

Deep Learning for Object Detection, Classification and Segmentation



Dr. Mohan Raj,

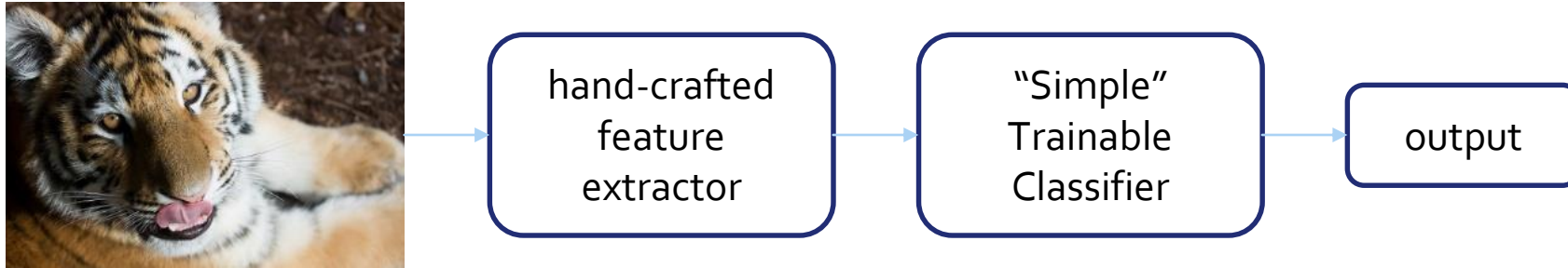
Data Scientist,
HCL Technologies,
Chennai.

Twitter - @mohanrajphd

Agenda

- Need for Deep Learning
- Image Classification
- Various CNN architecture for Image Classification
- Object Detection
- Various CNN architectures for Object Detection
- Image Segmentation
- Demo

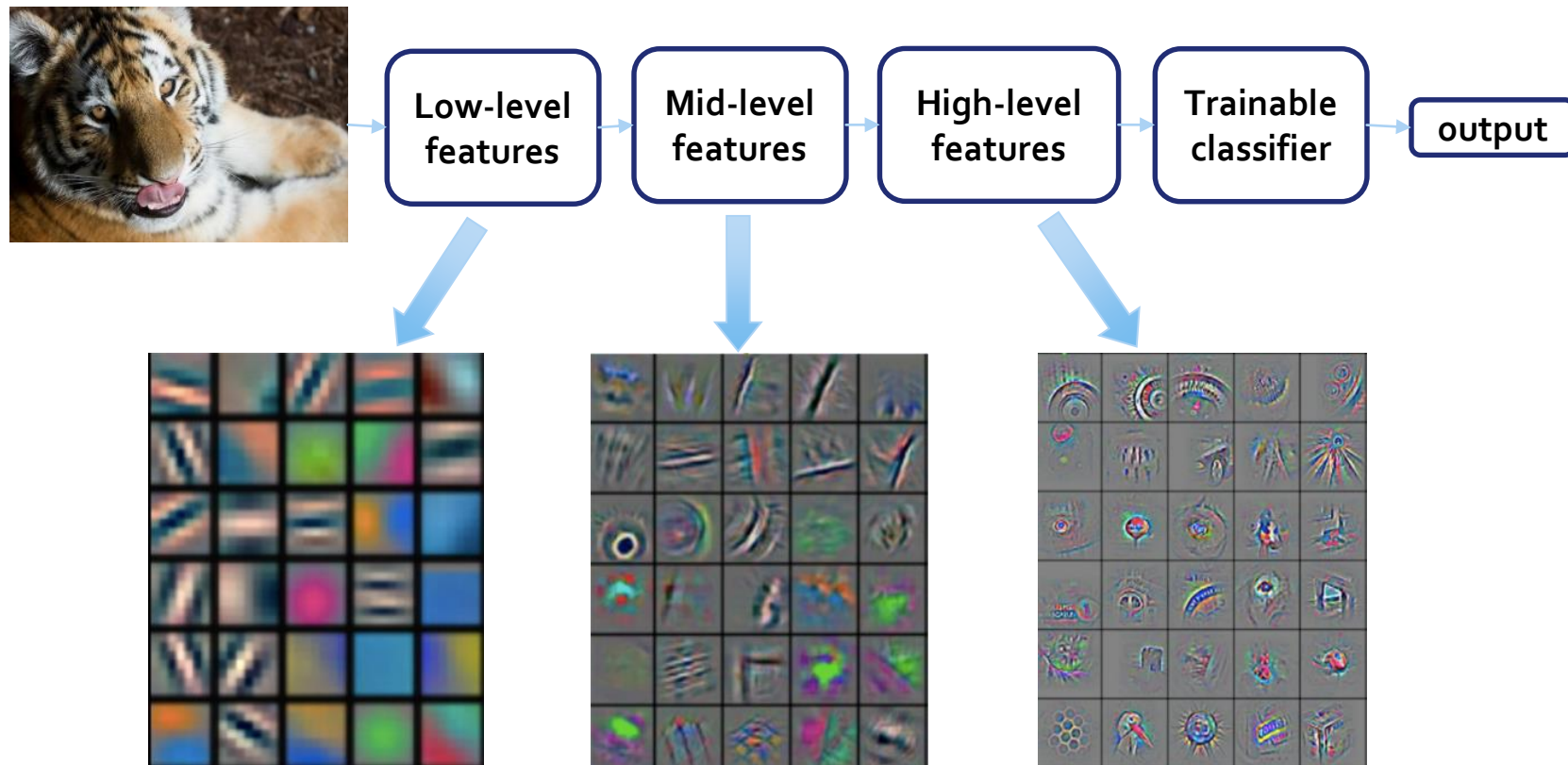
NEED FOR DEEP LEARNING



- Traditional pattern recognition models use hand-crafted features and relatively simple trainable classifier.
- This approach has the following limitations:
 - It is very tedious and costly to develop hand-crafted features
 - The hand-crafted features are usually highly dependents on one application, and cannot be transferred easily to other applications

DEEP LEARNING

- Deep learning (a.k.a. representation learning) seeks to learn rich hierarchical representations (i.e. features) automatically through multiple stage of feature learning process.



Feature visualization of convolutional net trained on ImageNet
(Zeiler and Fergus, 2013)

IMAGE CLASSIFICATION

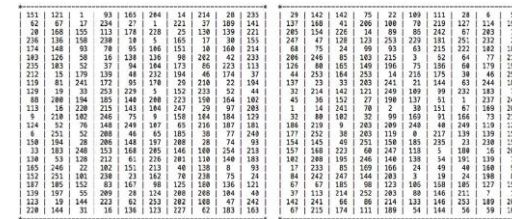
Image classification is the task of assigning a label to an image from a predefined set of categories.

- Let's assume the set of possible categories are:
categories = {cat, dog, panda}
- Classification algorithm assign multiple labels to the image via probabilities, such as
 - dog: 95%
 - cat: 4%
 - panda: 1%.



