

YOLO [You only look once]

[used for advance object detection]

- It is mostly used in real time object detection.
- Ex :- Traffic cameras.

use of YOLO → ① Different way of coding

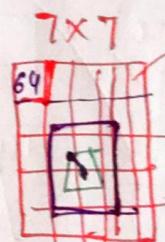
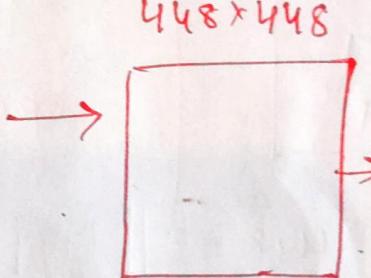
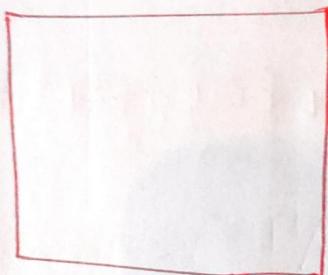
② Real time object Detection

③ One shot Training by using Regression

④ Losses used (BB, Object \rightarrow classes)

YOLO V1 Architecture Algorithm

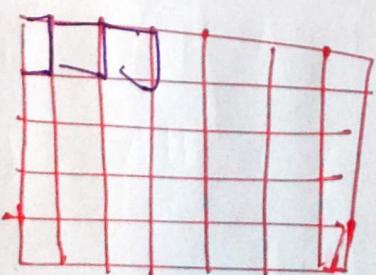
1500 x 700



for each box
one cell
covered

$\Delta x, \Delta y, \Delta w, \Delta h$

- ① Resize the image to (448×448) and $\frac{448}{7} = 64$
- ② divide in 8×8 as (7×7) .



$7 \times 7 \text{ BB}$

\rightarrow for each cell Output \rightarrow one hot encode

$$[\Delta x, \Delta y, \Delta w, \Delta h, c_i] + P(\text{2 classes})$$

$$5 + 20 \leftarrow (\text{parameters})$$

$$= 25 \text{ Parameters}$$

	Δx	Δy	Δw	Δh	ei	f	P^1 (20 classes)
A_1	0	0	0	0	0	0	(0.00-0.20)
A_2	0	0	0	0	0	0	0

A_3

A_{gt}

0.31	0.7	0.1	0.5	1	$[0.00 \ 1.00 \ 0]$
--------	-------	-------	-------	-----	---------------------

A_{obj}

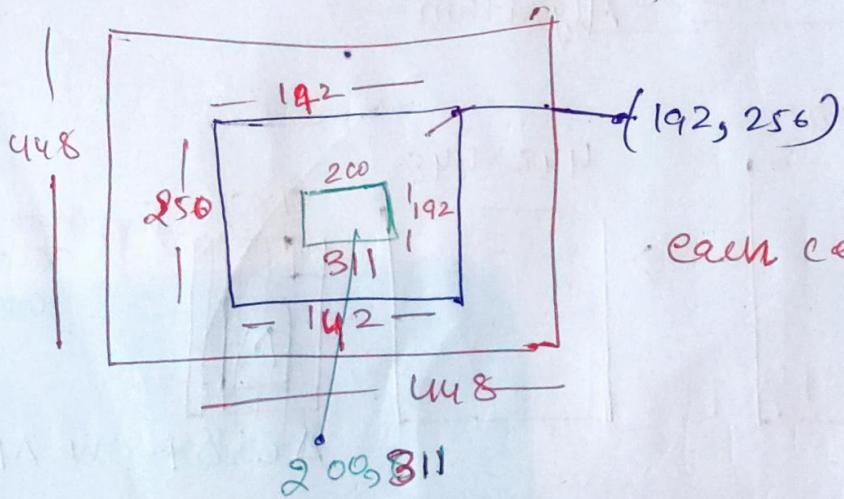
0.4	0.8	0.6	0.7	1	$[0.0 \ 0.0 \ 0.1]$
-------	-------	-------	-------	-----	---------------------

A_{pred}

0	0	0	0	0	- - -
---	---	---	---	---	-------

A

Ground Truth



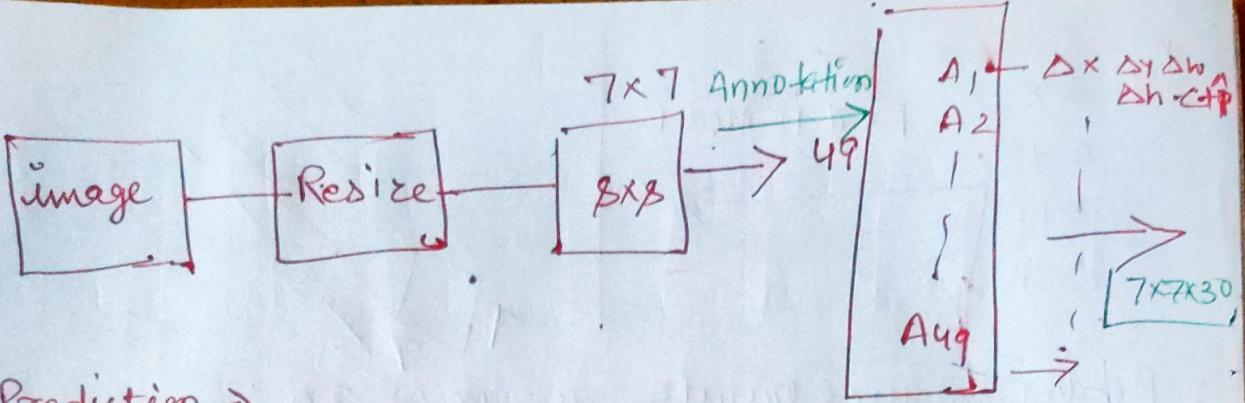
each cell grid = 64

$$\Delta x = \frac{200 - 192}{64} = x_1$$

$$\Delta y = \frac{311 - 256}{64} = y_1$$

$$\Delta w = \frac{142}{448}$$

$$\Delta h = \frac{142}{448}$$



Prediction →

Take 2 class / output

$$\Delta x \Delta y \Delta w \Delta h \cdot \underbrace{\Delta c_1 - \Delta c_2}_{C = \text{confidence score}} + p_1 p_2 p_3$$

Prediction

① Image is in $7 \times 7 \times 30$ format

② $\boxed{5} + \boxed{||| - 20 :}$
 $(0.1 \ 0.2 \ 0.8 \dots)$
 probability

- Training
- ① Take image as input
 - ② Resize image into (448×448)
 - ③ Convert image into 7×7 format as (8×8)
 - ④ Make grid cell as $20 + S(\Delta x \Delta y \Delta w \Delta h \Delta c)$

