

201900

Muhammad Hassan Mukhtar

Assignment 02

Artificial Neural Network

1.

In You Only Look Once there are several anchor boxes also known defined as priors boxes or default boxes.

V1 2 anchor boxes, No grid formation (single scale)

V2 5 anchor boxes, 13×13 grid (at each scale)

V3 3 anchor boxes (at 3 different scales)
 13×13 , 26×26 , 52×52 grids (at each scale)

V4 Same as V3 Mosaic Training
combine multiple grids

V5 same as V4 and improved anchor box initialization.

V6 Introduced "anchor free" approach
uses Key point Detection and Non maximum Separation

V7 Introduced "Auto Assign Labels"
"dual-teaching"
In dual-teachers model V7 learn from 2 teachers one is rigid and other is easy in such manner Knowledge distilled and learning is Consistance.

V8

Same as
Previous models

Anchor free

16x16, 32x32, 64x64 grid (at each scale)

Mosaic Training (combine multiple grid)

Improved auto assign label (IoU)

Spatial attention

- ① Apply Average Pooling and Max pooling along channel axis and concat them.
- ② On concat feature descriptors you apply Conv layers. which encodes where to emphasize or suppress.
- ③ Then aggregate channel info into 2D map by using two pooling operators.

$$M_s(F) = \sigma(f_7 \times 7([AvgPool(F), MaxPool(F)]))$$

$$M_s(F) = \sigma(f_7 \times 7([F_{avg}, F_{max}]))$$

$\sigma \Rightarrow$ Sigmoid function $f_7 \times 7 \Rightarrow$ Conv operation with filter 7x7

New use cases :- Instance Segmentation.
Keypoint Detection.

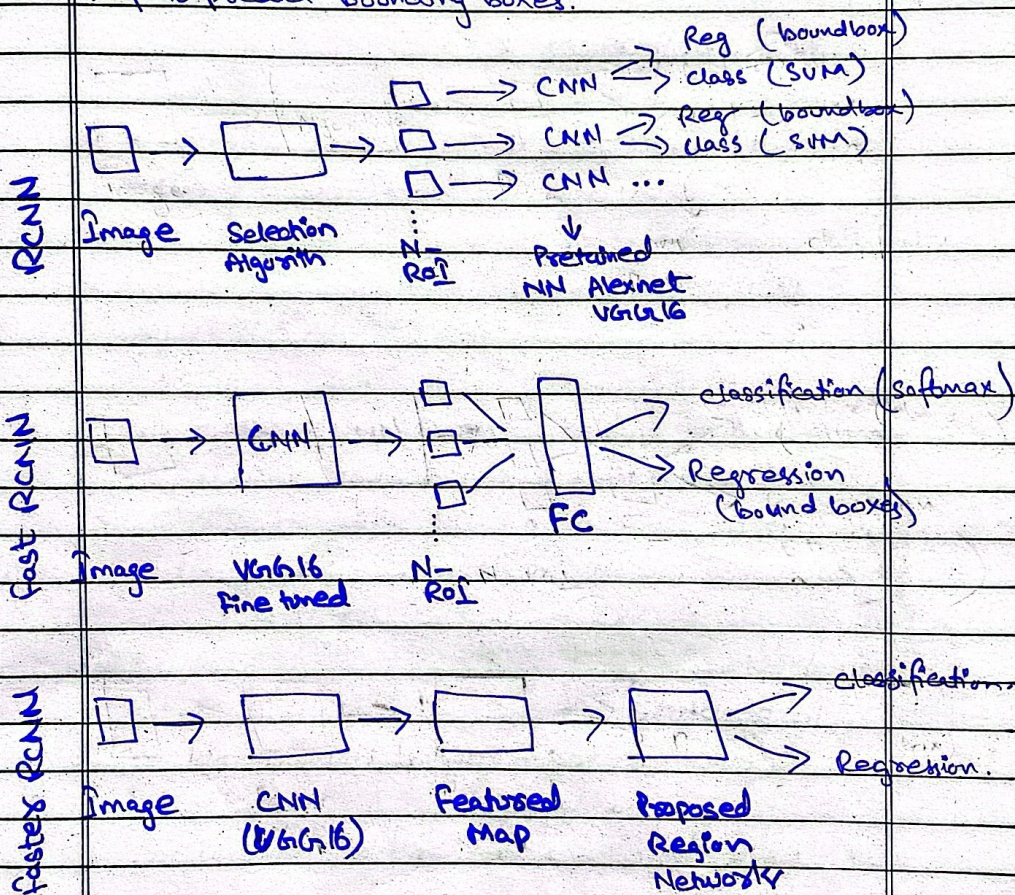
2.

