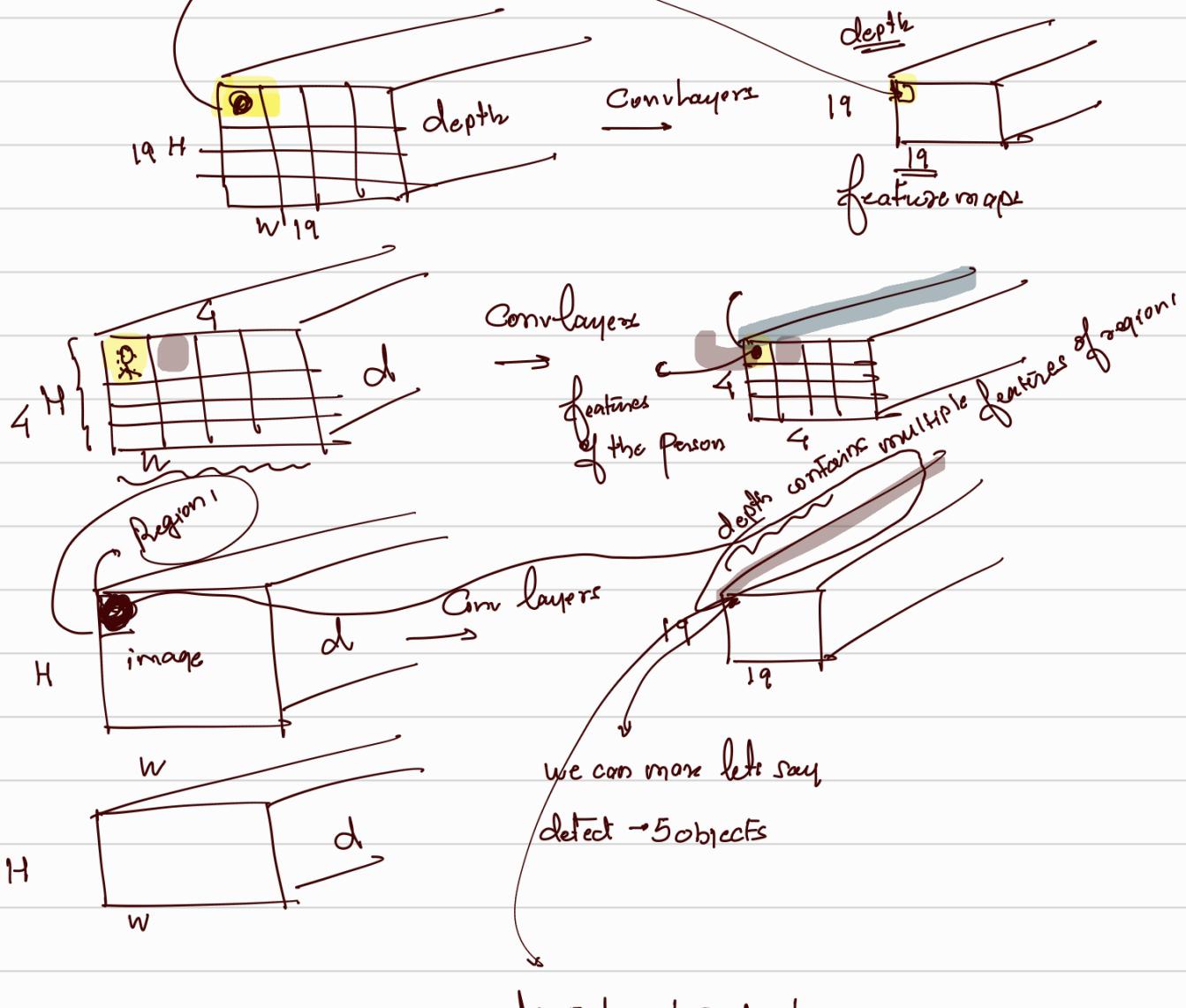


YOLO (You only look once) (Yolov2)

① Object detection algorithms: R-CNN, Fast R-CNN, Faster R-CNN



Within the same network we detect both the region proposals and the bounding boxes

