

YOLO (You Only Look Once)

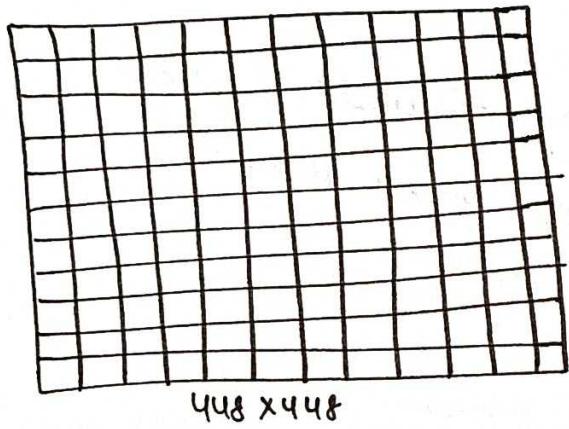
YOLO → Real time of detection

↳ one shot → regression

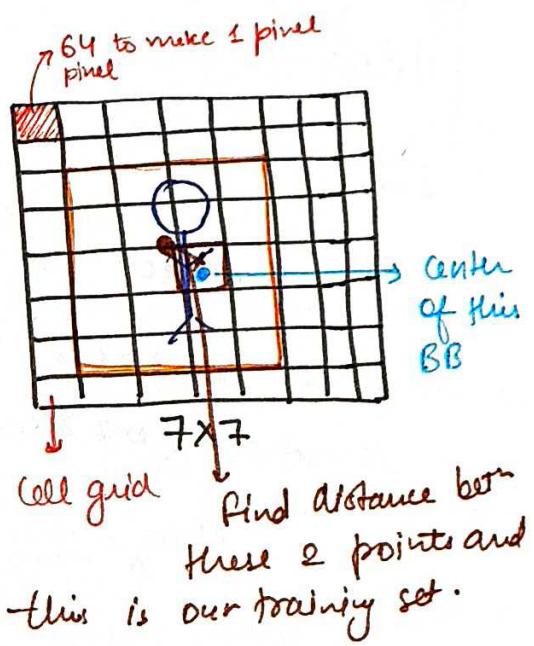
↳ losses (Obj, BB, classif)

Step-1:

Resize the image to 448×448



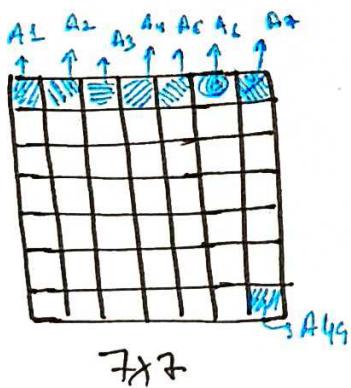
convert



$$\frac{448}{64} = 7$$

Image $\rightarrow 448 \times 448$

↓
convert into 7×7



$\frac{7 \times 7}{64}$

A_1

A_2

A_3

A_4

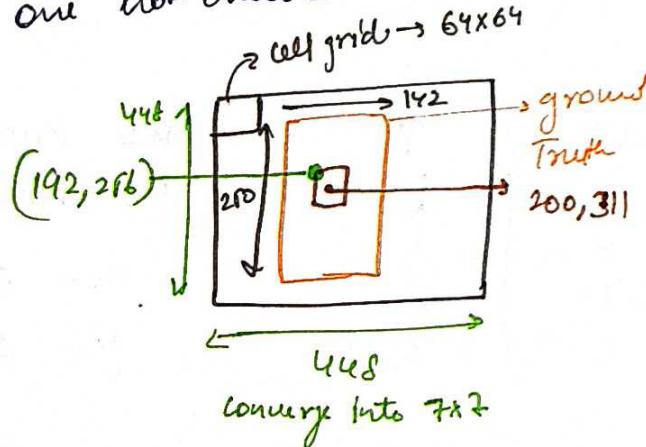
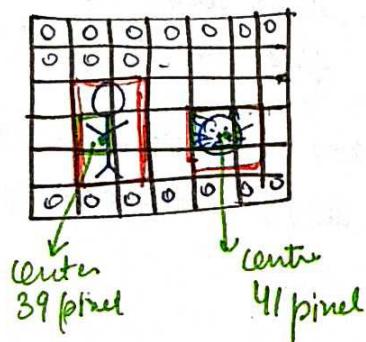
⋮

A_{49}

	Δx	Δy	Δw	Δh	c_i (object score)	\vec{P} (20 class)
A ₁	0	0	0	0	0	(0000...0) _{16k}
A ₂	0	0	0	0	0	0
⋮						
A ₃₉	0.31	0.2	0.1	0.5	1	[0001000] Person
A ₄₀	0.4	0.8	0.6	0.2	1	[0100000] Cat
⋮					0	[0000000]
A ₄₁	0	0	0	0	0	

Chart ↗

\vec{P} (20 class) = Pascal data with 20 classes and all classes one hot encoded



$$\Delta w_z = 142/448 = w_1$$

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

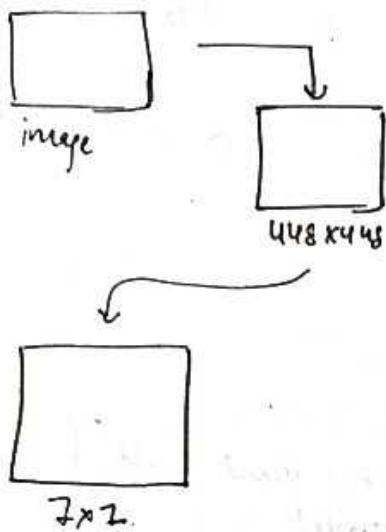
$$\Delta h = 250/448 = h_1$$

$$\Delta y = \frac{311 - 286}{64} = y_1$$

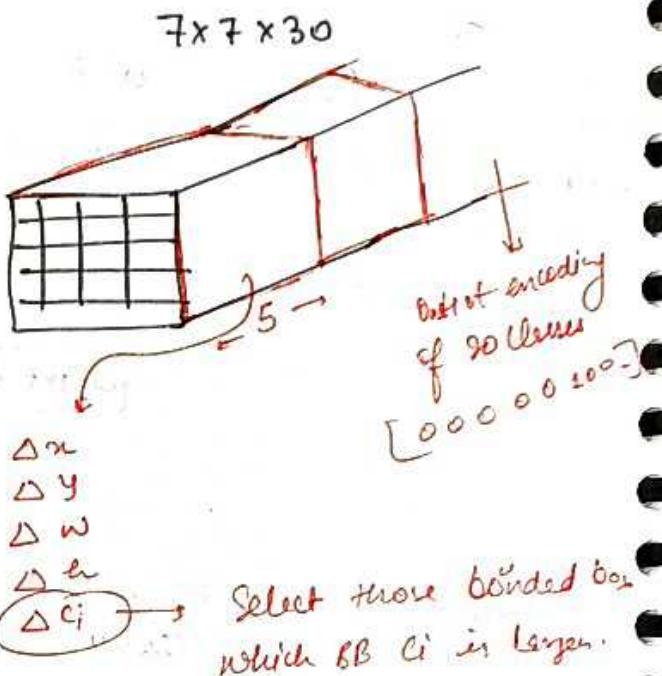
* These values are in upper chart.

Architecture

Training



Prediction



$$\begin{matrix} \Delta x \\ \Delta y \\ \Delta w \\ \Delta h \end{matrix} \left\{ \begin{matrix} x \\ y \\ w \\ h \end{matrix} \right.$$

we know Δx , Δy , Δw and Δh
so, we can easily find out x, y, w, h

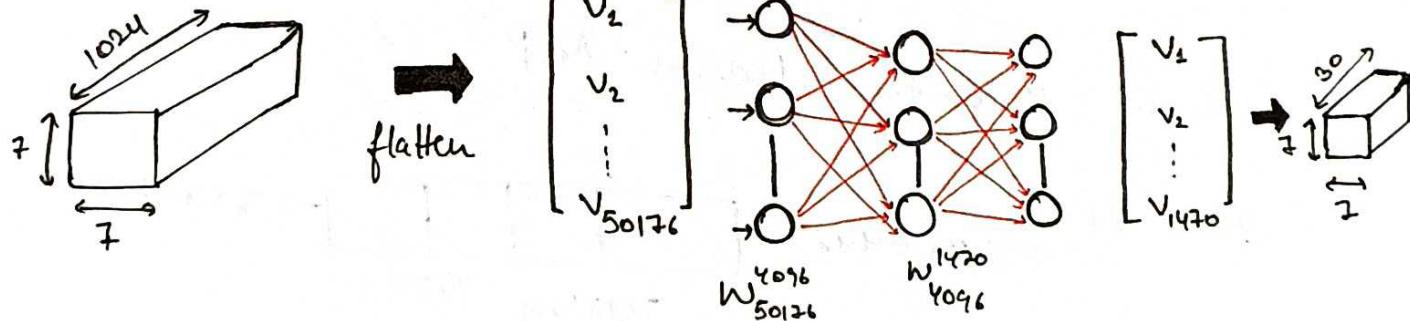
$$x = \Delta x + 64 + (192 - 200), \quad w = \Delta w \times 448 \\ y = \Delta y + 64 + (256 - 311), \quad h = \Delta h \times 448$$

$(x, y, h, w) \rightarrow$ Grid cell
Top left corner

Structure

Type	Size	Filter	Stride	Output
Conv.	7x7x3	64	2	224x224x64
max pool	2x2			112x112x64
Conv.	3x3x64	192	1	112x112x192
max pool	2x2			56x56x192

Conv.	$1 \times 1 \times 192$	128	1	$56 \times 56 \times 128$	(117)
Conv.	$3 \times 3 \times 128$	256	1	$56 \times 56 \times 256$	
Conv.	$1 \times 1 \times 256$	256	1	$56 \times 56 \times 256$	
Conv.	$3 \times 3 \times 256$	512	1	$56 \times 56 \times 512$	
maxpool	2×2			$28 \times 28 \times 512$	
4x [Conv.]	$1 \times 1 \times 512$	256	1	$28 \times 28 \times 256$	
	$3 \times 3 \times 256$	512	1	$28 \times 28 \times 512$	
Conv.	$1 \times 1 \times 512$	512	1	$28 \times 28 \times 512$	
Conv.	$3 \times 3 \times 512$	1024	1	$28 \times 28 \times 1024$	
maxpool	2×2			$14 \times 14 \times 1024$	
2x [Conv.]	$1 \times 1 \times 1024$	512	1	$14 \times 14 \times 512$	
	$3 \times 3 \times 512$	1024	1	$14 \times 14 \times 1024$	
Conv.	$3 \times 3 \times 1024$	1024	1	$14 \times 14 \times 1024$	
Conv.	$3 \times 3 \times 1024$	1024	2	$14 \times 14 \times 1024$	
Conv.	$3 \times 3 \times 1024$	1024	1	$14 \times 14 \times 1024$	
Conv.	$3 \times 3 \times 1024$	1024	1	$14 \times 14 \times 1024$	



YOLO V2

YOLO V2 → Accuracy

→ Increase No. of Objects

→ Increase Speed

Modification and Improvement in YOLO V2 over YOLO V1

1. Map grew from 64 to 78.

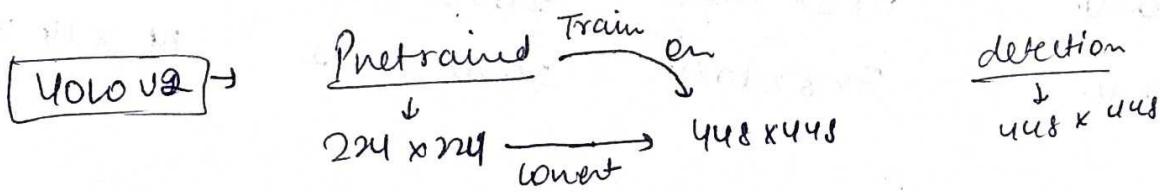
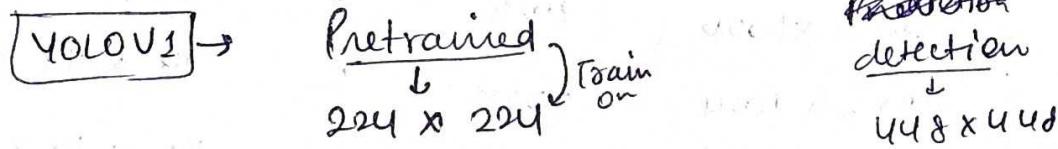
(A) Apply Batch Normalization to all CNN layers.

→ Improve map by 2%.

→ Solved overfitting

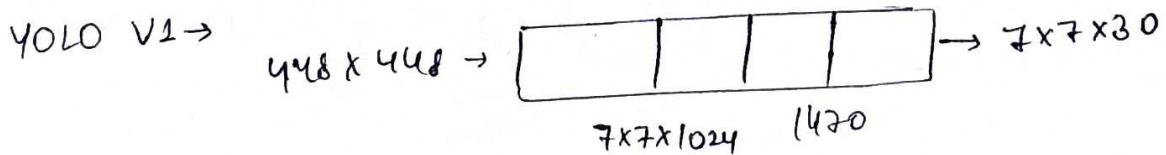
→ regularization effects

(B) Higher Resolution Classifier



→ Improve Map by 4%.

(C) High Resolution Feature Map



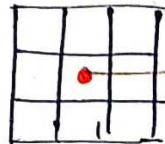
YOLOV2 →

$$448 \times 448 \quad | \quad \rightarrow 14 \times 14 \times 1024$$

remove
Max pool
from YOLOv1

but if we do
even x even x even
this is not good

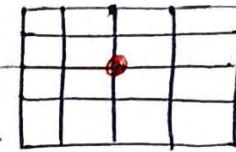
We can do oddxodd.



