# Departamento de Eletrónica, Telecomunicações e Informática

# Complements of Machine Learning

## Lecture 5 : Deep NN – Mobile Net & Object detection

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

# Outline

## *Part 1*

1. **MobileNet - MobileNet v1 ; MobileNet v2**

2. **EfficientNet**

## *Part 2*

1. **Object detection**

2. **Sliding windows detection algorithm**

universidade
de aveiro

# *Part 1: MobileNet & EfficientNet*
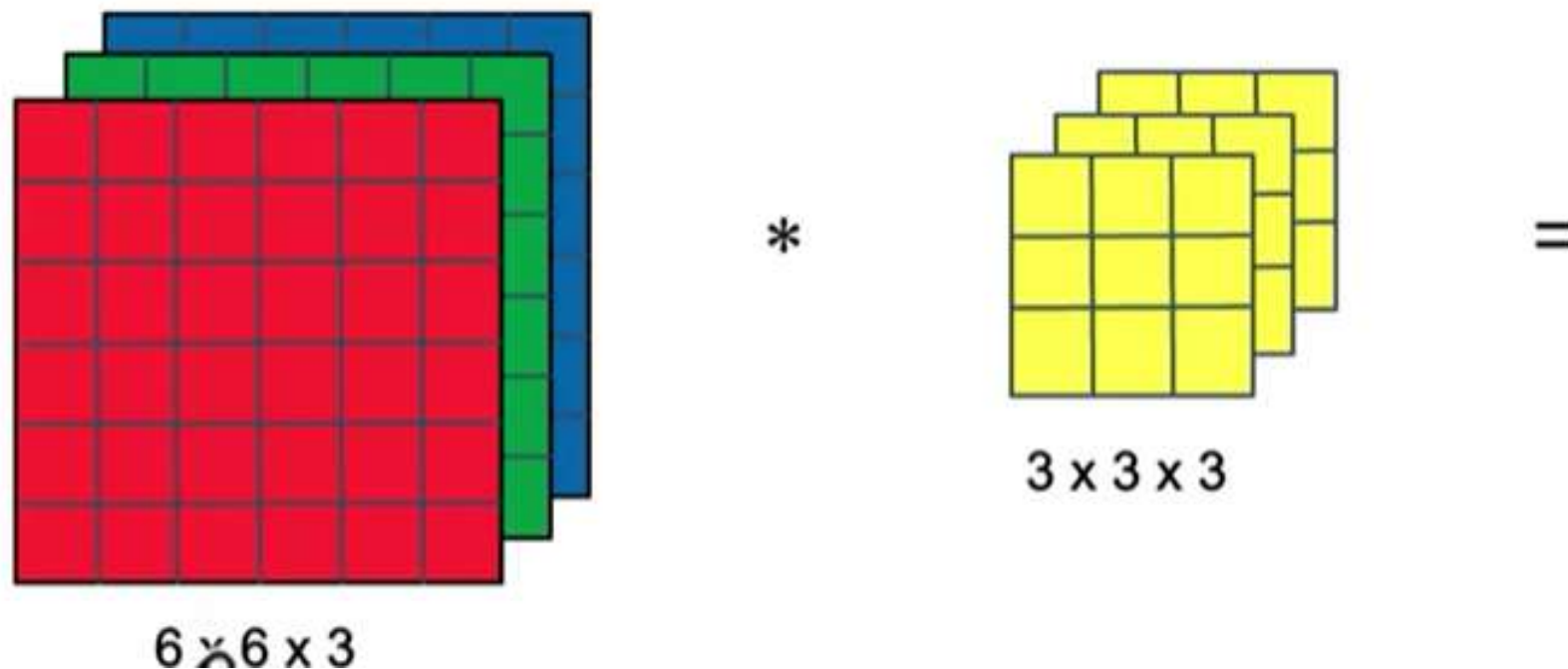
# Motivation for MobileNets



Low computational cost and low memory at deployment

Useful for mobile and embedded vision applications

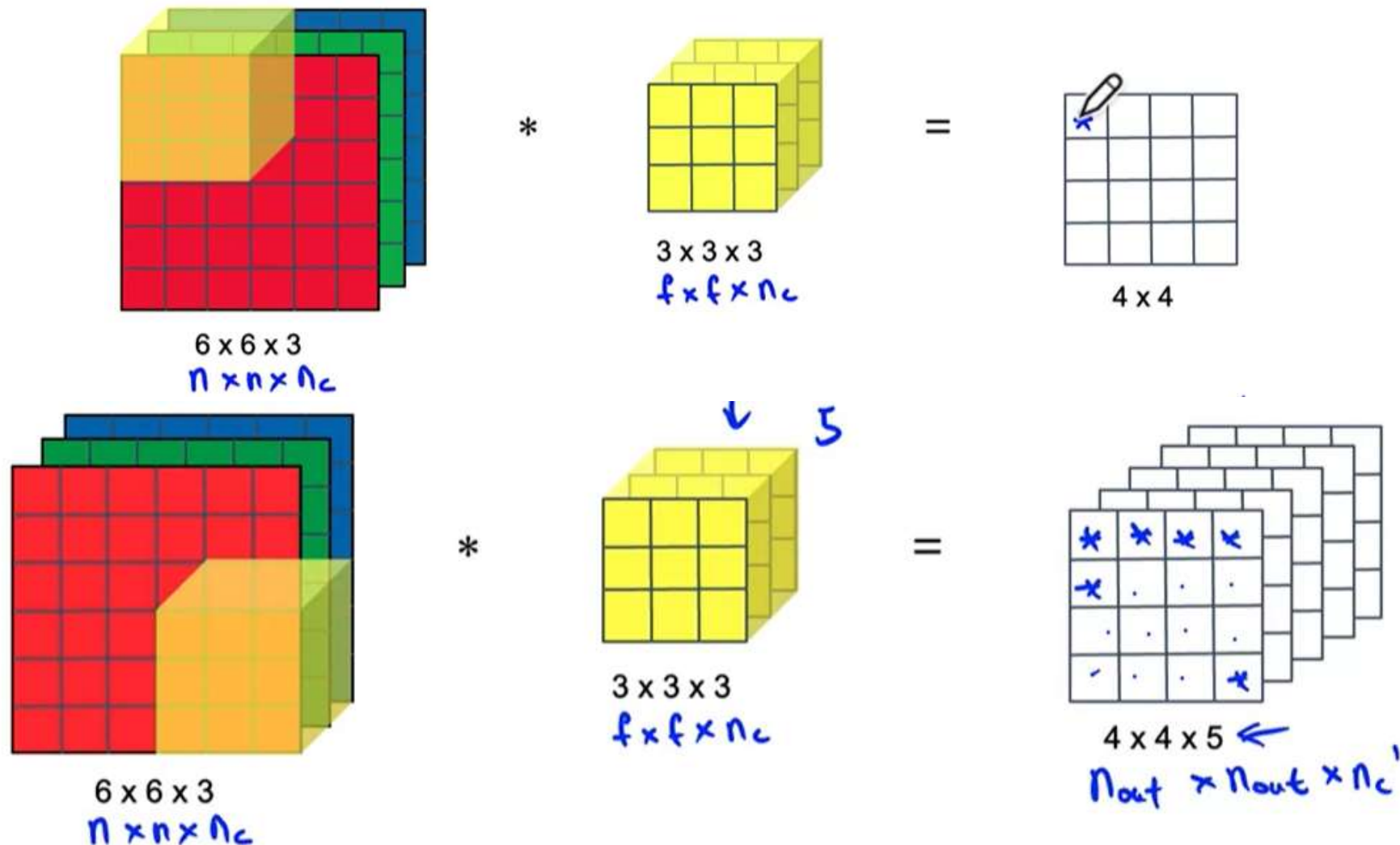Key idea: normal vs depth-wise separable convolutions

# Normal Convolution



6 x 6 x 3    *    3 x 3 x 3    =

**Stride =1, no padding**

**Computational cost  for applying 5 filters = ?**

# Rule of Normal Convolution



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

$*$

$3 \times 3 \times 3$
$f \times f \times n_c$

$=$

$4 \times 4$

$\downarrow \quad 5$

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

$*$

$3 \times 3 \times 3$
$f \times f \times n_c$
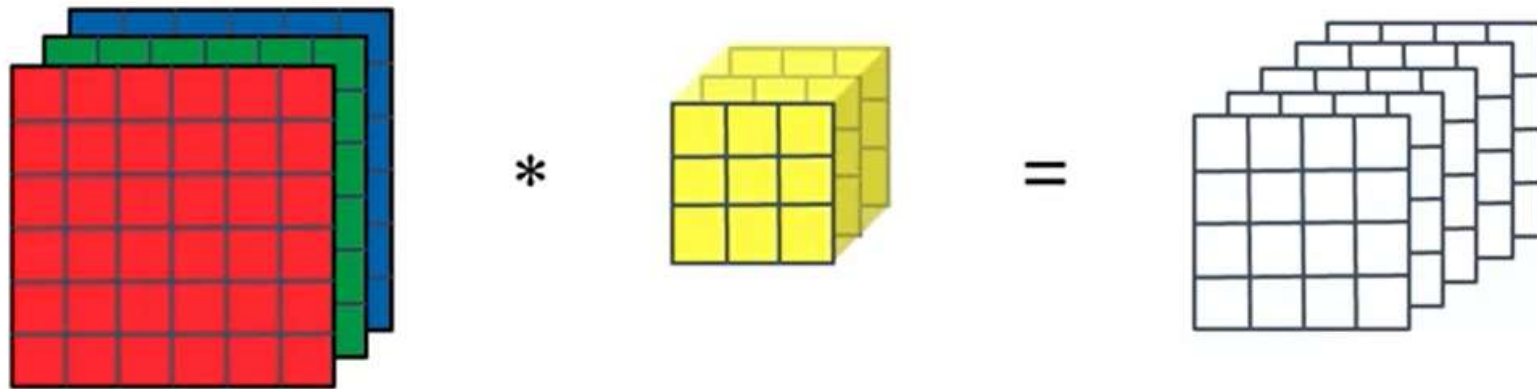
$=$