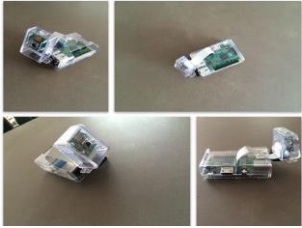
SEMANTIC GAP

- It should be fairly trivial for us to tell the difference between the two photos – there is clearly **a cat** on the left and **a dog** on the right. But all a computer sees is **two big matrices of pixels** (bottom).
- The **semantic gap** is the difference between how a human perceives the contents of an image versus how an image can be represented in a way a computer can understand the process.
- Visual examination of the two photos above can reveal the difference between the two species of an animal. But in reality, the computer has no idea there are animals in the image.
- We might describe the image as follows:
 - Spatial**: The sky is at the top of the image and the sand/ocean are at the bottom.
 - Color**: The sky is dark blue, the ocean water is a lighter blue than the sky, while the sand is tan.
 - Texture**: The sky has a relatively uniform pattern, while the sand is very coarse.



Feature extraction is the process of taking an input image, applying an algorithm, and obtaining a feature vector (i.e., a list of numbers) that quantifies our image.

REAL-TIME CHALLENGES



The object can be **oriented/rotated** in multiple dimensions with respect to how the object is photographed and captured.



The image on the left was photographed with standard **overhead lighting**. The image on the right was captured with very **little lighting**. We are still examining the same coffee cup — but based on the lighting conditions the cup looks dramatically different.



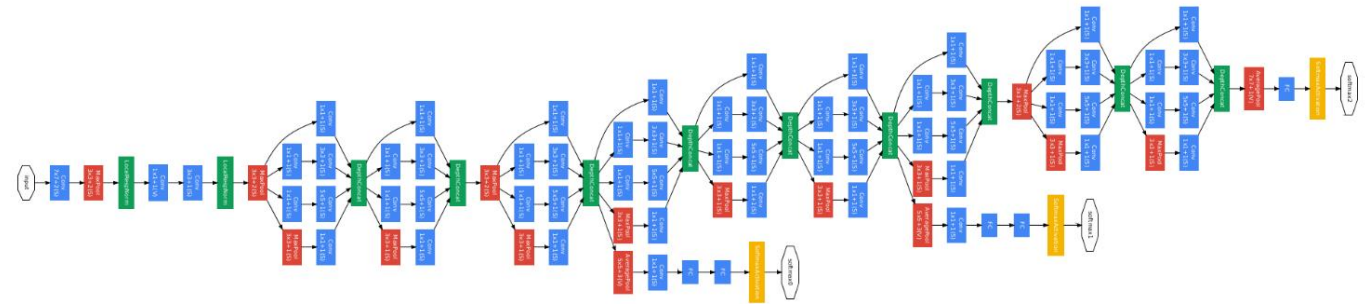
The same venti coffee will look dramatically different when it is **photographed up close** and when it is **captured farther way**.



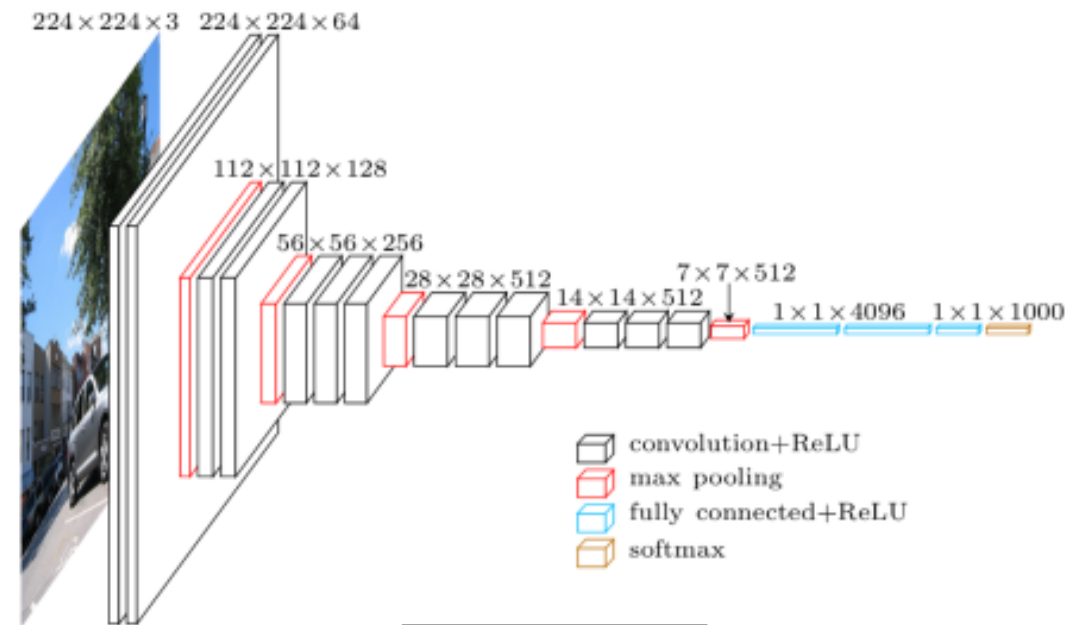
On the left we have a picture of a dog. The right we have a picture of the same dog, but notice how the dog is resting underneath the covers, **occluded from our view**.

STATE-OF-THE-ART CNN FOR IMAGE CLASSIFICATION

VGG16	VGG19	OverFeat	GoogLeNet	ResNet50
image	image	image	image	image
conv-64	conv-64	conv-96	conv-64	conv-64
conv-64	conv-64	maxpool	maxpool	maxpool
maxpool	maxpool			
		conv-256	conv-192	conv2_x
conv-128	conv-128	maxpool	maxpool	conv-64
conv-128	conv-128			conv-64 x 3
maxpool	maxpool	conv-512	inception-256	conv-256
		conv-1024	inception-480	
conv-256	conv-256	conv-1024	maxpool	conv3_x
conv-256	conv-256	maxpool		conv-128
conv-256	conv-256		inception-512	conv-128 x 4
maxpool	conv-256	FC-3072	inception-512	conv-512
	maxpool	FC-4096	inception-512	
conv-512	conv-512	FC-1000	inception-528	conv4_x
conv-512	conv-512	softmax	inception-832	conv-256
	conv-512		maxpool	conv-256 x 6
conv-512	conv-512		inception-832	conv-1024
maxpool	conv-512		inception-1024	
	maxpool		avgpool	conv5_x
conv-512	conv-512		dropout-1024	conv-512
conv-512	conv-512		FC-1000	conv-512 x 3
conv-512	conv-512		softmax	conv-2048
maxpool	conv-512			
	maxpool			avgpool
FC-4096	FC-4096			FC-1000
FC-4096	FC-4096			softmax
FC-1000	FC-1000			
softmax	softmax			