Never perfect
Center

Solution

This is centre
and we don't
want this centre



So, we use $\rightarrow 13 \times 13 \times 1024$

Type	Size	filters	stride	outputs
------	------	---------	--------	---------

Same as YOLOV1

4x	Conv. Conv.	$1 \times 1 \times 512$ $3 \times 3 \times 256$	256 512	1 1	$28 \times 28 \times 256$ $28 \times 28 \times 512$
	Conv. Conv.	$1 \times 1 \times 512$ $3 \times 3 \times 512$	512 1024	1 1	$28 \times 28 \times 512$ $28 \times 28 \times 1024$

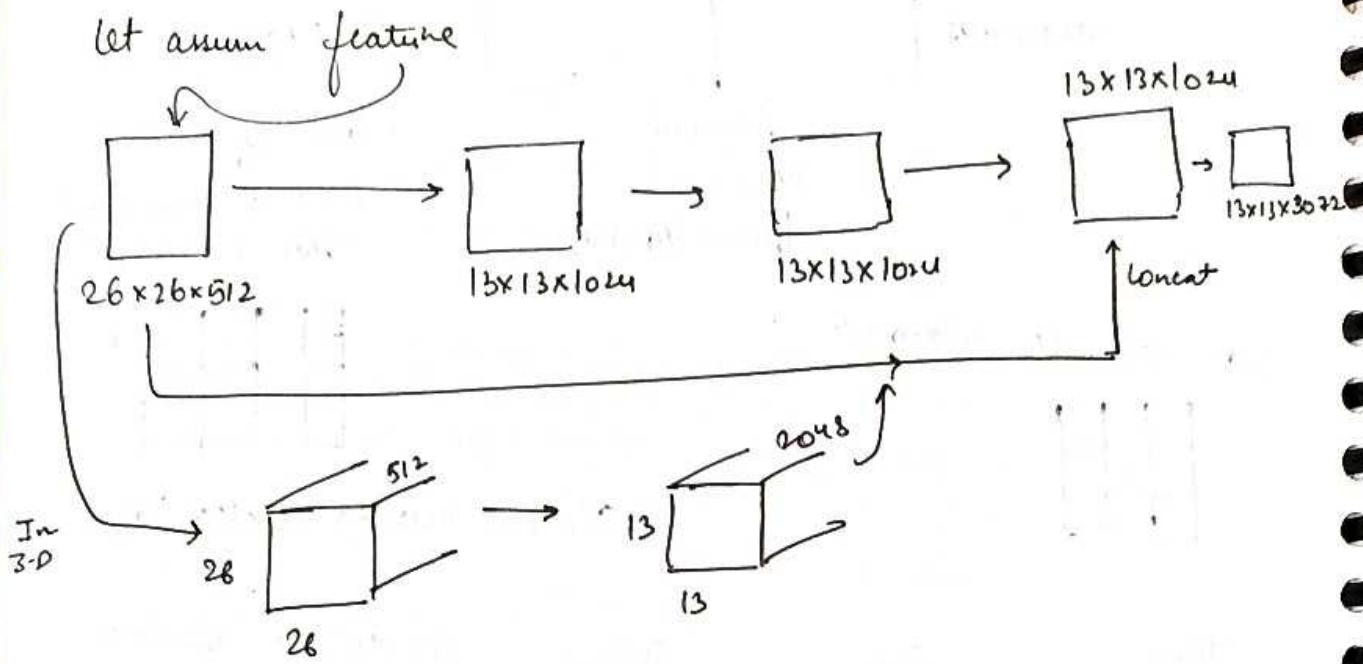
max pool

2x	conv. conv.	$1 \times 1 \times 1024$ $3 \times 3 \times 512$	512 1024	1 1	$28 \times 28 \times 512$ $28 \times 28 \times 1024$
	conv.	$3 \times 3 \times 1024$	1024	1	$28 \times 28 \times 1024$
	conv.	$3 \times 3 \times 1024$	1024	2	$14 \times 14 \times 1024$
	conv.	$3 \times 3 \times 1024$	1024	1	$14 \times 14 \times 1024$
	conv.	$3 \times 3 \times 1024$	1024	1	$14 \times 14 \times 1024$

Remove in YOLOV2

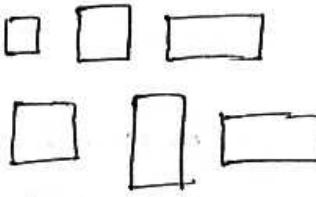
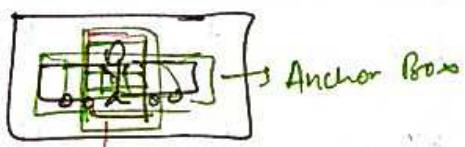
These values
are changed

④ Add Pass Through Layer



⑤ Bounding Box Improvement

→ Use Anchor Box →



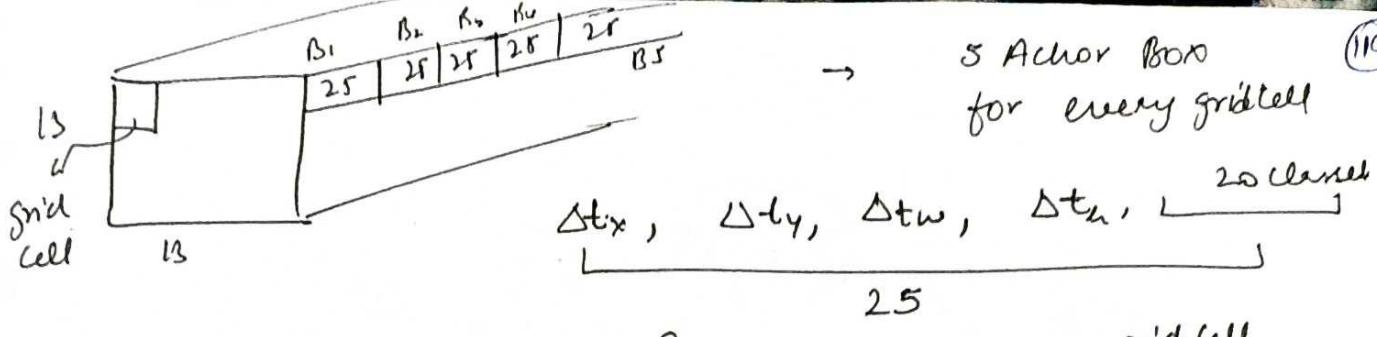
different
ratios
and scale

YOLOv1 detect person but hard to detect bus. So Anchor Box use to detect different shape as well.

→ In Yolo Faster RNN → we decide No. of Anchor Box

→ In Yolo v2 → K-Mean Cluster decide different shape Anchor Box.

5 Diff Modification which improve the model drastically from map of 64x to 781.



5 Anchor Boxes
for every gridcell

(119)

$\Delta t_x, \Delta t_y, \Delta t_w, \Delta t_h, \underbrace{25 \text{ classes}}$

Total Params = $125 \rightarrow 1 \text{ grid cell}$

YOLO V2 Params = $13 \times 13 \times 125$

YOLO V3

YOLO V3

using Anchor Box
3 prediction of different Scale

Scale 1

(Large)

Image size

~~Small~~ Anchor Box
Large

Scale 2

(Medium)

Image size

Medium Anchor Box

Scale 3

Small

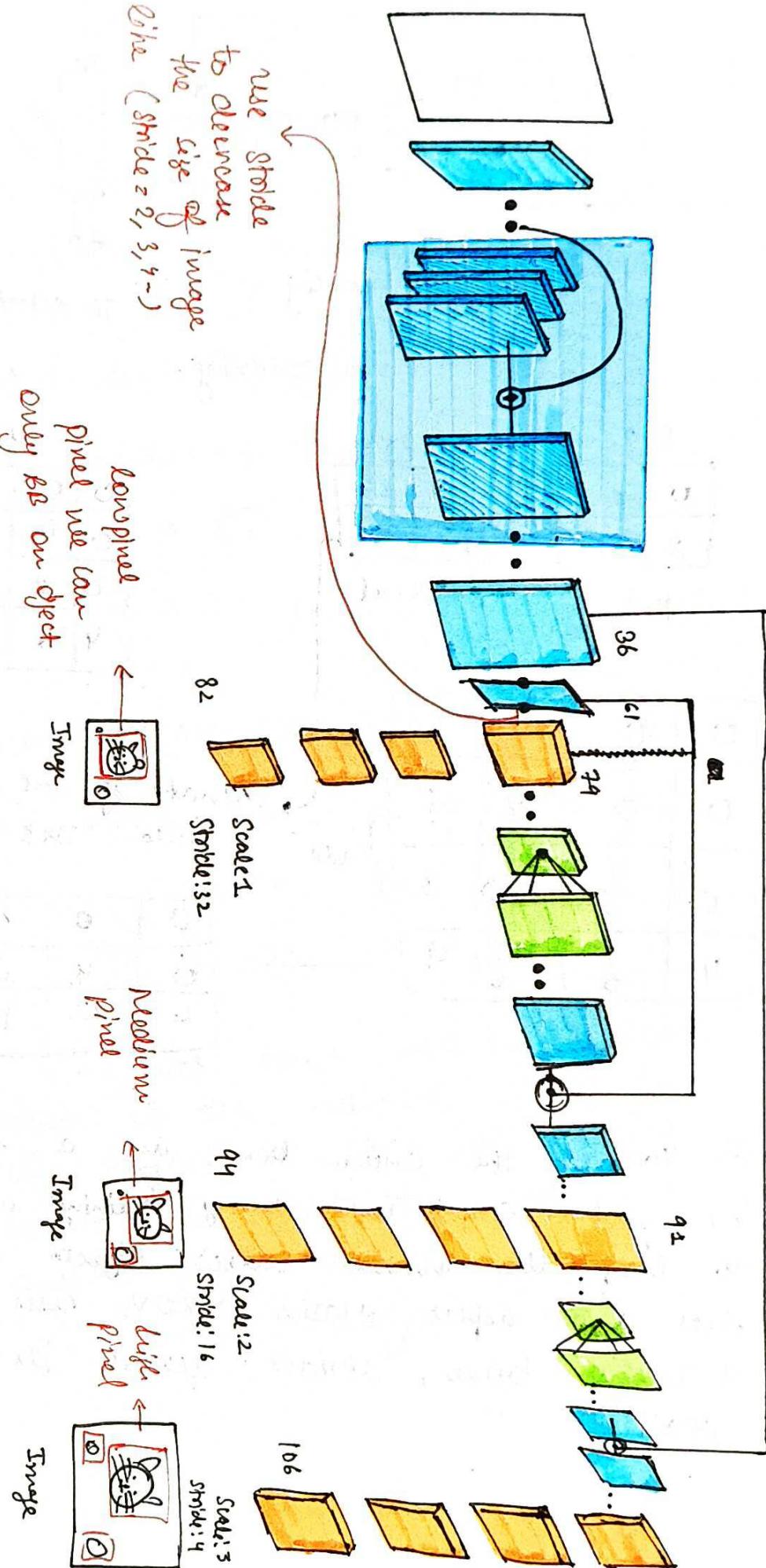
Image size

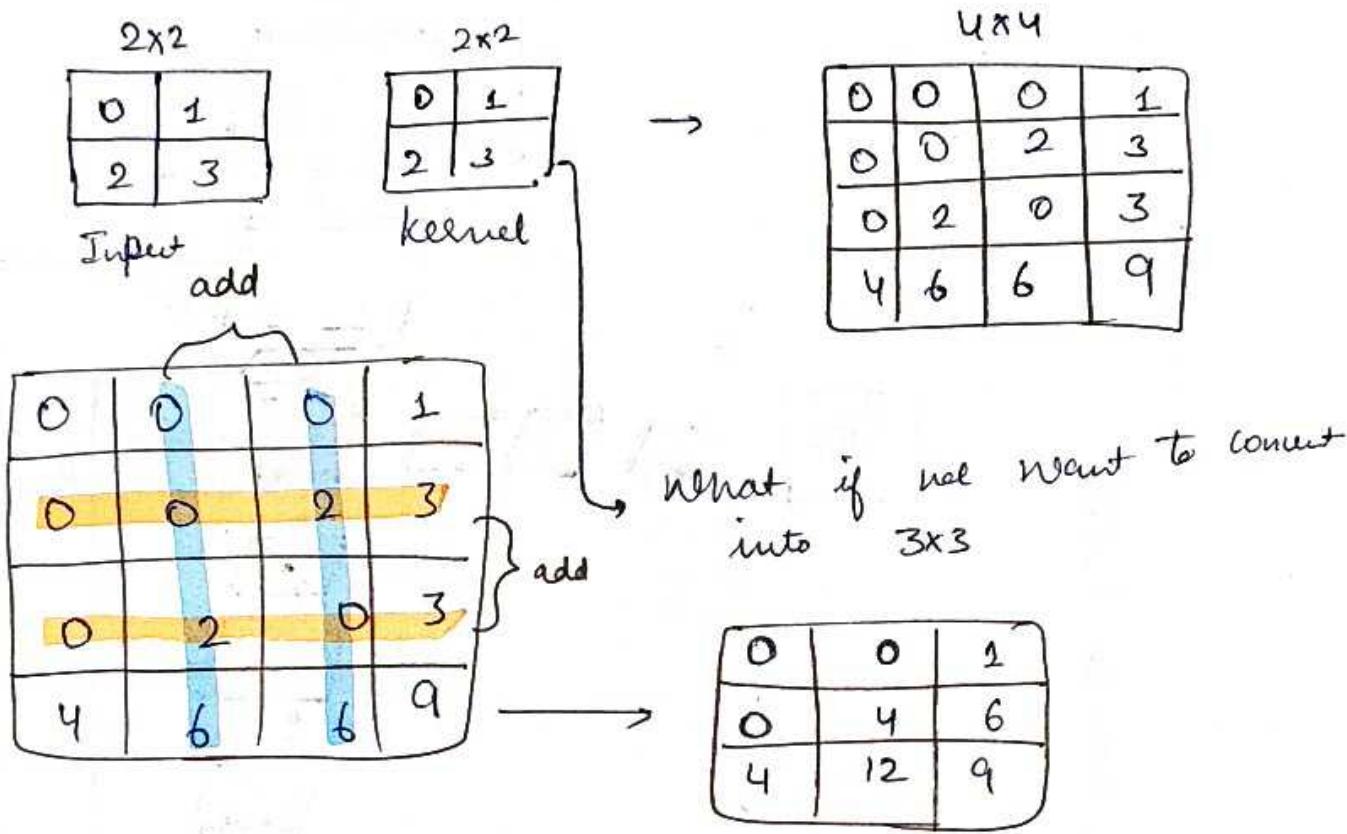
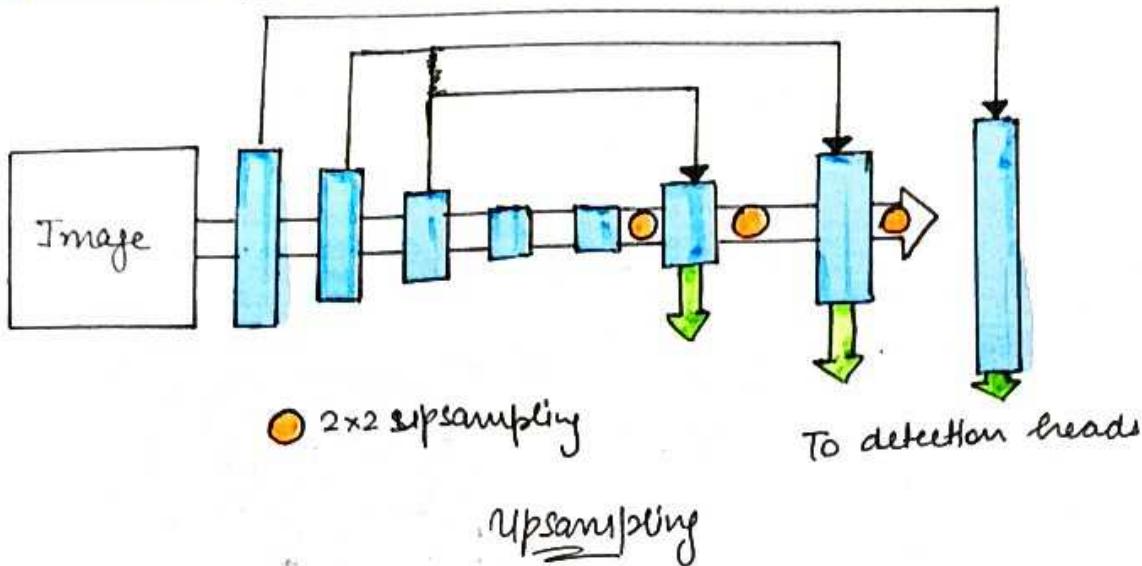
Small AB

In YOLO V2 architecture use → DarkNet (19) → 19 layer

In YOLO V3 architecture use → DarkNet (53) → 53 Layer

YOLO v3 network Architecture





In YOLOv3, the anchor boxes are a set of predefined bounding box shapes used during training and prediction to help the model detect objects of various sizes and aspect ratios. YOLOv3 uses a total of 9 anchor boxes, divided across its three different scales.