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

**3rd filter dimension = 3rd input dimension** (e.g. $n_c = 3$)

**Computational cost = #_filter parameters * #_filters positions * #_filters**

**2160        =        (3x3x3)        *        (4x4)        *        5**

universidade de aveiro

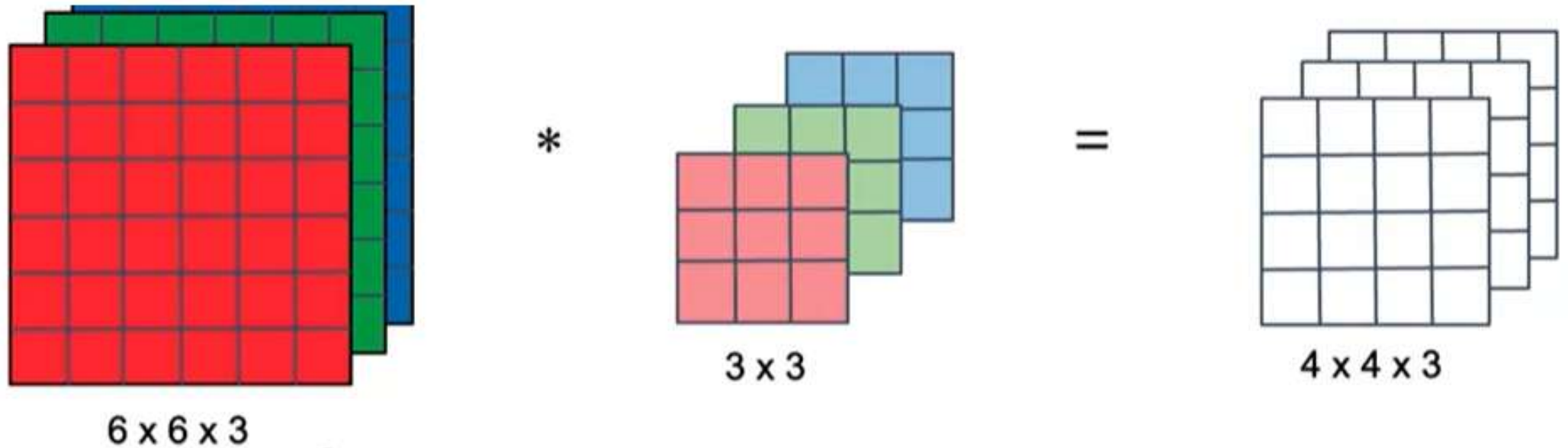# Normal vs Depthwise Separable Convolution



**Depthwise separable convolution has two steps:**

1) Depthwise convolutions
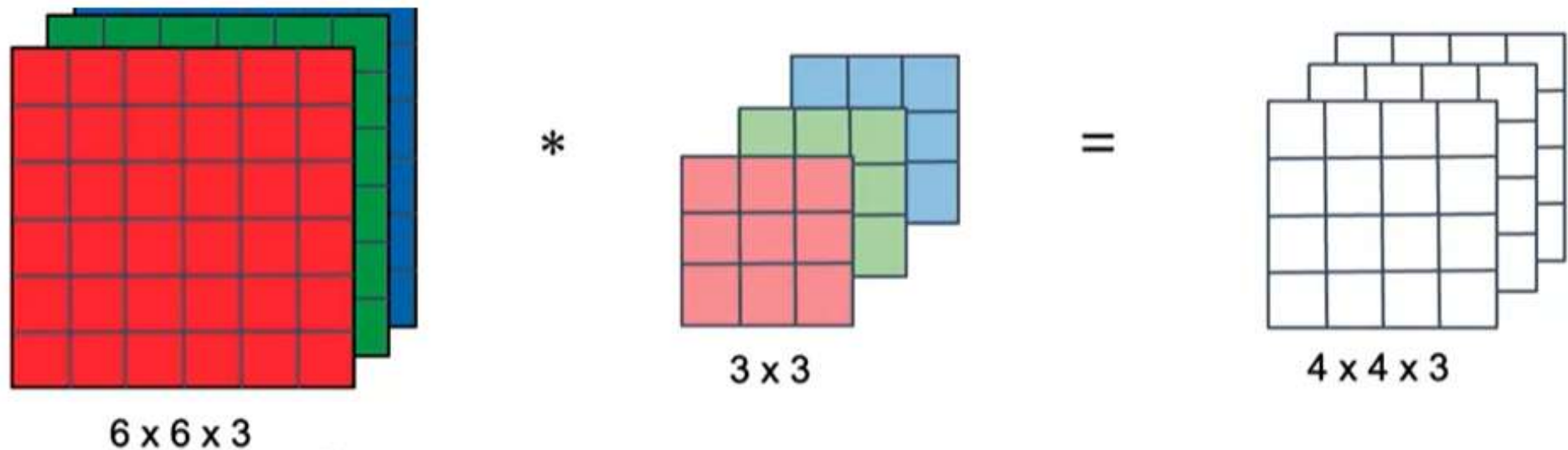2) Pointwise convolutions

# Step 1: Depthwise Convolution



6 x 6 x 3 * 3 x 3 = 4 x 4 x 3

**Depthwise convolution rule:**

 - Each filter has 2 dimensions (height and width)
 -  # of filters= 3$^{rd}$ dimension of the input volume (input channels)

**Computational cost (number of computations)  = ?**

# Step 1: Depthwise Convolution



$6 \times 6 \times 3$     *     $3 \times 3$     =     $4 \times 4 \times 3$

**Depthwise convolution rule:**

 - Each filter has 2 dimensions (height and width)
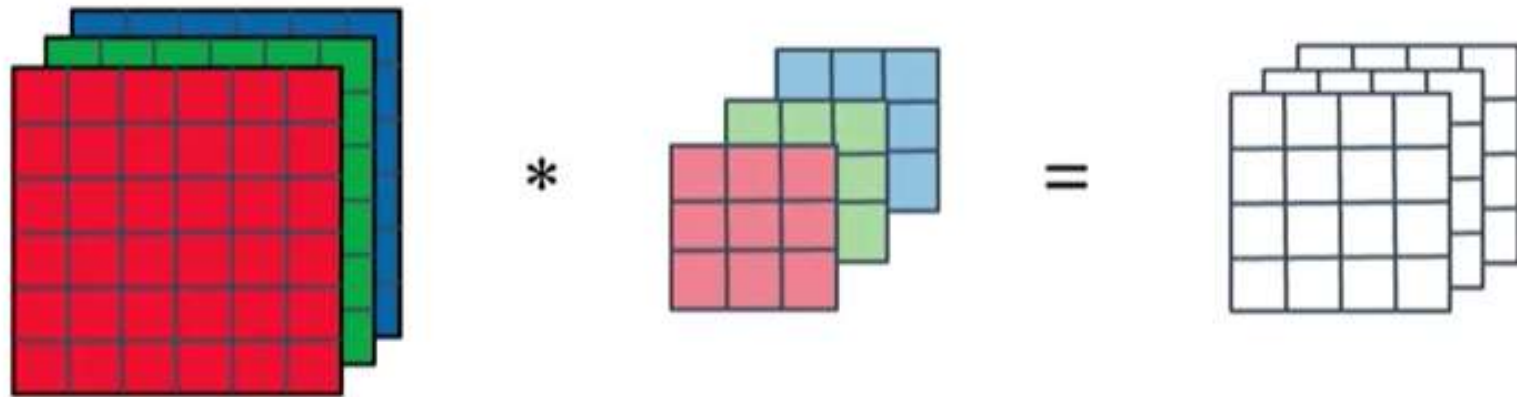-  # of filters= 3$^{rd}$ dimension of the input (input channels)

