Departamento de Eletrónica, Telecomunicações e Informática

# Complements of Machine Learning

## Lecture 4 : Deep NN - Residual Networks

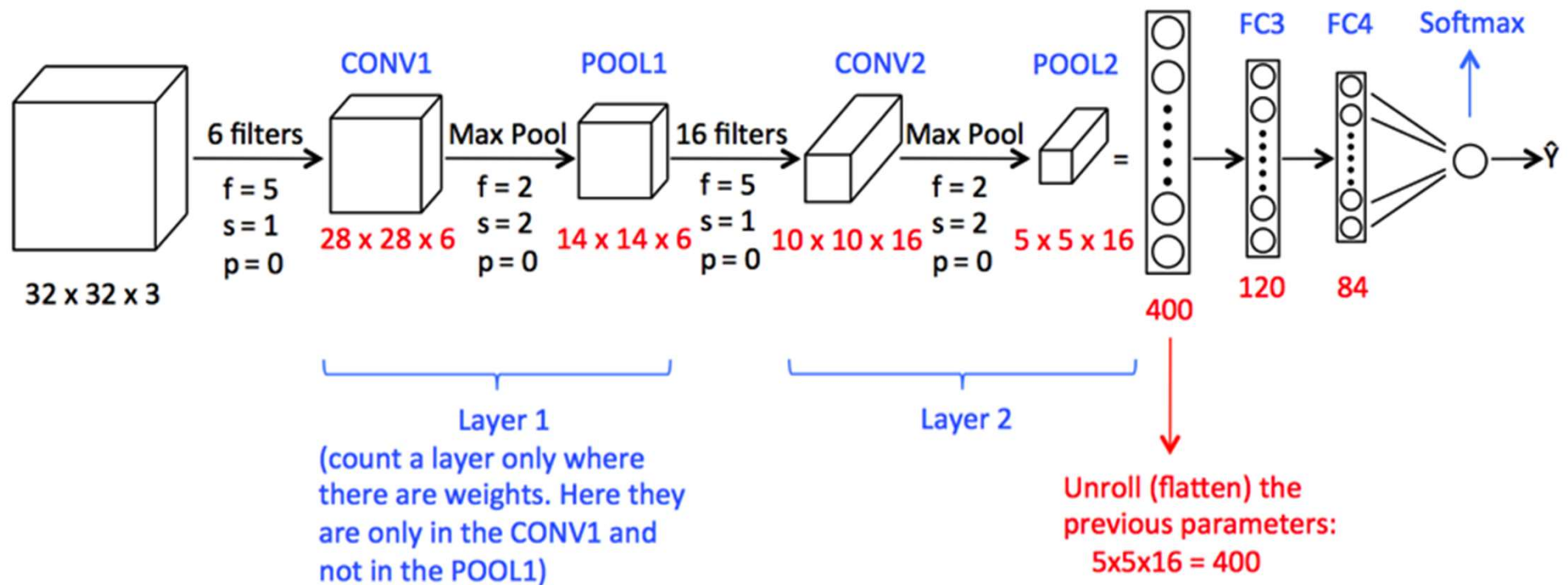### Petia Georgieva
(petia@ua.pt)

Petia Georgieva
(petia@ua.pt)

ieeta

universidade
de aveiro

# Outline

Deep NN case studies:

1. Residual Networks (ResNet)

2. 1x1 convolution

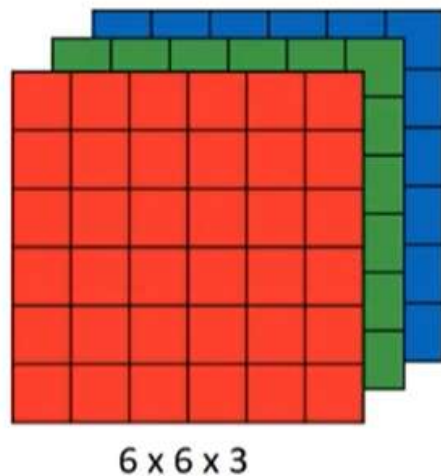3. Inception Network

universidade
de aveiro

# CNN: LeNet5*



**\*LeCun et al., 1998, "Gradient-based learning applied to document recognition".** Original LeNet5 applied to handwritten digit recognition (grey scale images). Avg pooling, no padding, softmax classifier; ReLU and sigmoid/tanh neurons in the Fully Connected (FC) layers.
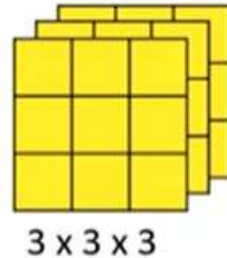
The activation function is present after the convolution, even if it is not drawn on the CNN diagram.

**General trend**: CNNs start with large image, then height and width gradually decrease as it goes deeper in the network, where as the number of channels increase.
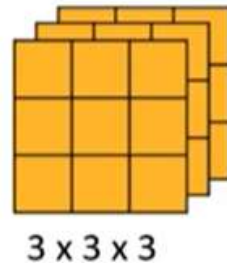
# Example of one Layer