(121)

The anchor box sizes in YOLOv3 are as follows.

For the first scale (largest scale):

Anchor box 1: width = 10, height = 13

Anchor box 2: width = 16, height = 30

Anchor box 3: width = 33, height = 23

For the second scale (Medium scale)

Anchor box 1: width = 30, height = 61

" " 2: " = 62, " = 45

" " 3: " = 59, " = 119

For the Third Scale (Smallest scale):

Anchor box 1: width = 116, height = 90

" " 2: " = 156, " = 198

" " 3: " = 873, " = 326

These anchor box sizes were determined based on the analysis of the dataset used for training YOLOv3. They cover range of object commonly found in dataset and model learn to adjust the prediction on the basis of anchor box during training.

Object detection Rule
from any architecture

Image resizing
and
Pre-processing



Backbone
CNN
↓
Feature
creation

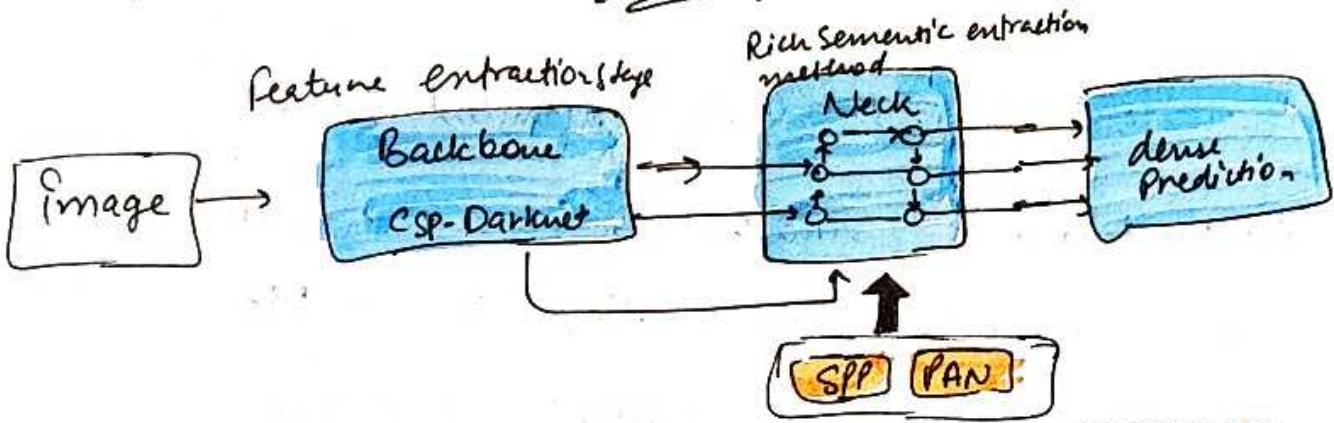
VGGNet/
DenseNet
etc.

Detector

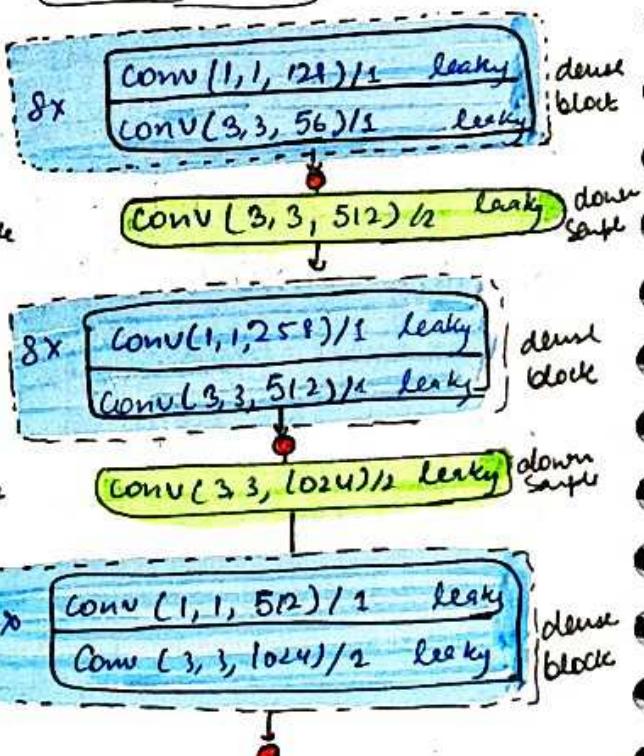
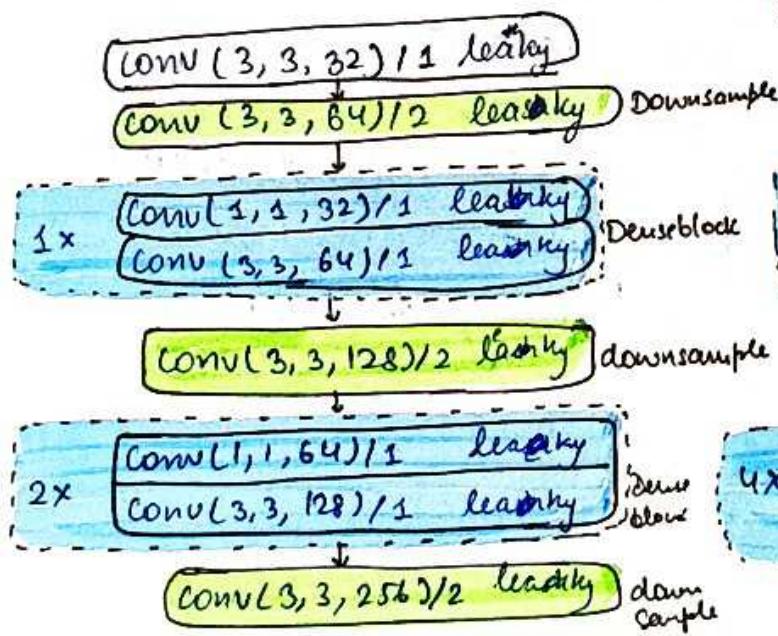


Neck

YOLO V4

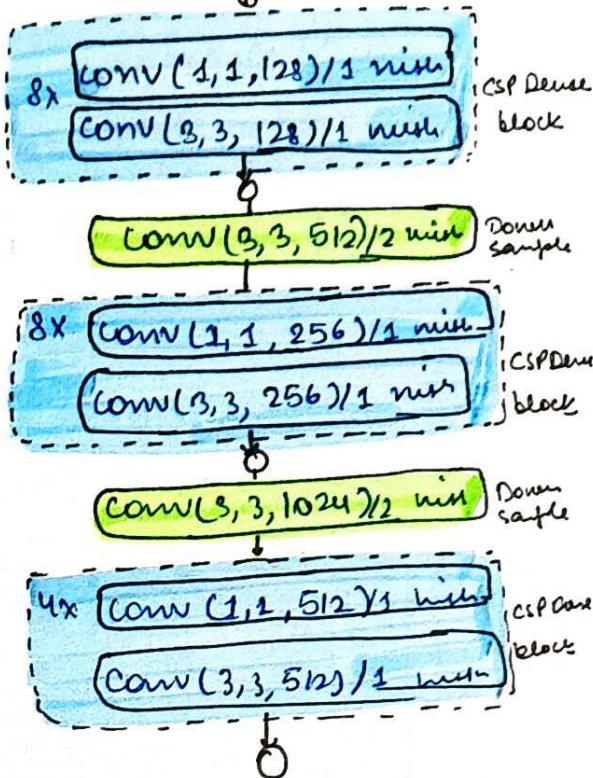
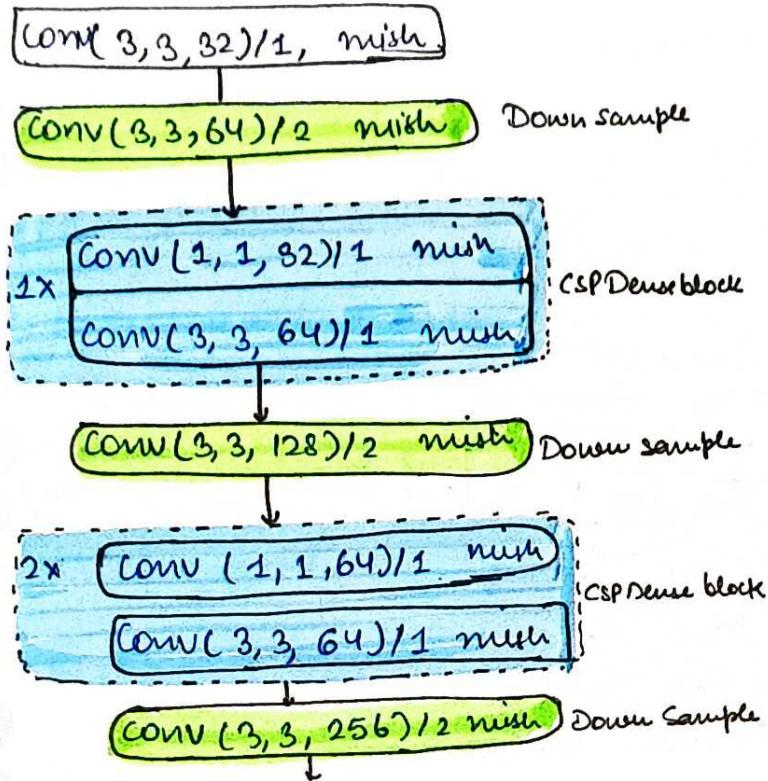


Darknet 53



CSP Darknet 53

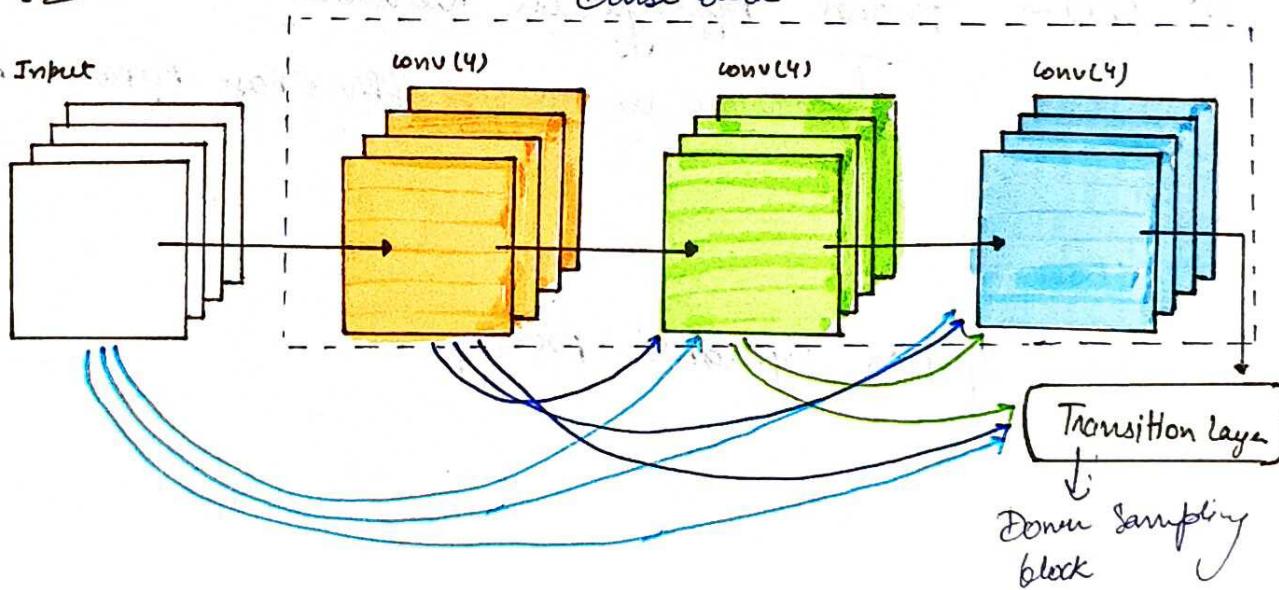
(122)



Difference betn Darknet and CSP Darknet

- ① Activation function change (LeakyRelu \rightarrow Mish)
- ② Dense block \rightarrow CSP - Dense block

Dense block

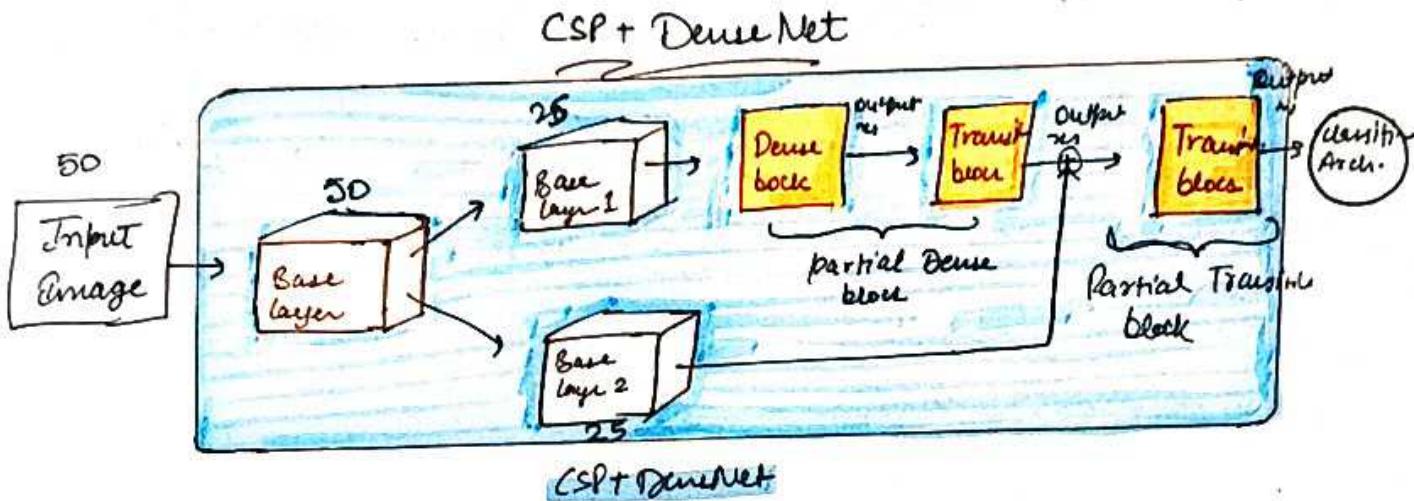


$$x_1 = w_1 * x_0$$

$$x_2 = w_2 * [x_0 \ x_1]$$

$$\vdots$$

$$x_k = w_k * [x_0 \ x_1 \ \dots \ x_{k-1}]$$



In YOLO V4 \rightarrow changes \rightarrow (i) Backbone CNN \rightarrow CSP + Darknet

