

Abdul Qayyum

Lecturer at University of Burgundy, France

- Postdoc in Electrical and Informatics Engineering
- PhD in Electrical & Electronics Engineering
- Masters in Electronics Engineering
- Bachelor in Computer Engineering

Collaborations & Expertise:



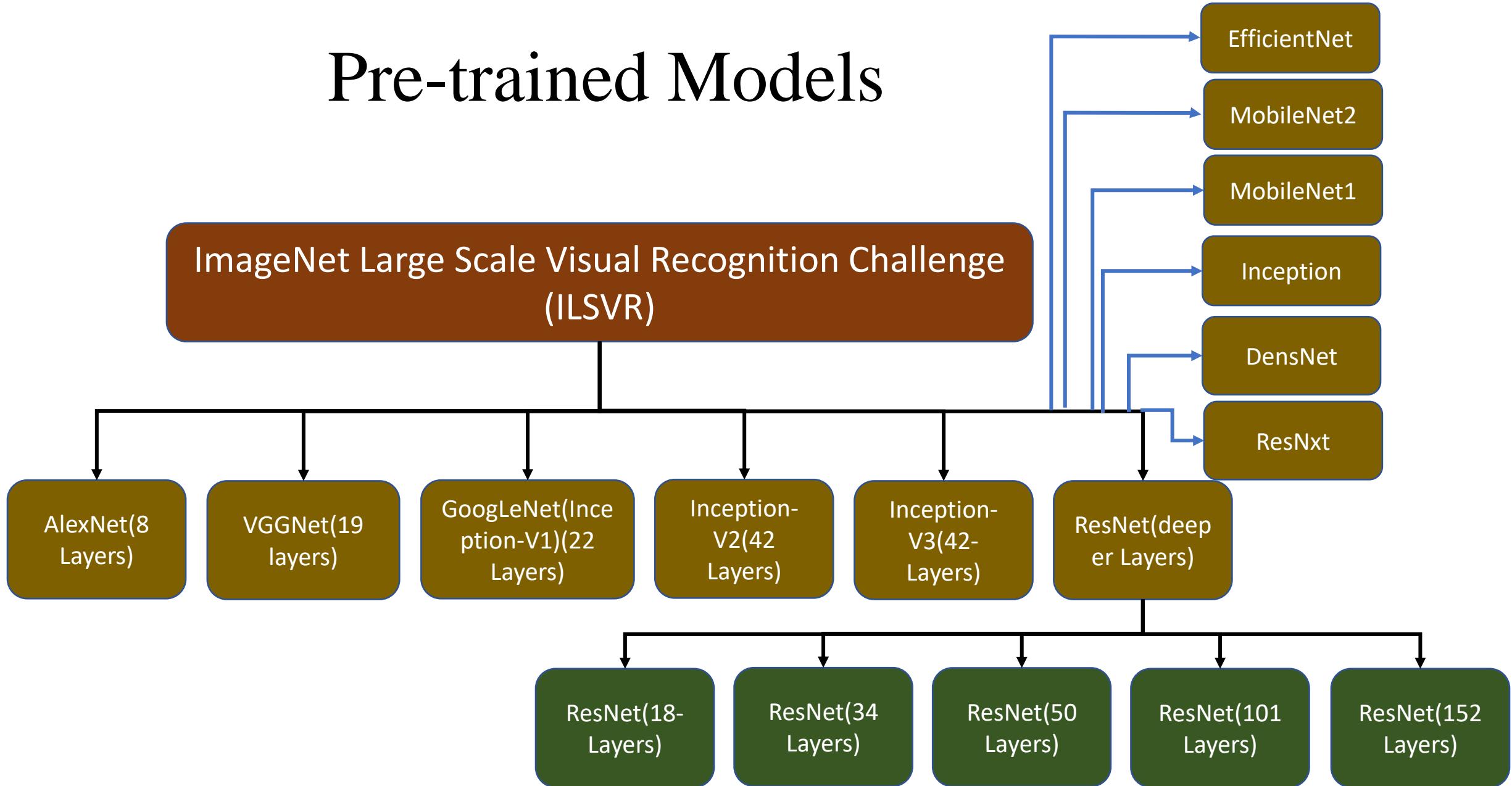
Topic: **Deep Learning**
Classification Models

Instructor: Abdul Qayyum, PhD

Class: MSCV

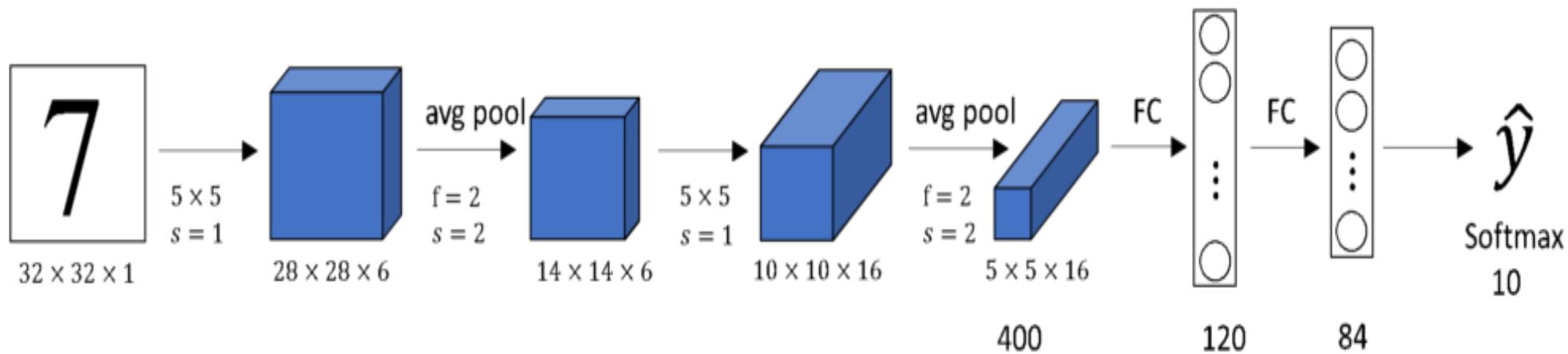
University of Burgundy, France

Pre-trained Models



LeNet

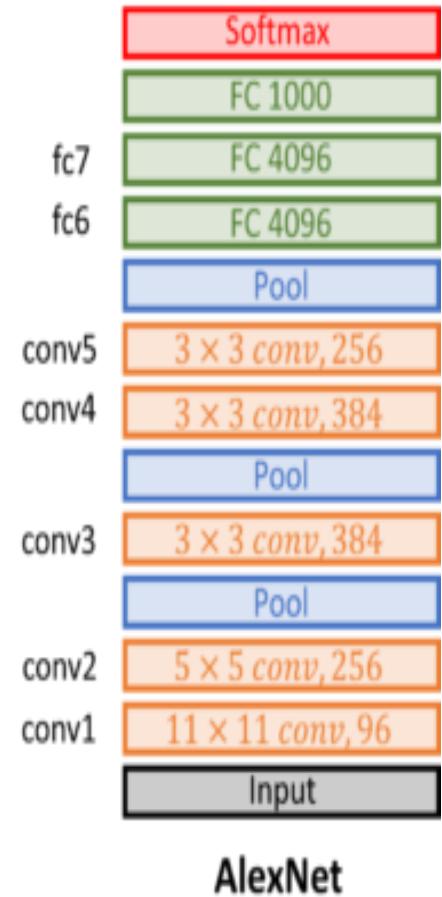
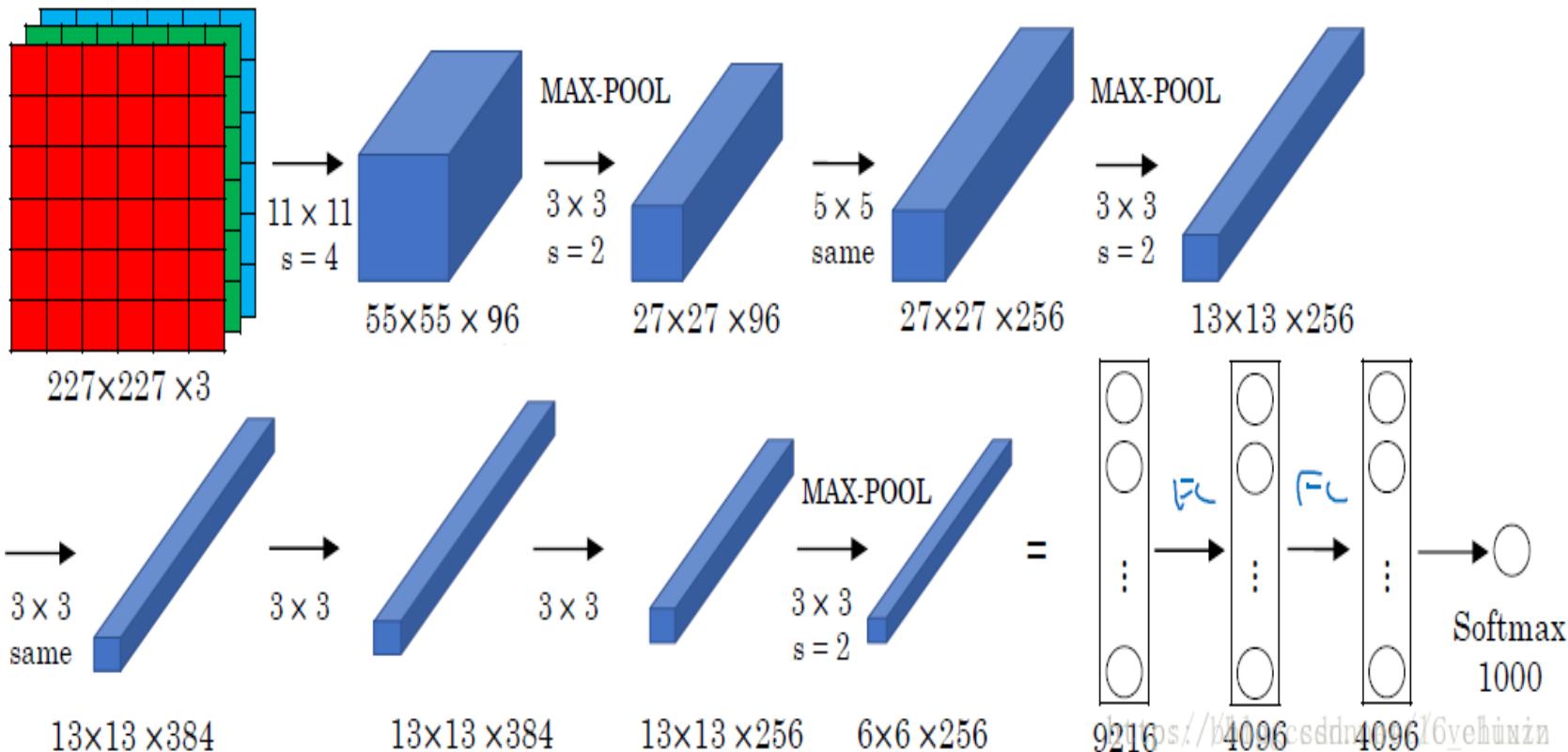
▪ LeNet(Layers)



AlexNet

AlexNet(8 Layers)

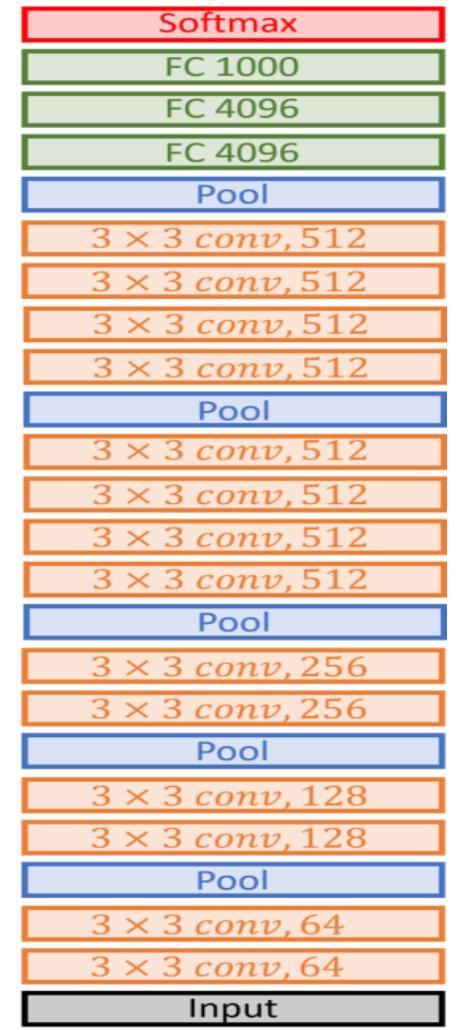
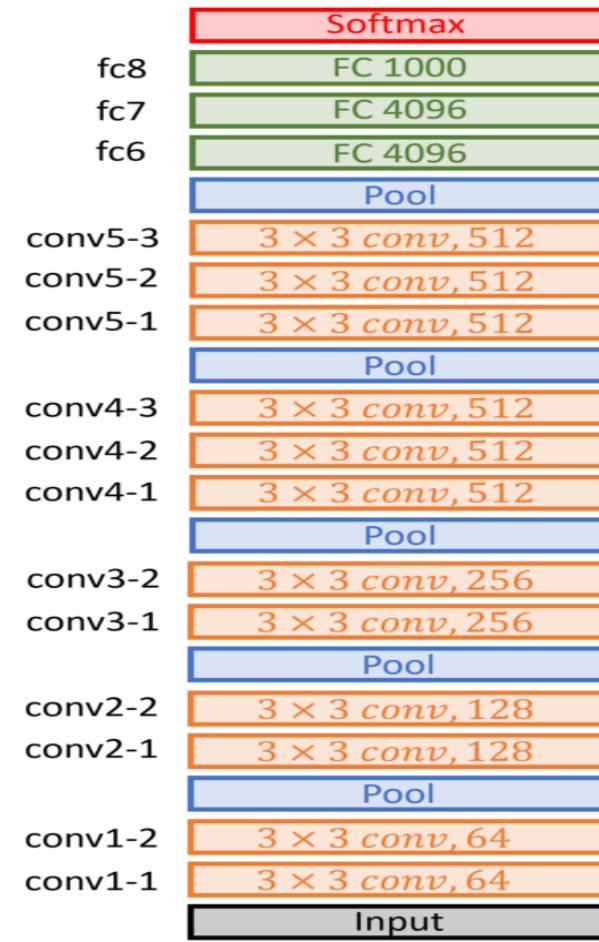
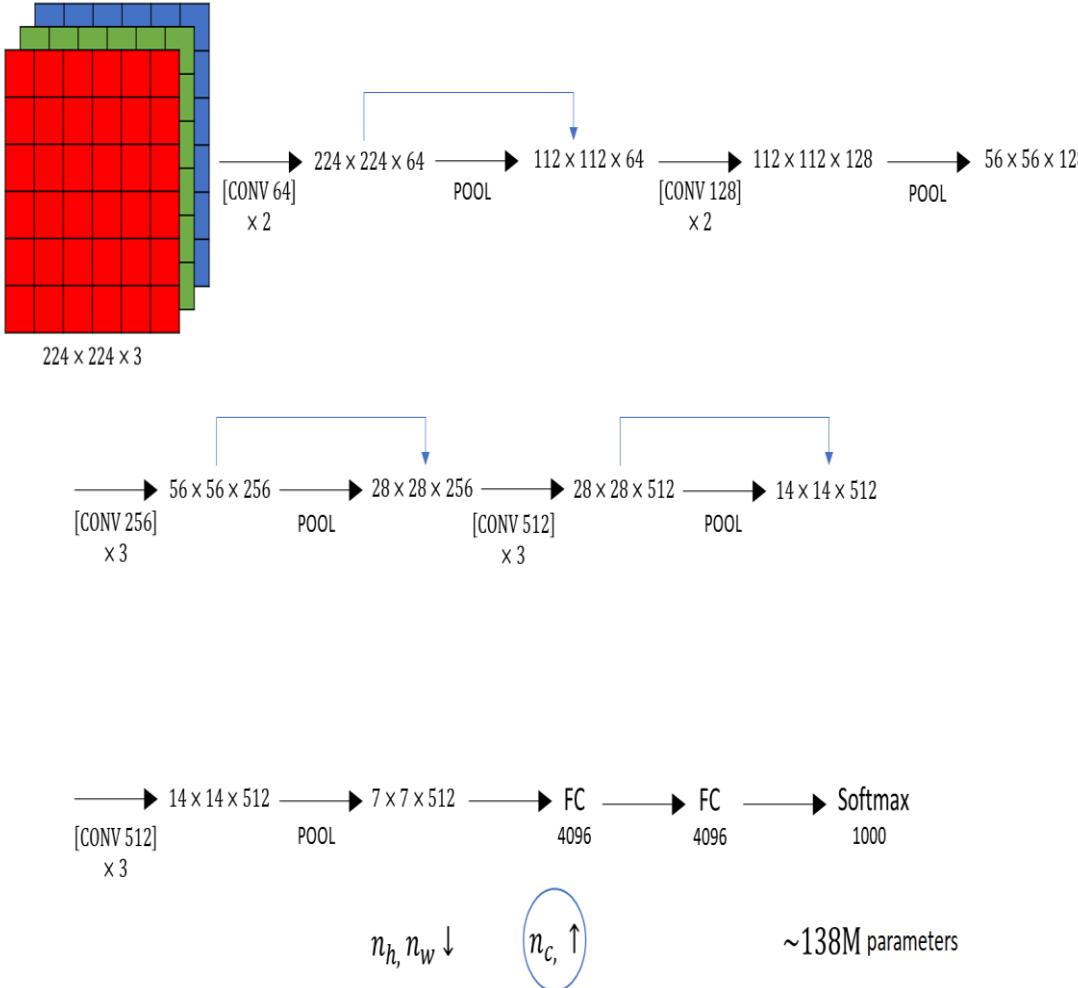
AlexNet



VGG Net

■ VGG16 and VGG19

CONV = 3×3 filter, $s=1$, same
MAX-POOL = 2×2 , $s=2$



Network in Network

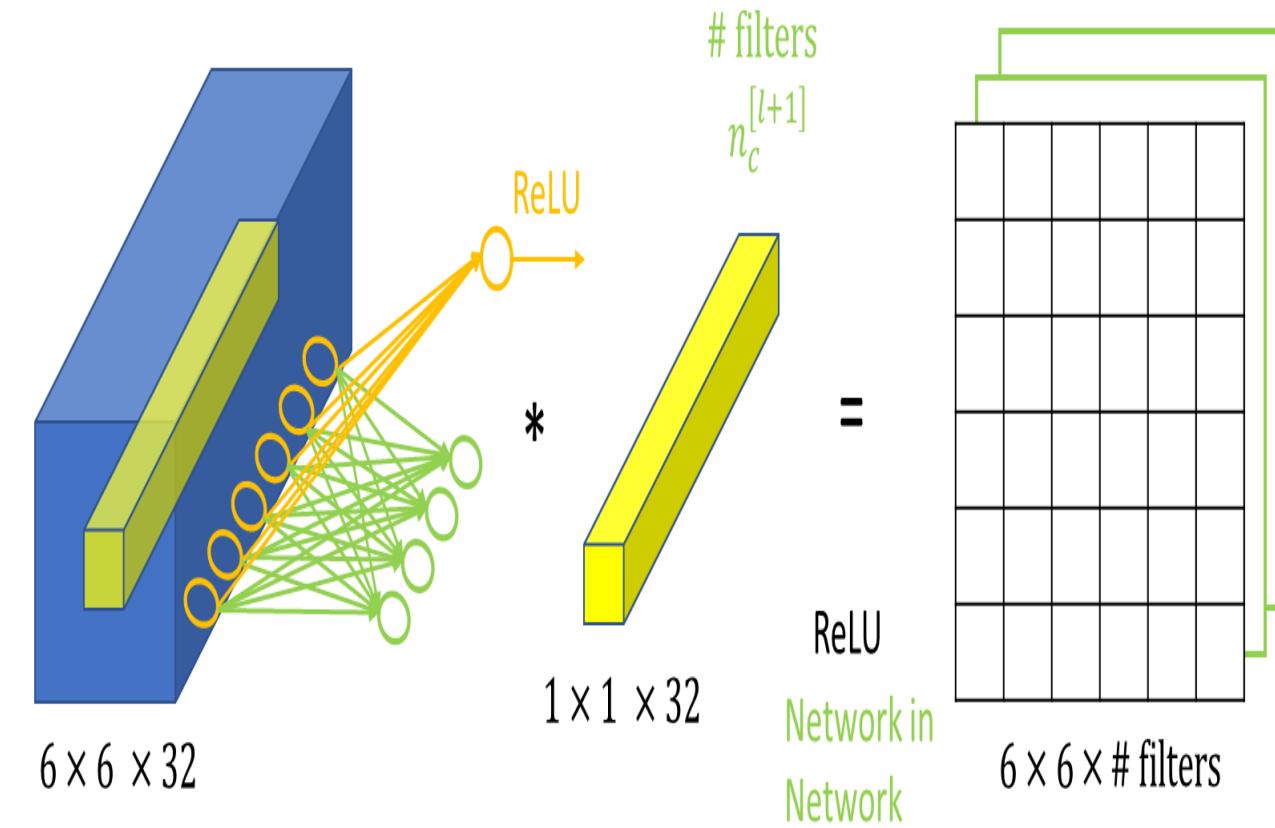
<http://datahacker.rs/introduction-into-1x1-convolution/>

▪ Network in Network – 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

$$\begin{array}{c} * \\ \text{ } \\ \text{ } \end{array} \quad \boxed{2} \quad = \quad \begin{array}{c} 2 \\ \text{ } \\ \text{ } \end{array}$$

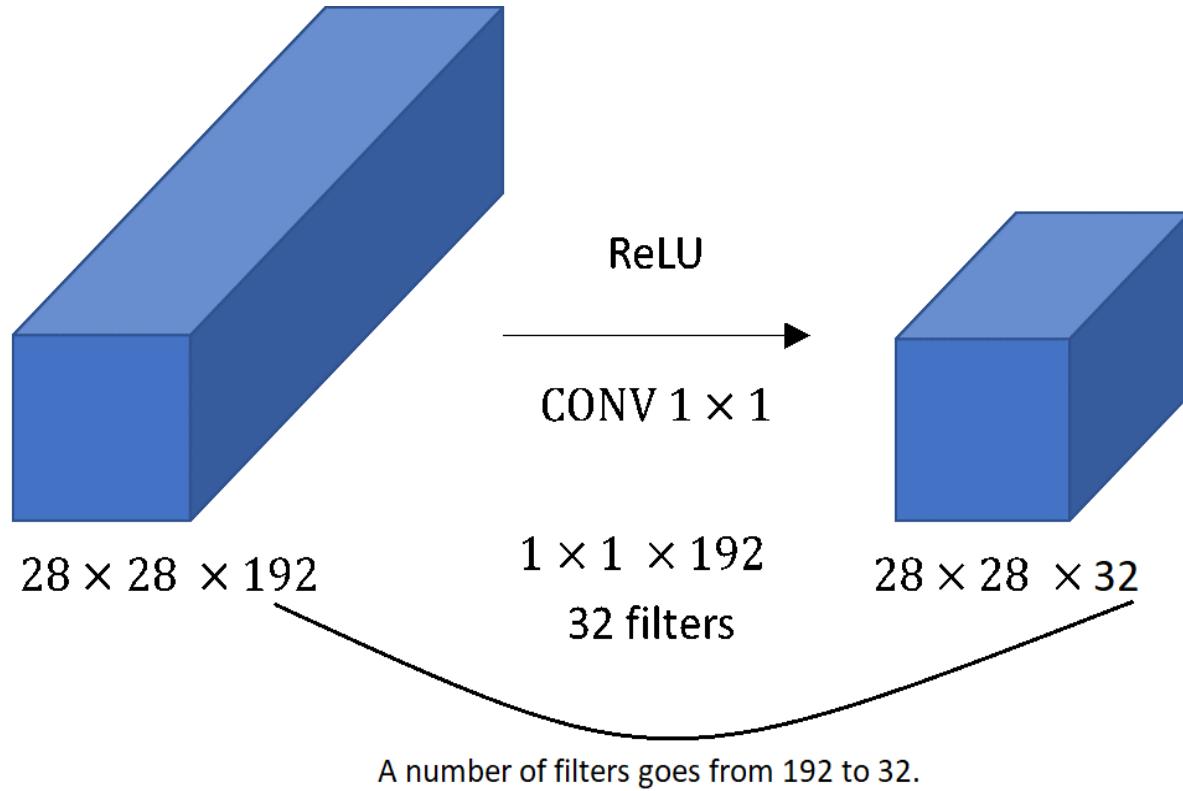
6×6 6×6



In particular a 1×1 convolution will look at each of the 36 different positions (6×6), and it will take the element-wise product between 32 numbers on the left and the 32 numbers in the filter. Then, a ReLU non-linearity will be applied. Looking at one out of the 36 positions, maybe one slice through this volume, we take these 32 numbers, multiply it by 1×1 slice through the volume, and we get a single number.

Network in Network

- Network in Network – 1×1 Convolutions

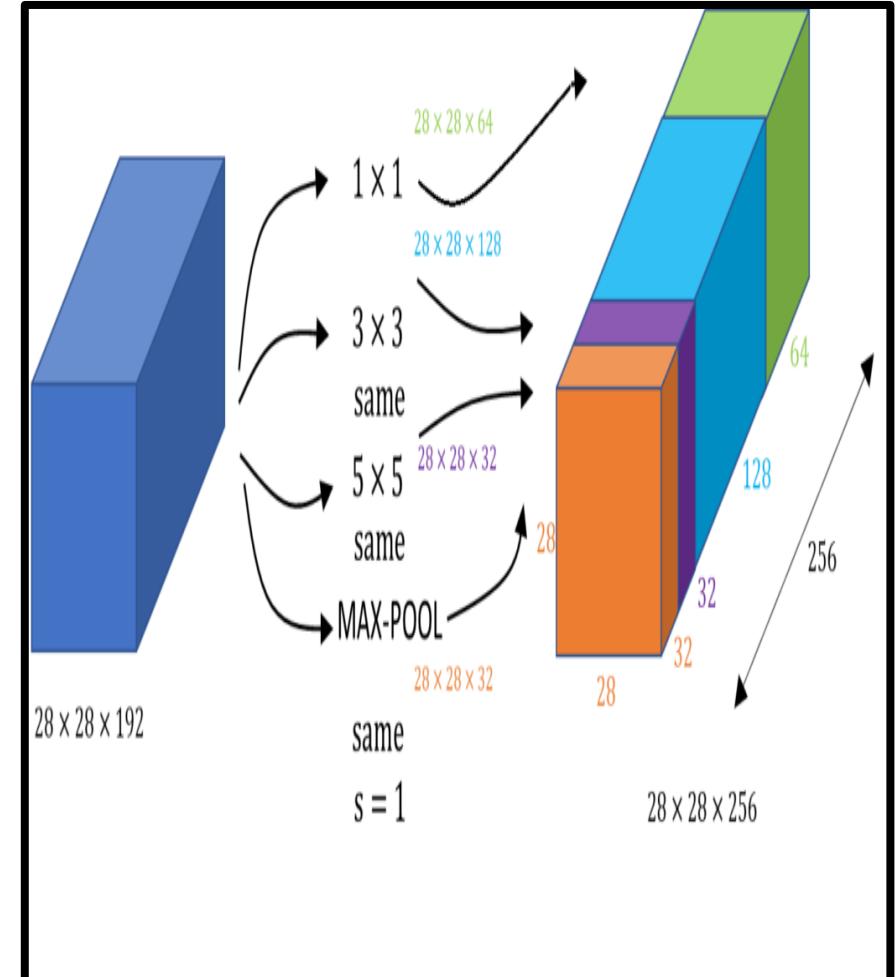
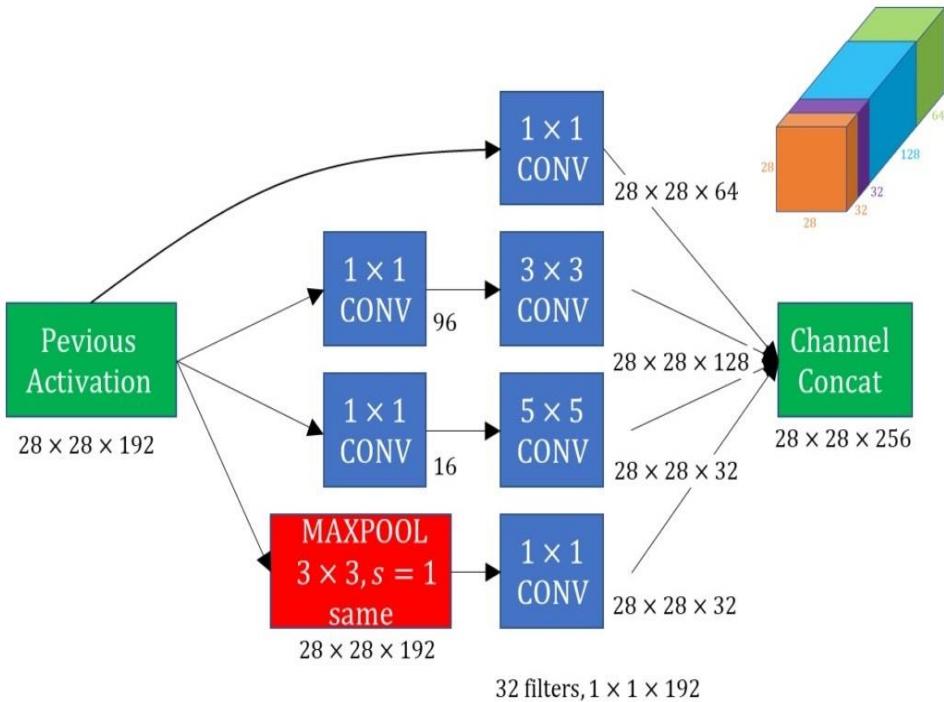


We can do is use 32 filters that are 1×1 , and technically each filter would be of dimension $1 \times 1 \times 192$ because the number of channels in our filter has to match the number of channels in input volume. But we use 32 filters and the output of this process will be a $28 \times 28 \times 32$ volume. This is a way to let us shrink nc (number of channels). We'll see later how this idea of 1×1 convolutions allows us to shrink the number of channels and therefore save on computations and networks.

Inception Module

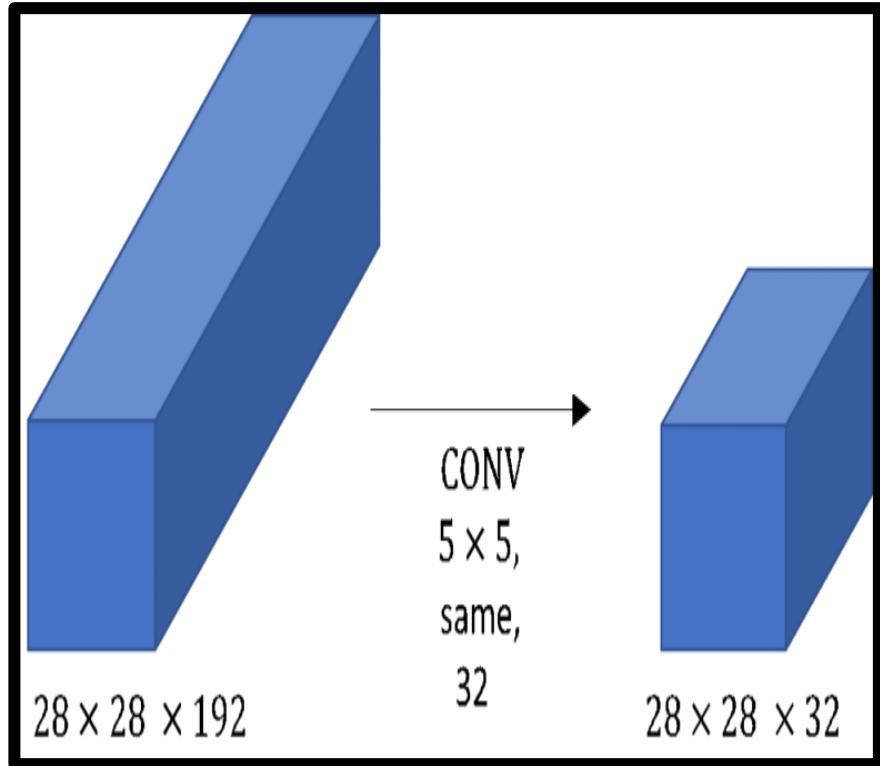
▪ Inception Module

An example of an Inception module

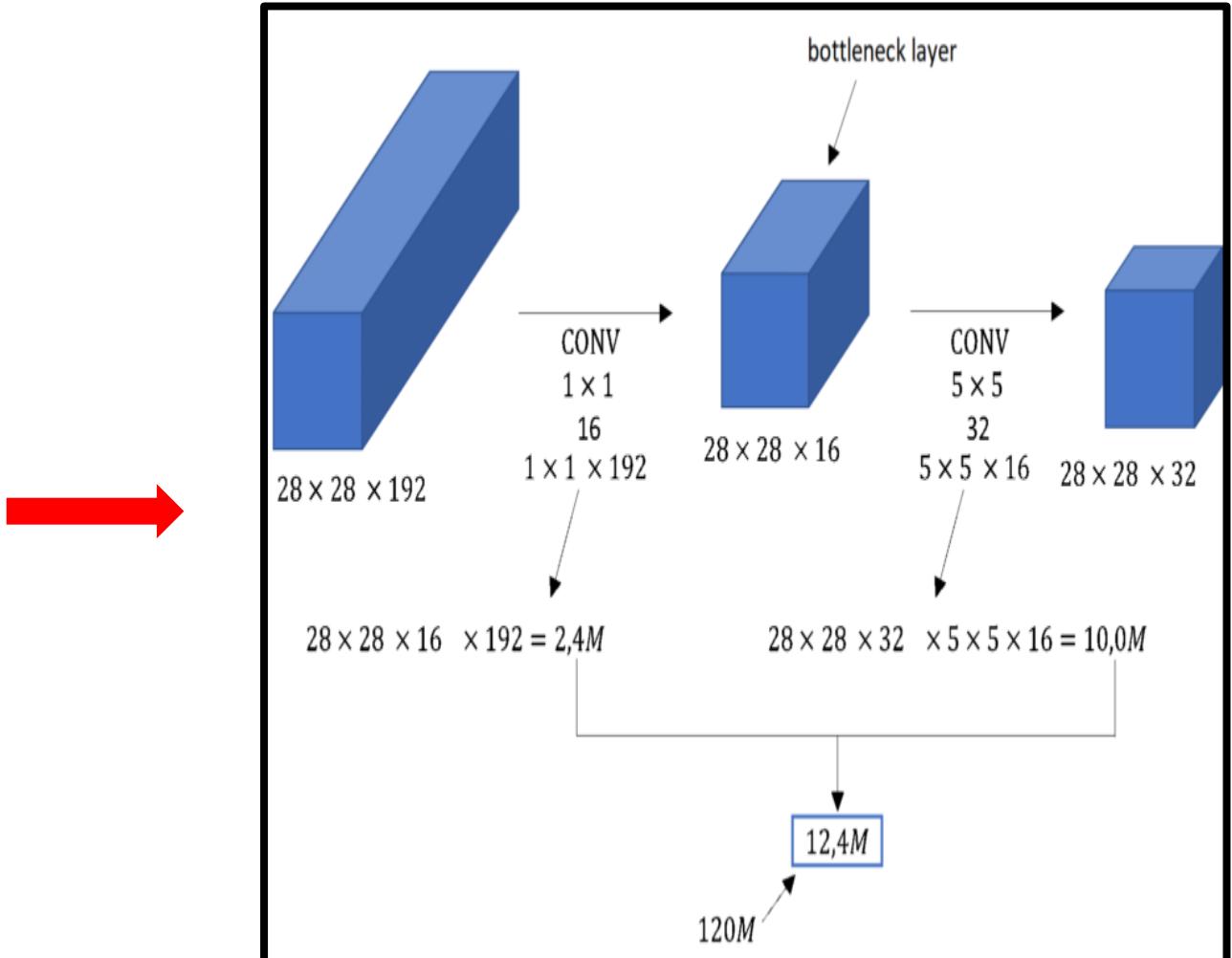


Inception Module

▪ Inception Module



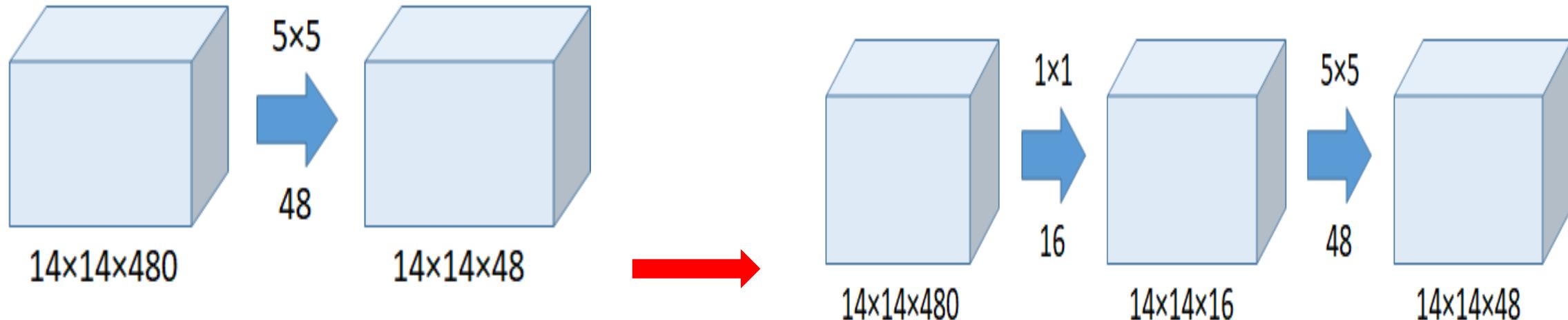
$$28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120M$$



Inception Module

▪ Inception Module

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

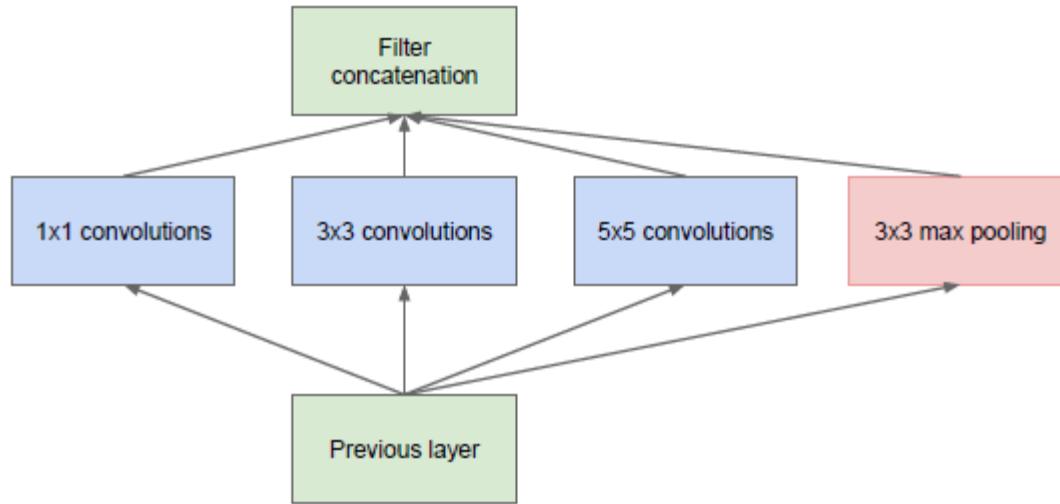


$$\text{Number of operations} = (14 \times 14 \times 48) \times (5 \times 5 \times 480) = 112.9M$$

Number of operations for 1×1 = $(14 \times 14 \times 16) \times (1 \times 1 \times 480) = 1.5M$
Number of operations for 5×5 = $(14 \times 14 \times 48) \times (5 \times 5 \times 16) = 3.8M$
Total number of operations = $1.5M + 3.8M = 5.3M$
which is much much smaller than 112.9M !!!!!!!!!!!!!!!

Optimization Algorithms

▪ Inception Module



(a) Inception module, naïve version

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

Previously, such as AlexNet, and VGGNet, conv size is fixed for each layer.

Now, **1×1 conv, 3×3 conv, 5×5 conv, and 3×3 max pooling** are done altogether for the previous input, and stack together again at output. When image's coming in, different sizes of convolutions as well as max pooling are tried. Then different kinds of features are extracted.