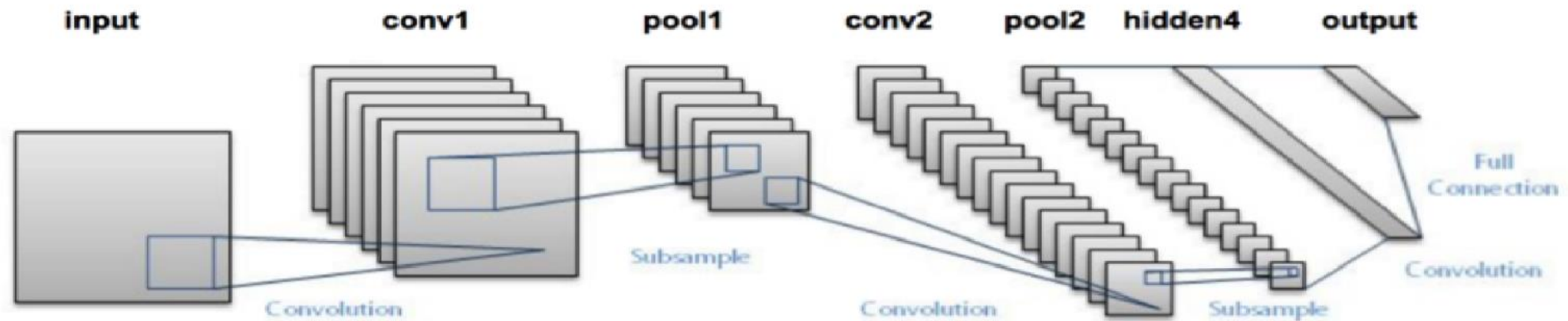
Google Inception-V3



VGG-16

LENET - RECOGNIZING HANDWRITTEN DIGITS

The **LeNet architecture** is a seminal work in the deep learning community, first introduced by **LeCun et al.** in their 1998 paper, Gradient-Based Learning Applied to Document Recognition.



WHY OBJECT DETECTION?

Image classification is the task of **assigning a label** to an image from a predefined set of categories.



To solve this problem, we can train a **multi-label classifier** which will predict the probabilities of both the classes (dog as well as cat).

However, we still don't know the **location** of cat or dog in the image.

OBJECT DETECTION

Predicting the **location of the object** in an image or video is called **object detection** or **localization**.

Classification



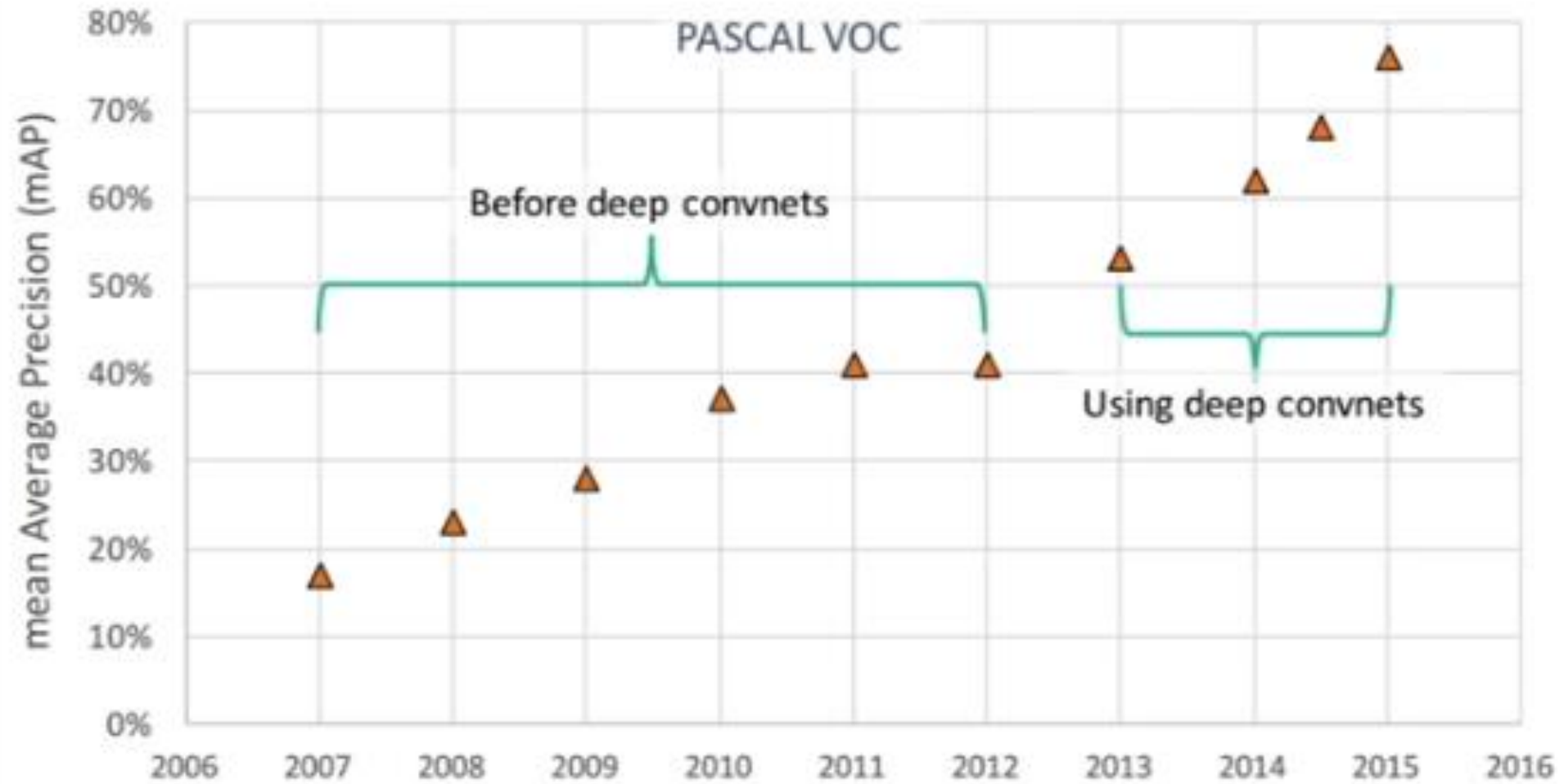
CAT

**Classification
+ Localization**



CAT

DEEP LEARNING FOR OBJECT DETECTION



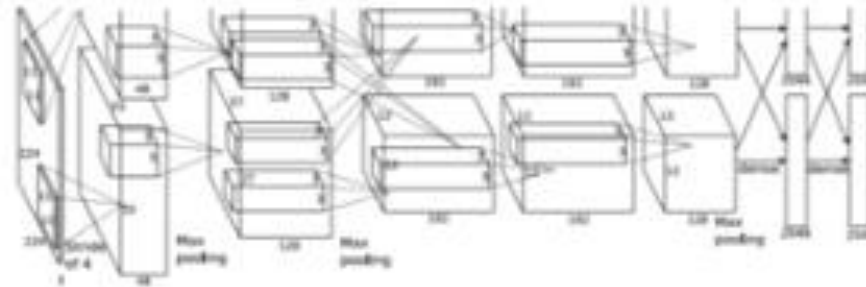
SLIDING WINDOWS - HOW IT WORKS?

- **Sliding window** is rectangular region of fixed width and height that “**slides**” across the image, from left-to-right and top-to-bottom.
- A sliding window slides from left-to-right and top-to-bottom across an input image taking **N pixel steps at a time**.
- The ROI at each step of the sliding window is **extracted and passed** into the feature extraction/object detection pipeline.