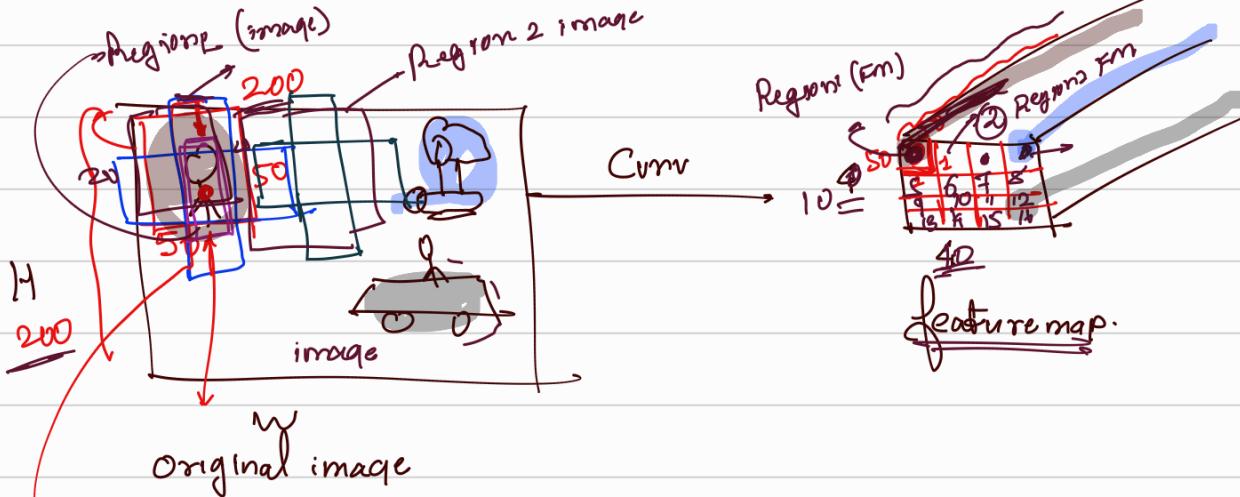


we can detect most 5 objects

\rightarrow Anchor Box1 \rightarrow Anchor Box5

\rightarrow Anchor Box2 \rightarrow Anchor Box4

\rightarrow Anchor Box3



At every 50×50 crop, we can attempt detect here say

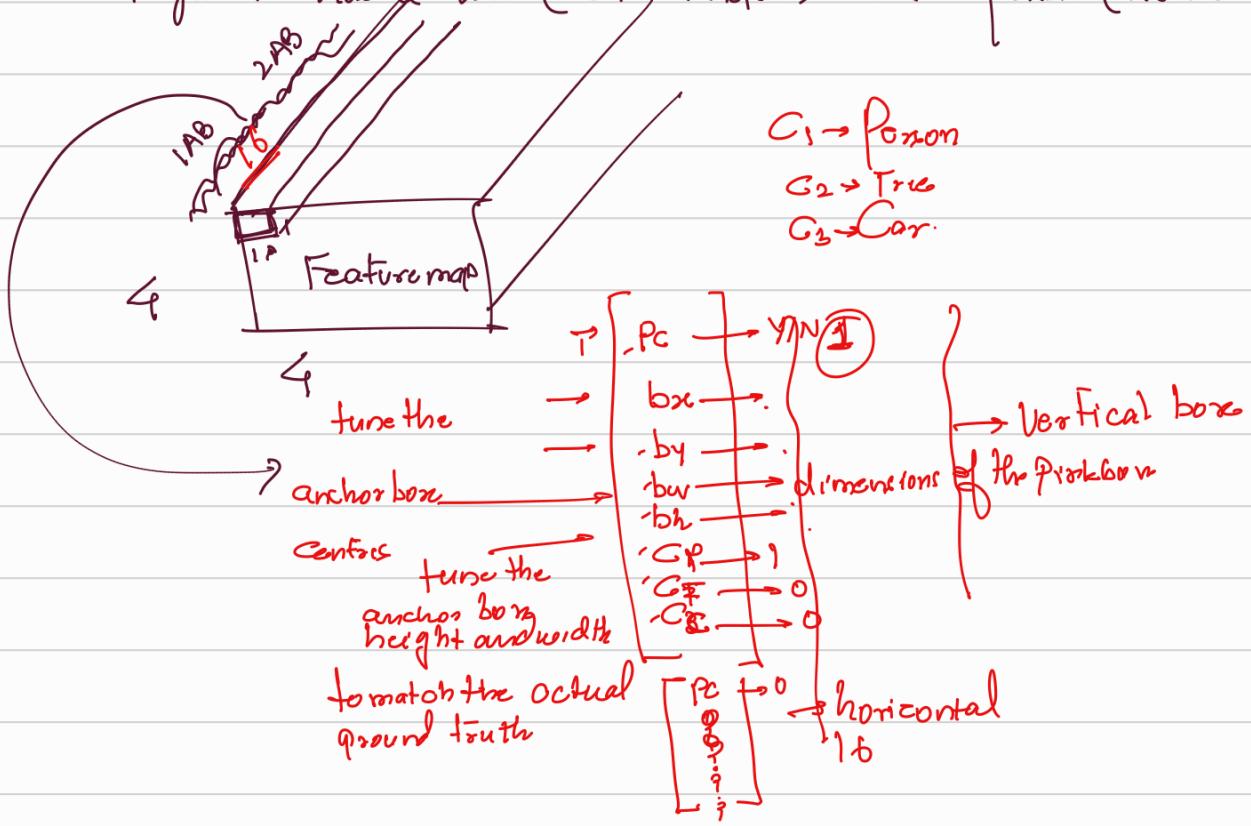
2 objects

- Horizontal object
- Vertical object

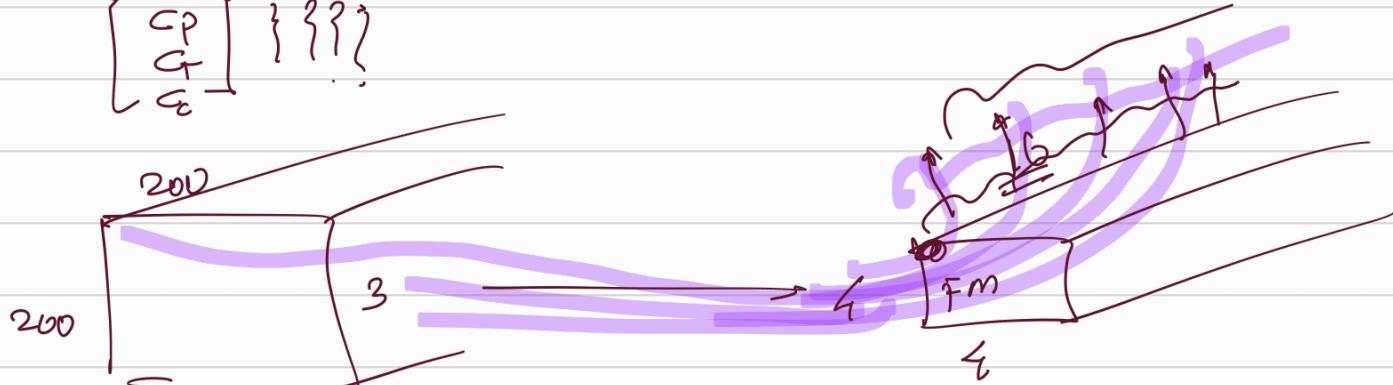
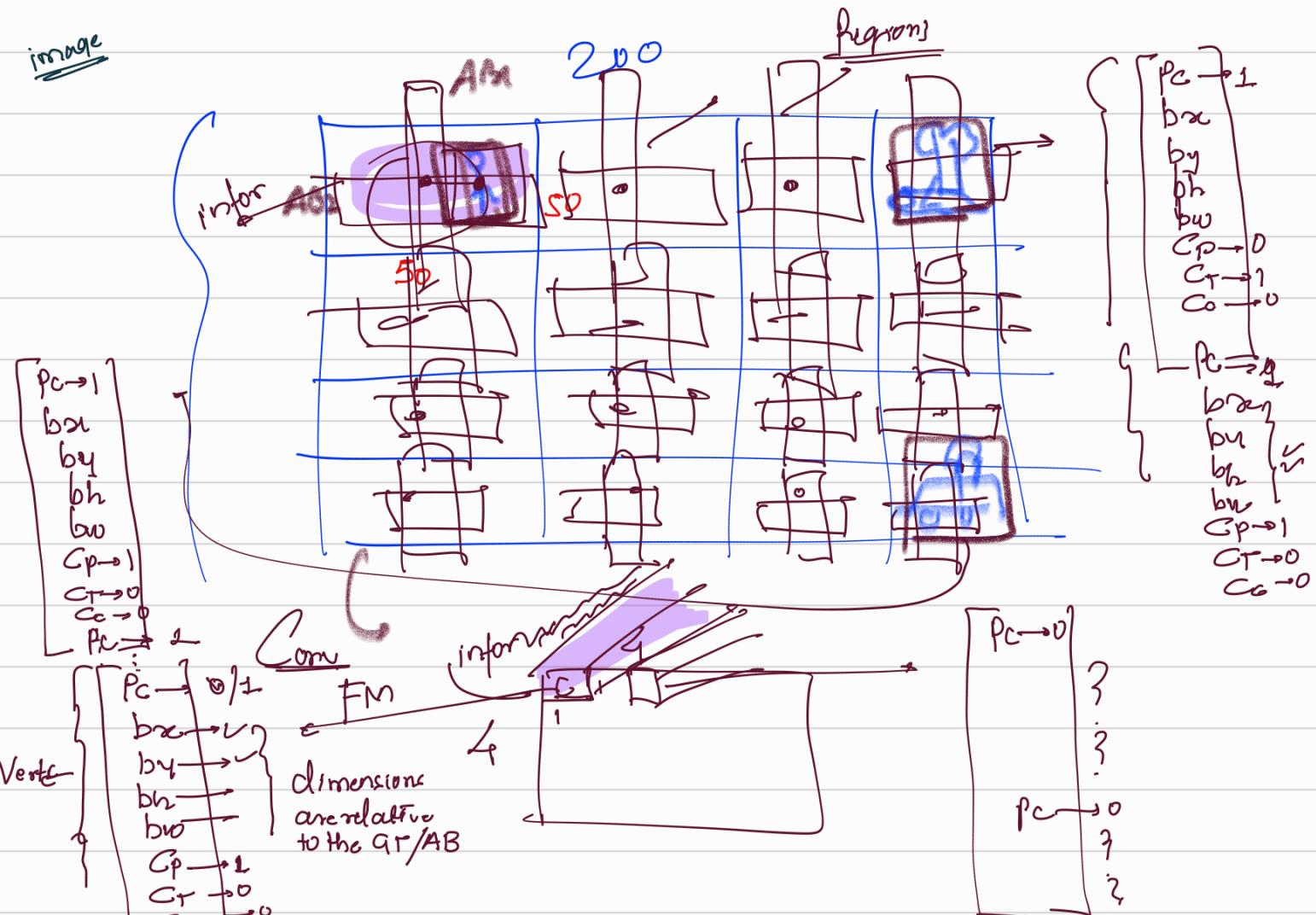
Which Anchor Box matches highest
IoU with the Pink Box