**Computational cost 1m =**
**#_filter parameters * #_filters positions * #_filters**

**Step 1:   432  =   (3x3)          *     (4x4)              *     3**

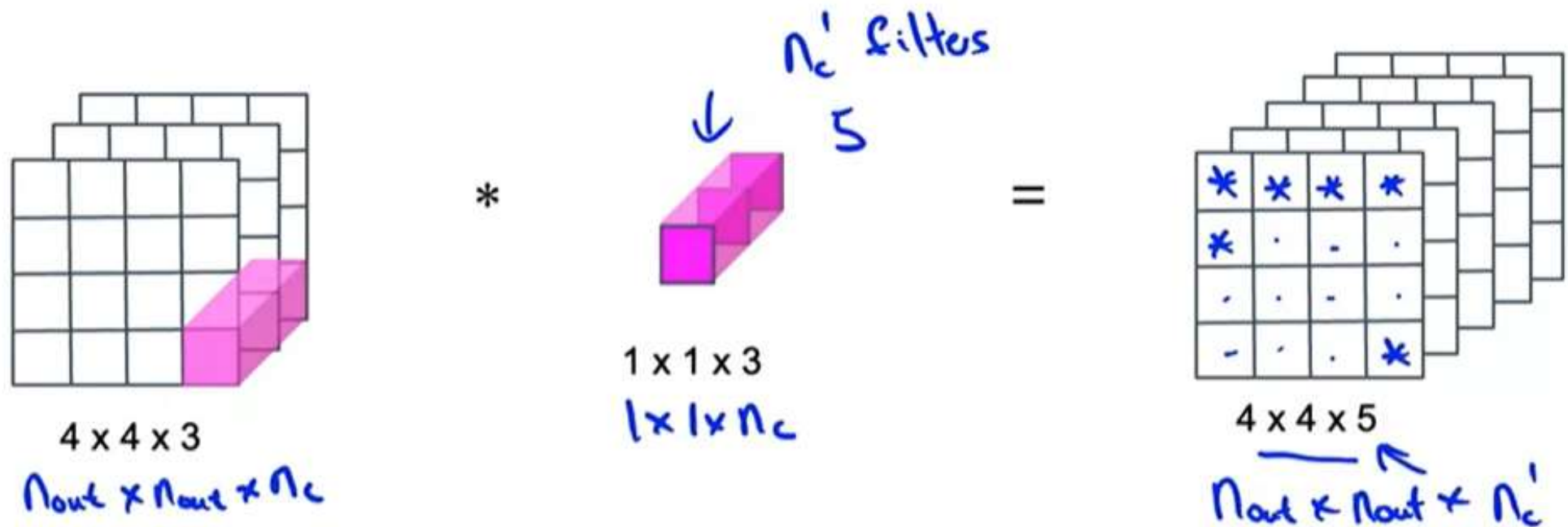universidade
de aveiro

# Step 2: Pointwise Convolution



Depthwise Convolution

Pointwise Convolution

# Step 2: Pointwise Convolution



$n_c'$ filters

5

$1 \times 1 \times 3$

$1 \times 1 \times n_c$

$4 \times 4 \times 3$

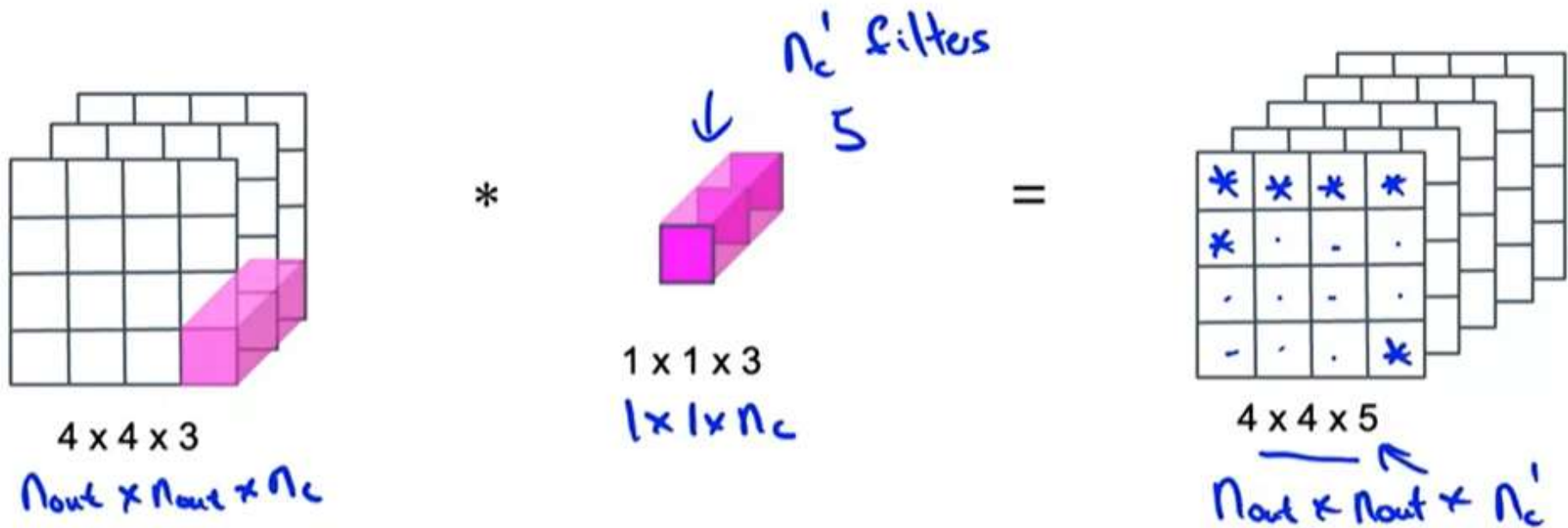$n_{out} \times n_{out} \times n_c$

$4 \times 4 \times 5$

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

**Computational cost 2 =**
   **#_filter parameters * #_filters positions * #_filters**

**Computational cost 2 = ?**

universidade
de aveiro

# Step 2: Pointwise Convolution



$n_c'$ filters
5

$1 \times 1 \times 3$
$1 \times 1 \times n_c$

$4 \times 4 \times 3$
$n_{out} \times n_{out} \times n_c$

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

**Computational cost 2 =**
        **#_filter parameters * #_filters positions * #_filters**

**Step 2:      240 =      (1x1x3)        *        (4x4)            *      5**

universidade
de aveiro

# Cost Summary

For the particular example:

Normal convolution: 2160

Depthwise separable convolution: 432+240=672

Ratio =672/2160=0.31  (31% as compared to Normal convolution)

In the paper *Andrew G. Howard at al., 2017, MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications.*

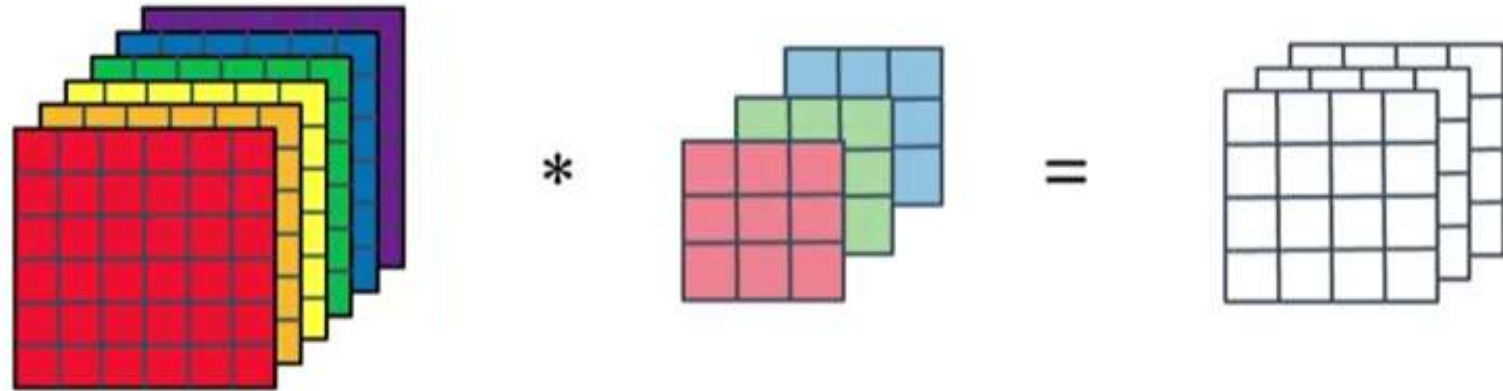General Ratio formula =  1/output channels +1/(filter_dimension^2 )

In our case: 1/5+1/(3^2)