Δx
 Δy
 Δw
 Δh
 Δc

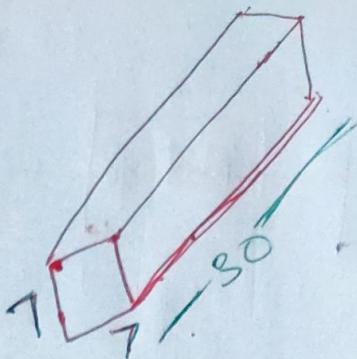
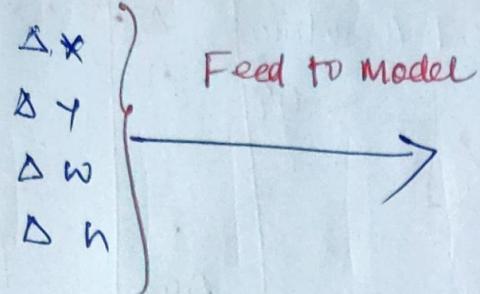
} helps to choose
Bounding Box

$$\left. \begin{array}{l} x \\ y \\ w \\ h \end{array} \right\} \rightarrow \left. \begin{array}{l} x \\ y \\ w \\ h \end{array} \right\} \quad \begin{aligned} x &= \Delta x * 64 + x_a \\ y &= \Delta y * 64 + y_a \end{aligned}$$

$(x_a, y_a) \Rightarrow$ Coordinates of grid cell
of top-left corner

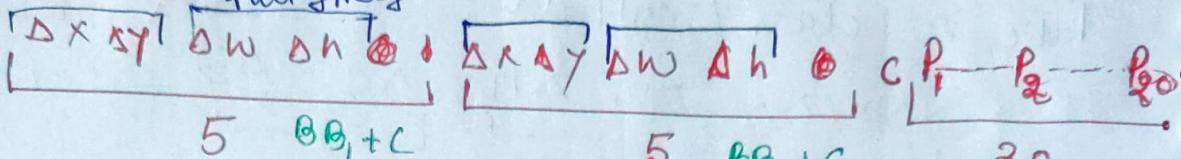
$$w = \Delta w * 448$$

$$h = \Delta h * 448$$



→ Prediction gives results as $7 \times 7 \times 30$

Unid cell
full image



→ Bounding Box added acc to Confidence value (c)

→ Which confidence value is large that Bounding Box

→ then Loss is find for accuracy of model.

YOLO V₁) → 24 conv layers + 2 fully connected layers

Mean Avg. Precision(MAP) of FasterRCNN = 70%.

MAP of YOLO V₁ = 68%.

MAP of YOLO V₂ = 78%.] high accuracy

* Disadvantages of YOLO (V₁) =>

(1) It can only detect 4 objects.

(2) Less accuracy as 68% than YOLO(V₂) & FasterRCNN.

(3) Fully connected layers is not present at the end.

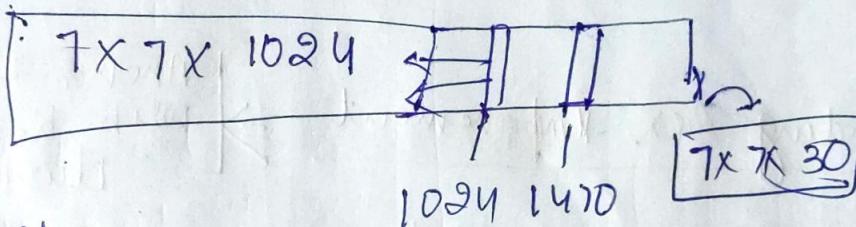
\Rightarrow Modification or improvement in YOLO v₂ over v₁

- ① MAP increases from 68 to 78.
- ② Apply Batch normalization to all CNN layers.
 - By BN, improve map by 2%.
 - overfitting Problem solved
 - Regularization effect came.
- ③ High Resolution classifier used. due to this MAP increases by 4%.

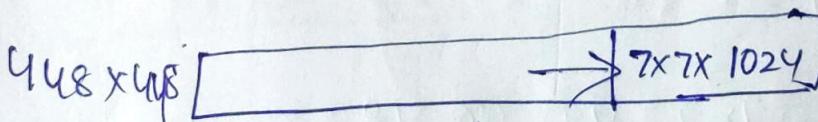


High Resolution Feature MAP \rightarrow

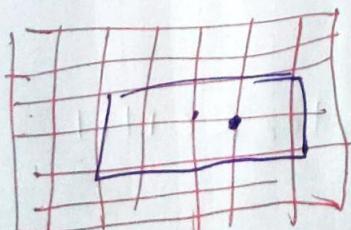
In YOLO v₁



YOLO v₂



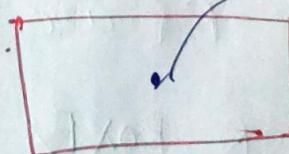
7x7



odd x odd

Center will
be at one
Particular
Point

14x14

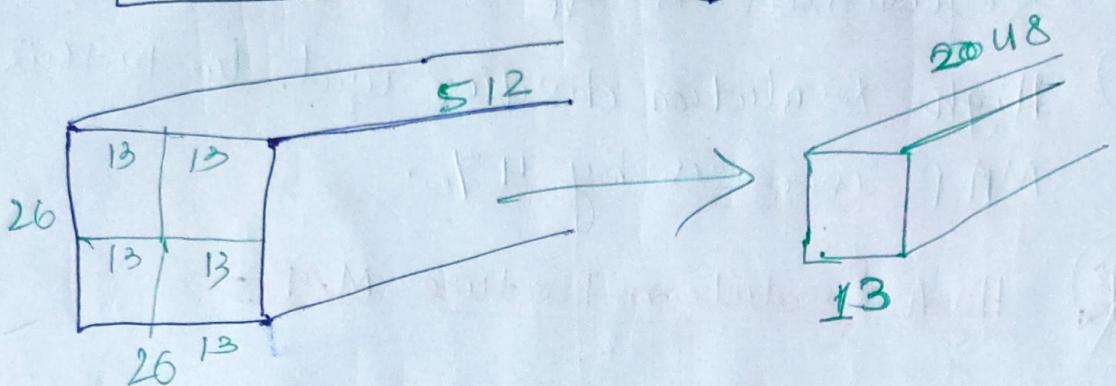
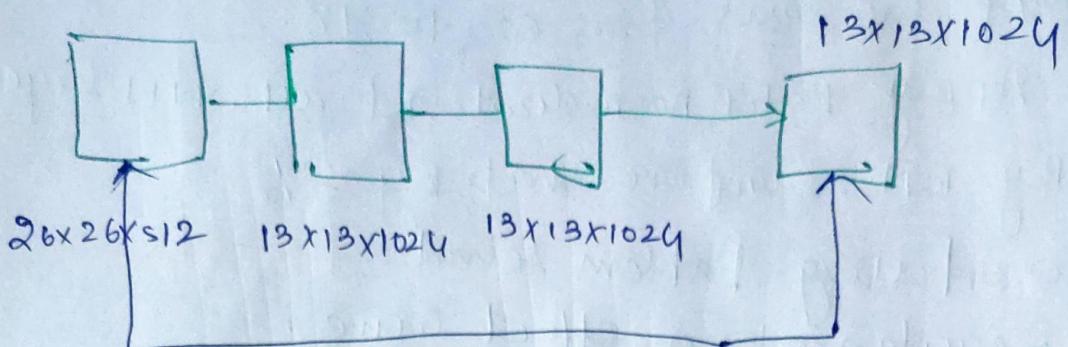


even x even

finding
center is
not easy
so, shift to
size as

13x13
13x13x 1024

④ Pass through Layer \Rightarrow



⑤ Bounding Box Improvement \Rightarrow [Anchor Boxes with K-Means clustering]