OBJECT DETECTION USING CNN

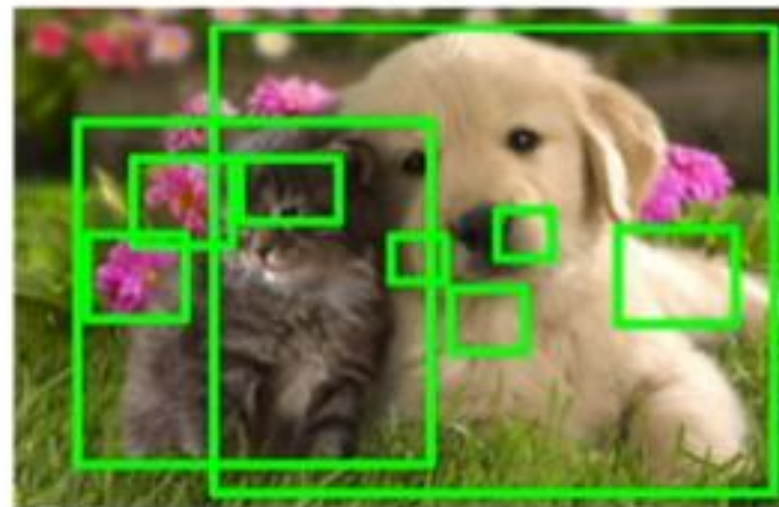
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? NO
Background? YES

REGION PROPOSALS

- Find “**blobby**” image regions that are **likely to contain objects**.
- Relatively fast to run; e.g. **Selective search** gives 1000 region proposals in a few seconds on CPU.



HISTORY OF OBJECT DETECTION

While there are many object detection methods in the computer vision literature, **two stand out amongst the others**:

- HOG + Linear SVM (Histogram of Oriented Gradients + Linear Support Vector Machine)
- Haar cascades

HISTOGRAM OF ORIENTED GRADIENTS (HOG)

- Normalizing the image prior to description.
- Computing gradients in both the x & y directions.
- Obtaining weighted votes in spatial & orientation cells.
- Contrast normalizing in the overlapping spatial cells.
- Collect all HOGs to form the final feature vector.

GRADIENT COMPUTATION



$$G_x = I \star D_x$$



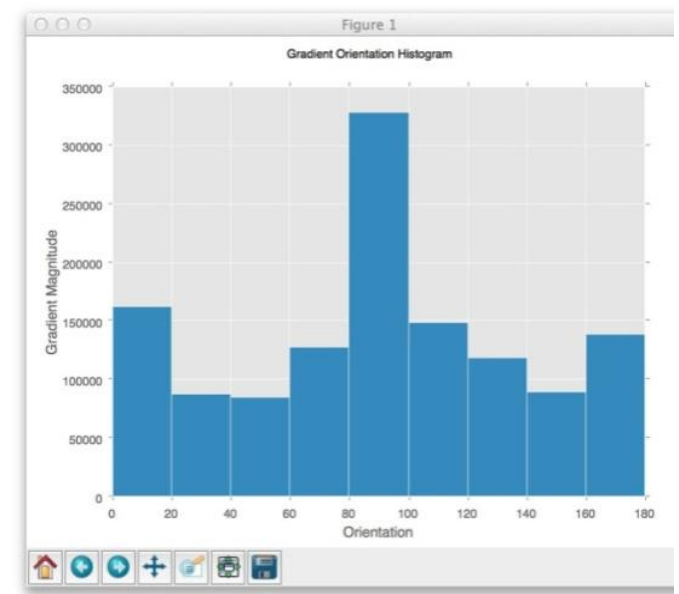
$$G_y = I \star D_y$$



$$|G| = \sqrt{G_x^2 + G_y^2}$$

$$\theta = \tan^{-1}\left(\frac{G_y}{G_x}\right)$$

WEIGHTED VOTES IN EACH CELL



CONTRAST NORMALIZATION OVER BLOCK

Block 1

Cell #1	Cell #2	Cell #3
Cell #4	Cell #5	Cell #6
Cell #7	Cell #8	Cell #8

HOG FEATURE VECTOR



HAAR CASCADES

- One of the most famous object detectors, **Rapid Object Detection** using a **Boosted Cascade of Simple Features**, by Viola and Jones (2004).
- **Pre-trained Haar cascades** are distributed with the OpenCV library, and are arguably the most used models for **face detection**.
- While Haar cascades are fast, they -
 - a. Tend to have a high false-positive detection rate.
 - b. Can miss objects entirely based on the parameters supplied to the cascade.

REGION-BASED CONVOLUTIONAL NEURAL NETWORK (R-CNN)

- **R-CNN** solves this problem by using an object proposal algorithm called **Selective Search** which reduces the number of bounding boxes that are fed to the classifier close to 2000 region proposals.
- Selective search uses **local cues** like texture, intensity, color to generate all the possible locations of the object.
- There are **3 important parts** in R-CNN :
 - a. Run Selective Search to generate probable objects.
 - b. Feed these patches to CNN, followed by SVM to predict the class of each patch.
 - c. Optimize patches by training bounding box regression separately.