my feature map has to know that in Region 1 FM

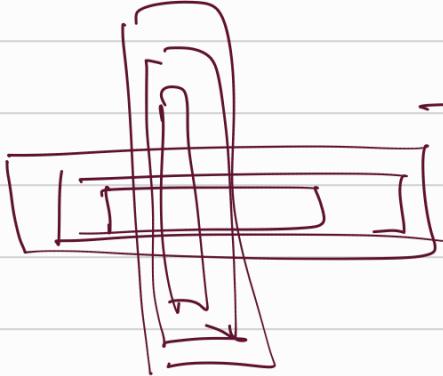
- we have a Person (Vertical object)
- Region 2 we do not have anything
- Region 3 does not have anything
- Region 4 has a tree (Vertical object) and a person (Horizontal)

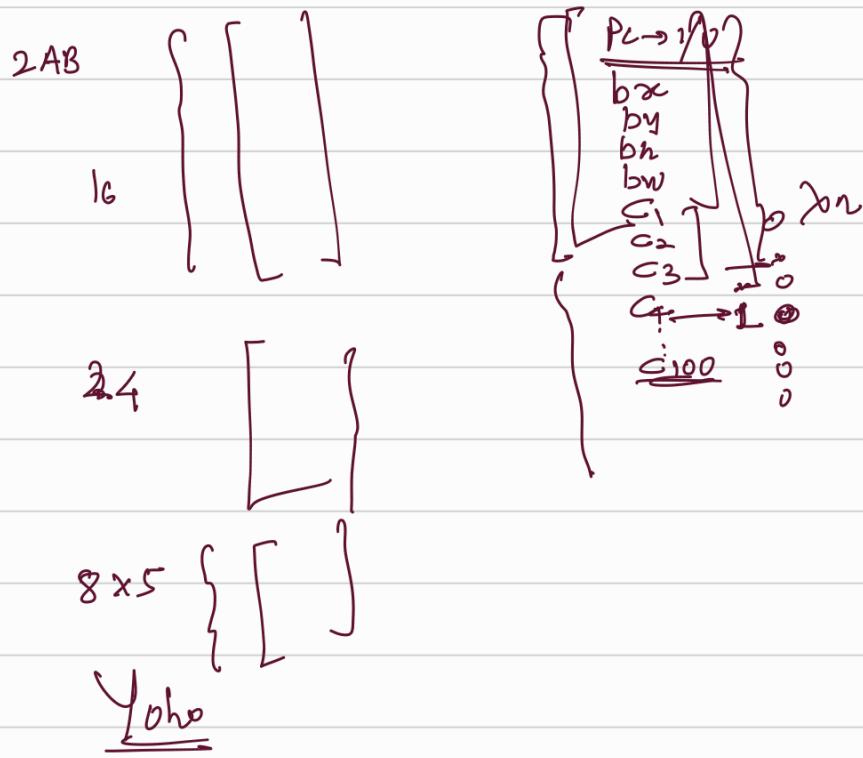


image



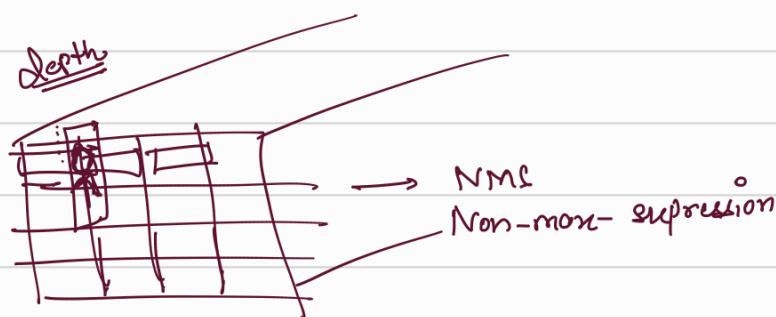
→ Number of anchor boxes
are decided.





Inference Time

image → Yolo model (Conv, MP, Conv, MP) →

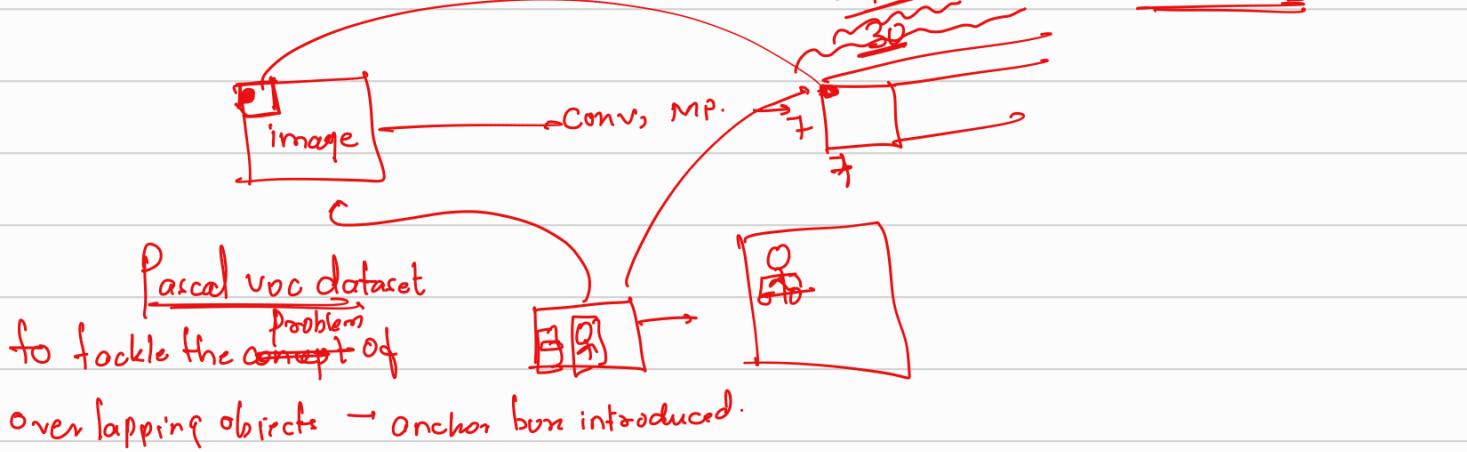


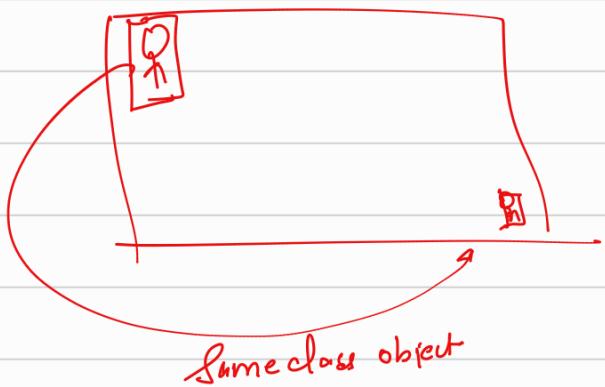
Anchor Boxes → how many classes
Kmeans clustering