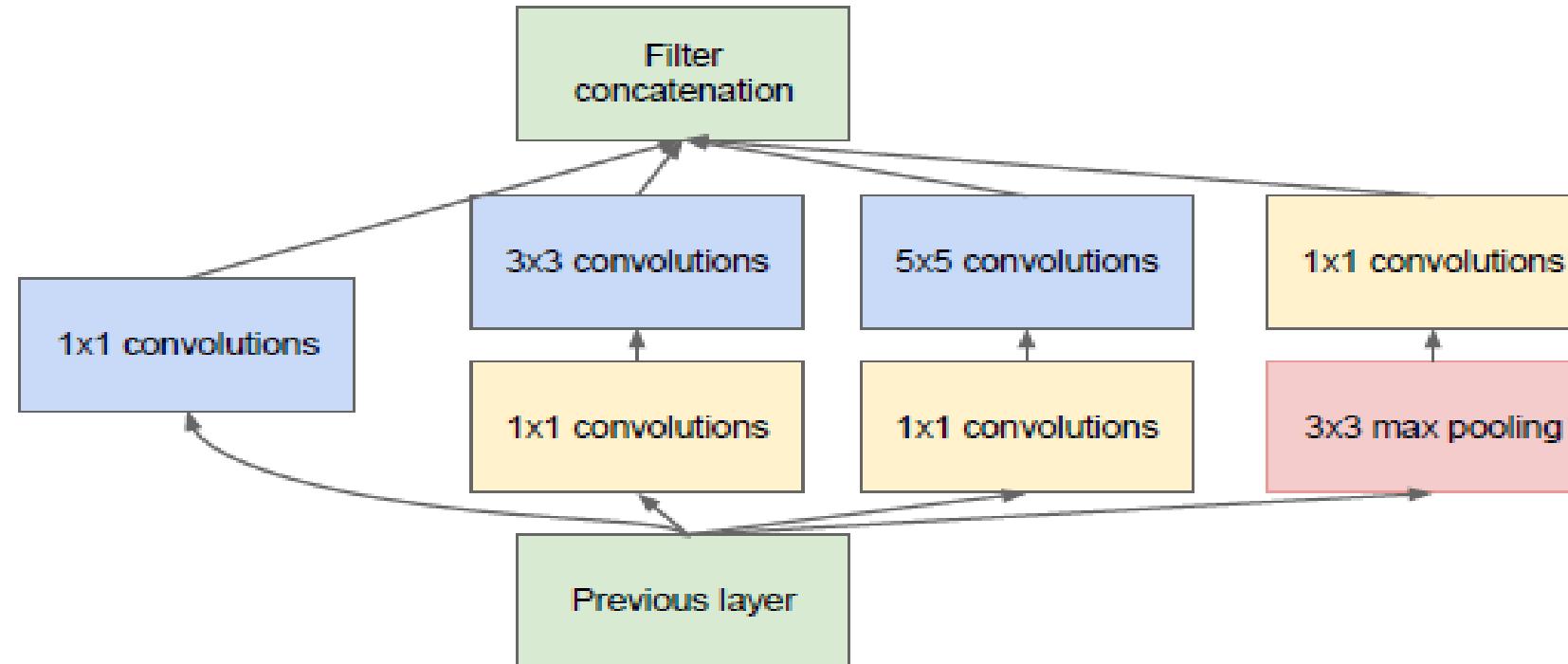
After that, all feature maps at different paths are concatenated together as the input of the next module.

However, without the 1×1 convolution as above, we can imagine how large the number of operation is!

Optimization Algorithms

▪ Inception Module

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>



(b) Inception module with dimensionality reduction

Optimization Algorithms

- **Inception Module**

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

“The Inception deep convolutional architecture was introduced as GoogLeNet in (Szegedy et al. 2015a), here named Inception-v1.

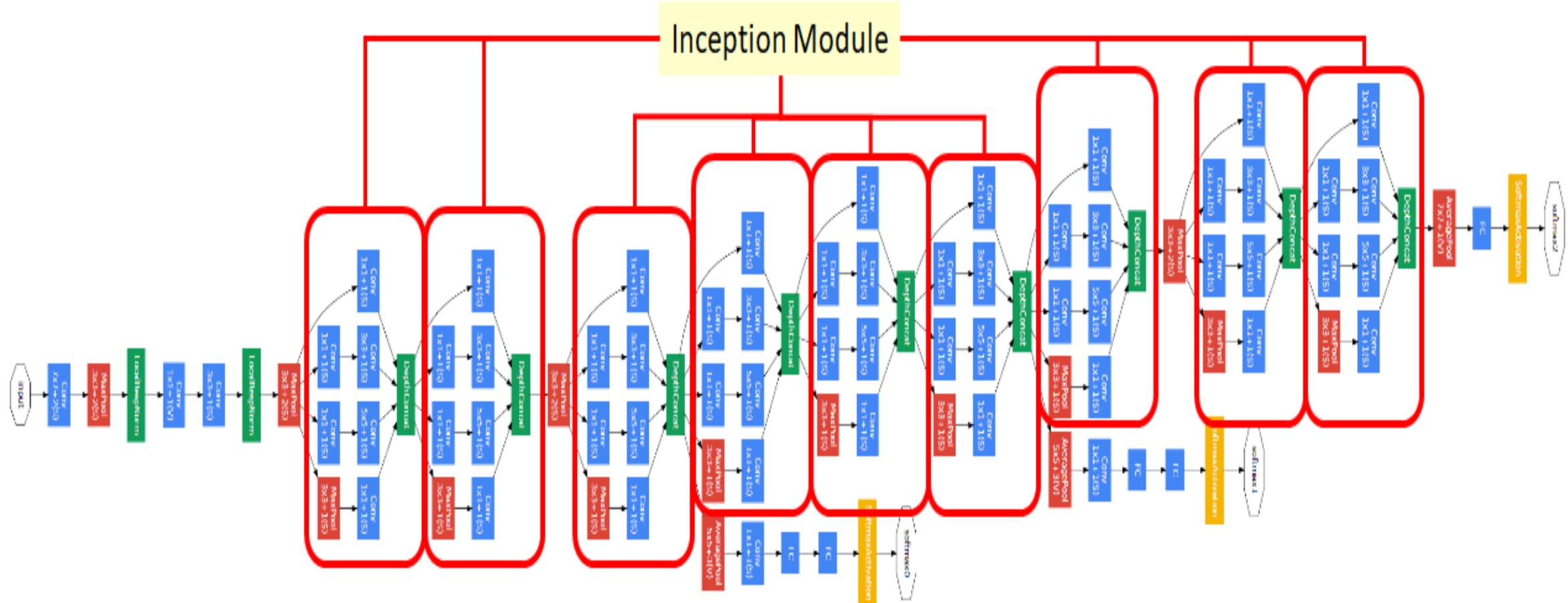
Later the Inception architecture was refined in various ways, first by the introduction of batch normalization (Ioffe and Szegedy 2015) (Inception-v2).

Later by additional factorization ideas in the third iteration (Szegedy et al. 2015b) which will be referred to as Inception-v3 in this report.”

GoogLeNet

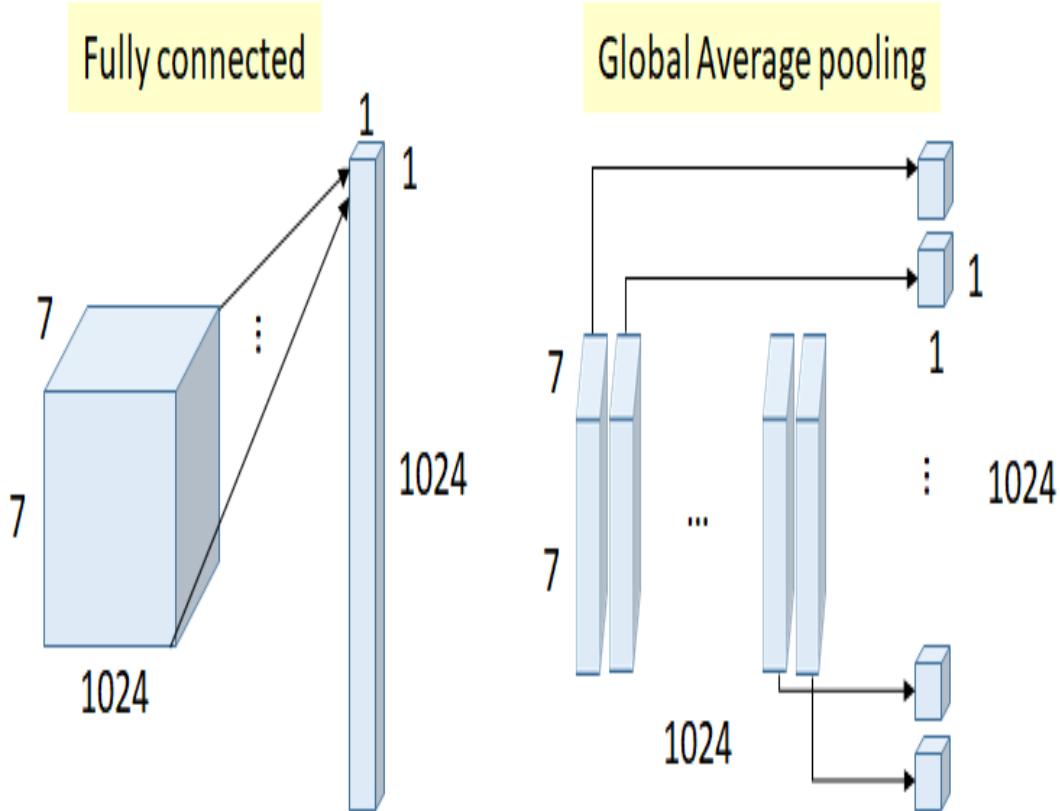
▪ GoogLeNet(Inception_V1)

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>



GoogleNet

▪ GoogleNet



Previously, fully connected (FC) layers are used at the end of network, such as in AlexNet.
All inputs are connected to each output.
Number of weights (connections) above = $7 \times 7 \times 1024 \times 1024 = 51.3M$

In GoogLeNet, global average pooling is used nearly at the end of network by averaging each feature map from 7×7 to 1×1 .
Number of weights = 0
And authors found that a move from FC layers to average pooling improved the top-1 accuracy by about 0.6%.

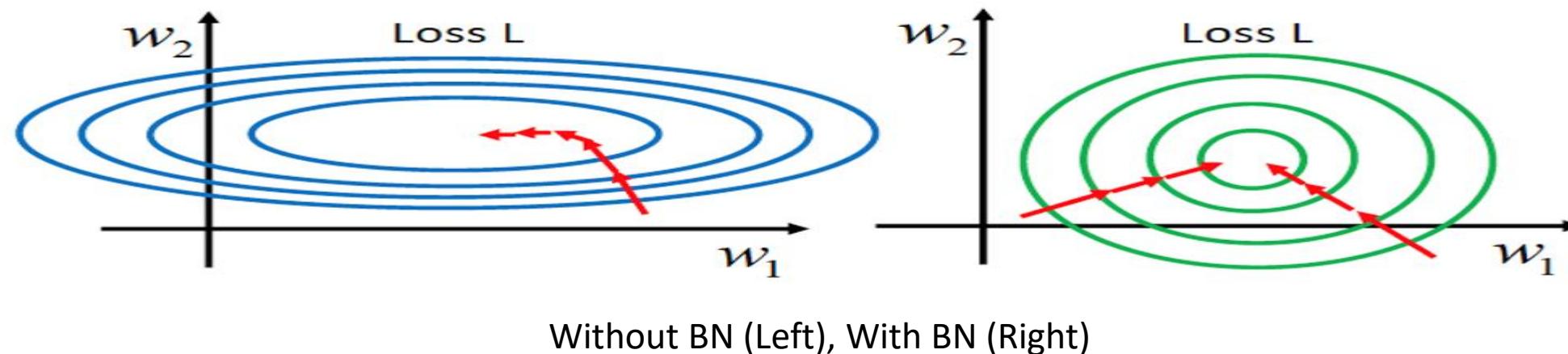
<https://medium.com/coinmonks/paper-review-of-googlenet-inception-v1-winner-of-ilsvrc-2014-image-classification-c2b3565a64e7>

Inception v2

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

- **Why we need Batch Normalization (BN)?**

- As we should know, the input X is multiplied by weight W and added by bias b and become the output Y at the next layer after an activation function F:
- $Y=F(W \cdot X+b)$
- Previously, F is sigmoid function which is easily saturated at 1 which easily makes the gradient become zero. As the network depth increases, this effect is amplified, and thus slow down the training speed.
- ReLU is then used as F, where $\text{ReLU}(x)=\max(x,0)$, to address the saturation problem and the resulting vanishing gradients. However, careful initialization, learning rate settings are required.



Inception v2

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

- **Why we need Batch Normalization (BN)?**
- It is advantageous for the distribution of X to remain fixed over time because a small change will be amplified when network goes deeper.
- BN can reduce the dependence of gradients on the scale of the parameters or their initial values. As a result,
- **Higher learning rate can be used.**
- **The need for Dropout can be reduced.**

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;
Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

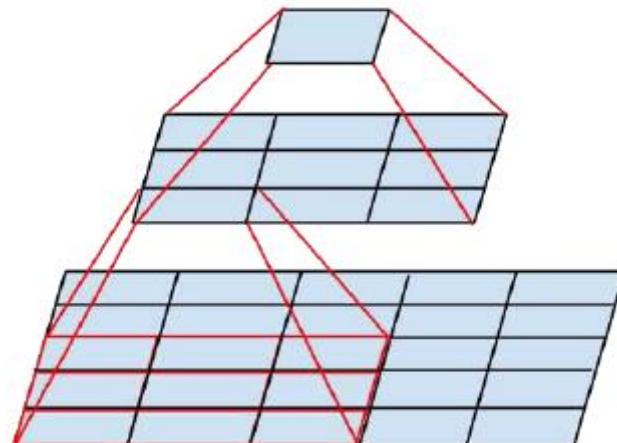
$Y=F(\text{BN}(W \cdot X+b))$
To have a more precise mean and variance, moving average is used to calculate the mean and variance.

Inception v2 and V3

- **Factorization Into Asymmetric Convolutions**
- Thus, the BN-Inception / Inception-v2 [6] is talking about batch normalization while Inception-v3 [1] is talking about factorization ideas.
- **Factorizing Convolutions**
- The aim of factorizing Convolutions is to reduce the number of connections/parameters without decreasing the network efficiency.

Factorization Into Smaller Convolutions

Two 3×3 convolutions replaces one 5×5 convolution as follows:



Two 3×3 convolutions replacing one 5×5 convolution

By using 1 layer of 5×5 filter, number of operations = $5 \times 5 = 25$

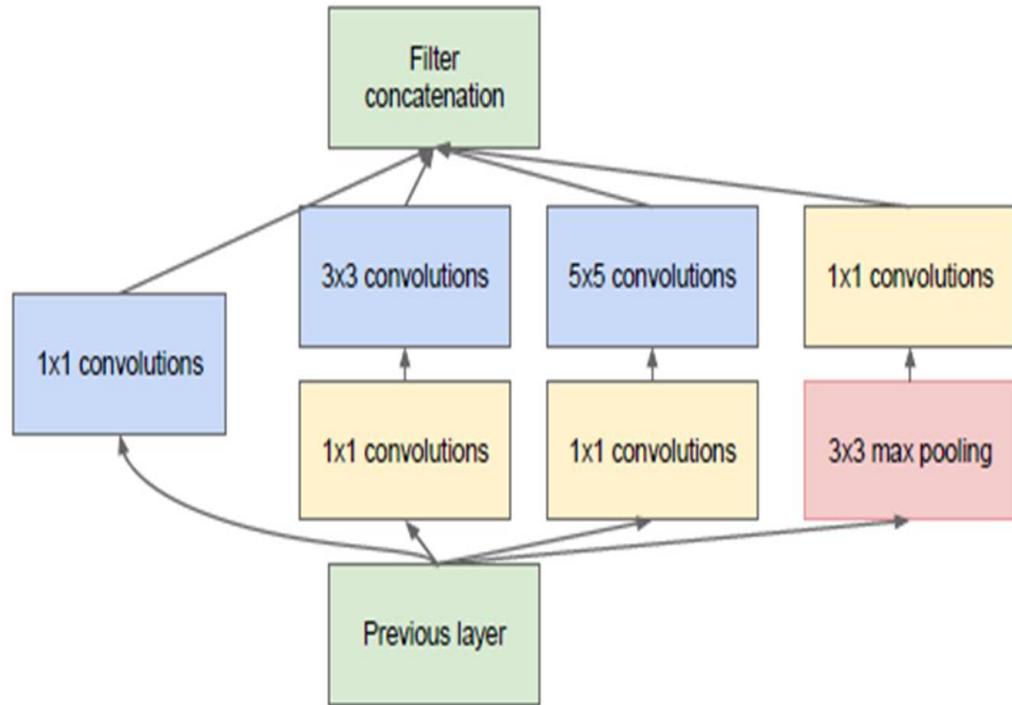
By using 2 layers of 3×3 filters, number of operations =
 $3 \times 3 + 3 \times 3 = 18$

Number of operations is reduced by 28%

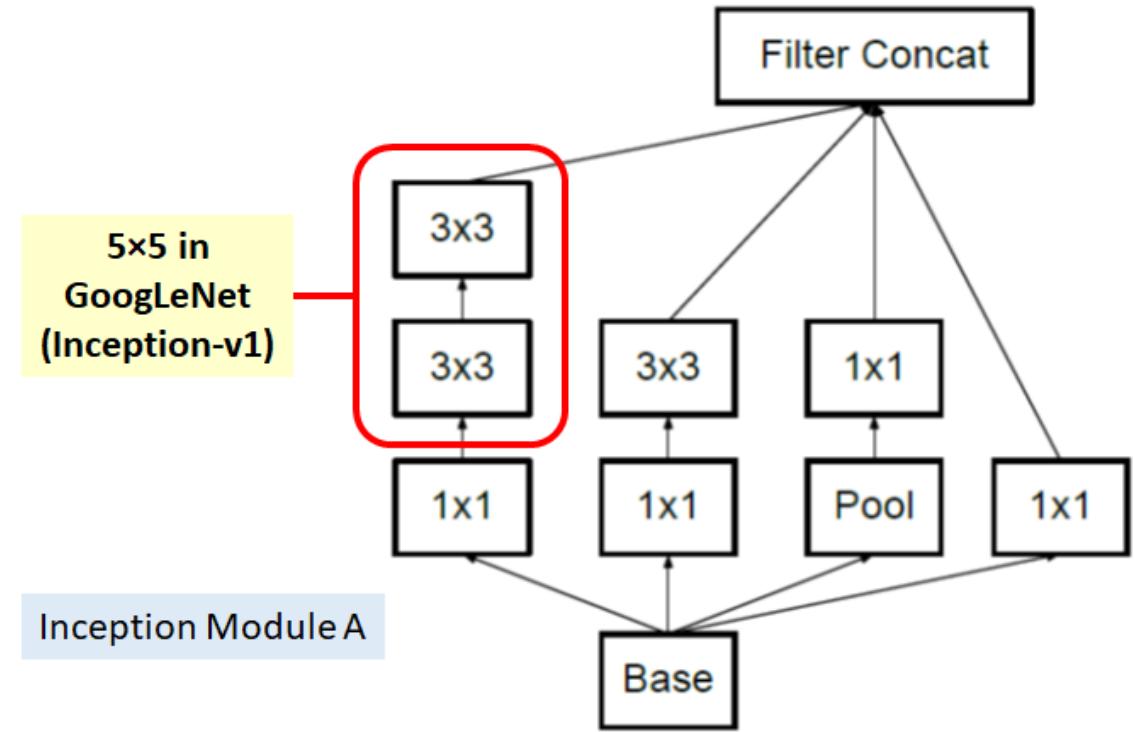
<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

Inception v2 and V3

- **Factorization Into Asymmetric Convolutions**
- With this technique, one of the new Inception modules (I call it Inception Module A here) becomes:



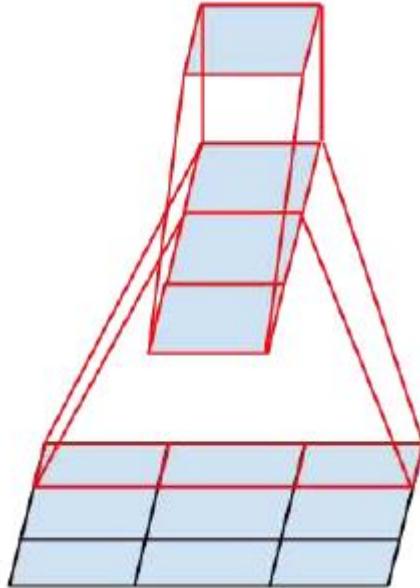
(b) Inception module with dimensionality reduction



<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

Inception v2 and V3

- **Factorization Into Asymmetric Convolutions**
- Factorization Into Asymmetric Convolutions
- One 3×1 convolution followed by one 1×3 convolution replaces one 3×3 convolution as follows:



By using 3×3 filter, number of operations = $3 \times 3 = 9$

By using 3×1 and 1×3 filters, number of operations = $3 \times 1 + 1 \times 3 = 6$

Number of operations is reduced by 33%

You may ask why we don't use two 2×2 filters to replace one 3×3 filter?

If we use two 2×2 filters, number of operations = $2 \times 2 \times 2 = 8$

Number of operations is only reduced by 11%

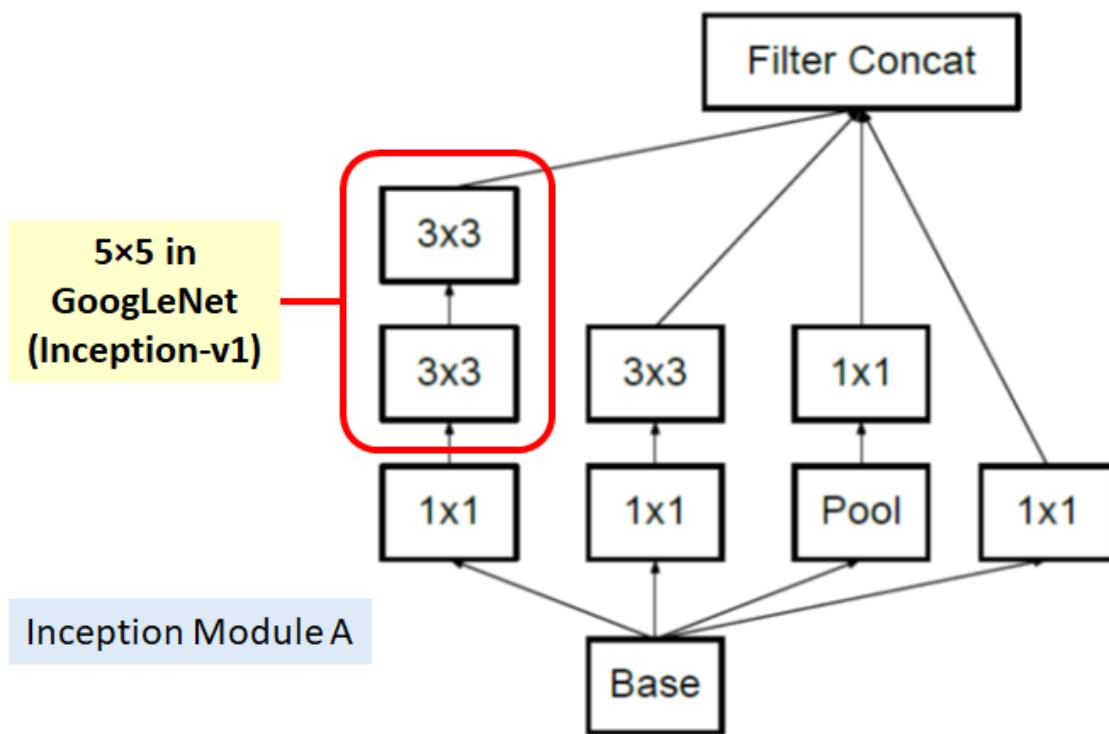
One 3×1 convolution followed by one 1×3 convolution replaces one 3×3 convolution

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

Inception v2 and V3

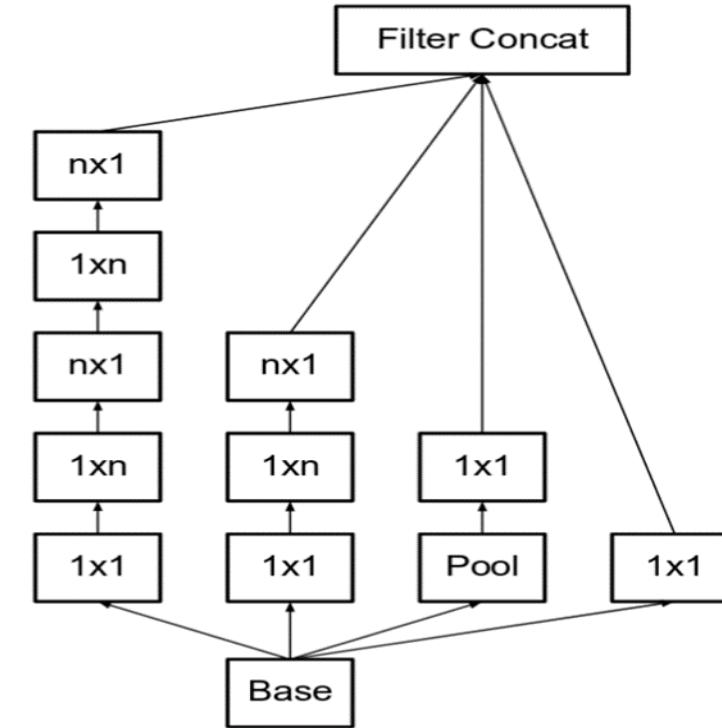
<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

- **Factorization Into Asymmetric Convolutions**
- With this technique, one of the new Inception modules (I call it Inception Module A here) becomes:



5x5 in
GoogLeNet
(Inception-v1)

Inception Module A

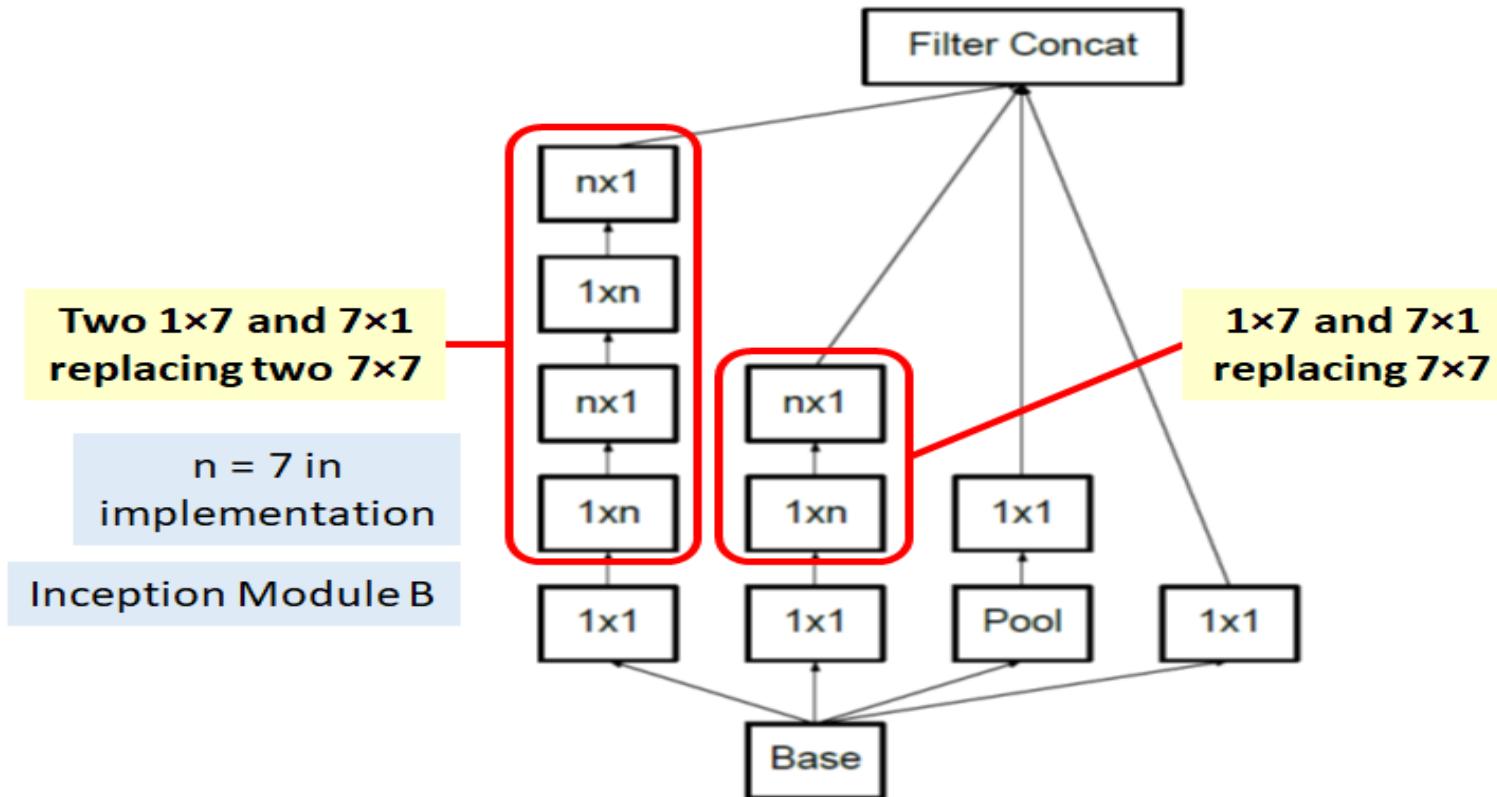


Here, put $n=3$ to obtain the equivalent of the previous image. The left-most 5x5 convolution can be represented as two 3x3 convolutions, which in turn are represented as 1×3 and 3×1 in series. (Source: Inception v2)

Inception v2 and V3

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

- **Factorization Into Asymmetric Convolutions**
- Factorization Into Asymmetric Convolutions
- With this technique, one of the new Inception modules (I call it Inception Module B here) becomes:

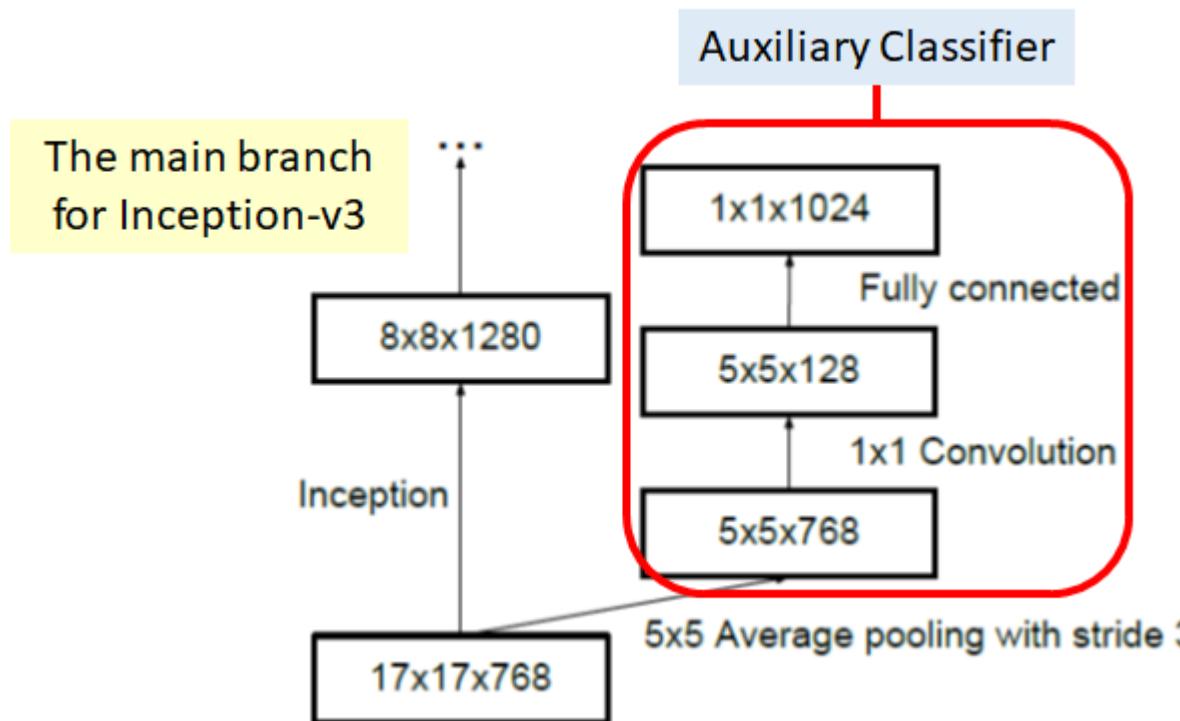


Inception Module B using asymmetric factorization(Inception-V3)

Inception v2 and V3

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

- **Inception v2 and V3**
- Auxiliary Classifiers were already suggested in GoogLeNet / Inception-v1 [4]. There are some modifications in Inception-v3.
- Only 1 auxiliary classifier is used on the top of the last 17×17 layer, instead of using 2 auxiliary classifiers. (The overall architecture would be shown later.)

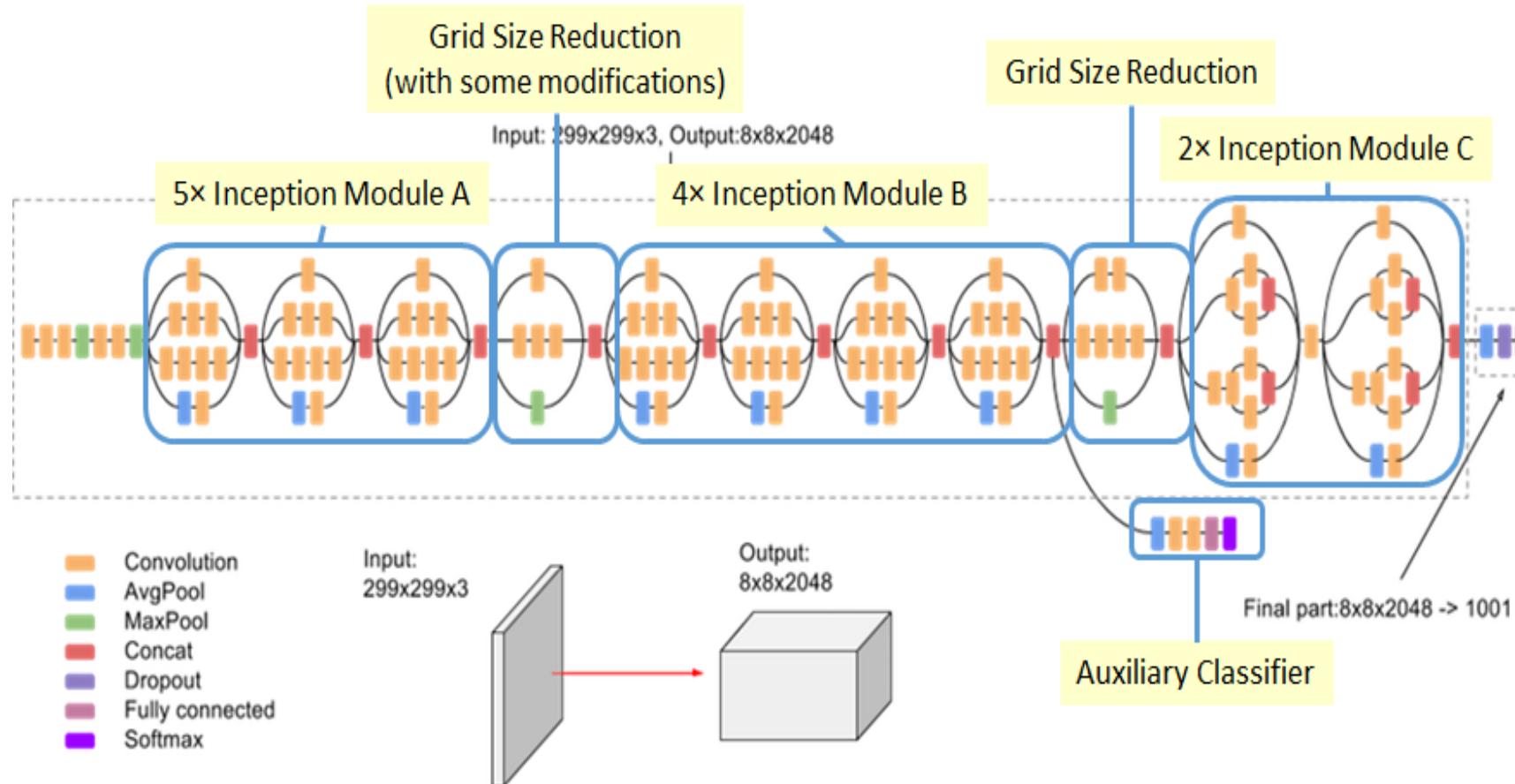


The purpose is also different. In GoogLeNet / Inception-v1 [4], auxiliary classifiers are used for having deeper network. In Inception-v3, auxiliary classifier is used as regularizer. So, actually, in deep learning, the modules are still quite intuitive.

InceptionV3

▪ InceptionV3

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>



With **42** layers deep, the computation cost is only about 2.5 higher than that of GoogLeNet [4], and much more efficient than that of VGGNet [3].