6 x 6 x 3

\* 3 x 3 x 3

\* 3 x 3 x 3

$$\lfloor \frac{n + 2p - f}{s} + 1 \rfloor \text{ by } \lfloor \frac{n + 2p - f}{s} + 1 \rfloor$$
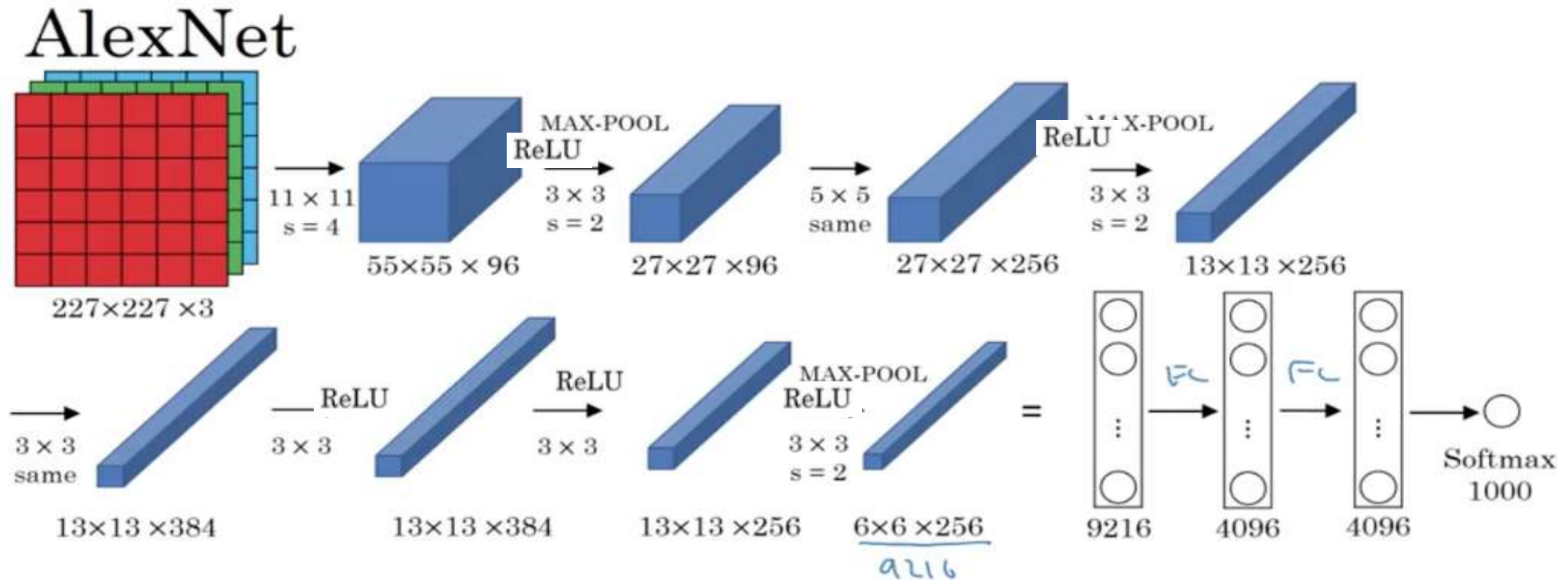
padding =0; stride =1:

Dimension of feature map (conv layer output) ?

How many trainable parameters has the layer ?

# AlexNet*



**\* Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton, 2012, ImageNet classification with deep convolutional neural networks.**
5 conv layers, 3 FC layers with softmax output, 60 million parameters in total.
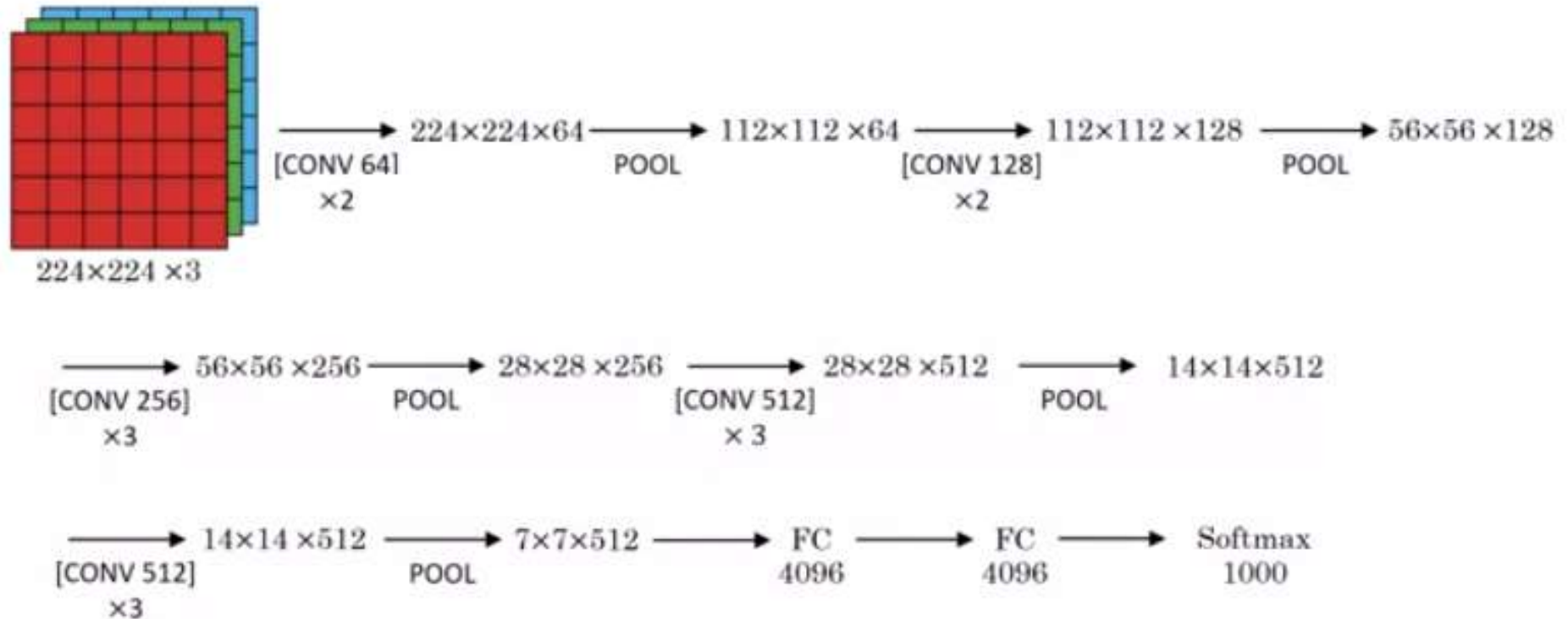AlexNet applied to ImageNet LSVRC-2010 dataset to recognize 1000 different classes.