\rightarrow here Region Proposal (RP) are approx 98 or

$$\rightarrow 7 \times 7 = 49 \text{ or } 98 \text{ or } 169 \text{ or } 338$$

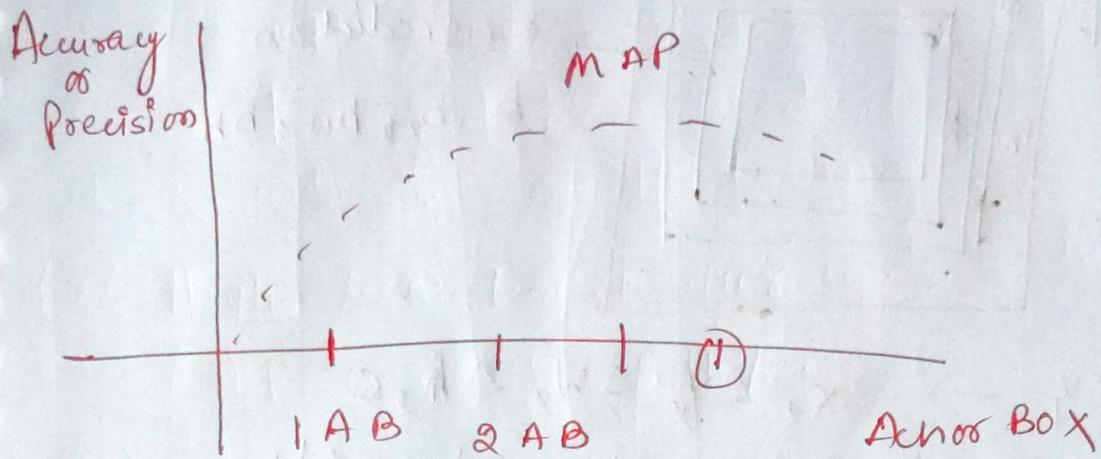
so, no. of RP are less here.

⑥ ANCHOR BOXES used in faster R-CNN

3 aspect ratio

- ANCHOR BOXES \rightarrow Every grid cell can have more than 1 object.

→ 9 anchor boxes present
so, K-Means clustering used on 5 Anchor Boxes for better result.



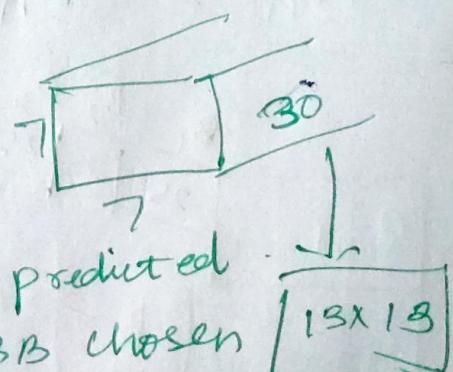
These 5 modifications which improved the model drastically from mAP of 68% to 78%.

* Predictions →

⑥ Different activation function used

* Predictions →

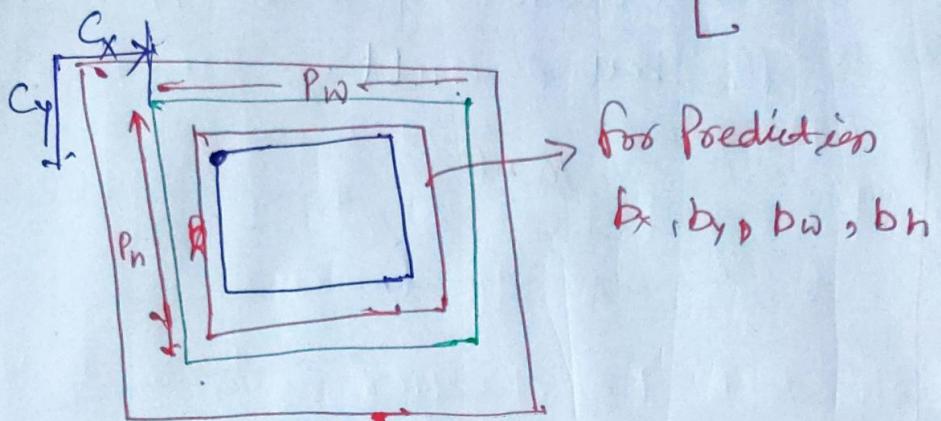
In YOLO v1 → offsets



2 BB predicted
and 1 BB chosen
as per confidence value.

→ YOLO v → Linear Activation fun used and range of O/P of (x,y) → $[-\infty, \infty]$.

YOLO V2 \Rightarrow Activation function + [Constraints]



YOLO V1 $\rightarrow \Delta x, \Delta y, \Delta w, \Delta h, c$

\Rightarrow YOLO V2 $\rightarrow \Delta t_x, \Delta t_y, \Delta t_w, \Delta t_h, c$

$$b_x = \sigma(\Delta t_x) + C_x$$

$$b_y = \sigma(\Delta t_y) + C_y$$

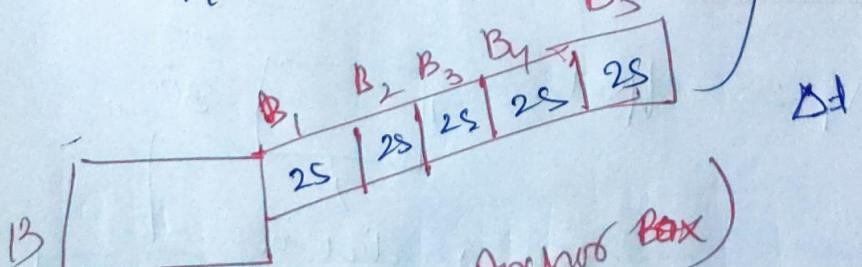
$$b_w = P_w \cdot e^{\Delta t_w}$$

$$b_h = b_n \cdot e^{\Delta t_h}$$

b_x, b_y, b_w, b_n

BB

targeted value



[\downarrow - Anchor Boxes used]