Region Based Conv Neural Network Family

RCNN uses selective search algorithm to propose Region of Interest (RoI) and then uses CNN to classify each RoI.

Fast RCNN uses shared convolution layers to whole image and all RoI rather than dealing with each separately (Reduces computational time)

Faster RCNN uses Region Proposed Network (RPN) to generate RoI which is faster than selective search algorithm. These networks also use feature map to predict bounding boxes.



Sliding Windows in CNN:

9 Anchors of different scales $\frac{1}{4}, \frac{1}{2}, 1, 2, 4$ are applied to feature map. Then fed to proposed Region Network.

Non-Maximum Suppression :

On output of RPN it is applied to remove duplication of detection.

Note:: Bounded Box with max confidence is chosen.

3.

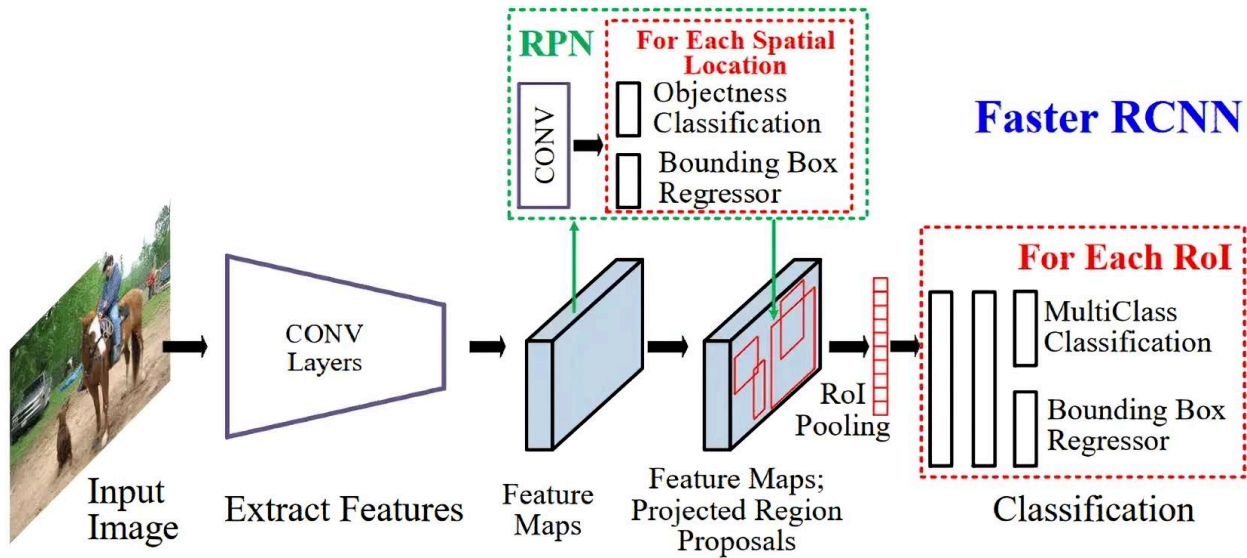
Key Differences :

1. RCNN uses selective search Algorithm while Yolo uses single neural network to predict bounding boxes.
2. RCNN has two stage approach (Region proposal and classification). Yolo is single stage.
3. RCNN prioritize ~~speed~~ accuracy over ~~accuracy~~ speed while Yolo prioritized speed over accuracy.
4. Yolo is used on live video stream while RCNN benefits for high accurate decision making on regions.

