



# *Computer Vision*

*COMPSCI 760  
2024 Semester 1*

*Olivier Graffeuille*

[ogra439@aucklanduni.ac.nz](mailto:ogra439@aucklanduni.ac.nz)



# Outline for Weeks 6-11

Week 6: Computer Vision – CNNs

Week 7: Computer Vision – Other vision tasks

Week 8: Transfer Learning

Week 9: Continual Learning

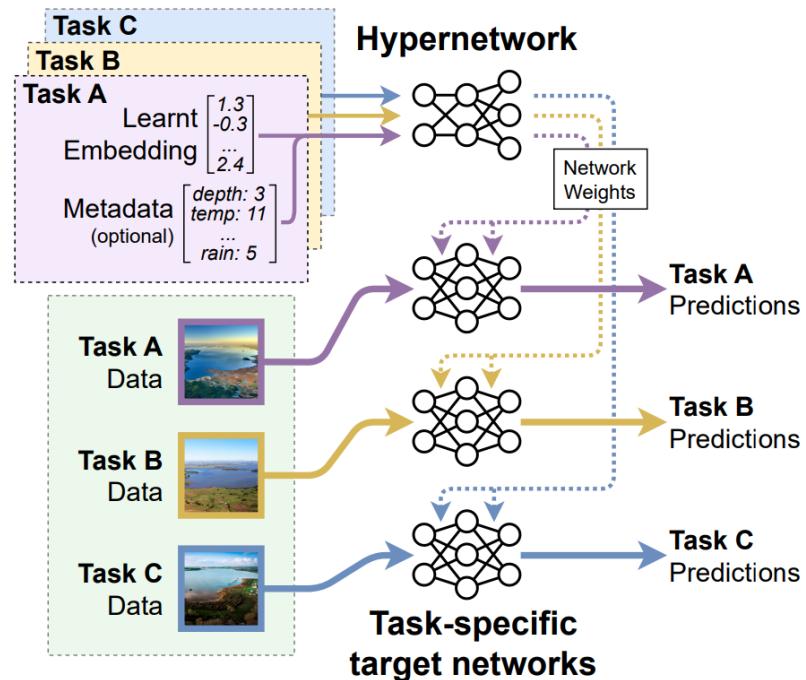
Week 10: Domain Generalisation

Week 11: Self-Supervised Learning

# About me



Taken from orbit in October 2011, the worst algae bloom that Lake Erie has experienced in decades. (ITTNASA)





## Background

# Disclaimer/Acknowledgment

The following slides are reusing some of the content from the Stanford CS230 and CS231n class:

<https://stanford.edu/~shervine/teaching/cs-230/>  
<https://cs231n.github.io/convolutional-networks/>

Some content is also reused from the CS230 Stanford course slides created by Andrew Ng, available at

<https://cs230.stanford.edu/syllabus>

# Outline

1. Background
2. Convolution
3. CNN components: convolutional layers, pooling layers
4. CNN architectures

# Outline

## ***1. Background***

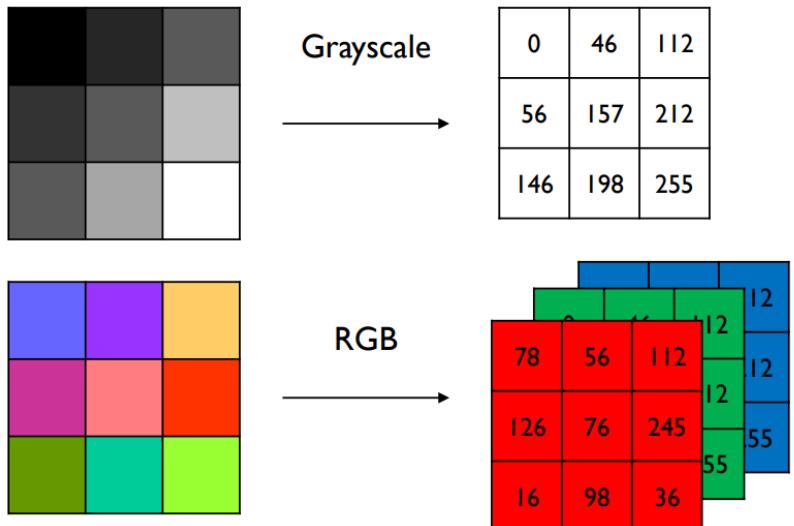
2. Convolution
3. CNN components: convolutional layers, pooling layers
4. CNN architectures

# Background

## Images as vectors

Neural networks usually take a vector of values as input

A digital image is generally represented as a 2D Matrix of pixels

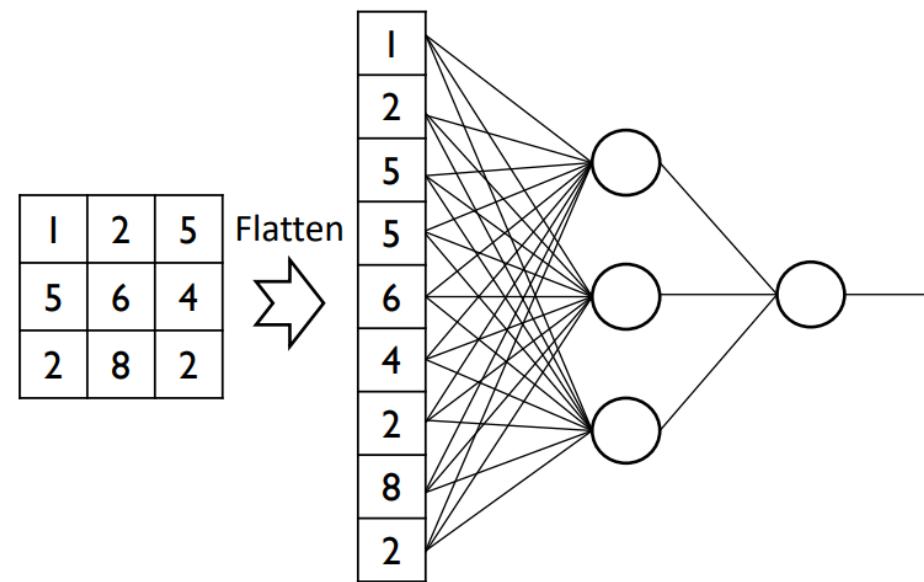


I
2
5
5
6
4
2
8
2

# Background

## Fully connected layers?

Why not just flatten an image into a vector, and use traditional NNs (with fully connected linear layers) to process an image?





## Background

# Fully connected layers?

1. Problem 1: The number of parameters would go quickly:
  1. A 512x512x3 image has 800k values
  2. A linear layer mapping from 800k input neurons to 800k hidden neurons would have 640 billion parameters (more than GPT-4?)



# Background

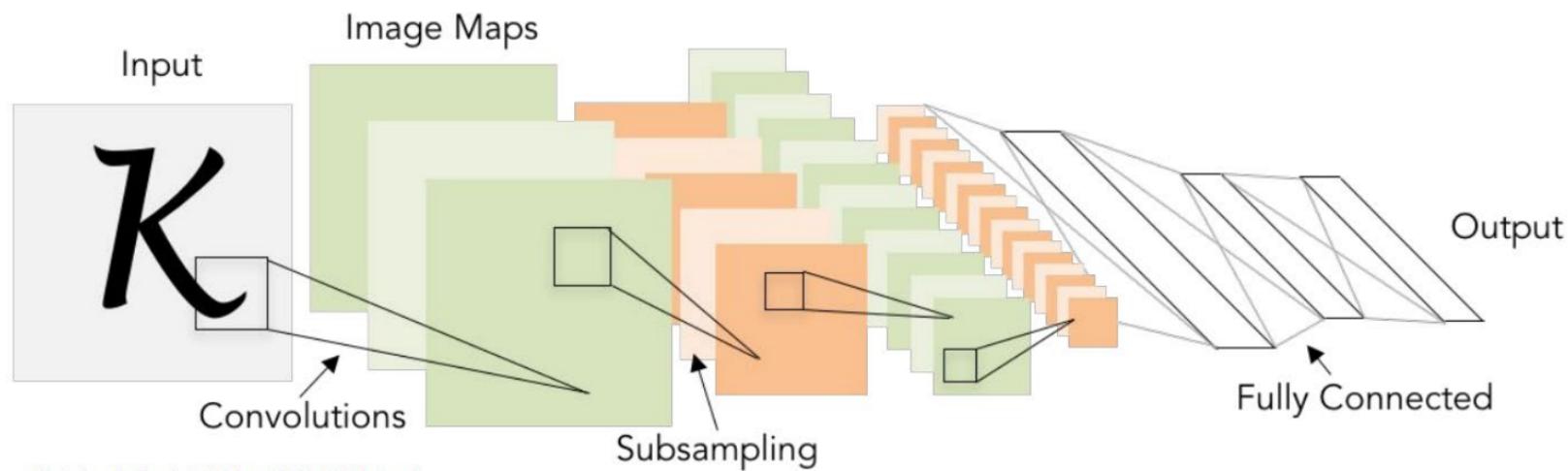
## Fully connected layers?

1. Problem 1: The number of parameters would go quickly:
  1. A 512x512x3 image has 800k values
  2. A linear layer mapping from 800k input neurons to 800k hidden neurons would have 640 billion parameters (more than GPT-4?)
2. Problem 2: Would lead to overfitting
  1. In a given dataset, potentially:  
if pixel #14,153 is greener than pixel #784,150 then the image is a cat
  2. The model's inductive bias is not suitable for images

# Background

## Images as vectors

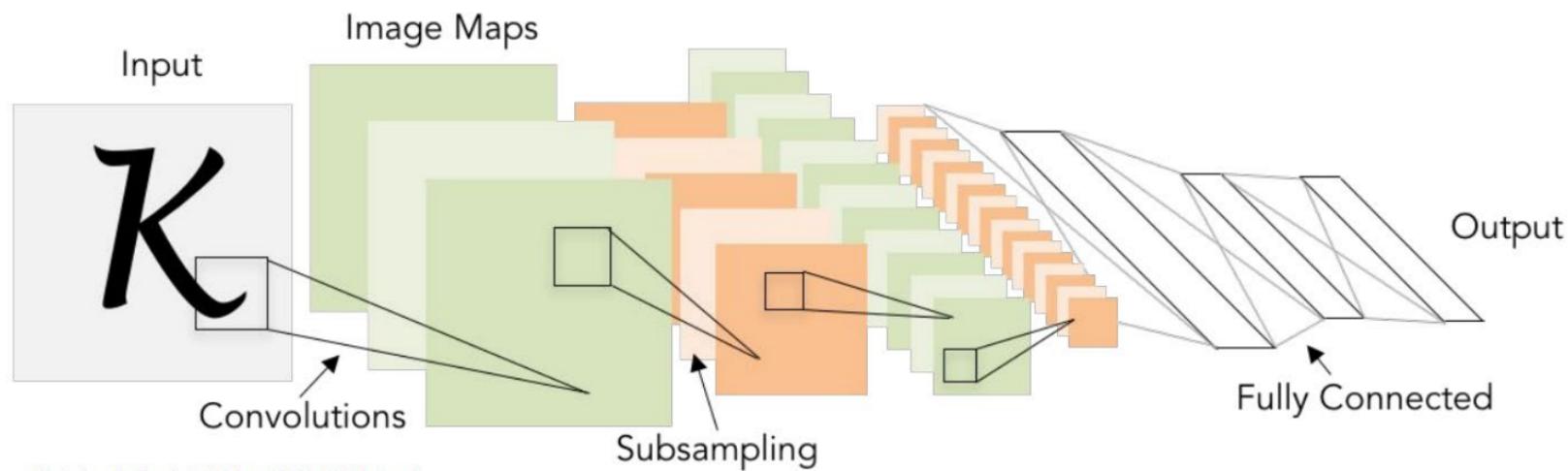
1. Instead, we should consider pixel locality



# Background

## Images as vectors

1. Main differences with a classic ANN:
  1. Explicit assumption that inputs are images.
  2. Introduce new types of layers: convolution and pooling layers





# Background

## Images as vectors

1. The benefits of CNNs:
  1. Make images easier to process
  2. Extract meaningful features from the image.
2. They are the go-to architecture for computer vision tasks\*

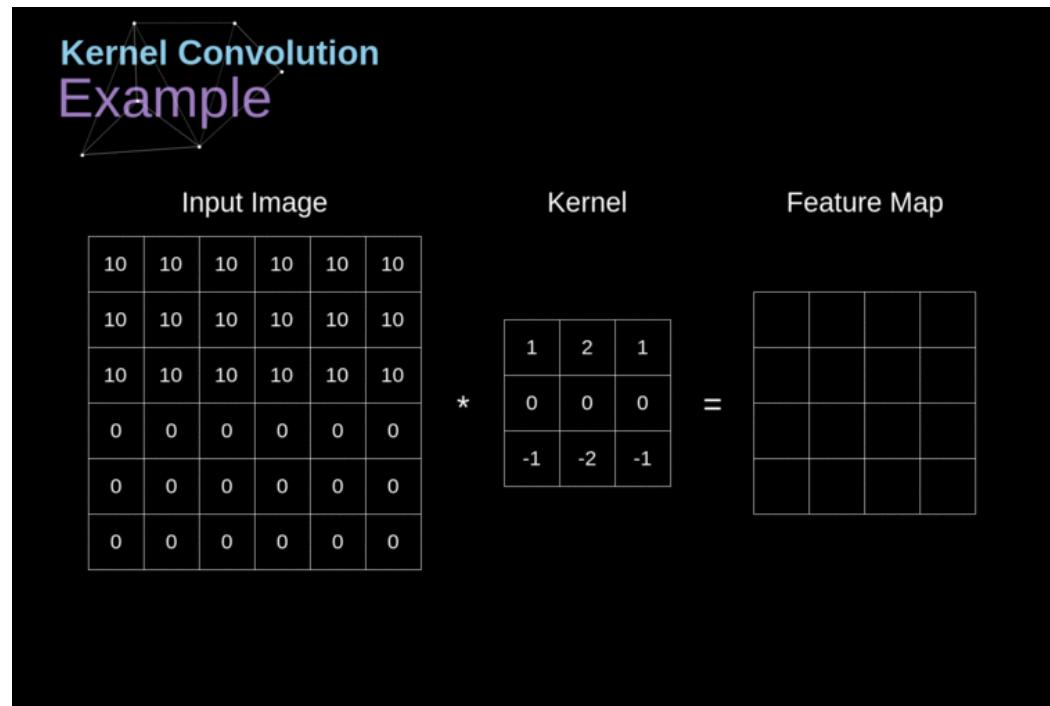
# Outline

1. Background
- 2. *Convolution***
3. CNN components: convolutional layers, pooling layers
4. CNN architectures

# Convolution

## Kernels

1. Convolution is common in image processing to filter images using **kernels**
  1. Edge detection



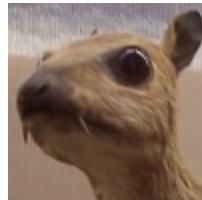


# Convolution

## Kernels

1. Convolution is common in image processing to filter images using **kernels**

1. Edge detection
2. Smoothing



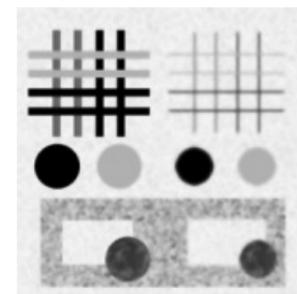
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	

# Convolution

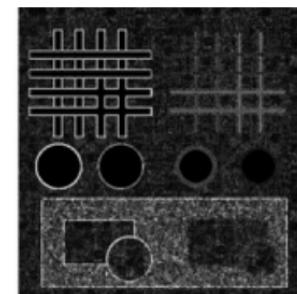
## Kernels

1. Convolution is common in image processing to filter images using **kernels**

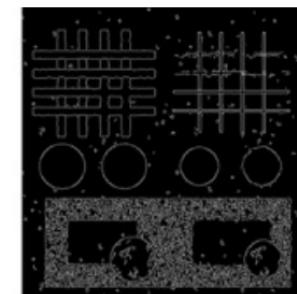
1. Edge detection
2. Smoothing
3. Other filters



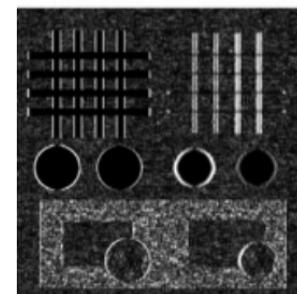
Original



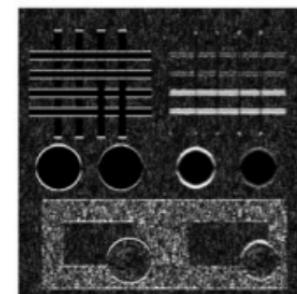
Laplacian



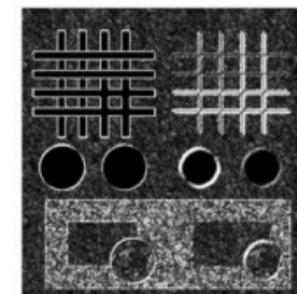
Canny



Sobel X



Sobel Y



Sobel X+Y

Sobel

-1	0	1
-2	0	2
-1	0	1

1	2	1
0	0	0
-1	-2	-1

Sobel X

Sobel Y

Laplacian

0	-1	0
-1	4	-1
0	-1	0

-1	-1	-1
-1	8	-1
-1	-1	-1



# Convolution

## Kernels

What filters should we use?

- Define them as neural network weights
- Learn them through backpropagation

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

\*

$w_1$	$w_2$	$w_3$
$w_4$	$w_5$	$w_6$
$w_7$	$w_8$	$w_9$

=


# Convolution

## Convolutional Layer

Sliding kernel (filter) over the image

Computationally:  
matrix multiplication + sum

Convolution of a  $f \times f$  filter  
with a  $n \times n$  image  
produces a  $n-f+1 \times n-f+1$  feature map

3 <sub>0</sub>	3 <sub>1</sub>	2 <sub>2</sub>	1	0
0 <sub>2</sub>	0 <sub>2</sub>	1 <sub>0</sub>	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

# Convolution

## Convolutional Layer

Sliding kernel (filter) over the image

Computationally:  
matrix multiplication + sum

Convolution of a  $f \times f$  filter  
with a  $n \times n$  image  
produces a  $n-f+1 \times n-f+1$  feature map

3 <sub>0</sub>	3 <sub>1</sub>	2 <sub>2</sub>	1	0
0 <sub>2</sub>	0 <sub>2</sub>	1 <sub>0</sub>	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

Problem: shrinking size, and border pixels  
don't contribute much to feature map

# Convolution

## Padding

0	0	0	0	0	0	0	0
0	I	I	I	0	0	0	0
0	0	I	I	I	0	0	0
0	0	0	I	I	I	0	0
0	0	0	I	I	0	0	0
0	0	I	I	0	0	0	0
0	0	0	0	0	0	0	0

padding p = 1

$$\begin{matrix} * & \begin{matrix} I & 0 & I \\ 0 & I & 0 \\ I & 0 & I \end{matrix} & = & \begin{matrix} 2 & 2 & 3 & I & I \\ I & 4 & 3 & 4 & 2 \\ I & 2 & 4 & 3 & 3 \\ I & 2 & 3 & 4 & I \\ 0 & 2 & 2 & I & I \end{matrix} \end{matrix}$$

Convolution of a  $f \times f$  filter with a  $n \times n$  image, with p padding produces a  $n+2p-f+1 \times n+2p-f+1$  feature map



# Convolution

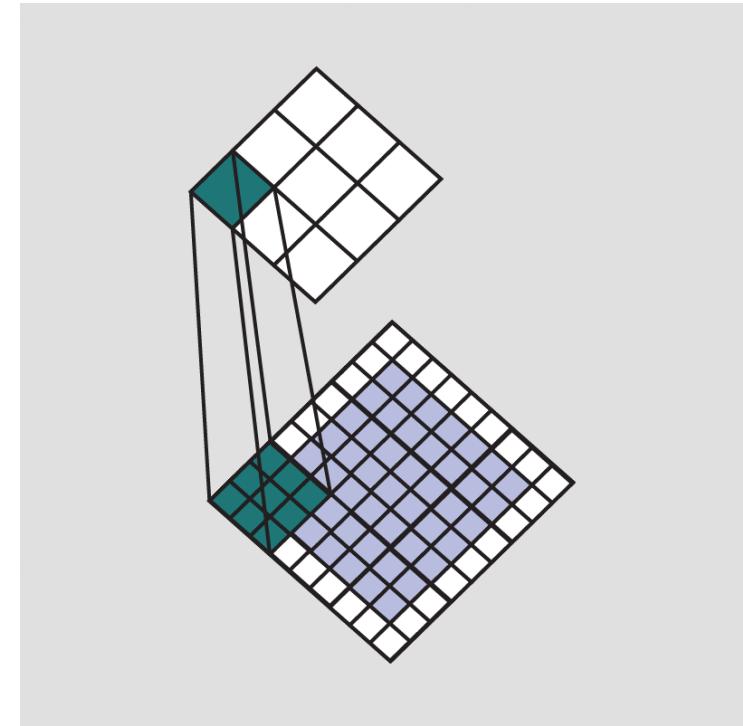
## Padding

- Valid padding (no padding)
  - Output:  $n-f+1 \times n-f+1$  feature map
- Same padding (padding to make output feature map the same size as the input)
  - Output:  $n+2p-f+1 \times n+2p-f+1$  feature map
  - What  $p$  should you use?  
$$n = n+2p-f+1$$
$$p = (f-1) / 2$$
 (depends on filter/kernel size!)
  - Filter size  $f$  is usually odd (3x3, 5x5...) to have integer  $p$