This paper convinced the Computer Vision (CV) community that DL really works and will have a huge impact not only in CV but also in speech/language processing.

# CNN: VGG16

## VGG - 16

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

224×224 ×3
→ [CONV 64] ×2 → 224×224×64 → POOL → 112×112 ×64 → [CONV 128] ×2 → 112×112 ×128 → POOL → 56×56 ×128

→ [CONV 256] ×3 → 56×56 ×256 → POOL → 28×28 ×256 → [CONV 512] × 3 → 28×28 ×512 → POOL → 14×14×512

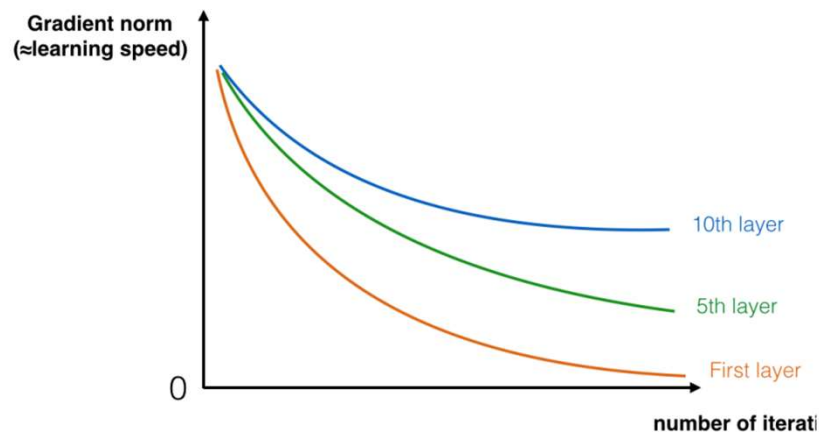→ [CONV 512] ×3 → 14×14 ×512 → POOL → 7×7×512 → FC 4096 → FC 4096 → Softmax 1000

**\* Karen Simonyan, Andrew Zisserman, 2015, Very Deep Convolutional Networks for Large-Scale Image Recognition**
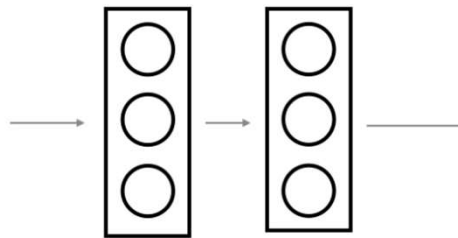VGG-16 has 16 layers (with weights), ***138 million parameters !!!***
Unified architecture: All conv layers: (3x3) filters, s=1, same; All MaxPool =2x2, s=2.
At each convolution the height and width go down by a factor of 2, the channels go up by a factor of 2.
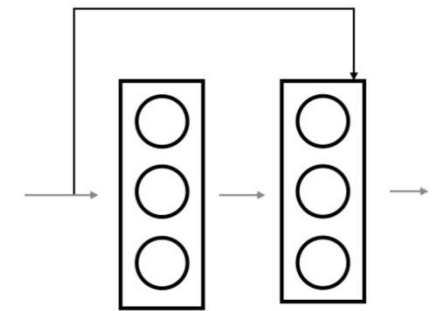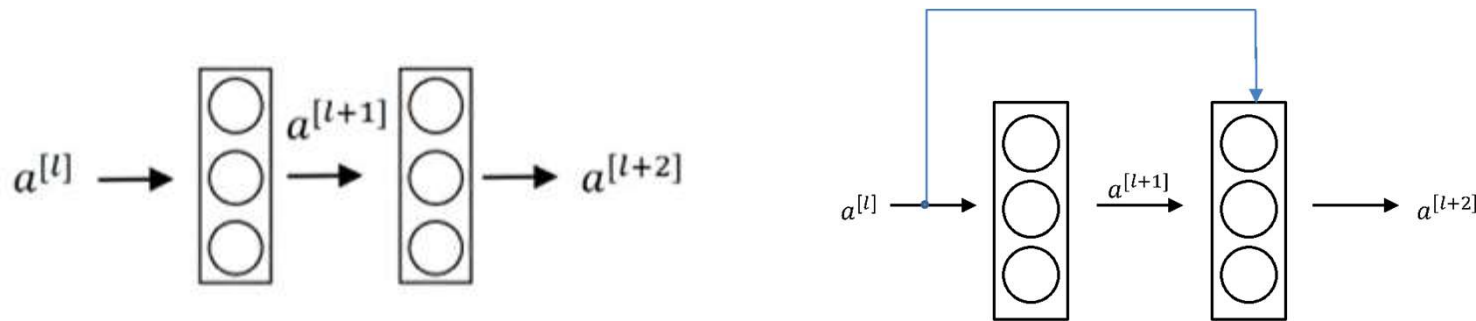
# Residual Networks



**Vanishing gradient**        **Residual block**

Very deep NNs are difficult to train because of vanishing or exploding gradient issues. During training, the magnitude of the gradient for the shallower (the first) layers decrease to 0 very rapidly as training proceeds, as shown on the figure.

**Skip connection (short cut)** takes the activation from one layer and feed it to another layer much deeper in the network.

Remark: Exploding gradient is easier to spot (NaN in computations – numerical overflow) => apply gradient clipping, i.e. limit the gradients up to some max value.