2D class probabilities [20]
→ 2 bounding boxes

Yolo-v1 model → RPN + Detection network

2 Bounding Boxes x_1, y_1, w_1, h_1 x_2, y_2, w_2, h_2
 P_{c_1} P_{c_2}





Yolo-v1 could not match the capability to predict:

- ② overlapping objects
- ① same objects different size

Yolo-v2

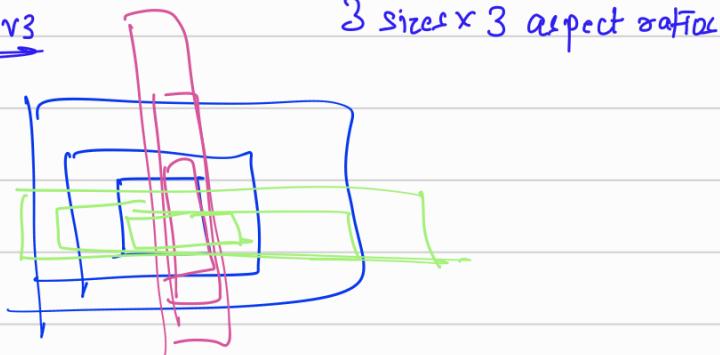
- ① Batch Normalization
- ② High Resolution images \rightarrow Yolo-v1 \rightarrow 24×224
 48×48 and
- ③ Anchor Boxes \rightarrow we define them by doing K-means
clustering on existing ground truths
of my training data
- ④ Feature map \rightarrow 13×13
- ⑤ 19 convolution layers and 5 more Pooling layers

V2 faster, accurate

$$\text{Yolo-v3} = 3 \times 3$$

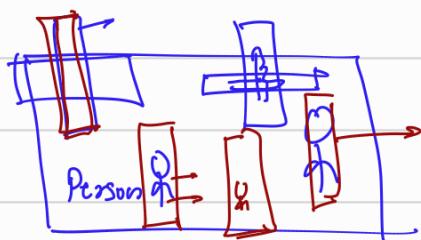


Yolo-v3



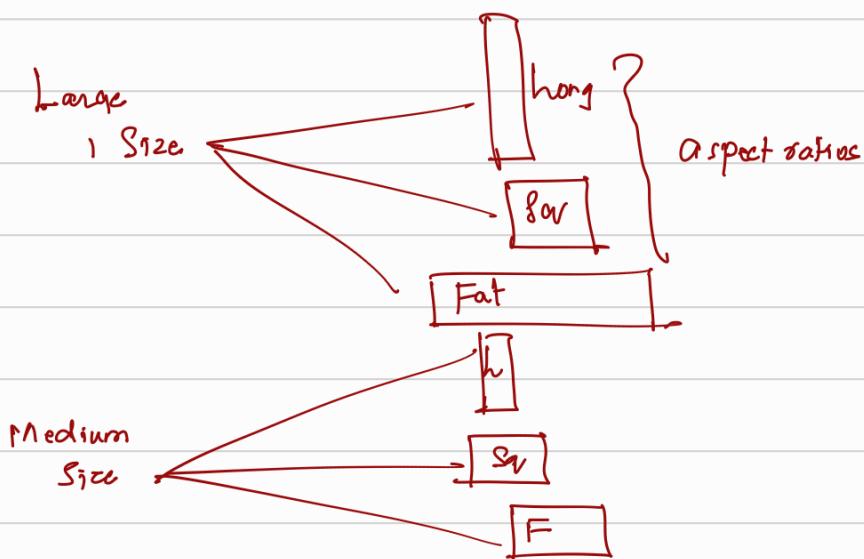
Anchor boxes → K-means Clustering

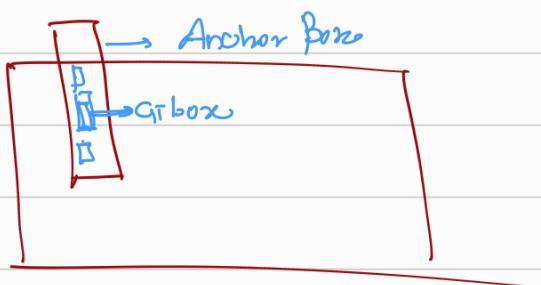
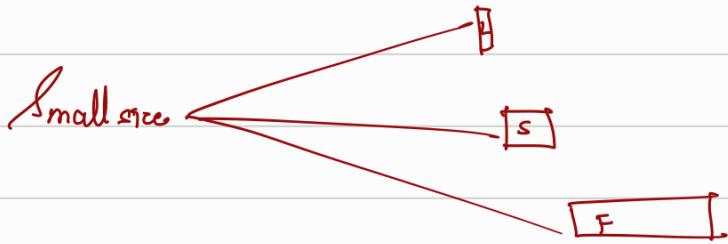
YOLO v3 model



Only one size for anchor boxes \Rightarrow

Yolo v3 → Sizes and aspect ratios for the anchor box

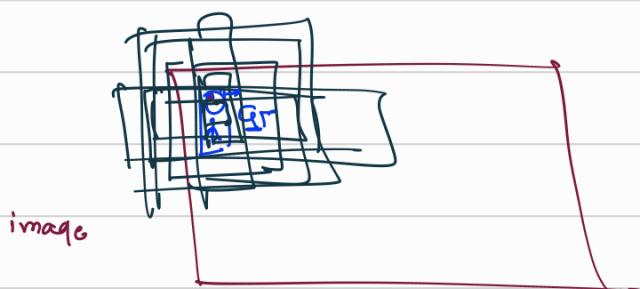
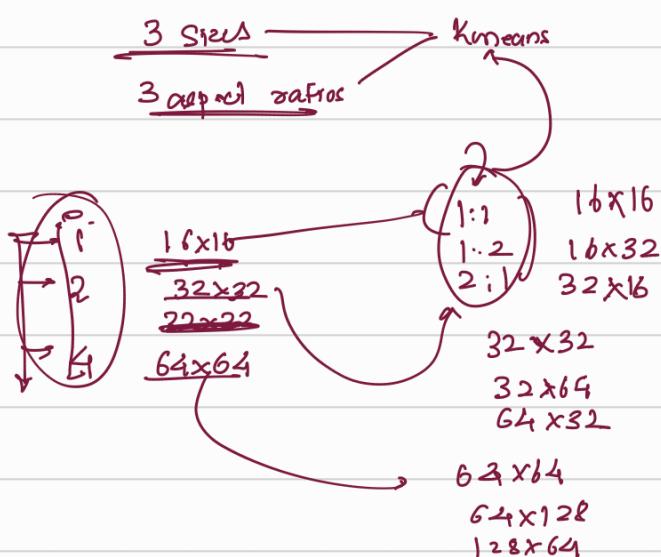