# Convolution

## Stride

- Strided convolutions reduce the size of the output



Convolution of  $f \times f$  filter on a  $n \times n$  image with padding  $p$  and stride  $s$  produces feature map of size  $\frac{n + 2p - f}{s} + 1 \times \frac{n + 2p - f}{s} + 1$



# Convolution

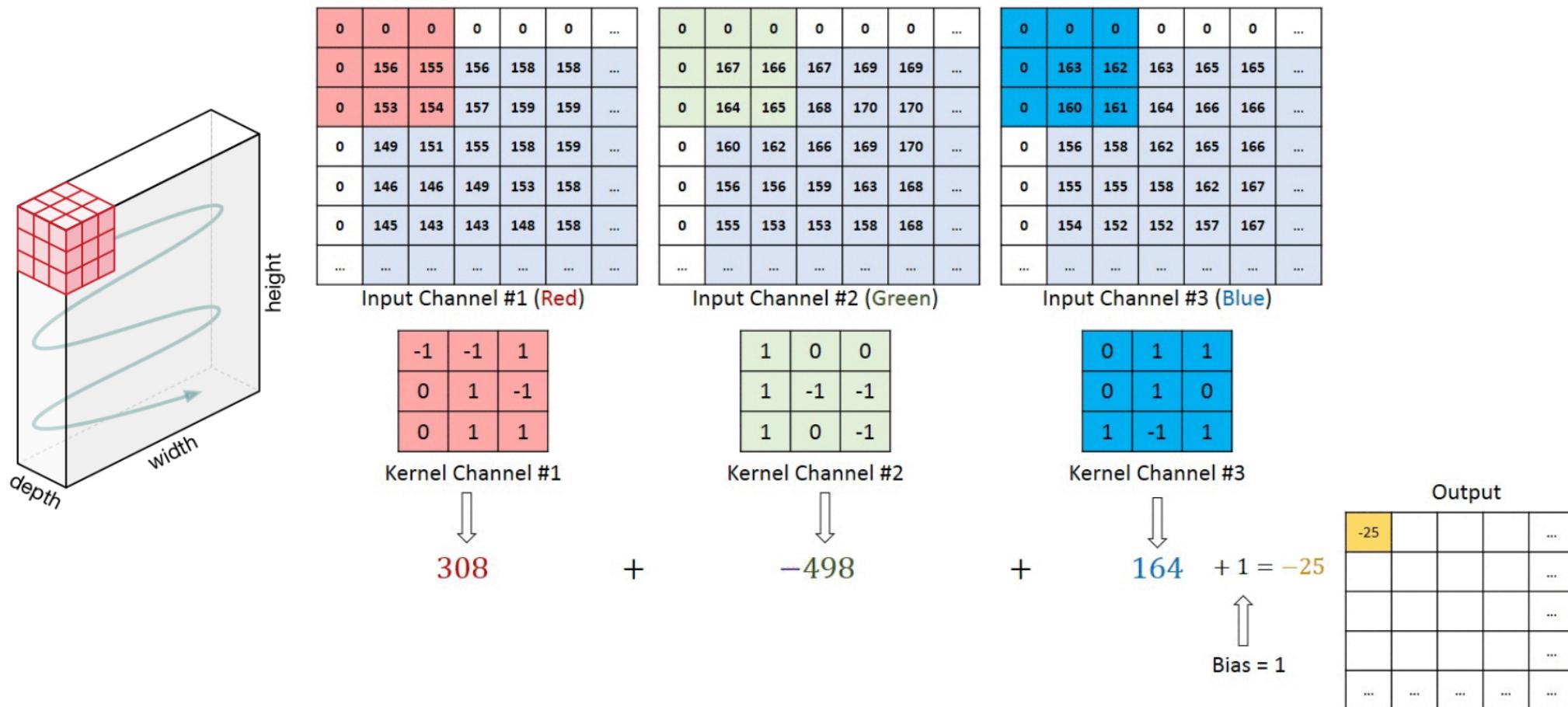
## Stride

Why use strided convolutions to reduce the size of the output?

- Reduces computational load
- Increases the receptive field of the features
- Alternative to pooling

# Convolution

## Images are 3D tensors

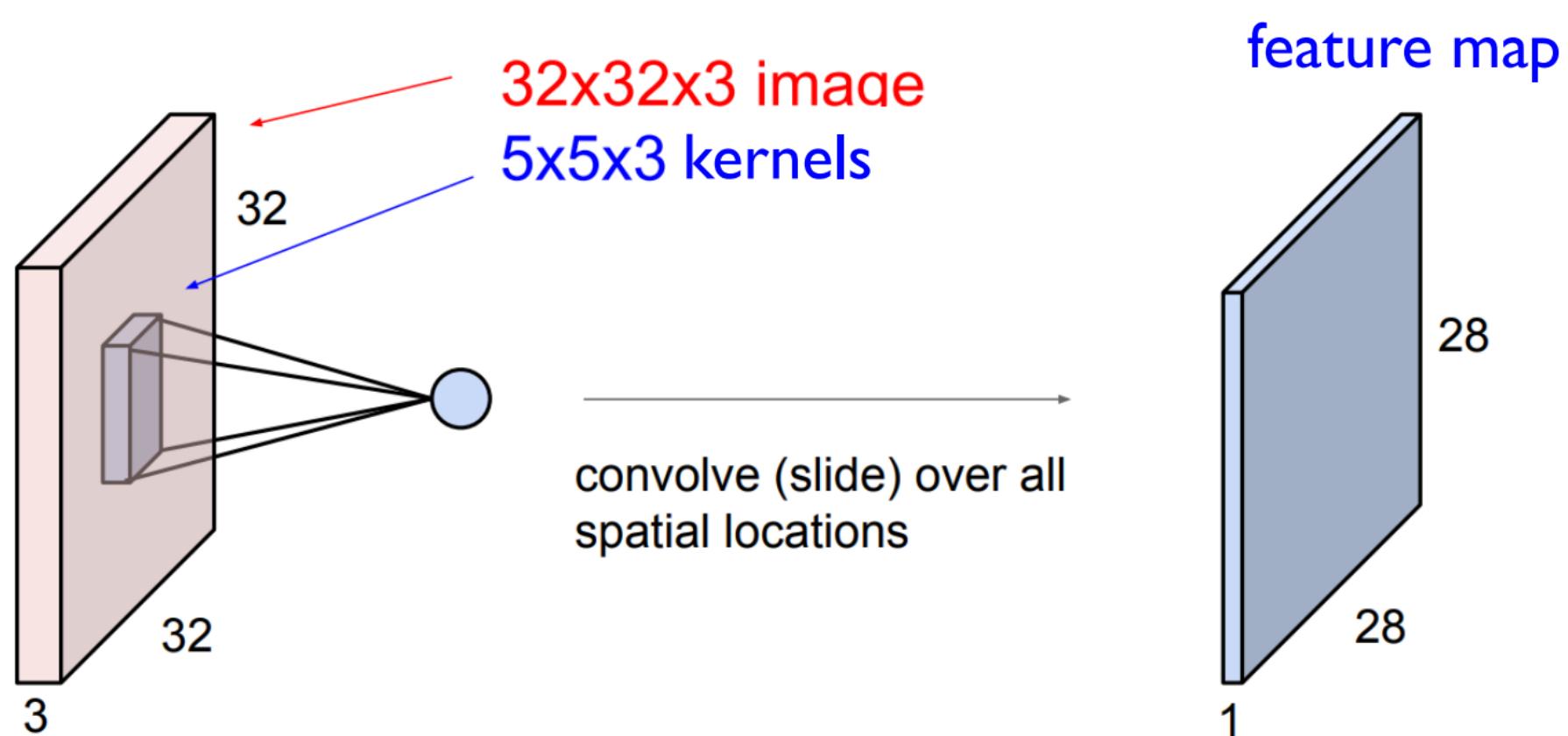


# Outline

1. Background
2. Convolution
- 3. *CNN components: convolutional layers, pooling layers***
4. CNN architectures

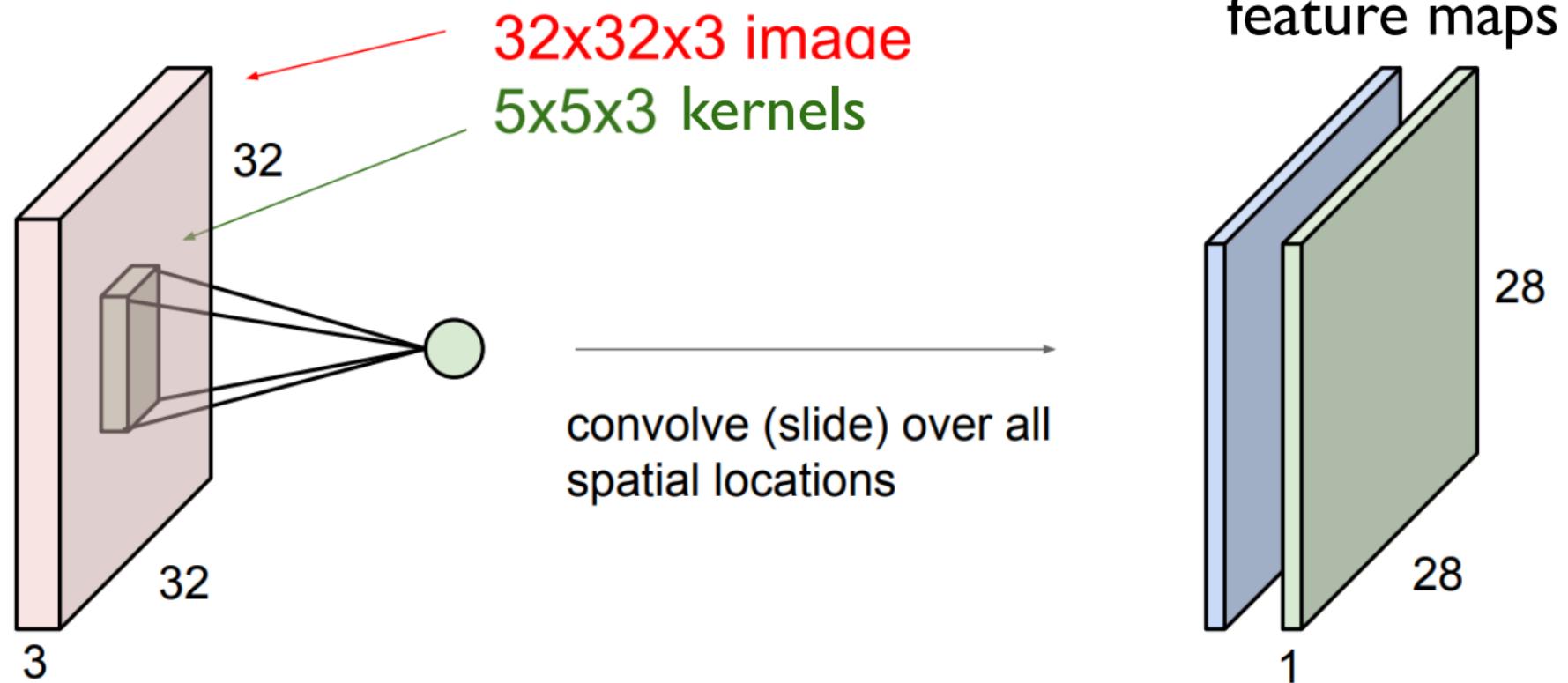
# Convolutional Layers

## Making a convolutional layer



# Convolutional Layers

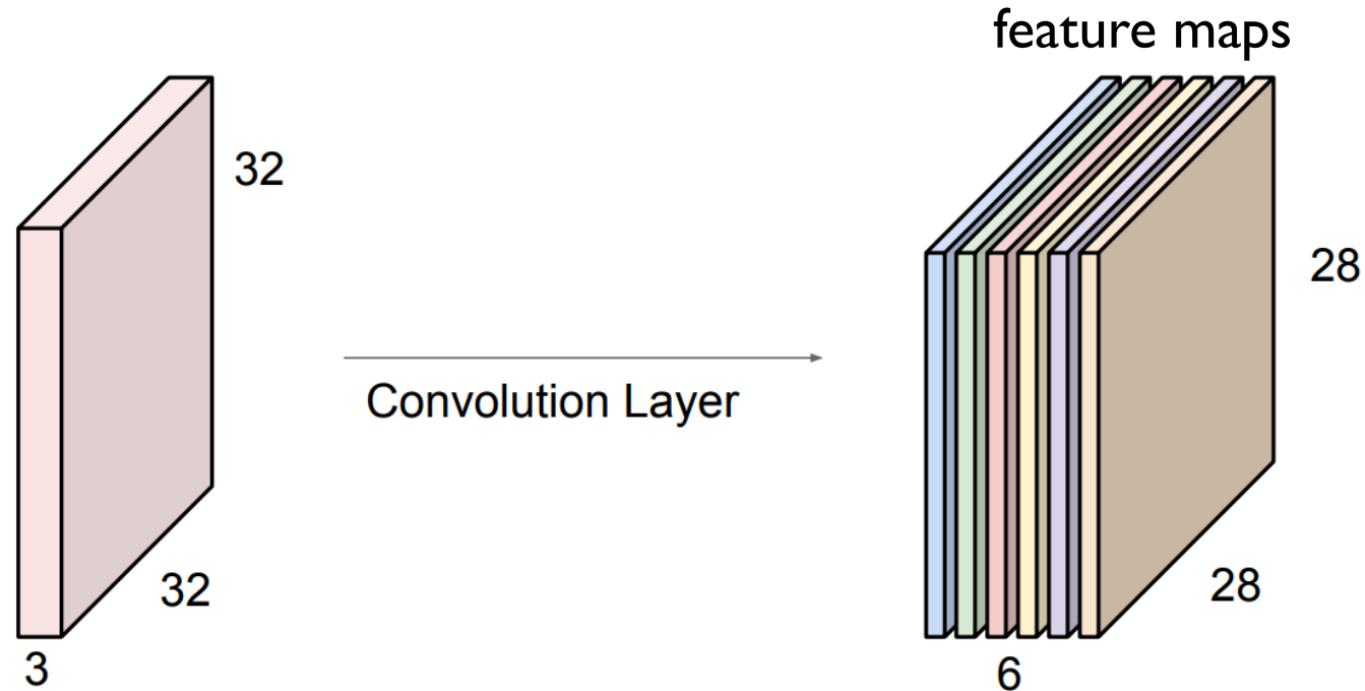
## Making a convolutional layer



# Convolutional Layers

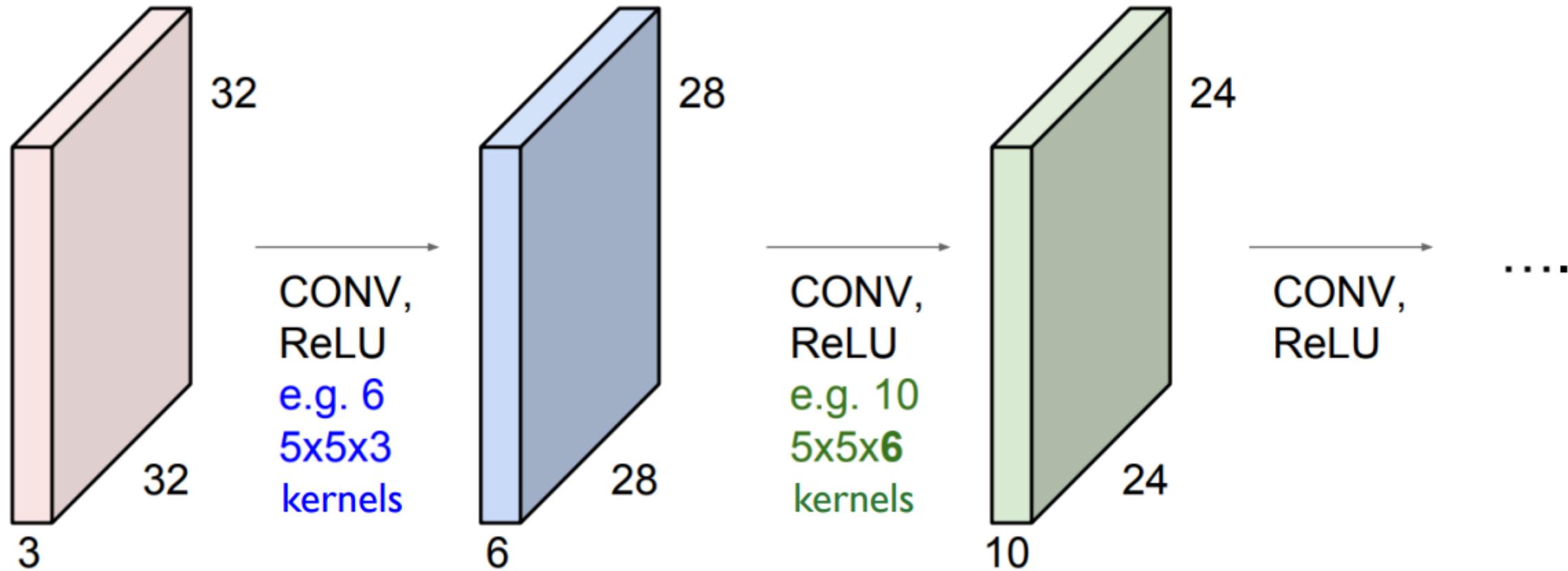
## Making a convolutional layer

- Stacking multiple kernels produces new 3D “image” (feature map)



# Convolutional Layers

## Multiple convolutional layer





# Convolutional Layers

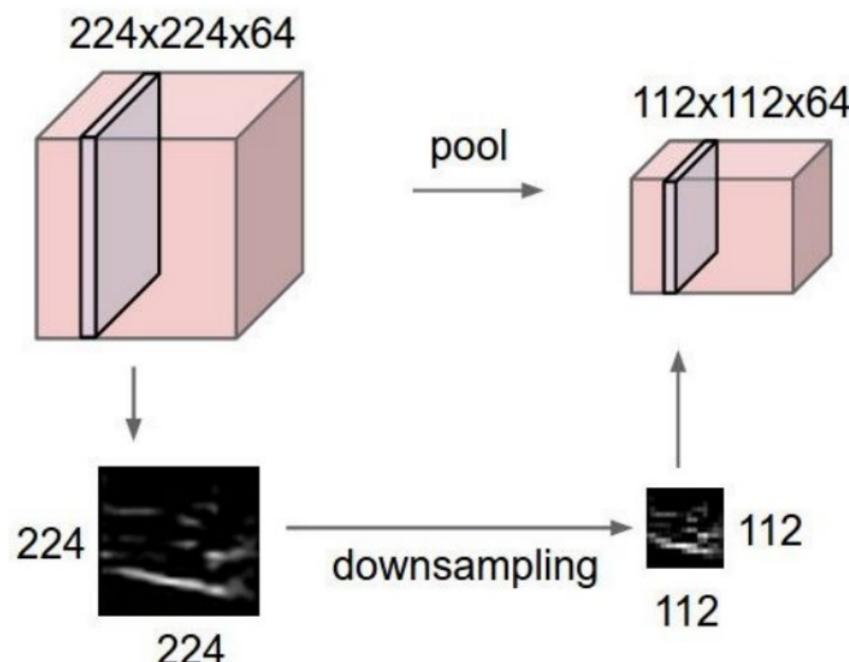
## 1x1 Convolutions

- Many models use 1x1 convolutions, which seems confusing
- On a 2D image, a 1x1 convolution is just a scaling factor  
 $\text{out} = \text{w} * \text{in}$
- But since images (and feature maps) are 3D, they are a layer-wise linear transformation

# Pooling Layers

## Pooling

- Makes representations smaller (similar to strides)
- Applies independently on each feature map (spatial aggregation)