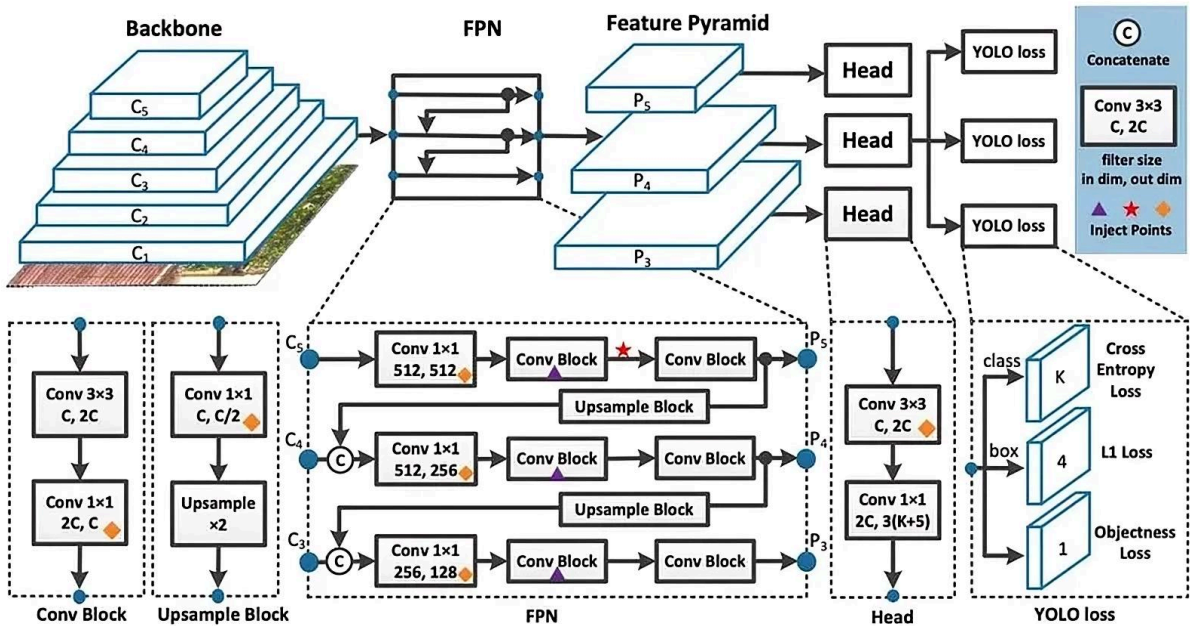
In structure

5. RCNN has Region Proposal Network and RoI pooling extra layers.
While Yolo uses one grid cell generation.

Architecture Difference



Faster RCNN



YOLO v8

Faster R-CNN:

- **Two-stage Detection:** Faster R-CNN follows a two-stage detection process. In the first stage, it proposes regions of interest (RoIs) using a region proposal network (RPN). In the second stage, it classifies these proposed regions.
- **Region Proposal Network (RPN):** The RPN generates region proposals by sliding a small network over the convolutional feature map.
- **Region of Interest Pooling:** After obtaining the region proposals, Faster R-CNN uses a process called region of interest pooling to extract fixed-size feature maps from these regions, which are then fed into a classifier.
- **Backbone Network:** Typically, Faster R-CNN uses a pre-trained convolutional neural network (CNN) such as VGG16, ResNet, or a similar architecture as its backbone network to extract features from the input image.

YOLO (You Only Look Once):

- **Single-stage Detection:** YOLO operates as a single-stage detector, meaning it performs both object localization and classification in a single pass through the network.
- **Grid-based Detection:** YOLO divides the input image into a grid and predicts bounding boxes and class probabilities directly from the grid cells.
- **Unified Framework:** YOLO predicts bounding boxes and class probabilities using a single neural network architecture, without needing separate region proposal and classification stages.
- **Loss Function:** YOLO uses a joint loss function that combines localization loss (for bounding box regression) and classification loss, optimizing both tasks simultaneously.
- **Feature Extraction:** YOLO has its own feature extraction network, typically composed of convolutional and pooling layers, which extracts features directly from the input image.