R-CNN ARCHITECTURE

R-CNN: *Regions with CNN features*

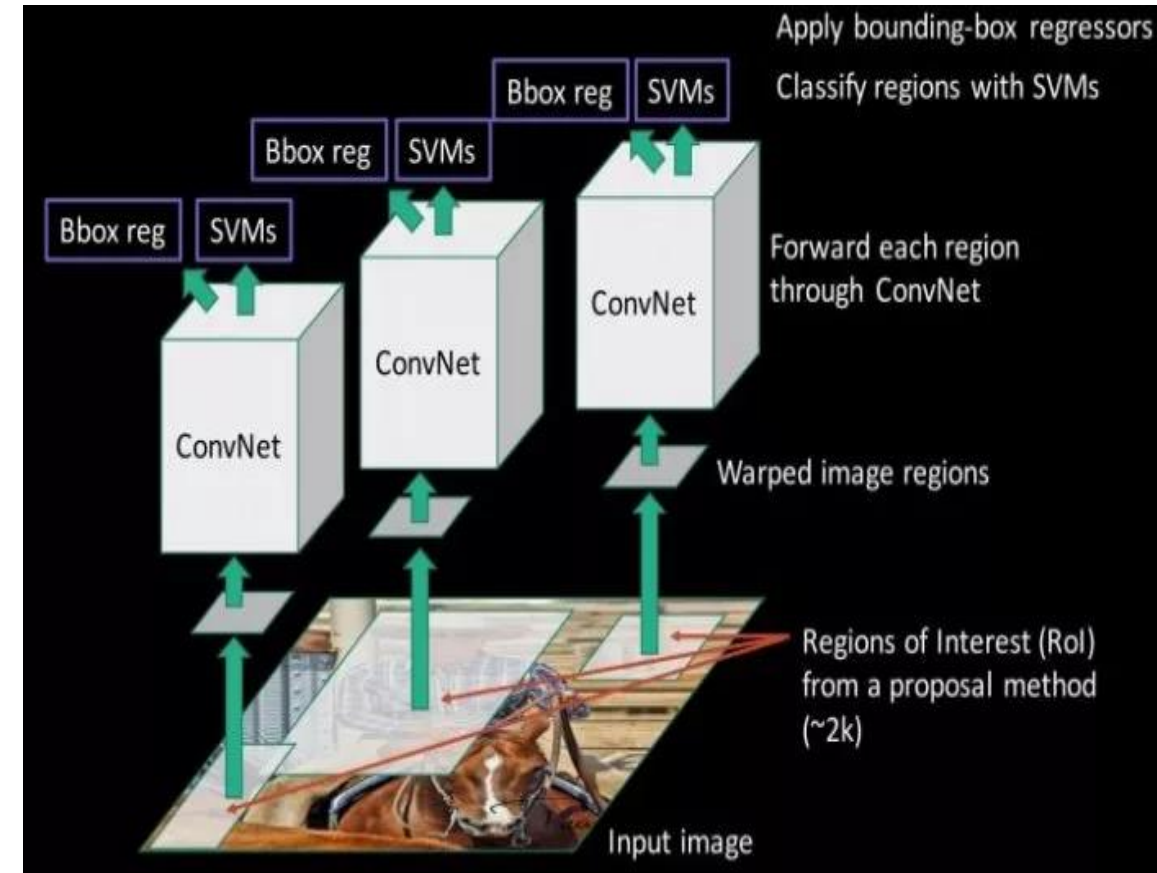
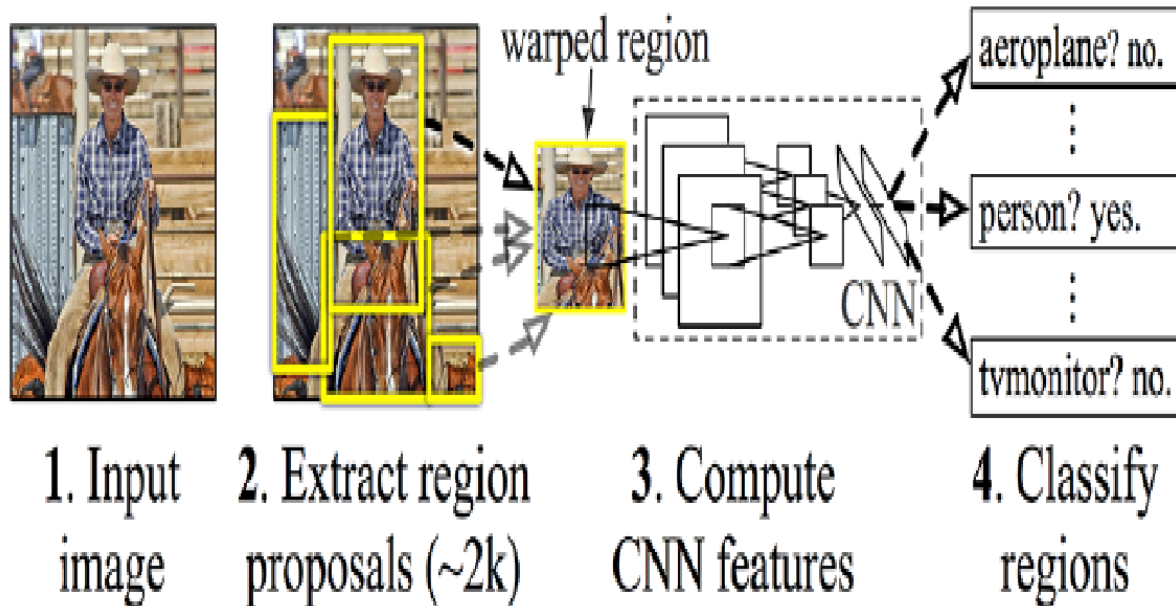
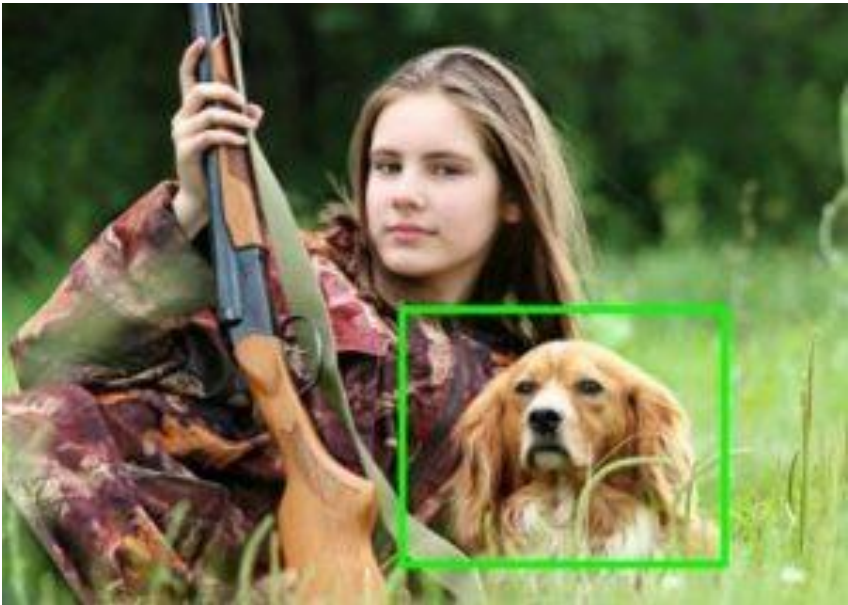
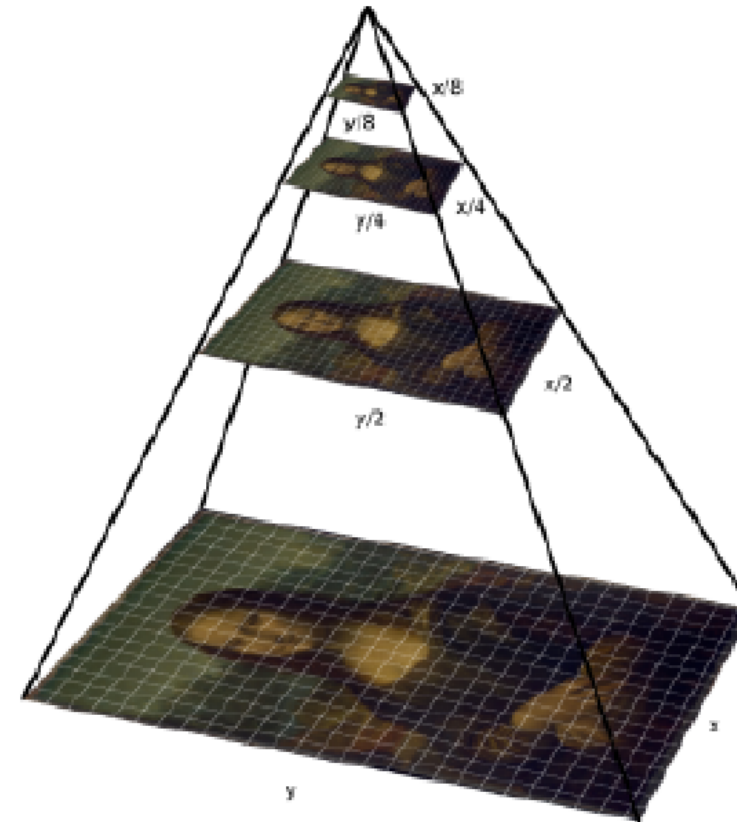
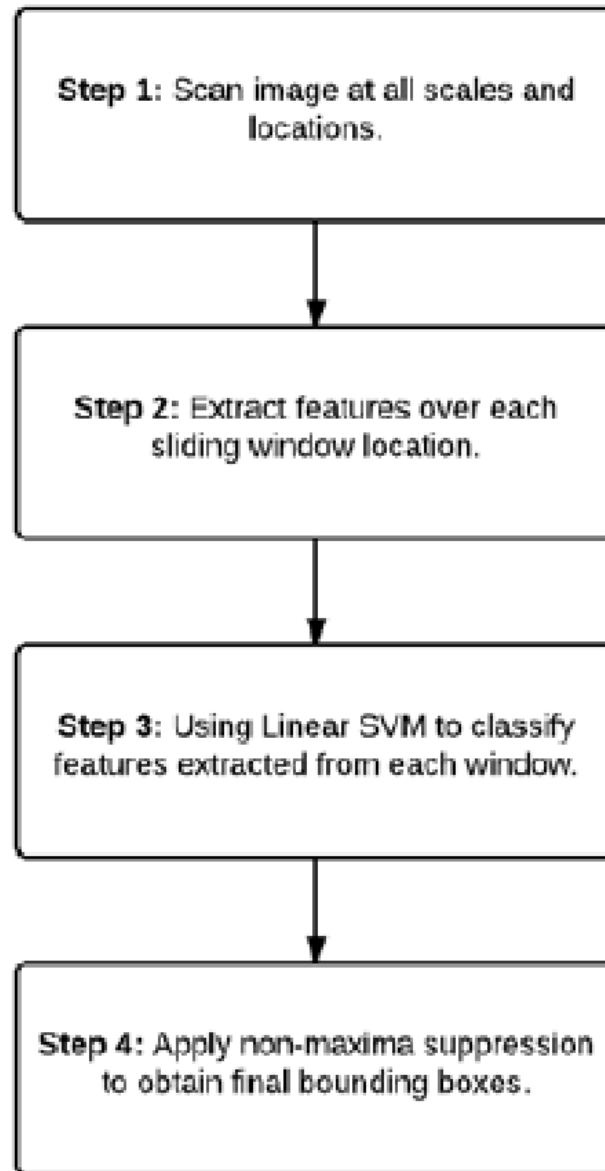