Inception v2

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>

▪ Inception v2

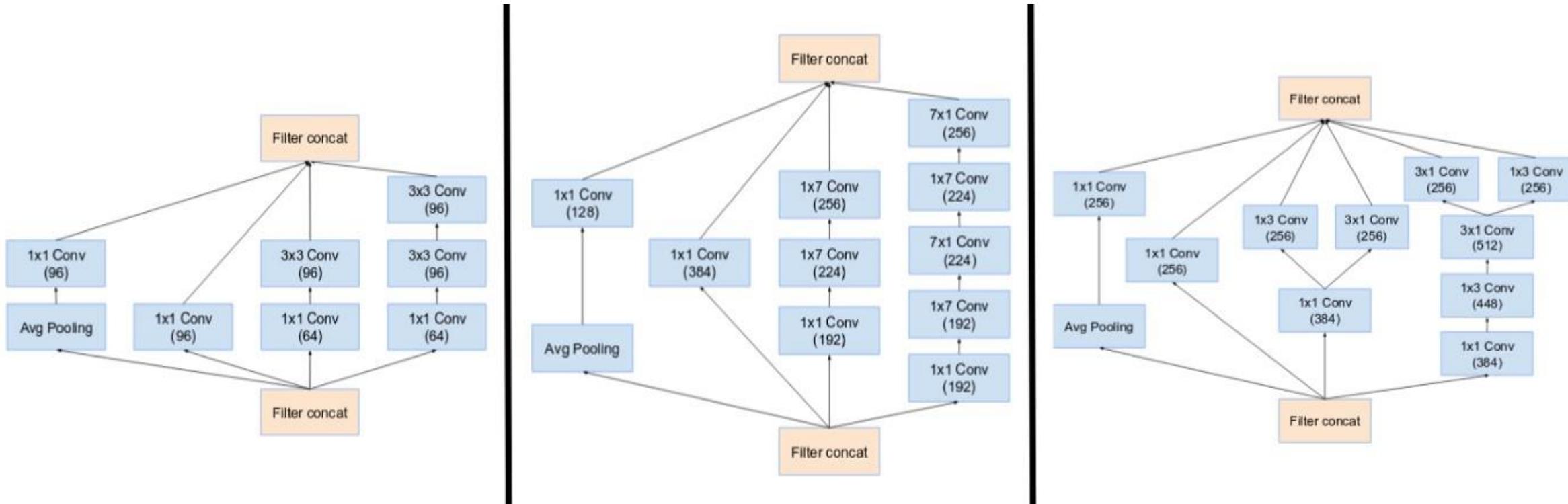
type	patch size/stride or remarks	input size
conv	$3 \times 3 / 2$	$299 \times 299 \times 3$
conv	$3 \times 3 / 1$	$149 \times 149 \times 32$
conv padded	$3 \times 3 / 1$	$147 \times 147 \times 32$
pool	$3 \times 3 / 2$	$147 \times 147 \times 64$
conv	$3 \times 3 / 1$	$73 \times 73 \times 64$
conv	$3 \times 3 / 2$	$71 \times 71 \times 80$
conv	$3 \times 3 / 1$	$35 \times 35 \times 192$
$3 \times$ Inception	As in figure 5	$35 \times 35 \times 288$
$5 \times$ Inception	As in figure 6	$17 \times 17 \times 768$
$2 \times$ Inception	As in figure 7	$8 \times 8 \times 1280$
pool	8×8	$8 \times 8 \times 2048$
linear	logits	$1 \times 1 \times 2048$
softmax	classifier	$1 \times 1 \times 1000$

Here, “figure 5” is module A, “figure 6” is module B and “figure 7” is module C. (Source: Incpetion v2)

Inception-V4

■ Inception-v4

<https://towardsdatascience.com/a-simple-guide-to-the-various-versions-of-the-inception-network-7fc52b863202>

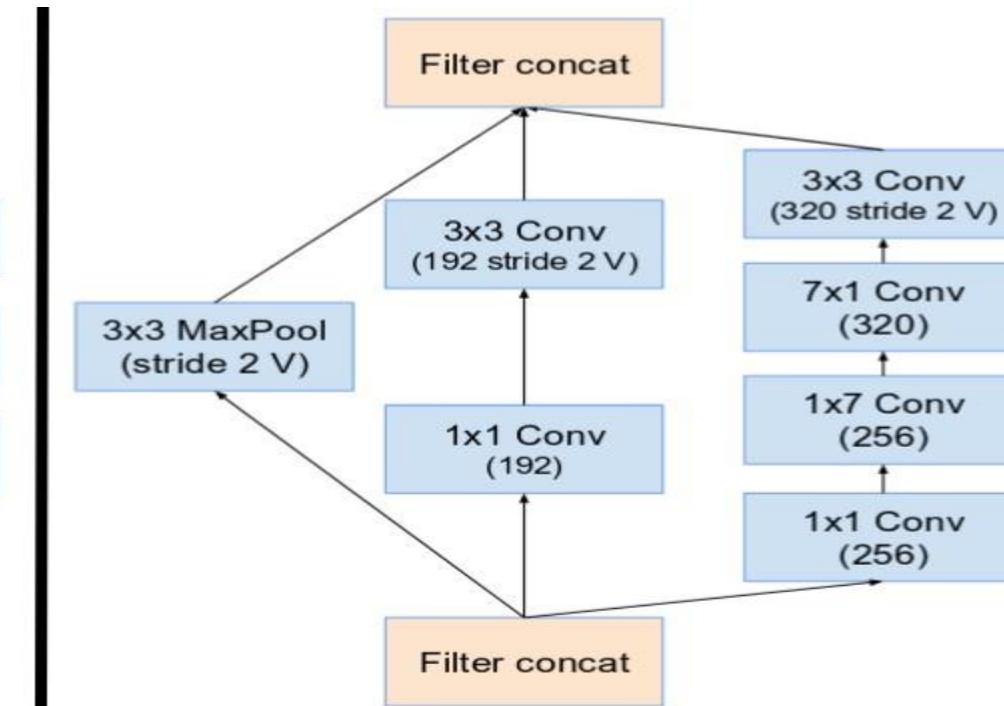
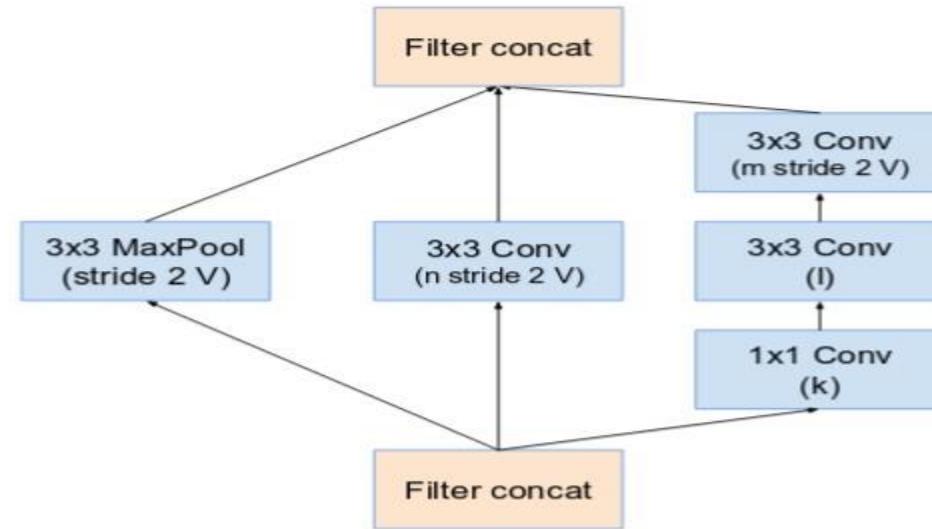


(From left) Inception modules A,B,C used in Inception v4. Note how similar they are to the Inception v2 (or v3) modules. (Source: Inception v4)

Inception-V4

▪ Inception-v4

<https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>



From Left) Reduction Block A (35x35 to 17x17 size reduction) and Reduction Block B (17x17 to 8x8 size reduction). Refer to the paper for the exact hyper-parameter setting (V,l,k). (Source: Inception v4)

Inception-V4

Inception-v4

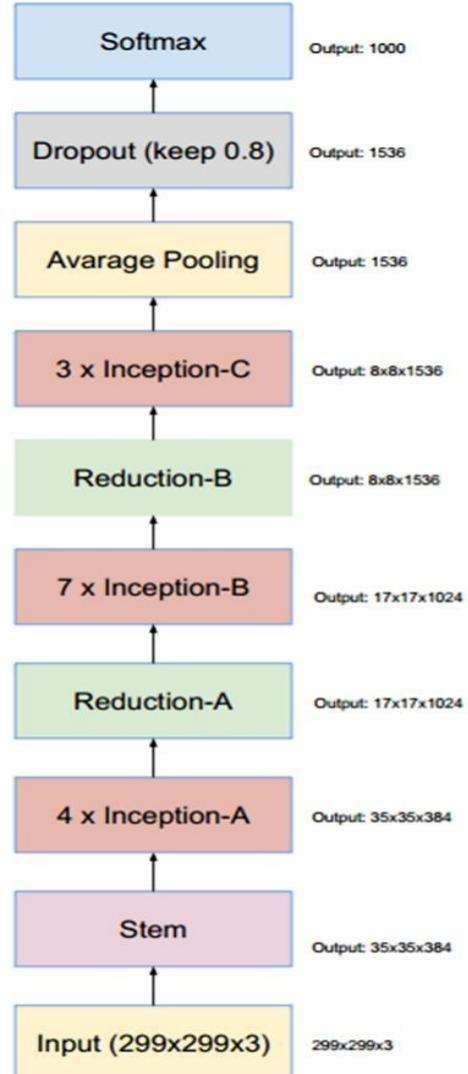


Figure 3. The schema for stem of the pure Inception-v4 and Inception-ResNet-v2 networks. This is the input part of those networks. Cf. Figures 9 and 15.

Figure 5. The schema for 17×17 grid modules of the pure Inception-v4 network. This is the Inception-B block of Figure 9.

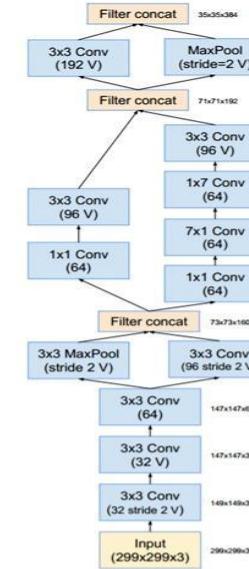


Figure 4. The schema for 35×35 grid modules of the pure Inception-v4 network. This is the Inception-A block of Figure 9.

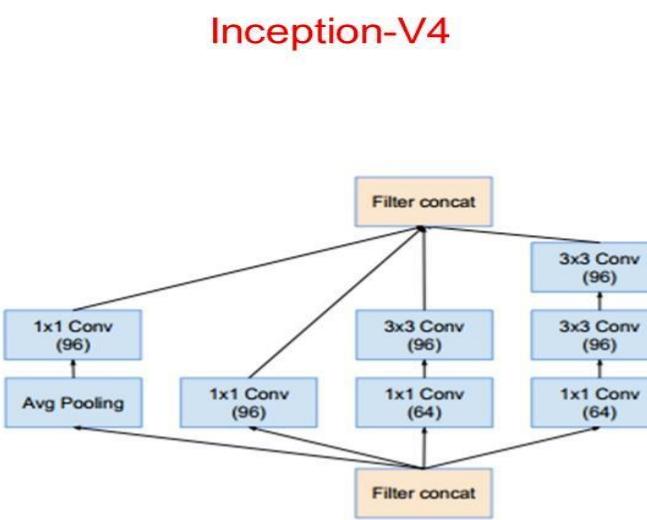


Figure 6. The schema for 8×8 grid modules of the pure Inception-v4 network. This is the Inception-C block of Figure 9.

Inception-V4

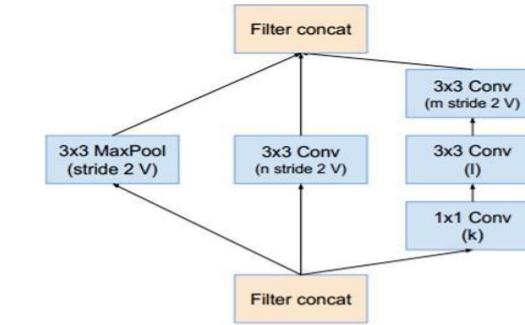


Figure 7. The schema for 35×35 to 17×17 reduction module. Different variants of this blocks (with various number of filters) are used in Figure 9 and 15 in each of the new Inception(-v4, -ResNet-v1, -ResNet-v2) variants presented in this paper. The k, l, m, n numbers represent filter bank sizes which can be looked up in Table II.

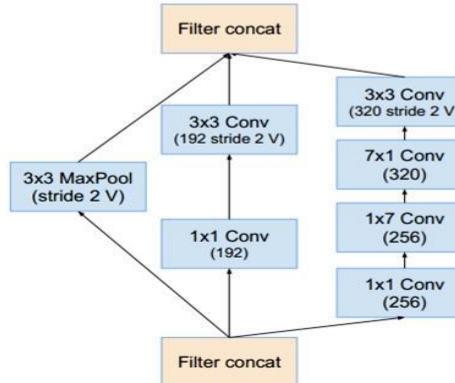
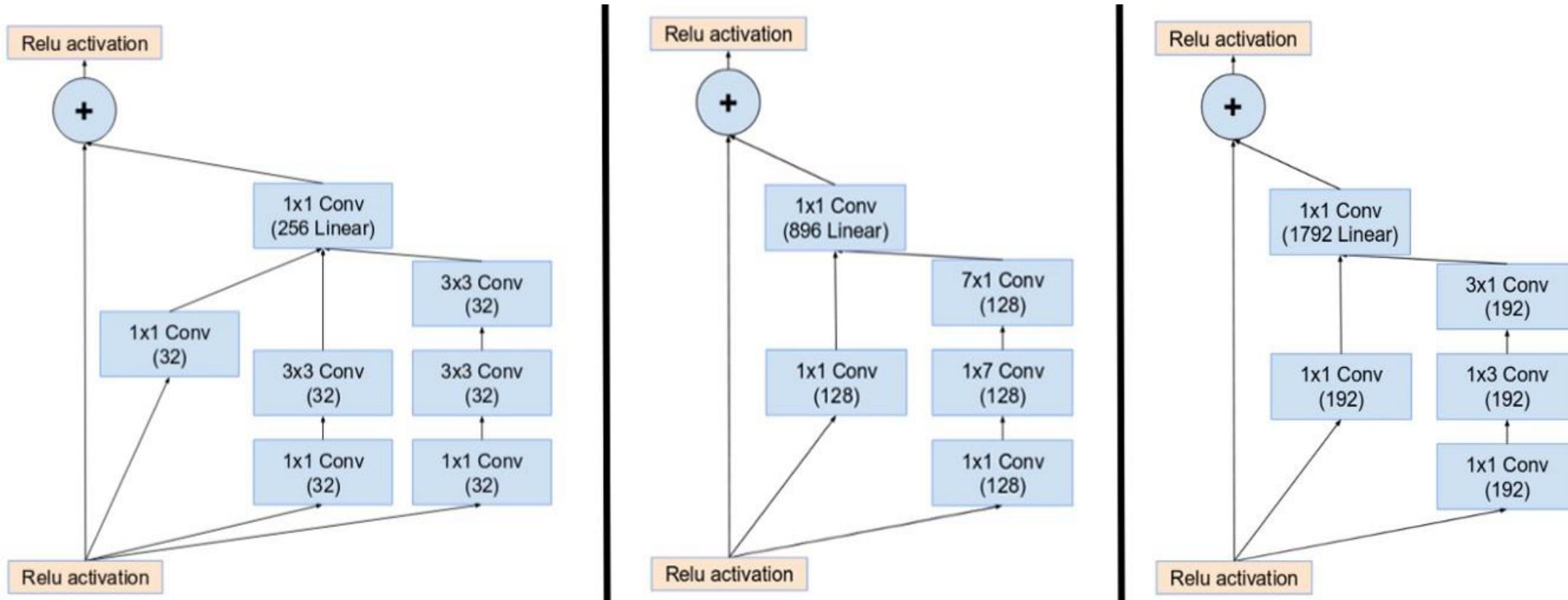


Figure 8. The schema for 17×17 to 8×8 grid-reduction module. This is the reduction module used by the pure Inception-v4 network in Figure 9.

Inception-ResNet-V1

▪ Inception-ResNet-V1

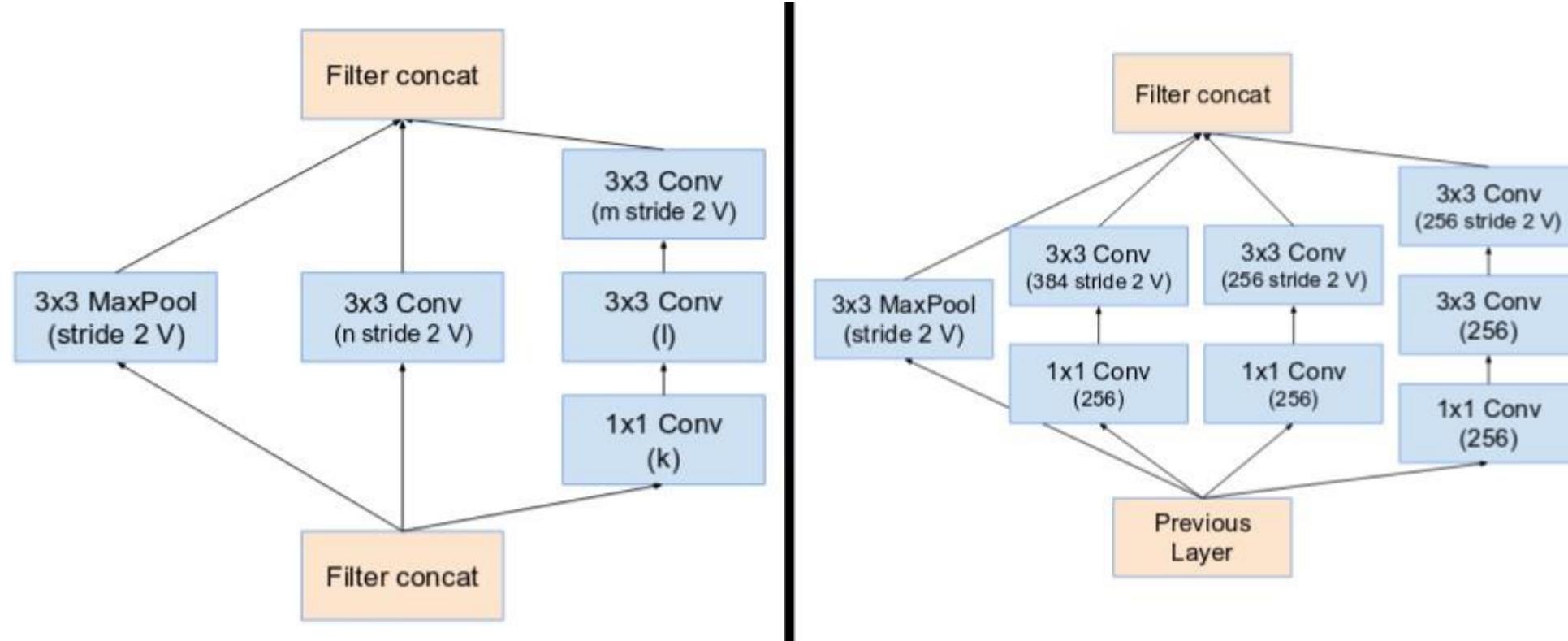


(From left) Inception modules A,B,C in an Inception ResNet. Note how the pooling layer was replaced by the residual connection, and also the additional 1x1 convolution before addition. (Source: Inception v4)

Inception-ResNet-V1

▪ Inception-ResNet-V1

<https://towardsdatascience.com/a-simple-guide-to-the-various-versions-of-the-inception-network-7fc52b863202>



(From Left) Reduction Block A (35x35 to 17x17 size reduction) and Reduction Block B (17x17 to 8x8 size reduction). Refer to the paper for the exact hyper-parameter setting (V,l,k). (Source: Inception v4)

Inception-ResNet-V1

Inception-ResNet-V1

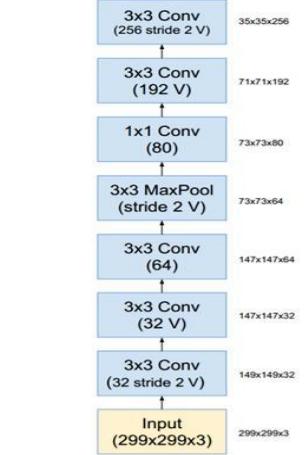
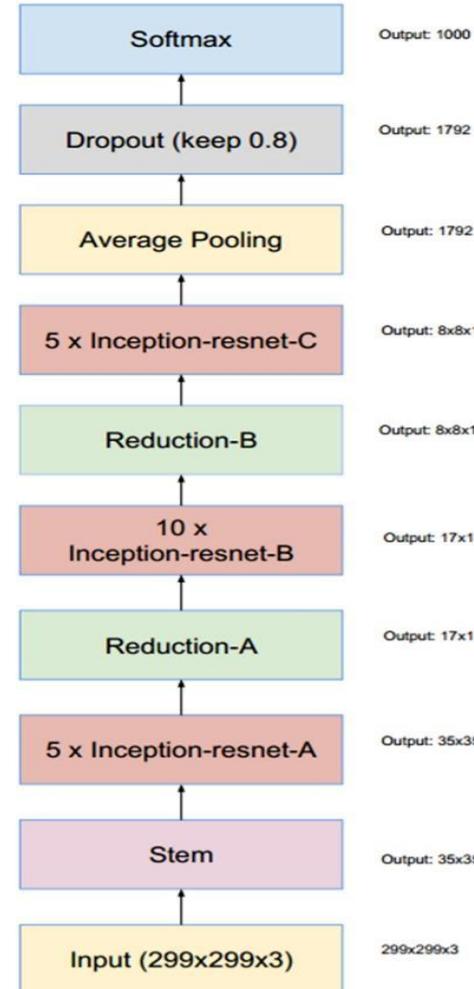


Figure 14. The stem of the Inception-ResNet-v1 network.

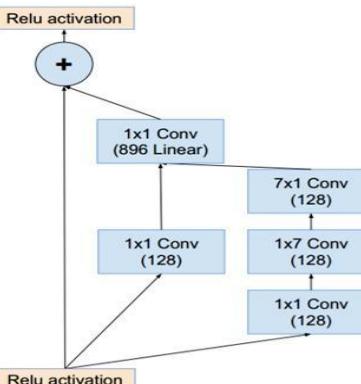


Figure 11. The schema for 17×17 grid (Inception-ResNet-B) module of Inception-ResNet-v1 network.

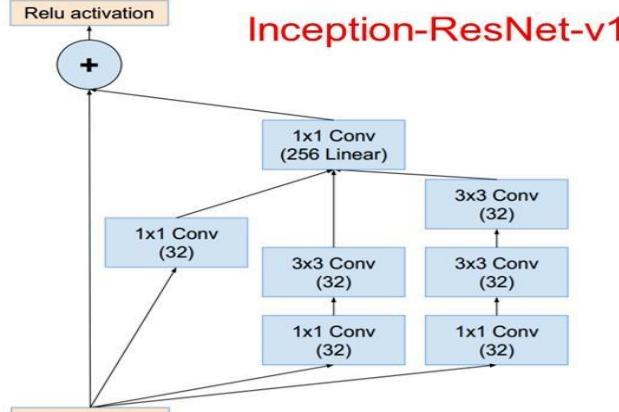


Figure 10. The schema for 35×35 grid (Inception-ResNet-A) module of Inception-ResNet-v1 network.

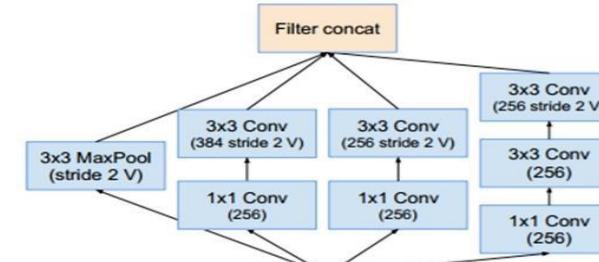


Figure 12. “Reduction-B” 17×17 to 8×8 grid-reduction module. This module used by the smaller Inception-ResNet-v1 network in Figure 15.

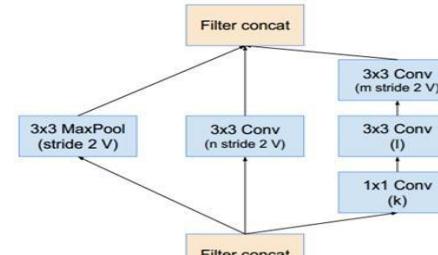


Figure 7. The schema for 35×35 to 17×17 reduction module. Different variants of this blocks (with various number of filters) are used in Figure 9 and 15 in each of the new Inception-v4, -ResNet-v1, -ResNet-v2 variants presented in this paper. The k, l, m, n numbers represent filter bank sizes which can be looked up in Table II.

Network	k	l	m	n
Inception-v4	192	224	256	384
Inception-ResNet-v1	192	192	256	384
Inception-ResNet-v2	256	256	384	384

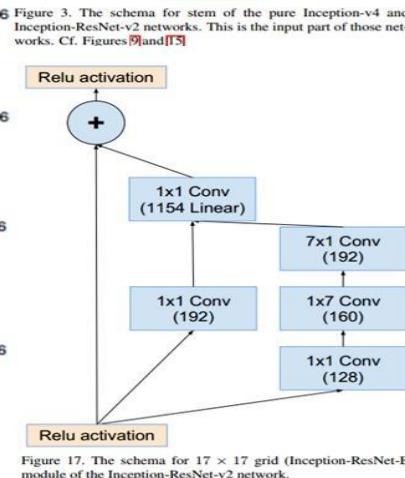
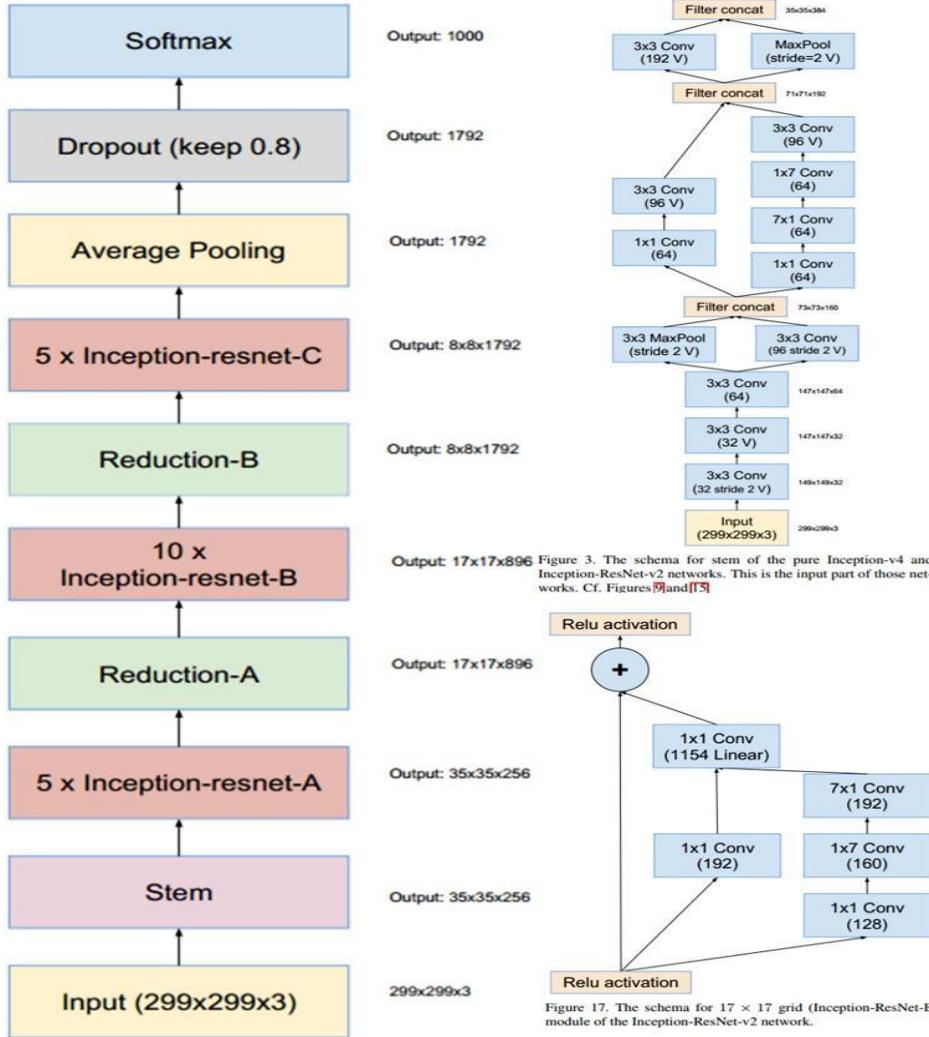
Table 1. The number of filters of the Reduction-A module for the three Inception variants presented in this paper. The four numbers in the columns of the paper parametrize the four convolutions of Figure 7.

<https://towardsdatascience.com/a-simple-guide-to-the-various-versions-of-the-inception-network-7fc52b863202>

<https://github.com/wewin189953602>

Inception-ResNet-V2

Inception-ResNet-V2



<https://towardsdatascience.com/a-simple-guide-to-the-various-versions-of-the-inception-network-7fc52b863202>

Inception-ResNet-v2

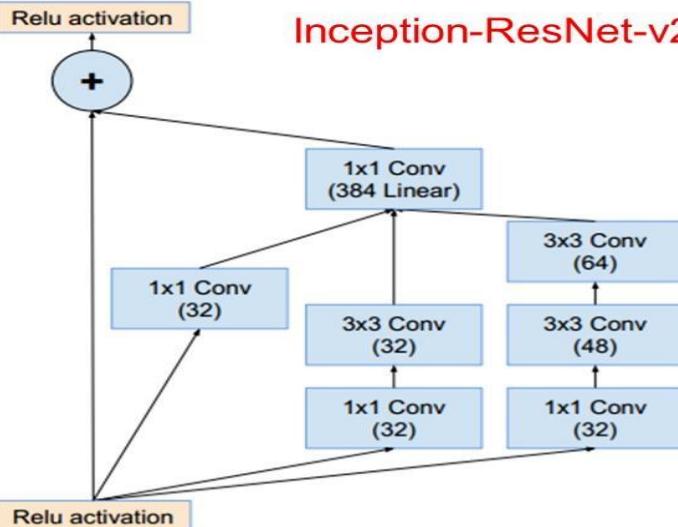


Figure 3. The schema for stem of the pure Inception-v4 and Inception-ResNet-v2 networks. This is the input part of those networks. Cf. Figures 9 and 15.

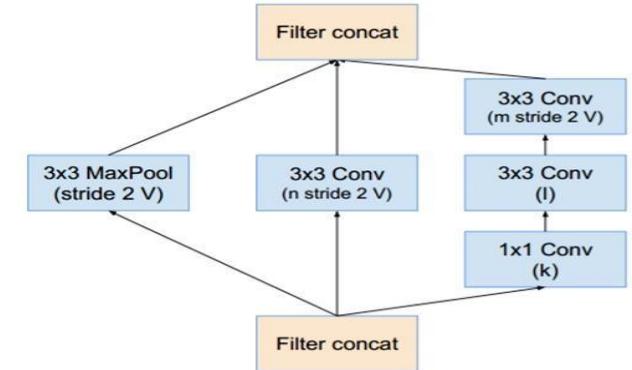


Figure 7. The schema for 35×35 to 17×17 reduction module. Different variants of this blocks (with various number of filters) are used in Figure 9 and 15 in each of the new Inception(-v4, -ResNet-v1, -ResNet-v2) variants presented in this paper. The k, l, m, n numbers represent filter bank sizes which can be looked up in Table II.

Network	k	l	m	n
Inception-v4	192	224	256	384
Inception-ResNet-v1	192	192	256	384
Inception-ResNet-v2	256	256	384	384

Table I. The number of filters of the Reduction-A module for the three Inception variants presented in this paper. The four numbers in the columns of the paper parametrize the four convolutions of Figure 7.

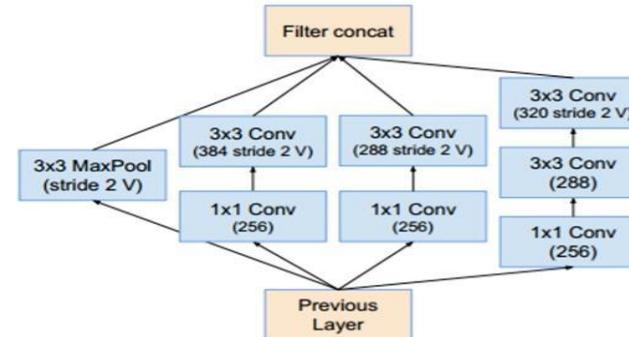


Figure 16. The schema for 35×35 grid (Inception-ResNet-A) module of the Inception-ResNet-v2 network.

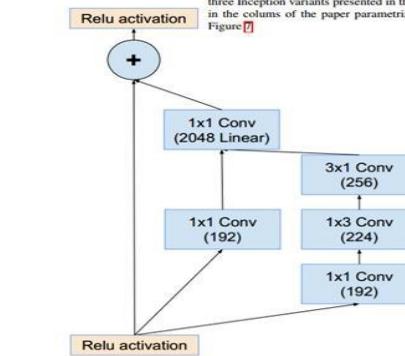


Figure 18. The schema for 17×17 to 8×8 grid-reduction module. Reduction-B module used by the wider Inception-ResNet-v1 network in Figure 15.

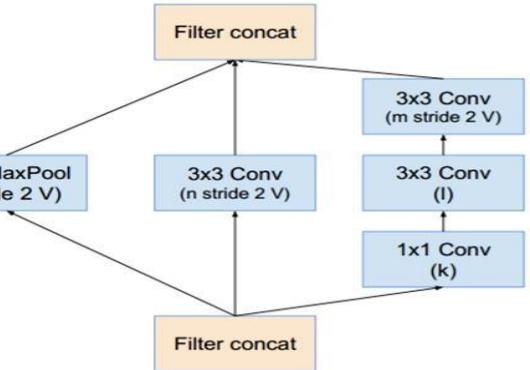


Figure 7. The schema for 35×35 to 17×17 reduction module. Different variants of this blocks (with various number of filters) are used in Figure 9 and 15 in each of the new Inception(-v4, -ResNet-v1, -ResNet-v2) variants presented in this paper. The k, l, m, n numbers represent filter bank sizes which can be looked up in Table II.

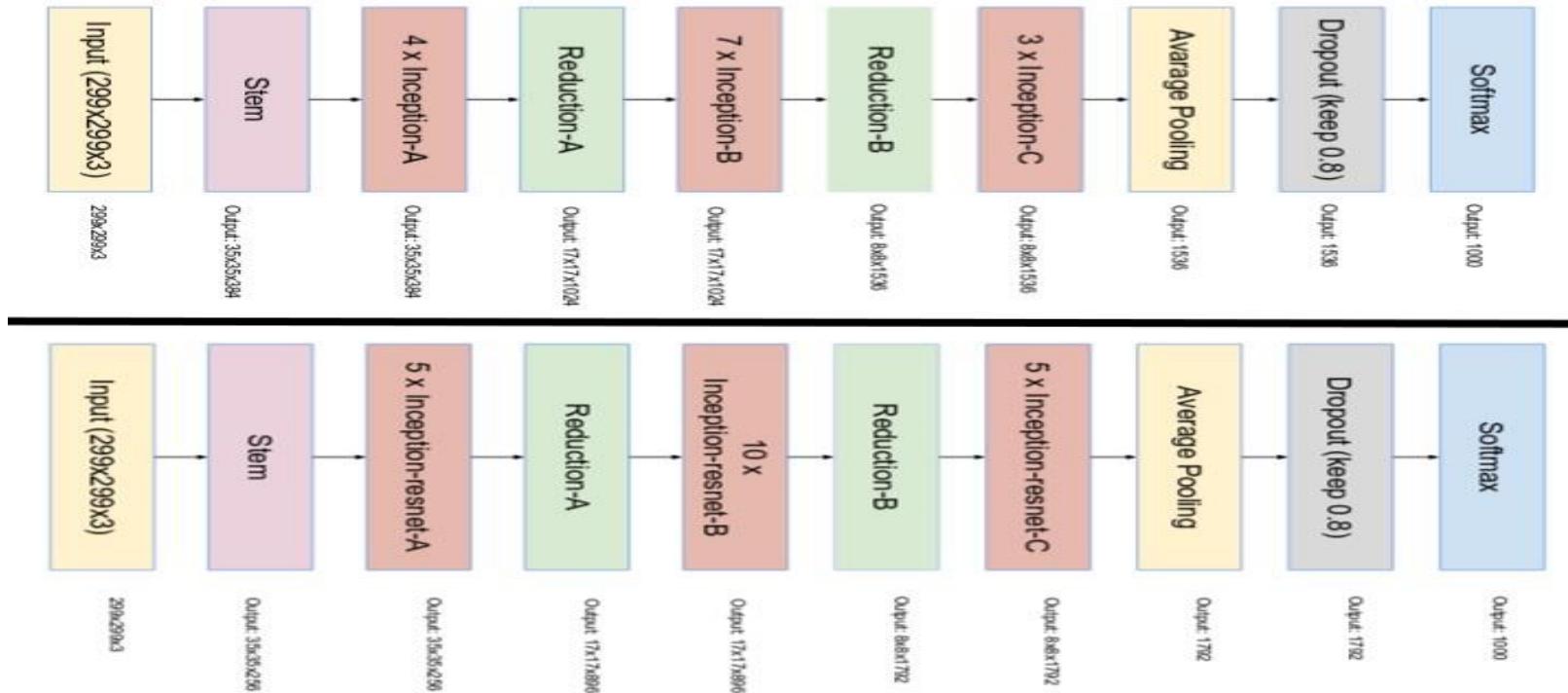
Network	k	l	m	n
Inception-v4	192	224	256	384
Inception-ResNet-v1	192	192	256	384
Inception-ResNet-v2	256	256	384	384

Table I. The number of filters of the Reduction-A module for the three Inception variants presented in this paper. The four numbers in the columns of the paper parametrize the four convolutions of Figure 7.

Inception-ResNet-V1 and V2

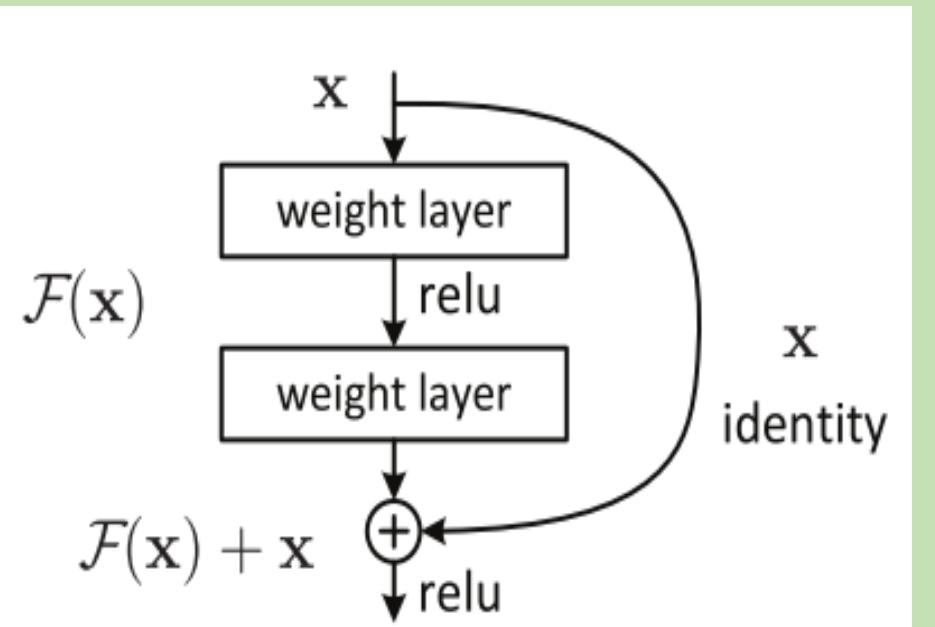
■ Inception-ResNet-V1 and V2

<https://towardsdatascience.com/a-simple-guide-to-the-versions-of-the-inception-network-7fc52b863202>

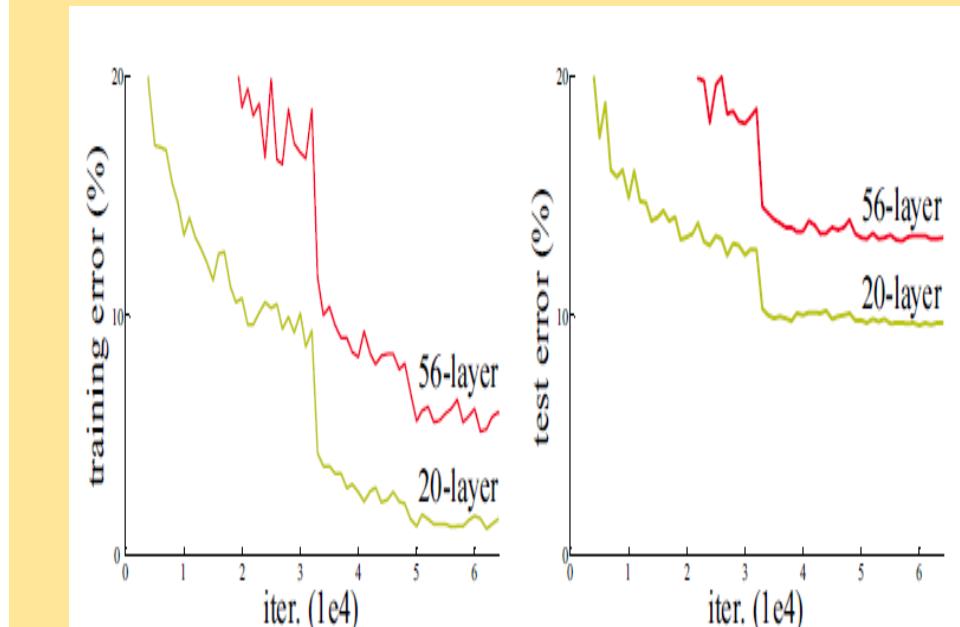


The top image is the layout of Inception v4. The bottom image is the layout of Inception-ResNet. (Source: Inception v4)

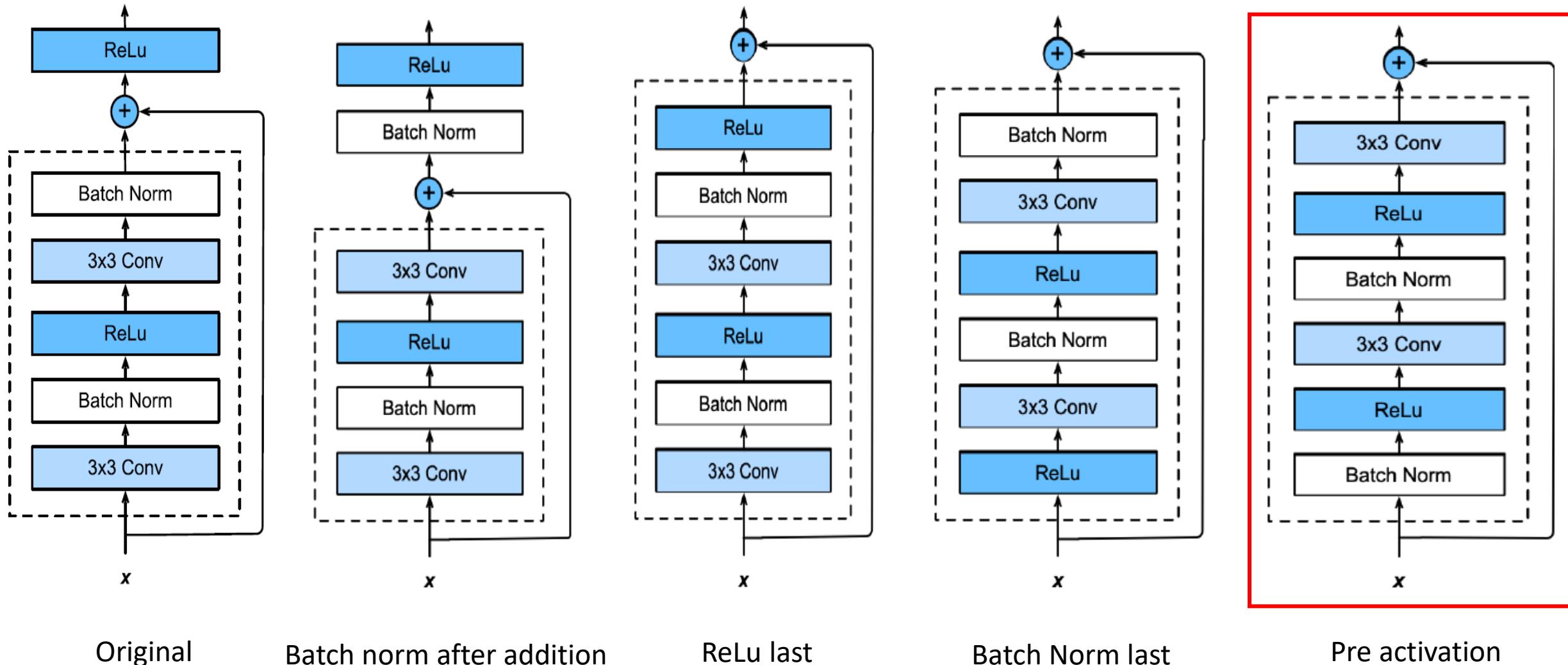
Residual Module



Comparison of Deeper Layers



Design Configuration



Original

Batch norm after addition

ReLU last

Batch Norm last

Pre activation

ResNet Model

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
				3×3 max pool, stride 2		
conv2_x	56×56	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
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