Typical ratio: 1/512 + 1/(3^2) => 0.11

Depthwise Separable Convolution about 10 times cheaper (less computations) than Normal convolution
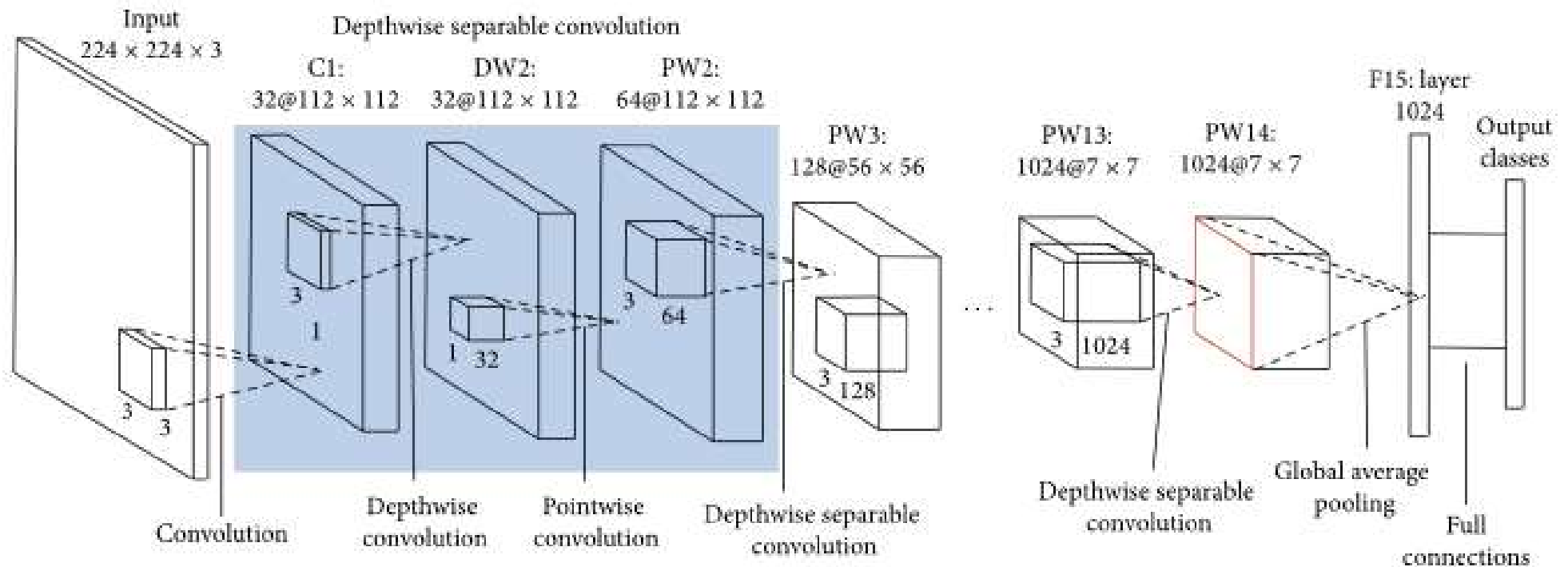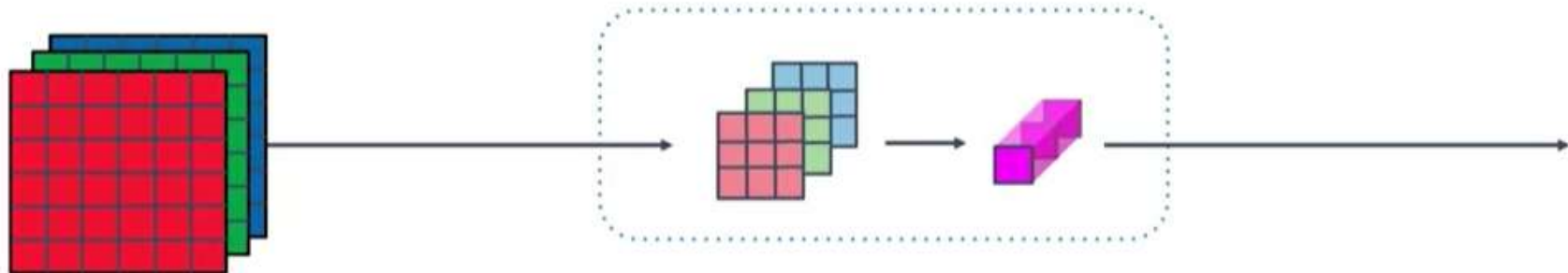
universidade
de aveiro

# Depthwise Separable Convolution



For more than 3 input channels (*n_input_channels*) we still use the same depthwise convolution icon, but will apply *n_input_channels* filters.
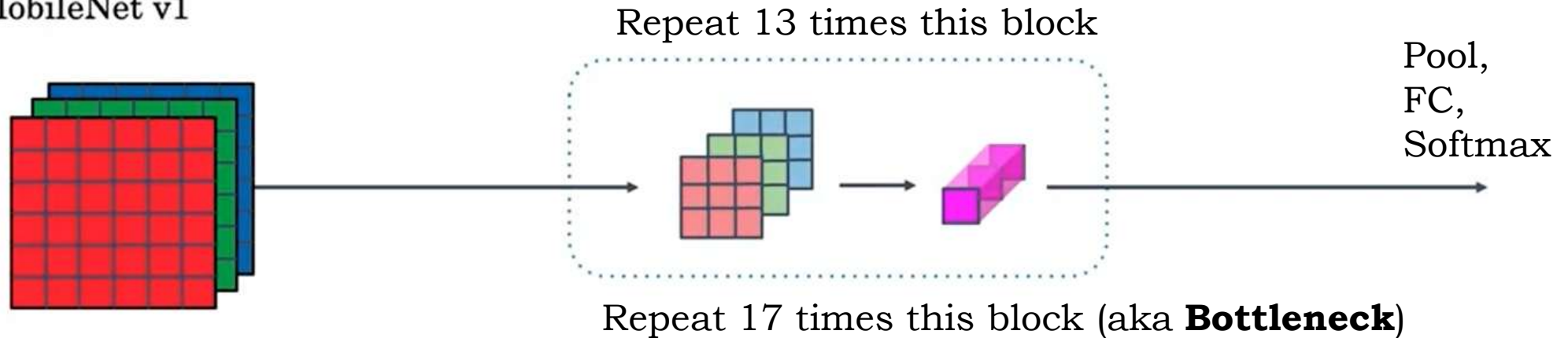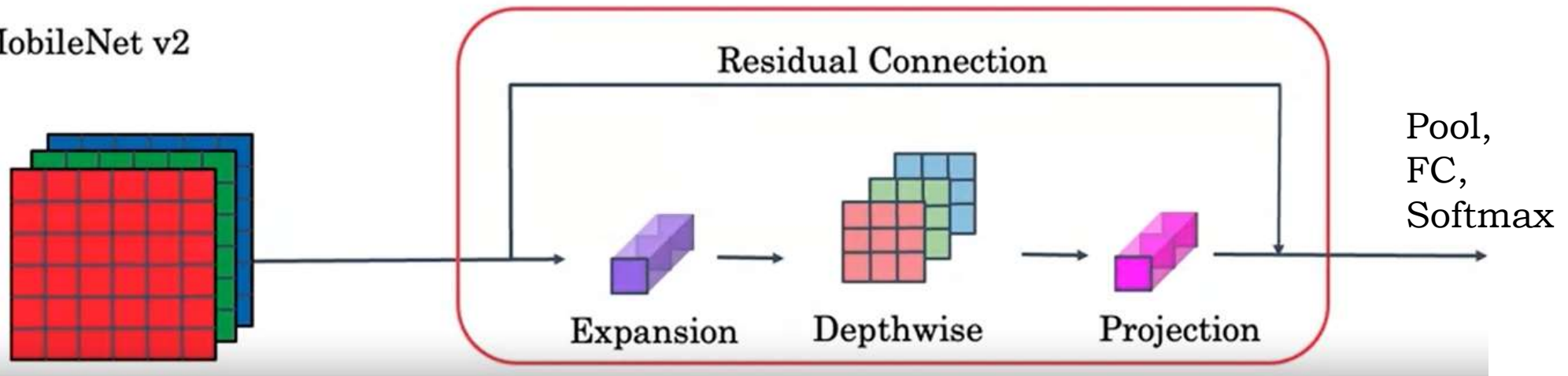
# MobileNet v1



Repeat 13 times this block

Pool,
FC,
Softmax

Input 224 × 224 × 3

Depthwise separable convolution

C1: 32@112 × 112

DW2: 32@112 × 112

PW2: 64@112 × 112

PW3: 128@56 × 56

PW13: 1024@7 × 7

PW14: 1024@7 × 7

F15: layer 1024

Output classes

Convolution

Depthwise convolution

Pointwise convolution

Depthwise separable convolution

Depthwise separable convolution

Global average pooling

Full connections

*Andrew G. Howard at al., 2017, MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications.*

# MobileNet v1 and v2 architectures

MobileNet v1

Repeat 13 times this block

Pool, FC, Softmax

MobileNet v2

Repeat 17 times this block (aka **Bottleneck**)

Residual Connection

Expansion   Depthwise   Projection

Pool, FC, Softmax

*Mark Sandler, Andrew Howard, et al. 2019, MobileNetV2: Inverted Residuals and Linear Bottlenecks.*

universidade de aveiro

16

# MobileNet v2 Bottleneck



**Expansion:** The input of the bottleneck block is a small volume (nxnx3) , apply 18 filters (i.e. typical factor of expansion 6) , we get big (expanded) volume (nxnx18).

**Depthwise:** Keep the dimension of the volumes by adding padding.

**Pointwise:** convolve 3 (1x1x18) dimensional filters, end up with nxnx3 output volume. In this last step, we project from big volume (nxnx18) down to smaller volume (nxnx3).

# MobileNet v2 – advantages

The bottleneck block with the expansion learns richer and more complex functions, while the size of the activations to pass from layer to layer is kept small => relatively low memory required.

MobileNet v2 can get  a better performance than MobileNet v1, and still use a modest amount of computing and memory resources.