\downarrow
Dense Block

+
Transition Block

(ii) Neck \rightarrow SPP + PAN

We use typically \rightarrow 2 feature Aggregation in YOLO V4

① PAN \rightarrow Partial Aggregation Network

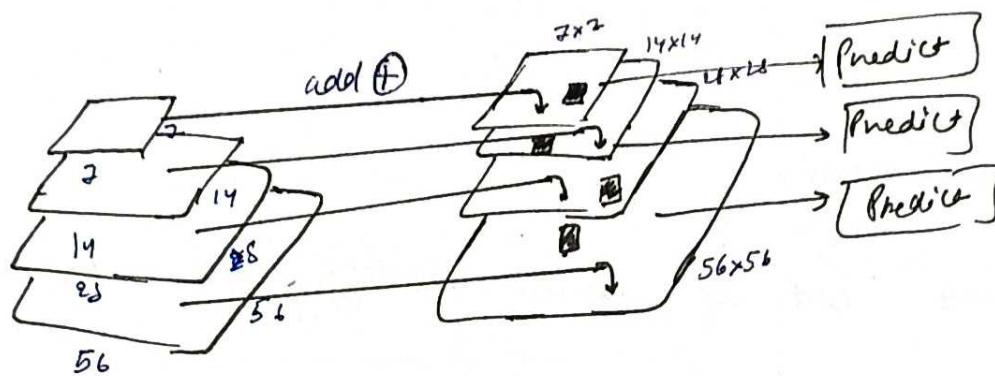
\hookrightarrow advance version of APN \rightarrow Fast Pyramid Network

② SPP

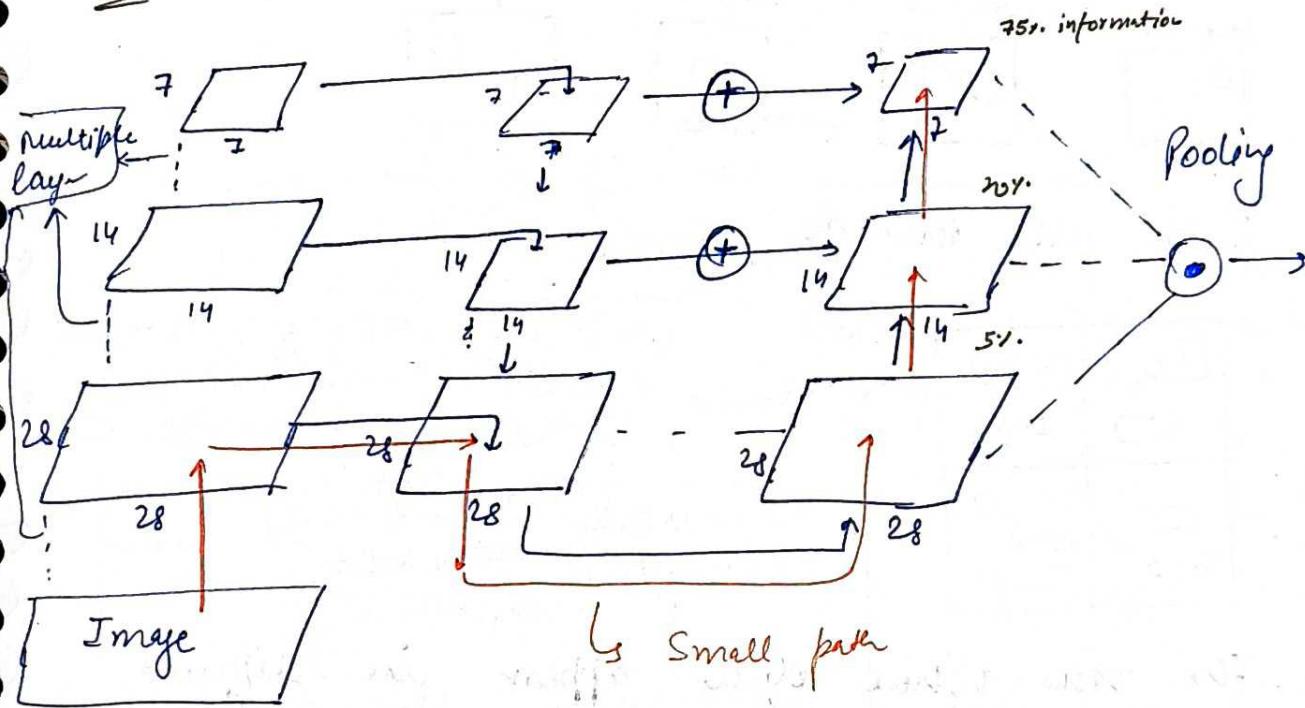
\hookrightarrow Spatial Pyramid Pooling

FPN → Feature Pyramid Network

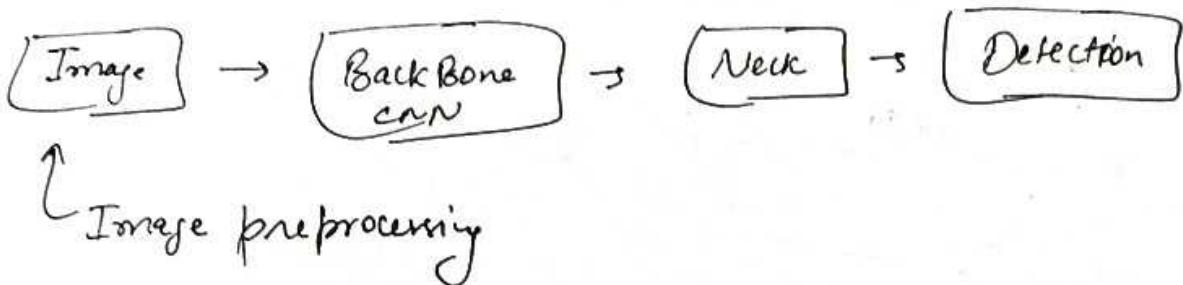
Used to detect small size object.



PAN → Path Aggregation Network



YOLO V5



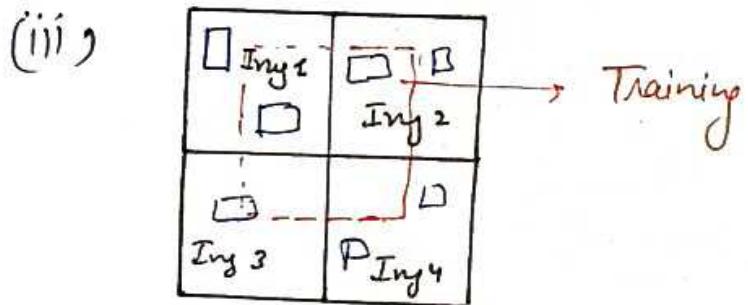
① MOSAIC → out of Content Identifying (Data Augmentation)

Steps:

(i) Take 4 image from the training set.



(ii) Resize all Image



Mosaic will help the model to recognise the objects that the model is not used to looking together.

All the cases where objects appear in different content.

MOSAIC used only Training Time.

→ Bag of Freebies

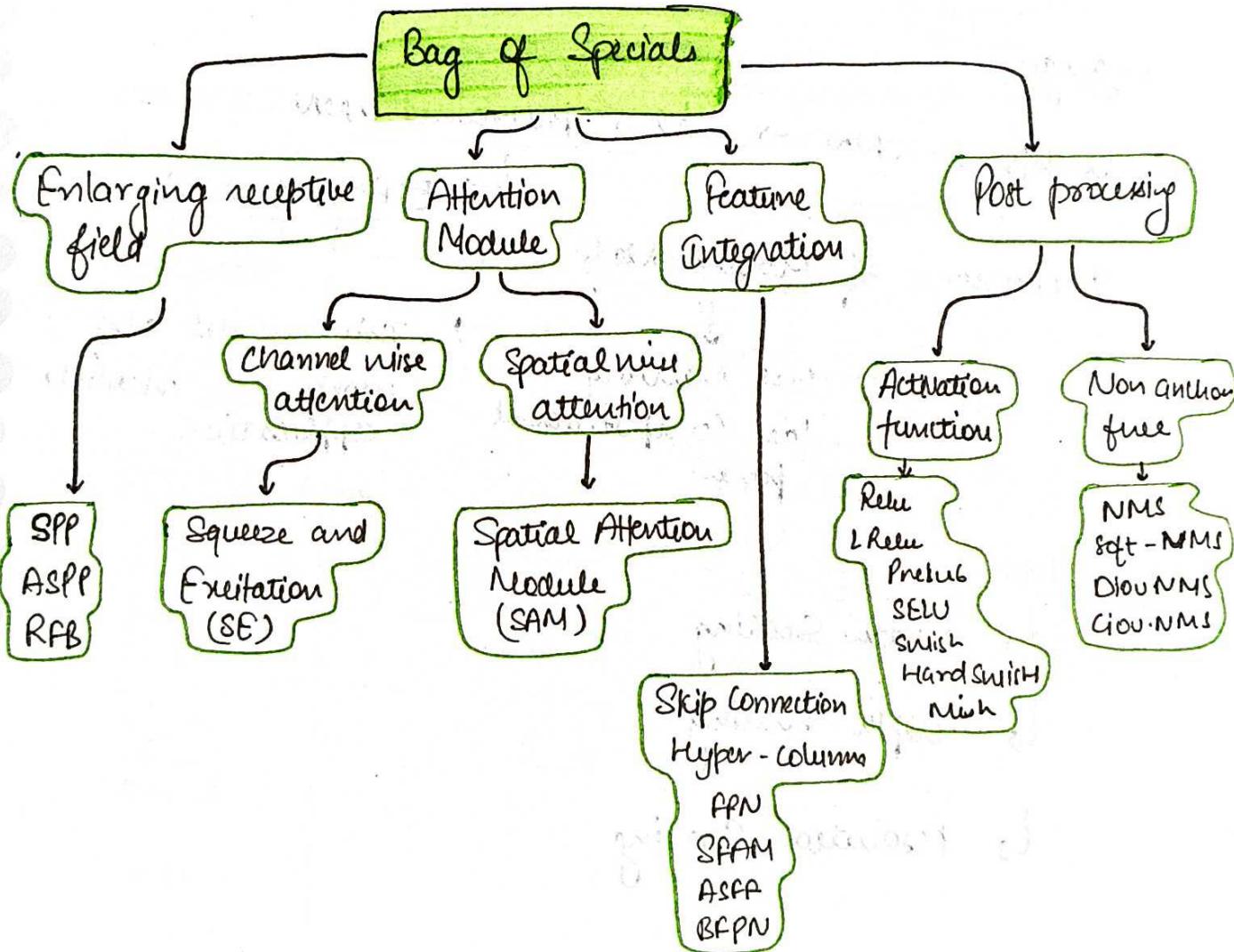
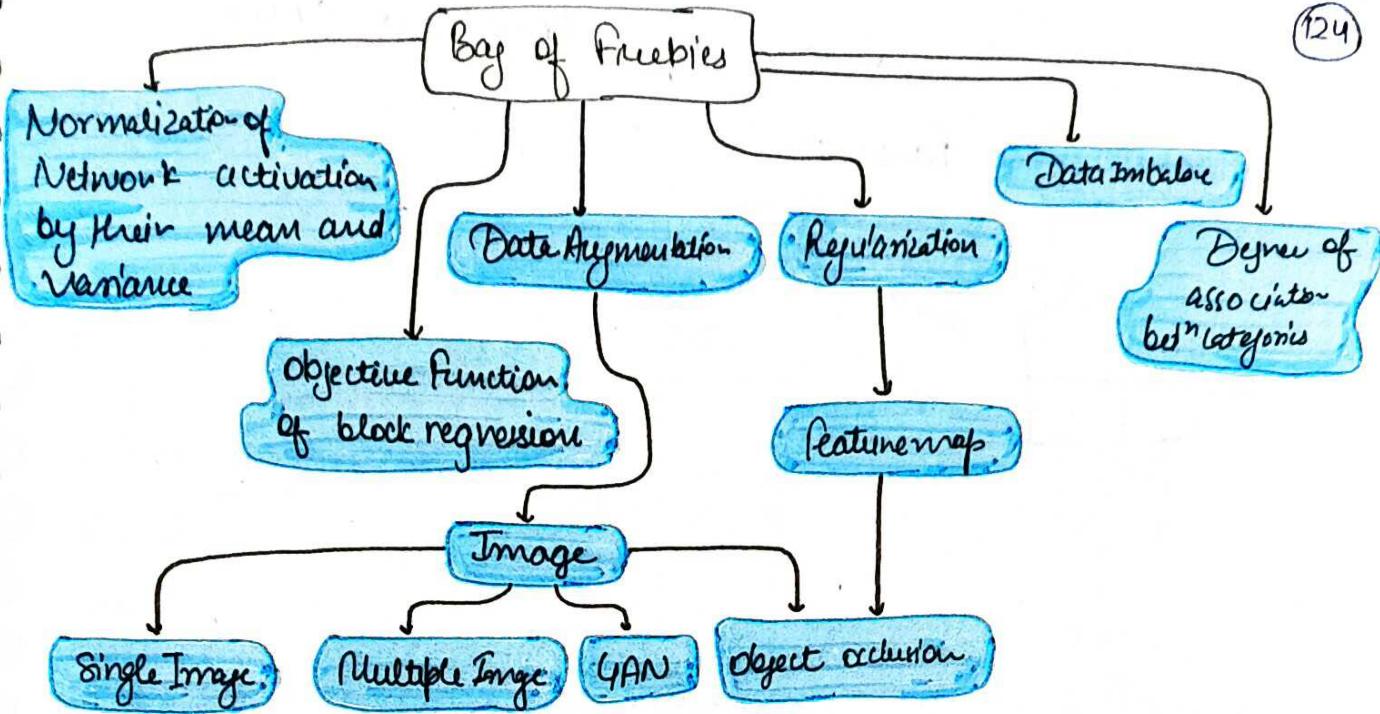
→ Bag of Specials