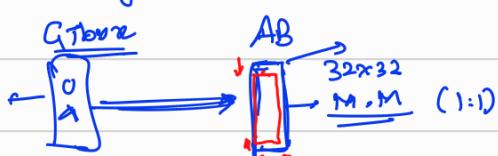
Yolov3 model \rightarrow 3×3 Anchor Boxes

3 sizes, 3 aspect ratios

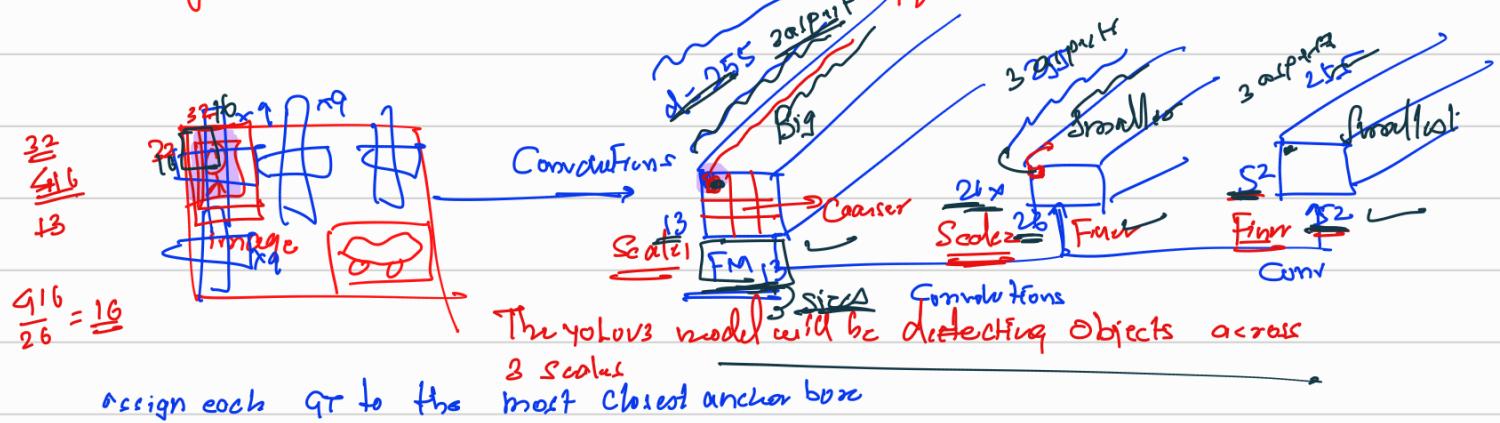
Every grid cell will have 9 anchor boxes to fit in the ground truth.



Assign the GT box to the most closest anchor box.

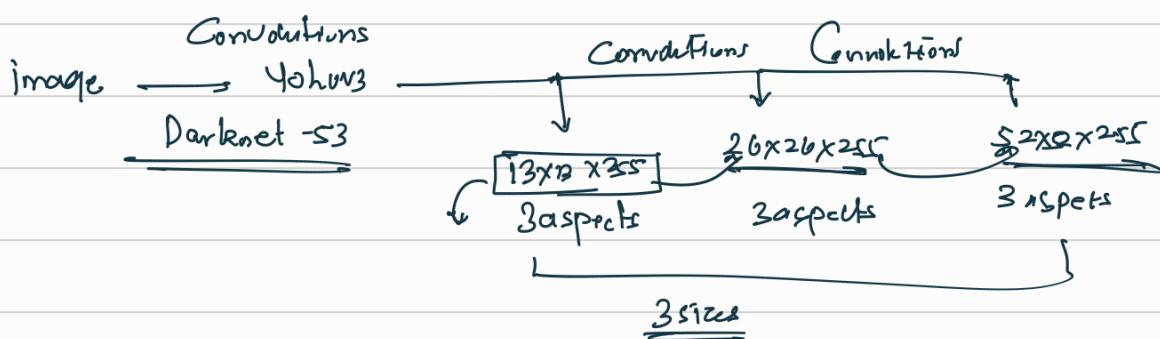
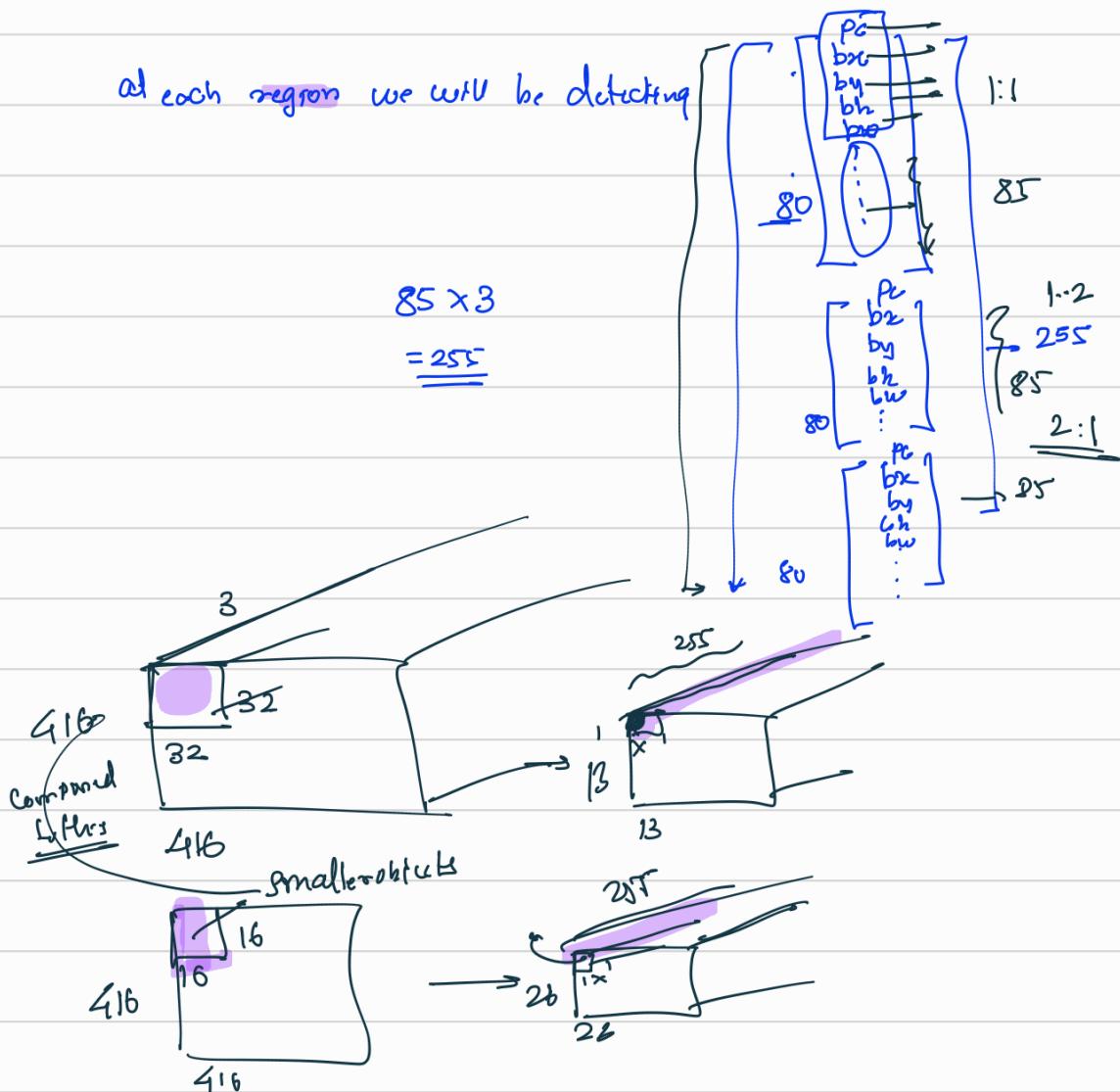