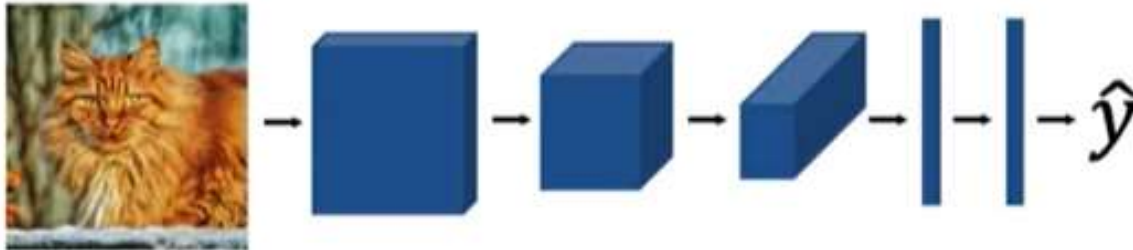*Mark Sandler, Andrew Howard, et al. 2019, MobileNetV2: Inverted Residuals and Linear Bottlenecks.*

*Tan & Le, 2020, EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks*

universidade
de aveiro
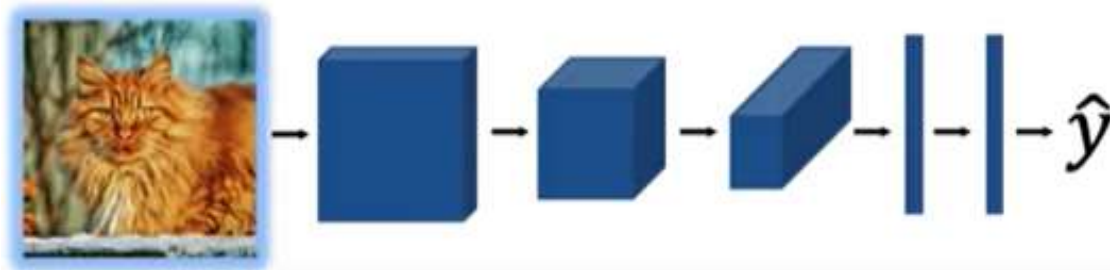
# Efficient Net

Major parameters to vary (up or down) to satisfy computacional limits:

**1) Image resolution (r)**
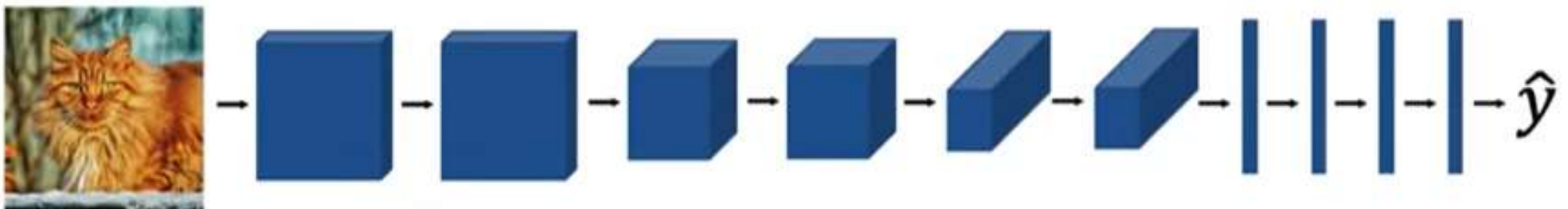
**2) NN Depth (d)**

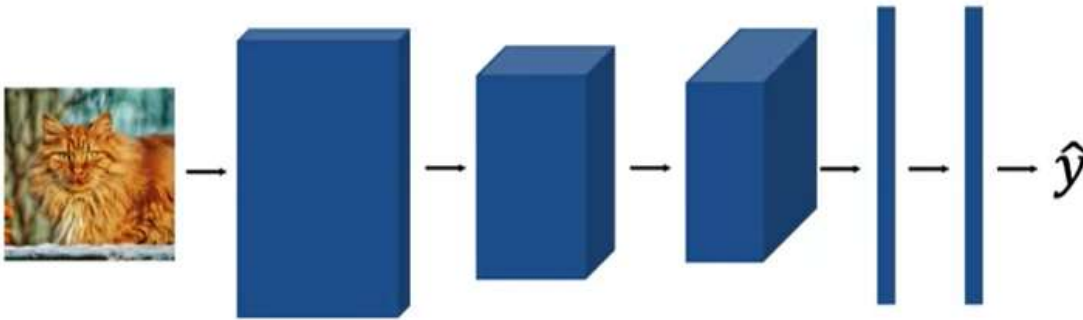Baseline



Higher Resolution



Deeper

# Efficient Net

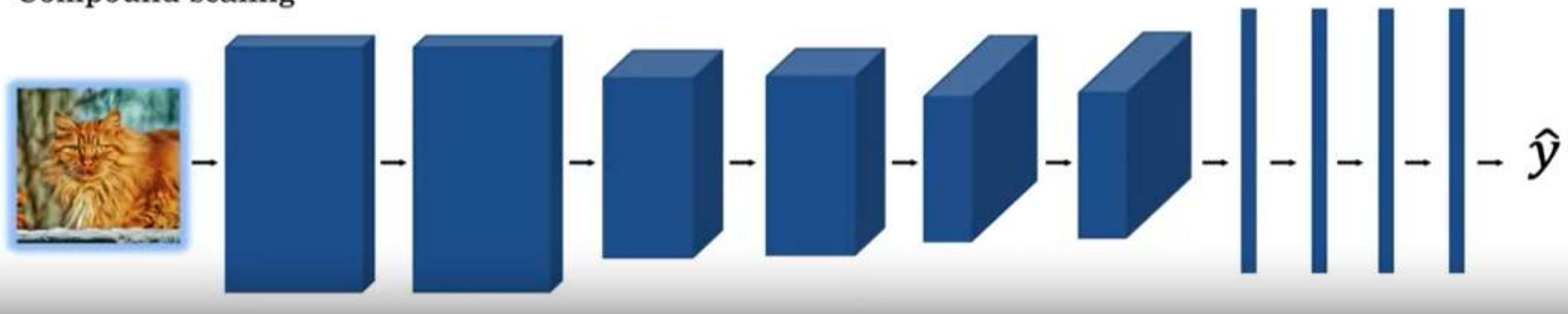Major parameters to vary (up or down) to satisfy computacional limits:

3) **Width of the layers (w)**

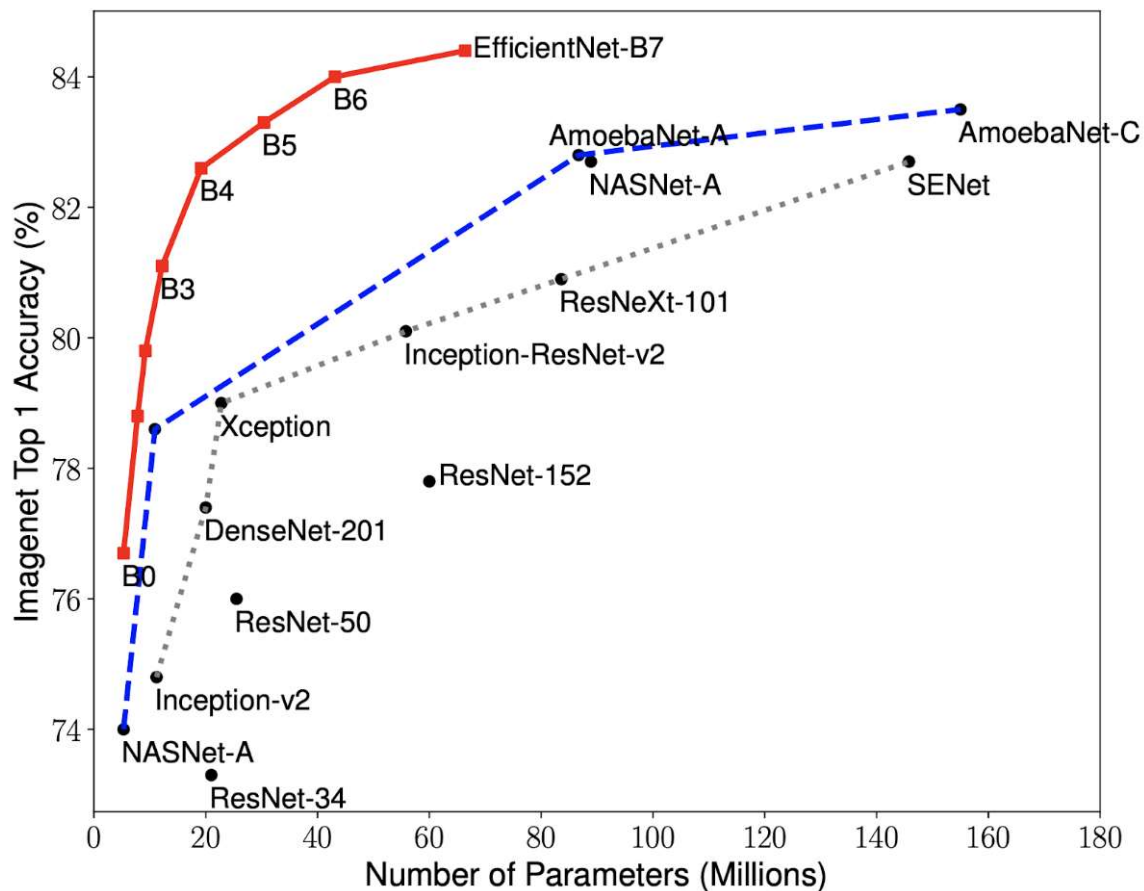4) **Compound scaling: Simultaneously scale up or down image resolution (r), depth (d), width (w).**