IMAGE PYRAMIDS

- An **image pyramid** is simply a **multi-scale representation** of an image. Using an image pyramid allows us to find objects in images at **different scales** of an image.
- A **sliding window** requires fixed spatial dimensions. If the object in the window is too large or small for the sliding window size, we can miss the detection.

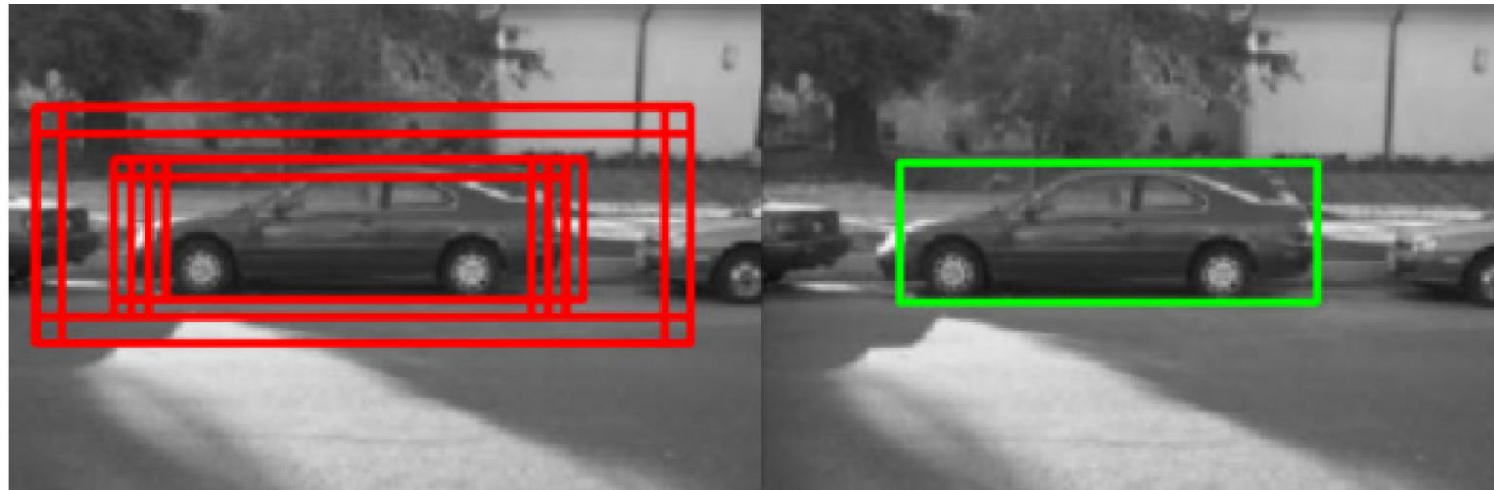


OBJECT DETECTION PIPELINE

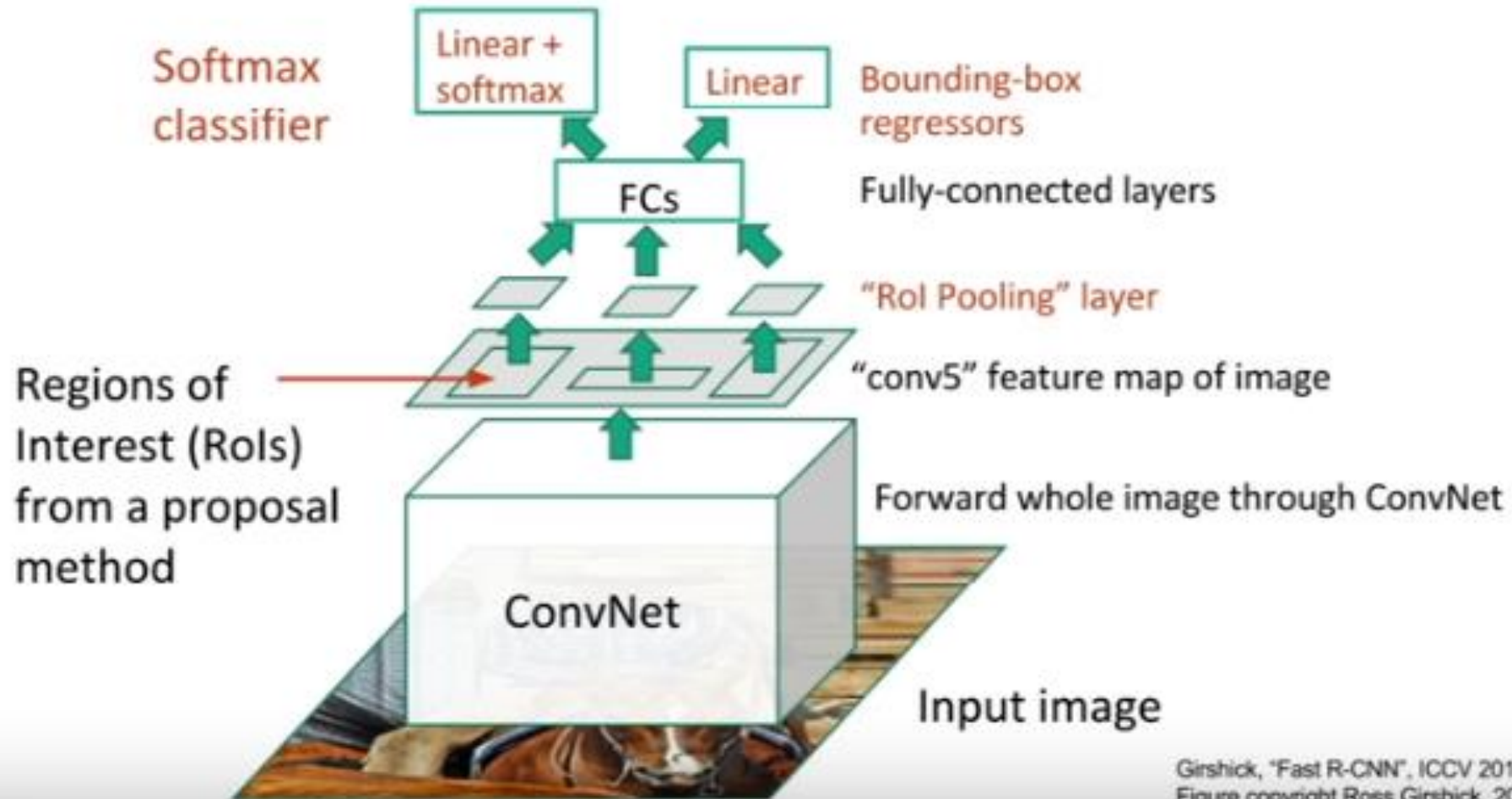


NON-MAXIMA SUPPRESSION

- When combined with image pyramids, this behavior implies that we will have **multiple bounding boxes** surrounding the object at multiple scales, even though there may only be one “true” object in the image.