# Residual Block



$$a^{[l+2]} = g\left(W^{[l+2]}a^{[l+1]} + b^{[l+2]}\right) - \quad \text{without skip connection}$$
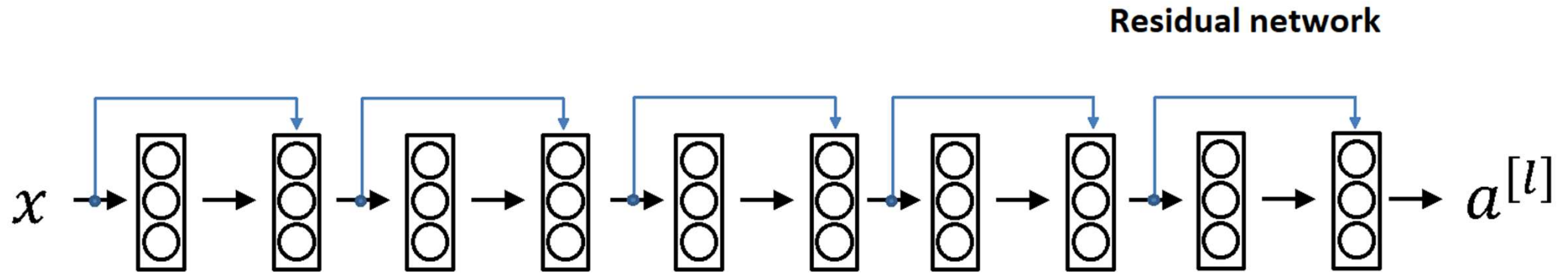
$$a^{[l+2]} = g\left(W^{[l+2]}a^{[l+1]} + b^{[l+2]} + a^{[l]}\right) - \quad \text{with identity skip connection}$$

Due to regularization (L2, drop out)  trainable parameters (*W, b*) in layer [*l*+2]
may tend to small values (close to 0) => next layers will stop learning.
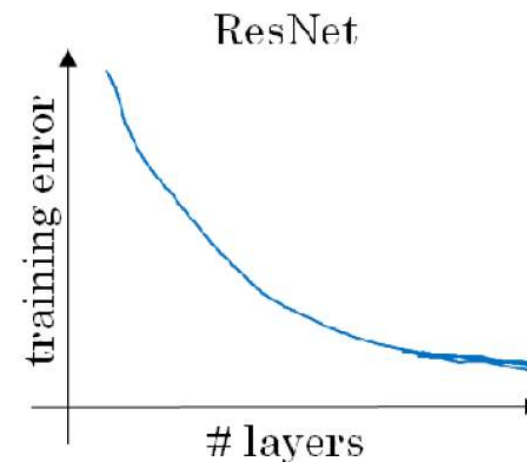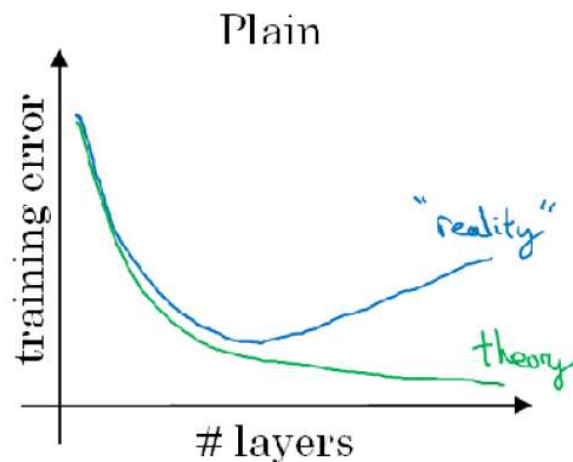With skip connection it will transmit learned patterns of layer *l*.