Wider



Compound scaling

universidade
de aveiro

# Efficient Net
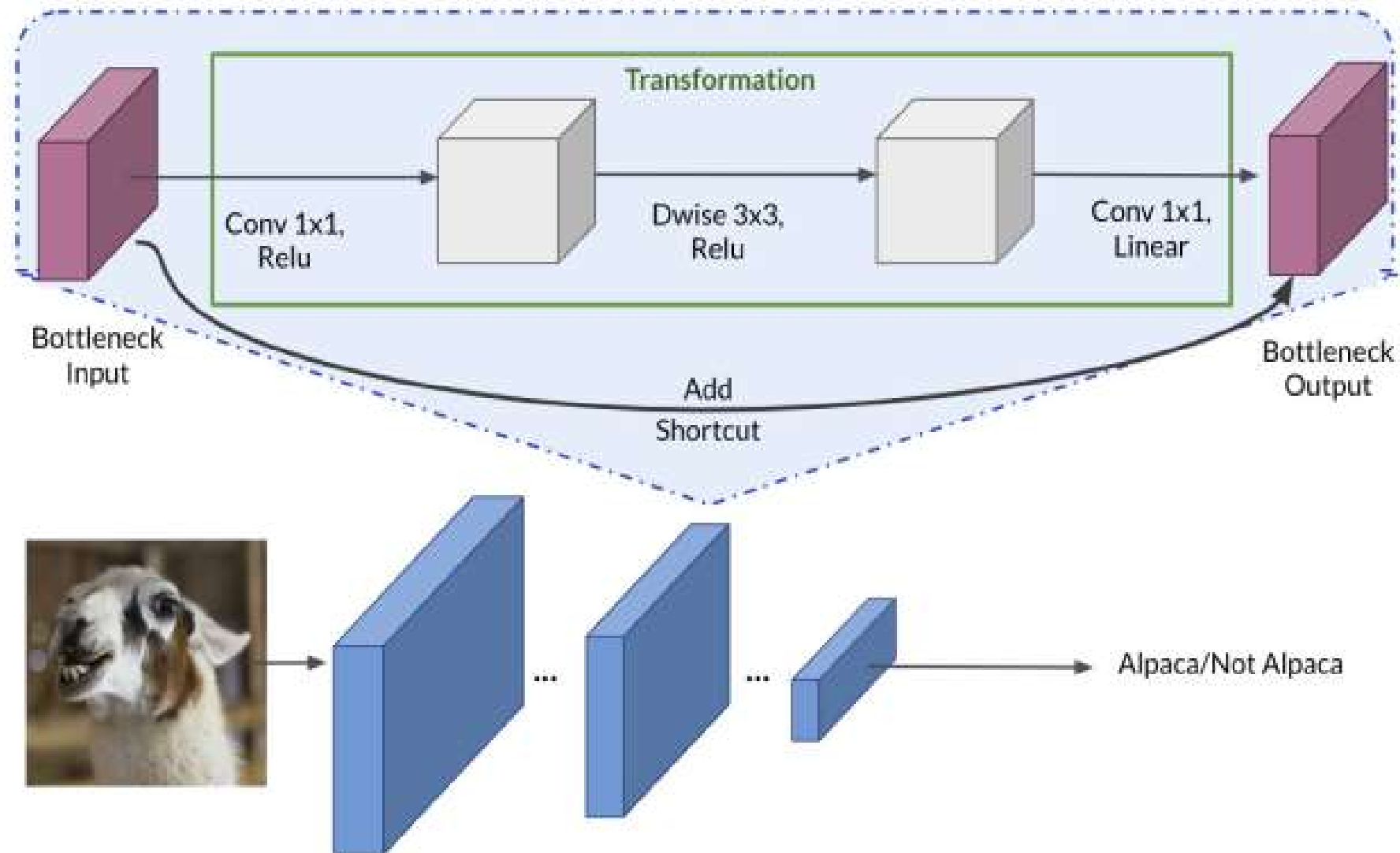
Suggestion: look at EfficientNet implementations, to choose a trade-off between r, d, w.

Unlike conventional practice that arbitrary scales these factors, the EfficientNet uniformly scales network width, depth, and resolution with a set of fixed scaling coeff.



EfficientNet: Improving Accuracy and Efficiency through AutoML and Model Scaling
*Tan, and Le, Google AI, 2019*

# Lab: MobileNetV2 for Transfer Learning



Use a pre-trained MobileNetV2 to build a binary classifier (Image of Alpaca (Lama) /Not Alpaca animal).
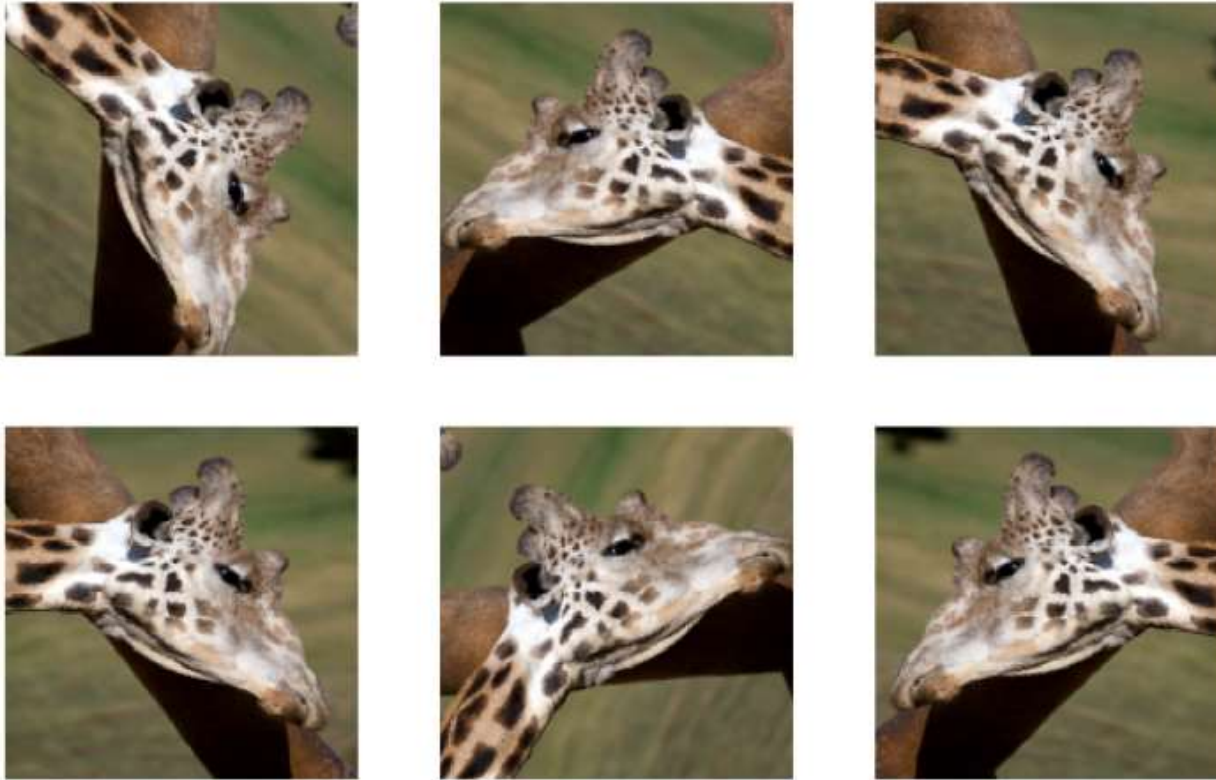
# Lab Data



**262 training images; 65 validation images**

# Lab: On-line Data Augmentation



**Data Prefetch :** CPU loads stream of images (e.g. 32) coming from the hard disc and generate distortions to form mini-batches that are passed to the training algorithm.

The two processes (data augmentation and training run in parallel.

```
Create a Sequential model composed of 2 layers
'''
data_augmentation = tf.keras.Sequential()
data_augmentation.add(RandomFlip("horizontal"))
data_augmentation.add(RandomRotation(0.2))
```

universidade
de aveiro

# Part 2: Object Detection

# Image classification/localization/detection

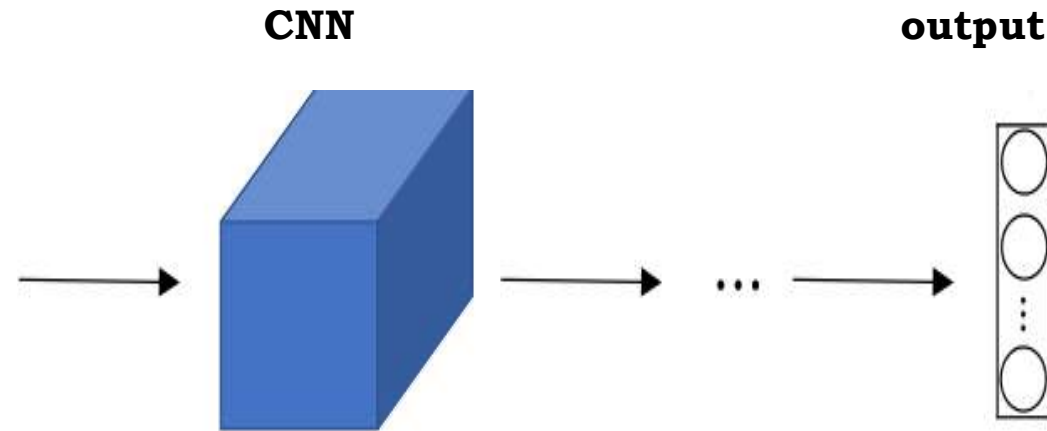| Image classification | Classification & Localization | Detection |
|---|---|---|
|  |  $b_x, b_y, b_h, b_w$ |  |
| **car** | **car + bounding box** | **many objects+ bounding boxes** |

**Image classification:** input a picture to CNN and the output is a class label (e.g. person, bike, car, background, etc.)

**Classification with localization**: the algorithm gives not only the class label of the object but also draws a bounding box (the coordinates) of its position in the image. Standard notation: (0,0) as the upper-left corner and (1,1) to be the lower-right corner. ($b_x$, $b_y$, $b_h$, $b_w$) describes the bounding box.