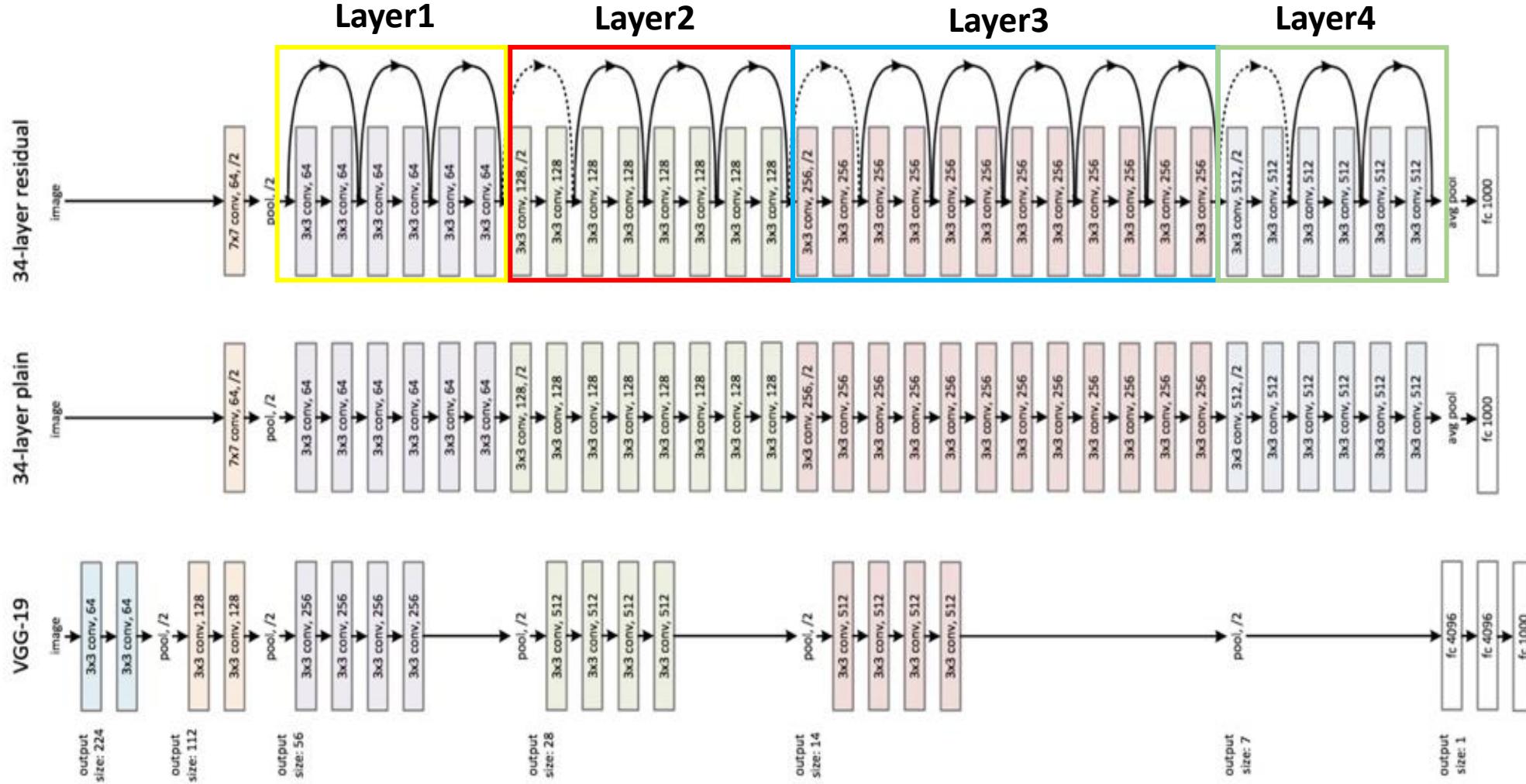
ResNet 18, 34

- Two 3x3 Conv layers
- Multiple repeated Residual Blocks
- **4 Major Layers**

ResNet50,101,152

- 3 Conv Layers
- Multiple repeated Residual Blocks
- **4 Major Layers**

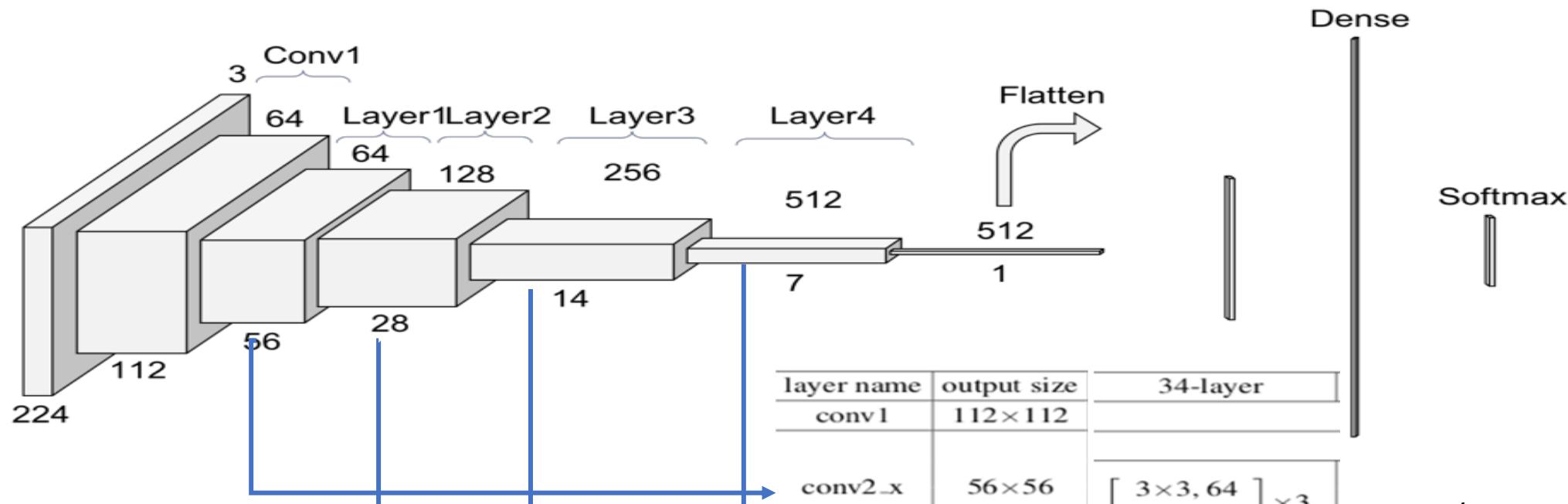
ResNet 34



<https://arxiv.org/pdf/1512.03385.pdf>

Figure 3. Example network architectures for ImageNet. **Left:** the VGG-19 model [41] (19.6 billion FLOPs) as a reference. **Middle:** a plain network with 34 parameter layers (3.6 billion FLOPs). **Right:** a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. **Table 1** shows more details and other variants.

ResNet 34

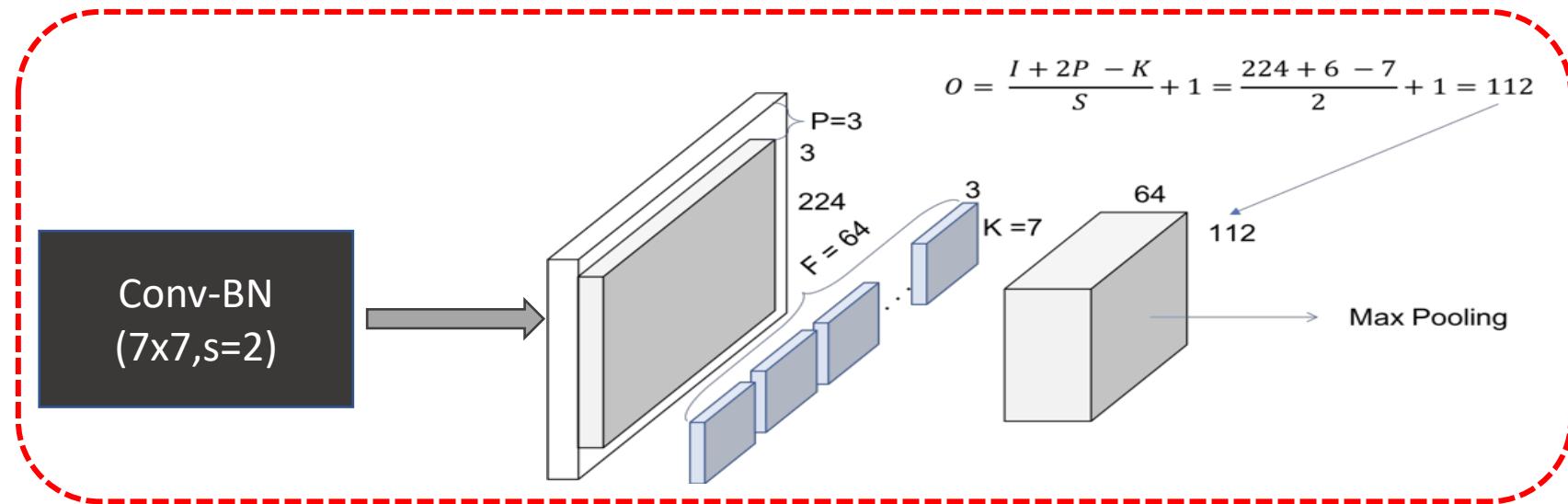
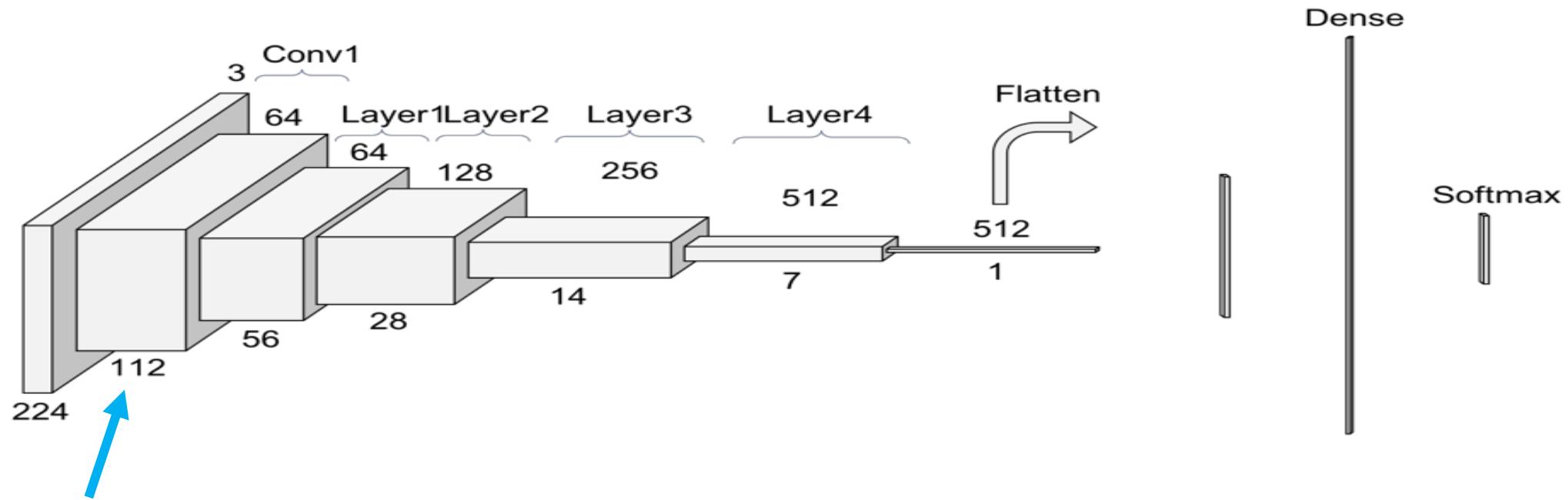


layer name	output size	34-layer
conv1	112×112	
conv2_x	56×56	$[3 \times 3, 64] \times 3$
conv3_x	28×28	$[3 \times 3, 128] \times 4$
conv4_x	14×14	$[3 \times 3, 256] \times 6$
conv5_x	7×7	$[3 \times 3, 512] \times 3$
	1×1	ave
FLOPs		3.6×10^9

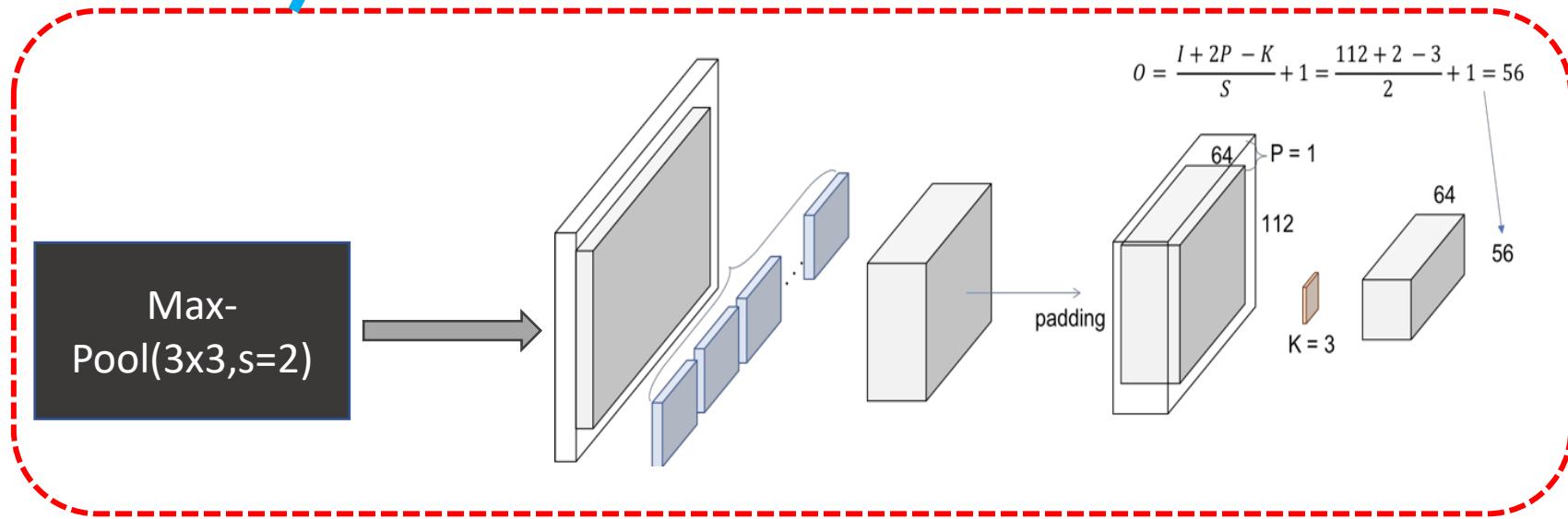
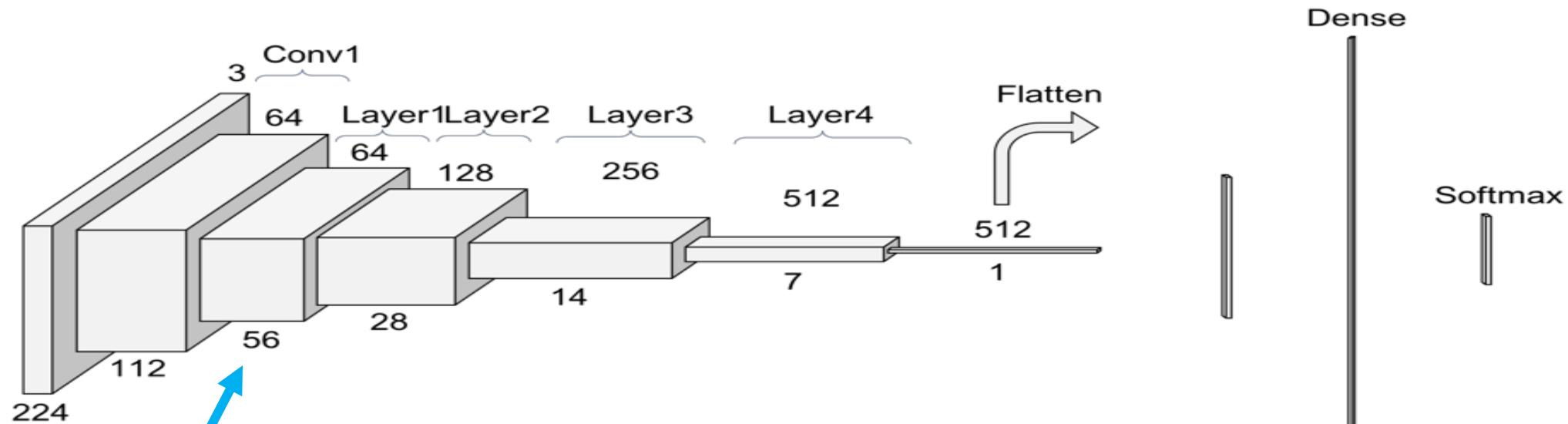
- Layer1(3 Residual Blocks)**
- Layer2 (4 Residual Blocks)**
- Layer3 (6 Residual Blocks)**
- Layer4 (3 Residual Blocks)**

Every Residual Block consisted of 2*(3x3 Conv Layers using skip connection) with different depth or feature maps (**64, 128, 256, 512**)

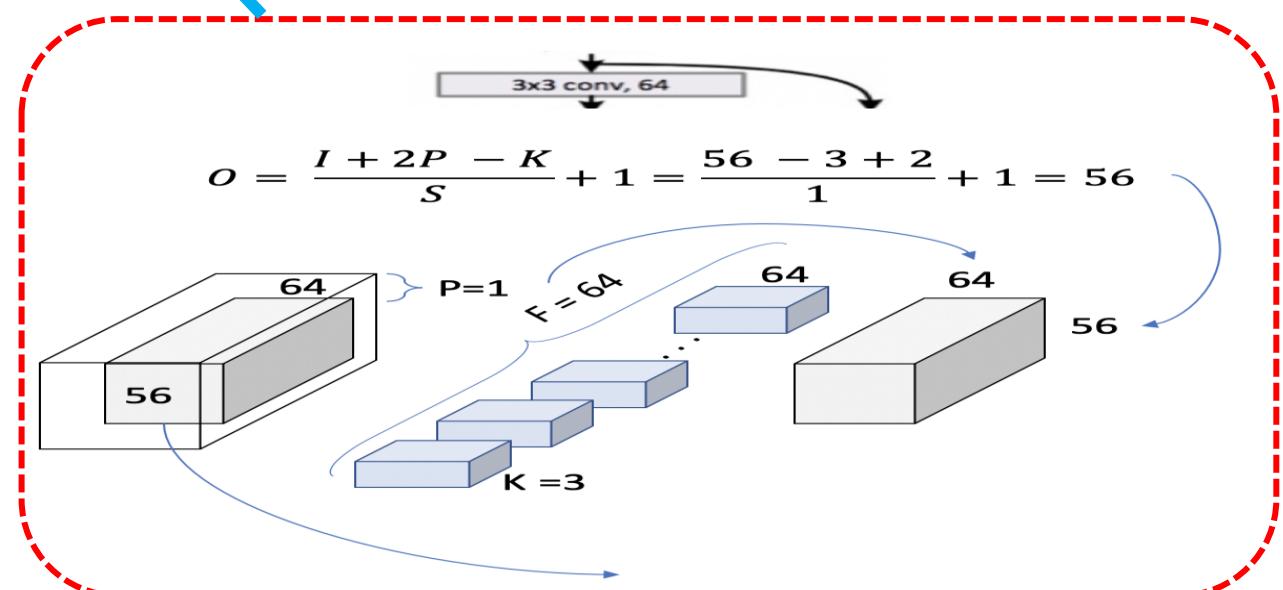
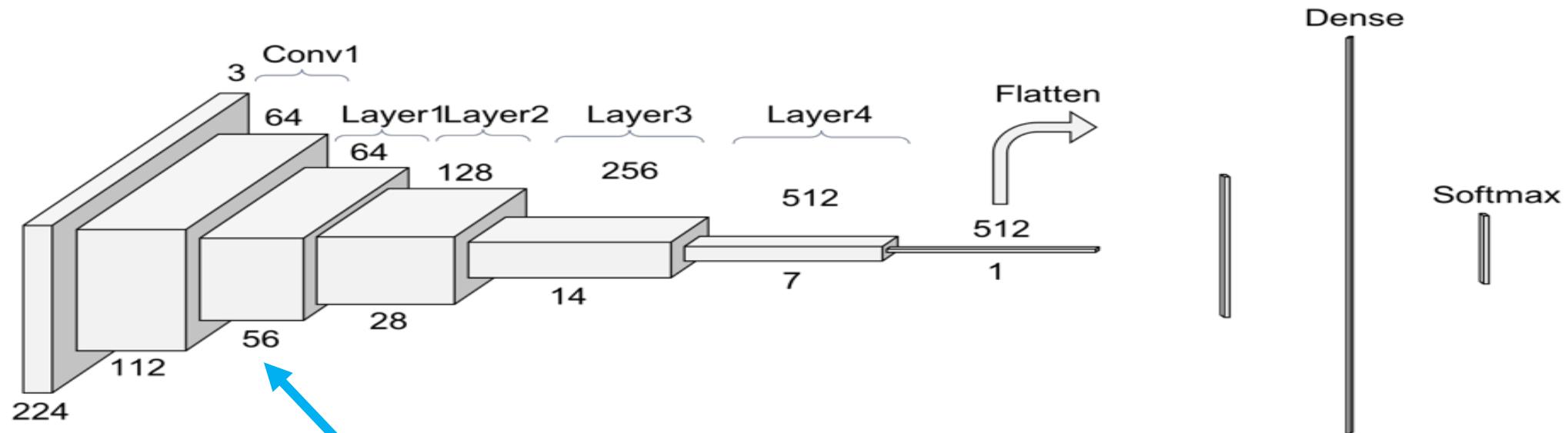
ResNet 34



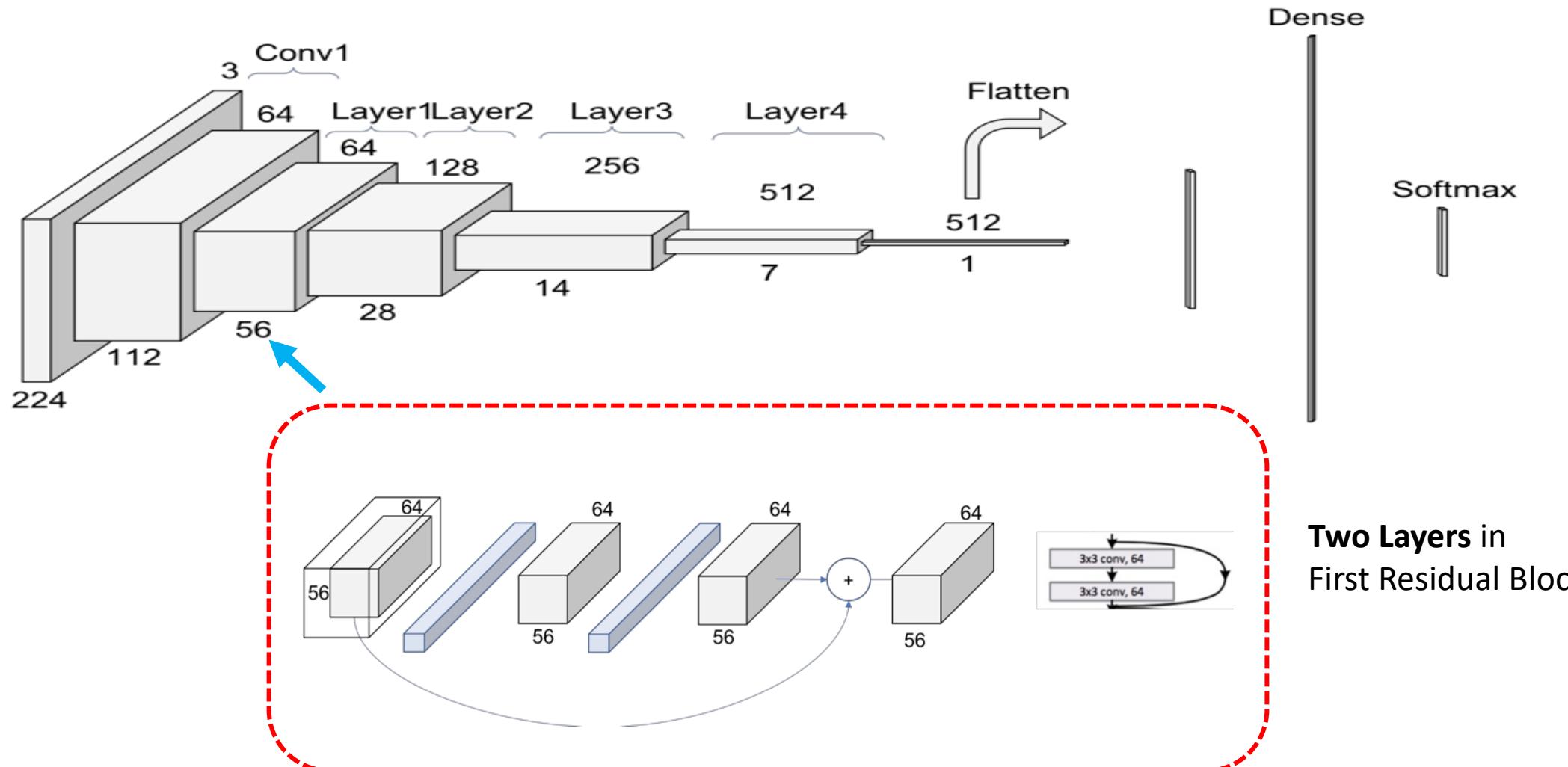
ResNet 34



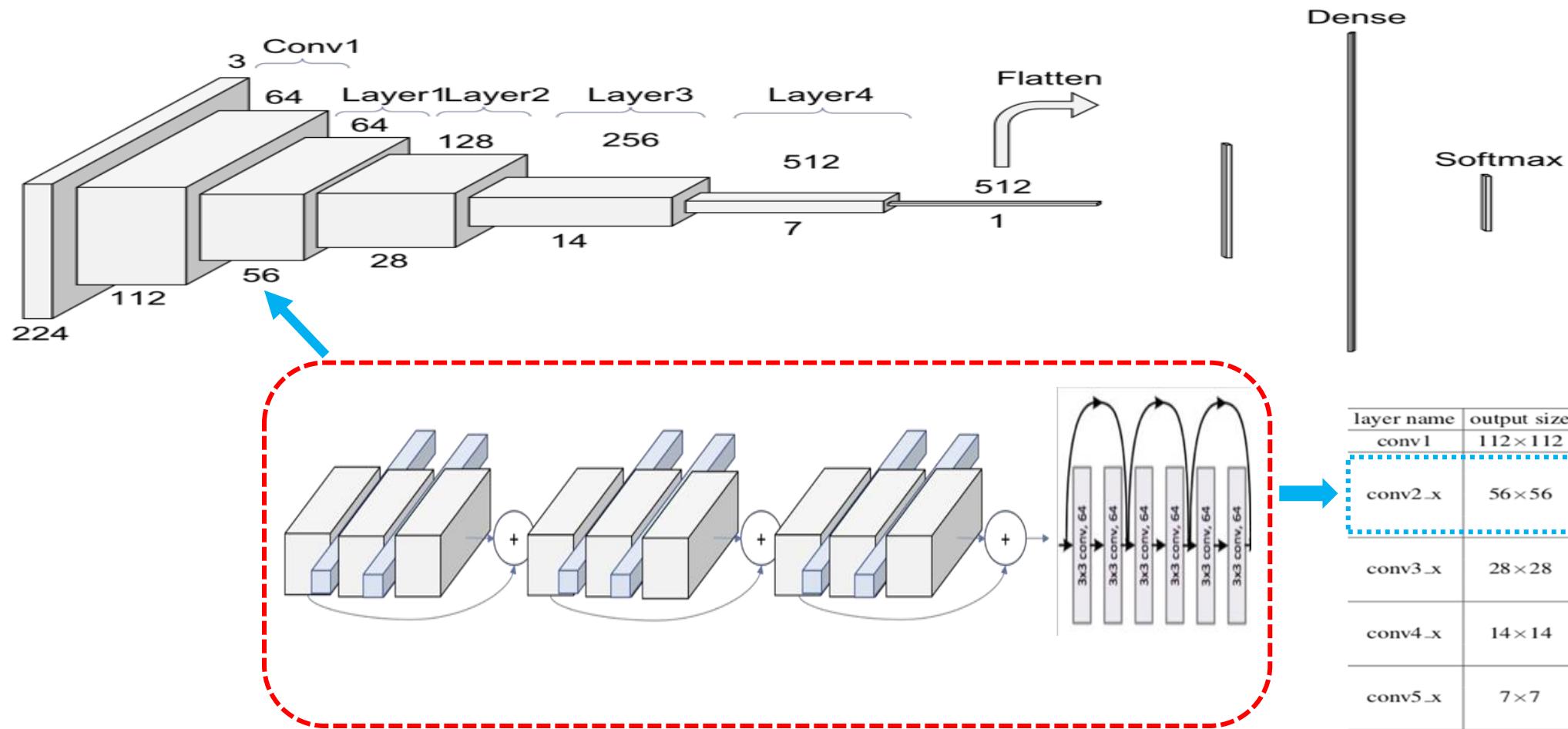
ResNet 34



ResNet 34



ResNet 34

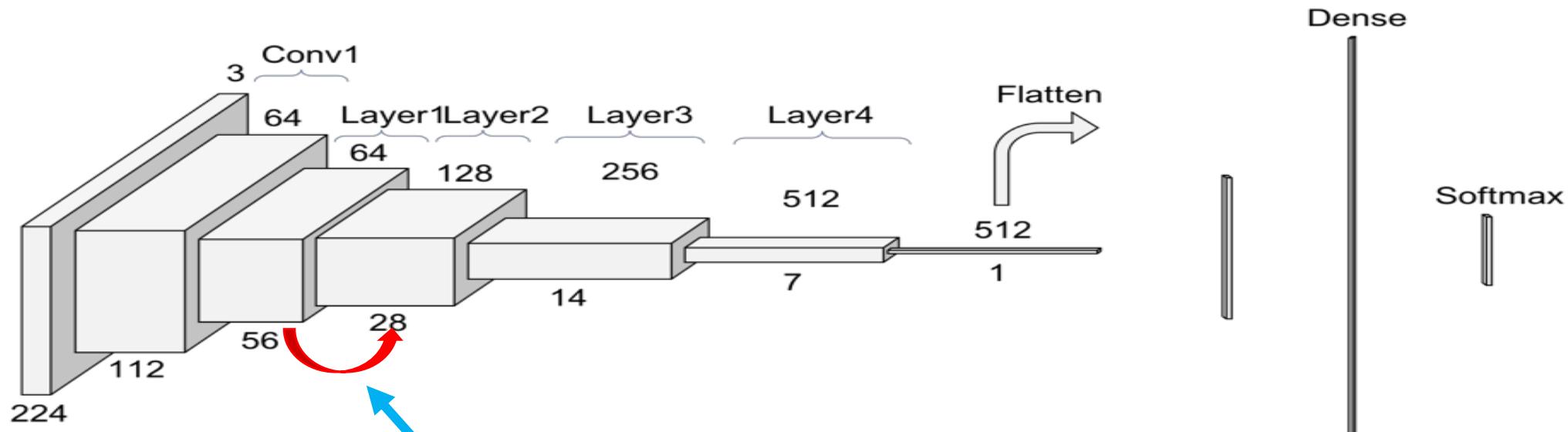


layer name	output size	34-layer
conv1	112×112	
conv2_x	56×56	$[3 \times 3, 64] \times 3$ $[3 \times 3, 64] \times 3$
conv3_x	28×28	$[3 \times 3, 128] \times 4$ $[3 \times 3, 128] \times 4$
conv4_x	14×14	$[3 \times 3, 256] \times 6$ $[3 \times 3, 256] \times 6$
conv5_x	7×7	$[3 \times 3, 512] \times 3$ $[3 \times 3, 512] \times 3$
		ave
FLOPs		3.6×10^9

First Layer (Layer1)

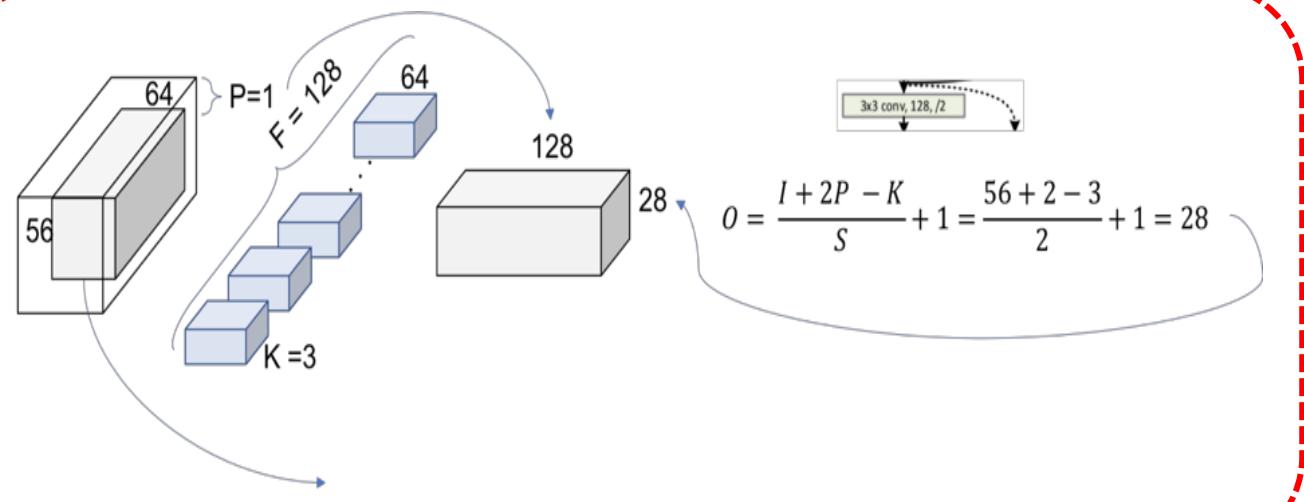
3 residual blocks (every Residual block consisted of two 3x3 Conv layer with skip connection)

ResNet 34

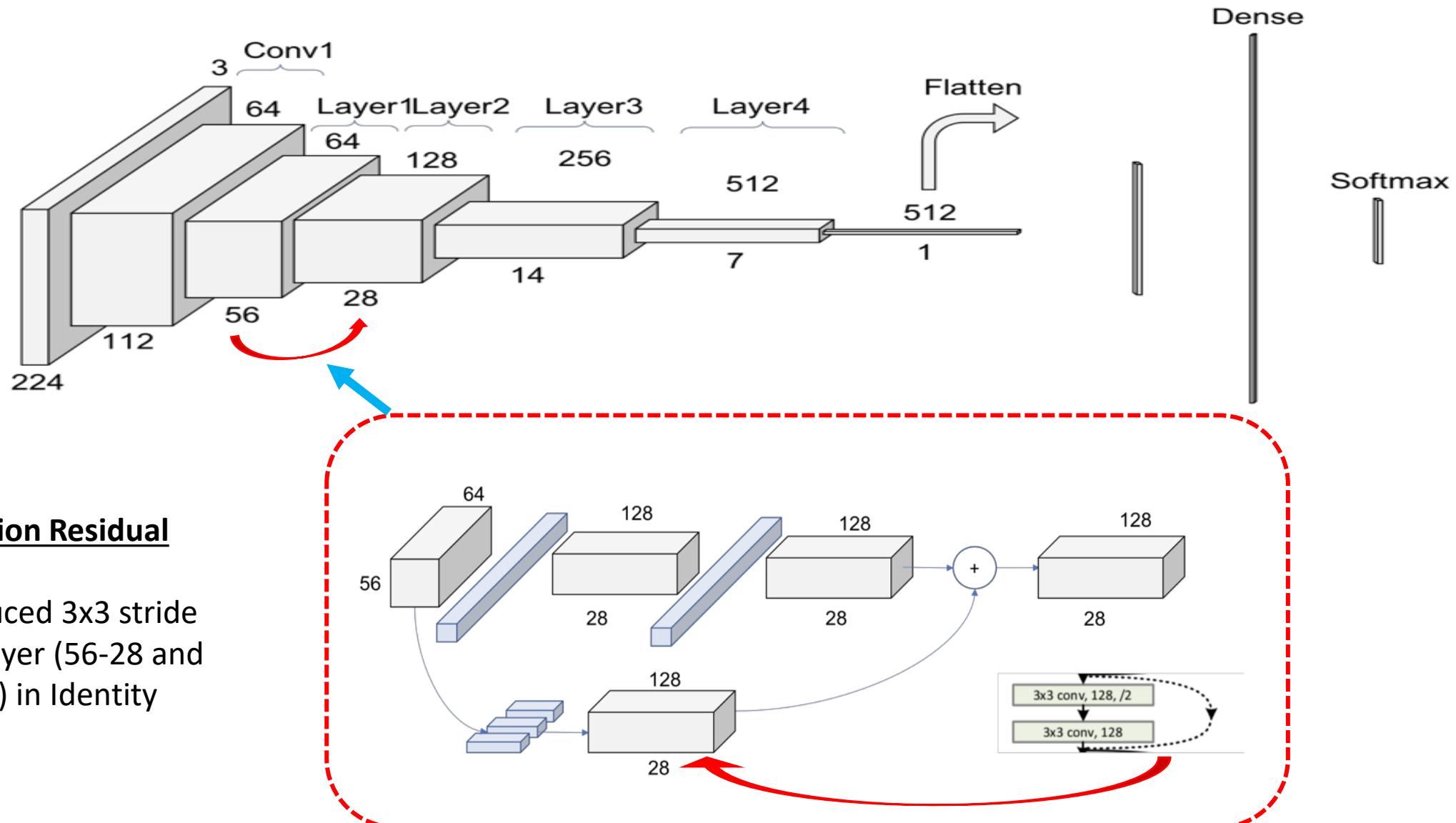


Stride Conv Layer

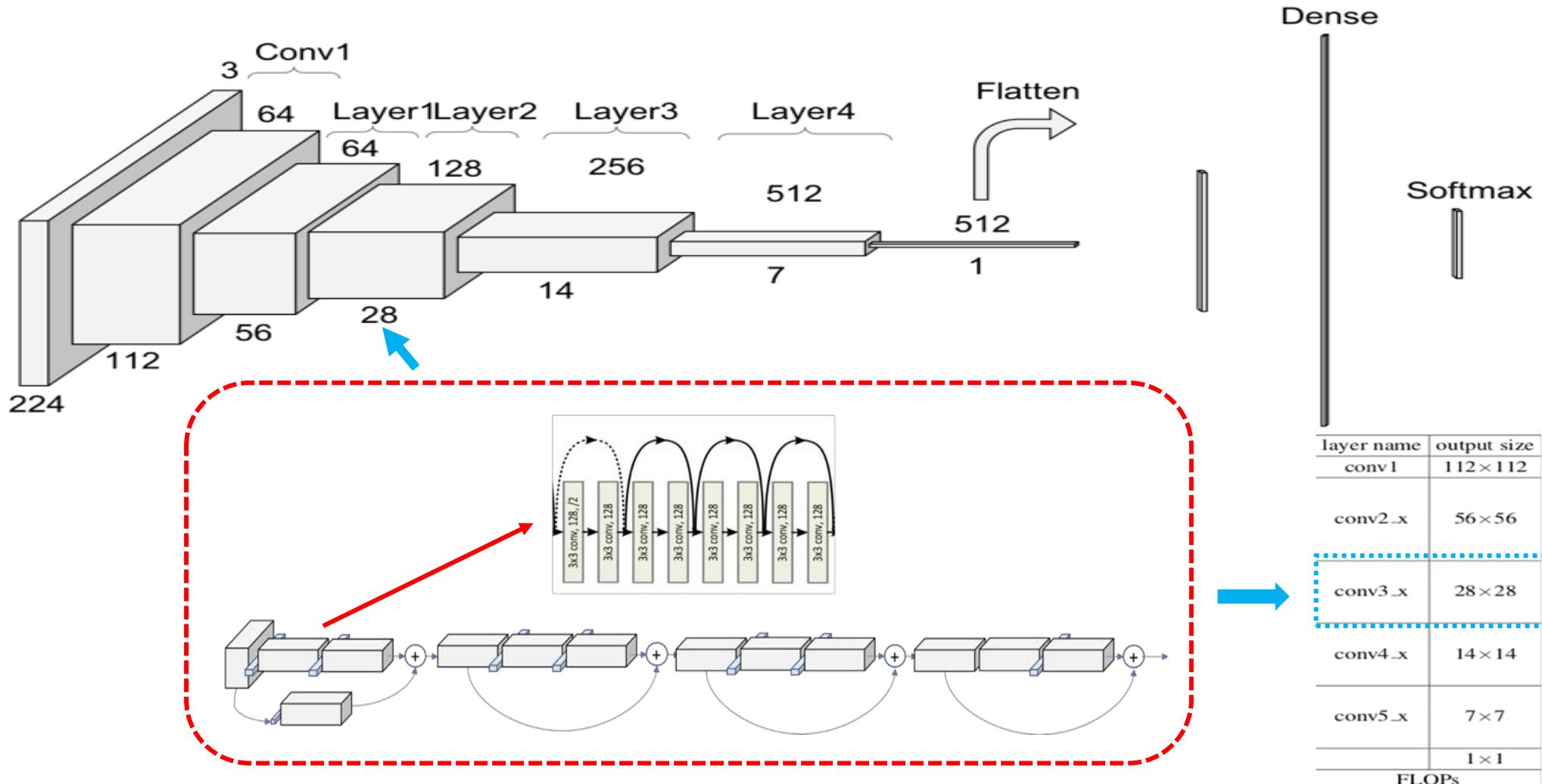
- 3x3 Conv layer with stride=2, P=1
- Down-sample the spatial dimension From 56-28,
- Depth from 64-128