# Residual Network vs Plain DNN

**Residual network**
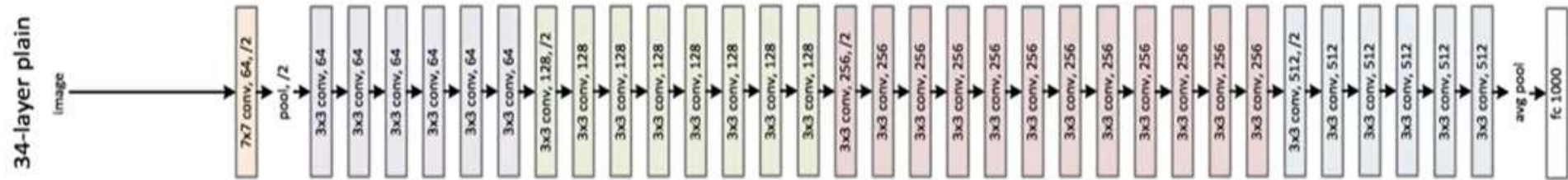


$$x \rightarrow \cdots \rightarrow a^{[l]}$$

*He, at all 2015, Deep residual networks for image recognition*

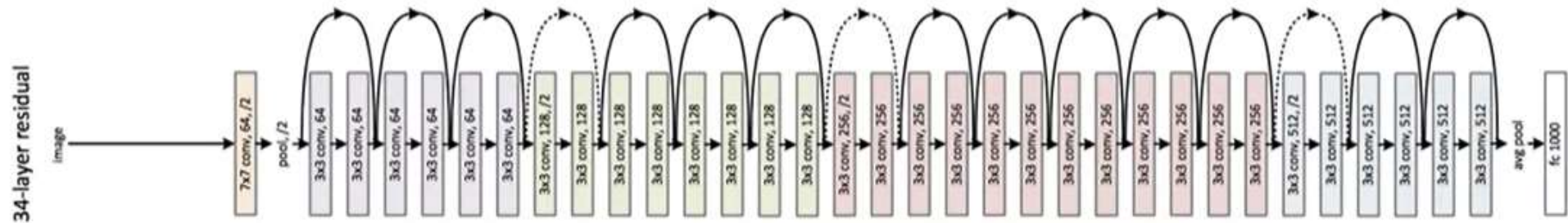Skip connection builds networks that enable to train very deep structures (over 100 layers).
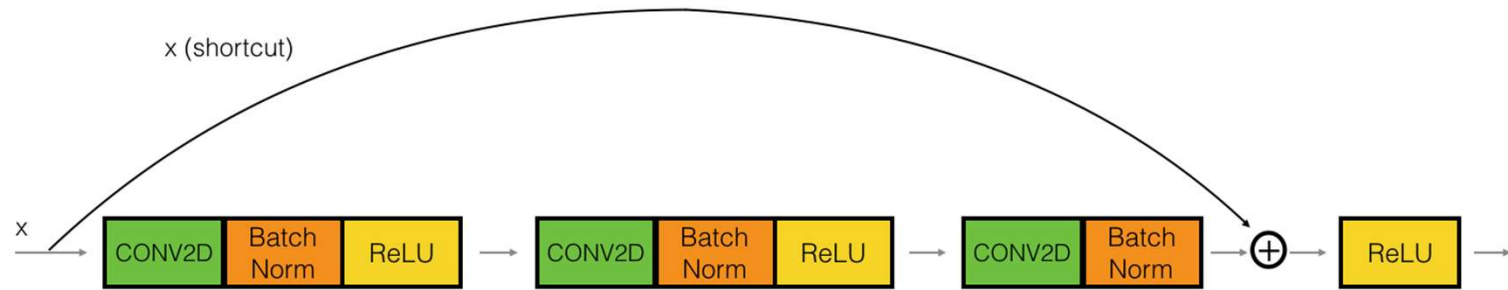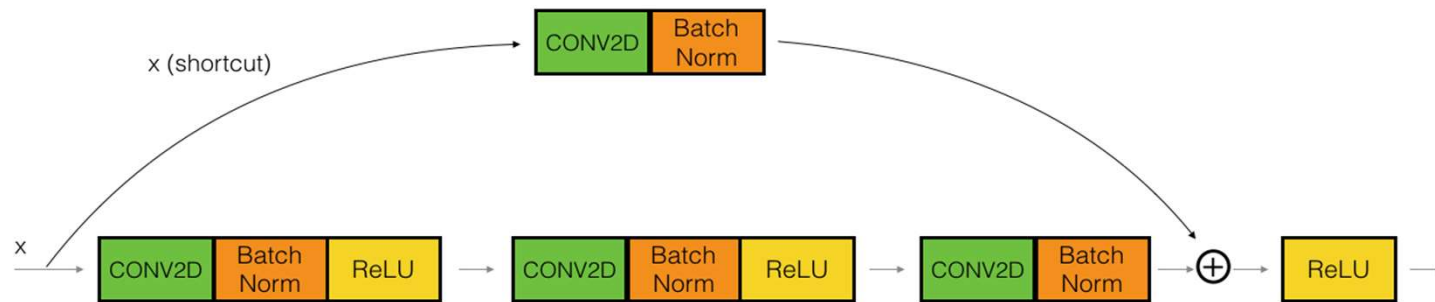
# ResNet on Images

# Lab: Residual blocks
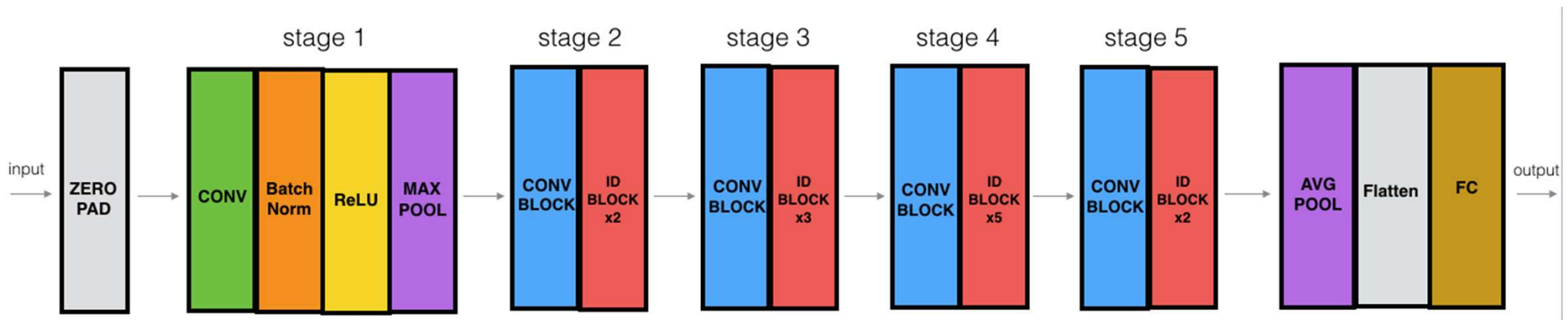
**Identity skip** connection /shortcut



**Convolution  skip** connection /shortcut



**ResNets** are built with such residual blocks.

# Lab: ResNet 50 model



ResNet-50 has 5 stages with Convolution and Identity blocks.

Each Conv BLOCK has 3 convolution layers and each identity block (ID BLOCK) also has 3 convolution layers.

Stage1(1)+Stage2(3+2*3)+ Stage3(3+3*3)+ Stage4 (3+5*3)+ Stage5 (3+2*3)+1FC =50

ResNet-50 has over 23 million trainable parameters !!!

universidade
de aveiro

# 1 by 1 Convolution



If we have 2D image and convolve with 1x1 filter it does not seem very usefull, just a miltiplication by a number.

But if we have a 3D volume (e.g.  6 x 6 x 32 channels),  convolution with 1x1x32 is useful.

# 1 by 1 Convolution (Network in Network)



1 by 1 convolution is useful for:

a) Learn complex non-linear relationships over a slice of the input volume.

b) Way to control (reduce, increase) number of output channels.

