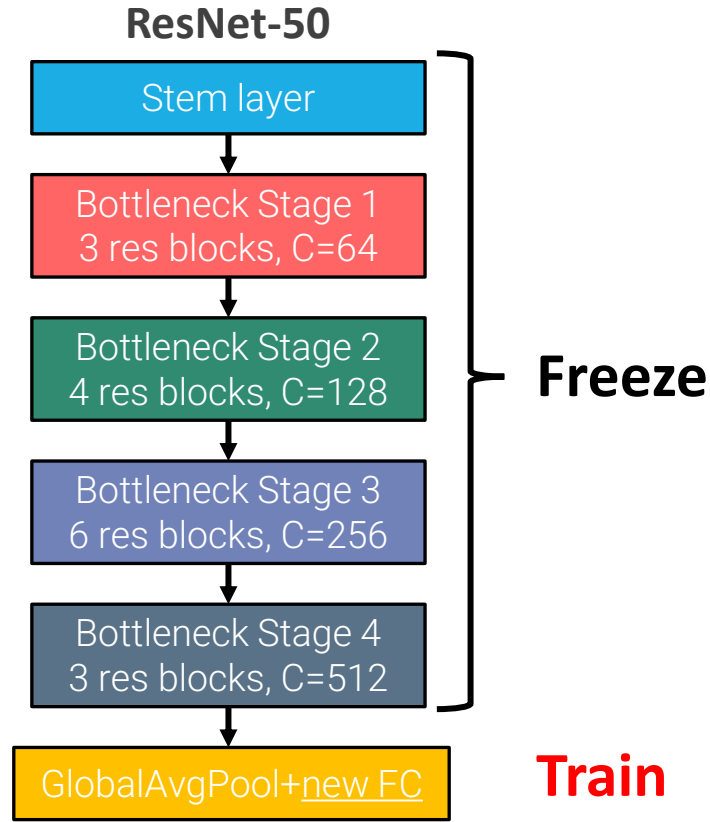
- To prevent **overfitting**, a large/deep (i.e. **high capacity**) neural network requires a large number of samples to effectively train its many parameters. But annotated (aka **supervised**) training data are expensive...what if in our scenario we only have a small training set?
- We may deploy a very effective two-steps approach referred to as **Transfer Learning**:
 1. **Pre-train** the network on a large dataset (e.g. ImageNet)
 2. **Fine-tune** the pre-trained network on the smaller, task-specific dataset.
- Typically, in the first step one relies on **standard architectures** (e.g. ResNet-50) and downloads the pretrained model from a public repository.

Transfer Learning (2)

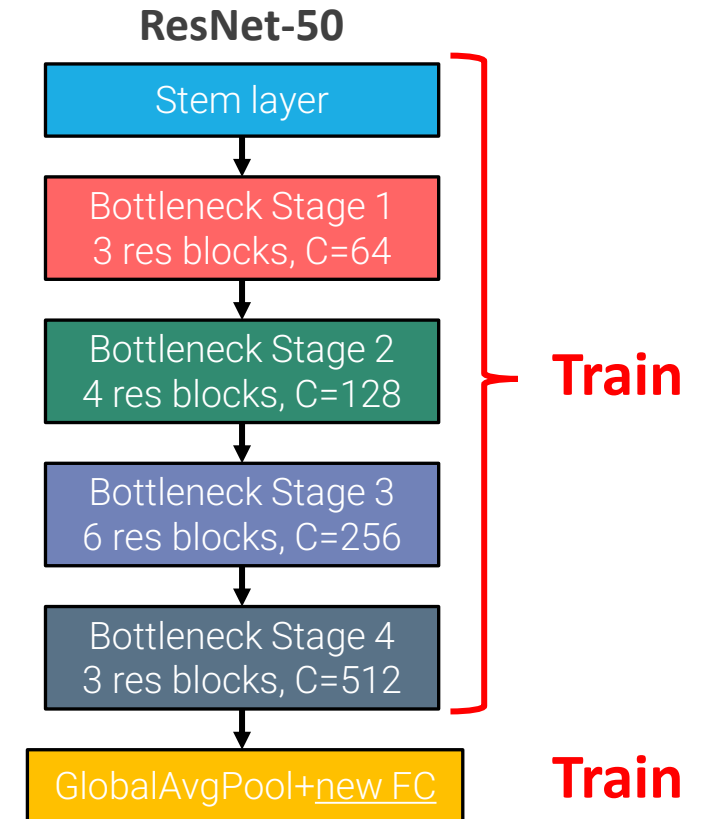
1. Pre-Trained Model



2.A Frozen Feature Extractor



2.B Train new Head and Feature Extractor



2.B: *warm-up* with a frozen feature extractor (like in **2.A**) then fine-tune the whole model with a very small learning rate. The initial layers may still be kept frozen as they typically learn general, low-level features.