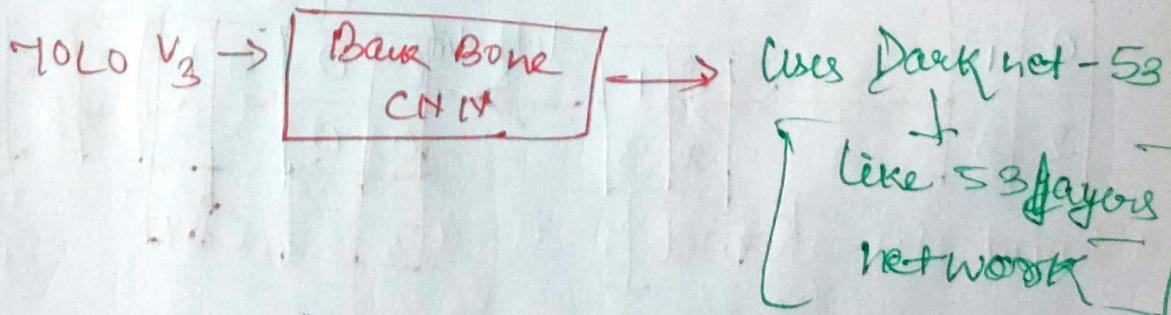
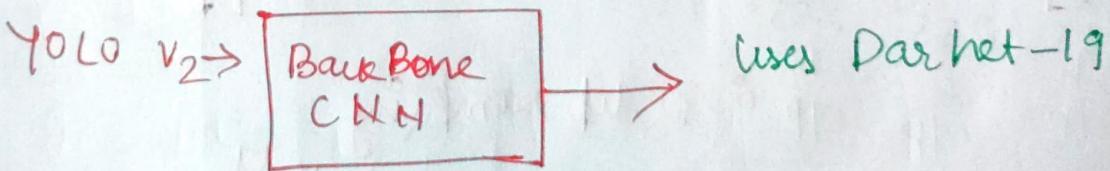
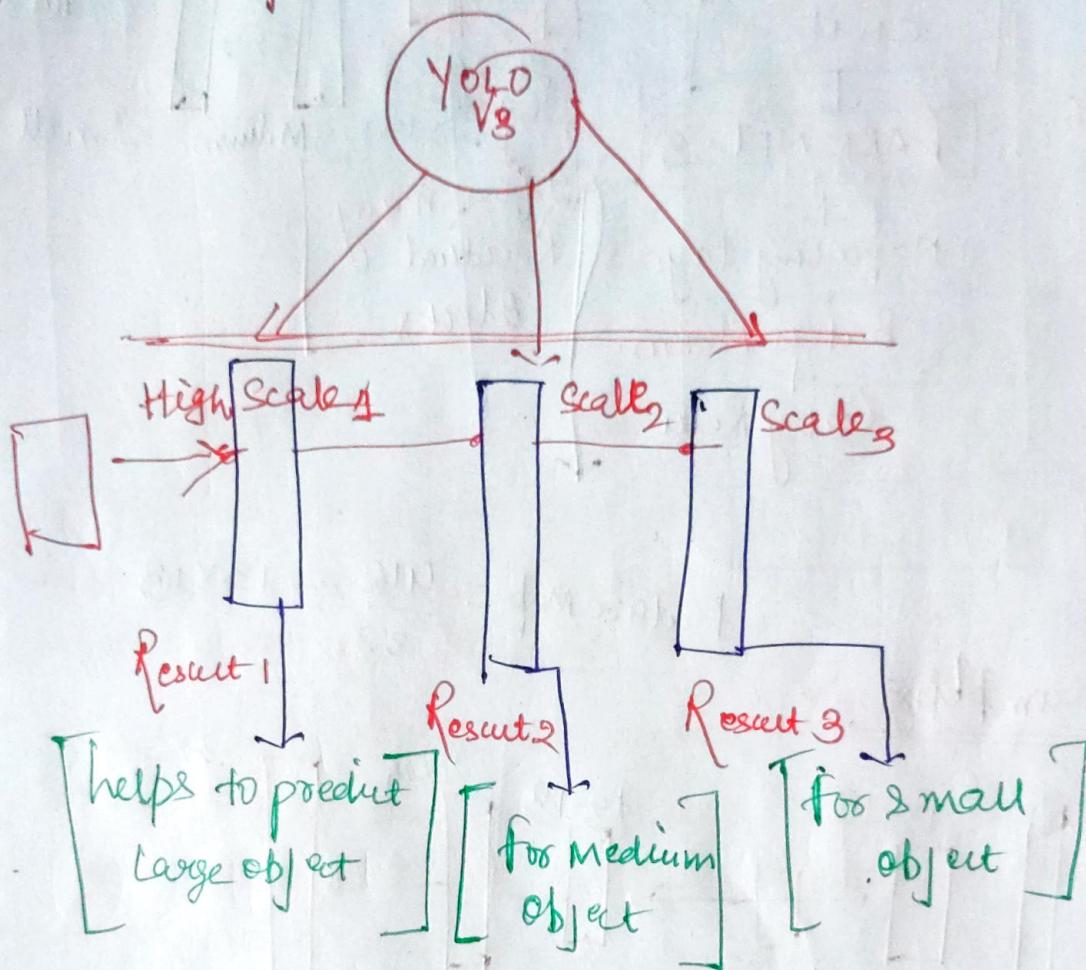
$\Delta t_x, \Delta t_y, \Delta t_w, \Delta t_h, p, c, [20]$

125

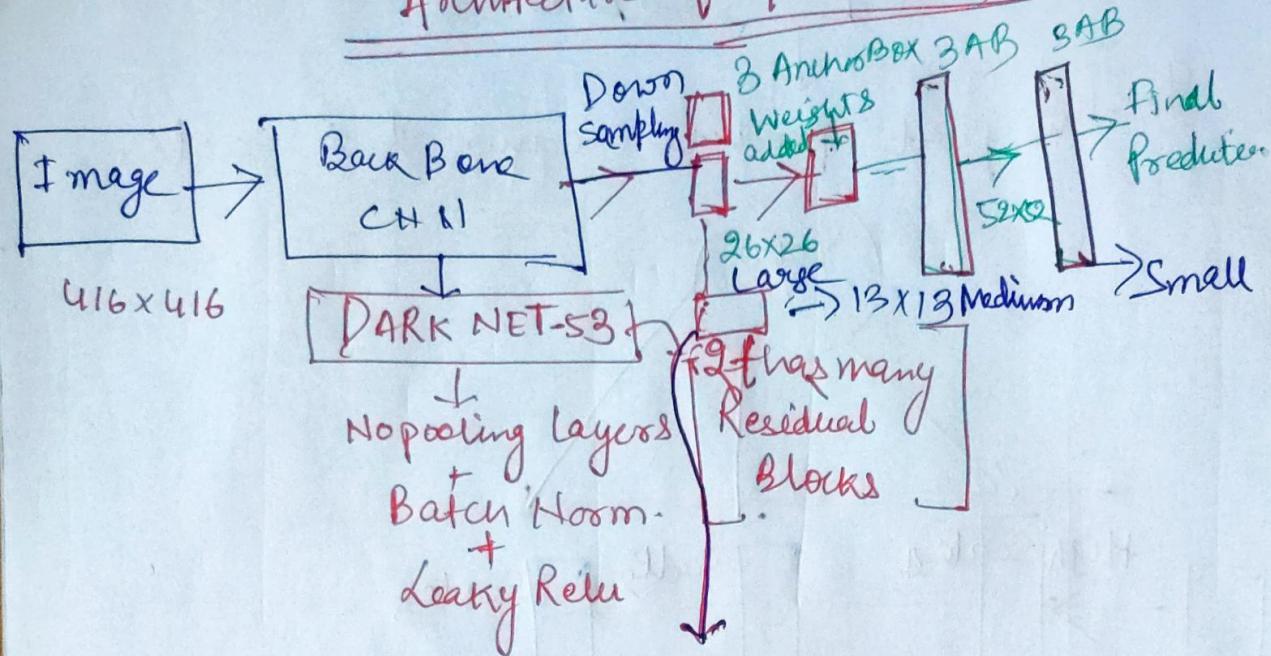
$[13 \times 13 \times 125]$

YOLO v₃

In this part prediction came of 3 parts.