*Paper: Lin at al., 2013, Network in Network*

# 1x1 Convolution –
# example of output channel reduction



**How to get the output volume ?**

$28 \times 28 \times 192$

$28 \times 28 \times 32$

universidade
de aveiro

# Inception Module

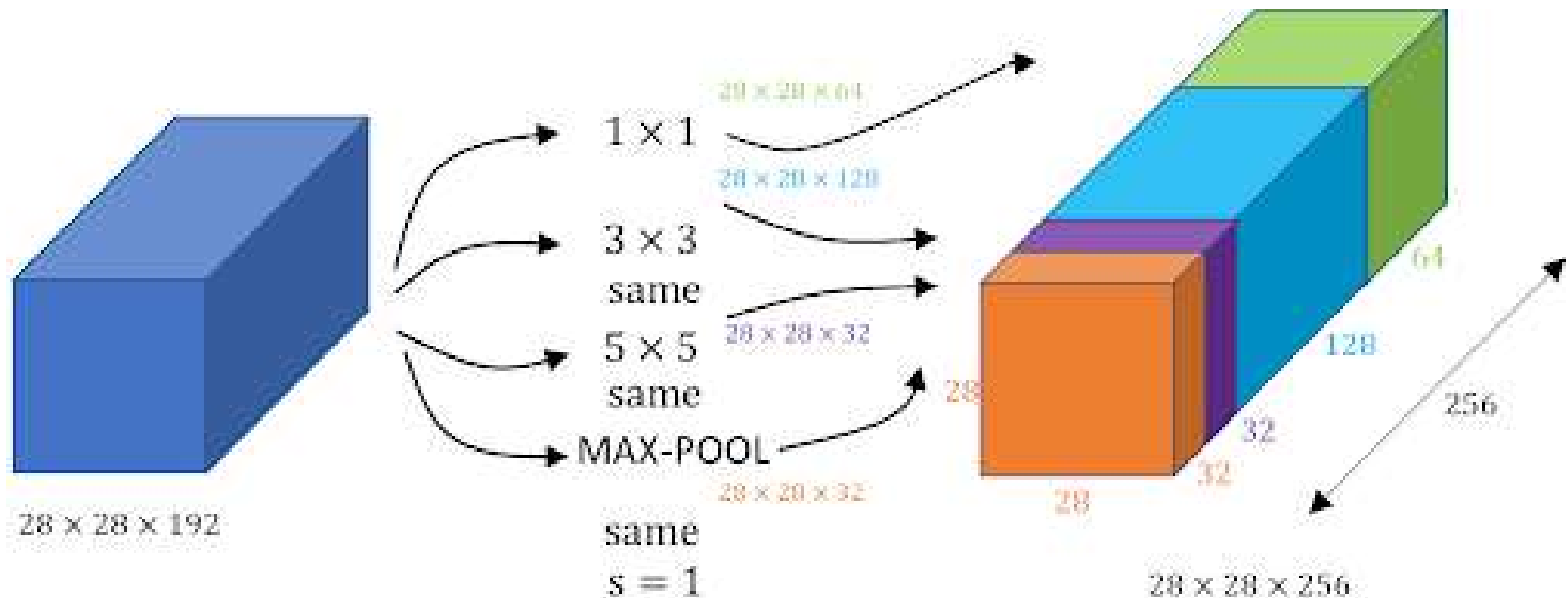When design a Conv Layer we may want to choose different filters (e.g. 1x1; 3x3; 5x5, or pooling payer).
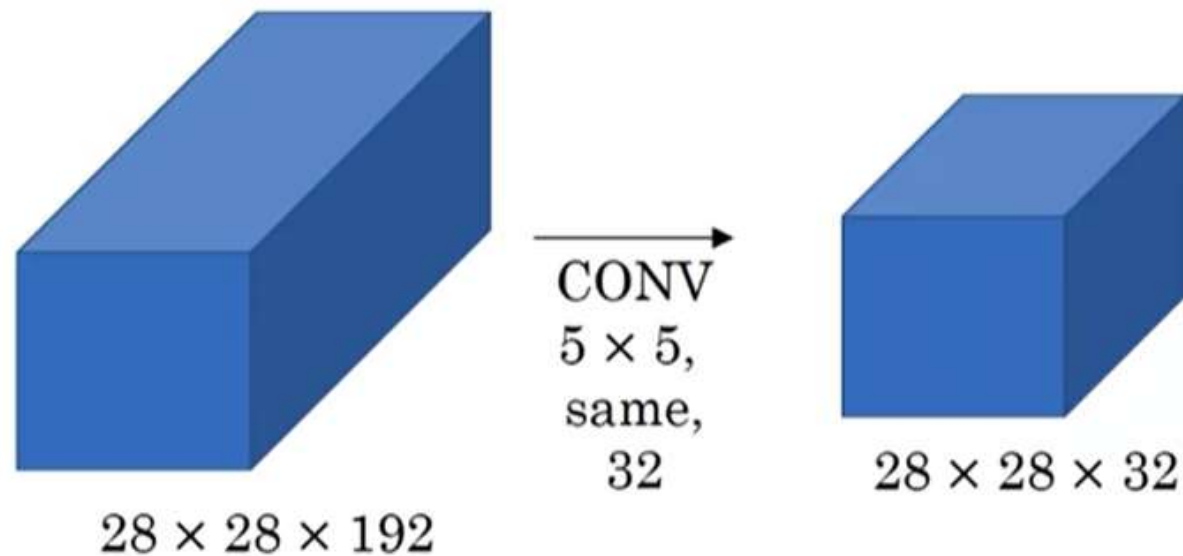In Inception network we may apply all of them and concatenate the outputs.

But such big/complex network will increase the computational cost.