Architecture Major changes

- * Increase complexity
- also Increase Accuracy

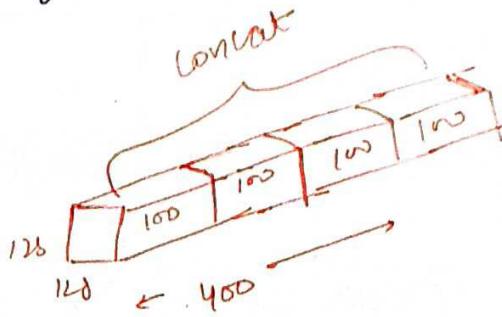
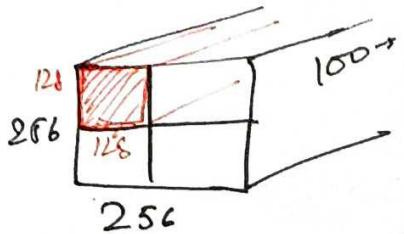
Training Time ↑↑

(Mosaic data augmentation)



Focus Layer → Just before back bone CNN

↳ Reduce the Parameter



YOLO V6

Map → 2.5x and Speed → 2.5x faster

Changes

1. Backbone CNN → New pretrained CNN

↳ Efficient NET

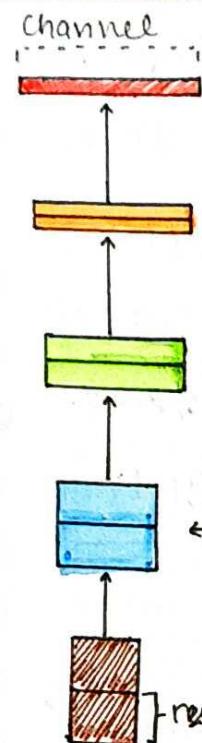
Successor of Mobile Net

↓
less Memory
less Computational power

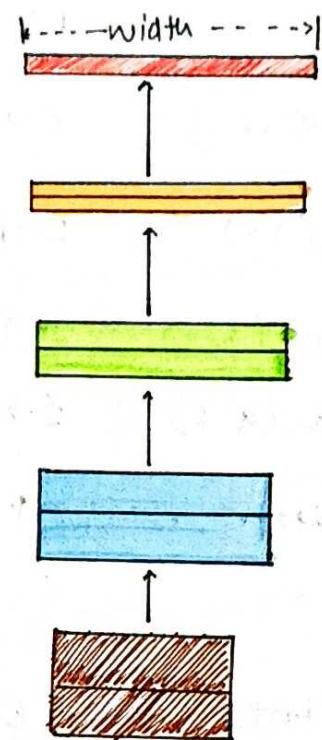
→ This model also work on mobile application.

Efficient Net

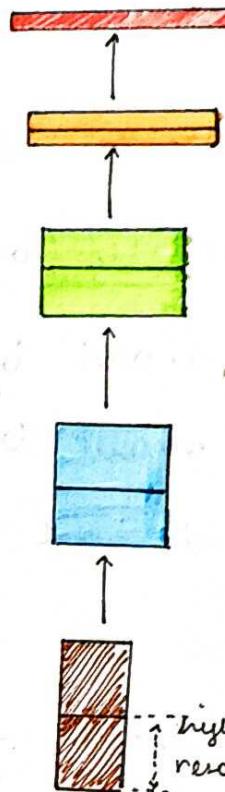
- ↳ width Scaling
- ↳ depth Scaling
- ↳ resolution Scaling



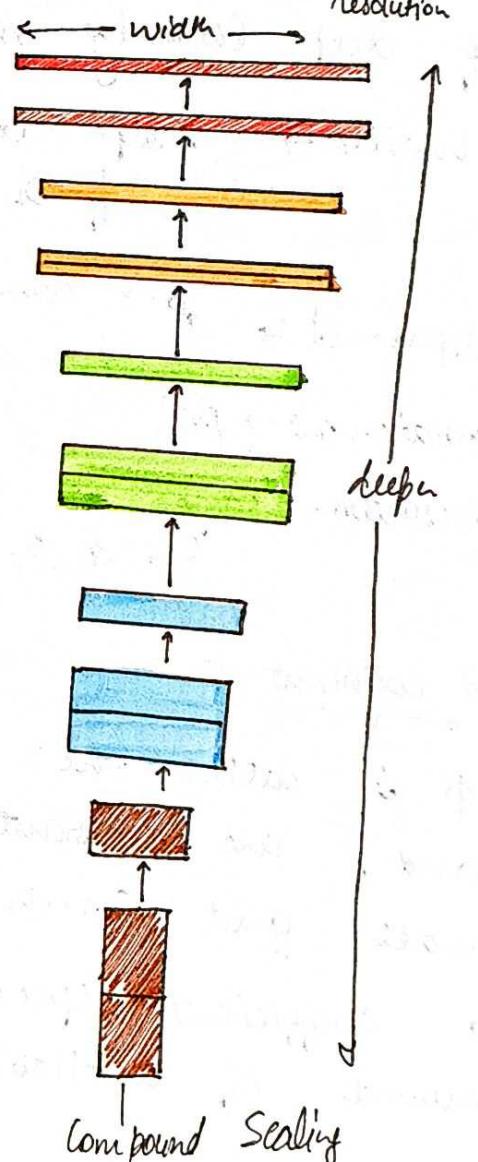
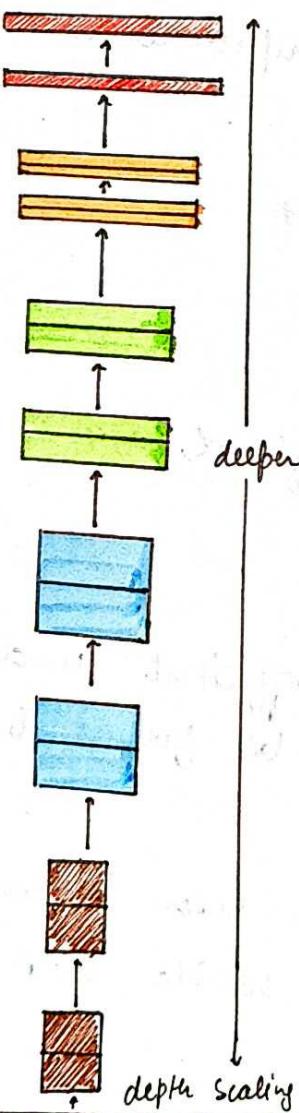
a) baseline



b) width scaling



d) depth Scaling



Compound Scaling

Models

Models	FLOPS	TOPS	Accuracy
Baseline model (EfficientNet-B0)	0.4B		77.3%
Scale model by depth ($d=4$)	1.8B		79.0%
Scale model by resolution ($r=2$)	1.9B		79.1%
Scale model by width ($w=2$)	3.8B		78.9%
<u>Compound Scale ($d=1.4$, $w=3.2$, $r=1.3$)</u>	<u>1.8B</u>		<u>81.1%</u>

Less parameters and high accuracy

Flops stands for "floating-point Operations Per Second" and it is a measure of computational performance used to evaluate the efficiency and computational cost of deep learning models.

Simple words \rightarrow helps to define computational cost of DL models.

depth = $d \rightarrow \alpha^\phi$ compound coefficient

width = $w \rightarrow \beta^\phi$

resolution = $r \rightarrow \gamma^\phi$
 α, β, γ is value of d, w, r

compound coefficient ϕ

where ϕ is called the compound coefficient and α, β and γ are constants that can be found by a small grid search.

ϕ is a coefficient specified by the user to control the amount of available resources. while α, β and γ

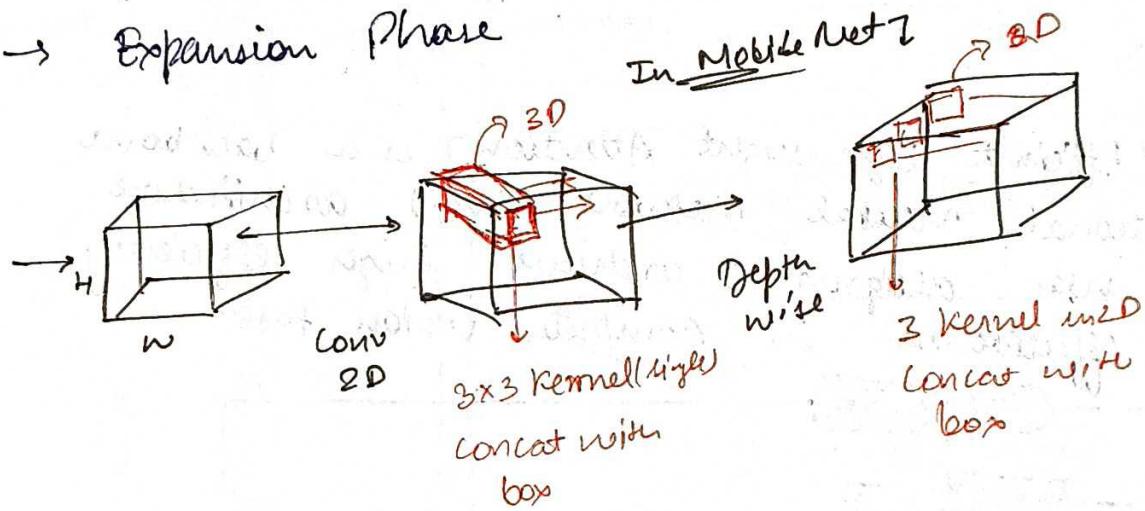
assign these resources to the network's depth, width and resolution respectively.

(126)

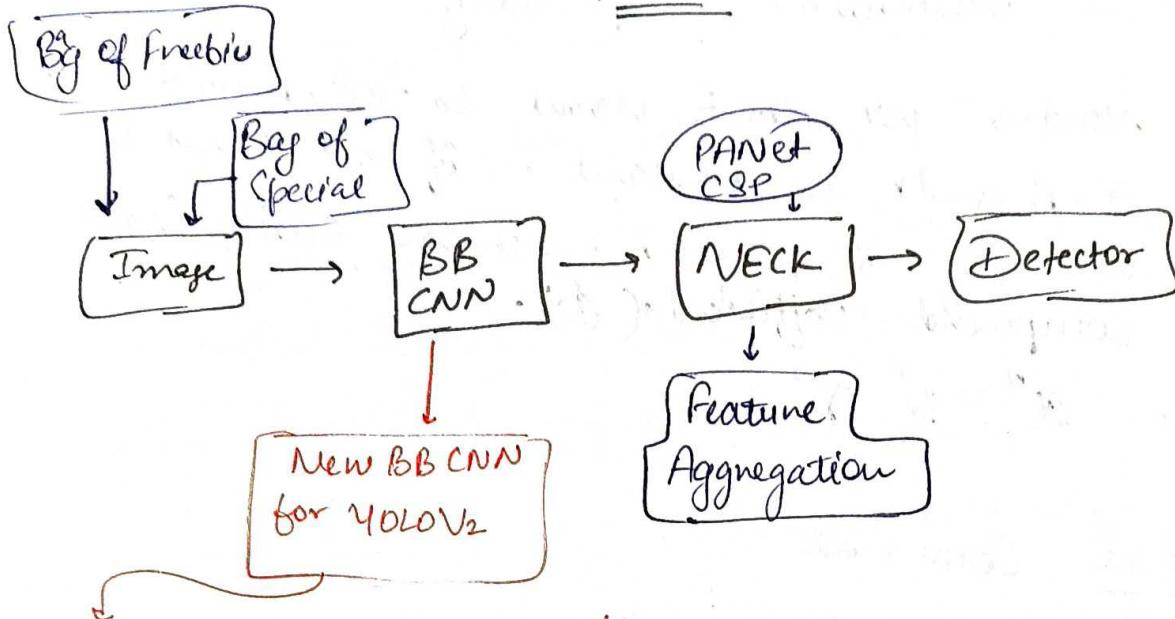
In simple word \rightarrow we don't want to change α, β and r .
Because α, β and r is constant. If we want to change value of α, β and r . we can easily change the of compound coefficient ϕ .
 $\phi \rightarrow \alpha^2, \beta^2, r^2$