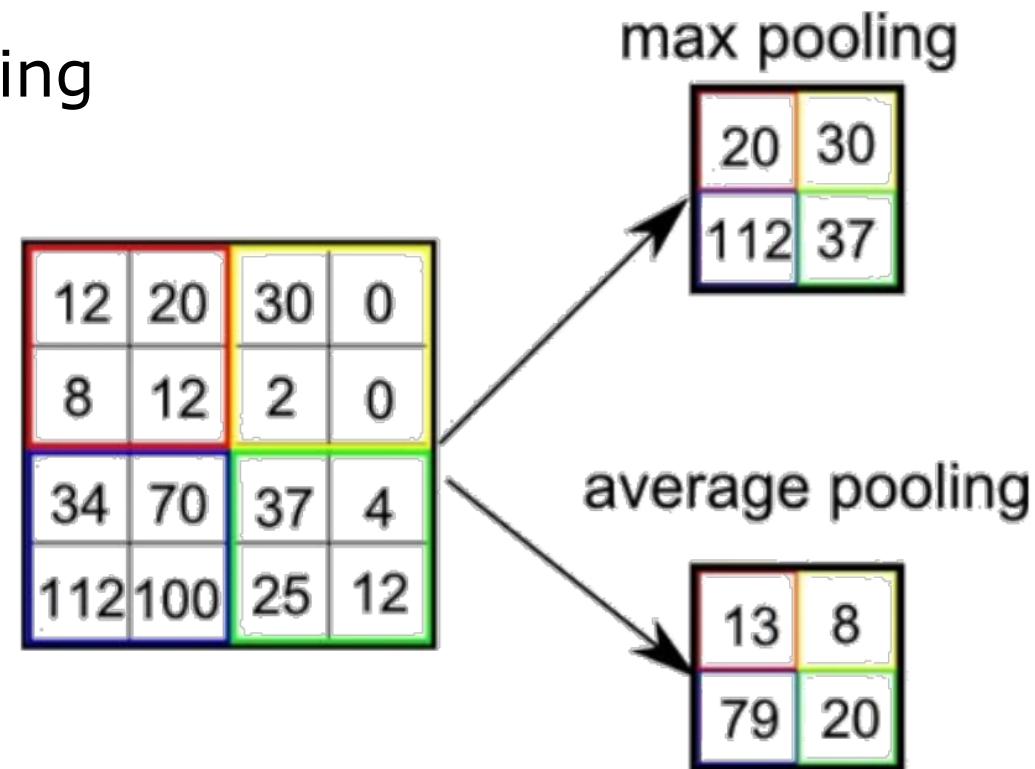
# Pooling Layers

## Pooling

- Two main types of pooling

kernel size  $f = 2$

stride  $s = 2$



- Max pooling is usually preferred (keeps the most significant features)



# Pooling Layers

## Hyperparameters

### Convolutional layer

f : kernel size

s : stride

p : padding

$n'_C$  : number of filters

If input size is  $n_H \times n_W \times n_C$  :

Output size of convolution will be:

$$\left\lfloor \frac{n_H+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W+2p-f}{s} + 1 \right\rfloor \times n'_C$$

### Pooling layer

f : kernel size

s : stride

Max or average pooling

Output size of pooling will be:

$$\left\lfloor \frac{n_H-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W-f}{s} + 1 \right\rfloor \times n_C$$

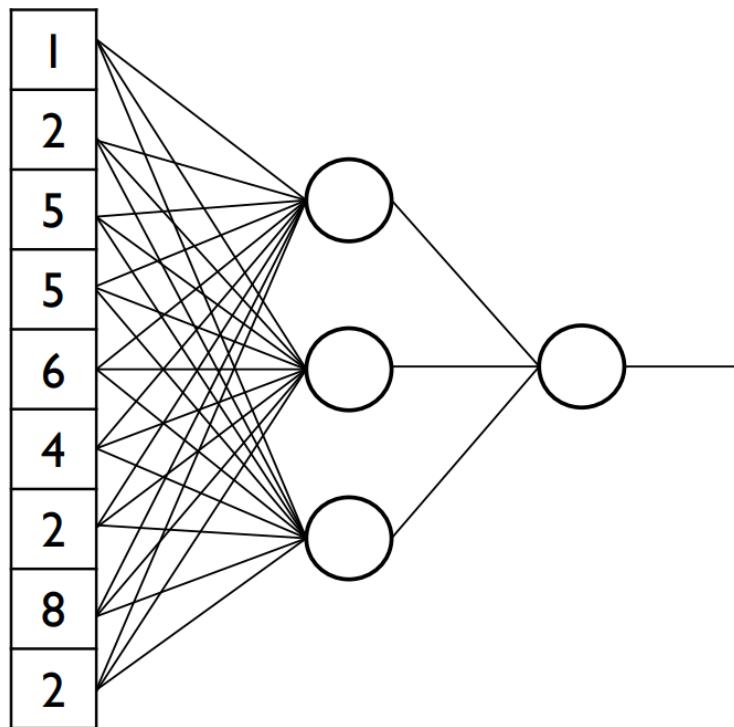
# Linear Layers

## Fully Connected Layer

For Classification:

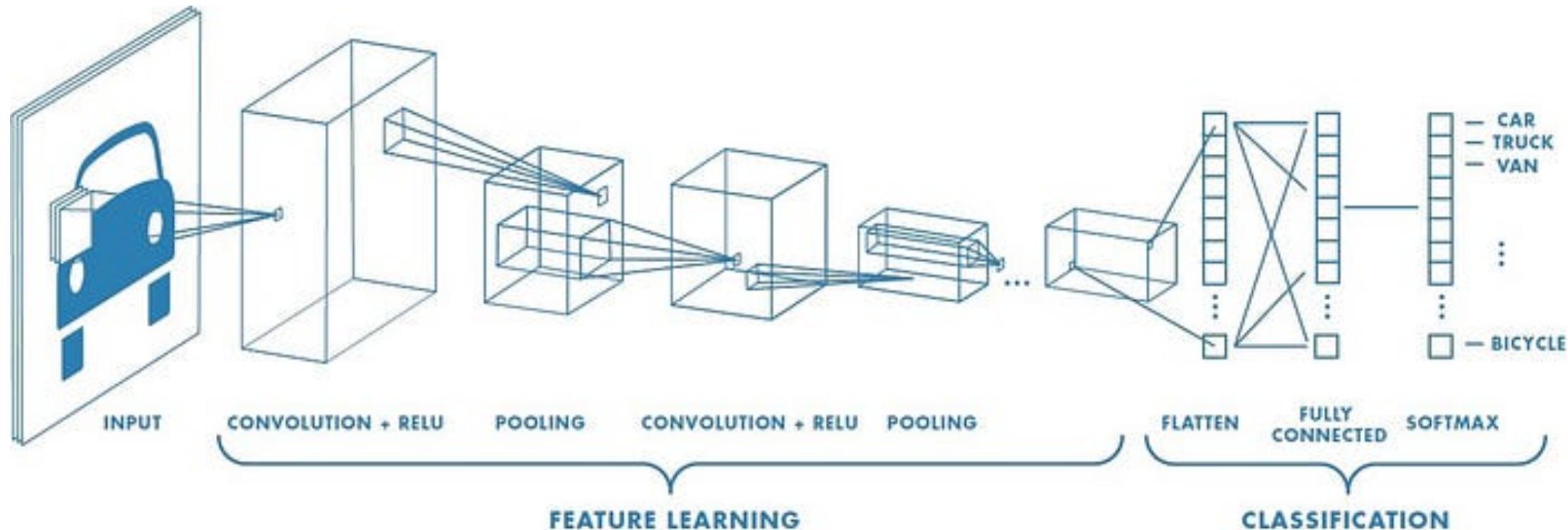
1	2	5
5	6	4
2	8	2

Flatten 



Feature map from the last convolution or pooling layer.

## Putting it all Together



# Outline

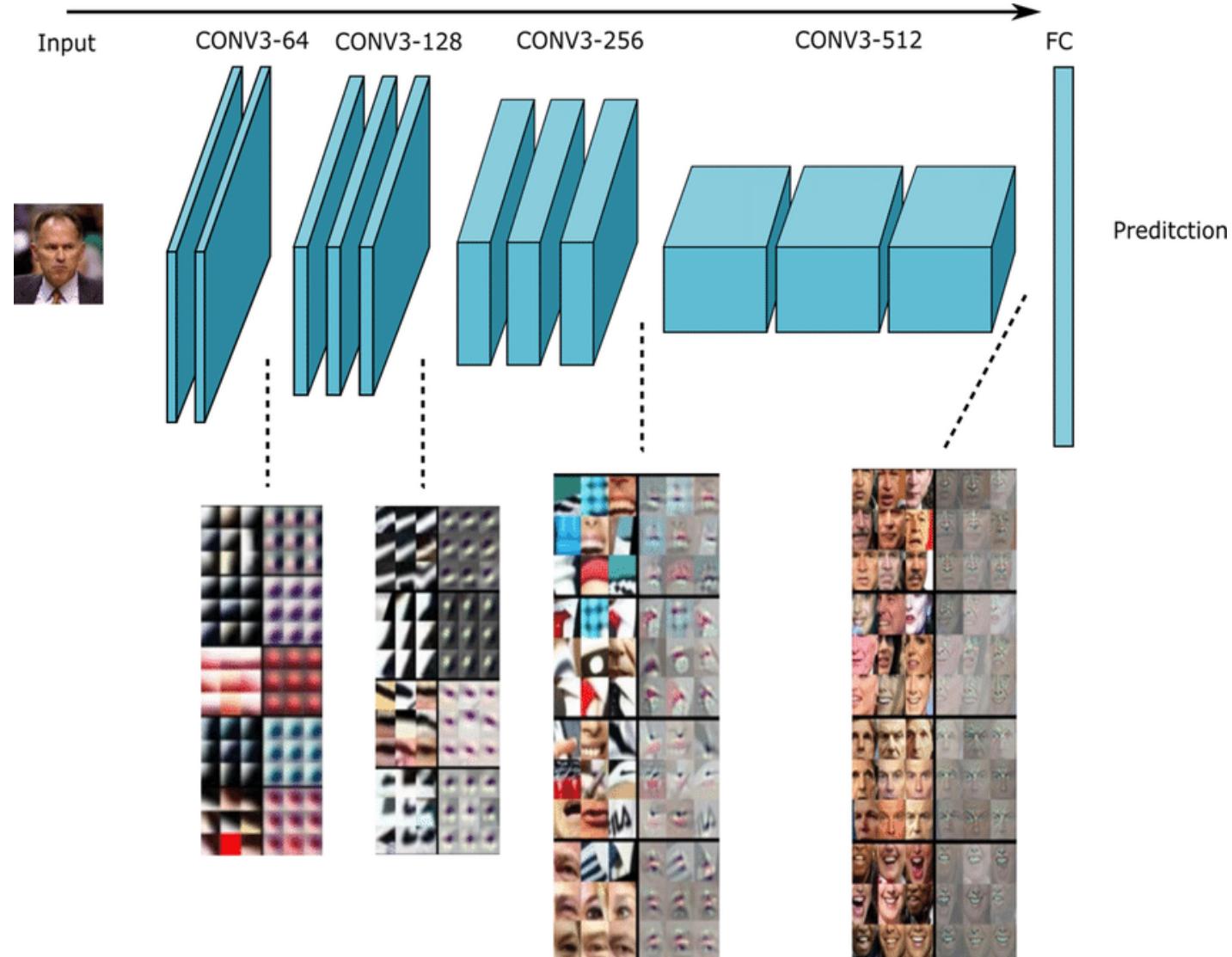
1. Background
2. Convolution
3. CNN components: convolutional layers, pooling layers
- 4. *CNN architectures***

# CNNs

## Deep Representations

Composition of convolutional layers allows these models to learn increasingly complex visual features

Biological vision systems do this



# Visualising Deep Representations

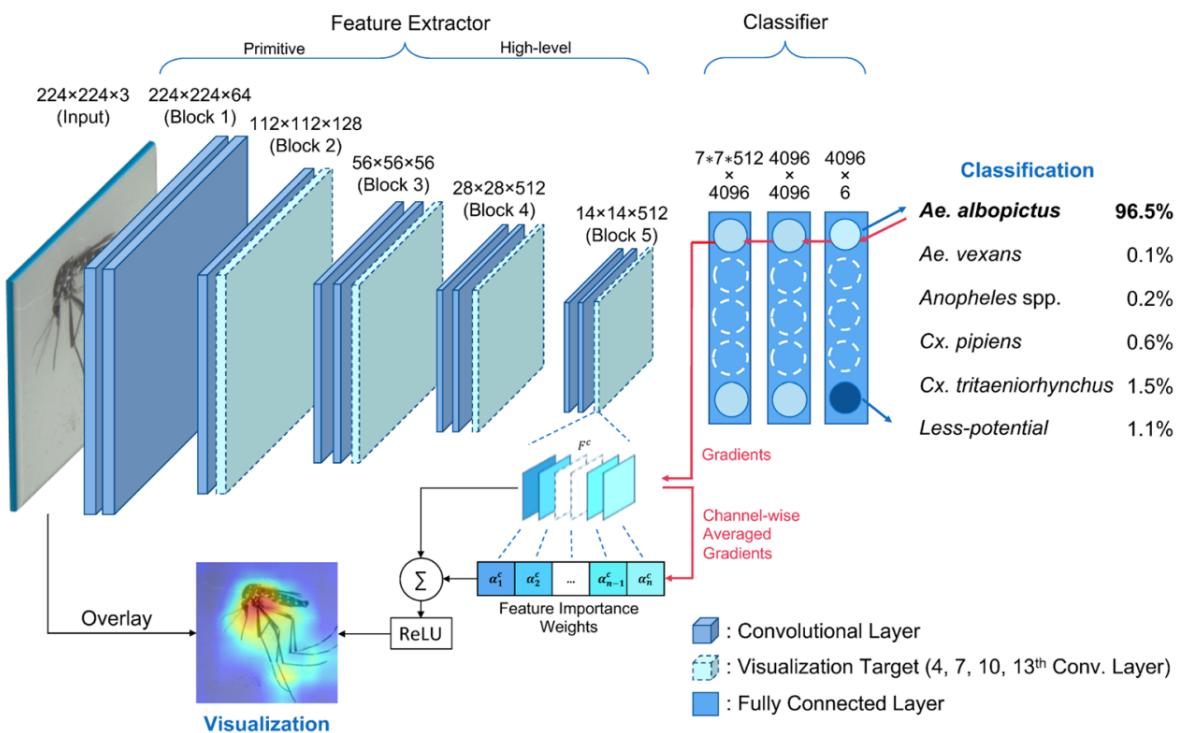


Figure 4. The visualization of discriminative regions in the 4-, 7-, 10- and 13-th convolution layers of VGG-16 when an *Aedes albopictus* input image is given.



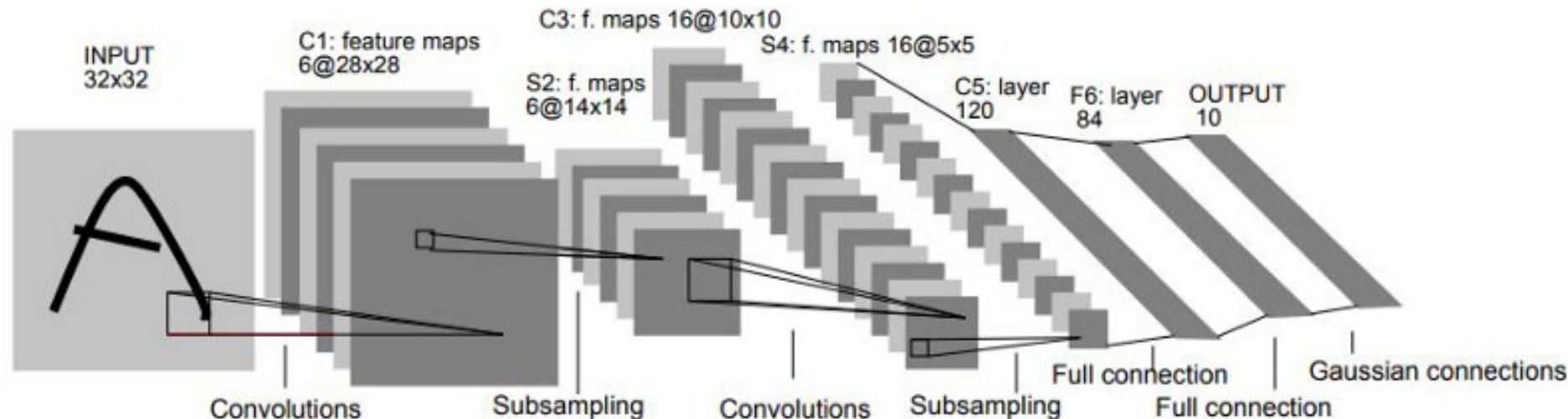
# CNN architectures

- CNN architectures are effectively a combination of the previous building blocks
  - This is largely an engineering problem
- Architectures are often effective across diverse domains

# CNN architectures

## LeNet-5 (1998)

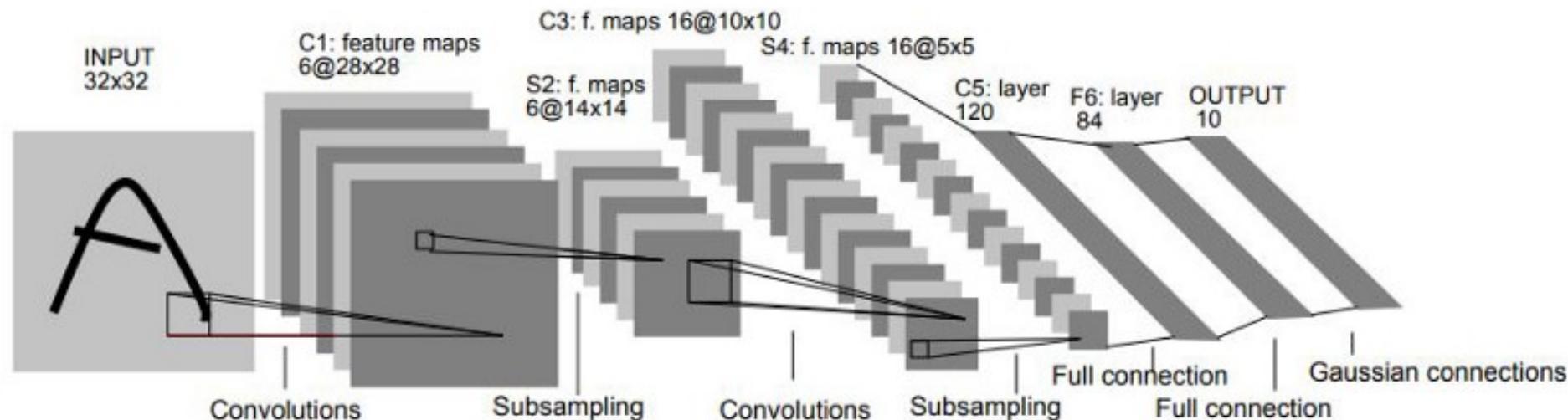
- Designed for handwritten digit recognition (MNIST)
- 60k parameters



# CNN architectures

## LeNet-5 (1998)

- Convolution kernels: kernel size  $f = 5$  ( $5 \times 5$ ), stride  $s = 1$
- Pooling layers: kernel size  $f = 2$  ( $2 \times 2$ ), stride  $s = 2$