train
during the learning we have to detect the above difference



YOLOv3 was trained on a dataset which has 80 classes.

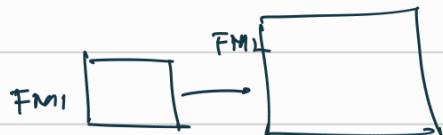
at each region we will be detecting



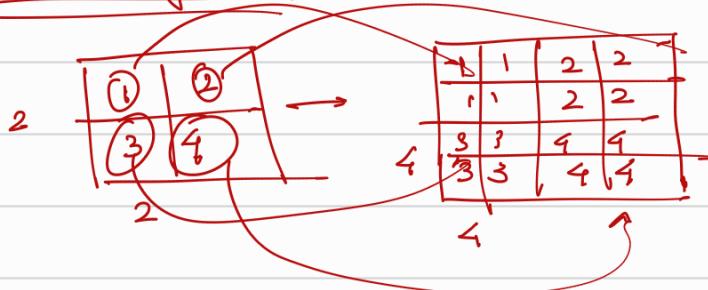
→ Convolutions → downsample

FM → Convolutions → Bigger FM

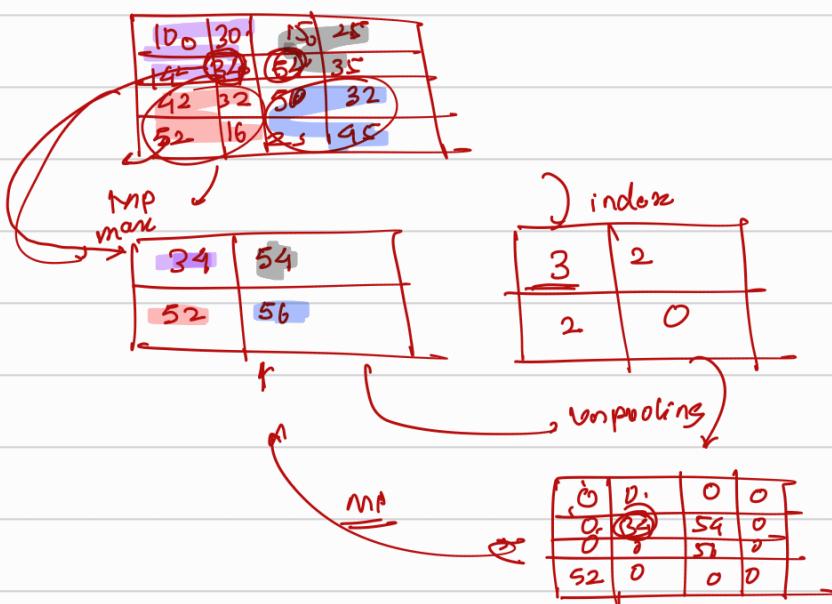
Upsampling (Transpose convolution)



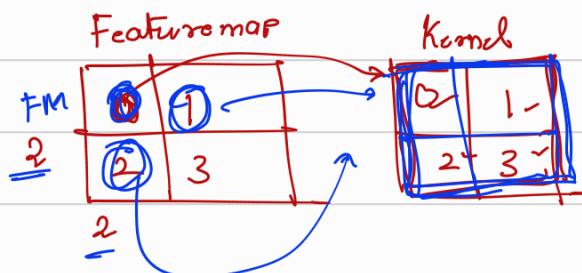
① Nearest neighbour



② Unpooling



③ Transpose convolution

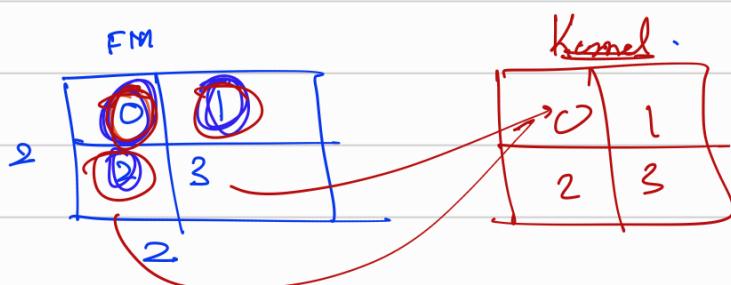


$s=1$

0	0	1
0	0+2	3
4	6	1

S=2

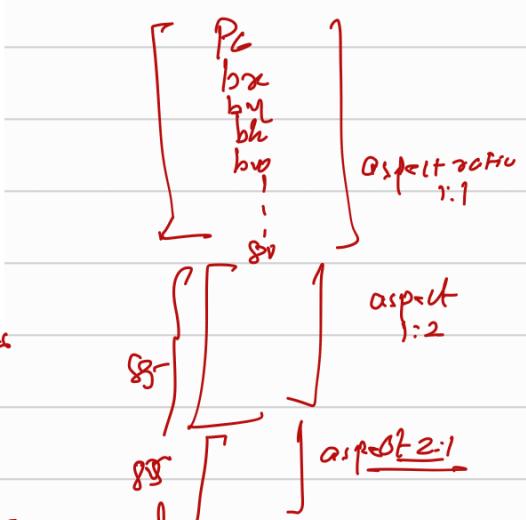
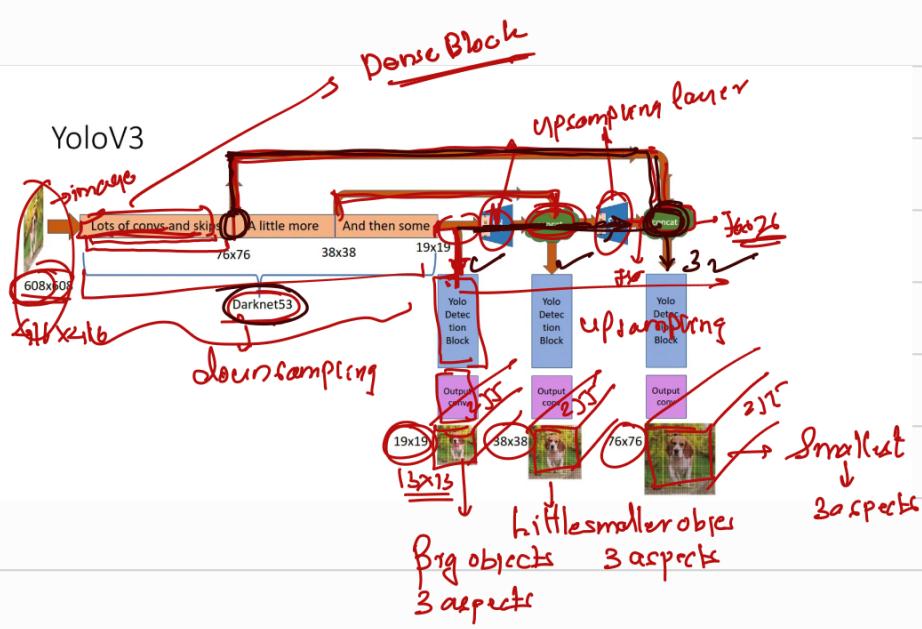
FM	Kernel		
0	0	0	1
0	0	2	3
0	2	0	3
4	6	8	9