Mobile Net Conv Block

- \rightarrow Squeeze and Excitation Phase
- \rightarrow Expansion Phase



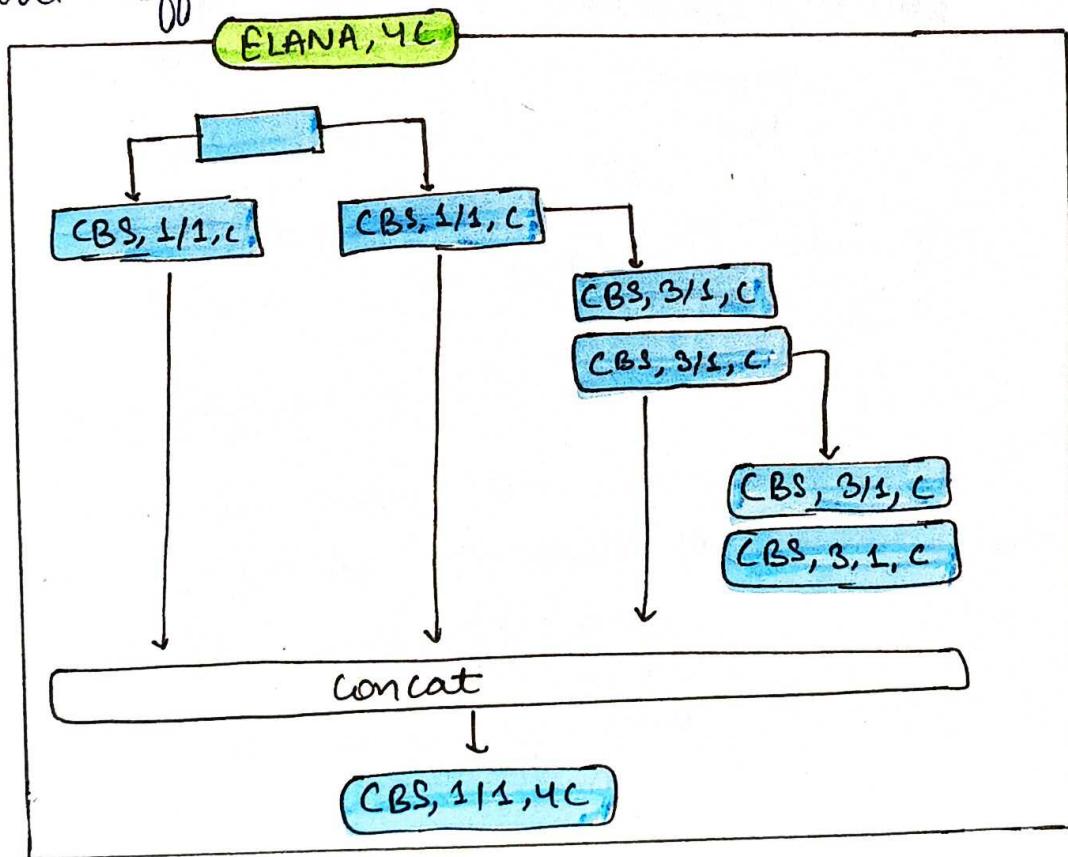
YOLO V₂

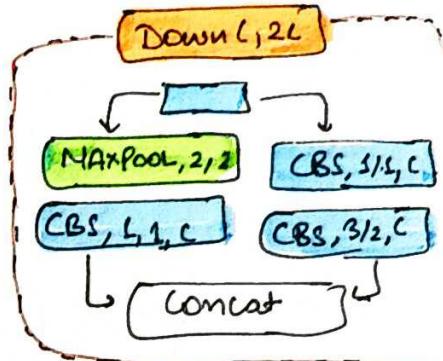
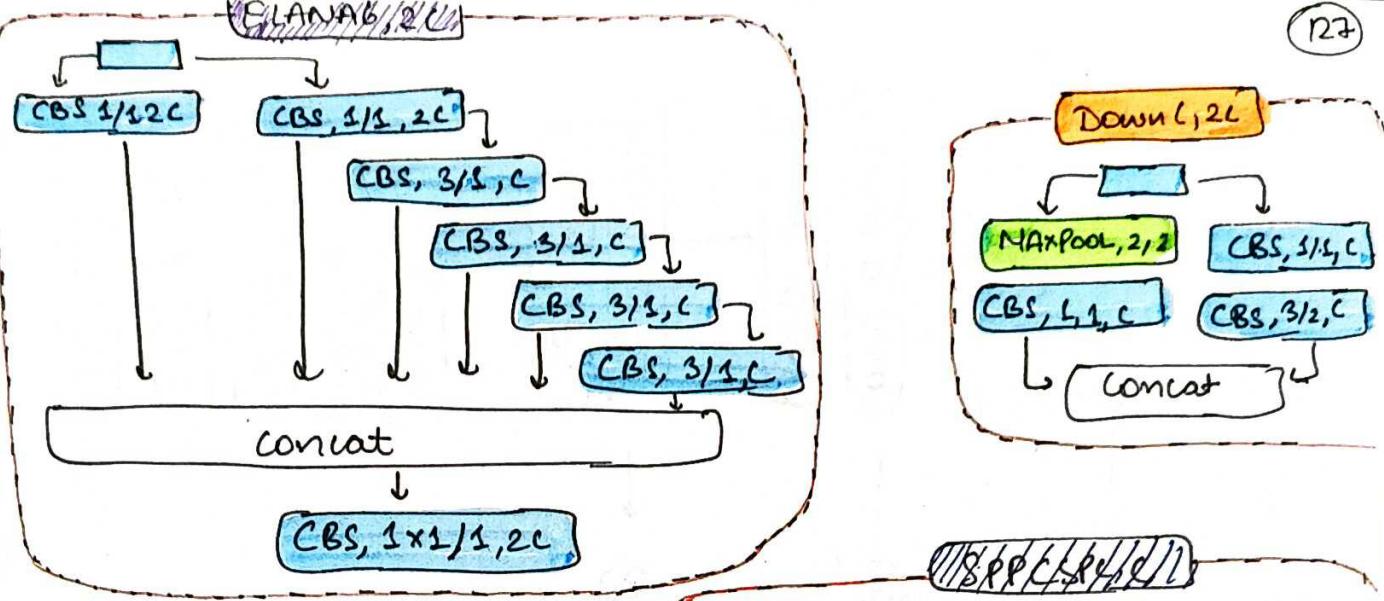


Past and Accurate than efficient Net

↳ Using Mechanism Attention → Ignore irrelevant things or information.

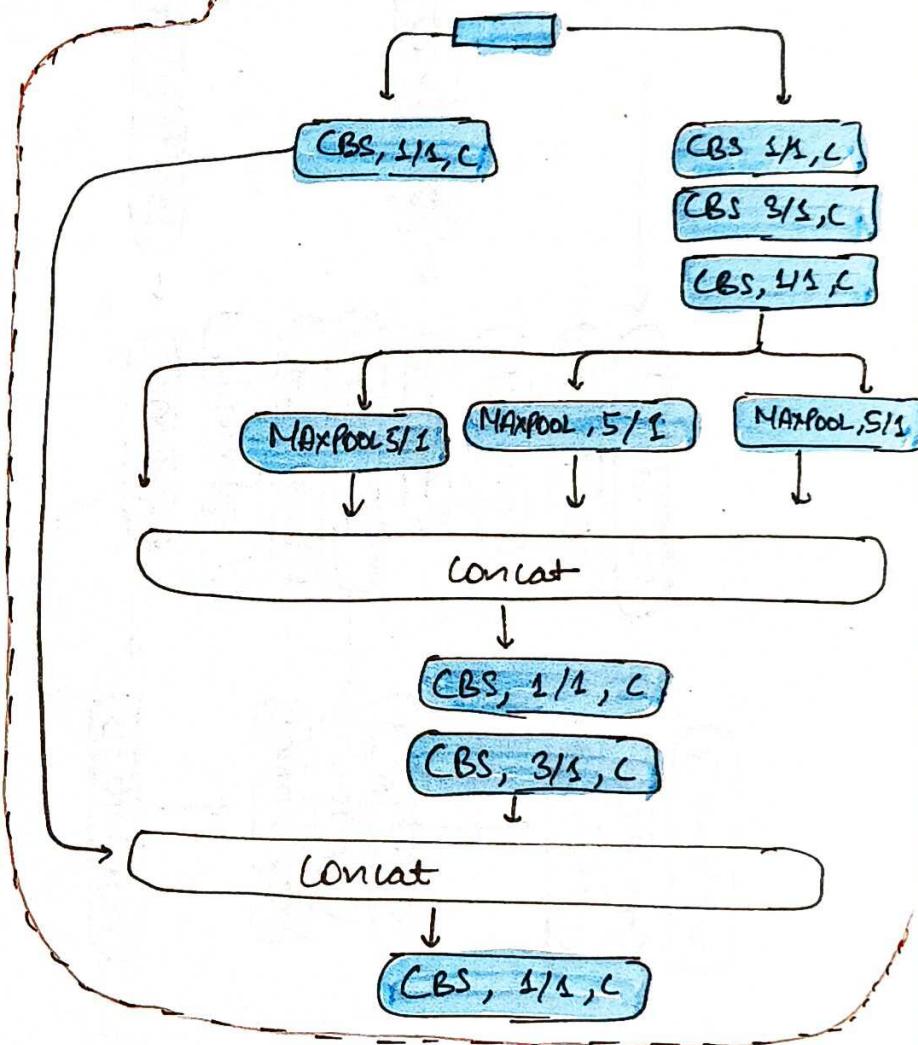
ELAN (Efficient Lightweight Attention) is a backbone convolutional neural network (CNN) architecture that was designed to achieve higher efficiency and effectiveness in computer vision task.

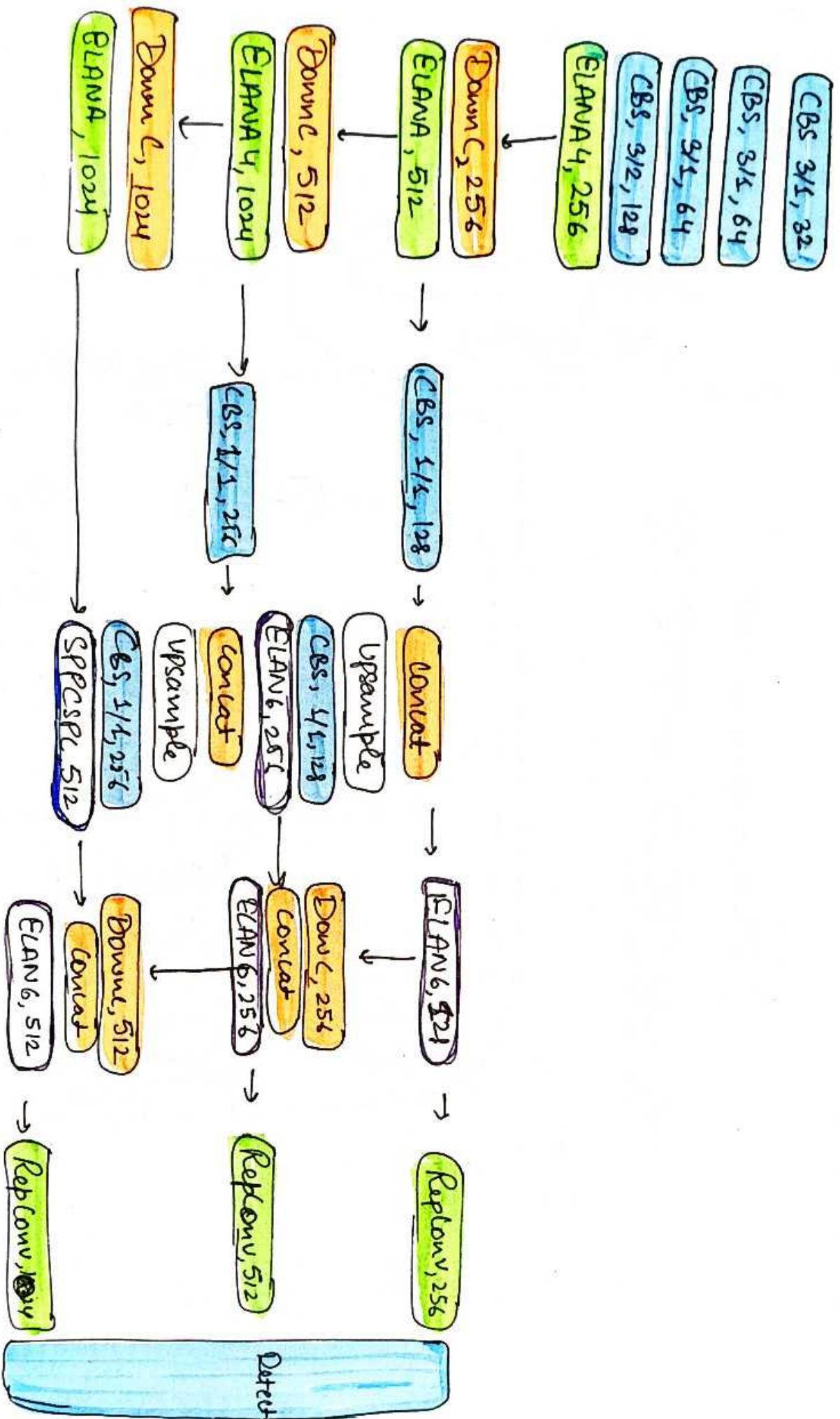




$$SILU = \alpha * \sigma(\alpha) = \alpha * \frac{1}{1 + e^{-2}}$$

$$\text{CBS, } k \times s, c = \text{conv } k \times h \times s - \text{BN} - \text{SILU}$$

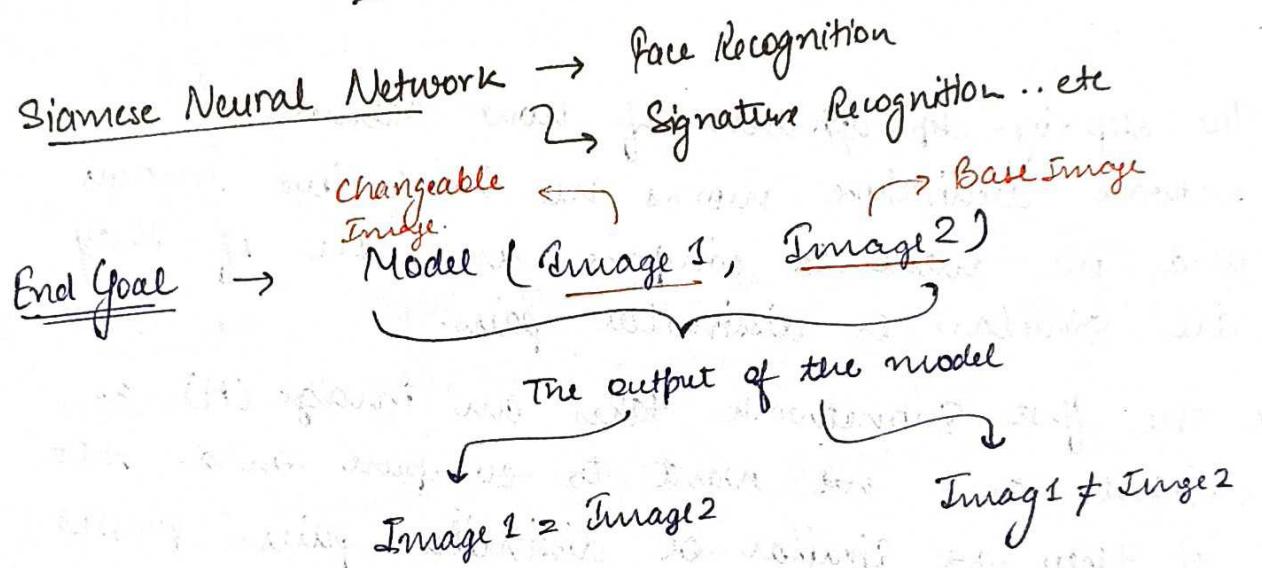




Roboflow

Roboflow basically is an online tool, using which we can ANNOTATE the image or our data and download the data in the required format.

SIAMESE NEURAL NETWORK

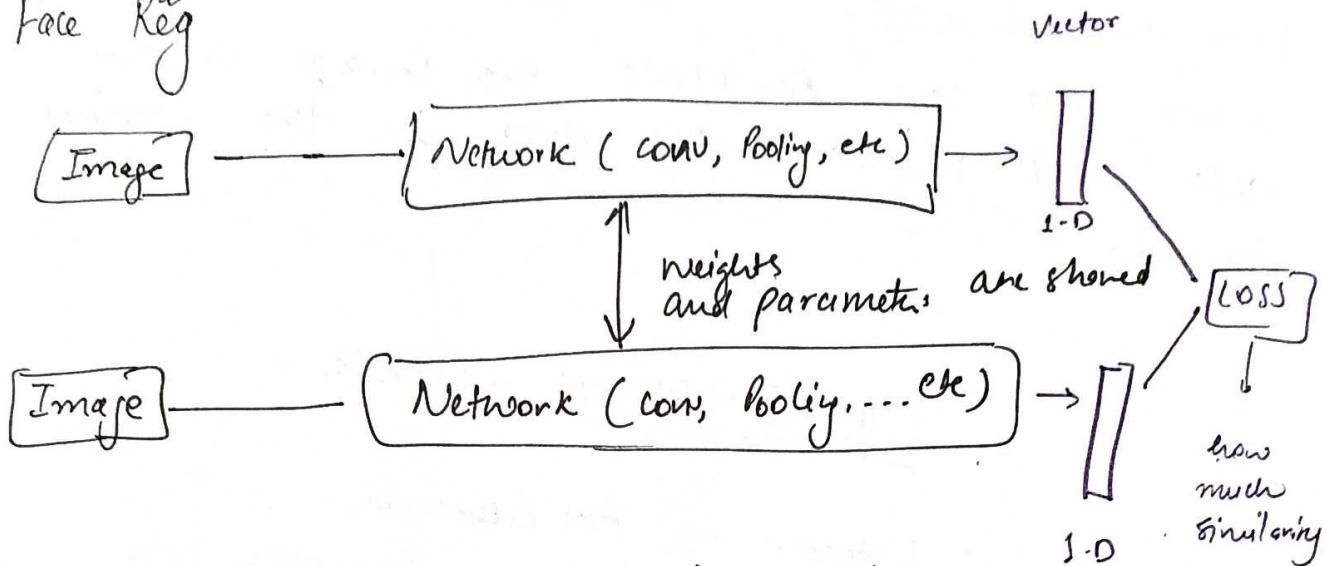


What we have to study?