*Paper: Szegedy at al., 2014, Going deeper with the convolutions*



28 × 28 × 192

1 × 1
28 × 28 × 64

3 × 3
same
28 × 28 × 128

5 × 5
same
28 × 28 × 32

MAX-POOL
28 × 28 × 32

same
s = 1

28 × 28 × 256

64
128
32
32
256
28
28

universidade
de aveiro

# Computational Cost - example



$28 \times 28 \times 192$ → CONV $5 \times 5$, same, 32 → $28 \times 28 \times 32$
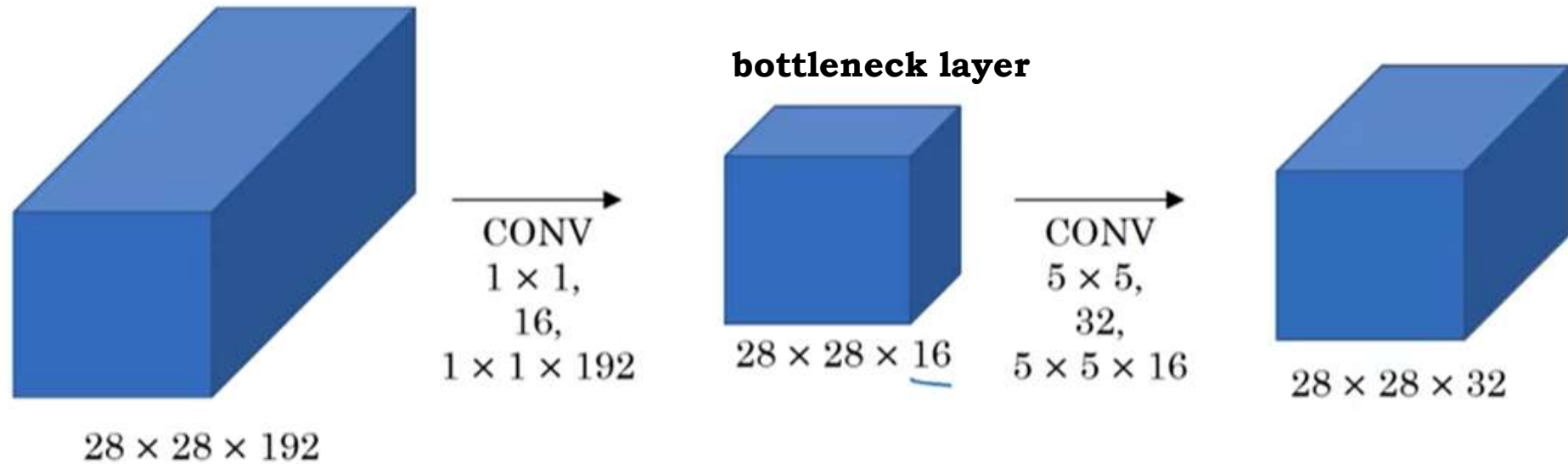
32 filters, each of them of size ?

padding ="same" (no height and width shrink)

**Computational Cost (how many computations) ?**

universidade de aveiro

# Comp. cost using 1x1 convolution -example

**bottleneck layer**

$28 \times 28 \times 192$
CONV $1 \times 1$, 16, $1 \times 1 \times 192$
$28 \times 28 \times 16$
CONV $5 \times 5$, 32, $5 \times 5 \times 16$
$28 \times 28 \times 32$

Input and output dimensions are still the same.

Shrink the big input volume into a much smaller intermediate volume (28x28x16) , known as **bottleneck layer** (the smallest part of the network).