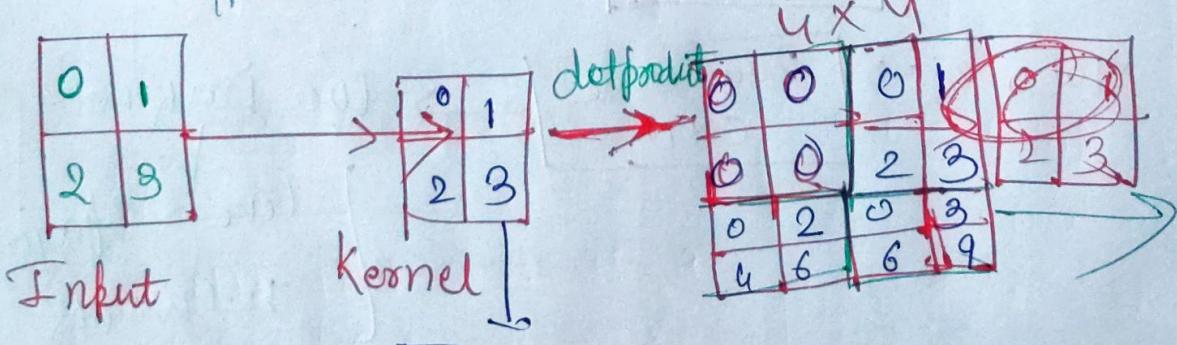
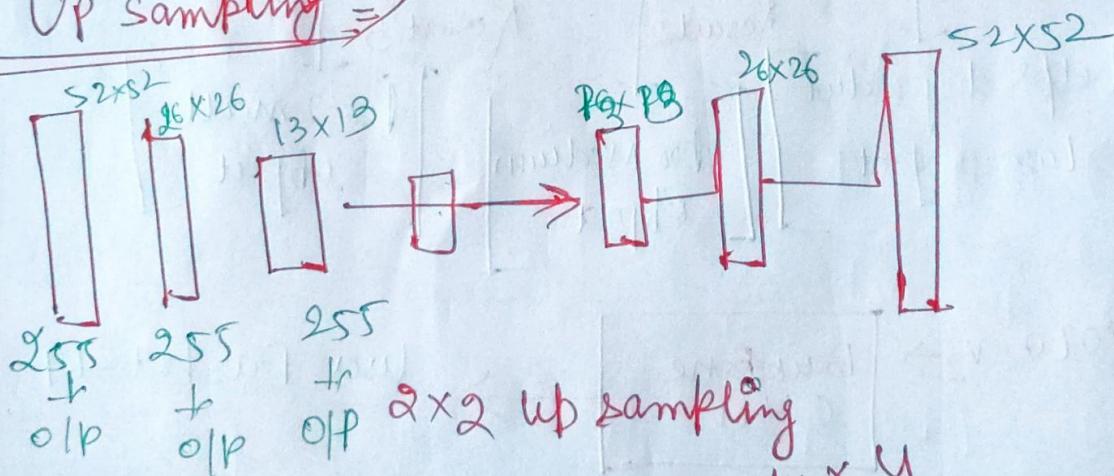


Architecture of YOLO V3



$$\text{Feature Map} = \frac{416}{32} = 13 \times 13$$

Up Sampling



Transpose conv \rightarrow Reverse conv
 ↓
 Subsampling

Add \rightarrow

$0 \ 0$	$0 \ 1$	$0 \ 0 \ 0 \ 1$
$0 \ 0$	$2 \ 3$	$0 \ 2 \ 2 \ 6$
$0 \ 2$	$0 \ 3$	$4 \ 6 \ 6 \ 9$

4×4

↓ add

$0 \ 0 \ 1$
$0 \ 4 \ 6$
$4 \ 12 \ 9$

$\downarrow 3 \times 3$

* In YOLO v3, All Anchor

Box sizes were predetermined

based on the analysis of the dataset used for training YOLO_{v3}

they cover a range of object Transpose conv
 commonly found in the dataset, and the model learns to adjust its predictions based on these anchor boxes during the training period.

* Predictions \Rightarrow

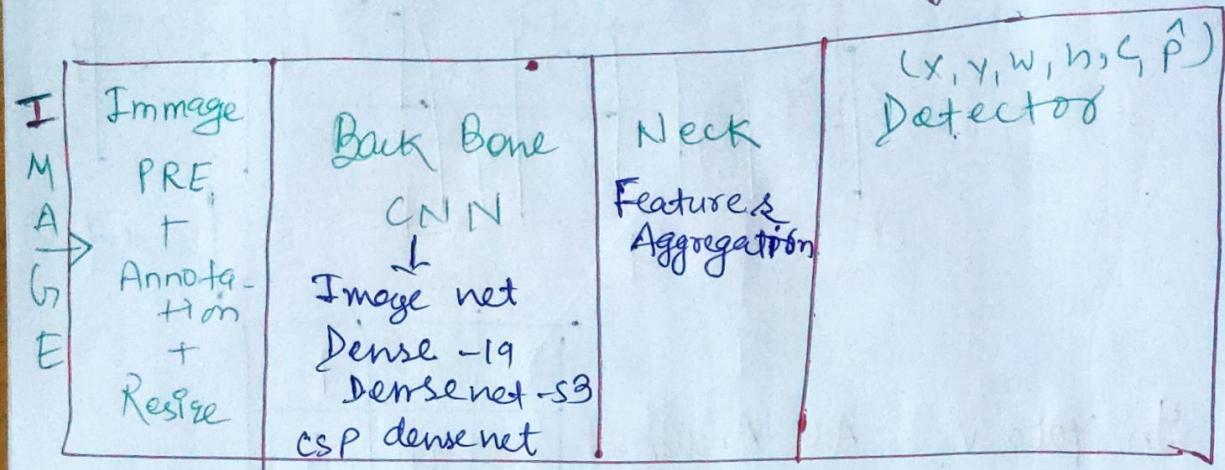
$$\frac{\Delta x \ \Delta y \ \Delta w \ \Delta h \ + c_i, \hat{P}}{8} \quad [80 \text{ classes}]$$