# Object classification with localization



**CNN**

**output**

Classes (e.g.):
1. person
2. Bikes
3. Car
4. Background (no object)

Output: ($p_c$ , $b_x$ , $b_y$ , $b_h$ , $b_w$ , c= [$c_1$, $c_2$....$c_{end}$ ])
$p_c$ – is there an object or not (1/0)

<= Image label: [1 , $b_x$ , $b_y$ , $b_h$ , $b_w$ , 0, 0, 1]

$b_x, b_y, b_h, b_w$

<= Image label: [0 , ? , ? , ? , ? . ?, ?; ?]
? – "don't care"

universidade
de aveiro

# Output/label vector

$$y = (p_c, b_x, b_y, b_h, b_w, c)$$

$p_c = 1$ : confidence of an object being present in the bounding box

$p_c$ $b_x$ $b_y$ $b_h$ $b_w$ ← 80 class probabilities →

universidade
de aveiro

# Landmark Detection

The idea of bounding boxes $(b_x, b_y, b_h, b_w)$ inspired the use of DL for e.g. emotion recognition from faces, person's pose detection.
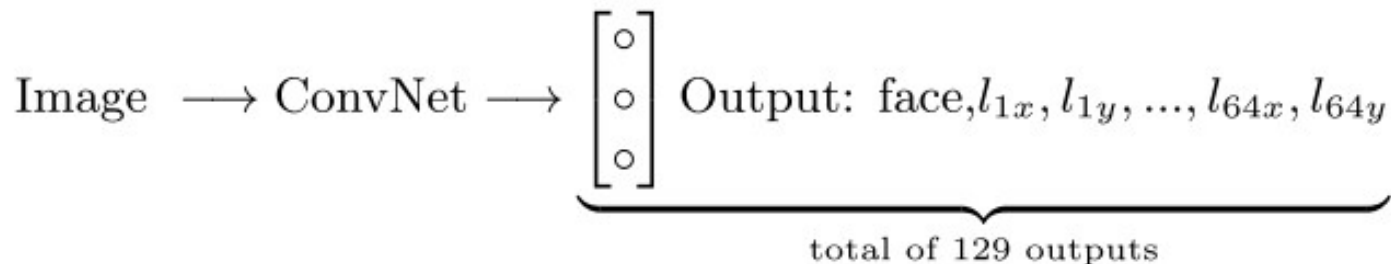
CNN can output the <u>important coordinates (called landmarks)</u> of targeted objects in the image. For example 64 chosen points on the face, or the body (key positions in the persons pose, e.g. the mid point of the chess, left/right shoulder, etc.).

Need to manually label all landmarks in the training data !!!

Consistent annotation over several images.
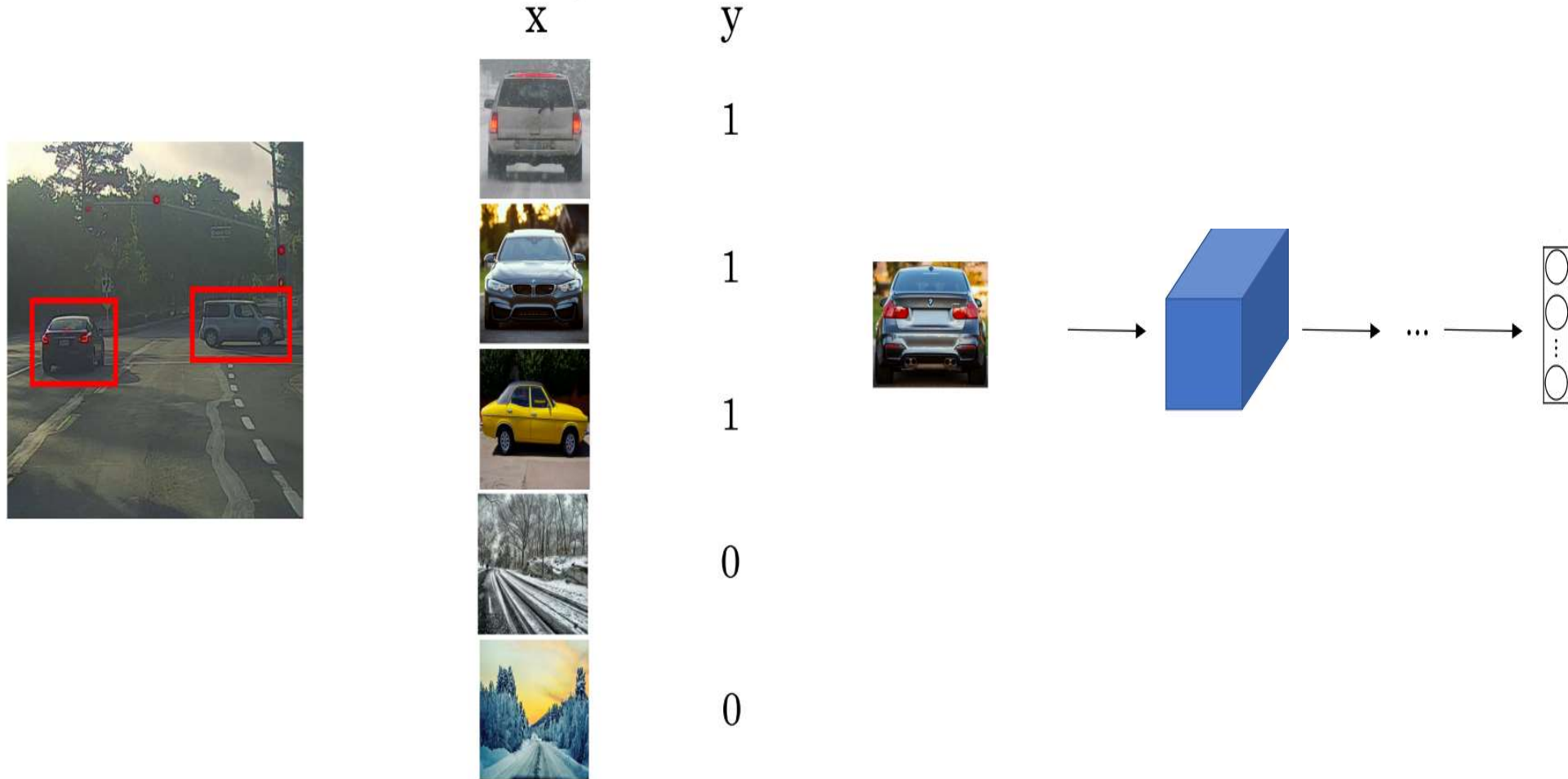


$b_x, b_y, b_h, b_w$

$$\text{Image} \longrightarrow \text{ConvNet} \longrightarrow \begin{bmatrix} \circ \\ \circ \\ \circ \\ \circ \end{bmatrix} \quad \text{Output: face}, l_{1x}, l_{1y}, ..., l_{64x}, l_{64y}$$

total of 129 outputs

# Object Detection algorithm

Training set:

x       y



1

1

1

0

0



Let's say we want to build a car detection algorithm.
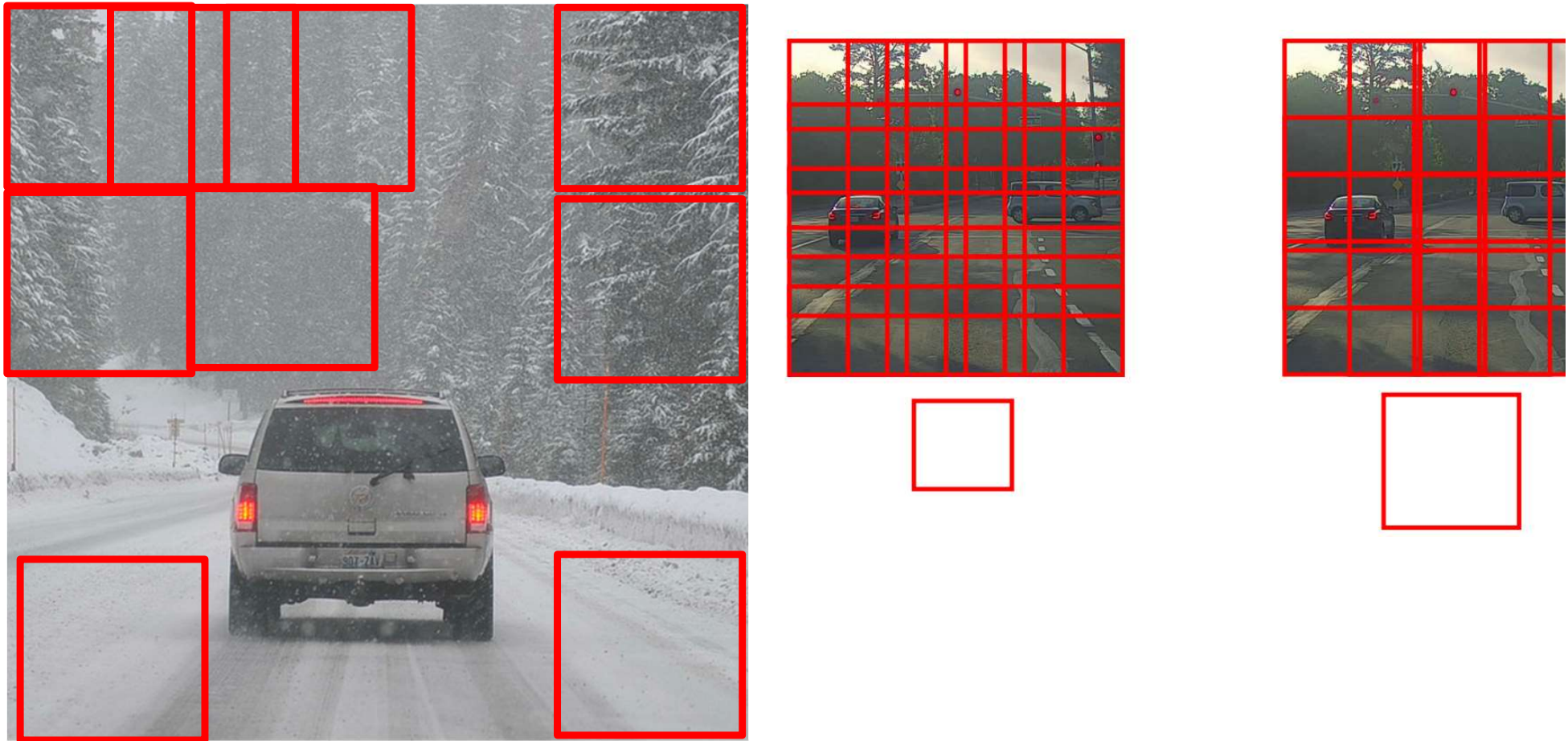First create a labelled training set.
Take a picture, cut out anything else that's not part of a car and get the car centred in pretty much the entire image (cropped examples of cars).
Train a CNN to output  y (0 or 1, is there a car or not).
The trained CNN is used in the Sliding Windows detection algorithm.
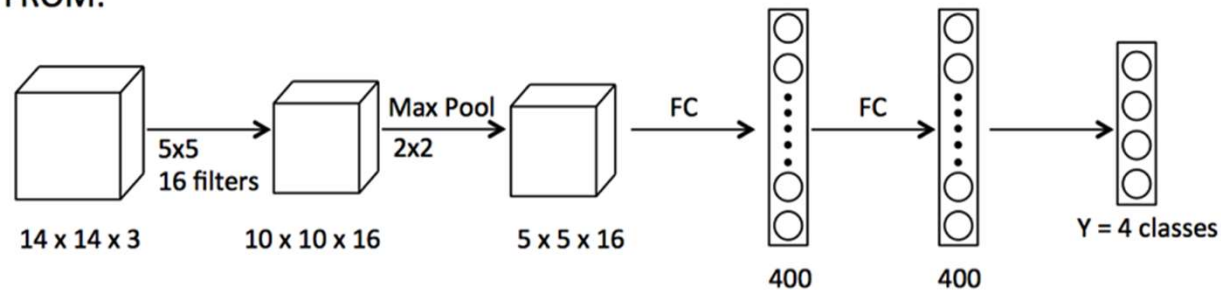
# Sliding windows detection algorithm



Pick a certain window size, input this sub-picture into <u>already trained CNN to detect objects</u>. Then shift the detection window to the right with one step (stride) and feed the new sub-picture into ConvNet. Go through every region (the stride needs to be small for the detection algorithm to run smoothly). Repeat the same with different sizes of the detection window in order to detect objects with different sizes of the picture.
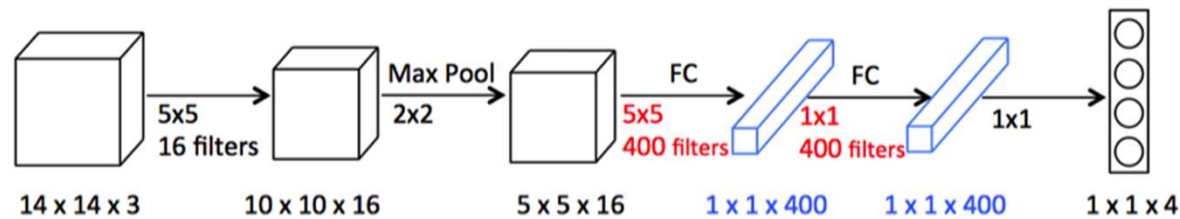
The Sliding Windows object detection algorithm has **infeasibly <u>high computationally cost.</u>**

**<u>Solution: Convolutional Implementation of sliding windows</u>**

# Turning FC layers into convolutional layers



The conv layers are the same in both implementations.
The difference is that the Fully Connected (FC) layers are implemented as convolutional layers.

**Ex.** *The first FC layer:* Let the input volume is 5x5x16, convolve it with 5x5x16 filter.
The outputs will be 1x1. Add 400 of these 5x5x16 filters, then the output dimension will be 1x1x400. Rather than seeing the FC as a set of 400 nodes, we view it as a 1x1x400 volume.
*The second FC layer:* implemented as convolution with 400 1x1 filters and will output 1x1x400 volume.
*The last layer:* is implemented as convolution with 4 (if we have 4 classes) 1x1 filters, followed by a softmax activation.

# Convolutional Implementation of Sliding Windows



The original sliding windows algorithm, takes one window (e.g. 14x14x3) and run it through the CNN, then take the next region and so on.
With the convolutional implementation instead of doing it sequentially, the entire image (16x16x3) is input and convolutionally make all the predictions at the same time by one forward pass through the CNN.

*Pierre Sermanet et al. , 2014, OverFeat: Integrated Recognition, Localization and Detection using Convolutional Networks*
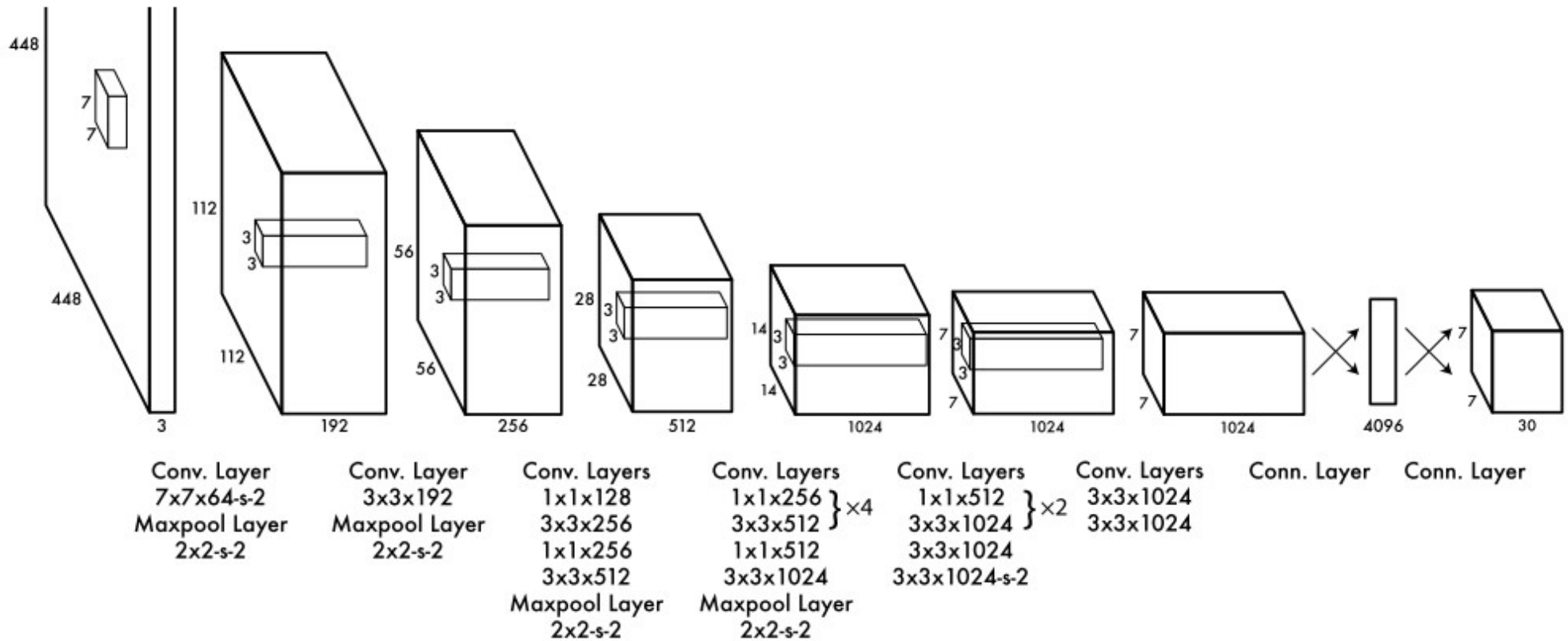
# Bounding Box Predictions



With the Sliding Windows algorithm, it may happen that none of the boxes really match perfectly with the position of the target object.

For example, we want to detect the car in this picture. The sliding windows (since we initially set the windows sizes and we don't not know how big the car is) may match only part of the car.

In some cases, the object may look like a rectangle instead of square.

**Solution: YOLO (You Only Look Once) algorithm** is a way to output more accurate bounding boxes.

# YOLO (You Only Look Once)



YOLO - CNN network for both classification and localising the object using bounding boxes.

24 convolutional layers + 2 fully connected layers.

Conv layers pretrained on ImageNet dataset.

*Redmon et al, 2015, "You Only Look Once: Unified, Real-Time Object Detection" (https://arxiv.org/abs/1506.02640)

Redmon & Farhadi, 2016 (https://arxiv.org/abs/1612.08242).