- Use case of
- SIAMESE Neural Network
- Architecture
- Loss function → (i) Triplet loss funcⁿ
(ii) contrastive loss

Siamese Neural Network

- Face Reg



The step-by-step approach of how siamese network architecture works we have two Images and we want to compare and see if they are similar or dissimilar pairs.

1. The first subnetwork takes an image (A) as input and we want to compare and see if they are similar or dissimilar pairs. pass through convolutional layers and fully connected layers, we get a vector representation of the image.
2. Again pass the second image (B) through a network that is exactly the same with the same weights and parameters.
3. Now we have two encodings $E(A)$ and $E(B)$ from the respective images, we can compare these two to know how similar the two images are. If the images are similar then the encodings will also be quite similar.

4. we will measure the distance b/w these two vectors and if the distance b/w these is small then the vectors are similar or of the same classes and if the distance b/w is large then the vectors are different from one another, based on the score.

Siamese

↳ we do not need lot of training data

↳ Time ↑↑

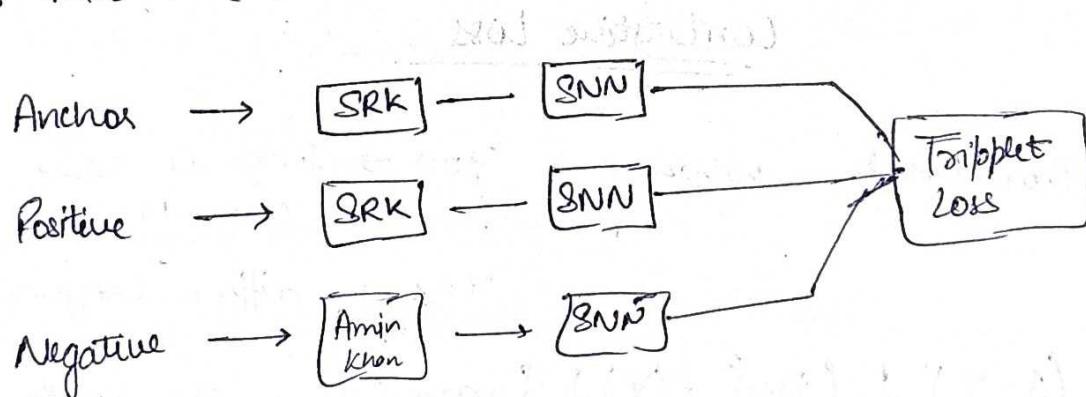
There are 2 losses

1. Triplet loss : Triplet loss is a loss function where a baseline is compared to a positive & a negative.

$$L(A, P, N) = \max (||f(A) - f(P)||^2 - ||f(A) - f(N)||^2 + \alpha)$$

- max → (A, N)

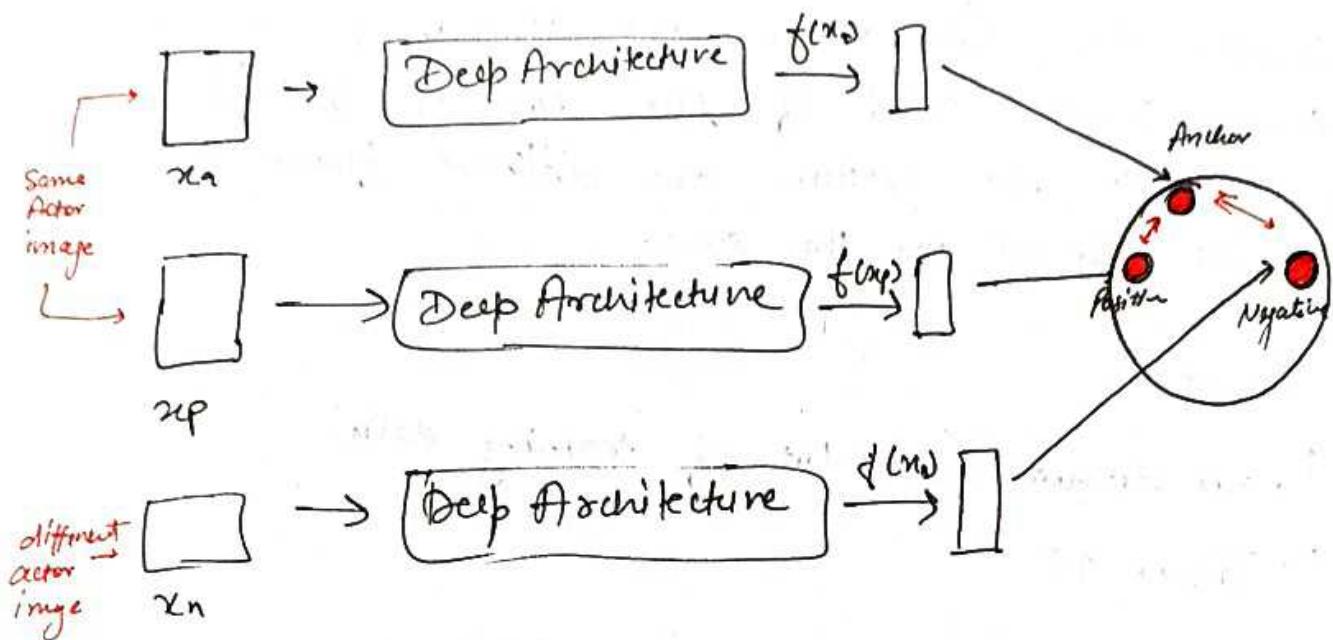
- min → (A, P)



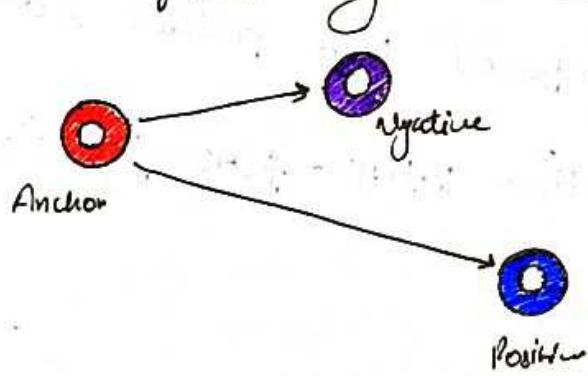
→ distance b/w Anchor and Positive is less.

→ distance b/w Anchor and Negative is high.

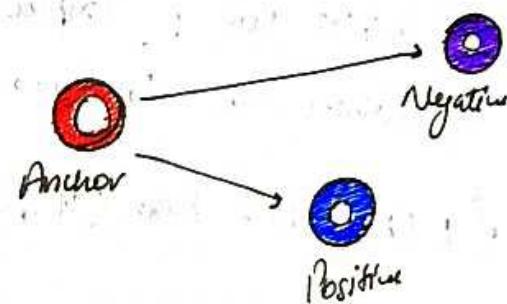
Triplet Loss with Keras and Tensorflow



Before Learning



After Learning



Contrastive Loss

$y \rightarrow$ True Label where

$y=0 \rightarrow$ Image of same categories

$y=1 \rightarrow$ different categories

$$\text{eqn} - (1-y) \frac{1}{2} (D_w)^2 + (y) \frac{1}{2} \left\{ \max(0, m - D_w) \right\}^2$$

margin (constant)

Case 1 $\Rightarrow \gamma = 0$

So, our eqn is

$$(1-\gamma) \frac{1}{2} (Dw)^2 + (\gamma) \frac{1}{2} \{ \max(0, m - Dw) \}$$

for case 1

Case 2 $\Rightarrow \gamma = 1$.

different

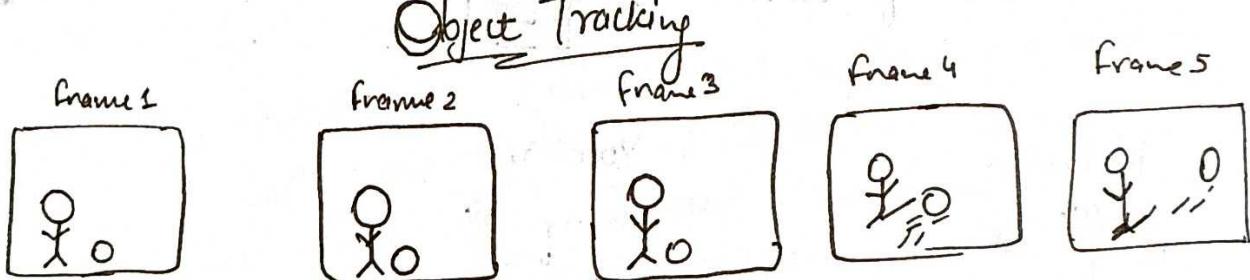
→ Where $Dw \rightarrow$ Euclidean distance
betⁿ 2 Images

- ⊕ In triplet loss we use Anchor, Positive and Negative

but in contrastive loss we use → Anchor, Positive
Negative
for different image

distance bet
A and N ↑↑

for similar image
distance bet
A, P is ↓↓



Video - 150 FPS

Real Time Object Tracking → Video

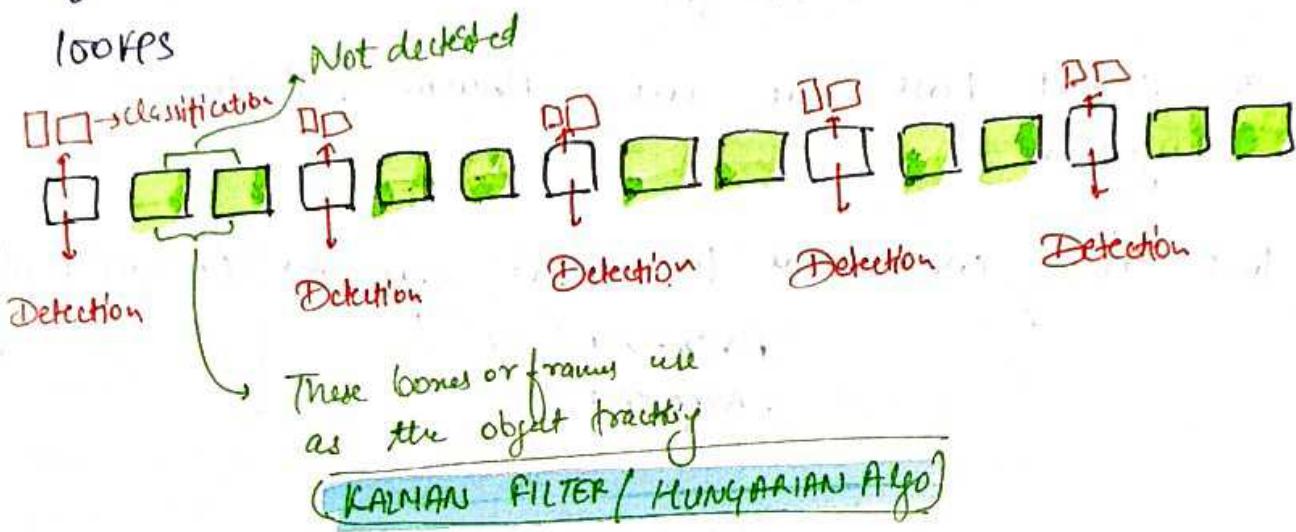
DeepSort → Object Tracking Algorithm

Object Tracking

Backbone → YOLO - V5

33.3 FPS

[Video] → helmets



$$IOU = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$

Kalman
Filter

Estimation

YOLO V5

Detection

Object Tracking → IOU (Estimation & Detection)

Object Tracking → Similarity (Estimation & Detection)
(SiamFC Network)

Code Explanation

Object Tracking Project

deebsort = DeepSort (cfg. DeepSort.REID_CKPT,
max_dist = cfg. DEEPSORT.MAX_DIST, min_confidence = cfg.
DEEPSORT.NMS_MAX_OVERLAP, max_iou_distance = cfg. DEEPSORT.
MAX_IOU_DISTANCE,
max_age = cfg. DEEPSORT.MAX_AGE, n_init = cfg. DEEPSORT.N_INIT
, nn_budget = cfg. DEEPSORT.NN_BUDGET,
use_cuda = True)

cfg. DEEPSORT.REID-CKPT: This is the path to the checkpoint file of a deep neural network model used for ReID (Re-Identification). ReID is used to generate feature embeddings that capture the unique characteristics of objects, enabling the algorithm to distinguish betn different objects based on their appearance.

max-dist: Maximum Squared Mahalanobis distance betn the feature embeddings of two detections in order to consider them as the same object in the tracking process.

min confidence: Minimum confidence threshold for a detection to be considered valid and used in the tracking process.

nms_max_overlap: Non-Maximum Suppression (NMS) threshold to remove overlapping bounding boxes from the detection outputs before tracking.

max_iou_distance: Maximum Intersection over Union (IoU) distance b/w two boundary boxes to be considered for association in the tracking process.

max_age: Maximum number of frames a track needs to be detected in order to keep alive without assigned detection.

n_init: Number of frames a track needs to be detected in order to be confirmed and included in the tracking process.

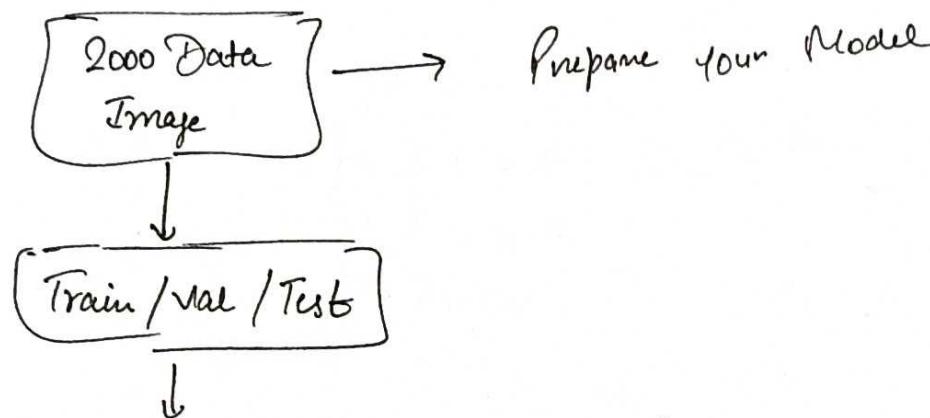
nn_budget: The maximum number of nearest neighbors that are considered when associating detection with tracks.

use_cuda: A boolean flag indicating whether to use CUDA (GPU) acceleration if available.