for each Anchor Box

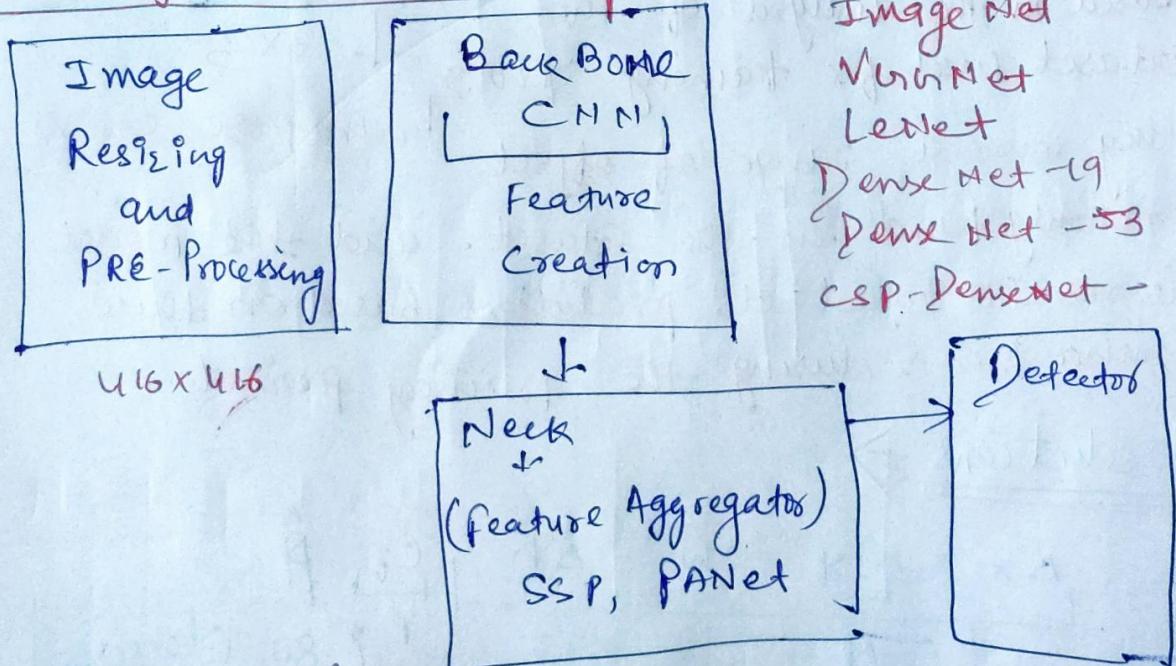
$$80 + 4 + 1 = 85 \text{ length output}$$

YOLO V4

It will be divided into 4 categories.



In object detection →



In YOLO (V3)

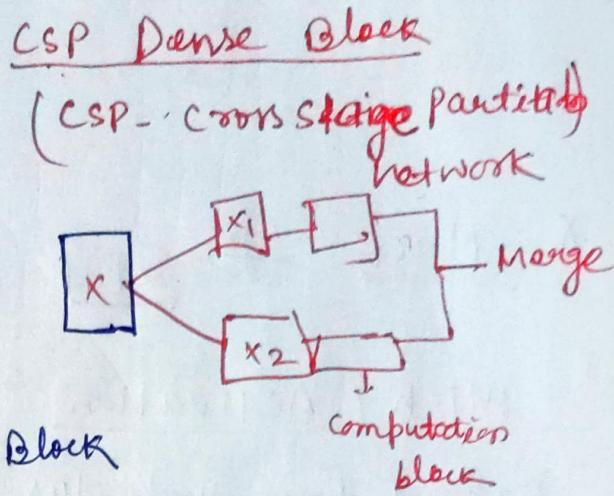
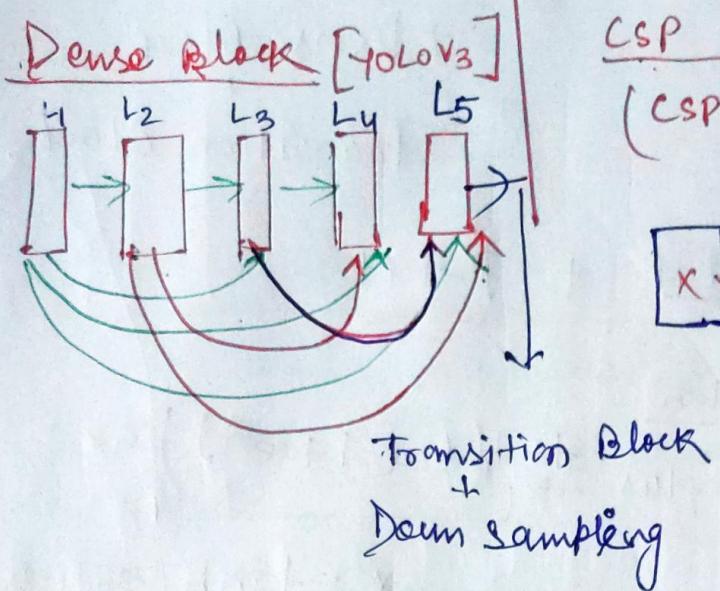
~~Darknet~~

A conv NN that acts as a backbone for the YOLO V3 object detection.

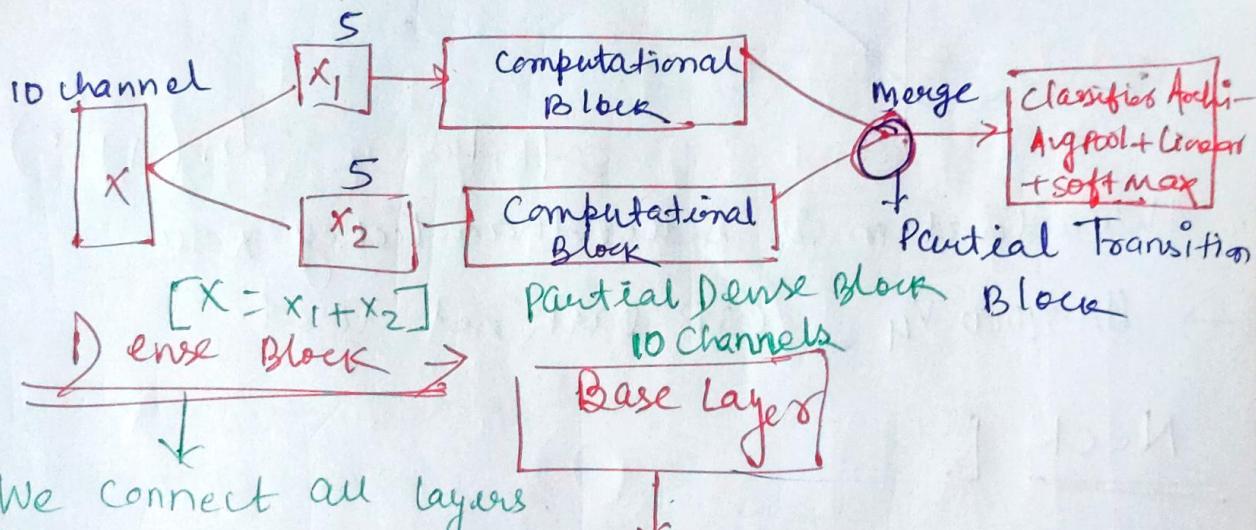
Darknet

A convo CNN and backbone for object detection.
 → Partition the feature map of the base layer into two parts then merge.

- In Dark net leaky
ReLU used
- Mish used as activation function



CSP Dense Block



We connect all layers.

(with matching feature-map size) directly with each other.

Dense Block

↓
Transition

↓
Down Sampling

In YOLO V4

* Backbone CNN \rightarrow CSP + Darknet
+ Dense Block
+ Transition Block

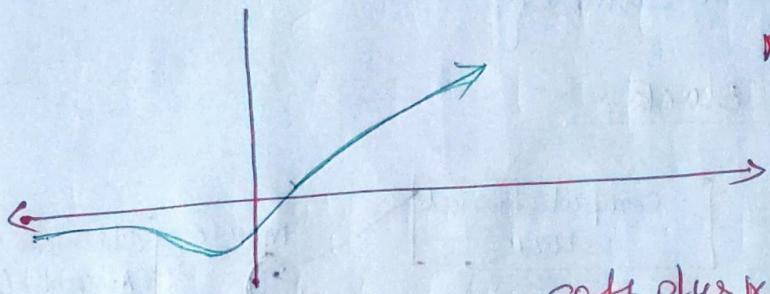
* Neck \rightarrow SPP + PAN

* Mish (Activation)

\rightarrow Derived from softplus

$$\tanh(\ln(1+e^x)) \times x$$

$$x * \tanh(\text{softplus}(x))$$

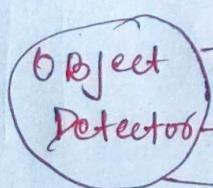


$$\text{softplus}(x) = \ln(1+e^x)$$

this activation is better than Leaky ReLU.