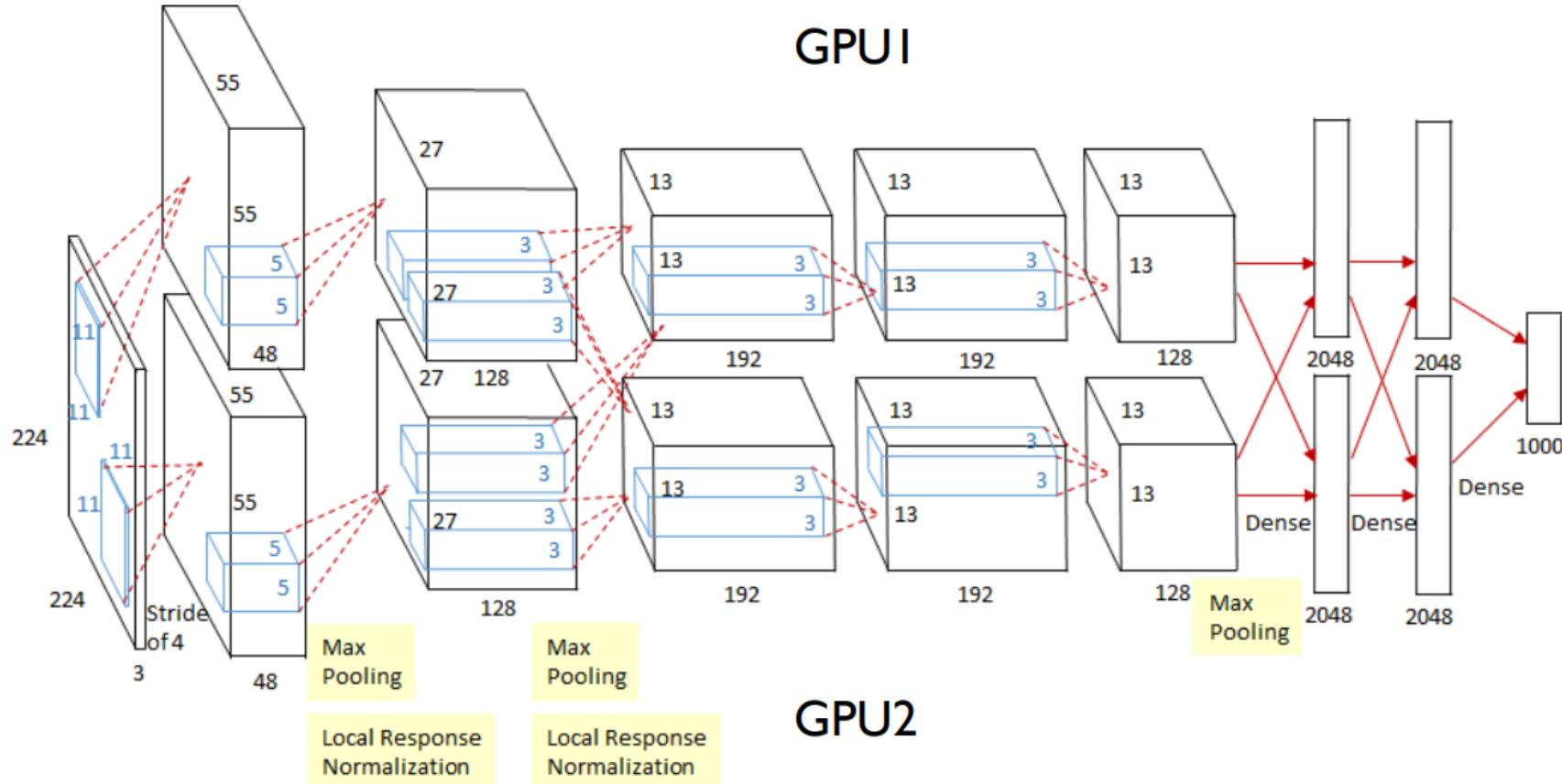


# CNN architectures

## AlexNet (2012)

Run on two GPUs due to memory constraints

60 million parameters





## AlexNet (2012)

### **Full (simplified) AlexNet architecture:**

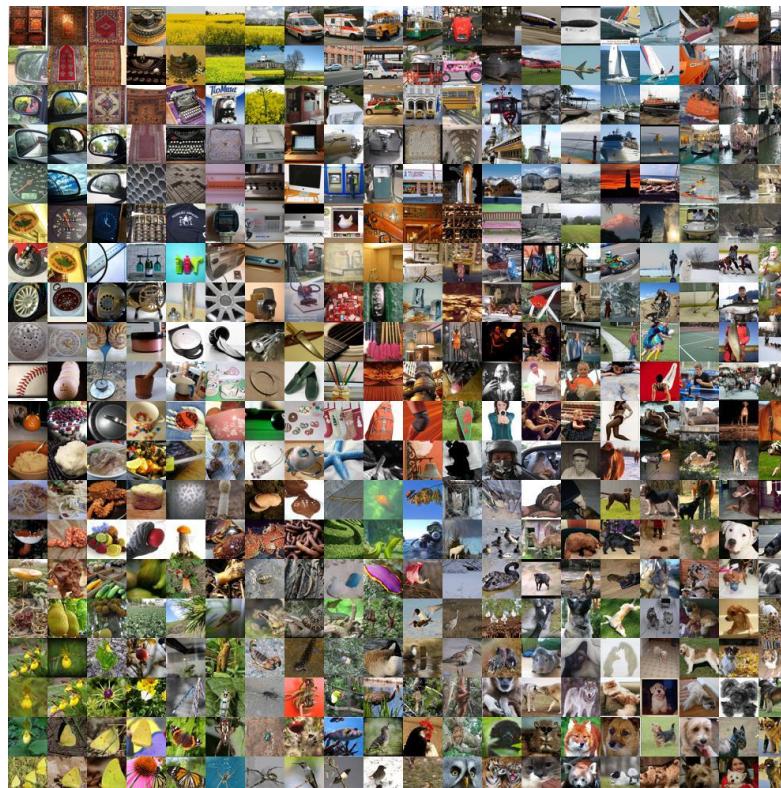
[227x227x3] INPUT  
[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0  
[27x27x96] MAX POOL1: 3x3 filters at stride 2  
[27x27x96] NORM1: Normalization layer  
[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2  
[13x13x256] MAX POOL2: 3x3 filters at stride 2  
[13x13x256] NORM2: Normalization layer  
[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1  
[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1  
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1  
[6x6x256] MAX POOL3: 3x3 filters at stride 2  
[4096] FC6: 4096 neurons  
[4096] FC7: 4096 neurons  
[1000] FC8: 1000 neurons (class scores)

### **Remarks/Hyperparameters:**

- First use of ReLU
- Use of normalisation layers (not common anymore)
- Heavy data augmentation (x 2048!)
- Dropout: 0.5
- Batch size: 128
- SGD Momentum: 0.9
- Learning rate: 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 regularisation weight decay: 5e-4

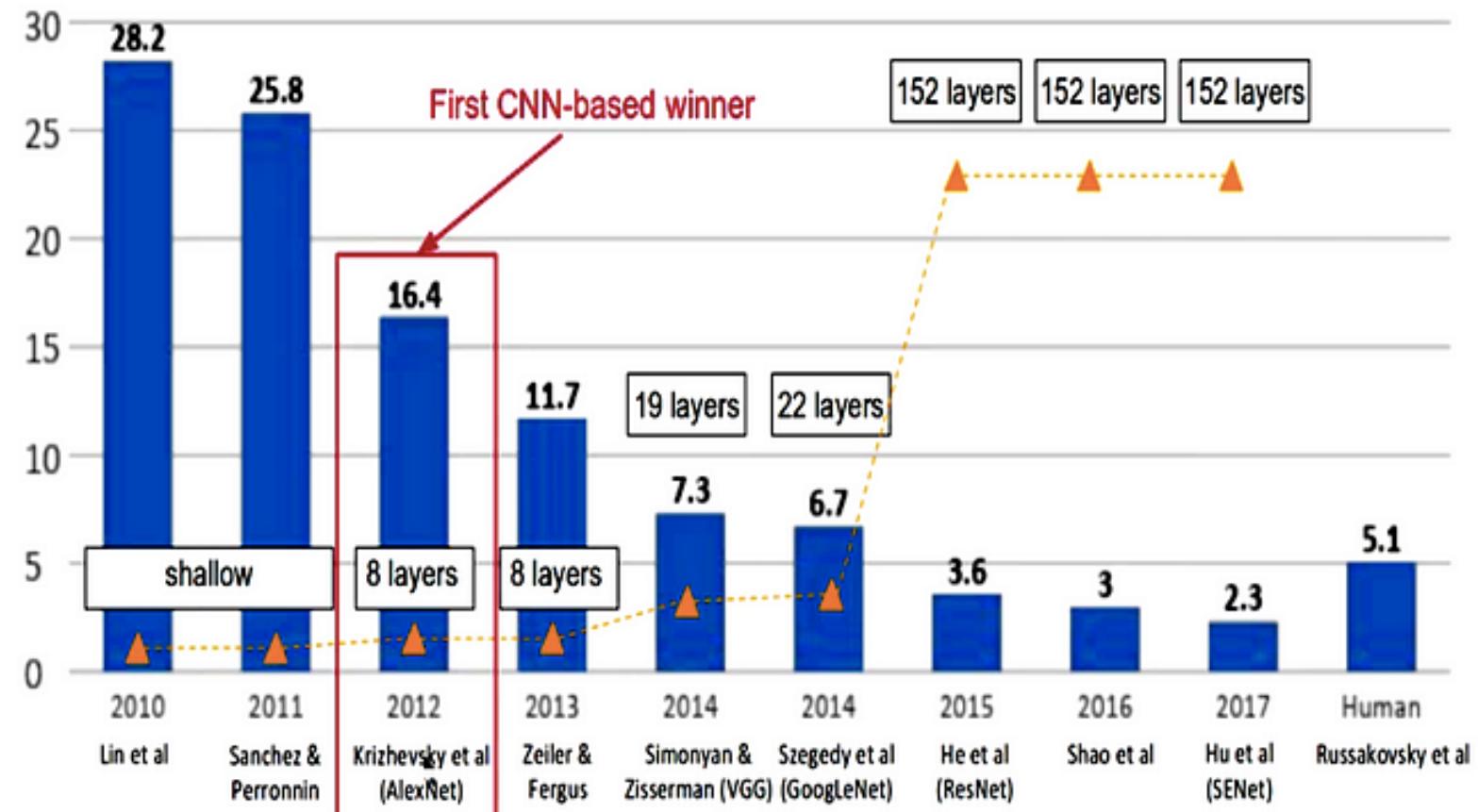
## ImageNet Classification Challenge

- Competition to achieve the highest accuracy on a visual classification dataset
- >1m images
- 1000 classes



# CNN architectures

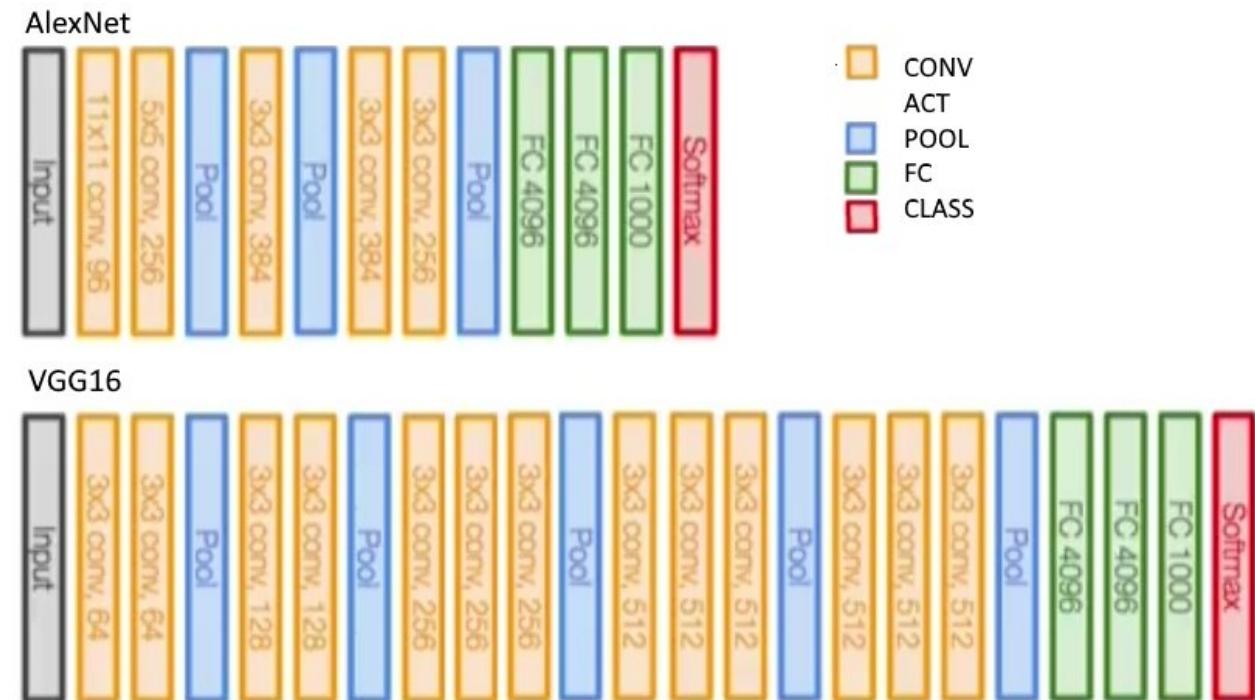
## ImageNet Classification Challenge



# CNN architectures

## VGGNet (2014)

- Small filters
  - Stack of small filters have the same effect than one larger one
- More layers
  - More non-linearities
  - More abstract feature representations



However, adding more layers than this would hurt performance...

# CNN architectures

## ResNet (2015)

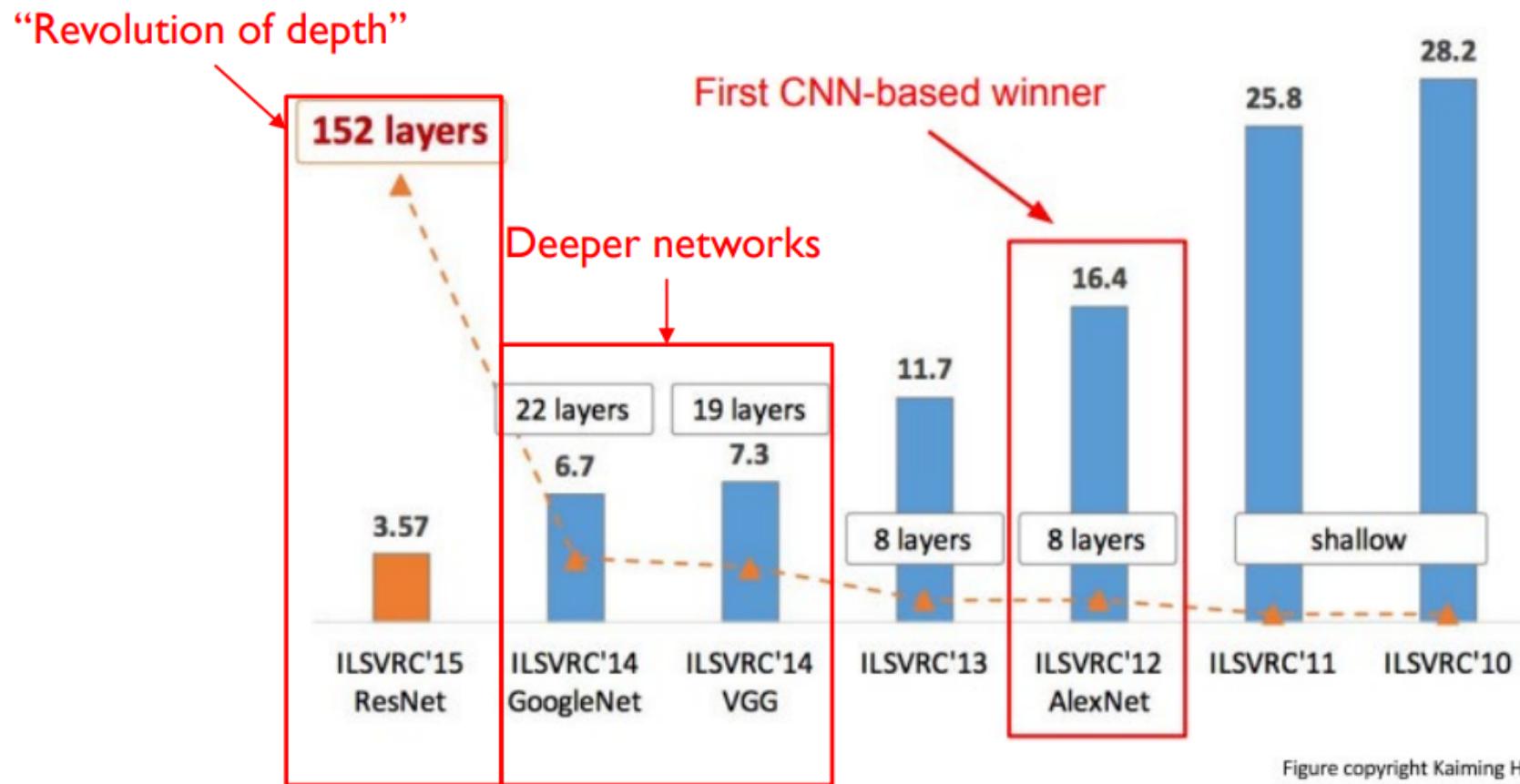
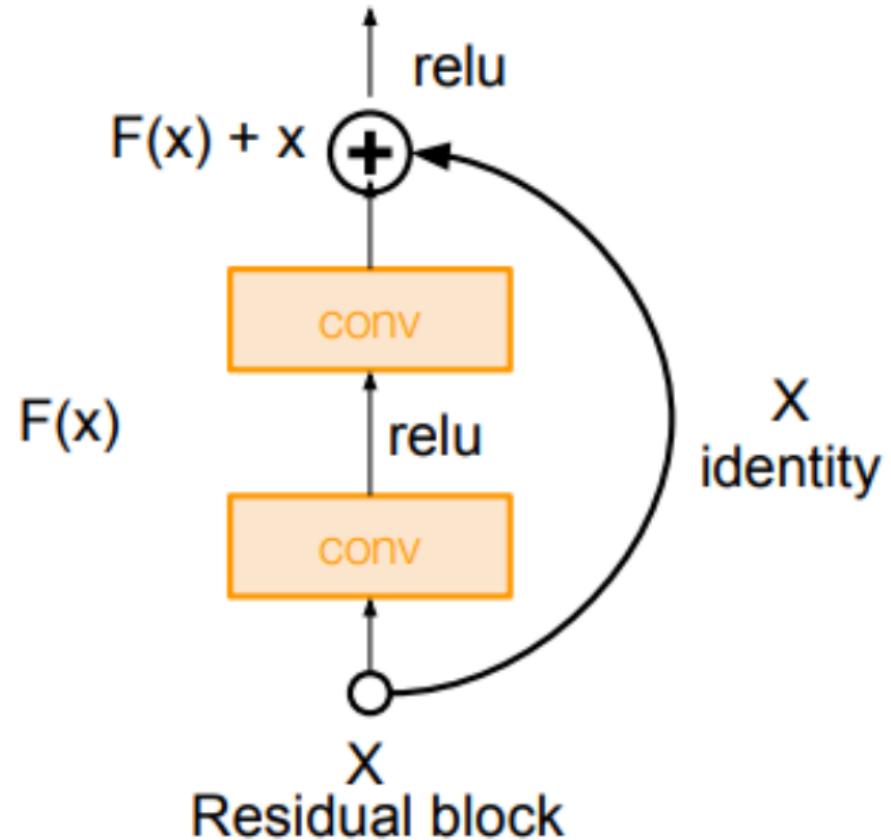


Figure copyright Kaiming He, 2016.