In summary, this line of code initializes an instance of the DEEP SORT tracker along with various configuration parameters. These parameters determine how detection and association with existing tracks, how track updates are performed, and how tracks are updated, how track updates are performed, and how tracks are updated for object tracking. The initialized deepsort object can then be used to perform object tracking on video frames using its tracking and updating methods.

Image Classification

Architecture



Transfer Learning → 95% Accuracy

↳ CNN

↳ VGG

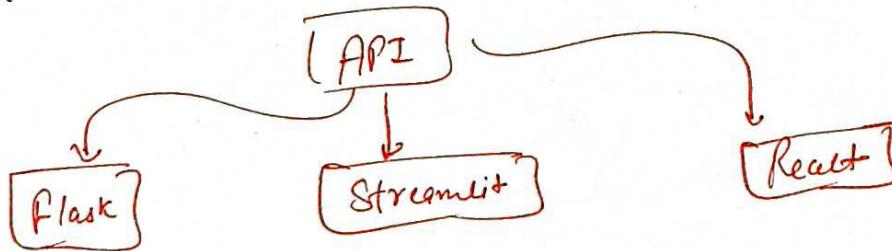
↳ ResNet

↳ LENET

↳ AlexNet

↳ DenseNet

↳ choose any one model and save it.
and use this same model for in application.



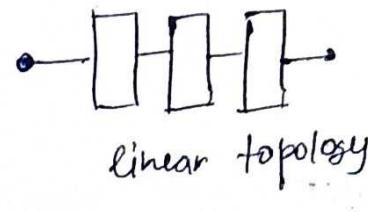
Functional Models

Problems with Sequential Model

ANN \leftrightarrow CNN

↳ Sequential

- 1 input
- 1 output
- linear

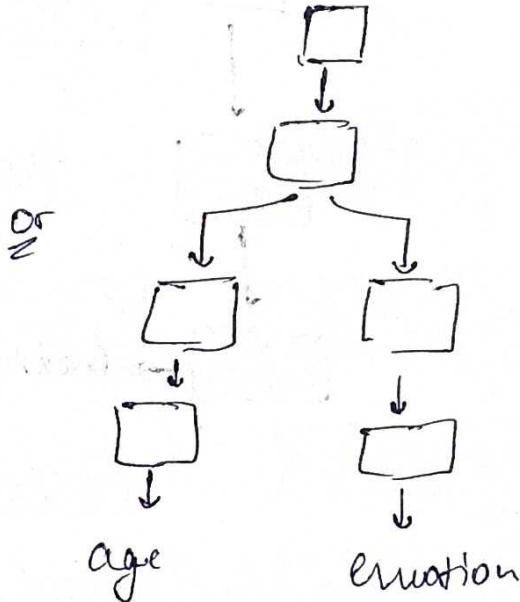
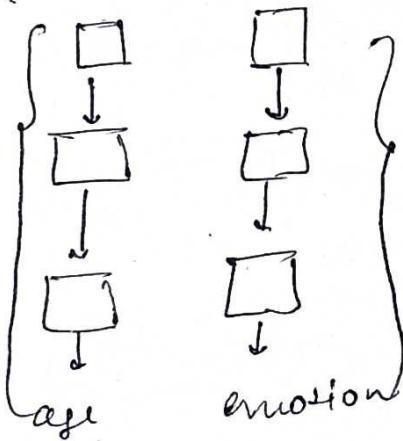


linear topology

Example → → Image → human → 25000 images
dataset faces

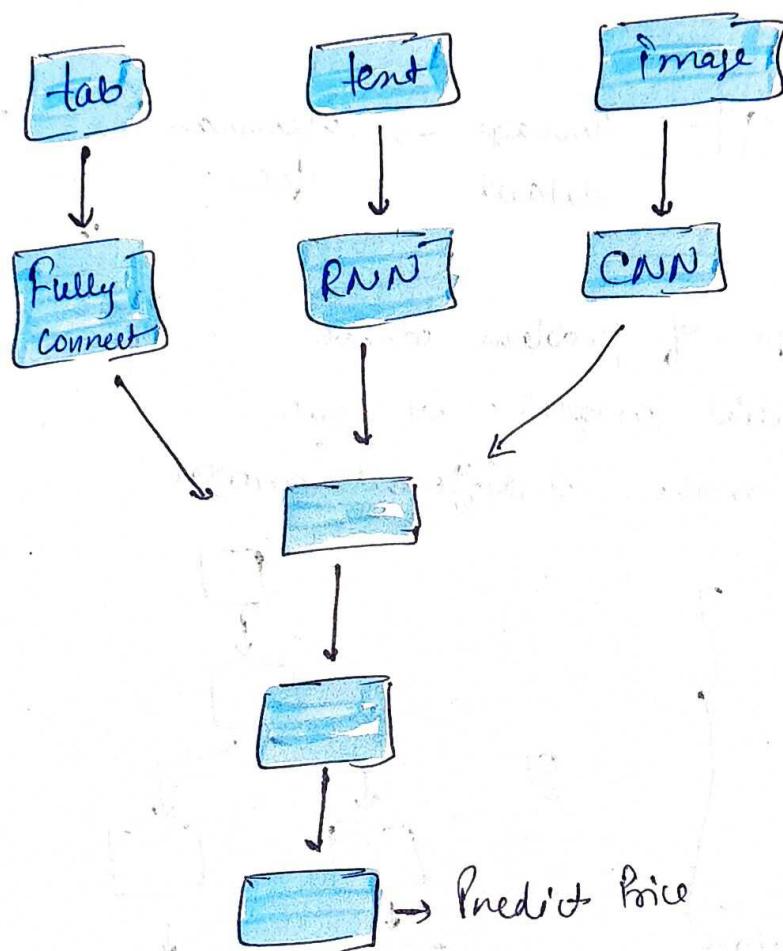
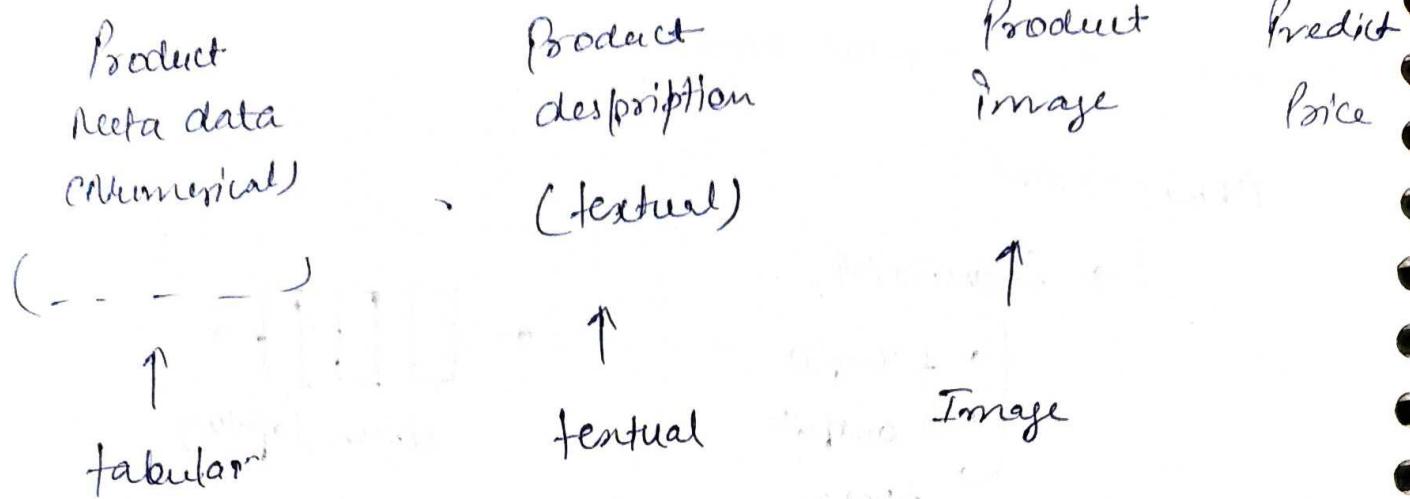
So, In this type of problem cannot solve sequential model. or we have to make 2 different model

age ? emotion ?
Sad happy angry



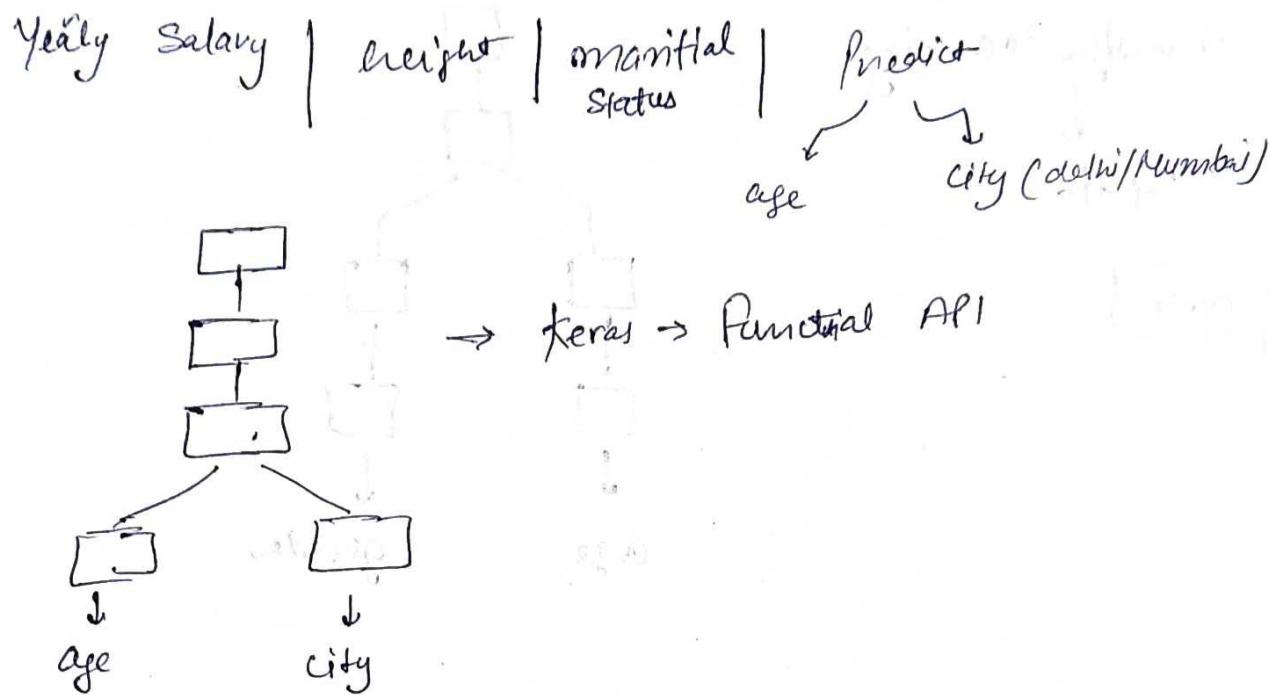
Here's sequential model fail.

Example → E-commerce website

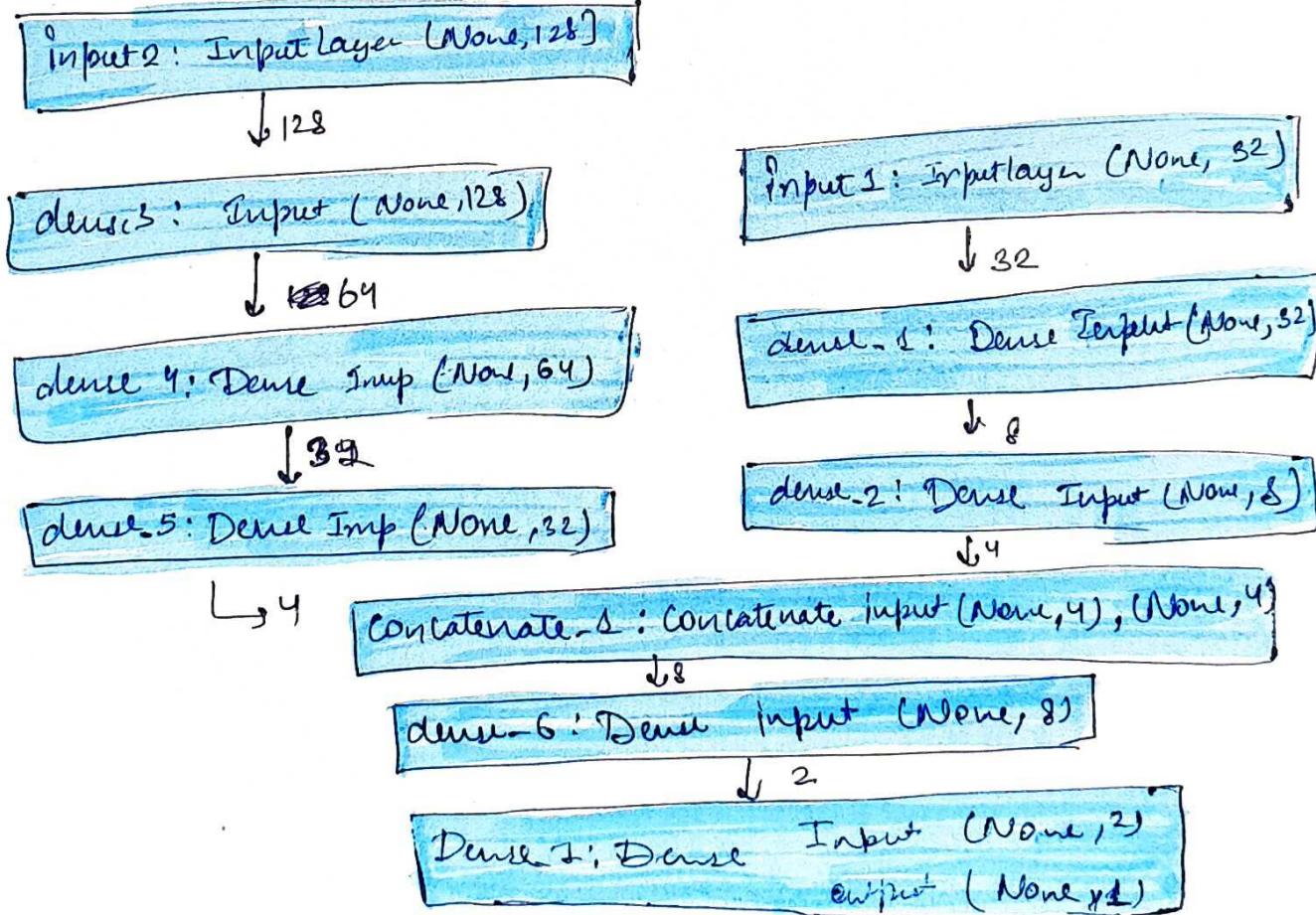


Let's apply functional API on some data.

A Simple Example



Multi Input Model



Multi Output Model

age/emotion → [8] → 20000 Images → gender
→ eye

Transfer Learning.

↓
VGG16

[code]