\rightarrow In YOLO V4, CSP Darknet used.

Neck [



BB - CNN - (Feature Map creators)
Neck
Feature Aggregator

Detector \rightarrow x, y, w, h, c, p

We use typically two feature aggregators in your YOLO - V4

(i) PAN | Partial Aggregation Networks)

→ It is modified Version of FPN.

(feature Pyramid Network)

(ii) SPP (Spatial pyramid pooling)

* FPN (Feature Pyramid Network) ⇒

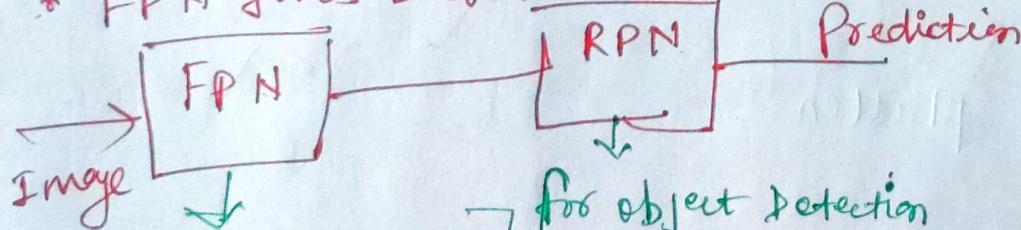
→ It is used for smaller object

→ FPN provides a top-down pathway to construct higher resolution layers from a semantic rich layer.

→ FPN is not an object detector by itself.

→ It is a feature extractor that works with object detectors.

• FPN gives better Result or MAP.

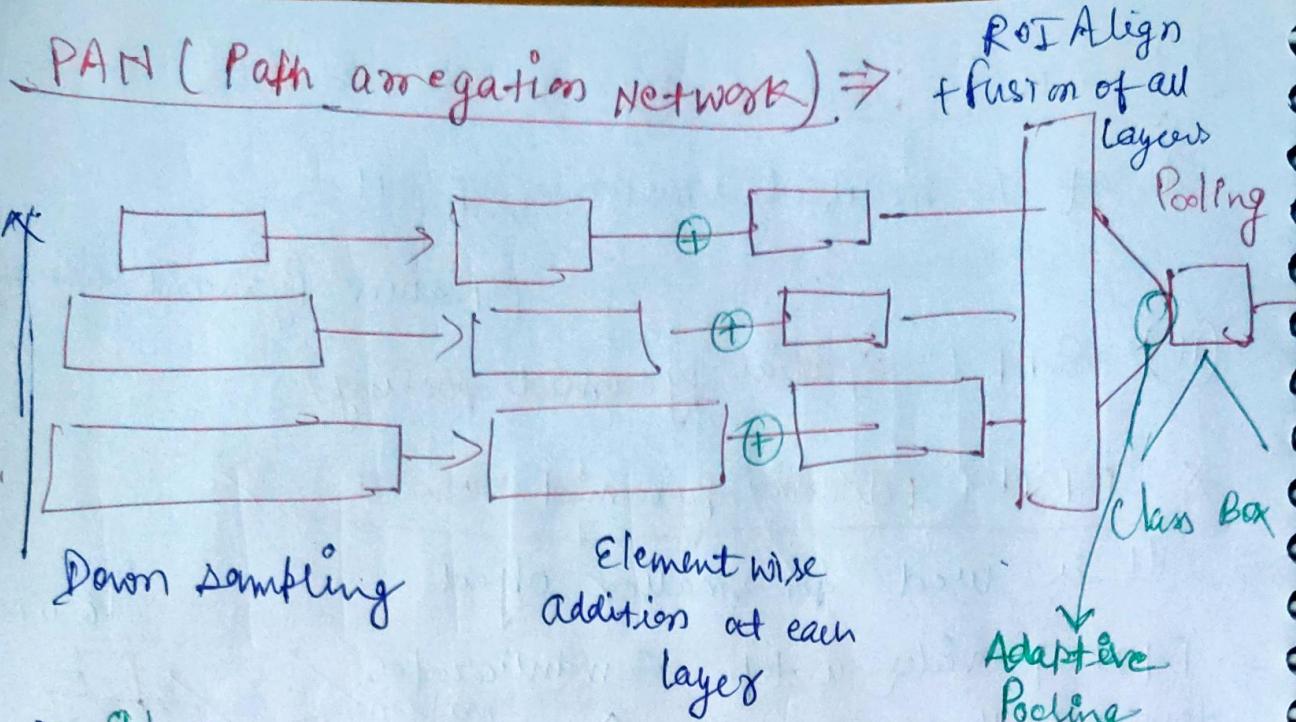


Extracts feature Map → It applies sliding window over the feature maps to

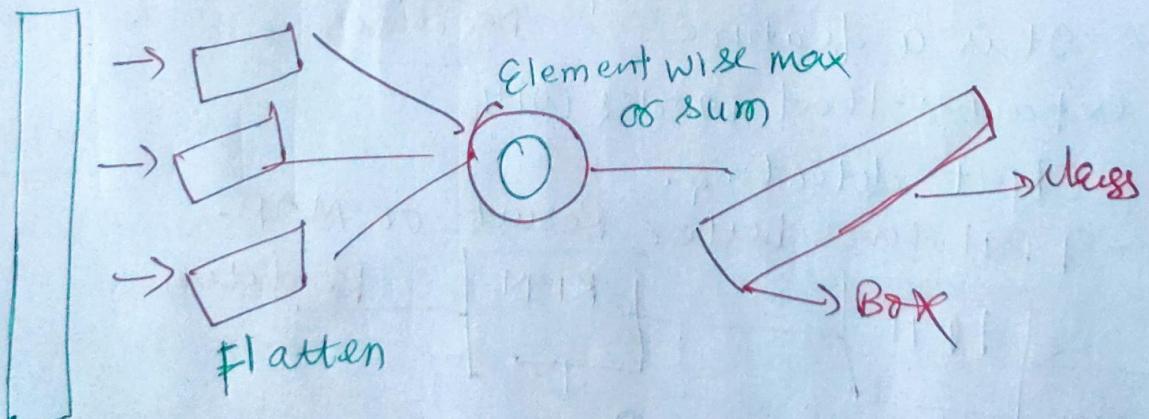
Make predictions on the objectness & object boundary box at each loc.

(feature Pyramid Network)

(feature Pyramid Network)



- \rightarrow It is mainly incorporated of instance segmentation by preserving spatial information.
- \rightarrow It has ability to preserve spatial information accurately which helps in proper localization of pixels for mask formation.