ResNet 34



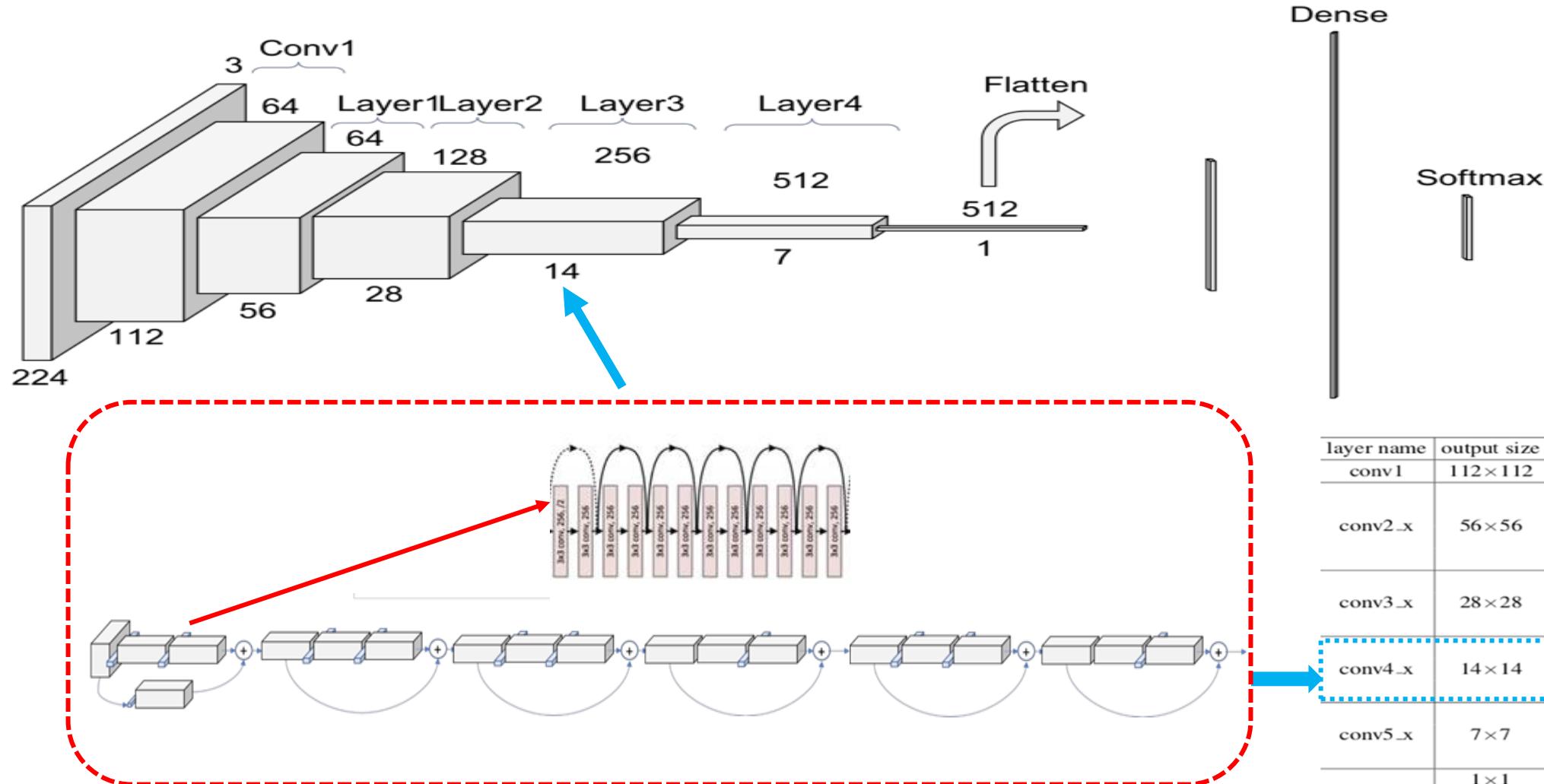
ResNet 34



Layer2

(1) Transition Residual block+(3)* Residual Blocks

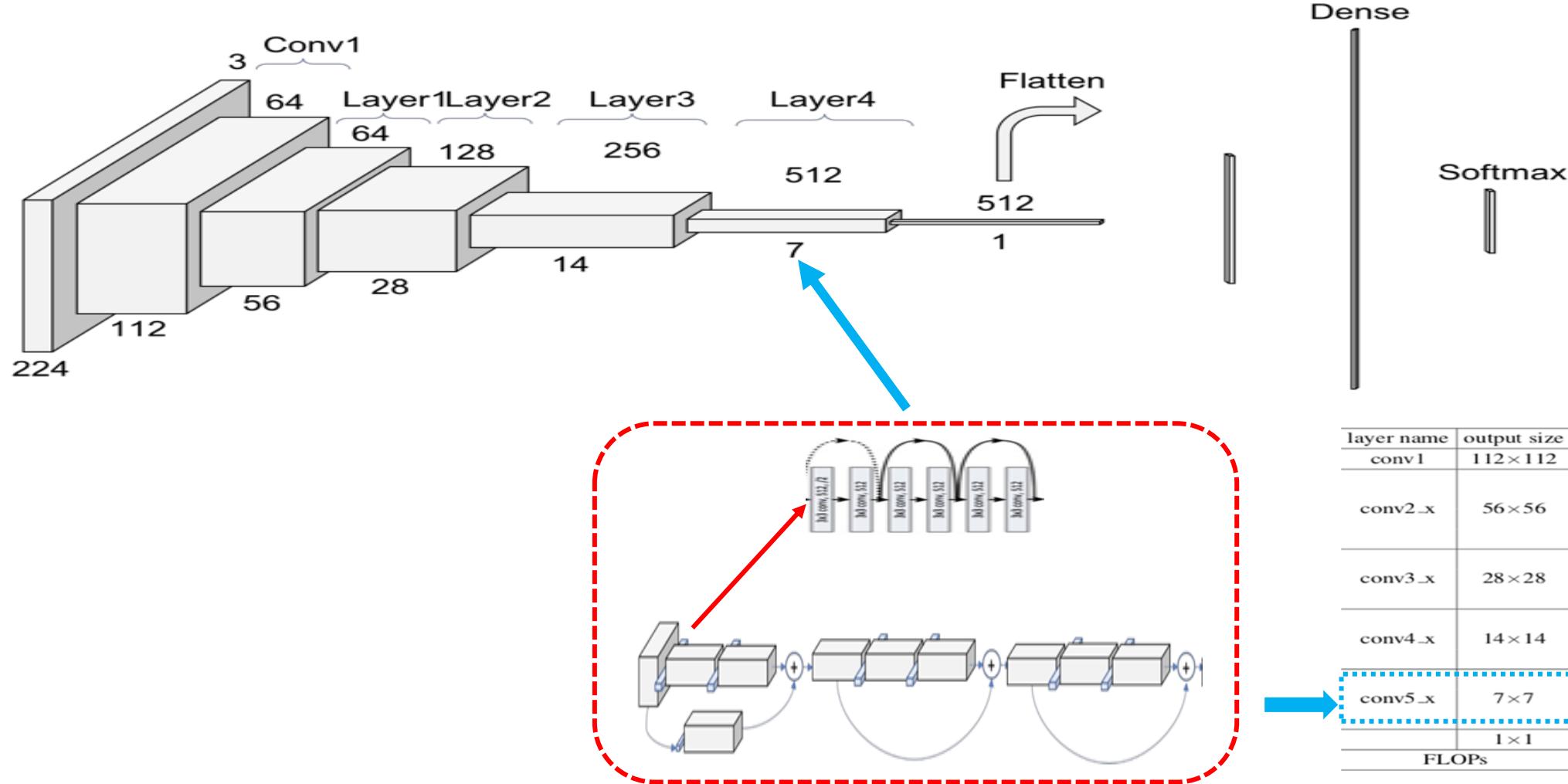
ResNet 34



Layer3

(1) Transition Residual block+(5)* Residual Blocks

ResNet 34

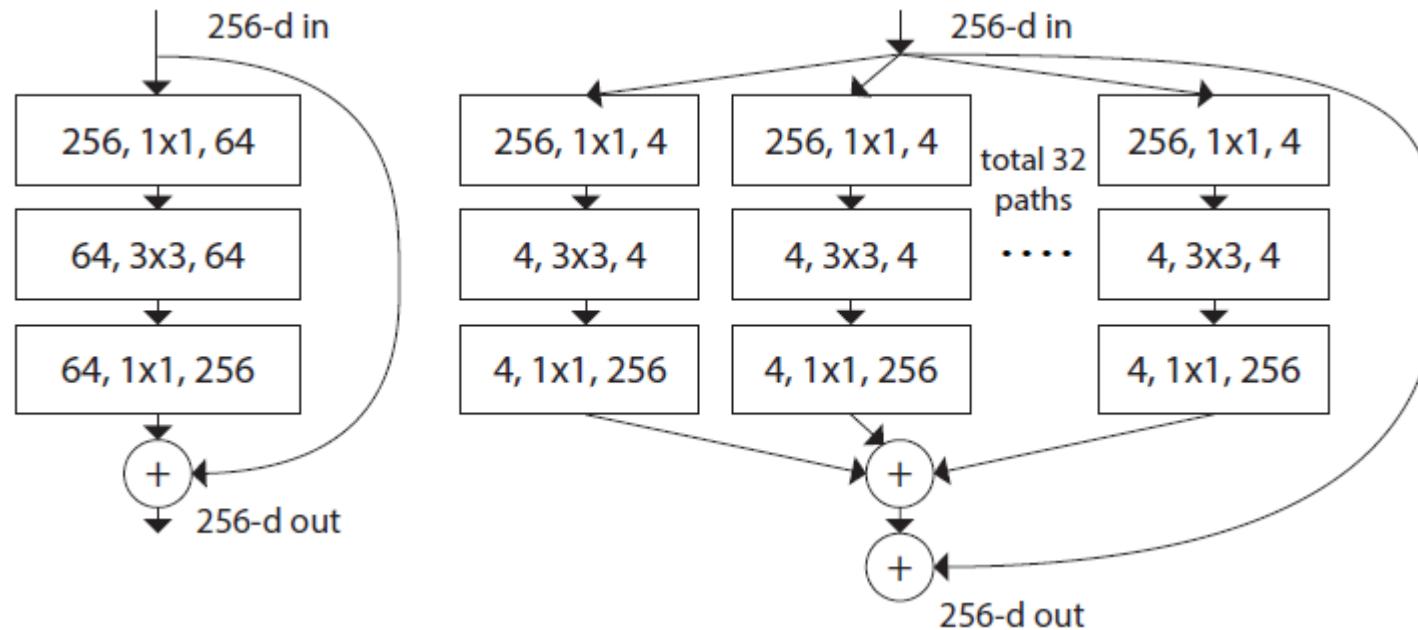


Layer4

(1)*Transition Residual block+(2)* Residual Blocks

ResNeXt

- **ResNeXt**
- ResNeXt — 1st Runner Up in ILSVRC 2016 (Image Classification)



Residual Block in ResNet (Left), A Block of ResNeXt with Cardinality = 32 (Right)

ResNeXt

- **ResNeXt**
- ResNeXt — 1st Runner Up in ILSVRC 2016 (Image Classification)

Detailed Architecture of ResNet-50 and ResNeXt-50 (32x4d)			
stage	output	ResNet-50	ResNeXt-50 (32x4d)
conv1	112×112	7×7, 64, stride 2	7×7, 64, stride 2
		3×3 max pool, stride 2	3×3 max pool, stride 2
conv2	56×56	$\begin{bmatrix} 1\times 1, 64 \\ 3\times 3, 64 \\ 1\times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 128 \\ 3\times 3, 128, C=32 \\ 1\times 1, 256 \end{bmatrix} \times 3$
		$\begin{bmatrix} 1\times 1, 128 \\ 3\times 3, 128 \\ 1\times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1\times 1, 256 \\ 3\times 3, 256, C=32 \\ 1\times 1, 512 \end{bmatrix} \times 4$
conv3	28×28	$\begin{bmatrix} 1\times 1, 256 \\ 3\times 3, 256 \\ 1\times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1\times 1, 512 \\ 3\times 3, 512, C=32 \\ 1\times 1, 1024 \end{bmatrix} \times 6$
		$\begin{bmatrix} 1\times 1, 512 \\ 3\times 3, 512 \\ 1\times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 1024 \\ 3\times 3, 1024, C=32 \\ 1\times 1, 2048 \end{bmatrix} \times 3$
	1×1	global average pool 1000-d fc, softmax	global average pool 1000-d fc, softmax
# params.		25.5×10⁶	25.0×10⁶
FLOPs		4.1×10⁹	4.2×10⁹

Number of Parameters (Proportional to FLOPs)					
$C \cdot (256 \cdot d + 3 \cdot 3 \cdot d \cdot d + d \cdot 256)$					
Different settings to maintain similar complexity					
cardinality C	1	2	4	8	32
width of bottleneck d	64	40	24	14	4
width of group conv.	64	80	96	112	128

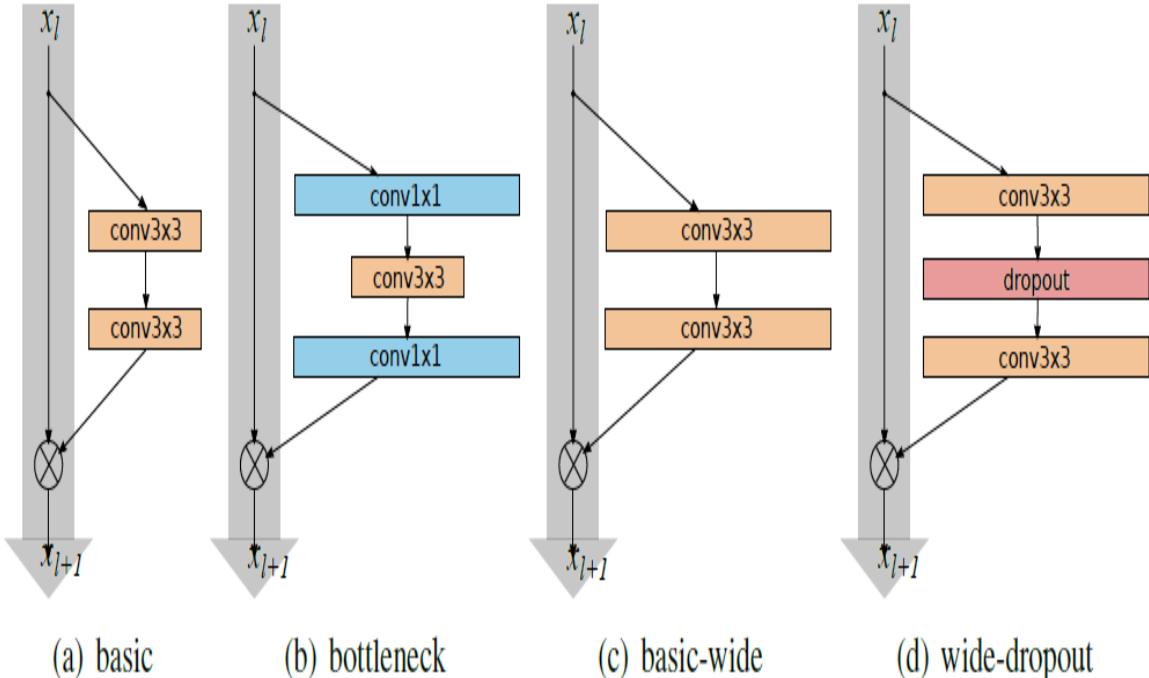
Comparisons under similar complexity		
	setting	top-1 error (%)
ResNet-50	1 × 64d	23.9
ResNeXt-50	2 × 40d	23.0
ResNeXt-50	4 × 24d	22.6
ResNeXt-50	8 × 14d	22.3
ResNeXt-50	32 × 4d	22.2
ResNet-101	1 × 64d	22.0
ResNeXt-101	2 × 40d	21.7
ResNeXt-101	4 × 24d	21.4
ResNeXt-101	8 × 14d	21.3
ResNeXt-101	32 × 4d	21.2

Detailed Architecture (Left), Number of Parameters for Each Block (Top Right), Different Settings to Maintain Similar Complexity (Middle Right), Ablation Study for Different Settings Under Similar Complexity (Bottom Right)

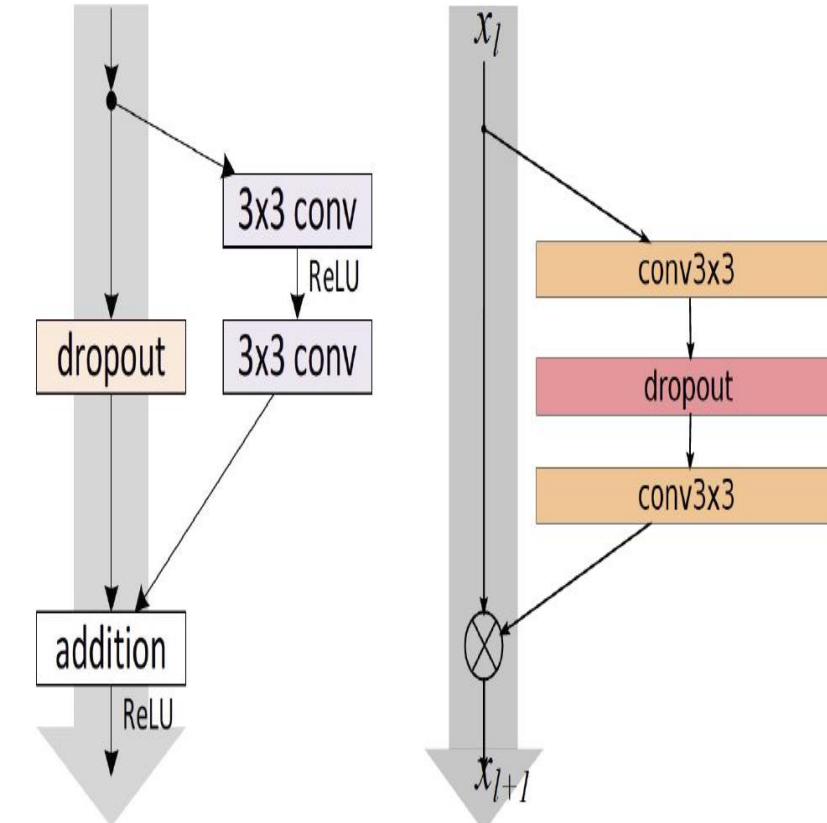
Wide Residual Networks

<https://towardsdatascience.com/review-wrns-wide-residual-networks-image-classification-d3feb3fb2004>

WRNs — Wide Residual Networks (Image Classification)



Various ResNet Blocks

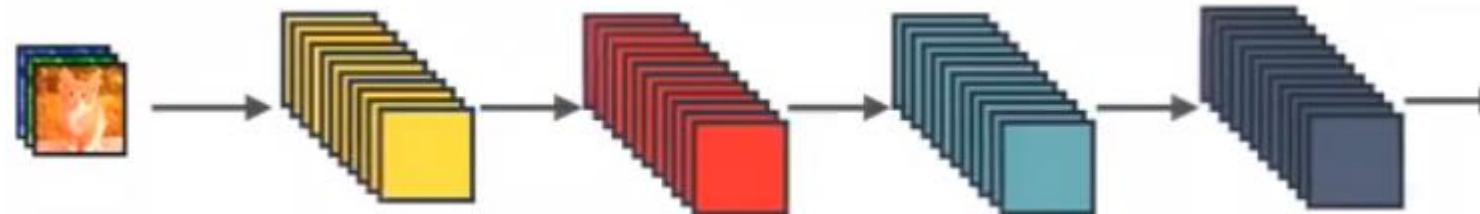


Dropout in Original ResNet (Left) and Dropout in WRNs (Right)

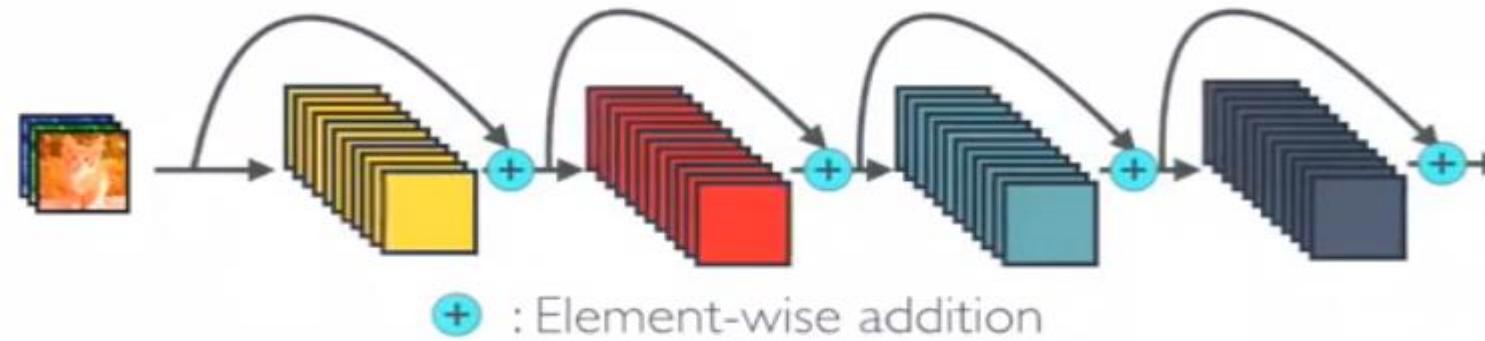
DenseNet

<https://towardsdatascience.com/review-densenet-image-classification-b6631a8ef803>

- **DenseNet —Dense Convolutional Network**



In Standard ConvNet, input image goes through multiple convolution and obtain high-level features.



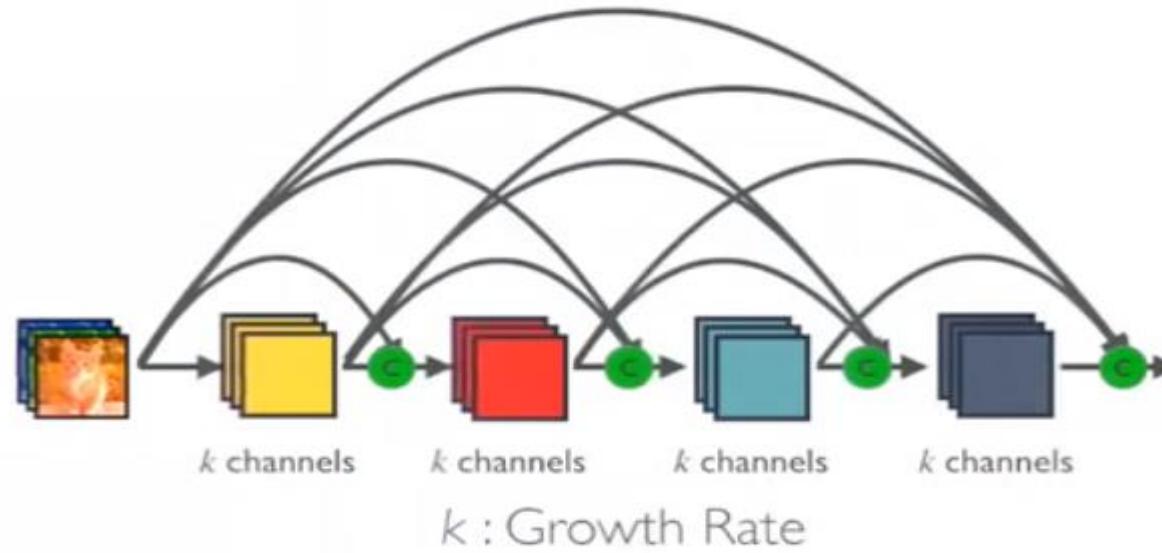
ResNet Concept

In ResNet, identity mapping is proposed to promote the gradient propagation. Element-wise addition is used. It can be viewed as algorithms with a state passed from one ResNet module to another one.

DenseNet

- **DenseNet —Dense Convolutional Network**

<https://towardsdatascience.com/review-densenet-image-classification-b6631a8ef803>



Dense Block in DenseNet with Growth Rate k

In DenseNet, each layer obtains additional inputs from all preceding layers and passes on its own feature-maps to all subsequent layers. Concatenation is used. Each layer is receiving a “collective knowledge” from all preceding layers.

Since each layer receives feature maps from all preceding layers, network can be thinner and compact, i.e. number of channels can be fewer. The growth rate k is the additional number of channels for each layer.

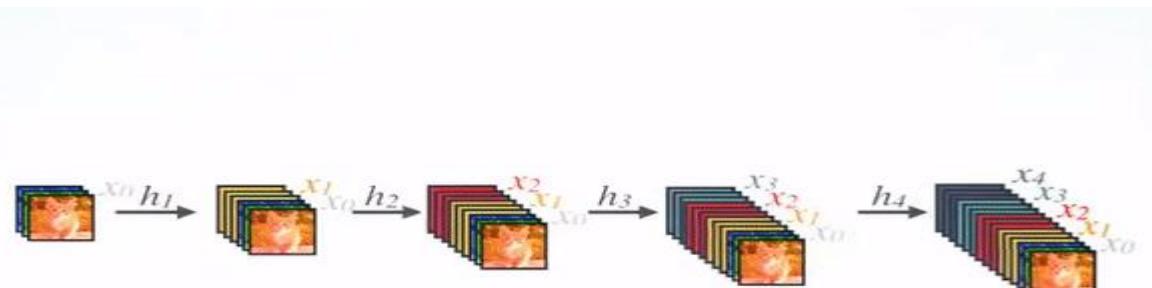
DenseNet

- **DenseNet —Dense Convolutional Network**

<https://towardsdatascience.com/review-densenet-image-classification-b6631a8ef803>



Concatenation during Forward Propagation

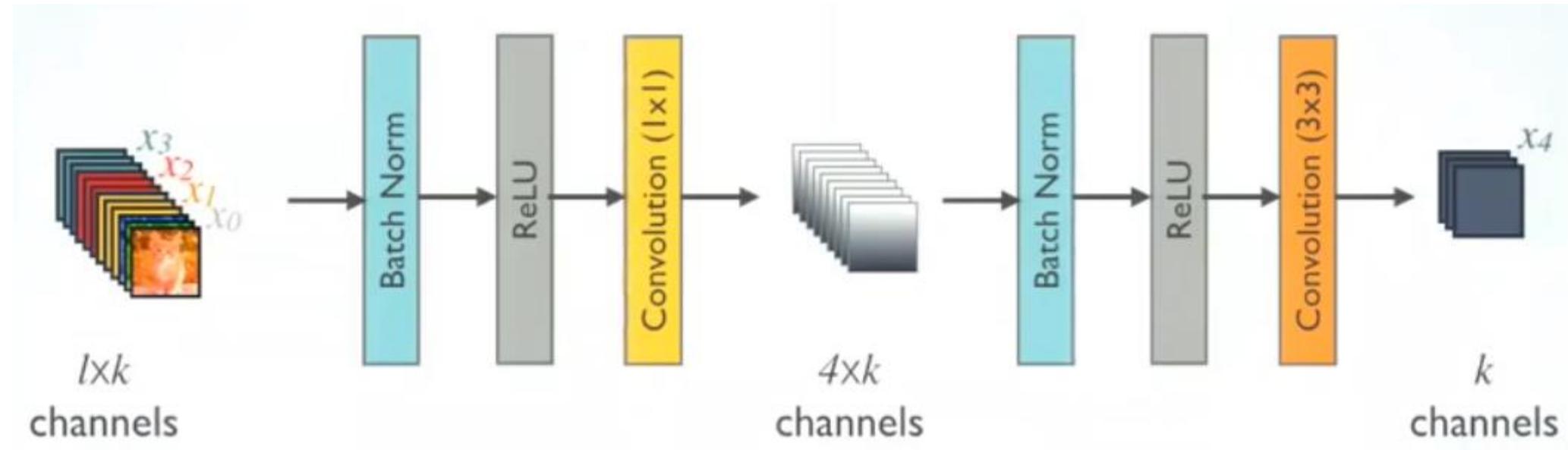


Composition Layer

DenseNet

- **DenseNet — Dense Convolutional Network**

<https://towardsdatascience.com/review-densenet-image-classification-b6631a8ef803>

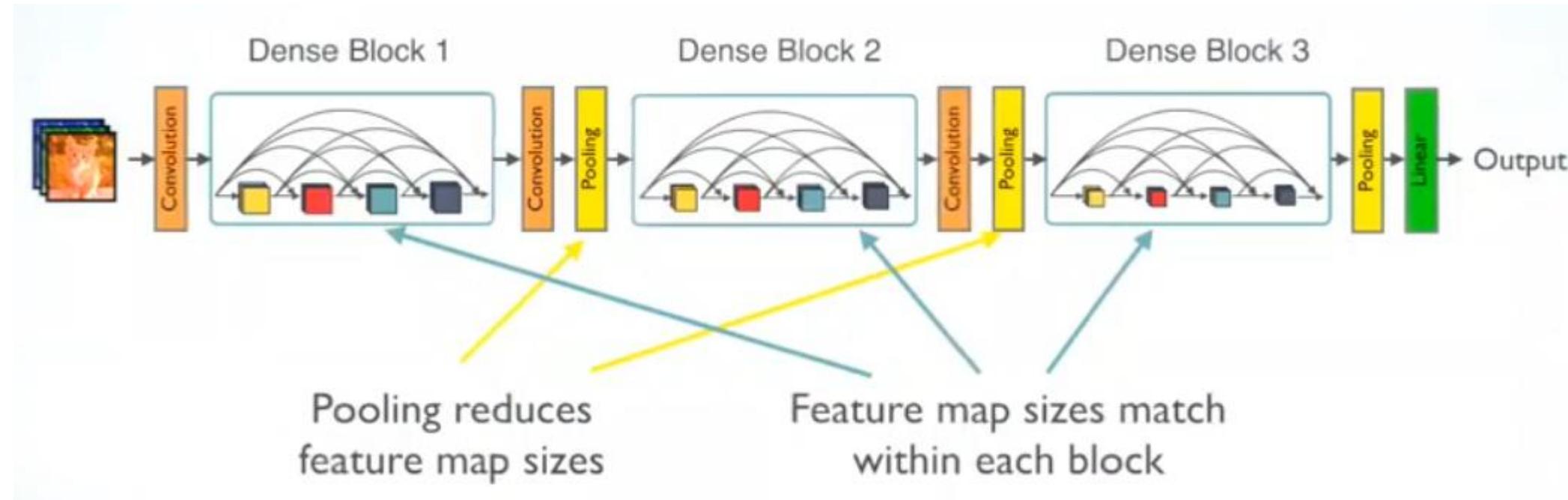


To reduce the model complexity and size, BN-ReLU- 1×1 Conv is done before BN-ReLU- 3×3 Conv.

DenseNet

<https://towardsdatascience.com/review-densenet-image-classification-b6631a8ef803>

▪ DenseNet —Dense Convolutional Network



1×1 Conv followed by 2×2 average pooling are used as the transition layers between two contiguous dense blocks.

Feature map sizes are the same within the dense block so that they can be concatenated together easily.

At the end of the last dense block, a global average pooling is performed and then a softmax classifier is attached.

DenseNet

<https://towardsdatascience.com/review-densenet-image-classification-b6631a8ef803>

▪ DenseNet —Dense Convolutional Network

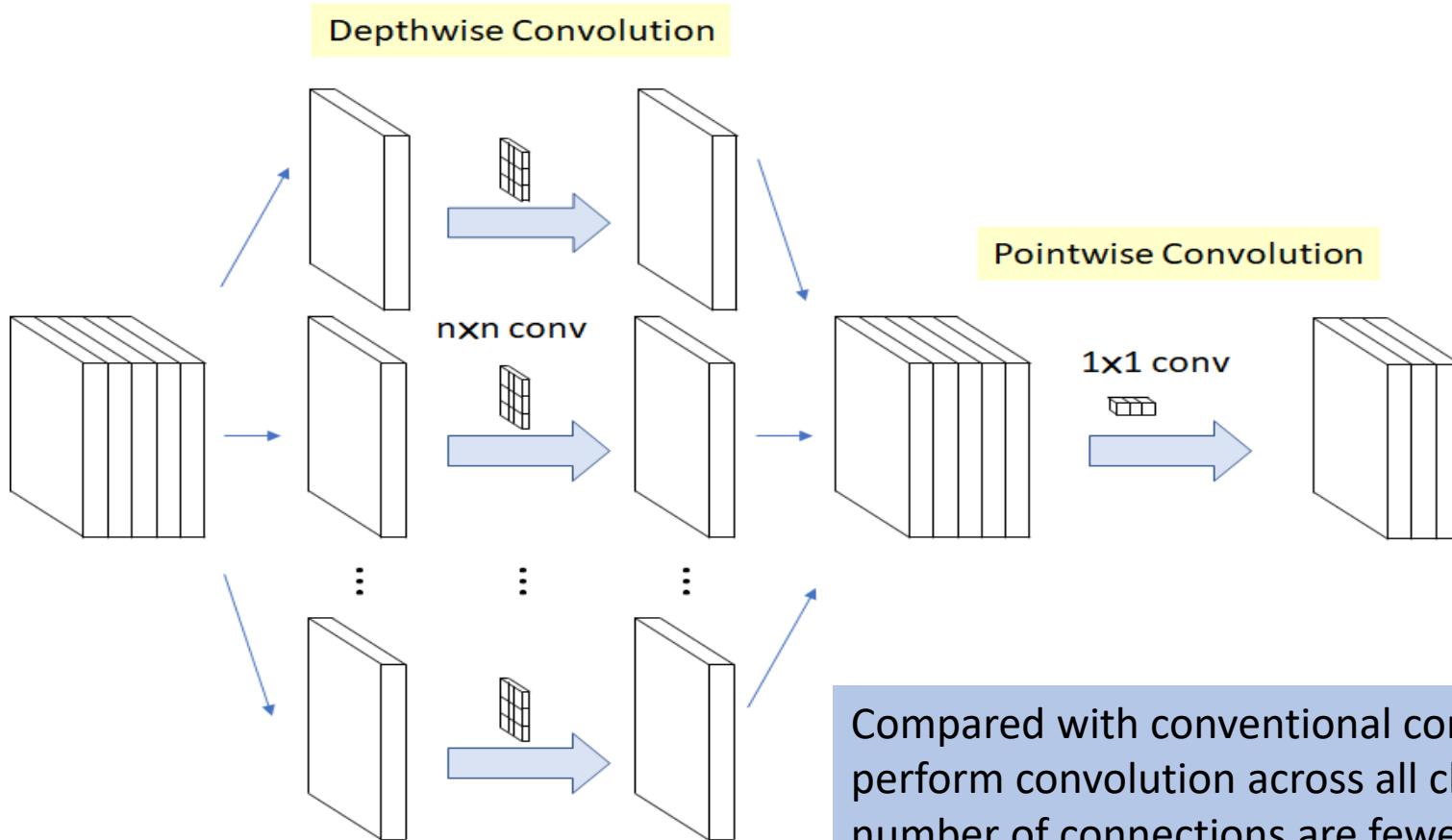
Layers	Output Size	DenseNet-121	DenseNet-169	DenseNet-201	DenseNet-264
Convolution	112×112		7×7 conv, stride 2		
Pooling	56×56		3×3 max pool, stride 2		
Dense Block (1)	56×56	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$
Transition Layer (1)	56×56		1×1 conv		
	28×28		2×2 average pool, stride 2		
Dense Block (2)	28×28	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$
Transition Layer (2)	28×28		1×1 conv		
	14×14		2×2 average pool, stride 2		
Dense Block (3)	14×14	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 24$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 48$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 64$
Transition Layer (3)	14×14		1×1 conv		
	7×7		2×2 average pool, stride 2		
Dense Block (4)	7×7	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 16$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 48$
Classification Layer	1×1		7×7 global average pool		
			1000D fully-connected, softmax		

Table 1: DenseNet architectures for ImageNet. The growth rate for all the networks is $k = 32$. Note that each “conv” layer shown in the table corresponds the sequence BN-ReLU-Conv.

Xception

- **Xception — With Depthwise Separable.**
- Original Depthwise Separable Convolution

<https://towardsdatascience.com/review-xception-with-depthwise-separable-convolution-better-than-inception-v3-image-dc967dd42568>



The original depthwise separable convolution is the depthwise convolution followed by a pointwise convolution.

Depthwise convolution is the channel-wise $n \times n$ spatial convolution.

Pointwise convolution actually is the 1×1 convolution to change the dimension.

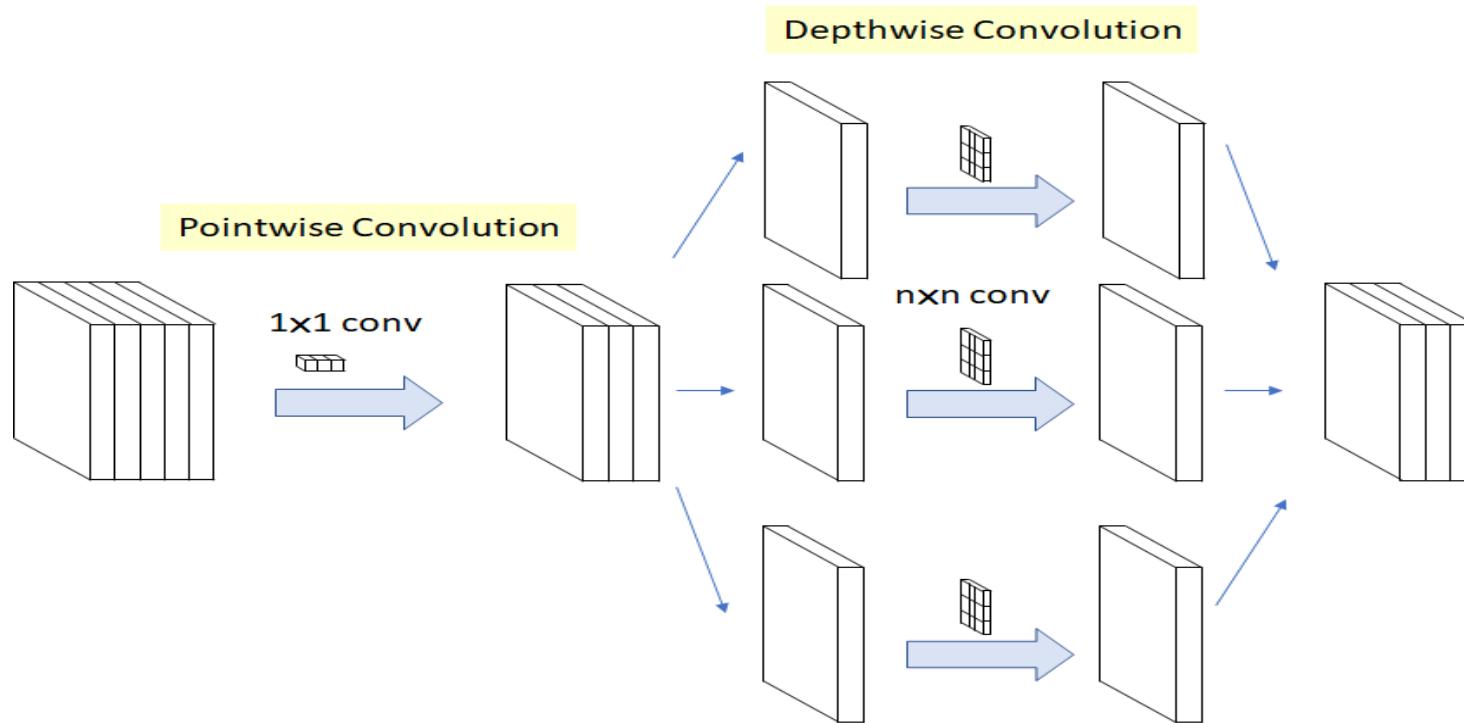
Compared with conventional convolution, we do not need to perform convolution across all channels. That means the number of connections are fewer and the model is lighter.

Xception

- **Xception — With Depthwise Separable.**
- Modified Depthwise Separable Convolution in Xception

<https://towardsdatascience.com/review-xception-with-depthwise-separable-convolution-better-than-inception-v3-image-dc967dd42568>

The modified depthwise separable convolution is the **pointwise convolution followed by a depthwise convolution**.



This is claimed to be unimportant because when it is used in stacked setting, there are only small differences appeared at the beginning and at the end of all the chained inception modules.