Output map:

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

Transpose convolution is basically for upsampling your
feature map



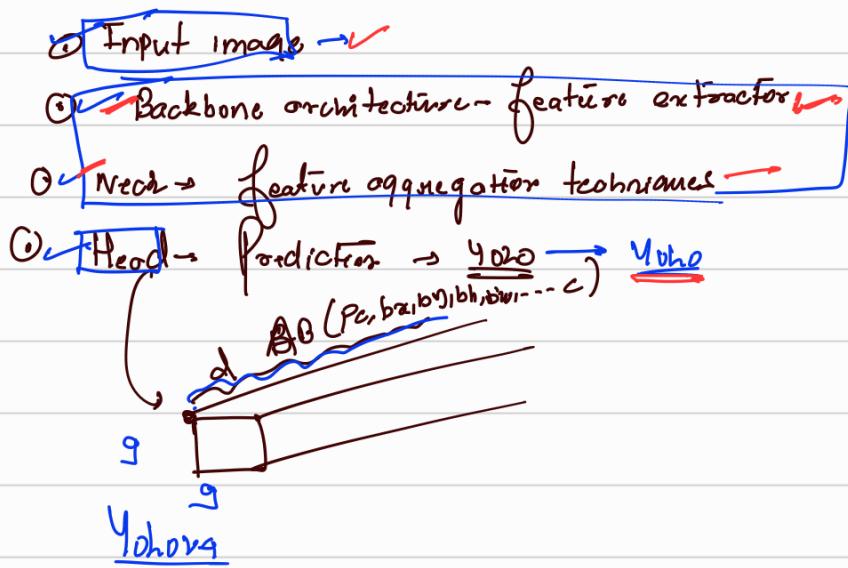
YoloV3 architecture is nothing but a series of convolutions applied along with certain skip connections to generate output at

19×19 (big object) \rightarrow upsample \rightarrow conv $\rightarrow 38 \times 38 \rightarrow$ upsample $\rightarrow 76 \times 76$

Yolo-v4 ... v7

Yolov4 model

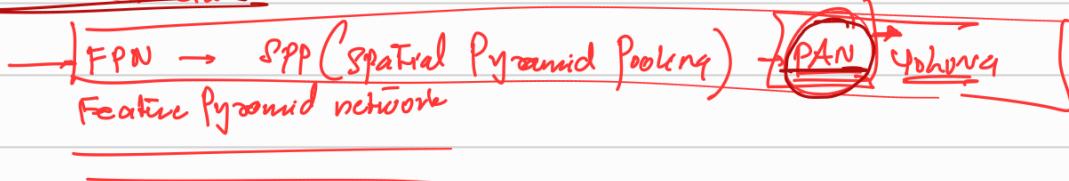
Building Components of an object detection technique:



Backbone architecture

- ① DenseNet, CspNet, CSPDarknet-53 → Yolov4

Neck architecture:



Dense Net (Backbone architecture)

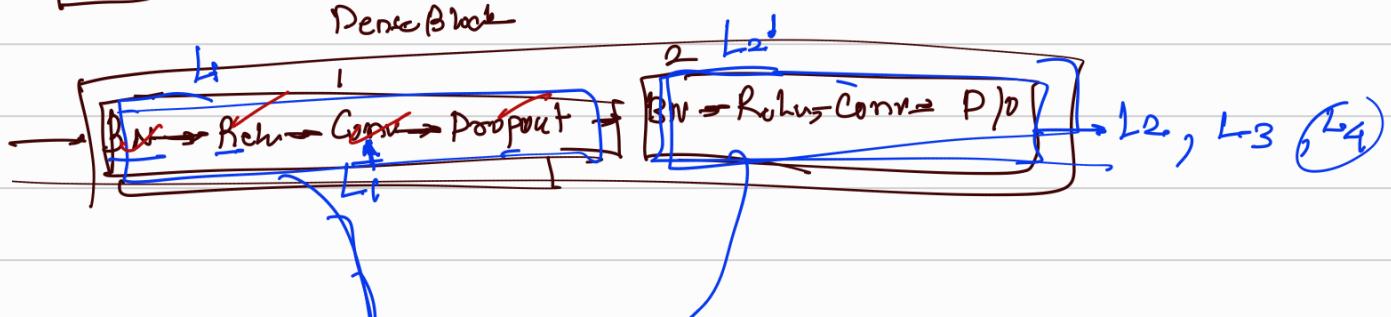


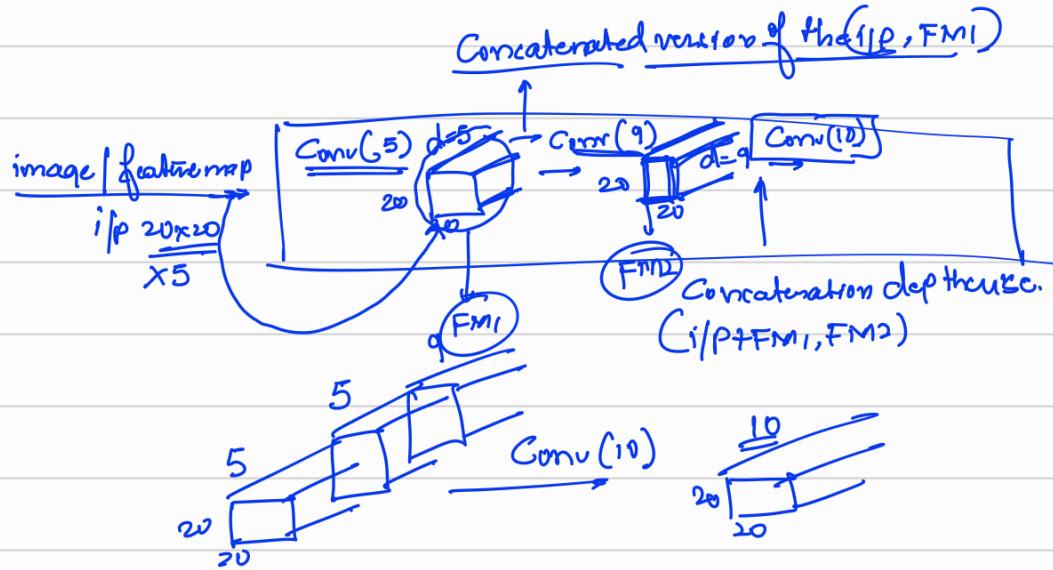
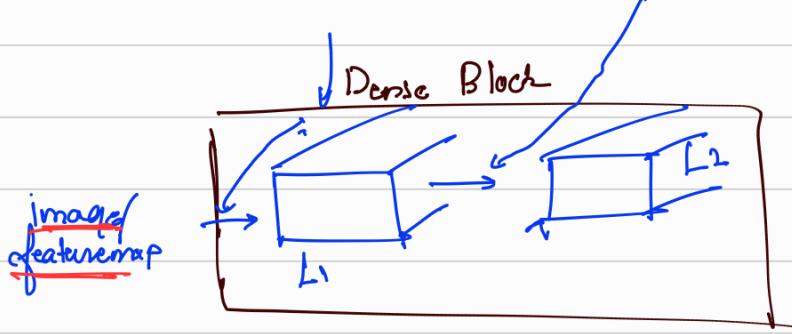
→ Transition block