- To handle the removal of overlapping bounding boxes (that refer to the same object) we can apply **non-maxima suppression (NMS)**.
- **NMS** works by computing the ratio of overlap between bounding boxes, then suppressing (i.e., removing) bounding boxes that have significant overlap.



FAST R-CNN

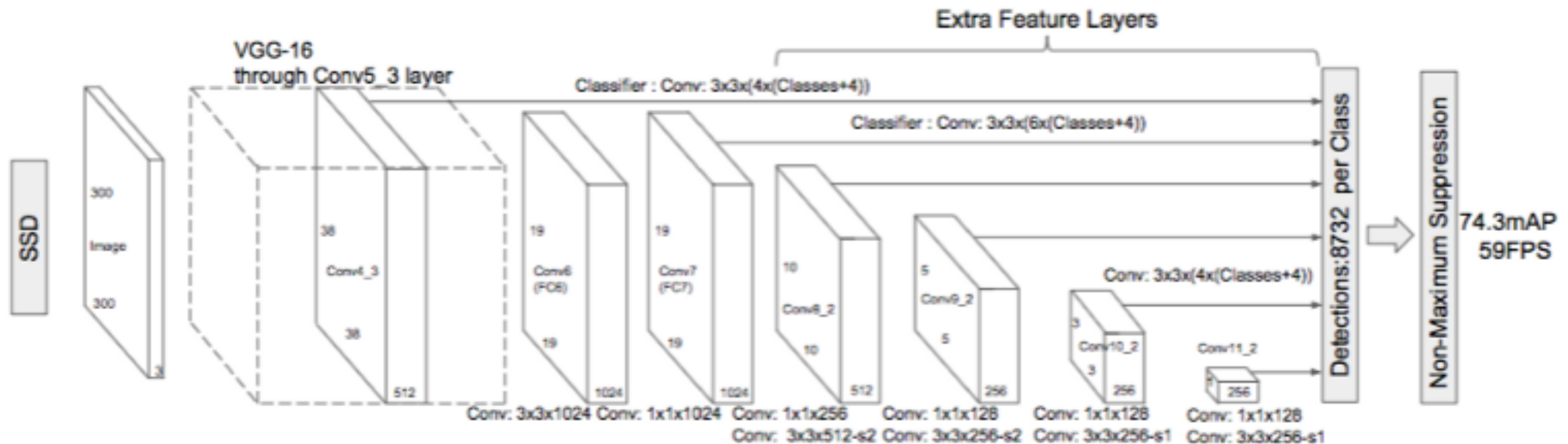


Girshick, "Fast R-CNN", ICCV 2015.

Figure copyright Ross Girshick, 2015. Reproduced with permission.

SINGLE-SHOT DETECTOR (SSD)

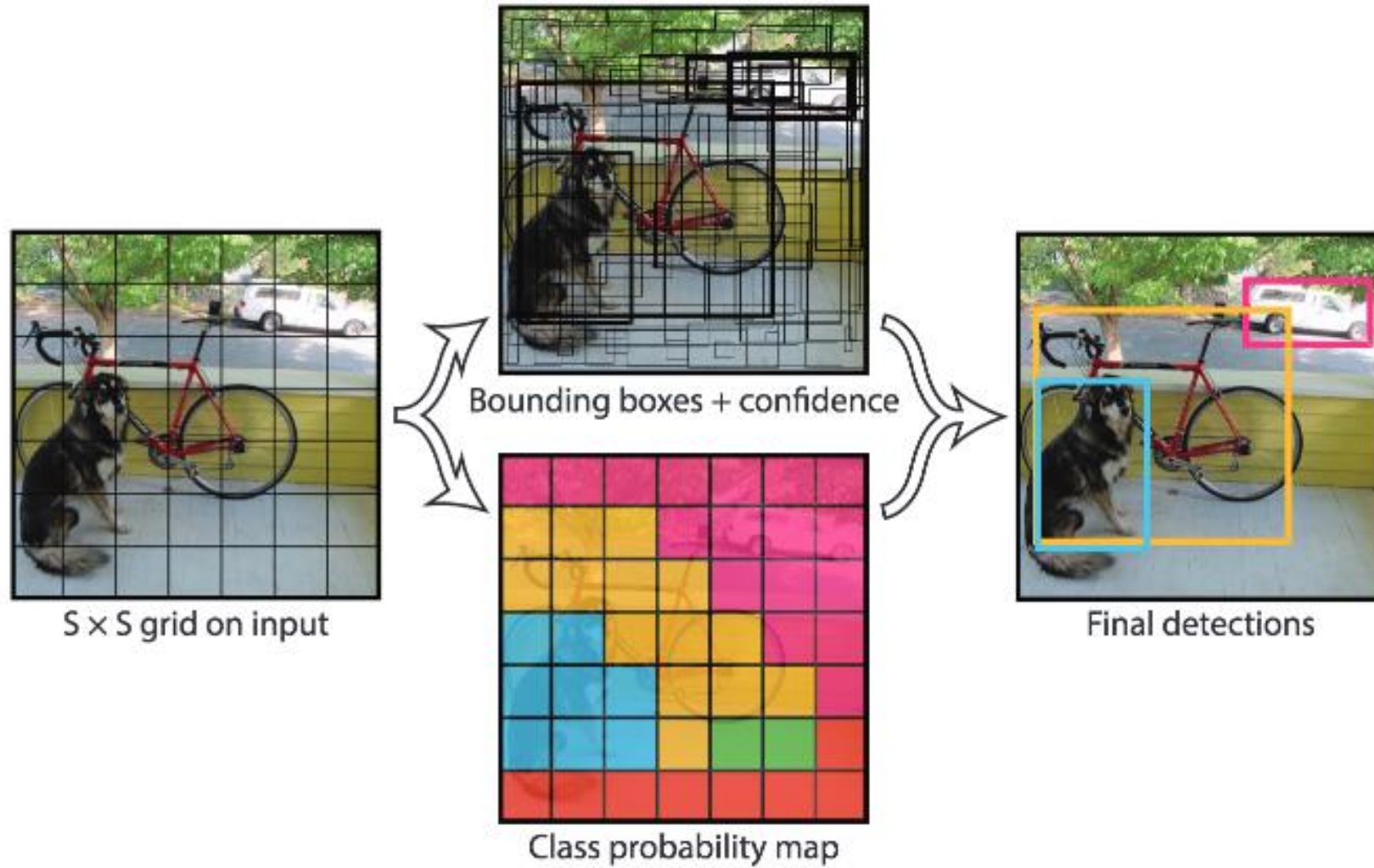
- **SSD** architecture builds on the **VGG-16** architecture by removing the fully connected layers.
- A set of *auxiliary* convolutional layers are added for extracting the features at multiple scales and progressively decrease the size of the input to each subsequent layer.