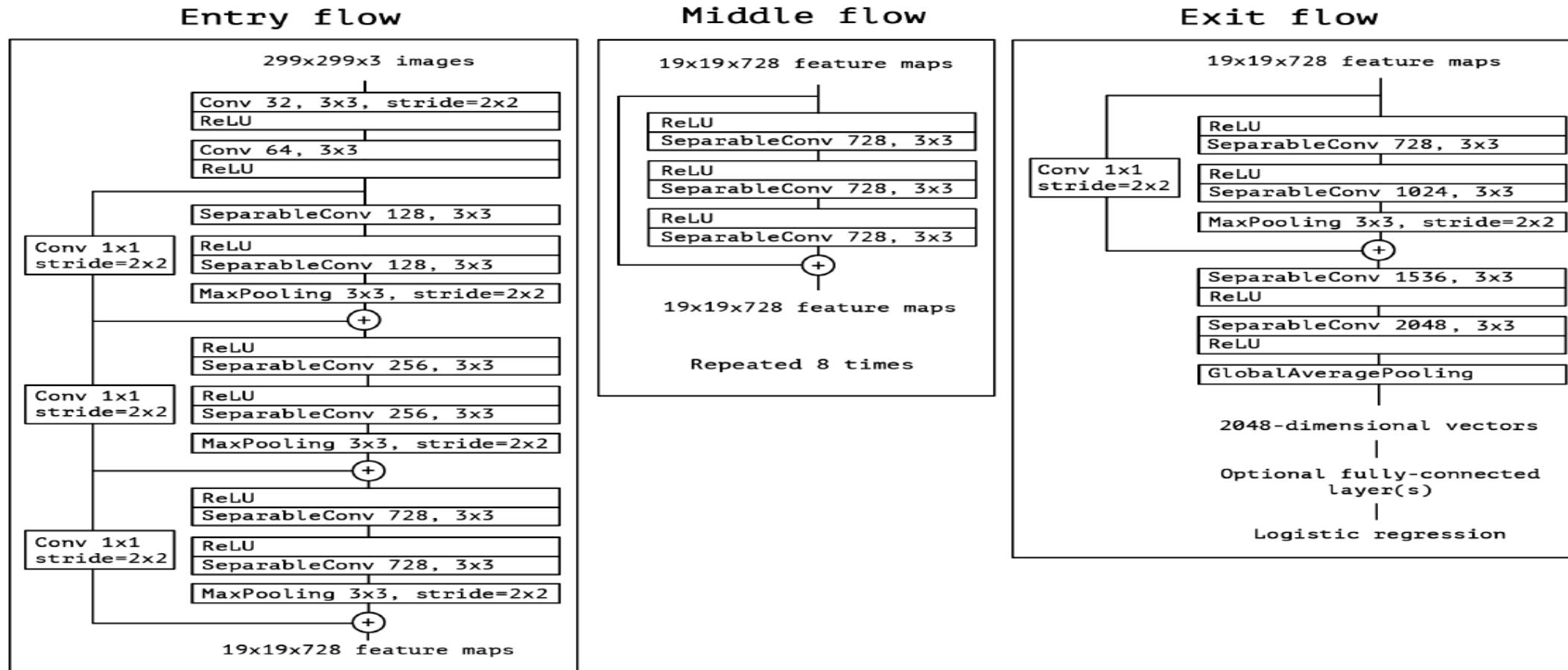
The Presence/Absence of Non-Linearity: In the original Inception Module, there is non-linearity after first operation. In Xception, the modified depthwise separable convolution, there is NO intermediate ReLU non-linearity.

The Modified Depthwise Separable Convolution used as an Inception Module in Xception, so called “extreme” version of Inception module (n=3 here).

Xception

<https://towardsdatascience.com/review-xception-with-depthwise-separable-convolution-better-than-inception-v3-image-dc967dd42568>

- Xception — With Depthwise Separable.
- Modified Depthwise Separable Convolution in Xception

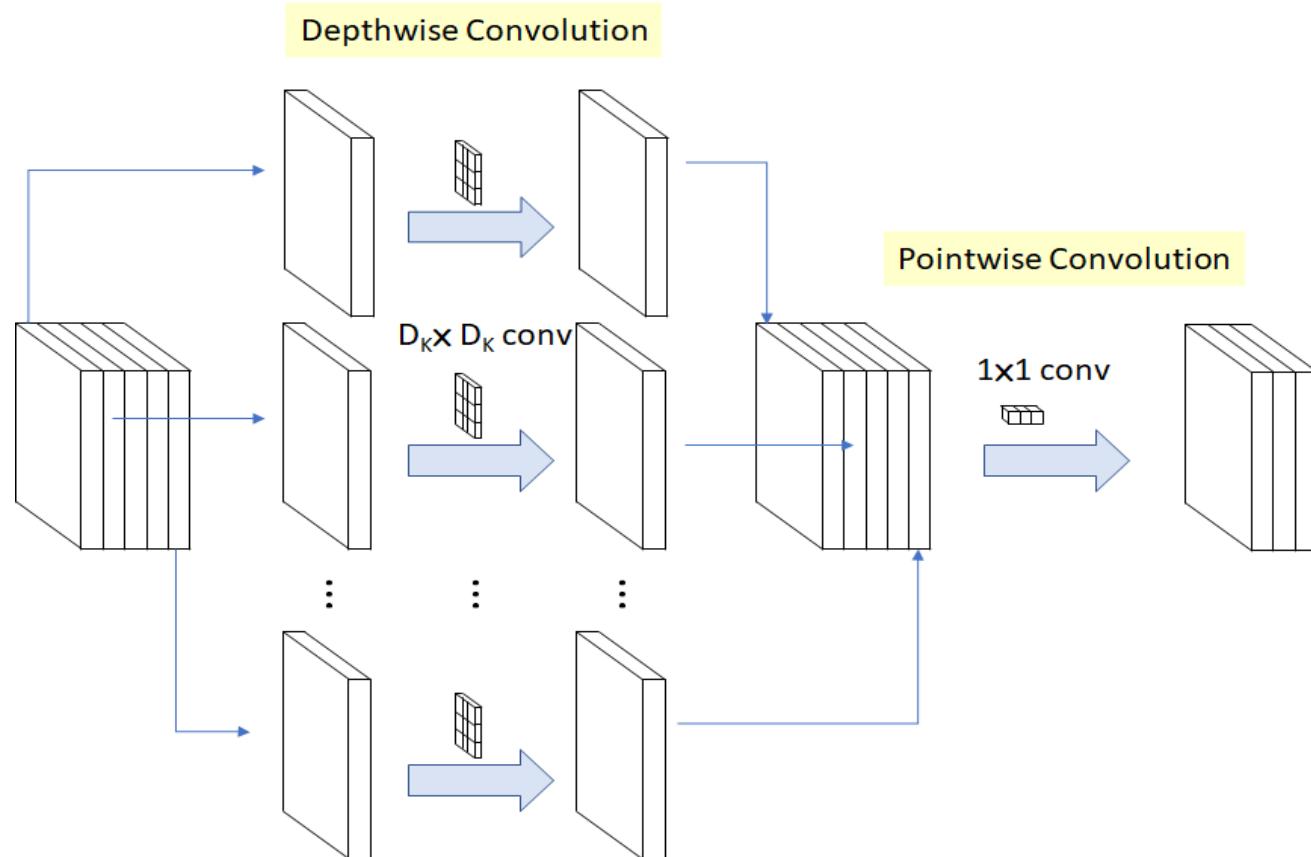


Overall Architecture of Xception (Entry Flow > Middle Flow > Exit Flow)

MobileNetV1

- **MobileNetV1**
- MobileNetV1 — Depthwise Separable Convolution (Light Weight Model)

<https://towardsdatascience.com/review-mobilenetv1-depthwise-separable-convolution-light-weight-model-a382df364b69>



Depthwise convolution is the channel-wise $D_K \times D_K$ spatial convolution.
Suppose in the figure , we have 5 channels, then we will have 5 $D_K \times D_K$ spatial convolution.
Pointwise convolution actually is the 1×1 convolution to change the dimension.

MobileNetV1

<https://towardsdatascience.com/review-mobilenetv1-depthwise-separable-convolution-light-weight-model-a382df364b69>

- **MobileNetV1**
- MobileNetV1 — Depthwise Separable Convolution (Light Weight Model)

With above operation, the operation cost is: $D_K \cdot D_K \cdot M \cdot D_F \cdot D_F + M \cdot N \cdot D_F \cdot D_F$

Depthwise Separable Convolution Cost: Depthwise Convolution Cost (Left), Pointwise Convolution Cost (Right)
where M: Number of input channels, N: Number of output channels, DK: Kernel size, and DF: Feature map size.

For standard convolution, it is: $D_K \cdot D_K \cdot M \cdot D_F \cdot D_F + M \cdot N \cdot D_F \cdot D_F$

Standard Convolution Cost $D_K \cdot D_K \cdot M \cdot N \cdot D_F \cdot D_F$

Thus, the computation reduction is:

$$\frac{D_K \cdot D_K \cdot M \cdot D_F \cdot D_F + M \cdot N \cdot D_F \cdot D_F}{D_K \cdot D_K \cdot M \cdot N \cdot D_F \cdot D_F} = \frac{1}{N} + \frac{1}{D_K^2}$$

Depthwise Separable Convolution Cost / Standard Convolution Cost

When $DK \times DK$ is 3×3 , 8 to 9 times less computation can be achieved, but with only small reduction in accuracy.

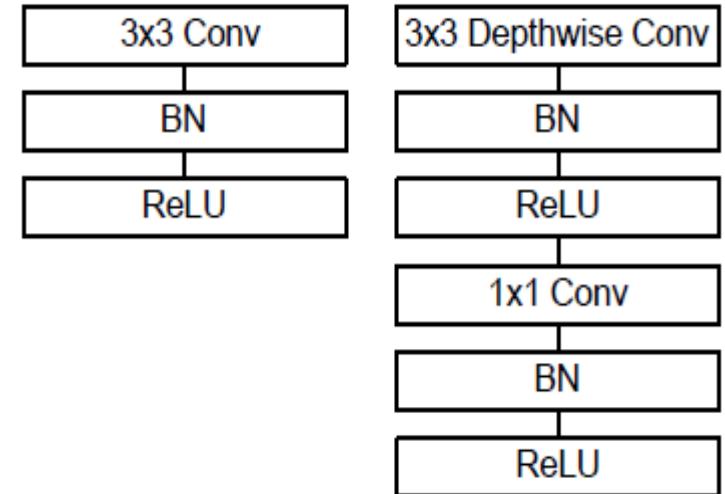
MobileNetV1

<https://towardsdatascience.com/review-mobilenetv1-depthwise-separable-convolution-light-weight-model-a382df364b69>

- **MobileNetV1**
- MobileNetV1 — Depthwise Separable Convolution (Light Weight Model)

Table 1. MobileNet Body Architecture

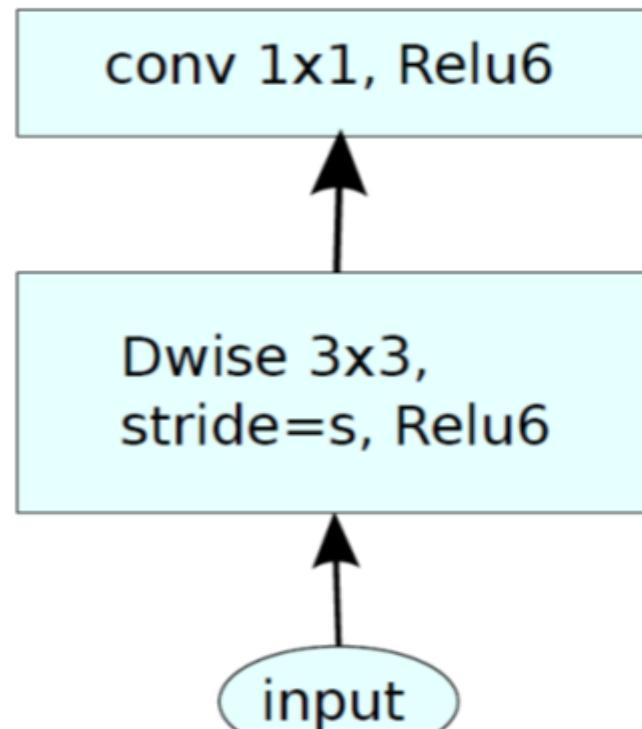
Type / Stride	Filter Shape	Input Size
Conv / s2	$3 \times 3 \times 3 \times 32$	$224 \times 224 \times 3$
Conv dw / s1	$3 \times 3 \times 32$ dw	$112 \times 112 \times 32$
Conv / s1	$1 \times 1 \times 32 \times 64$	$112 \times 112 \times 32$
Conv dw / s2	$3 \times 3 \times 64$ dw	$112 \times 112 \times 64$
Conv / s1	$1 \times 1 \times 64 \times 128$	$56 \times 56 \times 64$
Conv dw / s1	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 128$	$56 \times 56 \times 128$
Conv dw / s2	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 256$	$28 \times 28 \times 128$
Conv dw / s1	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 256$	$28 \times 28 \times 256$
Conv dw / s2	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 512$	$14 \times 14 \times 256$
5× Conv dw / s1	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 512$	$14 \times 14 \times 512$
Conv dw / s2	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 1024$	$7 \times 7 \times 512$
Conv dw / s2	$3 \times 3 \times 1024$ dw	$7 \times 7 \times 1024$
Conv / s1	$1 \times 1 \times 1024 \times 1024$	$7 \times 7 \times 1024$
Avg Pool / s1	Pool 7×7	$7 \times 7 \times 1024$
FC / s1	1024×1000	$1 \times 1 \times 1024$
Softmax / s1	Classifier	$1 \times 1 \times 1000$



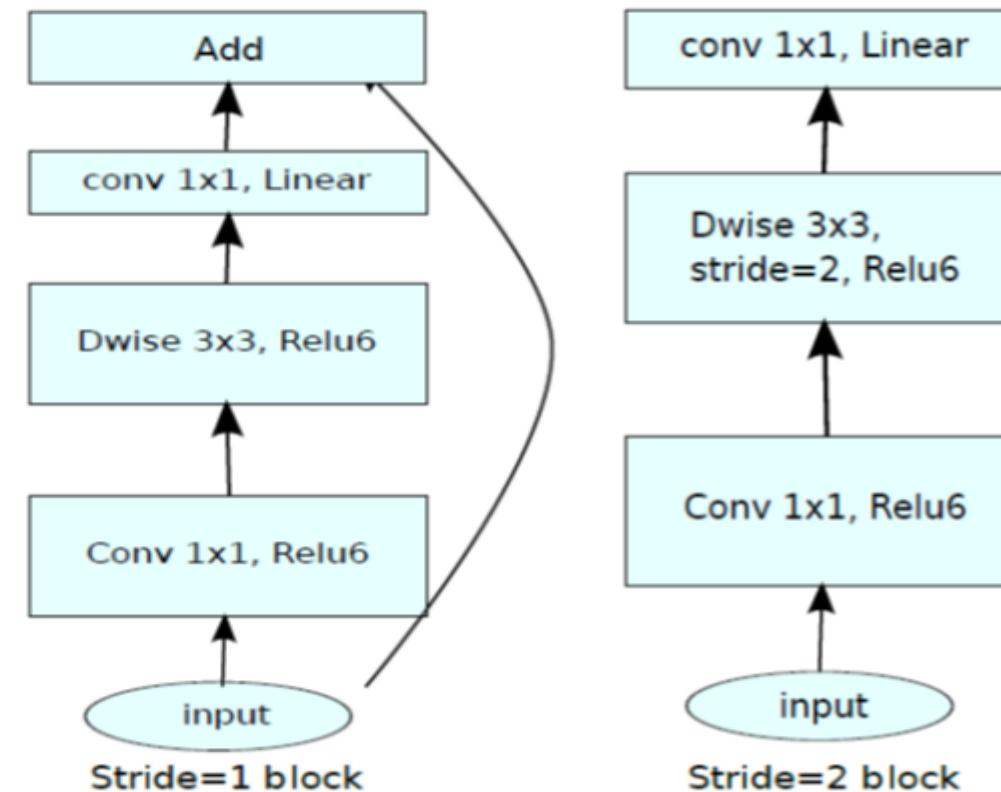
Standard Convolution (Left), Depthwise separable convolution (Right) With BN and ReLU

MobileNetV2

- MobileNetV2
- MobileNetV2 — Light Weight Model (Image Classification)



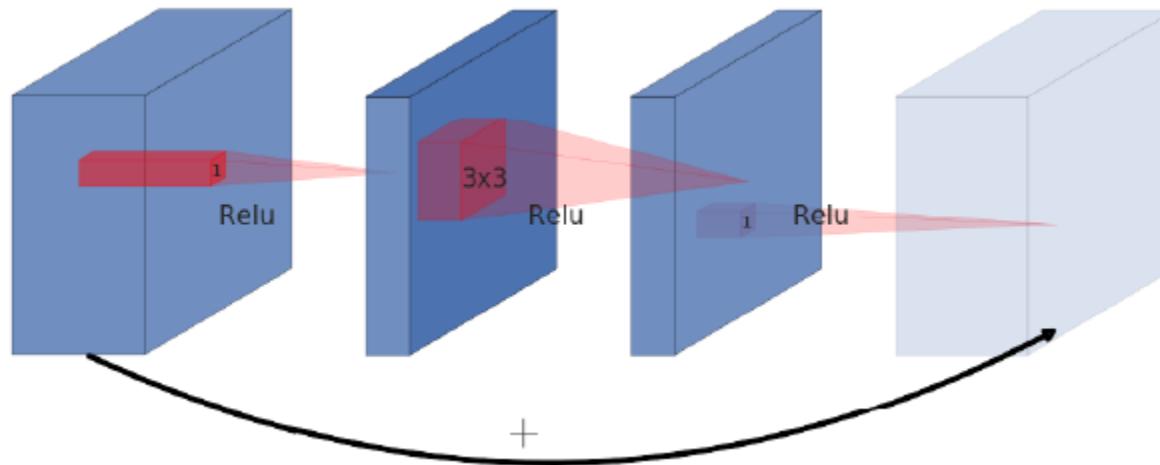
MobileNetV1



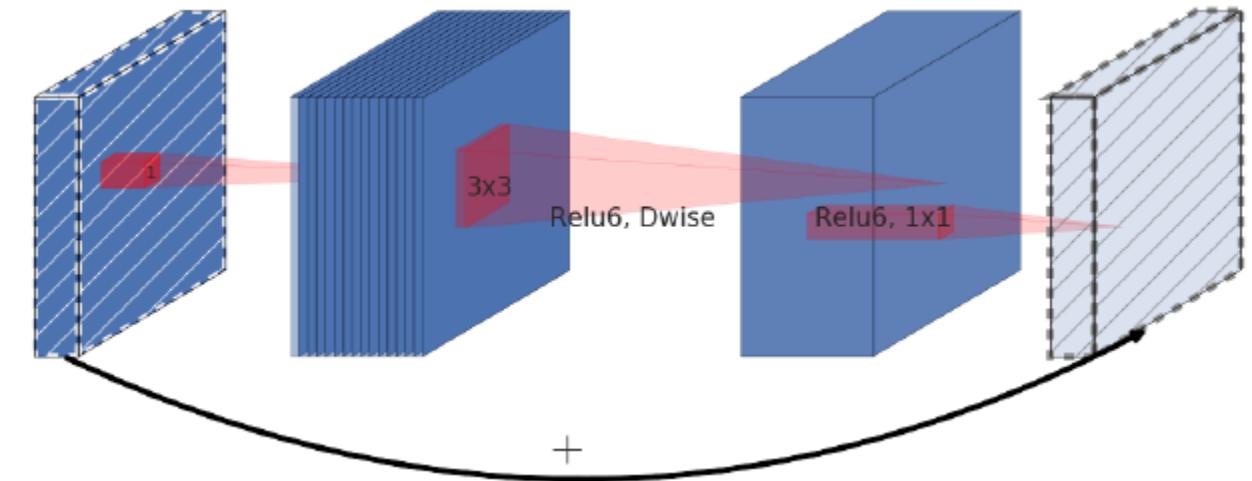
MobileNetV2

MobileNetV2

- MobileNetV2
- MobileNetV2 — Light Weight Model (Image Classification)



A residual block connects wide layers with a skip connection while layers in between are narrow



An inverted residual block connects narrow layers with a skip connection while layers in between are wide

MobileNetV2

- **MobileNetV2**
- **MobileNetV2 — Light Weight Model (Image Classification)**

In MobileNetV1, there are 2 layers.

The first layer is called a depthwise convolution, it performs lightweight filtering by applying a single convolutional filter per input channel.

The second layer is a 1×1 convolution, called a pointwise convolution, which is responsible for building new features through computing linear combinations of the input channels.

MobileNetV2

In MobileNetV2, there are two types of blocks. One is residual block with stride of 1.

Another one is block with stride of 2 for downsizing.

There are 3 layers for both types of blocks.

This time, the first layer is 1×1 convolution with ReLU6.

The second layer is the depthwise convolution.

The third layer is another 1×1 convolution but without any non-linearity. It is claimed that if ReLU is used again, the deep networks only have the power of a linear classifier on the non-zero volume part of the output domain.

MobileNetV2

- **MobileNetV2**
- **MobileNetV2 — Light Weight Model (Image Classification)**

Input	Operator	<i>t</i>	<i>c</i>	<i>n</i>	<i>s</i>
$224^2 \times 3$	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	-	-

where t: expansion factor,

c: **number of output channels**,

n: **repeating number**,

s: **stride. 3x3 kernels are used for spatial convolution.**

In typical, the primary network (width multiplier 1, 224×224), has a computational cost of 300 million multiply-adds and uses 3.4 million parameters. (Width multiplier is introduced in MobileNetV1.)

The performance trade offs are further explored, for input resolutions from 96 to 224, and width multipliers of 0.35 to 1.4. The network computational cost up to 585M M Adds, while the model size vary between 1.7M and 6.9M parameters.

MobileNetV2

▪ MobileNetV2

Input	Operator	<i>t</i>	<i>c</i>	<i>n</i>	<i>s</i>
$224^2 \times 3$	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	-	<i>k</i>	-	-

```
def residual_block(x, squeeze=16, expand=64):
    m = Conv2D(squeeze, (1,1), activation='relu')(x)
    m = Conv2D(squeeze, (3,3), activation='relu')(m)
    m = Conv2D(expand, (1,1), activation='relu')(m)
    return Add()([m, x])
```

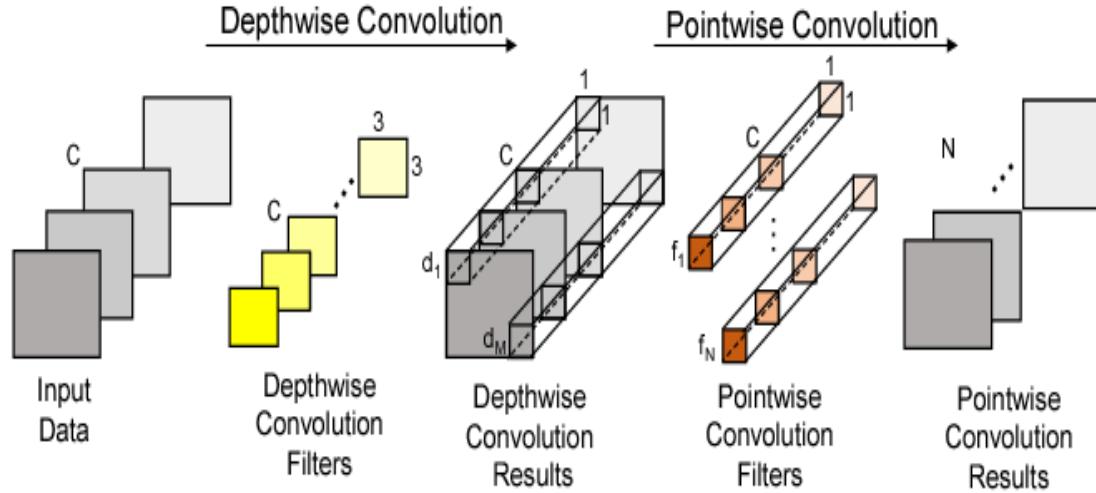
```
def bottleneck_block(x, expand=64, squeeze=16):
    m = Conv2D(expand, (1,1))(x)
    m = BatchNormalization()(m)
    m = Activation('relu6')(m)
    m = DepthwiseConv2D((3,3))(m)
    m = BatchNormalization()(m)
    m = Activation('relu6')(m)
    m = Conv2D(squeeze, (1,1))(m)
    m = BatchNormalization()(m)
    return Add()([m, x])
```

```
def inverted_residual_block(x, expand=64, squeeze=16):
    m = Conv2D(expand, (1,1), activation='relu')(x)
    m = DepthwiseConv2D((3,3), activation='relu')(m)
    m = Conv2D(squeeze, (1,1), activation='relu')(m)
    return Add()([m, x])
```

EfficientNet Models

▪ Rethinking Model Scaling for Convolutional Neural Networks

(a)



(b)

$$\begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_N \end{bmatrix} \times \begin{bmatrix} d_1 & d_2 & \cdots & d_M \end{bmatrix} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_N \end{bmatrix}$$

Filter Matrix
(Pointwise Convolution Filters)

Data Matrix
(Depthwise Convolution Results)

Results Matrix
(Pointwise Convolution Results)

<https://heartbeat.fritz.ai/reviewing-efficientnet-increasing-the-accuracy-and-robustness-of-cnns-6aaf411fc81d>

Depthwise Convolution + Pointwise Convolution: Divides the original convolution into two stages to significantly reduce the cost of calculation, with a minimum loss of accuracy.

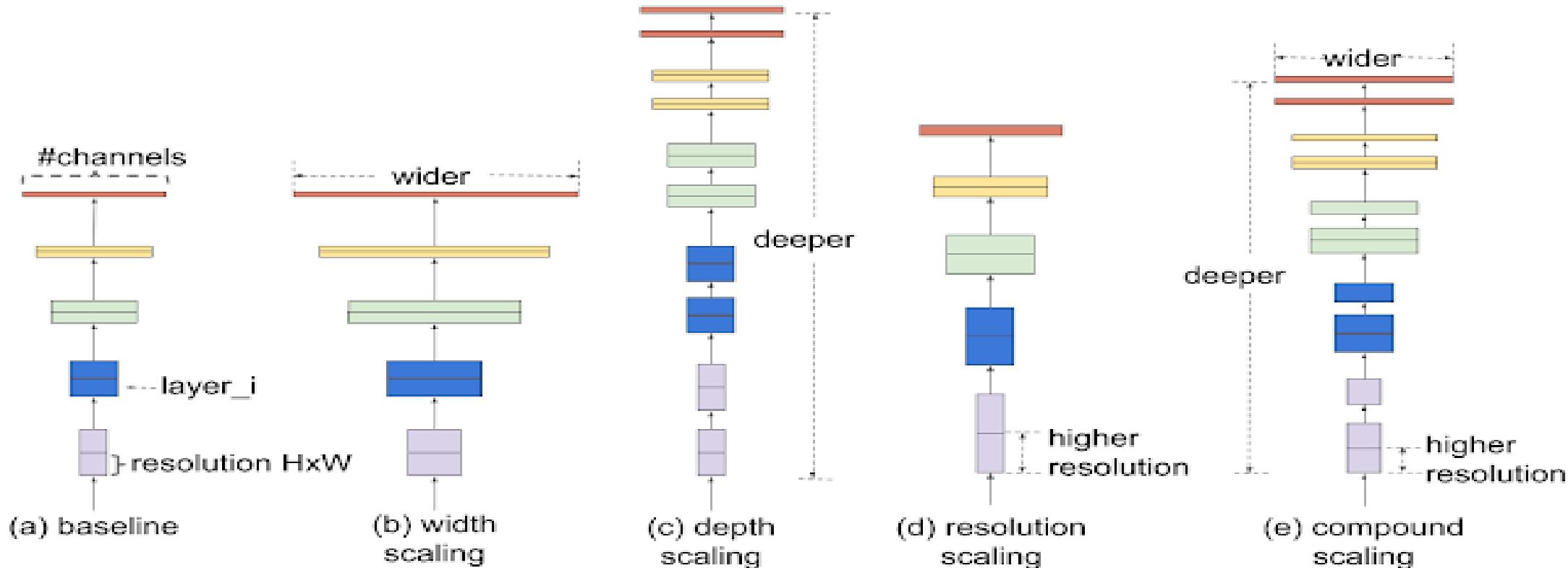
Inverse Res: The original ResNet blocks consist of a layer that squeezes the channels, then a layer that extends the channels. In this way, it links skip connections to rich channel layers. In **MBConv**, however, blocks consist of a layer that first extends channels and then compresses them, so that layers with fewer channels are skip connected.

Linear bottleneck: Uses linear activation in the last layer in each block to prevent loss of information from ReLU.

EfficientNet Models

▪ EfficientNet-B0-B7

<https://heartbeat.fritz.ai/reviewing-efficientnet-increasing-the-accuracy-and-robustness-of-cnns-6aaf411fc81d>



Model Scaling. (a) is a baseline network example; (b)-(d) are conventional scaling that only increases one dimension of network width, depth, or resolution. (e) is our proposed compound scaling method that uniformly scales all three dimensions with a fixed ratio.

EfficientNet Models

▪ EfficientNet-B0-B7

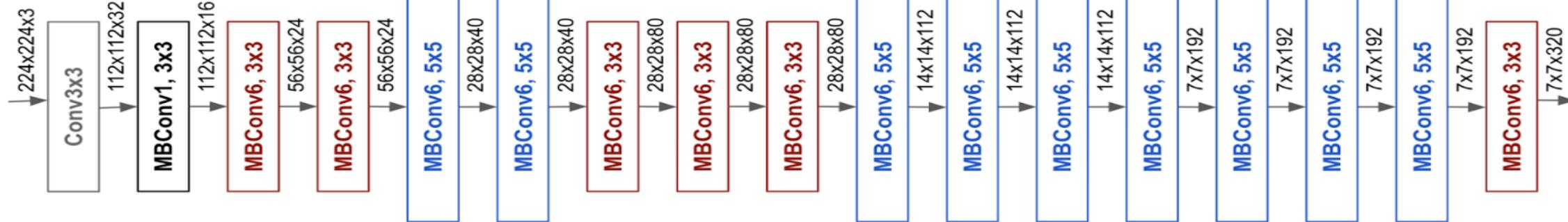
Stage i	Operator $\hat{\mathcal{F}}_i$	Resolution $\hat{H}_i \times \hat{W}_i$	#Channels \hat{C}_i	#Layers \hat{L}_i
1	Conv3x3	224×224	32	1
2	MBConv1, k3x3	112×112	16	1
3	MBConv6, k3x3	112×112	24	2
4	MBConv6, k5x5	56×56	40	2
5	MBConv6, k3x3	28×28	80	3
6	MBConv6, k5x5	14×14	112	3
7	MBConv6, k5x5	14×14	192	4
8	MBConv6, k3x3	7×7	320	1
9	Conv1x1 & Pooling & FC	7×7	1280	1

Step 1:

Building a basic model is called EfficientNet-B0.

MBConv is used with MobileNet's inverted Res bottlenecks.

Basic network structure of EfficientNet-B0



A basic block representation of EfficientNet-B0

EfficientNet Models

<https://towardsdatascience.com/review-mobilenetv2-light-weight-model-image-classification-8febb490e61c>

▪ EfficientNet-B0-B7

Step 1:

Building a basic model is called EfficientNet-B0.

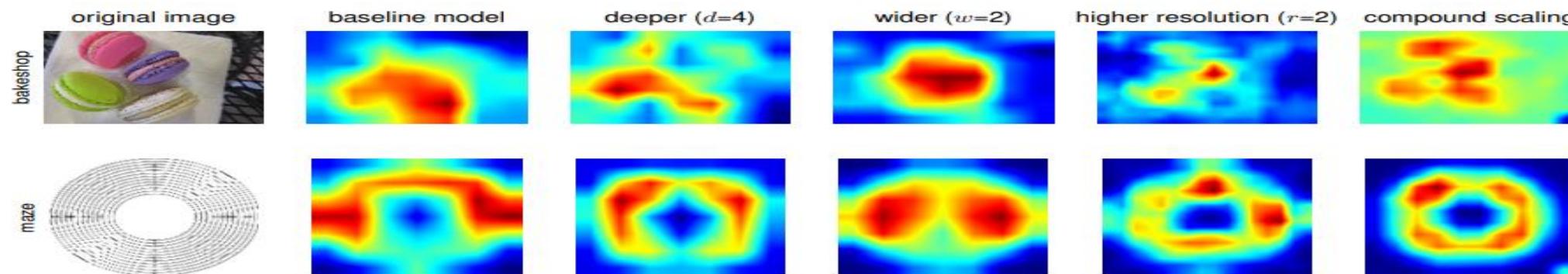
MBConv is used with MobileNet's inverted Res bottlenecks.

Step 2:

$\varphi = 1$ and grid search for α , β , and γ to scale from B0 to B1.

Step 3:

α, β, γ set. Thus, for scaling from B2 to B7, φ is selected between $2 \sim 7$. Below, you can see class activation maps for models with different scaling methods.



Class Activation Map (CAM) (Zhou et al., 2016) for Models with different scaling methods- Our compound scaling method allows the scaled model (last column) to focus on more relevant regions with more object details.

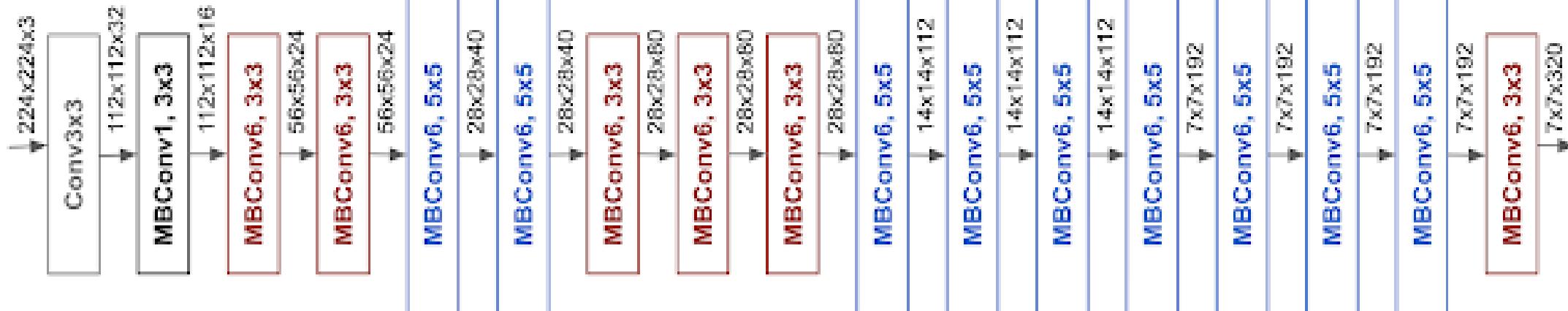
EfficientNet Models

<https://towardsdatascience.com/review-mobilenetv2-light-weight-model-image-classification-8febb490e61c>

■ EfficientNet-B0-B7

Table 1. EfficientNet-B0 baseline network – Each row describes a stage i with \hat{L}_i layers, with input resolution $\langle \hat{H}_i, \hat{W}_i \rangle$ and output channels \hat{C}_i . Notations are adopted from equation 2.

Stage i	Operator $\hat{\mathcal{F}}_i$	Resolution $\hat{H}_i \times \hat{W}_i$	#Channels \hat{C}_i	#Layers \hat{L}_i
1	Conv3x3	224×224	32	1
2	MBConv1, k3x3	112×112	16	1
3	MBConv6, k3x3	112×112	24	2
4	MBConv6, k5x5	56×56	40	2
5	MBConv6, k3x3	28×28	80	3
6	MBConv6, k5x5	14×14	112	3
7	MBConv6, k5x5	14×14	192	4
8	MBConv6, k3x3	7×7	320	1
9	Conv1x1 & Pooling & FC	7×7	1280	1

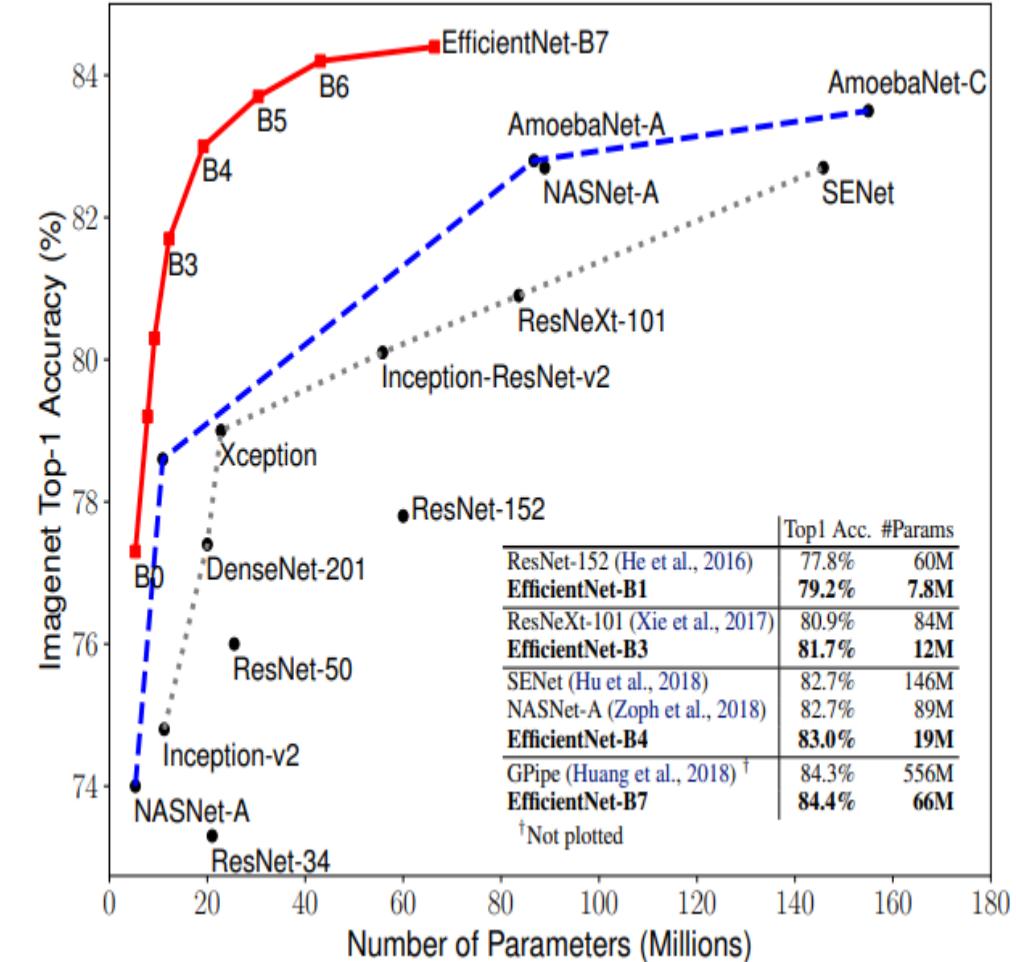


EfficientNet Models

▪ EfficientNet-B0-B7

	Comparison to best public-available results						Comparison to best reported results					
	Model	Acc.	#Param	Our Model	Acc.	#Param(ratio)	Model	Acc.	#Param	Our Model	Acc.	#Param(ratio)
CIFAR-10	NASNet-A	98.0%	85M	EfficientNet-B0	98.1%	4M (21x)	Gpipe	99.0%	556M	EfficientNet-B7	98.9%	64M (8.7x)
CIFAR-100	NASNet-A	87.5%	85M	EfficientNet-B0	88.1%	4M (21x)	Gpipe	91.3%	556M	EfficientNet-B7	91.7%	64M (8.7x)
Birdsnap	Inception-v4	81.8%	41M	EfficientNet-B5	82.0%	28M (1.5x)	GPipe	83.6%	556M	EfficientNet-B7	84.3%	64M (8.7x)
Stanford Cars	Inception-v4	93.4%	41M	EfficientNet-B3	93.6%	10M (4.1x)	DAT	94.8%	-	EfficientNet-B7	94.7%	-
Flowers	Inception-v4	98.5%	41M	EfficientNet-B5	98.5%	28M (1.5x)	DAT	97.7%	-	EfficientNet-B7	98.8%	-
FGVC Aircraft	Inception-v4	90.9%	41M	EfficientNet-B3	90.7%	10M (4.1x)	DAT	92.9%	-	EfficientNet-B7	92.9%	-
Oxford-IIIT Pets	ResNet-152	94.5%	58M	EfficientNet-B4	94.8%	17M (5.6x)	GPipe	95.9%	556M	EfficientNet-B6	95.4%	41M (14x)
Food-101	Inception-v4	90.8%	41M	EfficientNet-B4	91.5%	17M (2.4x)	GPipe	93.0%	556M	EfficientNet-B7	93.0%	64M (8.7x)
Geo-Mean					(4.7x)						(9.6x)	

<https://towardsdatascience.com/review-mobilenetv2-light-weight-model-image-classification-8febb490e61c>



Transfer Learning DL models

▪ Transfer Learning

ImageNet Challenge

IMAGENET

- 1,000 object classes (categories).
- Images:
 - 1.2 M train
 - 100k test.