# CNN architectures

## ResNet (2015)

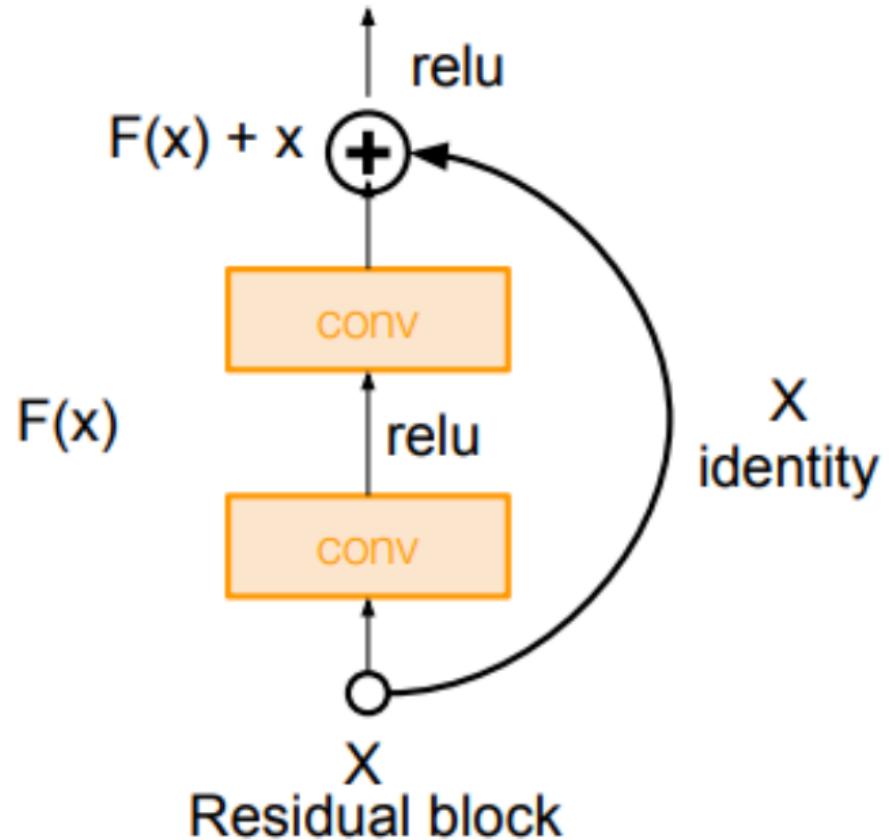
- Main novelty: residual connections
- Architecture consists of residual blocks



# CNN architectures

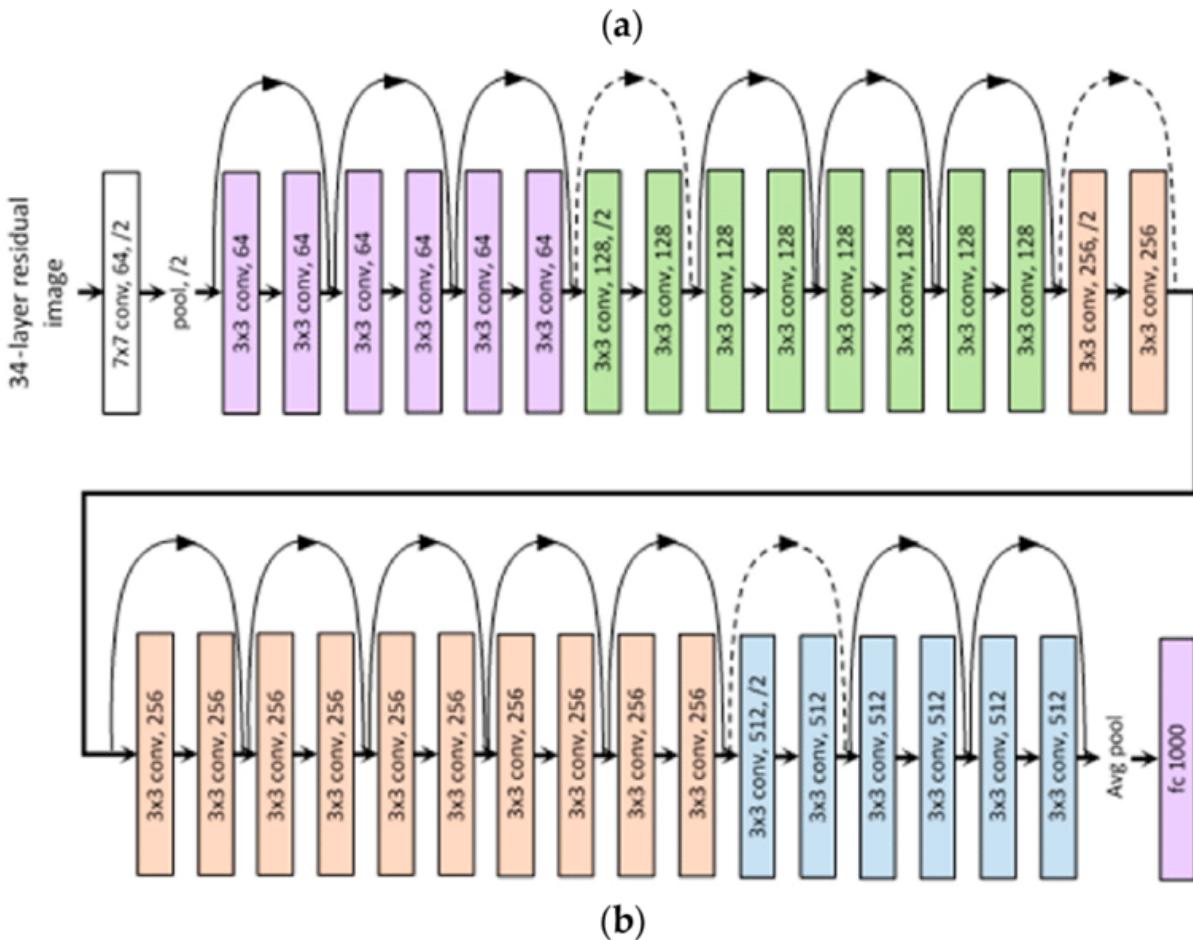
## ResNet (2015)

- Why? Deeper networks!
- Training deep networks is hard
  - Vanishing gradient problem
  - Composition of many functions is difficult to optimise
- Residual blocks offer a skip/shortcut
  - We only need to learn deviation from the identity
- Commonly used in most deep architectures



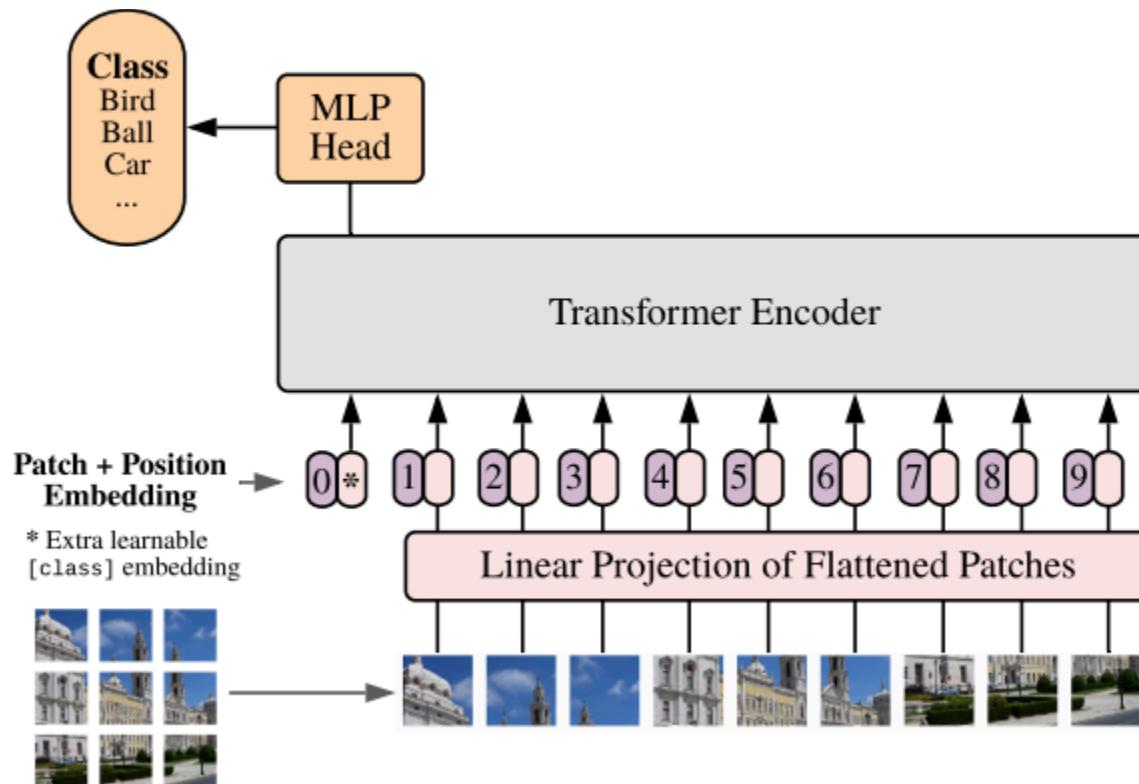
# CNN architectures

## ResNet (2015)



# Vision Transformers (2020)

- Not a CNN:  
Original implementation  
avoids convolutions
- Attention and positional  
encoding allows interactions  
across entire image, without  
needing to build a receptive  
field like CNNs





# *Computer Vision 2*

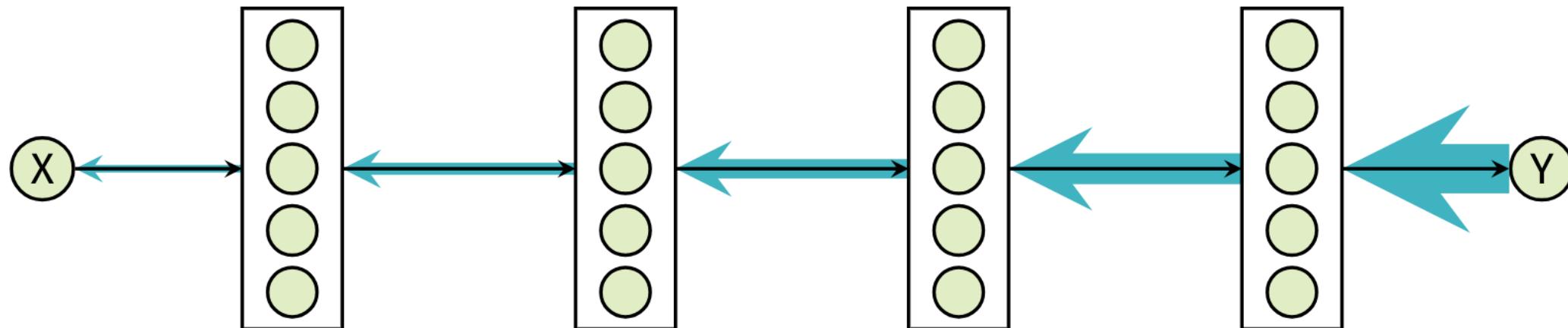
*COMPSCI 760  
2024 Semester 1*

*Olivier Graffeuille*

[ogra439@aucklanduni.ac.nz](mailto:ogra439@aucklanduni.ac.nz)

# Clarifications from Last Lecture

## Vanishing Gradients

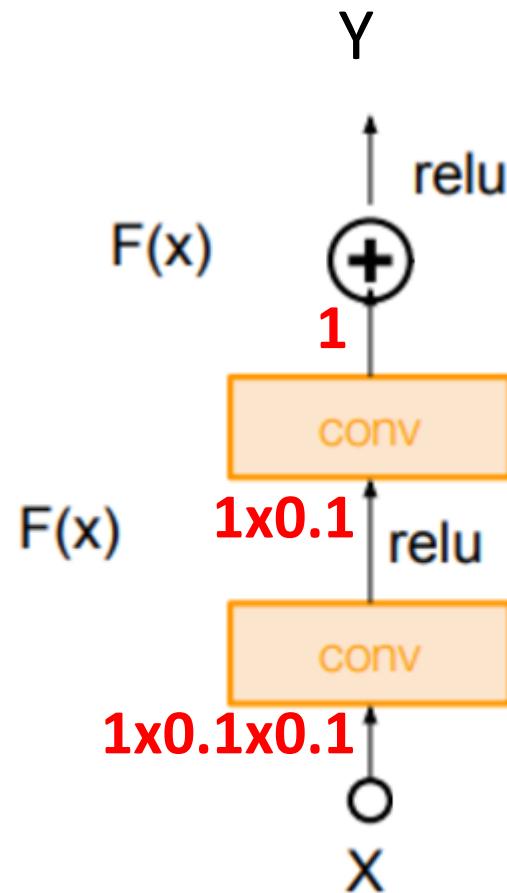


$$\delta^{(1)} = \varphi'(A^{(1)}) \circ W^{(0)T} \left( \varphi'(A^{(2)}) \circ W^{(1)T} \left( \varphi'(A^{(3)}) \circ W^{(2)T} \left( \varphi'(A^{(4)}) \circ (Y^{(3)} - Y) \right) \right) \right)$$

- A long sequence of mathematical operations is unstable

# Clarifications from Last Lecture

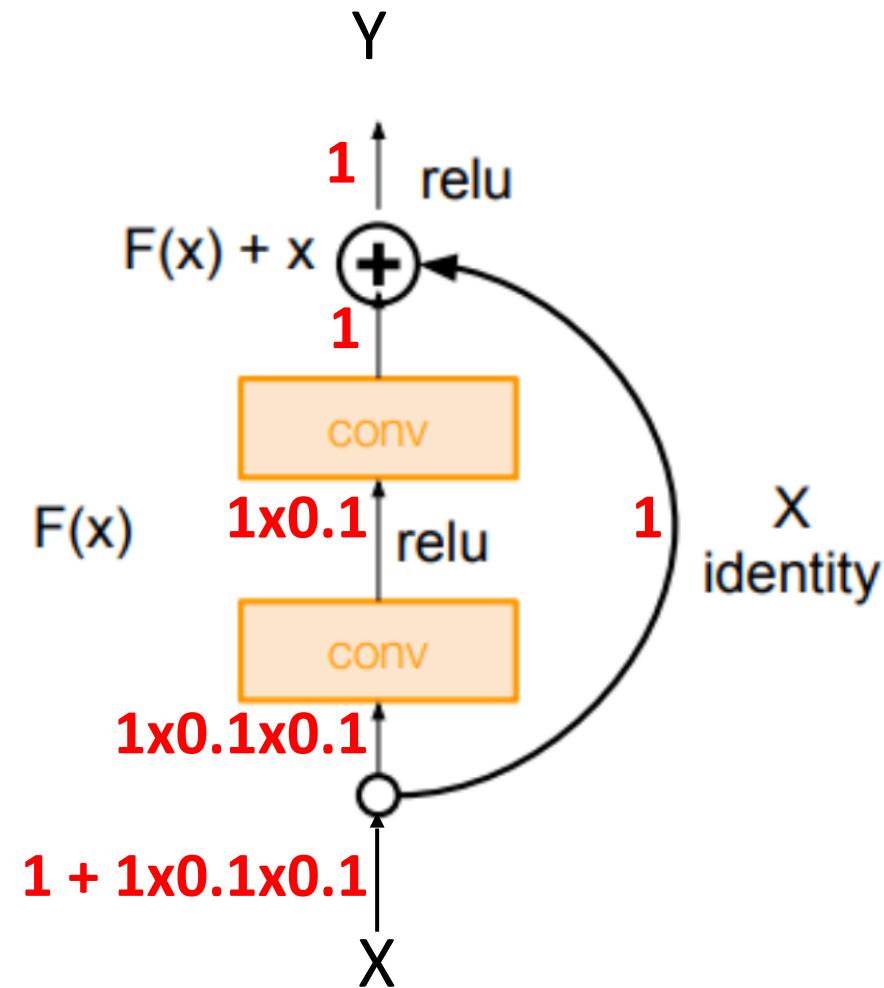
## Regular Deep Networks



# Clarifications from Last Lecture

## Residual Networks

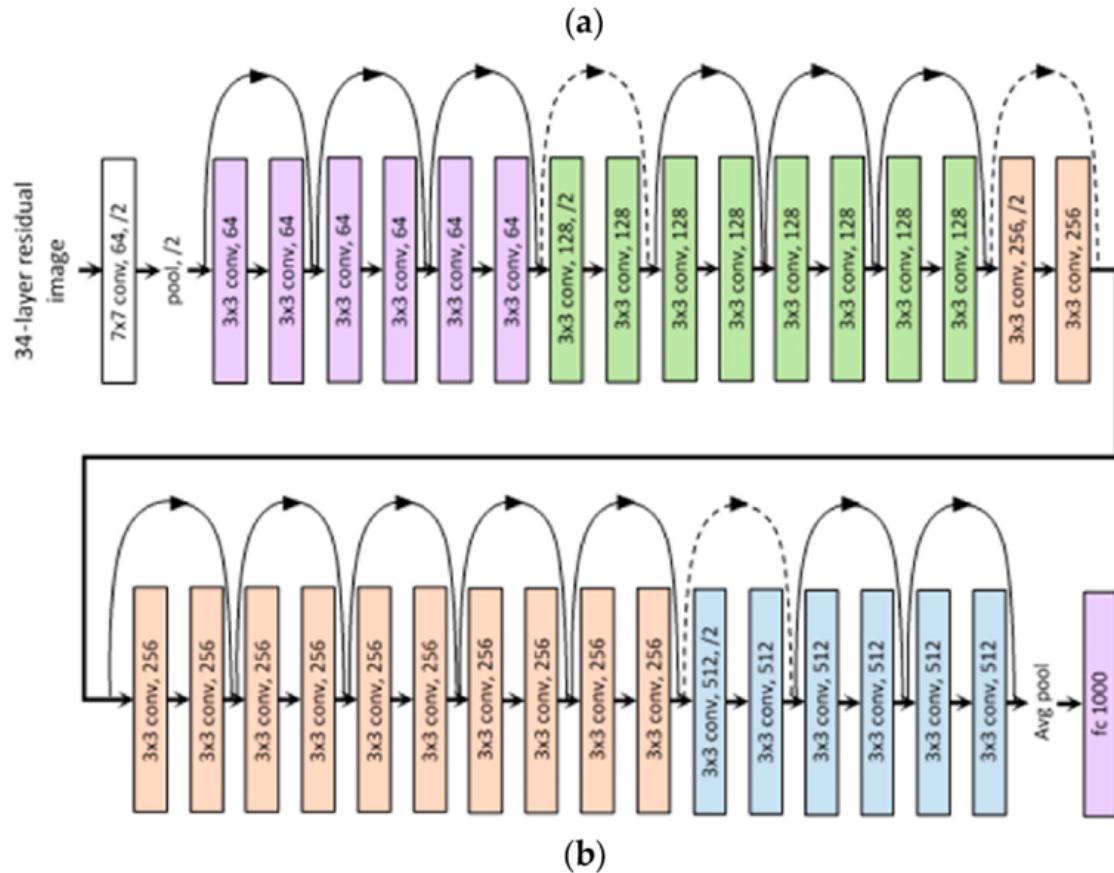
- Avoids vanishing gradients



# Clarifications from Last Lecture

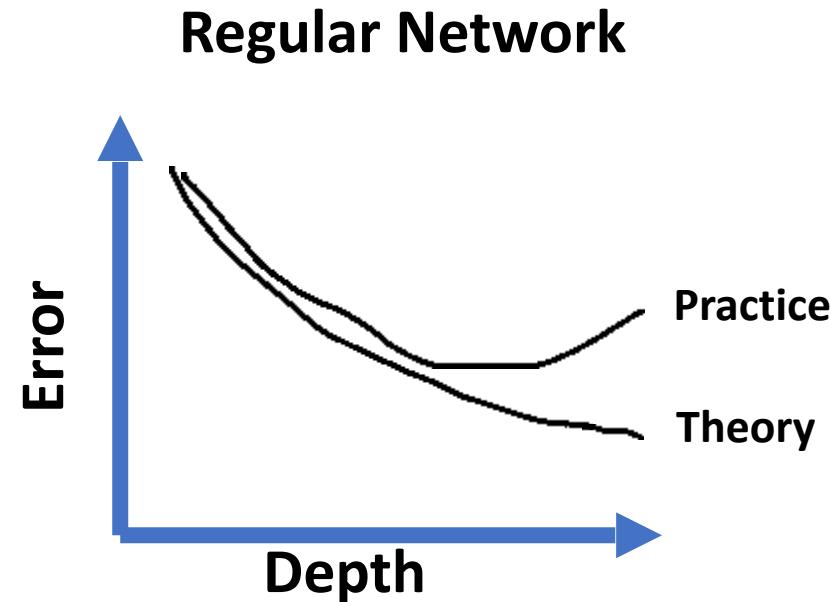
## ResNet (2015)

- Adding new layers, at worst, is like adding the identity function (no change)
  - These new layers can learn a function which is better than the identity



# Clarifications from Last Lecture

## Depth vs. Performance



- Regular CNNs have enough **capacity** to learn accurate models
- However, we aren't able to train them accurately
  - Vanishing gradients
  - Identity function start



## Clarifications from Last Lecture

# Vision Transformers vs. CNNs

### CNNs

- Requires large datasets
- Similar to human vision

### Vision transformers

- Requires **HUGE** datasets
- Very different from human vision
- Useful for vision/language interaction (foundation VLMs)

# Outline

- Computer vision tasks
  - Classification, localisation, landmark detection
- Object Detection
  - Sliding window method
  - Bounding box prediction (YOLO)
  - Anchor boxes
  - Putting it all together

# Outline

- ***Computer vision tasks***
  - ***Classification, localisation, landmark detection***
- Object Detection
  - Sliding window method
  - Bounding box prediction (YOLO)
  - Anchor boxes
  - Putting it all together

# Acknowledgement/Disclaimer

The following slides are taken from Andrew's Ng Coursera Convolutional Neural Networks lectures, part of the Deep Learning Specialization.