**Computational Cost** (using 1x1 convolutions):

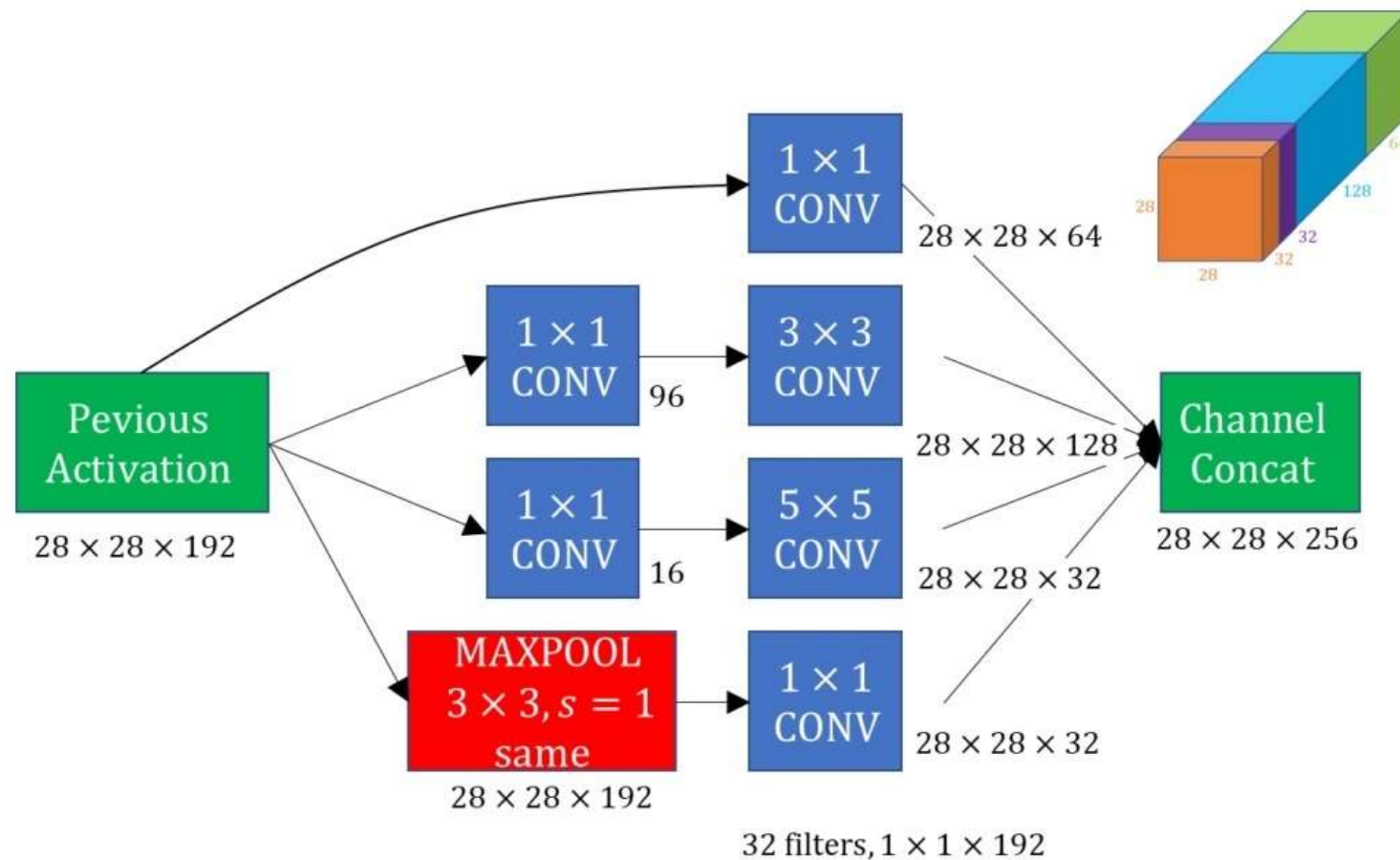**Stage1 (Input-Bottleneck connection):** ? computations
+
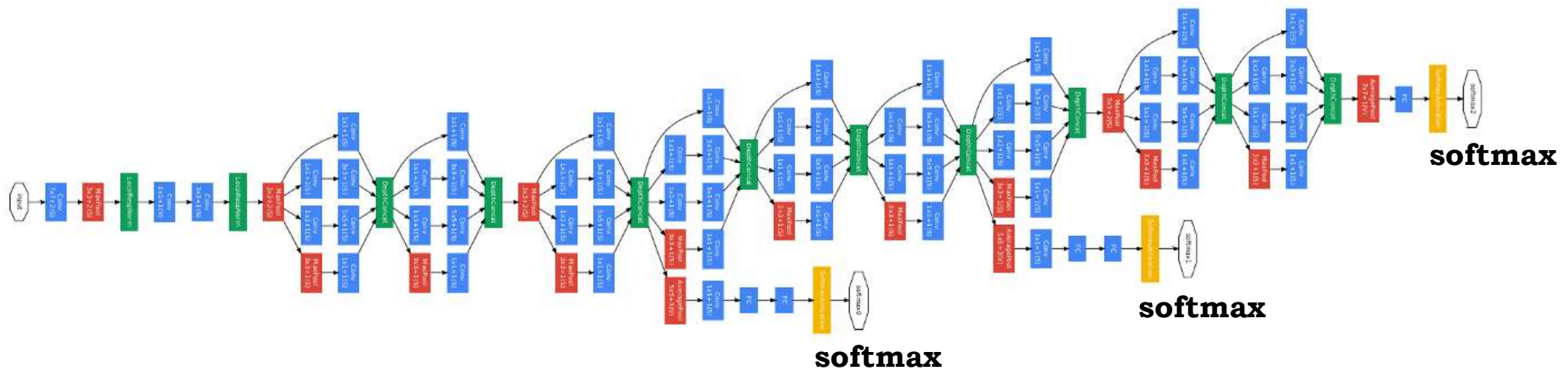**Stage2 (Bottleneck –Output connection):** ? computations

universidade
de aveiro

# INCEPTION MODULE

Inception Network <mark>puts all these modules together</mark>, using <mark>1x1 convolutions to save computations</mark>.



An example of an Inception module

# INCEPTION NETWORK (GoogLeNet)
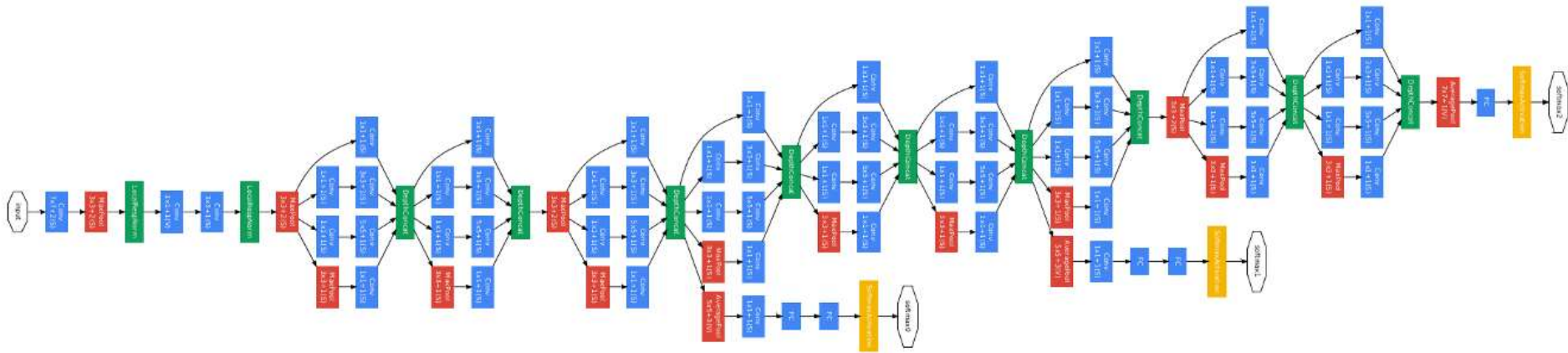


**softmax**

**softmax**

**softmax**

Ref. Szegedy et al 2014, "Going deeper with Convolutions".  Around 25 mln. of parameters.

Many inception blocks repeated at different positions of the network.

The last layers are fully connected followed by softmax layer trying to predict what is the output label.

**Side branches** take inf from previous (early) layers and try to make predictions.
The earlier features are not that bad to make predictions.
This has a regularization effect, combat the over fitting.

# INCEPTION NETWORK – several architectures



**INCEPTION NETWORK = GoogLeNet =Inception_v1**
Other versions are available, Inception v2, v3, v4, etc., they all are based on the idea of inception module.

universidade
de aveiro