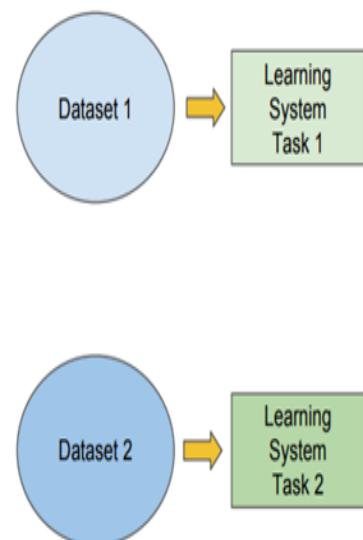
<https://towardsdatascience.com/a-comprehensive-hands-on-guide-to-transfer-learning-with-real-world-applications-in-deep-learning-212bf3b2f27a>

Traditional ML

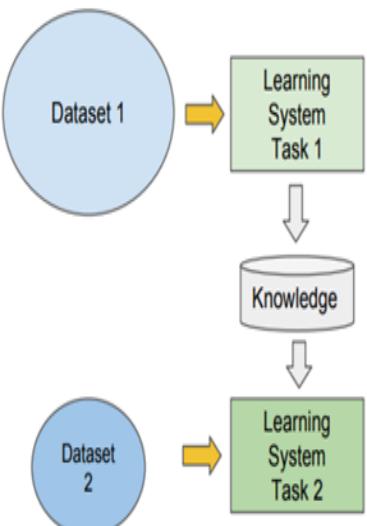
vs

Transfer Learning

- Isolated, single task learning:
 - Knowledge is not retained or accumulated. Learning is performed w.o. considering past learned knowledge in other tasks



- Learning of a new tasks relies on the previous learned tasks:
 - Learning process can be faster, more accurate and/or need less training data



Transfer Learning DL models

▪ Transfer Learning

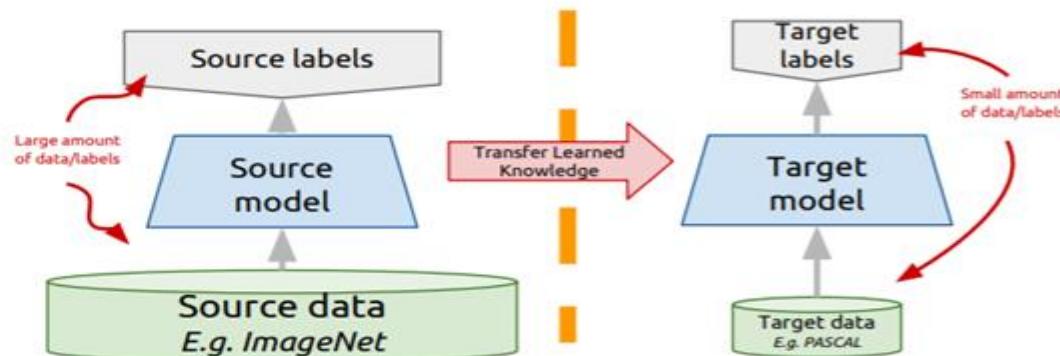
Transfer learning: idea

Instead of training a deep network from scratch for your task:

- Take a network trained on a different domain for a different **source task**
- Adapt it for your domain and your **target task**

Variations:

- Same domain, different task
- Different domain, same task



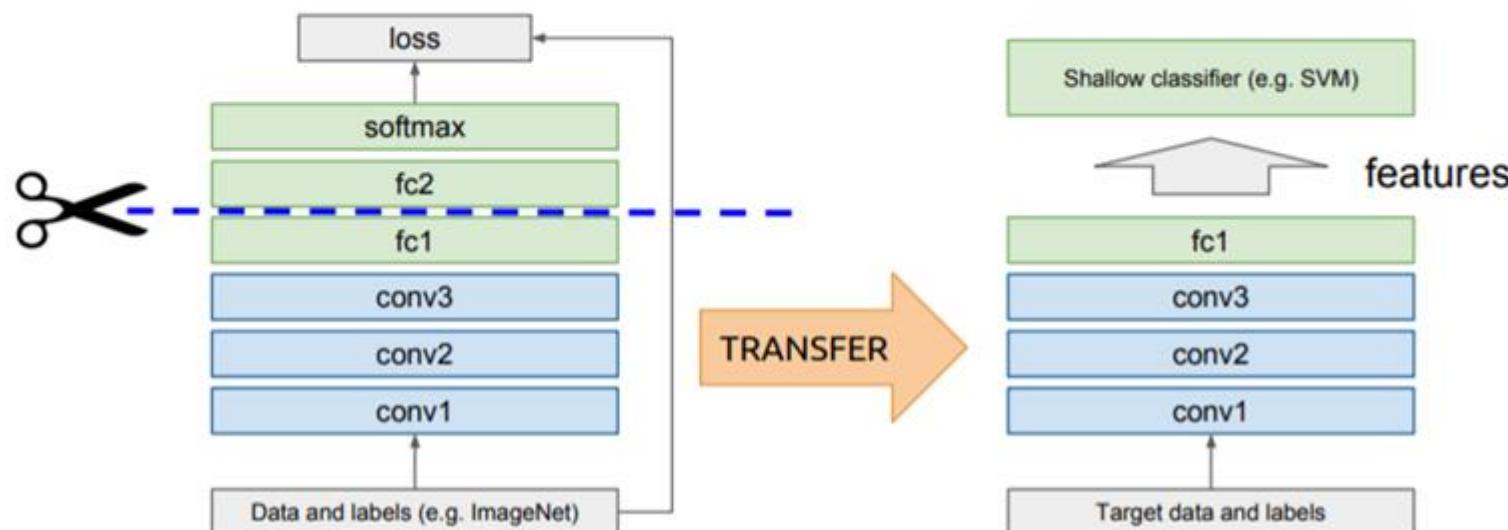
<https://towardsdatascience.com/a-comprehensive-hands-on-guide-to-transfer-learning-with-real-world-applications-in-deep-learning-212bf3b2f27a>

Transfer Learning DL models

▪ Off-the-shelf Pre-trained Models as Feature Extractors

Idea: use outputs of one or more layers of a network trained on a different task as generic feature detectors. Train a new shallow model on these features.

Assumes that $D_S = D_T$

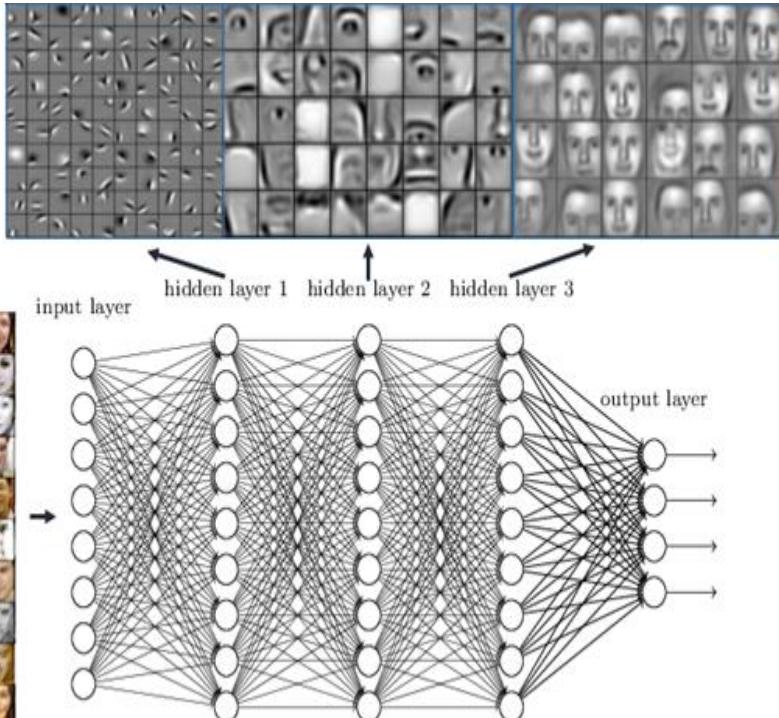


<https://towardsdatascience.com/a-comprehensive-hands-on-guide-to-transfer-learning-with-real-world-applications-in-deep-learning-212bf3b2f27a>

Transfer Learning DL models

▪ Fine Tuning Off-the-shelf Pre-trained Models

Deep neural networks learn hierarchical feature representations



Fine-tuning: supervised domain adaptation

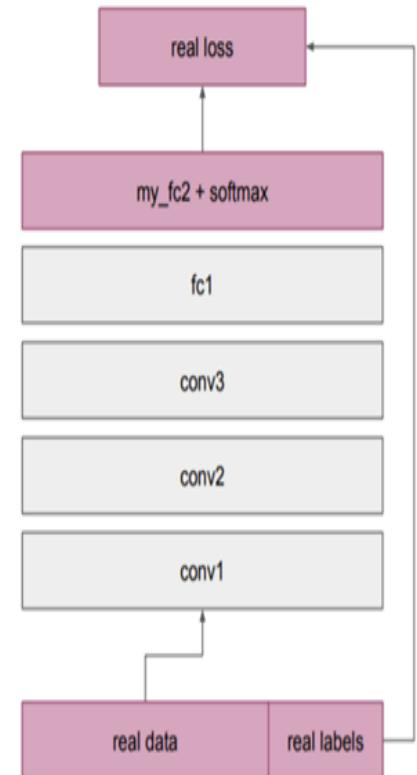
Train deep net on “nearby” task for which it is easy to get labels using standard backprop

- E.g. ImageNet classification
- Pseudo classes from augmented data
- Slow feature learning, ego-motion

Cut off top layer(s) of network and replace with supervised objective for target domain

Fine-tune network using backprop with labels for target domain until validation loss starts to increase

Aligns D_S with D_T



Transfer Learning DL models

▪ Fine Tuning Off-the-shelf Pre-trained Models

Freeze or fine-tune?

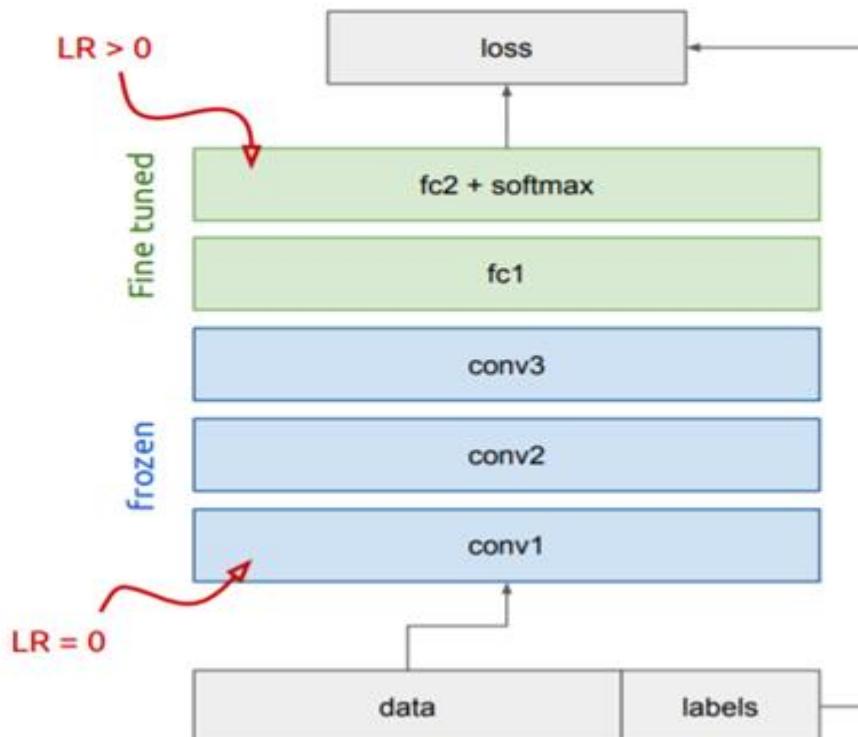
Bottom n layers can be frozen or fine tuned.

- **Frozen:** not updated during backprop
- **Fine-tuned:** updated during backprop

Which to do depends on target task:

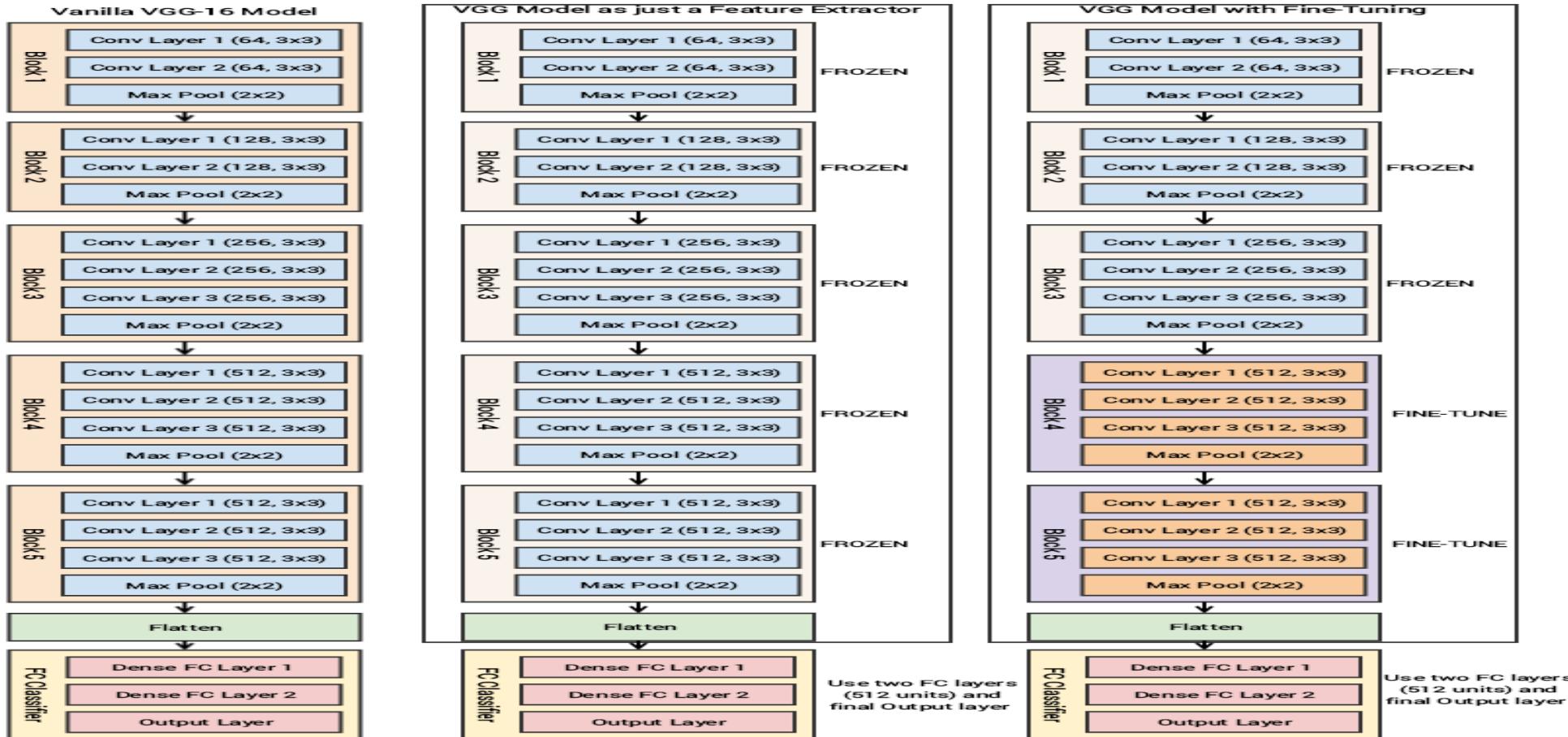
- **Freeze:** target task labels are scarce, and we want to avoid overfitting
- **Fine-tune:** target task labels are more plentiful

In general, we can set learning rates to be different for each layer to find a tradeoff between freezing and fine tuning



Transfer Learning DL models

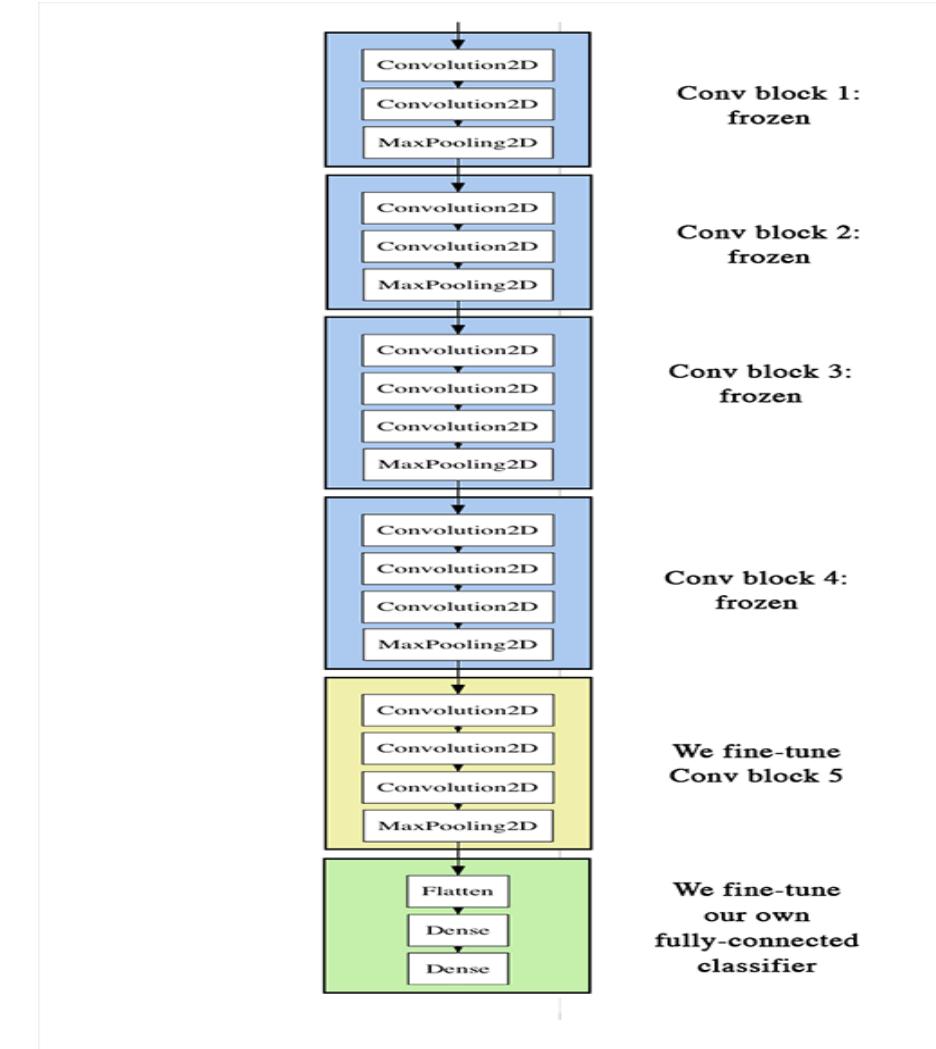
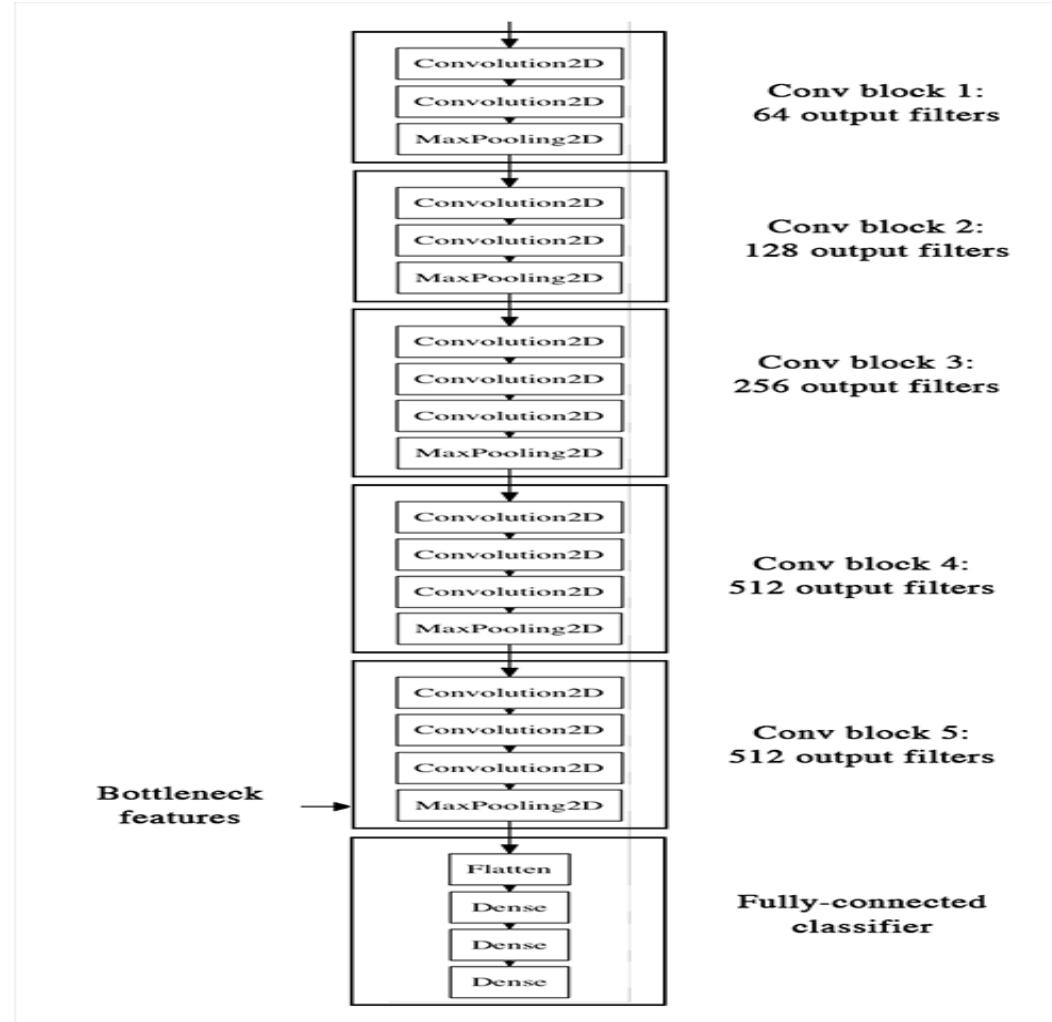
▪ Fine Tuning Off-the-shelf Pre-trained Models



<https://towardsdatascience.com/a-comprehensive-hands-on-guide-to-transfer-learning-with-real-world-applications-in-deep-learning-212bf3b2f27a>

Transfer Learning DL models

▪ Fine Tuning Off-the-shelf Pre-trained Models



Object Detection DL models

- Deep learning approach for object detection and classification

OverFeat

One of the first advances in using deep learning for object detection was OverFeat from NYU published in 2013. They proposed a multi-scale sliding window algorithm using Convolutional Neural Networks (CNNs).

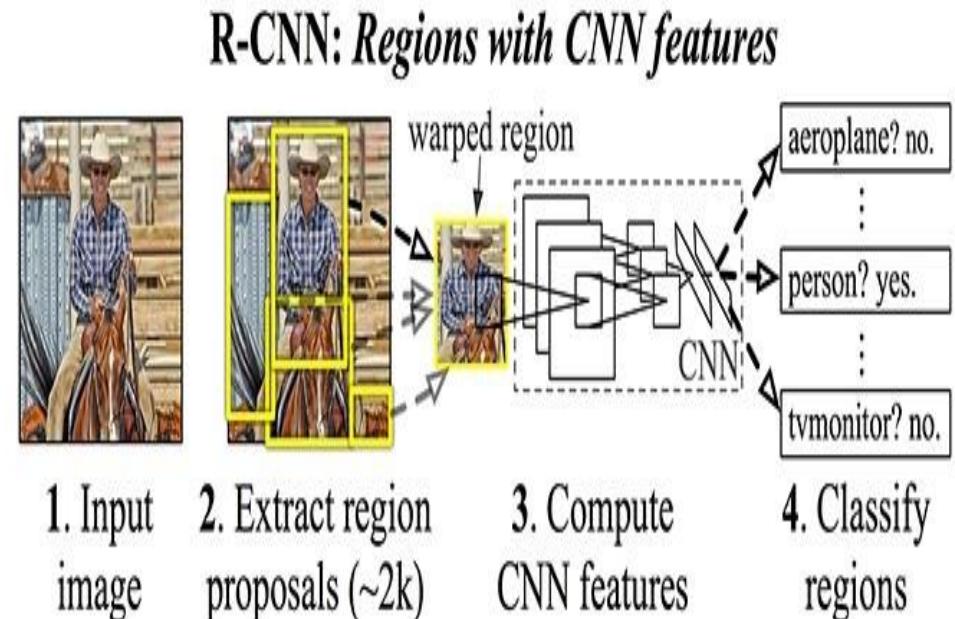
R-CNN

Quickly after OverFeat, Regions with CNN features or R-CNN from Ross Girshick, et al. at the UC Berkeley was published which boasted an almost 50% improvement on the object detection challenge. What they proposed was a three stage approach:

Extract possible objects using a region proposal method (the most popular one being Selective Search).

Extract features from each region using a CNN.

Classify each region with SVMs.



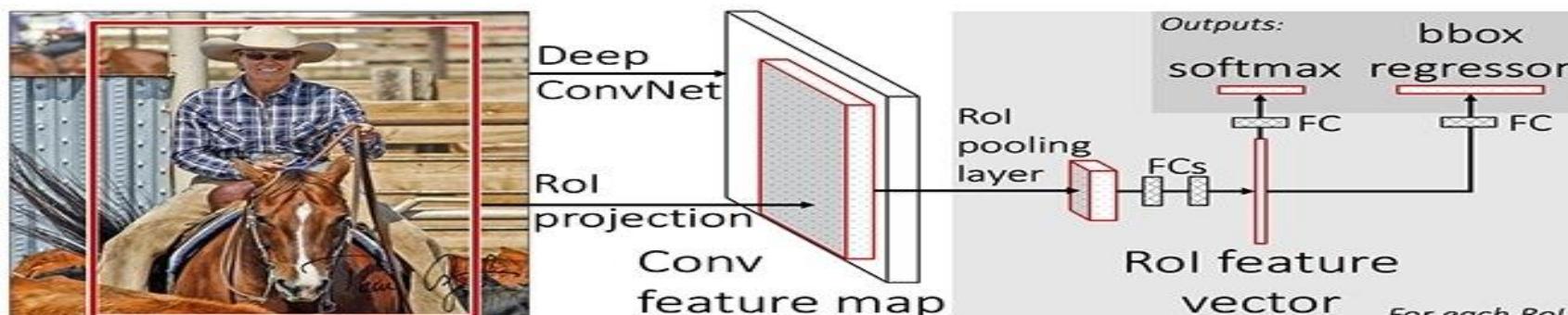
Object Detection DL models

- Deep learning approach for object detection and classification

Fast R-CNN

This approach quickly evolved into a purer deep learning one, when a year later Ross Girshick (now at Microsoft Research) published Fast R-CNN. **Similar to R-CNN, it used Selective Search to generate object proposals, but instead of extracting all of them independently and using SVM classifiers, it applied the CNN on the complete image and then used both Region of Interest (RoI) Pooling on the feature map with a final feed forward network for classification and regression.**

Not only was this approach faster, but having the RoI Pooling layer and the fully connected layers allowed the model to be end-to-end differentiable and easier to train. The biggest downside was that the model still relied on **Selective Search (or any other region proposal algorithm)**, which became the bottleneck when using it for inference.

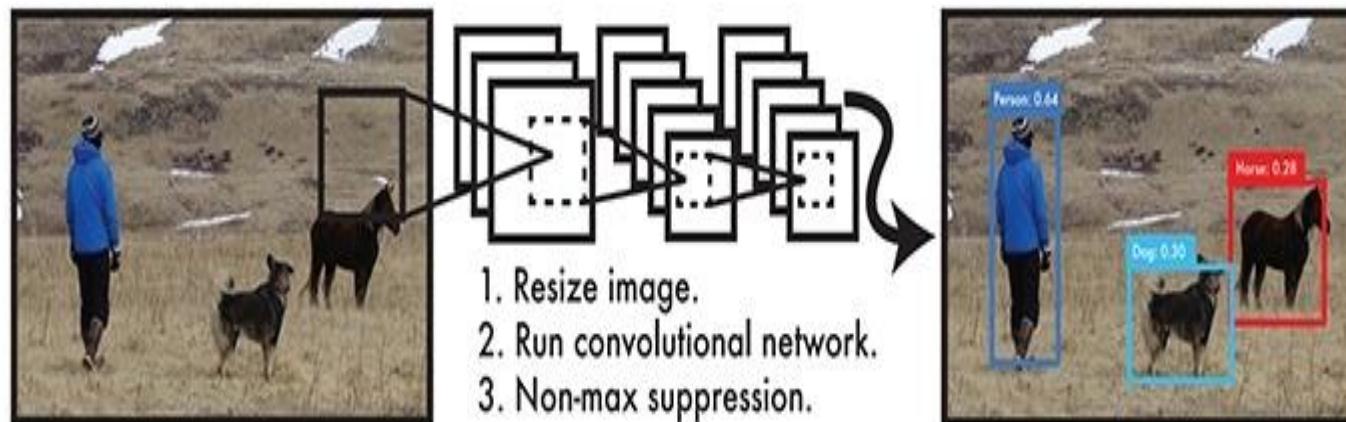


Object Detection DL models

- Deep learning approach for object detection and classification

YOLO

Shortly after that, [You Only Look Once: Unified, Real-Time Object Detection](#) (YOLO) paper published by Joseph Redmon (with Girshick appearing as one of the co-authors). YOLO proposed a simple convolutional neural network approach which has both great results and high speed, allowing for the first time real time object detection.



YOLO Architecture.

Redmon, Joseph, et al. "You only look once: Unified, real-time object detection." 2016.

Object Detection DL models

- **Deep learning approach for object detection and classification**

Faster R-CNN

Subsequently, Faster R-CNN authored by Shaoqing Ren (also co-authored by Girshick, now at Facebook Research), **the third iteration of the R-CNN series.**

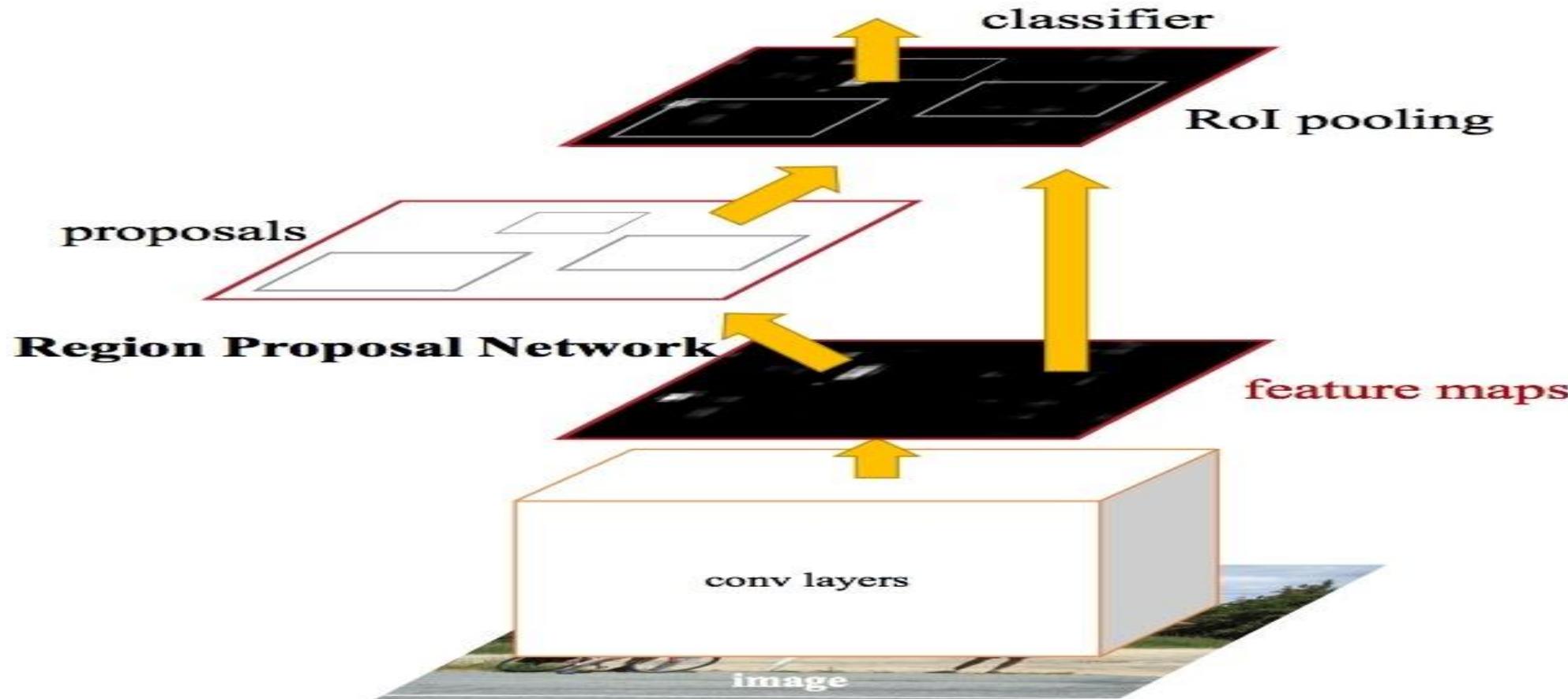
Faster R-CNN added what they called a Region Proposal Network (RPN), in an attempt to get rid of the Selective Search algorithm and make the model completely trainable end-to-end.

It has the task to output objects based **on an “objectness” score.**

These objects are used by the ROI Pooling and fully connected layers for classification.

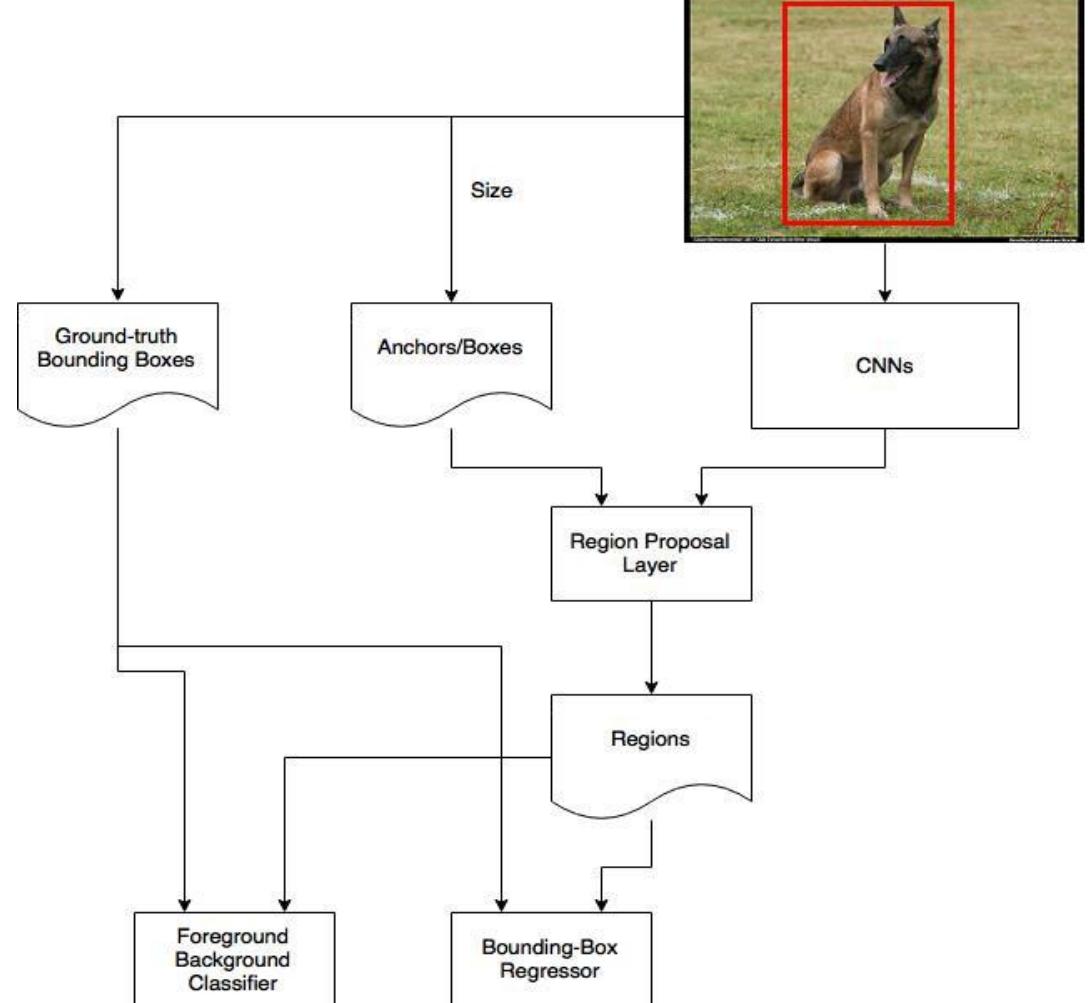
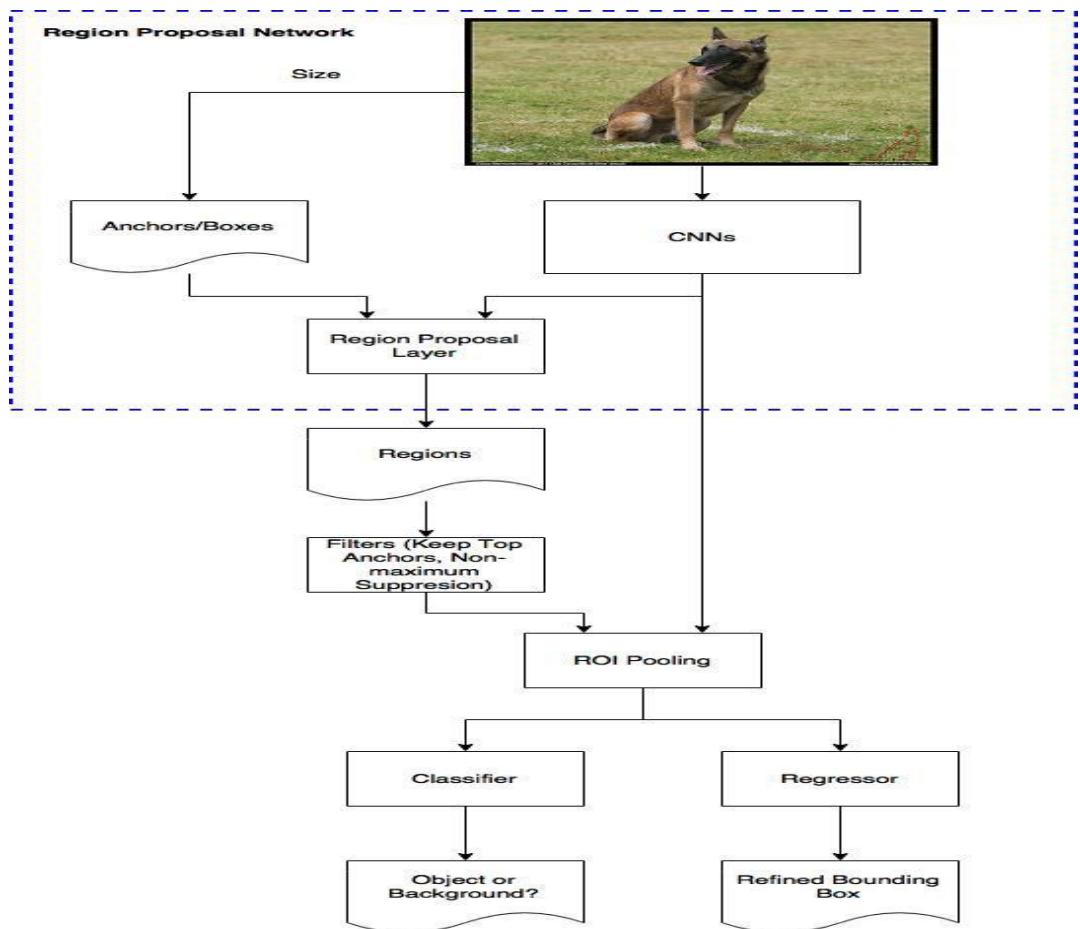
Object Detection DL models

- Deep learning approach for object detection and classification



Object Detection DL models

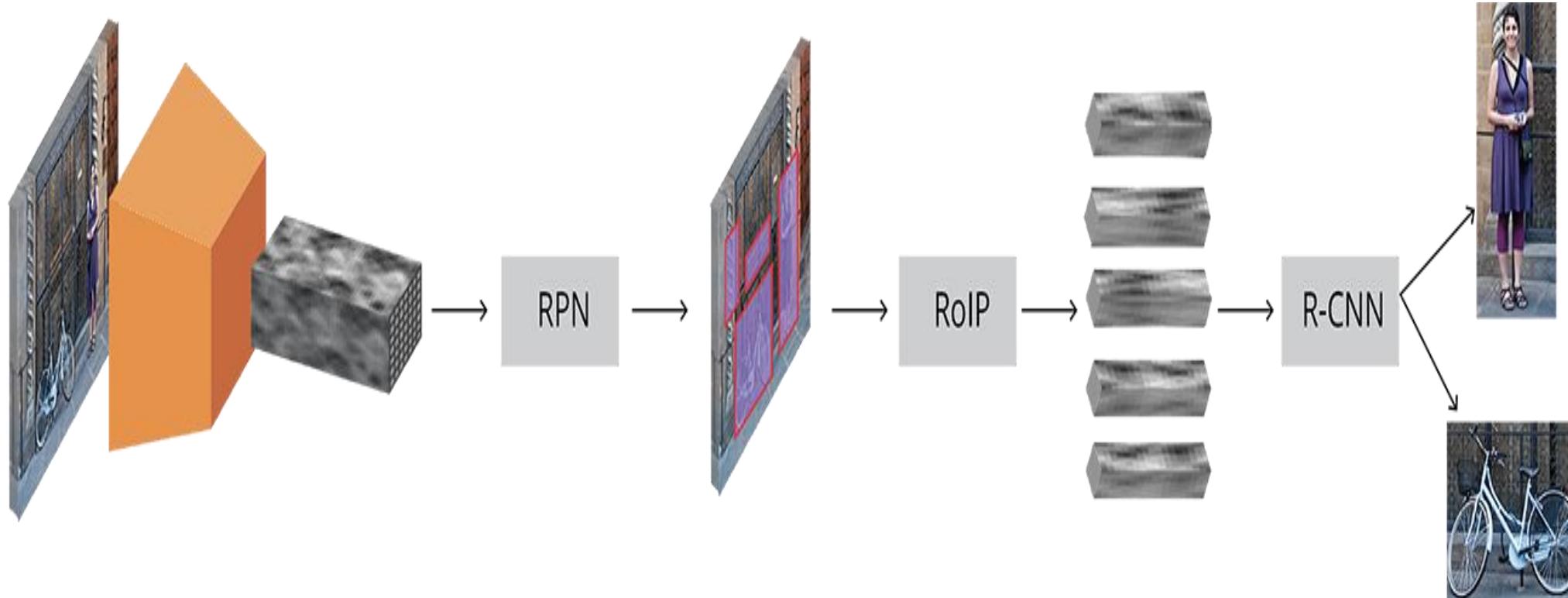
- Deep learning approach for object detection and classification