The source of the slides comes from [DeepLearning.AI](#)

## Copyright notice:

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see:

<https://creativecommons.org/licenses/by-sa/2.0/legalcode>



# Computer Vision Tasks

Query Image	Surface Normals	Eucl. Distance	Object Class.	Scene Class.	3D Curvature	Image Reshading	In-painting
			<p>Object Class. <b>Top 5 prediction:</b></p> <ul style="list-style-type: none"><li>sliding door</li><li>home theater, home theatre</li><li>studio couch, day bed</li><li>china cabinet, china closet</li><li>entertainment center</li></ul>	<p>Scene Class. <b>Top 2 prediction:</b></p> <ul style="list-style-type: none"><li>living room</li><li>television room</li></ul>			
Jigsaw puzzle	Colorization	2D Segm.	2.5D Segm.	Semantic Segm.	Cam. Pose (non-fixed)	Cam. Pose (fixed)	Autoencoding
Vanishing Points	2D Edges	3D Edges	2D Keypoints	3D Keypoints	Room Layout	Point Matching	Denoising



# Computer Vision Tasks

## Classification

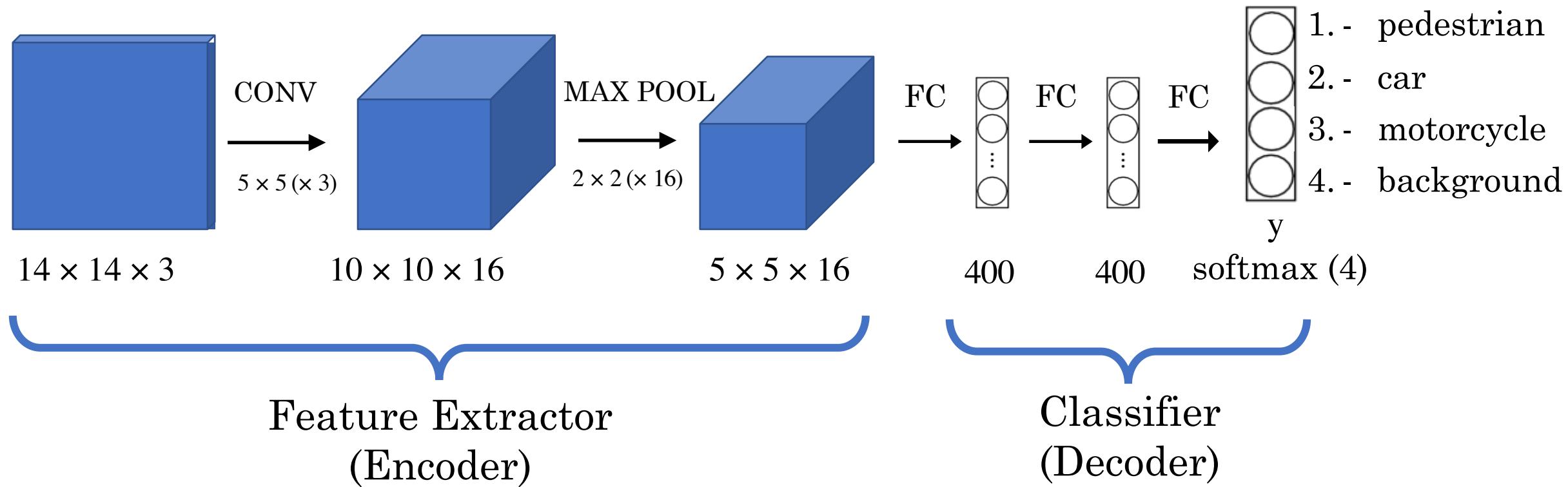
Image classification



« Car »

# Computer Vision Tasks

## Classification



# Computer Vision Tasks

## Classification + Localisation

Image classification



« Car »

Classification with  
localization

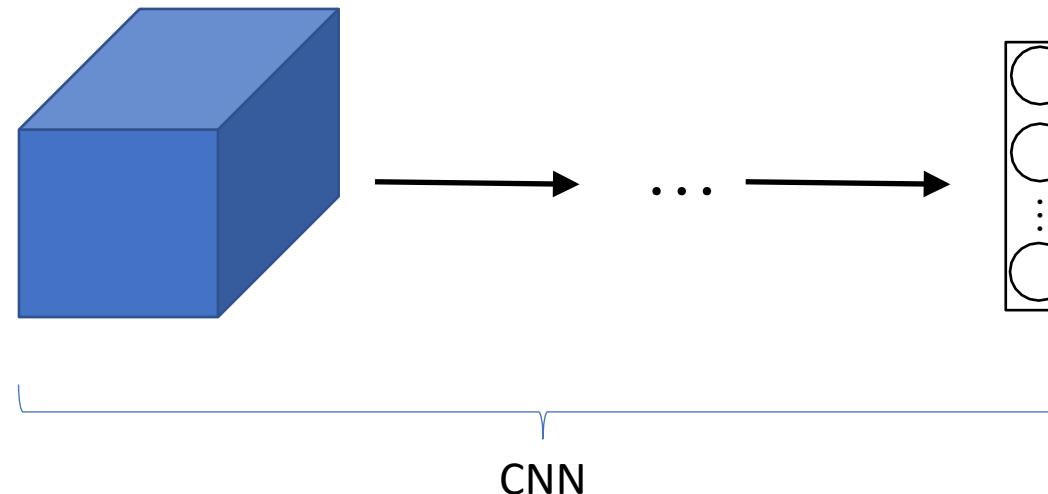
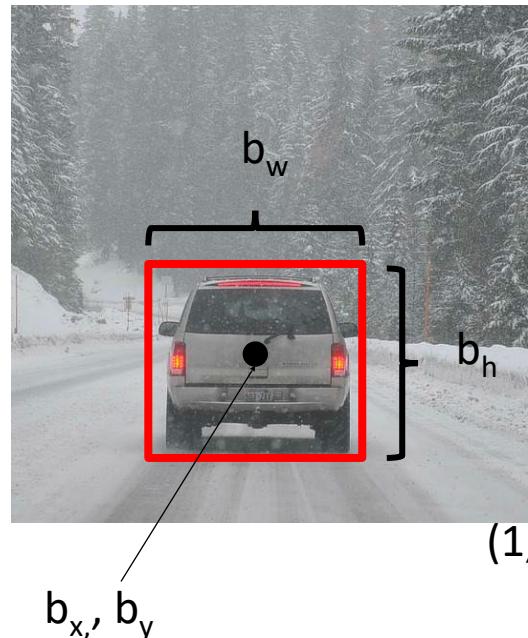


« Car » + localisation  
(e.g., ***bounding box***)

# Computer Vision Tasks

## Classification + Localisation

(0,0)



1. - pedestrian
2. - car
3. - motorcycle
4. - background

Softmax unit  
(e.g., 2 if there is  
a car)

Bounding box  
( $b_x, b_y, b_h, b_w$ )

# Computer Vision Tasks

Need to output  $b_x, b_y, b_h, b_w$ , class label (1-4)

- 1.- pedestrian
- 2.- car
- 3.- motorcycle
- 4.- background

Is there any object? (1 or 0 if none of the classes)

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \end{bmatrix}$$

If  $p_c = 1$   
Bounding box  
coordinates  
Class



1
0.5
0.7
0.3
0.4
0
1
0

0
?
?
?
?
?
?
?
?

Do not  
care



## Classification + Localisation

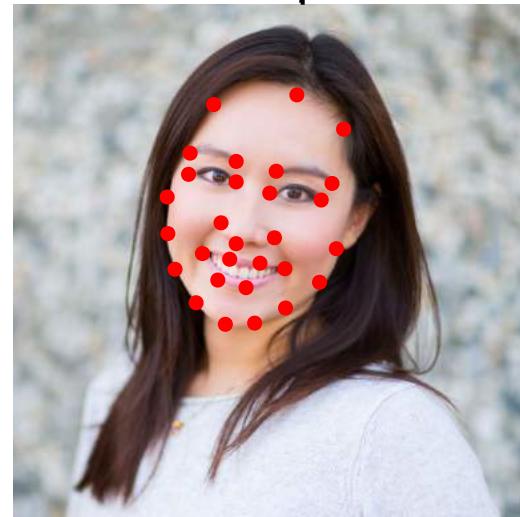
$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \quad L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2 & \text{if } p_c = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } p_c = 0 \end{cases}$$

# Computer Vision Tasks

## Landmark Detection



$b_x, b_y, b_h, b_w$



$l_{1x}, l_{1y}$

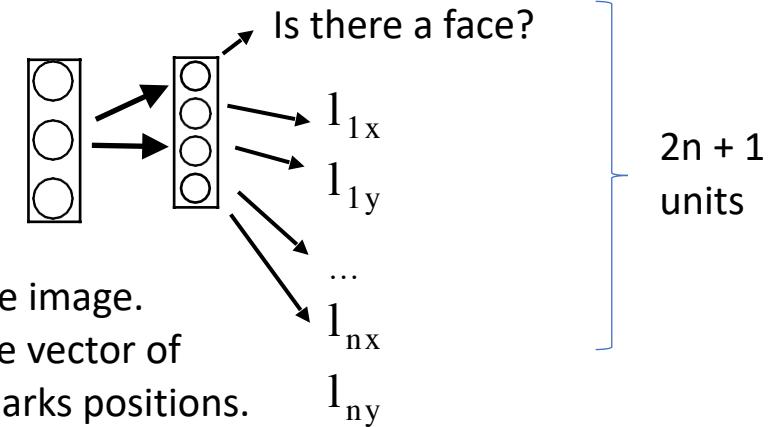
$l_{2x}, l_{2y}$

$l_{1x}, l_{1y}$

$l_{2x}, l_{2y}$

...

ConvNet



Can also be used for pose recognition.

# Computer Vision Tasks

## Object Detection

Image classification



« Car »

1 object

Classification with  
localization



« Car » + localisation (e.g.,  
bounding box)

Object  
Detection



Multiple objects + localisations  
(e.g., bounding boxes)

# Outline

- Computer vision tasks
  - Classification, localisation, landmark detection
- ***Object Detection***
  - ***Sliding window method***
  - Bounding box prediction (YOLO)
  - Anchor boxes
  - Putting it all together

# Object Detection

## Can we just use classification?

Training set:



x

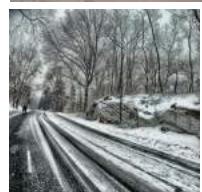


y

1



1



0



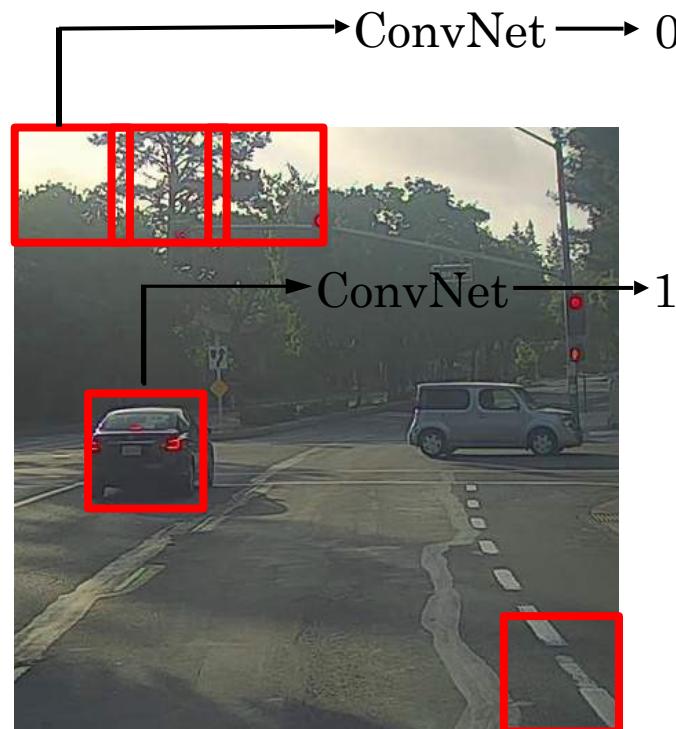
0



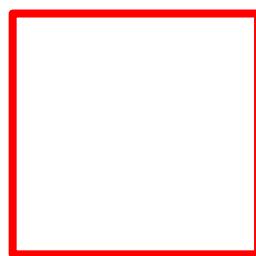
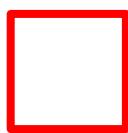
→ ConvNet → 1

# Object Detection

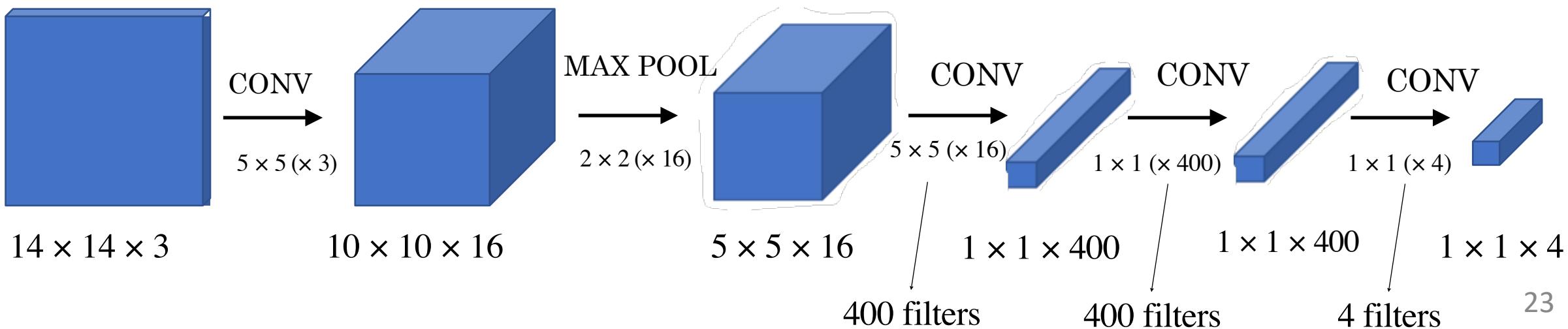
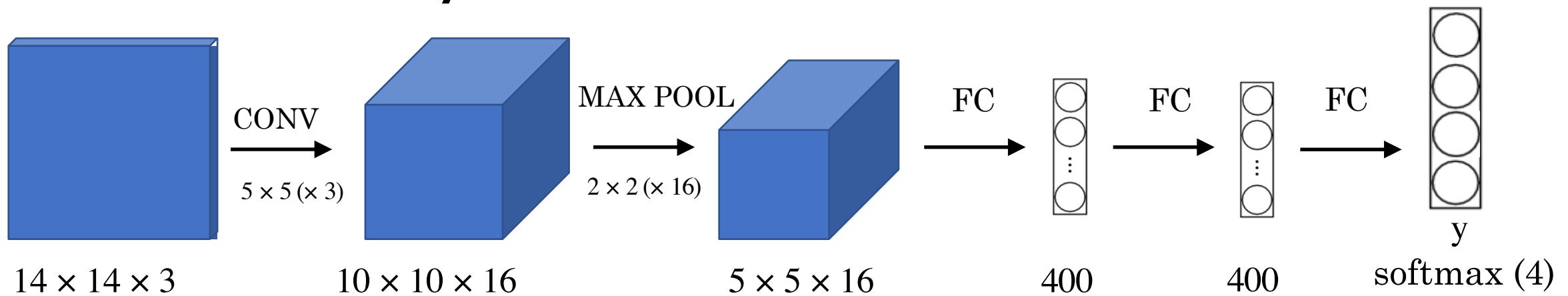
## Sliding Window Method



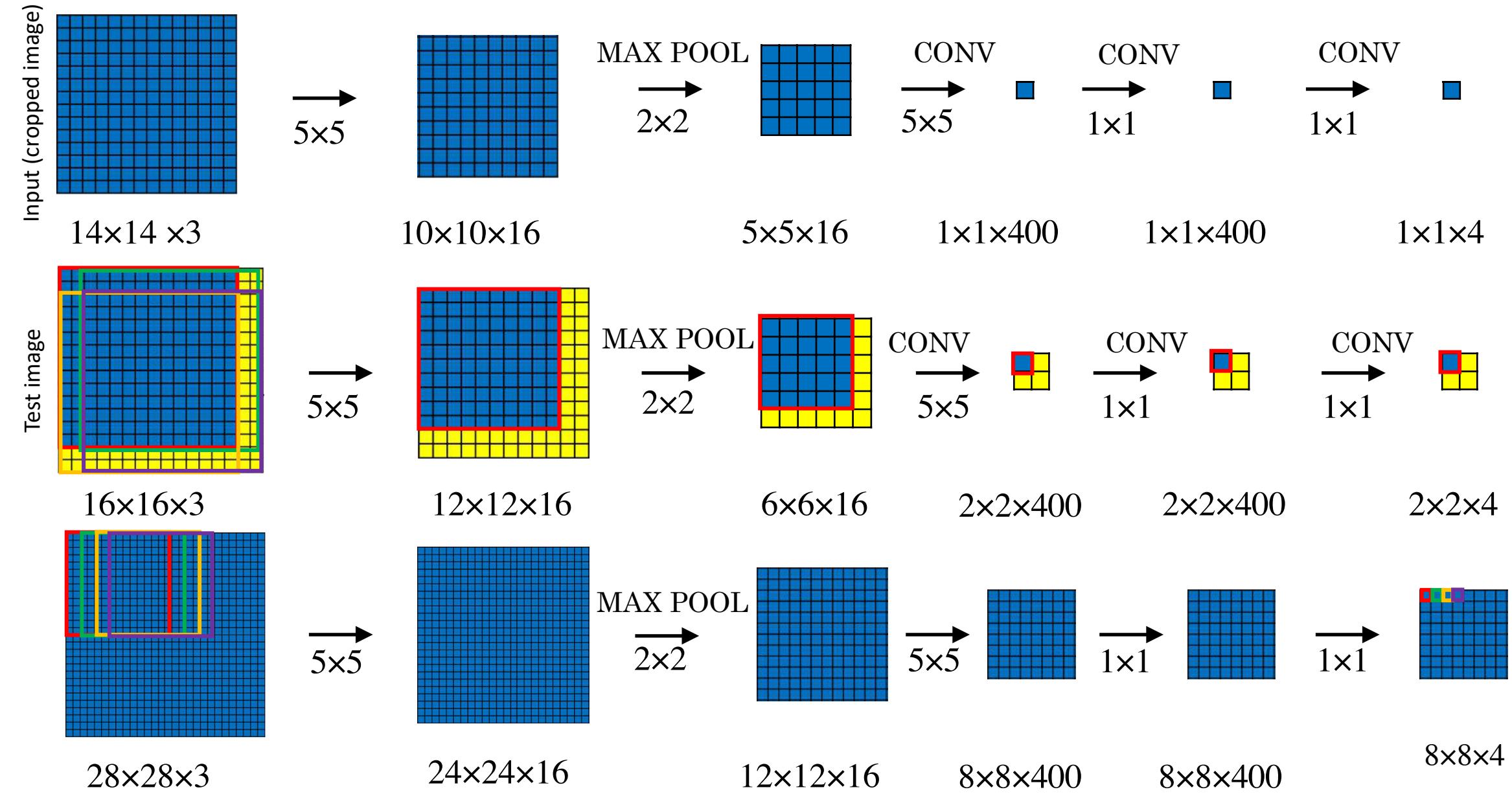
What are the problems  
with this approach?



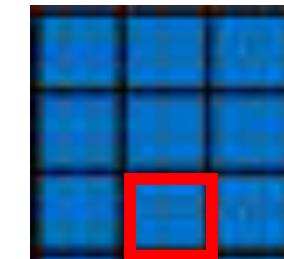
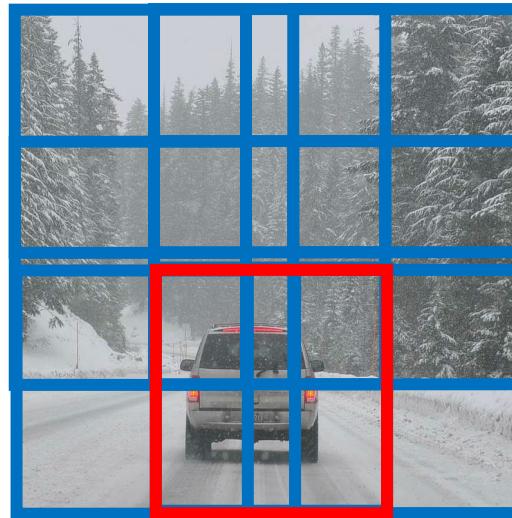
## Linear Layers as Convolutions



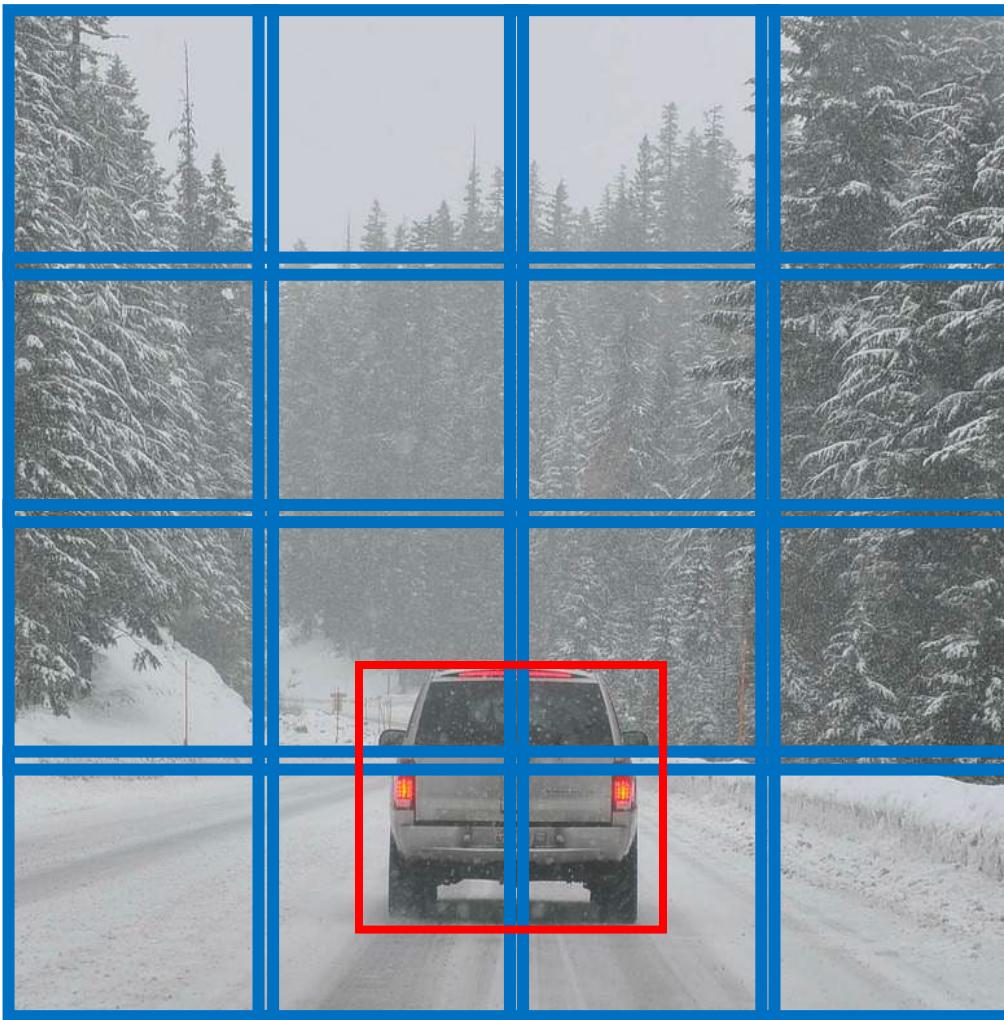
# Convolution implementation of sliding windows