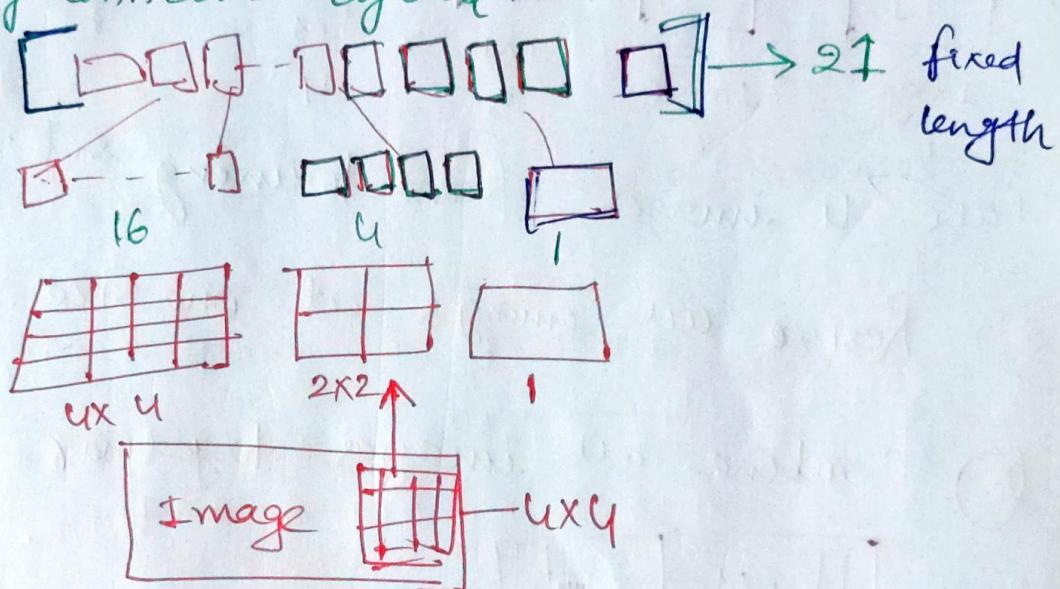


- \rightarrow In YOLO v4 Modified PAN used as we add conv. layers for concatenation.

→ In YOLOv4, Modified PAN uses concat in replace of adding.

SPP (spatial pyramid pooling)

- g_t is used \odot in replace of ^{max} pooling.
- g_t pools the features and generates fixed-length outputs. which are then fed to the fully connected layers.



→ g_t gives better contextual information through scaling

* YOLO v5

There are mainly two changes

- (i) Mosaic Data Augmentation Technique
- (ii) Focus layer

* MOSAIC Data Augmentation Technique

- It helps to design image that out of context object also detected.
- It is used to increase data.

①

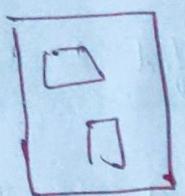


Image
1

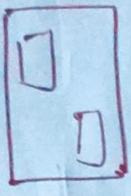
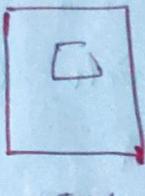
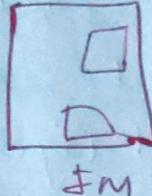


Image
2



IM
3

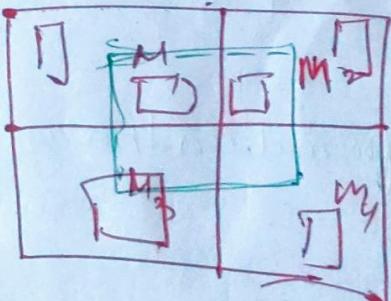


IM
4

Take 4 images from training set.

② Resize all images at one size.

③ Combine all images together



- Mosaic will help to recognise the objects that the model is not used to looking at together.
- All the cases where objects appear in different contexts.
- Mosaic used during training time.

* Bag of Freebies →

During training time increases and Model will not change.

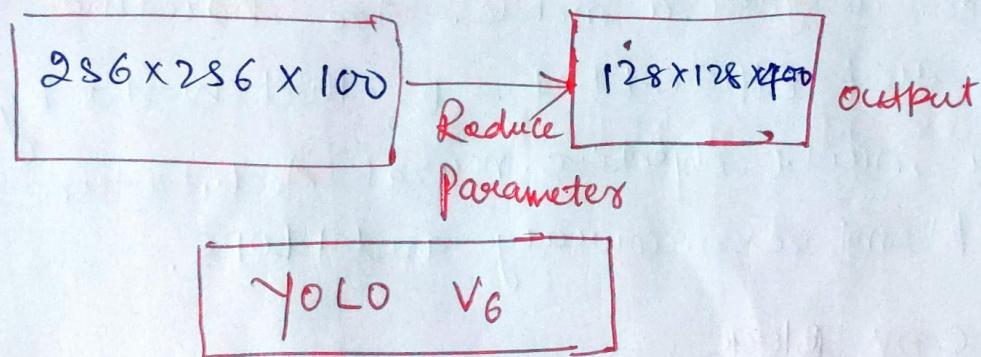
So, Most AC data aug. is this type.

* Bag of specials →

→ Major changes in architecture which increases the complexity & accuracy.

② Focus layer →

→ It is used to reduce no of parameters.



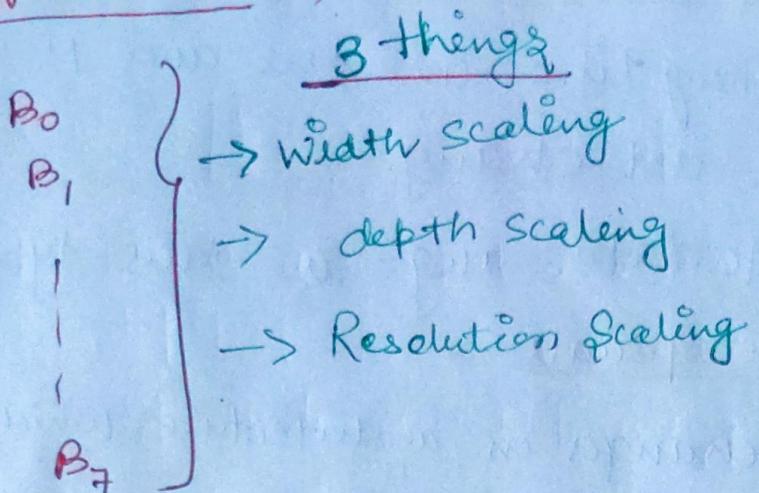
* Whenever we are building the model we will use the latest object detection.

which is based on accuracy & speed.

* Efficient Net used here as BackBone CNN.



* Efficient Net →



→ It expands the network's width, depth and resolution.

Hence, it is critical to have a good baseline network.

The authors designed a mobile-size baseline network called EfficientNet - B₀, that works by using a multi-objective neural architecture that optimizes accuracy and FLOPS.

* MB conv Block →

- ① Squeeze and Excitation Phase
- ② Expansion Phase.

YOLO V7

- New back Bone CNN used.
- this Version is Very accurate and fast.
- ELAN (Efficient lightWeight Attention)

→ It is a backbone CNN architecture that was designed to achieve high efficiency and efficiency in computer vision tasks.
It was introduced in the research paper titled "ELAN" for real time 'Object Detection' by Yeranho Cai, et al., which was published in 2020.

YOLO V7 [conv, Pooling, CBS, Concat, ELAN]
SPP, CSVC