is a collection of \approx set of layers

1 set of layer → Batch Norm, ReLU, Conv (3x3), Dropout

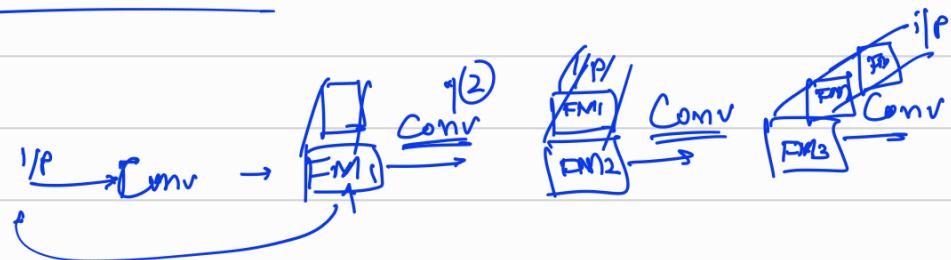
Dense Block





The dense block every turn of apply a [conv]
we apply always to the concatenated version

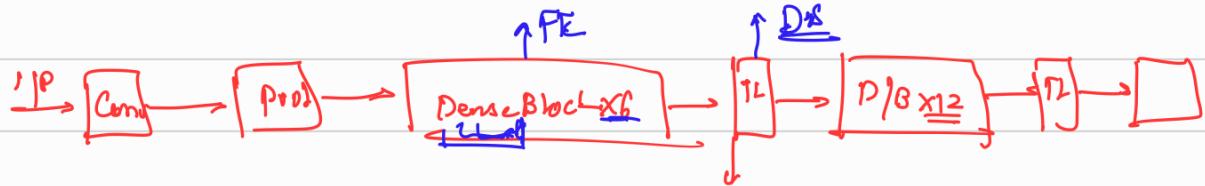
→ ^{ith}
[FMI-1, FMI-2, ..., i/A]



this is a very computationally expensive phenomenon

Cross stage Partial network

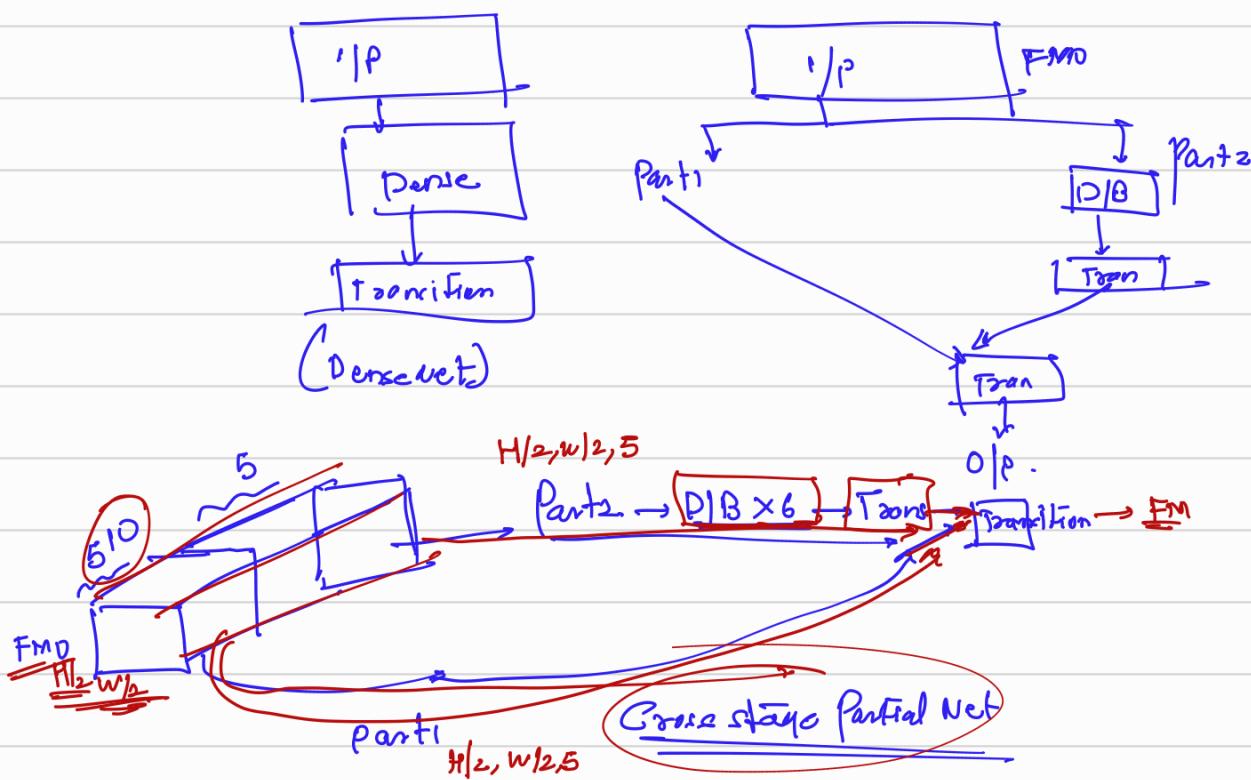
Transition Block



downsampling

$\text{TL} \rightarrow \text{BatchNorm, ReLU, Conv}(1 \times 1), \Delta/0, \text{Pool}(2 \times 2)$

Cross Stage Partial Network



Yohors

Mesh activation \rightarrow Yohors

Backbone Yohors \longrightarrow image $\longrightarrow \dots \text{dIB} \dots \text{dIS} \dots \text{Yohos}$

Backbone Yohors \longrightarrow image $\longrightarrow \text{dIB CSP} \xrightarrow{\text{L}} \text{CSP} \rightarrow \text{CSP}$
Mesh activation