<https://medium.com/@smallfishbigsea/faster-r-cnn-explained-864d4fb7e3f8>

Object Detection DL models

- Deep learning approach for object detection and classification



Complete Faster R-CNN architecture

<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

Object Detection DL models

- Deep learning approach for object detection and classification

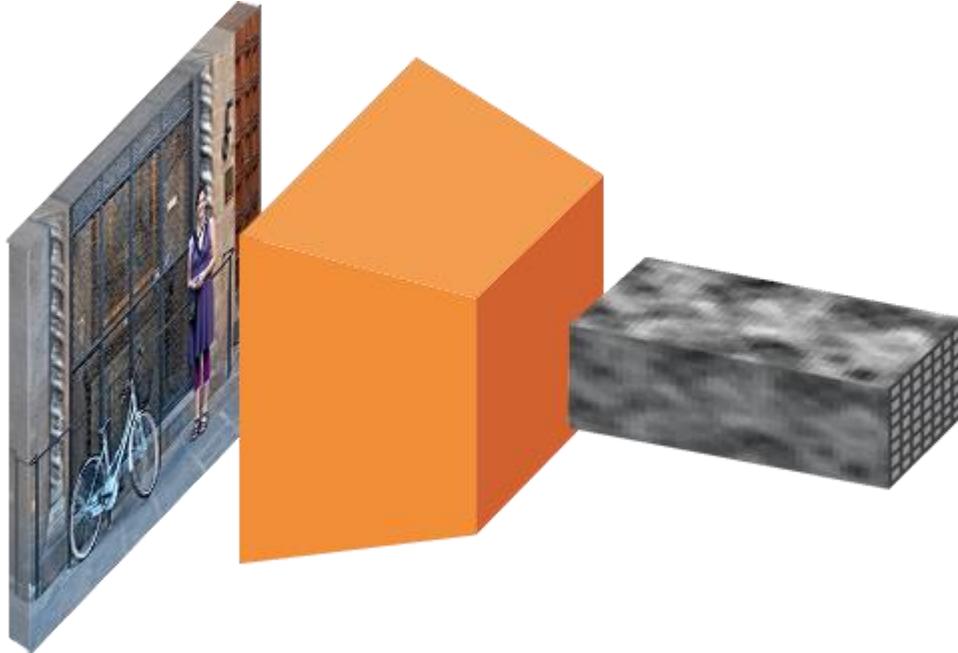
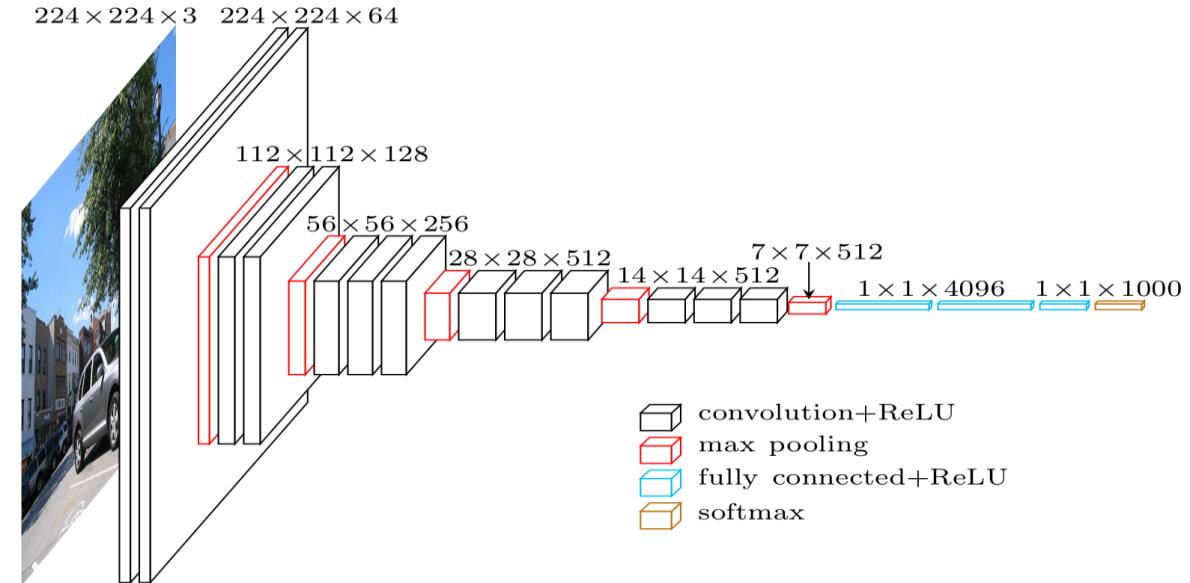


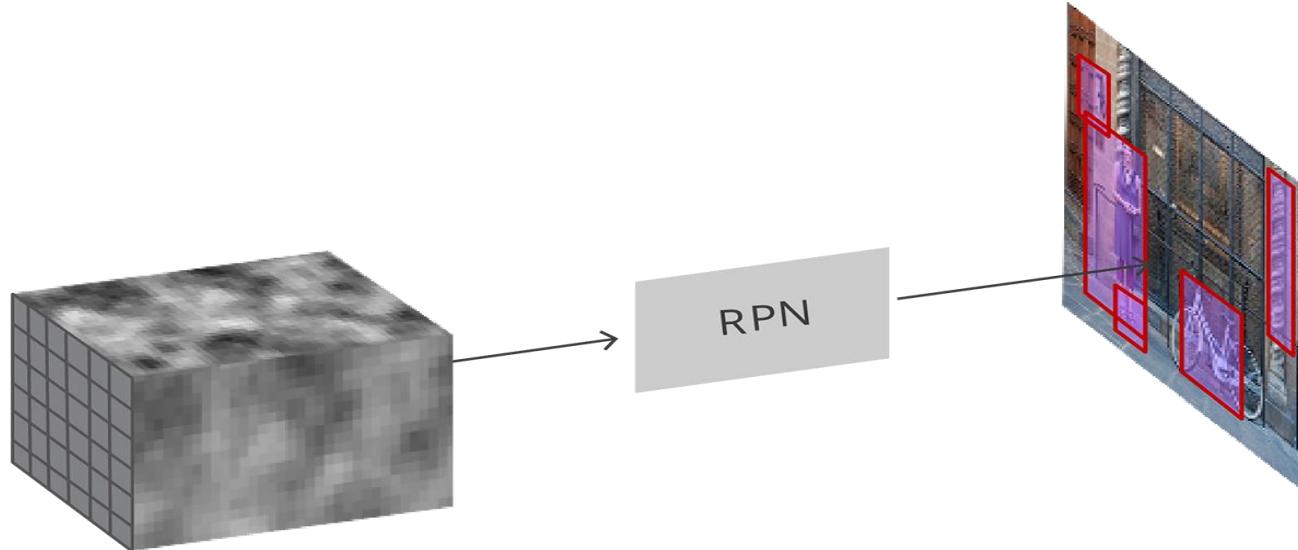
Image to convolutional feature map



VGG architecture

Object Detection DL models

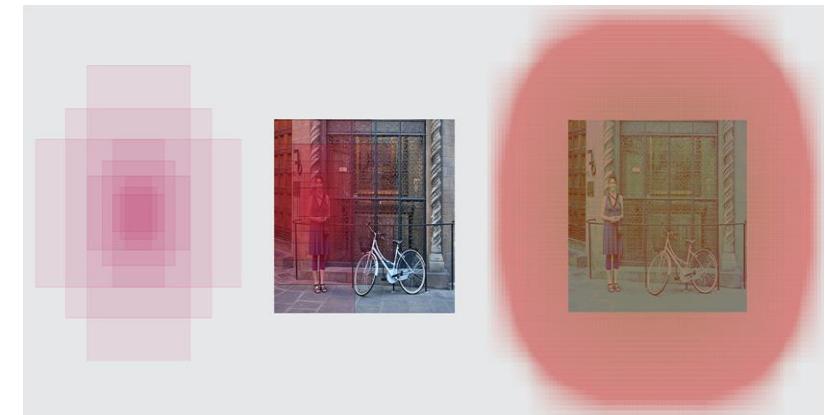
- Deep learning approach for object detection and classification



The RPN takes the convolutional feature map and generates proposals over the image



Anchor centers through the original image

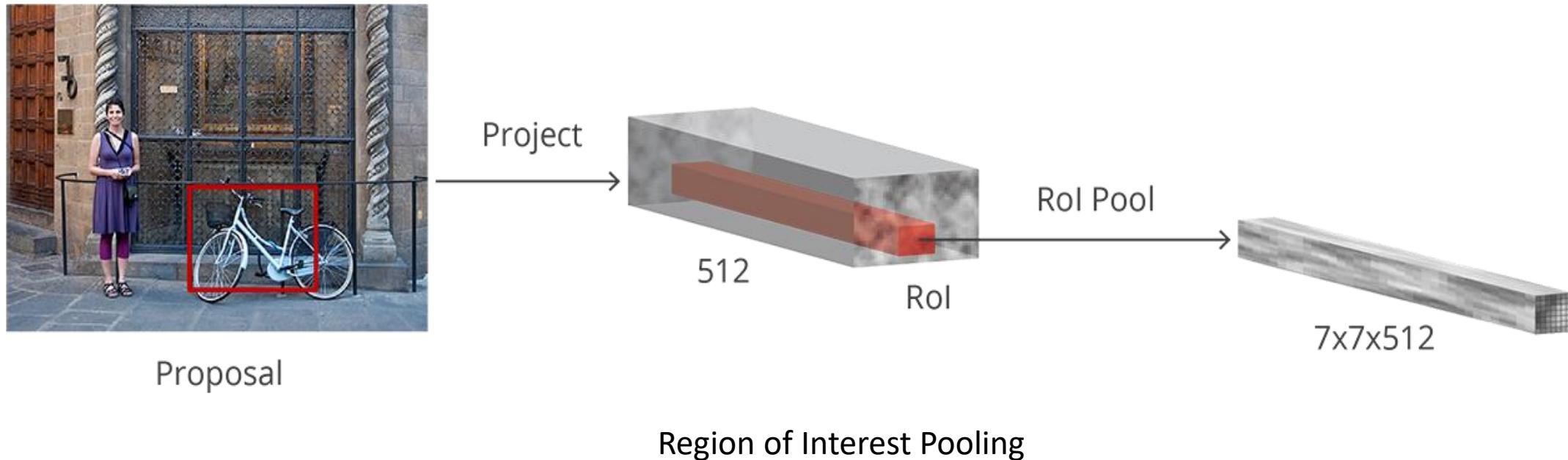


Left: Anchors, Center: Anchor for a single point, Right: All anchors

Object Detection DL models

- Deep learning approach for object detection and classification

Region of Interest Pooling



<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

Object Detection DL models

- Deep learning approach for object detection and classification

Region-based Convolutional Neural Network

Region-based convolutional neural network (R-CNN) is the final step in Faster R-CNN's pipeline. After getting a convolutional feature map from the image, using it to get object proposals with the RPN and finally extracting features for each of those proposals (via RoI Pooling),

we finally need to use these features for classification.

R-CNN tries to mimic the final stages of classification CNNs where a fully-connected layer is used to output a score for each possible object class.

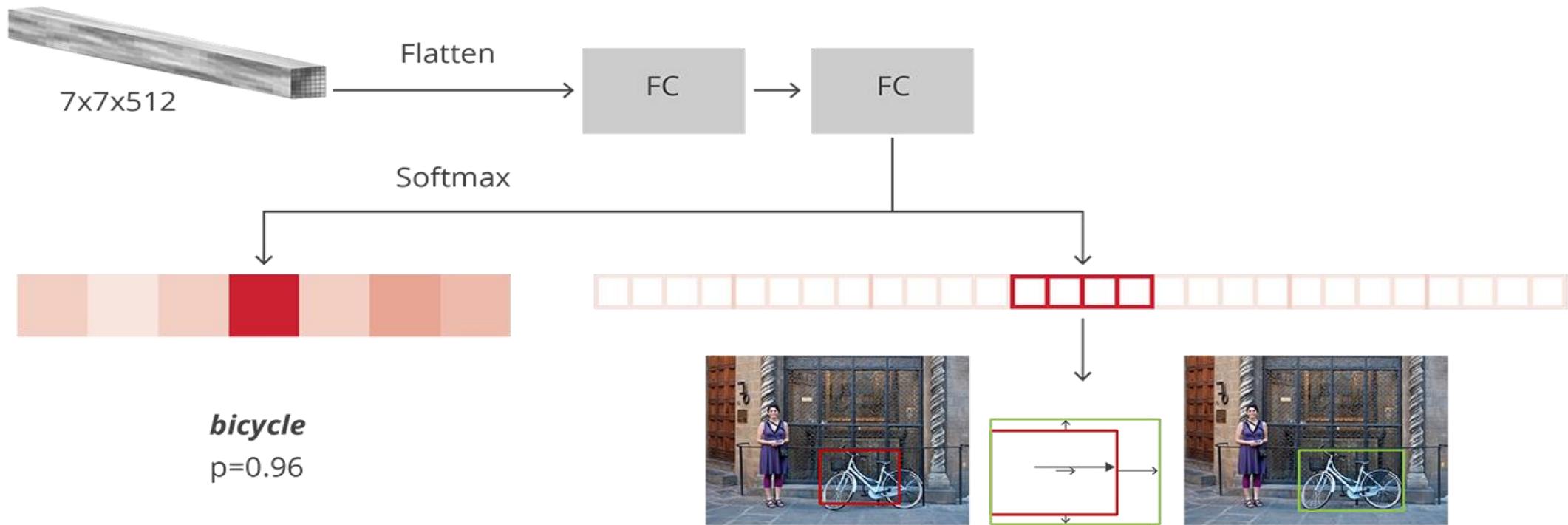
R-CNN has two different goals:

Classify proposals into one of the classes, plus a background class (for removing bad proposals).
Better adjust the bounding box for the proposal according to the predicted class.

Object Detection DL models

- Deep learning approach for object detection and classification

Region-based Convolutional Neural Network



bicycle
p=0.96

<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

Object Detection DL models

- Deep learning approach for object detection and classification

Region-based Convolutional Neural Network

In the original paper, Faster R-CNN was trained using a multi-step approach, training parts independently and merging the trained weights before a final full training approach. Since then, it has been found that doing end-to-end, joint training leads to better results.

After putting the complete model together we end up with 4 different losses, two for the RPN and two for R-CNN. We have the trainable layers in RPN and R-CNN, and we also have the base network which we can train (fine-tune) or not.

Object Detection DL models

- Deep learning approach for object detection and classification

SSD and R-FCN

Finally, there are two notable papers, Single Shot Detector (SSD) which takes on YOLO by using multiple sized convolutional feature maps achieving better results and speed, and Region-based Fully Convolutional Networks (R-FCN) which takes the architecture of Faster R-CNN but with only convolutional networks.

Keras models

▪ Fine Tuning Off-the-shelf Pre-trained Models

Model	Size	Top-1 Accuracy	Top-5 Accuracy	Parameters	Depth
Xception	Available models		0.790	0.945	22,910,480
VGG16		0.713	0.901	138,357,544	23
VGG19		0.713	0.900	143,667,240	26
ResNet50	98 MB	0.749	0.921	25,636,712	-
ResNet101	171 MB	0.764	0.928	44,707,176	-
ResNet152	232 MB	0.766	0.931	60,419,944	-
ResNet50V2	98 MB	0.760	0.930	25,613,800	-
ResNet101V2	171 MB	0.772	0.938	44,675,560	-
ResNet152V2	232 MB	0.780	0.942	60,380,648	-
InceptionV3	92 MB	0.779	0.937	23,851,784	159
InceptionResNetV2	215 MB	0.803	0.953	55,873,736	572
MobileNet	16 MB	0.704	0.895	4,253,864	88
MobileNetV2	14 MB	0.713	0.901	3,538,984	88
DenseNet121	33 MB	0.750	0.923	8,062,504	121
DenseNet169	57 MB	0.762	0.932	14,307,880	169
DenseNet201	80 MB	0.773	0.936	20,242,984	201
NASNetMobile	23 MB	0.744	0.919	5,326,716	-
NASNetLarge	343 MB	0.825	0.960	88,949,818	-
EfficientNetB0	29 MB	-	-	5,330,571	-
EfficientNetB1	31 MB	-	-	7,856,239	-
EfficientNetB2	36 MB	-	-	9,177,569	-
EfficientNetB3	48 MB	-	-	12,320,535	-
EfficientNetB4	75 MB	-	-	19,466,823	-
EfficientNetB5	118 MB	-	-	30,562,527	-
EfficientNetB6	166 MB	-	-	43,265,143	-
EfficientNetB7	256 MB	-	-	66,658,687	-

Keras models

- **Training Deep Learning Models in keras**

Available dataset preprocessing utilities

Image data preprocessing

- image_dataset_from_directory function
- load_img function
- img_to_array function
- ImageDataGenerator class
- flow method
- flow_from_dataframe method
- flow_from_directory method

Timeseries data preprocessing

- timeseries_dataset_from_array function
- pad_sequences function
- TimeseriesGenerator class

Text data preprocessing

- text_dataset_from_directory function
- Tokenizer class

Keras models

■ Training Deep Learning Models in keras

Example of using `.flow_from_directory(directory)`:

```
train_datagen = ImageDataGenerator(  
    rescale=1./255,  
    shear_range=0.2,  
    zoom_range=0.2,  
    horizontal_flip=True)  
test_datagen = ImageDataGenerator(rescale=1./255)  
train_generator = train_datagen.flow_from_directory(  
    'data/train',  
    target_size=(150, 150),  
    batch_size=32,  
    class_mode='binary')  
validation_generator = test_datagen.flow_from_directory(  
    'data/validation',  
    target_size=(150, 150),  
    batch_size=32,  
    class_mode='binary')  
model.fit(  
    train_generator,  
    steps_per_epoch=2000,  
    epochs=50,  
    validation_data=validation_generator,  
    validation_steps=800)
```

Keras models

▪ Training Deep Learning Models in keras

ImageDataGenerator class

```
tf.keras.preprocessing.image.ImageDataGenerator(  
    featurewise_center=False,  
    samplewise_center=False,  
    featurewise_std_normalization=False,  
    samplewise_std_normalization=False,  
    zca_whitening=False,  
    zca_epsilon=1e-06,  
    rotation_range=0,  
    width_shift_range=0.0,  
    height_shift_range=0.0,  
    brightness_range=None,  
    shear_range=0.0,  
    zoom_range=0.0,  
    channel_shift_range=0.0,  
    fill_mode="nearest",  
    cval=0.0,  
    horizontal_flip=False,  
    vertical_flip=False,  
    rescale=None,  
    preprocessing_function=None,  
    data_format=None,  
    validation_split=0.0,  
    dtype=None,  
)
```

Generate batches of tensor image data with real-time data augmentation.

Keras models

- **Training Deep Learning Models in keras**

flow_from_directory method

```
ImageDataGenerator.flow_from_directory(  
    directory,  
    target_size=(256, 256),  
    color_mode="rgb",  
    classes=None,  
    class_mode="categorical",  
    batch_size=32,  
    shuffle=True,  
    seed=None,  
    save_to_dir=None,  
    save_prefix="",  
    save_format="png",  
    follow_links=False,  
    subset=None,  
    interpolation="nearest",  
)
```

Takes the path to a directory & generates batches of augmented data.

Keras models

▪ **Training Deep Learning Models in keras**

Dataset preprocessing

Keras dataset preprocessing utilities, located at `tf.keras.preprocessing`, help you go from raw data on disk to a `tf.data.Dataset` object that can be used to train a model.

Here's a quick example: let's say you have 10 folders, each containing 10,000 images from a different category, and you want to train a classifier that maps an image to its category.

Your training data folder would look like this:

```
training_data/  
...class_a/  
.....a_image_1.jpg  
.....a_image_2.jpg  
...class_b/  
.....b_image_1.jpg  
.....b_image_2.jpg  
etc.
```

You may also have a validation data folder `validation_data/` structured in the same way.

You could simply do:

Keras models

■ Training Deep Learning Models in keras

Dataset preprocessing

Keras dataset preprocessing utilities, located at `tf.keras.preprocessing`, help you go from raw data on disk to a `tf.data.Dataset` object that can be used to train a model.

You could simply do:

```
from tensorflow import keras
from tensorflow.keras.preprocessing.image import image_dataset_from_directory

train_ds = image_dataset_from_directory(
    directory='training_data/',
    labels='inferred',
    label_mode='categorical',
    batch_size=32,
    image_size=(256, 256))
validation_ds = image_dataset_from_directory(
    directory='validation_data/',
    labels='inferred',
    label_mode='categorical',
    batch_size=32,
    image_size=(256, 256))

model = keras.applications.Xception(weights=None, input_shape=(256, 256, 3),
classes=10)
model.compile(optimizer='rmsprop', loss='categorical_crossentropy')
model.fit(ds, epochs=10, validation_data=validation_ds)
```

Keras models

■ Training Deep Learning Models in keras

Optimizers

Usage with `compile()` & `fit()`

An optimizer is one of the two arguments required for compiling a Keras model:

```
from tensorflow import keras
from tensorflow.keras import layers

model = keras.Sequential()
model.add(layers.Dense(64, kernel_initializer='uniform', input_shape=(10,)))
model.add(layers.Activation('softmax'))

opt = keras.optimizers.Adam(learning_rate=0.01)
model.compile(loss='categorical_crossentropy', optimizer=opt)
```

You can either instantiate an optimizer before passing it to `model.compile()`, as in the above example, or you can pass it by its string identifier. In the latter case, the default parameters for the optimizer will be used.

```
# pass optimizer by name: default parameters will be used
model.compile(loss='categorical_crossentropy', optimizer='adam')
```

Keras models

- **Training Deep Learning Models in keras**

Available optimizers

- [SGD](#)
- [RMSprop](#)
- [Adam](#)
- [Adadelta](#)
- [Adagrad](#)
- [Adamax](#)
- [Nadam](#)
- [Ftrl](#)

Keras models

<https://keras.io/api/applications/>

■ Training Deep Learning Models in keras

Available metrics

Accuracy metrics

- [Accuracy class](#)
- [BinaryAccuracy class](#)
- [CategoricalAccuracy class](#)
- [TopKCategoricalAccuracy class](#)
- [SparseTopKCategoricalAccuracy class](#)

Probabilistic metrics

- [BinaryCrossentropy class](#)
- [CategoricalCrossentropy class](#)
- [SparseCategoricalCrossentropy class](#)
- [KLDivergence class](#)
- [Poisson class](#)

Regression metrics

- [MeanSquaredError class](#)
- [RootMeanSquaredError class](#)
- [MeanAbsoluteError class](#)
- [MeanAbsolutePercentageError class](#)
- [MeanSquaredLogarithmicError class](#)
- [CosineSimilarity class](#)
- [LogCoshError class](#)

Classification metrics based on True/False positives & negatives

- [AUC class](#)
- [Precision class](#)
- [Recall class](#)
- [TruePositives class](#)
- [TrueNegatives class](#)
- [FalsePositives class](#)
- [FalseNegatives class](#)
- [PrecisionAtRecall class](#)
- [SensitivityAtSpecificity class](#)
- [SpecificityAtSensitivity class](#)

Image segmentation metrics

- [MeanIoU class](#)

Hinge metrics for "maximum-margin" classification

- [Hinge class](#)
- [SquaredHinge class](#)
- [CategoricalHinge class](#)

Keras models

<https://keras.io/api/applications/>

■ Training Deep Learning Models in keras

Usage with `compile()` & `fit()`

The `compile()` method takes a `metrics` argument, which is a list of metrics:

```
model.compile(  
    optimizer='adam',  
    loss='mean_squared_error',  
    metrics=[  
        metrics.MeanSquaredError(),  
        metrics.AUC(),  
    ]  
)
```

Metric values are displayed during `fit()` and logged to the `History` object returned by `fit()`. They are also returned by `model.evaluate()`.

Note that the best way to monitor your metrics during training is via [TensorBoard](#).

To track metrics under a specific name, you can pass the `name` argument to the metric constructor:

```
model.compile(  
    optimizer='adam',  
    loss='mean_squared_error',  
    metrics=[  
        metrics.MeanSquaredError(name='my_mse'),  
        metrics.AUC(name='my_auc'),  
    ]  
)
```

Keras models

<https://keras.io/api/applications/>

■ Training Deep Learning Models in keras

Losses

The purpose of loss functions is to compute the quantity that a model should seek to minimize during training.

Probabilistic losses

- [BinaryCrossentropy class](#)
- [CategoricalCrossentropy class](#)
- [SparseCategoricalCrossentropy class](#)
- [Poisson class](#)
- [binary_crossentropy function](#)
- [categorical_crossentropy function](#)
- [sparse_categorical_crossentropy function](#)
- [poisson function](#)
- [KLDivergence class](#)
- [kl_divergence function](#)

Regression losses

- [MeanSquaredError class](#)
- [MeanAbsoluteError class](#)
- [MeanAbsolutePercentageError class](#)
- [MeanSquaredLogarithmicError class](#)
- [CosineSimilarity class](#)
- [mean_squared_error function](#)
- [mean_absolute_error function](#)
- [mean_absolute_percentage_error function](#)
- [mean_squared_logarithmic_error function](#)
- [cosine_similarity function](#)
- [Huber class](#)
- [huber function](#)
- [LogCosh class](#)
- [log_cosh function](#)

Hinge losses for "maximum-margin" classification

- [Hinge class](#)
- [SquaredHinge class](#)
- [CategoricalHinge class](#)
- [hinge function](#)
- [squared_hinge function](#)
- [categorical_hinge function](#)

Keras models

<https://keras.io/api/applications/>

■ Training Deep Learning Models in keras

Losses

The purpose of loss functions is to compute the quantity that a model should seek to minimize during training.

Usage of losses with `compile()` & `fit()`

A loss function is one of the two arguments required for compiling a Keras model:

```
from tensorflow import keras
from tensorflow.keras import layers

model = keras.Sequential()
model.add(layers.Dense(64, kernel_initializer='uniform', input_shape=(10,)))
model.add(layers.Activation('softmax'))

loss_fn = keras.losses.SparseCategoricalCrossentropy()
model.compile(loss=loss_fn, optimizer='adam')
```

All built-in loss functions may also be passed via their string identifier:

```
# pass optimizer by name: default parameters will be used
model.compile(loss='sparse_categorical_crossentropy', optimizer='adam')
```

Keras models

<https://keras.io/api/applications/>

■ Training Deep Learning Models in keras

Losses

The purpose of loss functions is to compute the quantity that a model should seek to minimize during training.

Callbacks API

A callback is an object that can perform actions at various stages of training (e.g. at the start or end of an epoch, before or after a single batch, etc).

You can use callbacks to:

- Write TensorBoard logs after every batch of training to monitor your metrics
- Periodically save your model to disk
- Do early stopping
- Get a view on internal states and statistics of a model during training
- ...and more

Usage of callbacks via the built-in `fit()` loop

You can pass a list of callbacks (as the keyword argument `callbacks`) to the `.fit()` method of a model:

```
my_callbacks = [
    tf.keras.callbacks.EarlyStopping(patience=2),
    tf.keras.callbacks.ModelCheckpoint(filepath='model.{epoch:02d}-
{val_loss:.2f}.h5'),
    tf.keras.callbacks.TensorBoard(log_dir='./logs'),
]
model.fit(dataset, epochs=10, callbacks=my_callbacks)
```

The relevant methods of the callbacks will then be called at each stage of the training.

Keras models

<https://keras.io/api/applications/>

■ **Training Deep Learning Models in keras**

Losses

The purpose of loss functions is to compute the quantity that a model should seek to minimize during training.

Available callbacks

- [Base Callback class](#)
- [ModelCheckpoint](#)
- [TensorBoard](#)
- [EarlyStopping](#)
- [LearningRateScheduler](#)
- [ReduceLROnPlateau](#)
- [RemoteMonitor](#)
- [LambdaCallback](#)
- [TerminateOnNaN](#)
- [CSVLogger](#)
- [ProgbarLogger](#)

Keras models

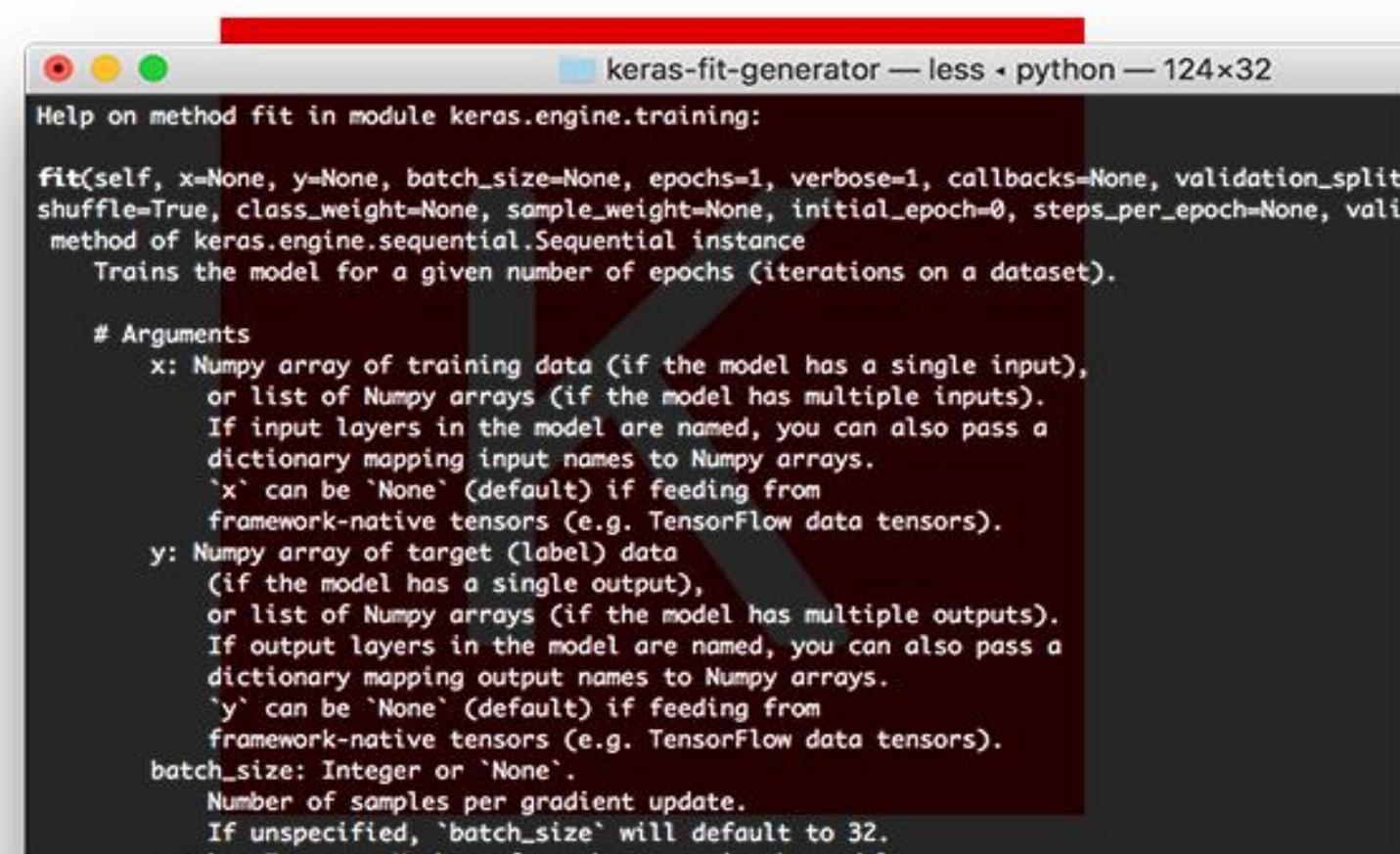
■ Training Deep Learning Models in keras

https://www.pyimagesearch.com/2018/12/24/how-to-use-keras-fit-and-fit_generator-a-hands-on-tutorial/

When to use Keras' fit, fit_generator, and train_on_batch functions?

Keras provides three functions that can be used to train your own deep learning models:

- 1..fit
- 2..fit_generator
- 3..train_on_batch



The screenshot shows a terminal window titled "keras-fit-generator — less - python — 124x32". The content is the help documentation for the `fit` method in the `keras.engine.training` module. It describes the `fit` method as a class method of `keras.engine.sequential.Sequential` instances, which trains the model for a given number of epochs (iterations on a dataset). It details the arguments `x`, `y`, and `batch_size`, explaining their types and behaviors.

```
Help on method fit in module keras.engine.training:  
  
fit(self, x=None, y=None, batch_size=None, epochs=1, verbose=1, callbacks=None, validation_split=  
      shuffle=True, class_weight=None, sample_weight=None, initial_epoch=0, steps_per_epoch=None, vali  
      method of keras.engine.sequential.Sequential instance  
      Trains the model for a given number of epochs (iterations on a dataset).  
  
# Arguments  
    x: Numpy array of training data (if the model has a single input),  
       or list of Numpy arrays (if the model has multiple inputs).  
       If input layers in the model are named, you can also pass a  
       dictionary mapping input names to Numpy arrays.  
       'x' can be 'None' (default) if feeding from  
       framework-native tensors (e.g. TensorFlow data tensors).  
    y: Numpy array of target (label) data  
       (if the model has a single output),  
       or list of Numpy arrays (if the model has multiple outputs).  
       If output layers in the model are named, you can also pass a  
       dictionary mapping output names to Numpy arrays.  
       'y' can be 'None' (default) if feeding from  
       framework-native tensors (e.g. TensorFlow data tensors).  
    batch_size: Integer or 'None'.  
       Number of samples per gradient update.  
       If unspecified, 'batch_size' will default to 32.
```

Keras models

■ Training Deep Learning Models in keras

Let's start with a call to

.fit

:How to use Keras fit and

fit_generator (a hands-on tutorial)

model.fit(trainX, trainY, batch_size=32, epochs=50)

Here you can see that we are supplying our training data (trainX) and training labels (trainY).

We then instruct Keras to allow our model to train for 50 epochs with a batch size of 32

.The call to

.fit

1. is making two primary assumptions here:Our *entire* training set can fit into RAM

2. There is *no* data augmentation going on (i.e., there is no need for Keras generators)

Instead, our network will be trained on the raw data.

The raw data itself will fit into memory — we have no need to move old batches of data out of RAM and move new batches of data into RAM.

Furthermore, we will not be manipulating the training data on the fly using data augmentation.

https://www.pyimagesearch.com/2018/12/24/how-to-use-keras-fit-and-fit_generator-a-hands-on-tutorial/

Keras models

https://www.pyimagesearch.com/2018/12/24/how-to-use-keras-fit-and-fit_generator-a-hands-on-tutorial/

■ Training Deep Learning Models in keras

For small, simplistic datasets it's perfectly acceptable to use Keras' `.fit` function.

These datasets are often not very challenging and do not require any data augmentation.

However, real-world datasets are rarely that simple:

- Real-world datasets are often *too large to fit into memory*.
- They also tend to be *challenging*, requiring us to perform data augmentation to avoid overfitting and increase the ability of our model to generalize.

In those situations we need to utilize Keras' `.fit_generator`

function:How to use Keras fit and fit_generator (a hands-on tutorial)

```
# initialize the number of epochs and batch size
```

```
EPOCHS = 100
```

```
BS = 32
```

```
# construct the training image generator for data augmentation
```

```
aug = ImageDataGenerator(rotation_range=20, zoom_range=0.15,  
width_shift_range=0.2, height_shift_range=0.2, shear_range=0.15,  
horizontal_flip=True, fill_mode="nearest")
```

```
# train the network
```

```
H = model.fit_generator(aug.flow(trainX, trainY, batch_size=BS),  
validation_data=(testX, testY), steps_per_epoch=len(trainX) // BS,  
epochs=EPOCHS)
```

Keras models

■ Training Deep Learning Models in keras

Here we start by first initializing the number of epochs we are going to train our network for along with the batch size. We then initialize `aug`, a Keras `ImageDataGenerator` object that is used to apply data augmentation, randomly translating, rotating, resizing, etc. images on the fly. Performing data augmentation is a form of regularization, enabling our model to generalize better. However, applying data augmentation implies that our training data is no longer “static” — the data is constantly changing.

Each new batch of data is randomly adjusted according to the parameters supplied to

`ImageDataGenerator`. Thus, we now need to utilize Keras’

`.fit_generator` function to train our model. As the name suggests, the `.fit_generator` function assumes there is an underlying function that is *generating* the data for it.

Internally, Keras is using the following process when training a model with `.fit_generator`

1. Keras calls the generator function supplied to `.fit_generator` (in this case, `aug.flow`).

2. The generator function yields a batch of size BS to the `.fit_generator` function.

3. The `.fit_generator` function accepts the batch of data, performs backpropagation, and updates the weights in our model.

4. This process is repeated until we have reached the desired number of epochs.

You’ll notice we now need to supply a `steps_per_epoch` parameter when calling `.fit_generator` (the `.fit` method had no such parameter).

https://www.pyimagesearch.com/2018/12/24/how-to-use-keras-fit-and-fit_generator-a-hands-on-tutorial/

Keras models

- **Training Deep Learning Models in keras**

https://www.pyimagesearch.com/2018/12/24/how-to-use-keras-fit-and-fit_generator-a-hands-on-tutorial/

The `train_on_batch`

function accepts a *single batch of data*, performs backpropagation, and then updates the model parameters.

The batch of data can be of arbitrary size (i.e., it does not require an explicit batch size to be provided).

The data itself can be generated however you like as well. This data could be raw images on disk or data that has been modified or augmented in some manner.

You'll typically use the `.train_on_batch` function when you have *very explicit* reasons for wanting to maintain your own training data iterator, such as the data iteration process being extremely complex and requiring custom code. If you find yourself asking if you need the `.train_on_batch` function then in all likelihood you probably don't.

Transfer Learning DL models

▪ Fine Tuning Off-the-shelf Pre-trained Models

TORCHVISION.MODELS

The models subpackage contains definitions of models for addressing different tasks, including: image classification, pixelwise semantic segmentation, object detection, instance segmentation, person keypoint detection and video classification.

Classification

The models subpackage contains definitions for the following model architectures for image classification:

- [AlexNet](#)
- [VGG](#)
- [ResNet](#)
- [SqueezeNet](#)
- [DenseNet](#)
- [Inception v3](#)
- [GoogLeNet](#)
- [ShuffleNet v2](#)
- [MobileNet v2](#)
- [ResNeXt](#)
- [Wide ResNet](#)
- [MNASNet](#)

Transfer Learning DL models

- **Fine Tuning Off-the-shelf Pre-trained Models**
- You can construct a model with random weights by calling its constructor:

```
import torchvision.models as models  
resnet18 = models.resnet18()  
alexnet = models.alexnet()  
vgg16 = models.vgg16()  
squeezenet = models.squeezezenet1_0()  
densenet = models.densenet161()  
inception = models.inception_v3()  
googlenet = models.googlenet()  
shufflenet = models.shufflenet_v2_x1_0()  
mobilenet = models.mobilenet_v2()  
resnext50_32x4d = models.resnext50_32x4d()  
wide_resnet50_2 = models.wide_resnet50_2()  
mnasnet = models.mnasnet1_0()
```

Transfer Learning DL models

- **Fine Tuning Off-the-shelf Pre-trained Models**
- We provide pre-trained models, using the PyTorch torch.utils.model_zoo. These can be constructed by passing pretrained=True:

```
import torchvision.models as models
resnet18 = models.resnet18(pretrained=True)
alexnet = models.alexnet(pretrained=True)
squeezenet = models.squeeze1_0(pretrained=True)
vgg16 = models.vgg16(pretrained=True)
densenet = models.densenet161(pretrained=True)
inception = models.inception_v3(pretrained=True)
googlenet = models.googlenet(pretrained=True)
shufflenet = models.shufflenet_v2_x1_0(pretrained=True)
mobilenet = models.mobilenet_v2(pretrained=True)
resnext50_32x4d = models.resnext50_32x4d(pretrained=True)
wide_resnet50_2 = models.wide_resnet50_2(pretrained=True)
mnasnet = models.mnasnet1_0(pretrained=True)
```