YOLO (You Only Look Once)

- **YOLO** divides each image into a **grid of $S \times S$** and each grid predicts **N bounding boxes and confidence**.
- The confidence reflects the accuracy of the bounding box and whether the bounding box actually contains an object (regardless of class).
- YOLO also predicts the classification score for each box for every class in training. It can combine both the classes to calculate the probability of each class being present in a predicted box.
- So, total $S \times S \times N$ boxes are predicted. However, most of these boxes have low confidence scores and if we set a threshold say 30% confidence.

YOLO OBJECT DETECTION

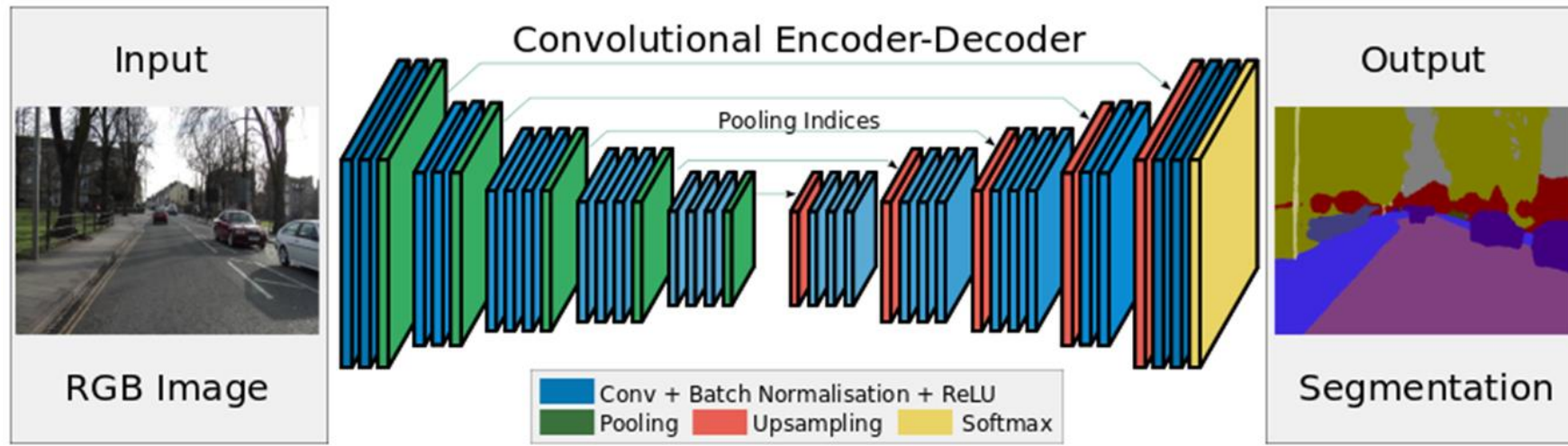


CNN FOR IMAGE SEGMENTATION

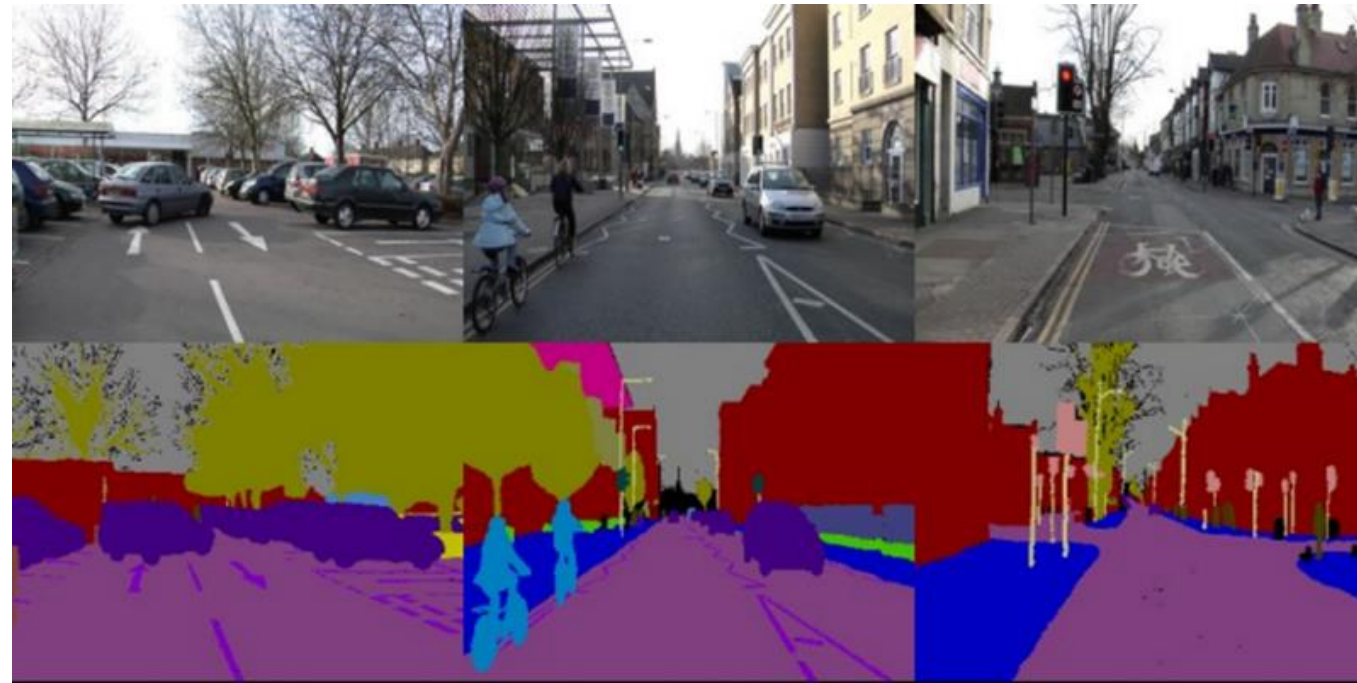
- Pixel level semantic segmentation of a road scene understanding.
- Assign a **label** to **every pixel** in an image.
- Object detection, medical imaging, machine vision, traffic control systems.



SEGNET



- Deep encoder-decoder architecture for multi-class pixelwise segmentation .
- Keras and Caffe library implementation.
- Indoor and outdoor scene understanding.
- 11 classes .





Demo



THANK YOU FOR YOUR KIND ATTENTION!

Twitter - @mohanrajphd