# Convolution implementation of sliding windows



# Problem: Inaccurate estimation



# Outline

- Computer vision tasks
  - Classification, localisation, landmark detection
- Object Detection
  - Sliding window method
  - *Bounding box prediction (YOLO)*
  - Anchor boxes
  - Putting it all together

# YOLO: You Only Look Once

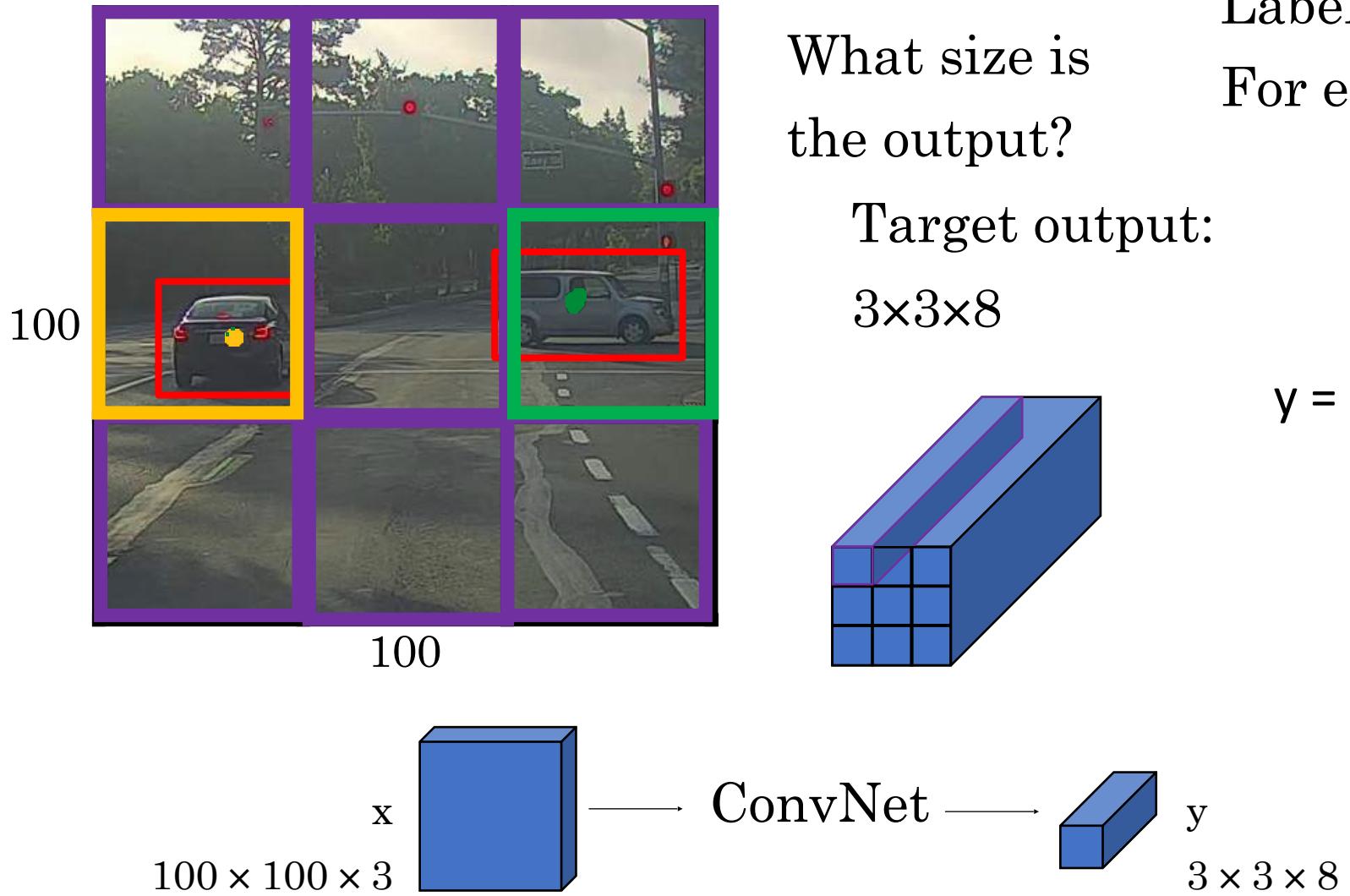
## YOLO: You Only Look Once (2015)

- This paper revolutionised object detection
- What do these two images have in common?

YOLO

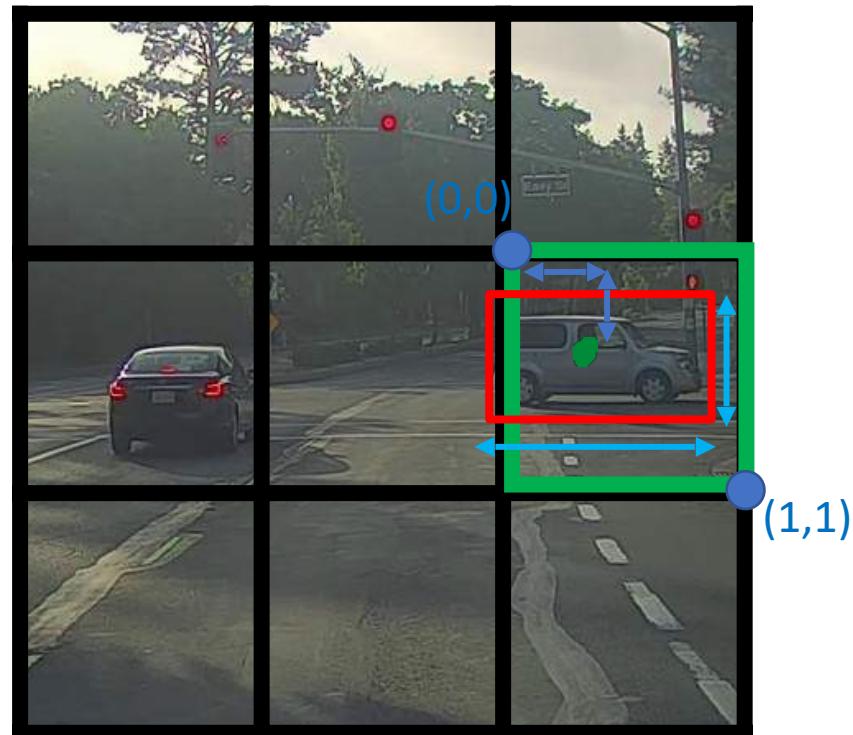


# YOLO: You Only Look Once



## YOLO: You Only Look Once

## Bounding Boxes for Object Detection



$$\begin{bmatrix} b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

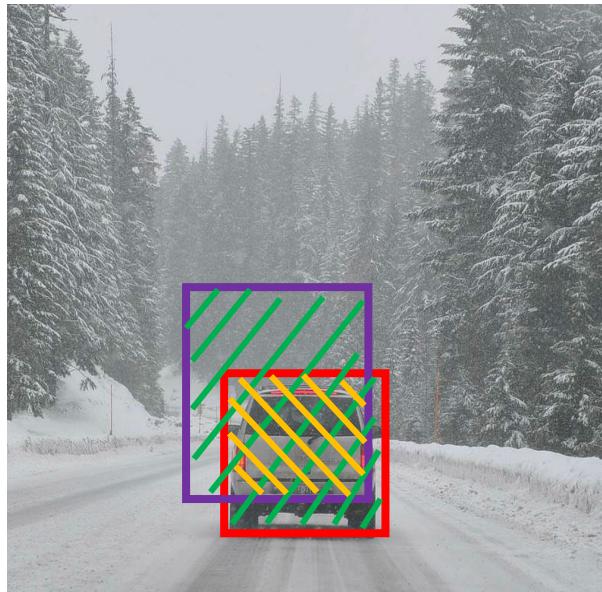
- x position of middle of bb to cell (0,0) = 0.4
- y position of middle of bb to cell (0,0) = 0.3
- height of bb as a fraction of cell size = 0.5
- width of bb as a fraction of cell size = 0.9

Between 0 and 1

Can be > 1

# YOLO: You Only Look Once

## Bounding Boxes: Evaluation



Intersection over union (IoU):

$$= \frac{\text{size of } \cap}{\text{size of } \cup}$$

“Correct” if  $\text{IoU} \geq 0.5$

More generally, IoU is a measure of the overlap between two bounding boxes.

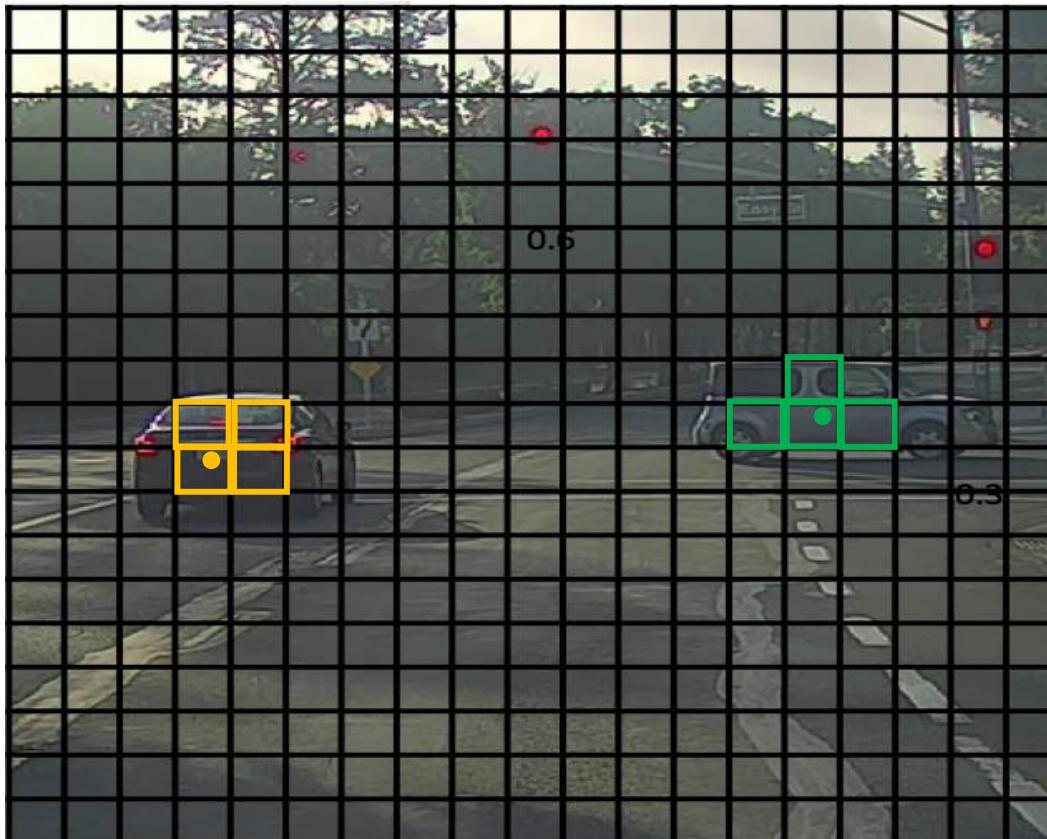
YOLO: You Only Look Once

# Bounding Boxes: Duplicates



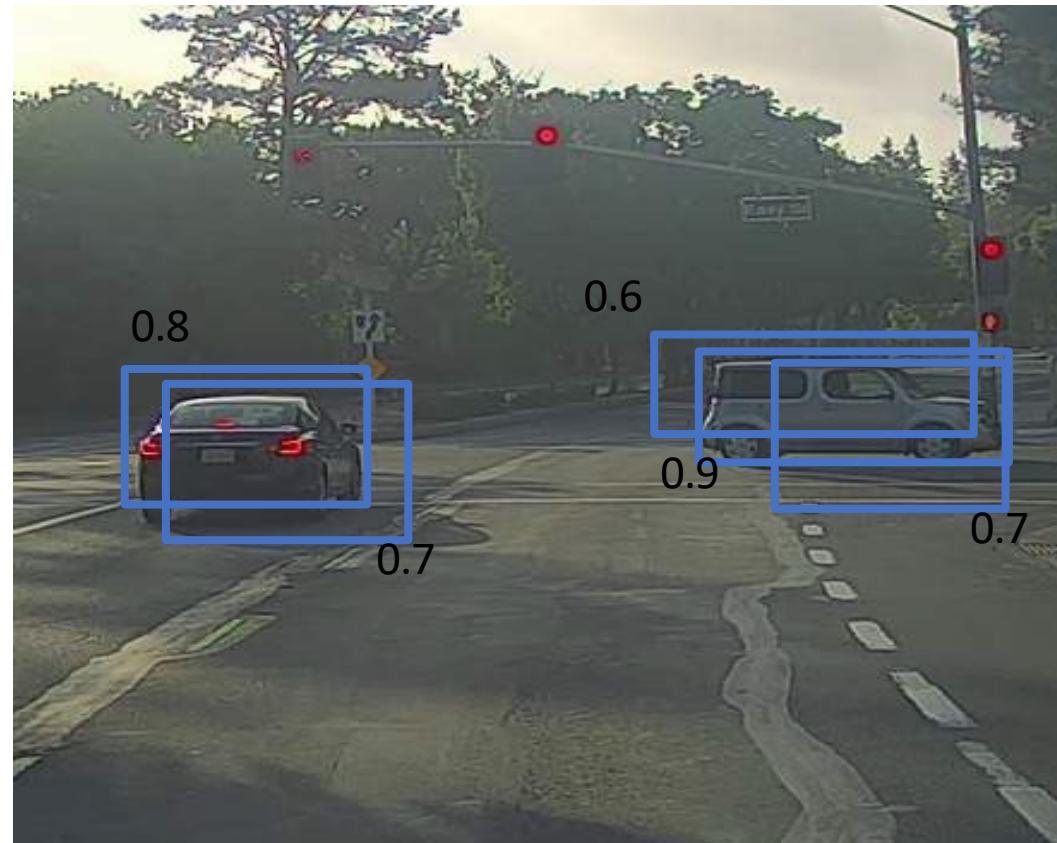
# YOLO: You Only Look Once

## Bounding Boxes: Duplicates



19 × 19

# YOLO: You Only Look Once



## Non-Max Suppression

Each output prediction is:

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \end{bmatrix}$$

Discard all boxes with  $p_c \leq 0.6^w$

While there are any remaining boxes:

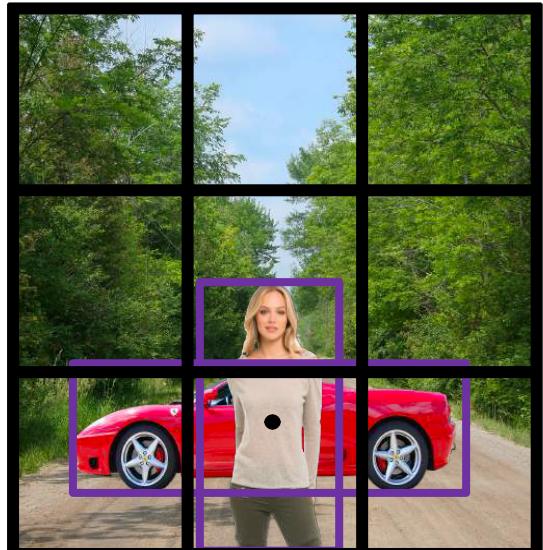
- Pick the box with the largest  $p_c$ .  
Output that as a prediction.
- Discard any remaining box with  $\text{IoU} \geq 0.5$  with the box output in the previous step.

# Outline

- Computer vision tasks
  - Classification, localisation, landmark detection
- Object Detection
  - Sliding window method
  - Bounding box prediction (YOLO)
  - *Anchor boxes*
  - Putting it all together

# YOLO: You Only Look Once

## Anchor Box Method



$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

$$y =$$

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

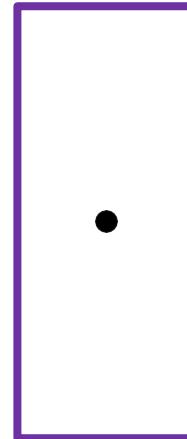
$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Anchor box 1  
(more similar to  
pedestrian's shape)

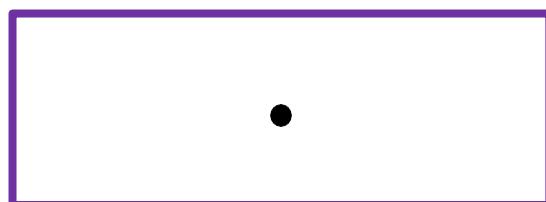
Anchor box 2  
(more similar to  
car's shape)

36

Anchor box 1:



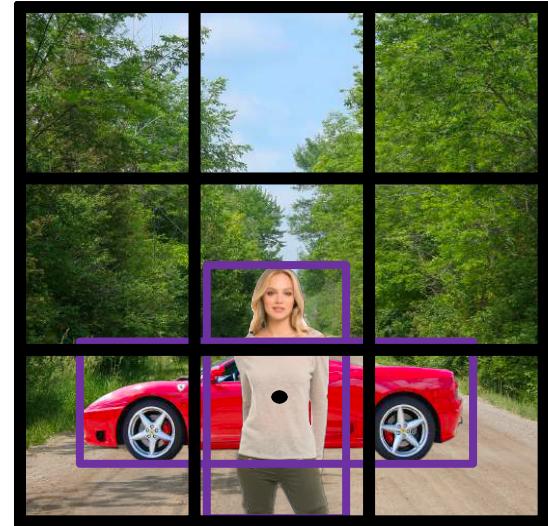
Anchor box 2:



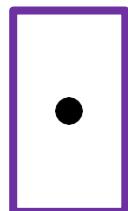
36

# YOLO: You Only Look Once

## Anchor Box: Example



Anchor box 1



Anchor box 2



$y =$

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 1 \\ 0 \\ 0 \\ 1 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ c_1 \\ 1 \end{bmatrix}$$

Car only?

$$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Anchor box 1

Anchor box 2



# YOLO: You Only Look Once

## Anchor Box Method

Previously:

Each object in training image is assigned to a grid cell that contains that object's midpoint.

Output (for each grid cell and 3 classes)

$y: 3 \times 3 \times 8$

With two anchor boxes:

Each object in training image is assigned to grid a cell that contains that object's midpoint, and anchor box for the grid cell with highest IoU.

Output (for each grid cell, 3 classes, and 2 anchor boxes)

$y: 3 \times 3 \times 16 \text{ (} 3 \times 3 \times 2 \times 8 \text{)}$



## Anchor Box Method

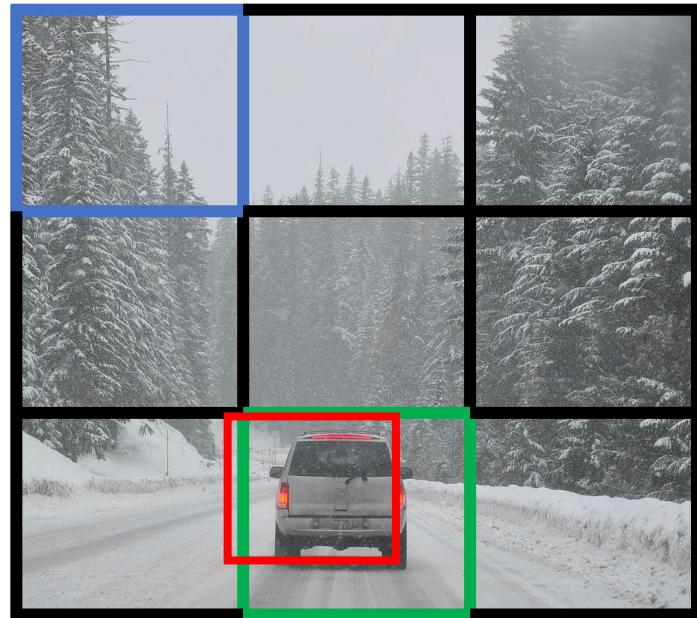
- Can detect multiple objects within a cell, but:
  - Only one object related to each anchor box per cell (usually enough)
- Allows the network to specialise into different object types

# Outline

- Computer vision tasks
  - Classification, localisation, landmark detection
- Object Detection
  - Sliding window method
  - Bounding box prediction (YOLO)
  - Anchor boxes
  - *Putting it all together*

# YOLO: You Only Look Once

## Training YOLO



$3 \times 3 \times 16$

y is  $3 \times 3 \times 2 \times 8$

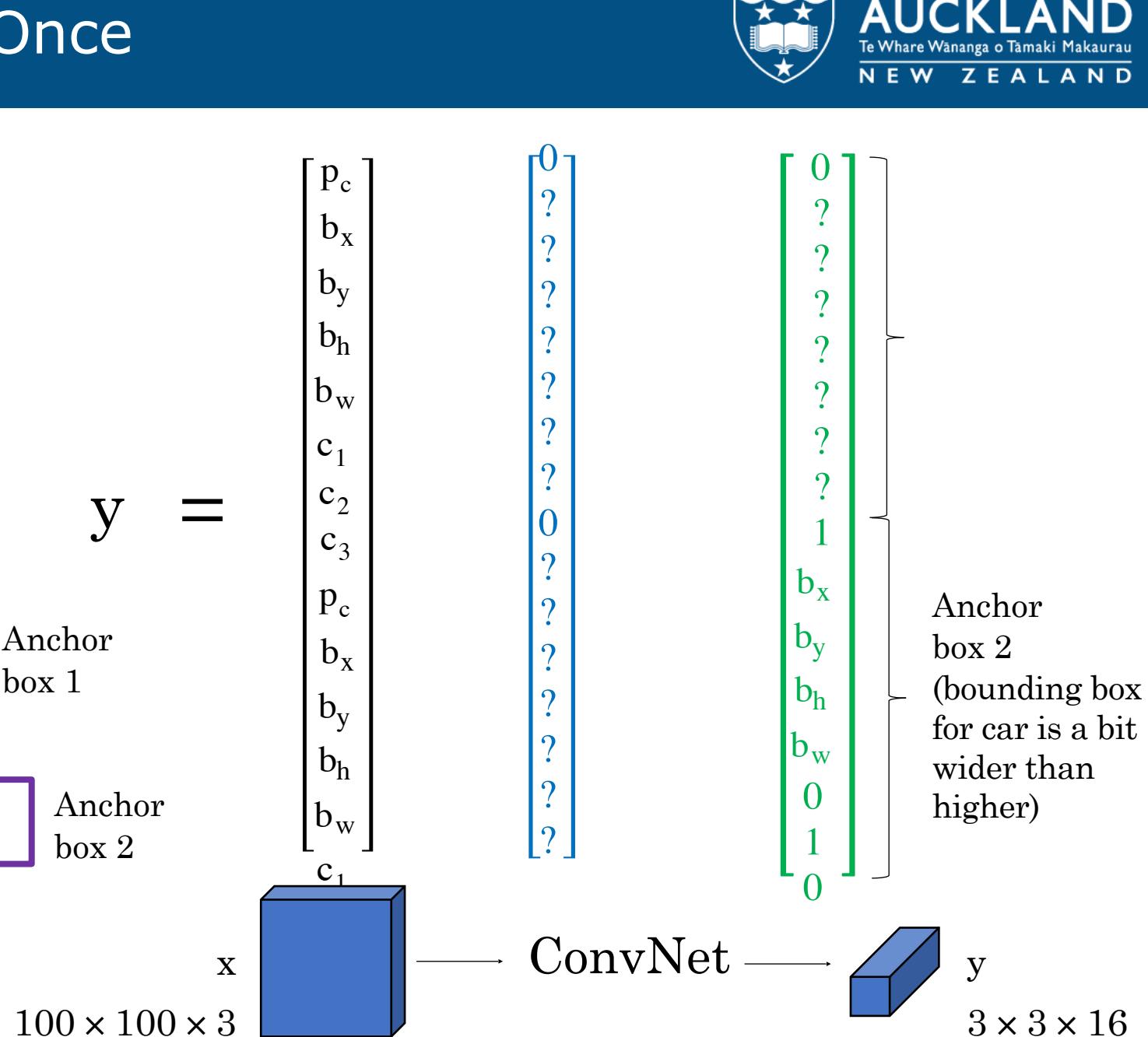
# anchor boxes

5 + 3 classes

y =

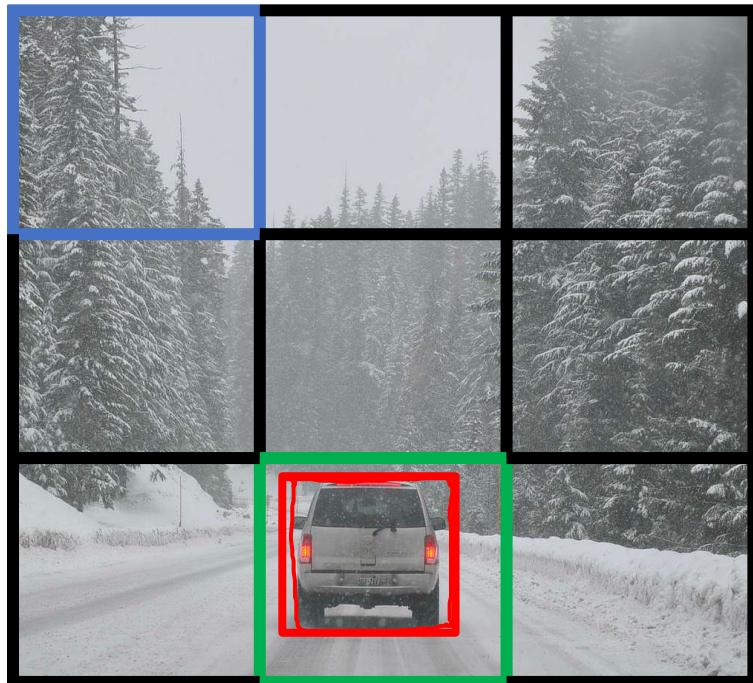
Anchor  
box 1

Anchor  
box 2

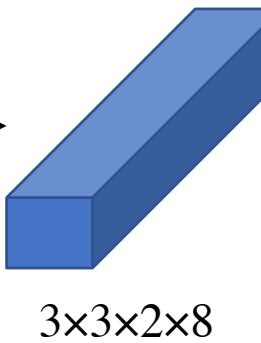


# YOLO: You Only Look Once

## Predicting with YOLO



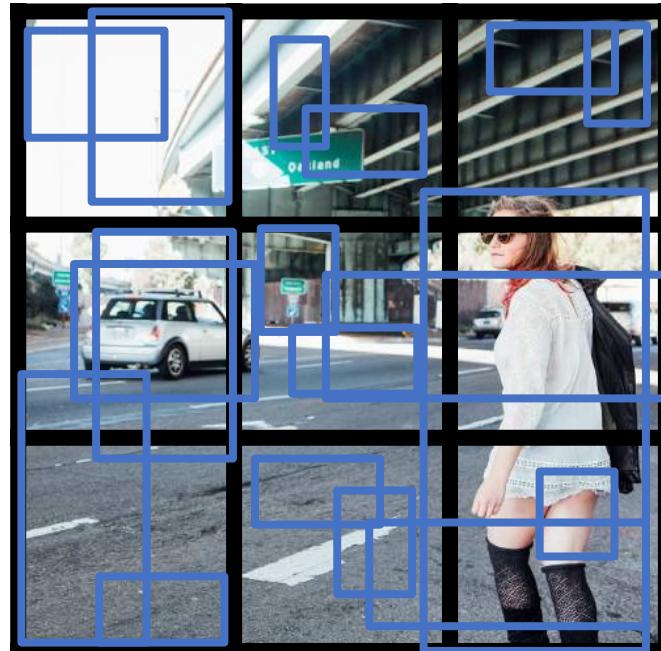
→ ... →



$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \begin{bmatrix} 0 \\ \dots \\ \dots \\ \dots \\ \dots \\ \dots \\ \dots \\ 0 \\ \dots \\ 0 \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \\ 42 \end{bmatrix}$$

# YOLO: You Only Look Once

## Predicting with YOLO: Non-Max Suppression



- For each grid cell, get 2 predicted bounding boxes.
- Get rid of low probability predictions.
- For each class (pedestrian, car, motorcycle) independently run non-max suppression to generate final predictions.



## YOLO Summary

- Object detection by estimating object position in a grid of cells
- Convolutional implementation for computational speed
- Anchor boxes to predict multiple objects per cell
- Non-max suppression to remove duplicate predictions



# YOLO: You Only Look Once

## YOLO Versions

- **YOLOv1:** Introduced a novel real-time object detection system that predicts bounding boxes and class probabilities directly from full images in one evaluation, significantly speeding up detection speeds.
- **YOLOv2:** Improved accuracy with anchor boxes, finer-grained features using a passthrough layer from earlier layers, and the ability to detect over 9000 object categories using a joint training method on both ImageNet and COCO datasets.
- **YOLOv3:** Enhanced detection at various scales by incorporating three different scales for prediction, and used a deeper network backbone (Darknet-53) for improved feature extraction.
- **YOLOv4:** Boosted accuracy and efficiency by integrating advanced techniques such as Mish activation, Cross-Stage Partial connections, and Self-Adversarial Training, while still maintaining real-time performance.
- **YOLOv8...**



# YOLO: You Only Look Once

## An interesting history

- The main author, Joseph Redmon, was a computer vision PhD student
- The first paper's name "YOLO" is a reference to a meme at the time
- The second paper was called "YOLO9000" and promised to predict "over 9000" classes, a reference to dragon ball Z



# YOLO: You Only Look Once

## An interesting history

- Joseph stopped working on computer vision for ethical reasons, e.g. military drones



**Joe Redmon**

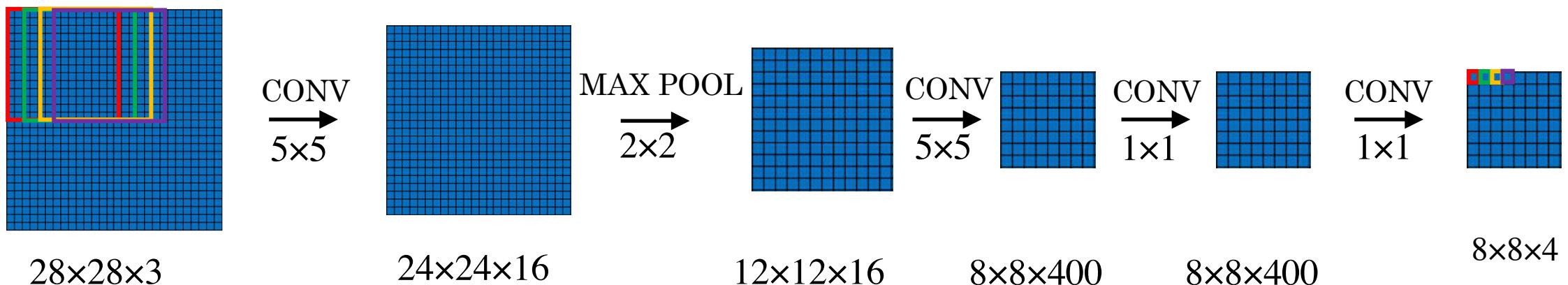
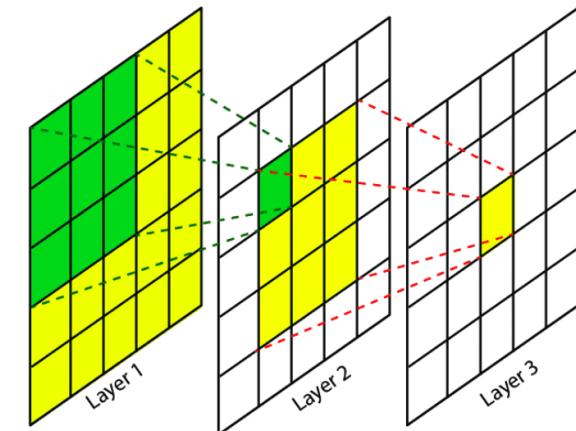
@pjreddie

I stopped doing CV research because I saw the impact my work was having. I loved the work but the military applications and privacy concerns eventually became impossible to ignore.[twitter.com/RogerGrosse/st...](https://twitter.com/RogerGrosse/status/1000000000000000000)

# Clarifications from Last Lecture

## Receptive Field

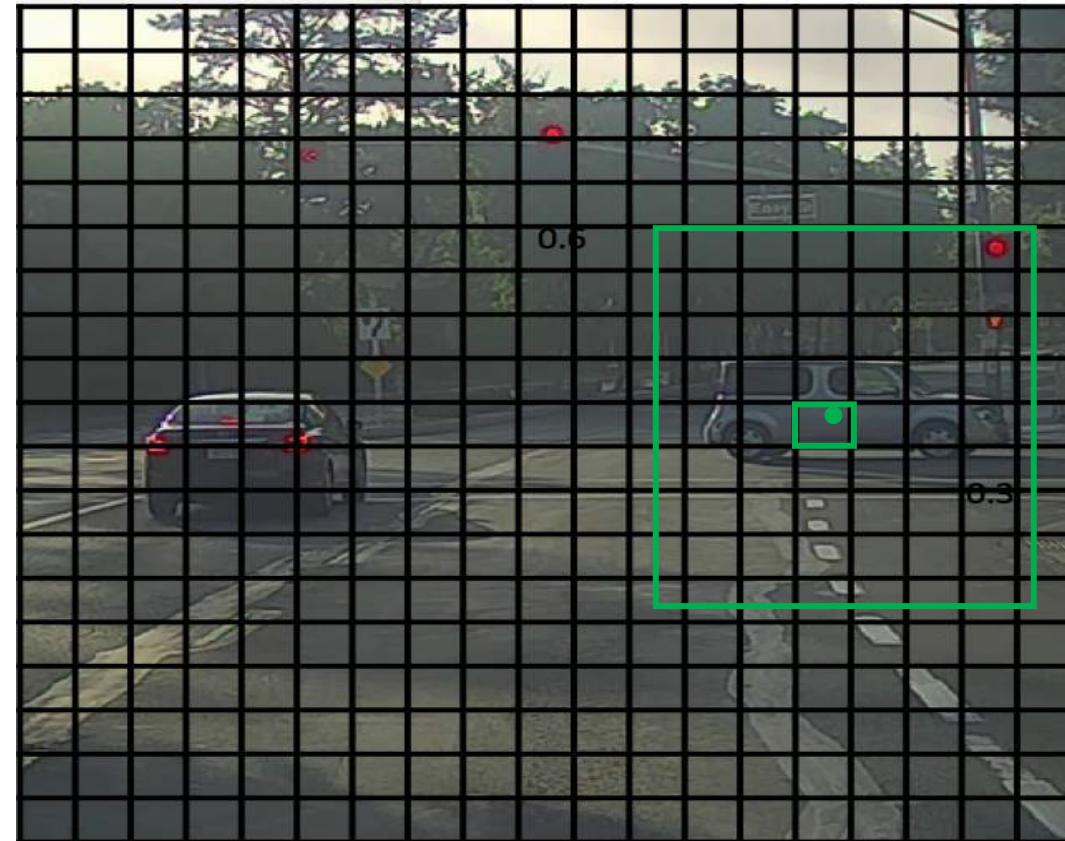
- The prediction within each feature/cell is the function of a larger patch of the input image aka
  - This is due to convolutional and pooling layer



# Clarifications from Last Lecture

## Cells in YOLO

- YOLO splits the image into cells
- Cells are a way of organising labels (object positions)
- The CNN can “see” more than the contents of a cell for each cell prediction

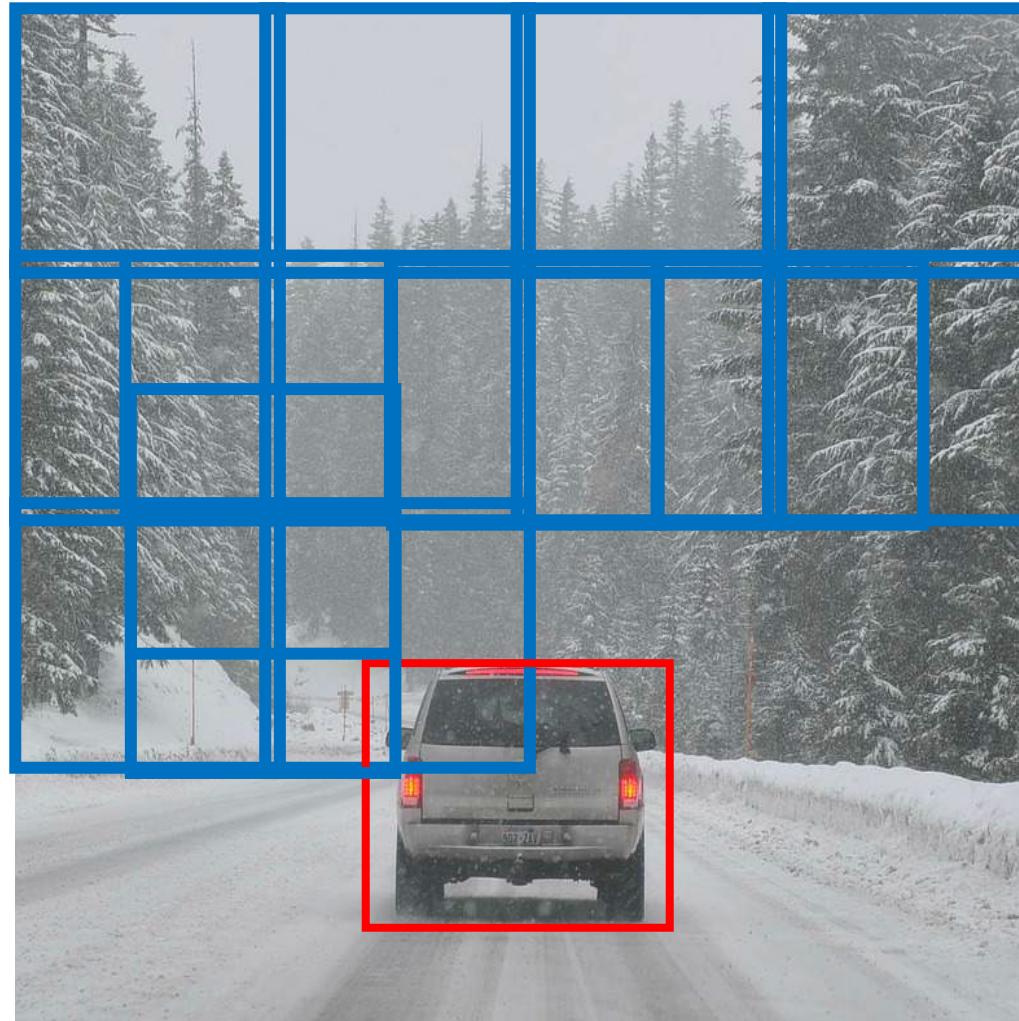


$19 \times 19$

# Clarifications from Last Lecture

## Cells in YOLO

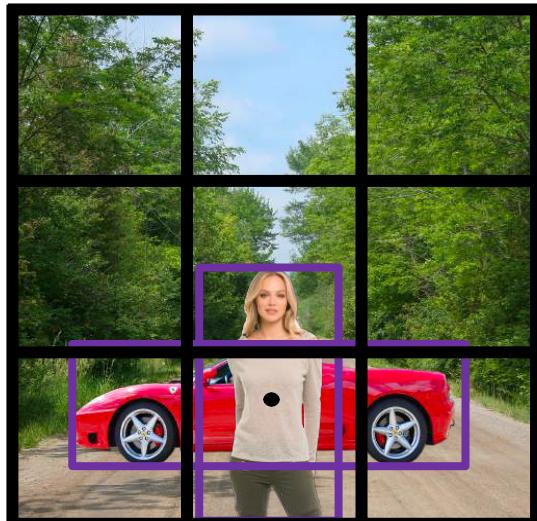
- Patches can overlap



# Clarifications from Last Lecture

# Anchor Boxes in YOLO

- They allow us to predict multiple objects per cell



$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \end{bmatrix}$

Which is correct?

or

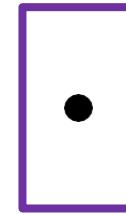
1	1
$b_x$	$b_x$
$b_y$	$b_y$
$b_h$	$b_h$
$b_w$	$b_w$
1	0
0	1
0	1
1	1
$b_x$	$b_x$
$b_y$	$b_y$
$b_h$	$b_h$
$b_w$	$b_w$
0	1
1	0
0	0

# Clarifications from Last Lecture

## Anchor Boxes in YOLO

- Anchor boxes are set before training
  - Manually or with clustering over data
- They serve as “reference point” shapes
  - Not directly linked to classes e.g. bounding box of a person will be different if they are standing vs. they are lying down

Anchor box 1



Anchor box 2

