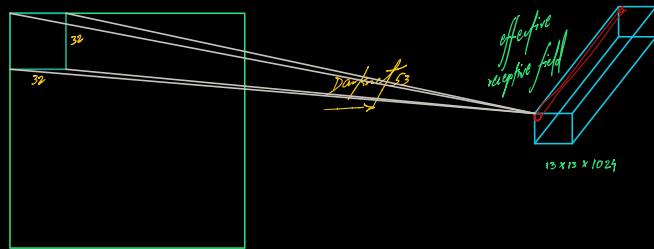


→ intermediate result  
of  $100 \times 100$

\* If there's only one object, then  $P^o$

↓  
probability of  
having an object  
of interest will be  
zero or close to  
zero

\* YOLO is designed to work  
on square images but  
we can always resize  
the input to be a square

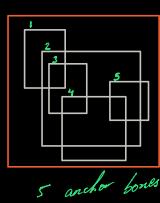


\* what if the object  
spans multiple receptive fields?

↓  
the output of NN has  
to work receptive field, the  
height weight can be anything

↓  
anchor boxes solve  
some of the issue

↑  
 $\text{range} \rightarrow 0^{+0.2}$  → it'll be at most one  
pixel away. Try not  
bb to be very close to  
anchor boxes



\* solution → have anchor boxes

\* left represent bb as a modification  
of anchor boxes

$$(\text{cx}, \text{cy}, \text{pw}, \text{ph})$$

The improved performance by 3x or 4x  
by using anchor boxes

sigmoid function

$$\sigma(cx) + cx$$

$$\sigma(cy) + cy$$

$$(\text{bx}, \text{by}, \text{bw}, \text{bh})$$

$$\begin{cases} \text{bx} \\ \text{by} \\ \text{bw} \\ \text{bh} \end{cases} = \begin{cases} \sigma(cx) + cx \\ \sigma(cy) + cy \\ \text{pw} \cdot e^{cw} \\ \text{ph} \cdot e^{ch} \end{cases}$$

→ means dropping on train after bouncing boxes

$K = 5$  anchor boxes

↓  
5 boxes in which there is a high likelihood of finding an object

↓  
this can act as "prior" if to find object more successfully.

# Default weights  
are used for not only  
loss but propagated to  
all the layers

\* left work with objectives now (?)  
prob. of finding an object in the first

society of deep learning b  
↑ loss for design

→ model learns all this implicitly since  
we penalize the model accordingly

tn | tg | tw | Th | P0 | P1 | P2 | ... | Pn

[if class = woman] → if we do find  
a woman, obj will  
put her in person  
or woman class. →  
with weight the  
not only one  
class

for every class we multiply

P0 with P1  
↓ general prob  
for finding

so we use squared loss  
as often  
doesn't work  
here

# loss function  
if we take  $P = 3$

12x12x255

∴ 12x12x3 possible losses

they v if a high parameter.  
found v if a low parameter.  
An obj works well → helps w.r.t. the  
refuses for weight

①  $\lambda_{coord} * \text{square loss on } tn, tg, tw & Th$   
+ log. loss from linear  
for each ② + log. loss for each  $P_i$   
↳ due to logistic regression  
bounding box containing an object in ground truth

single result will be independent  
of individual loss → 0 to 1  
↳ no object losses in ground truth → over here we  
don't care about anything  
else so we just want  
P0 to be close to zero or zero.

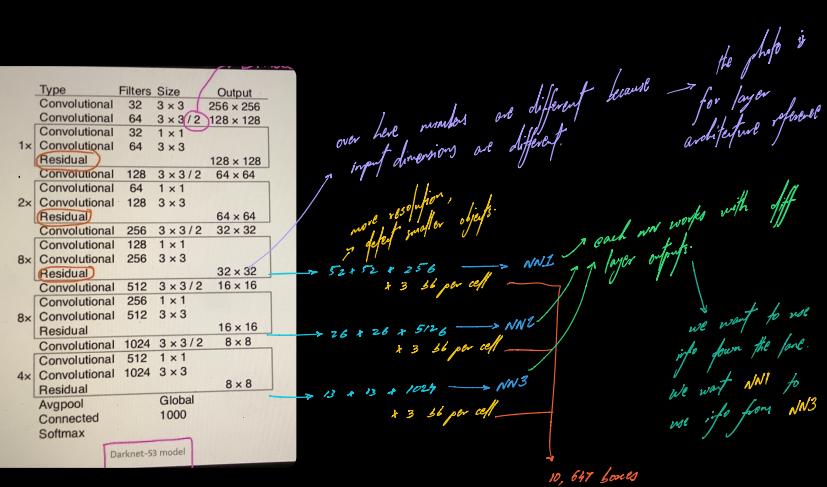
# We use pretrained Darknet  
but were not freezing the weights.

logits are of a probability  
square loss because red value

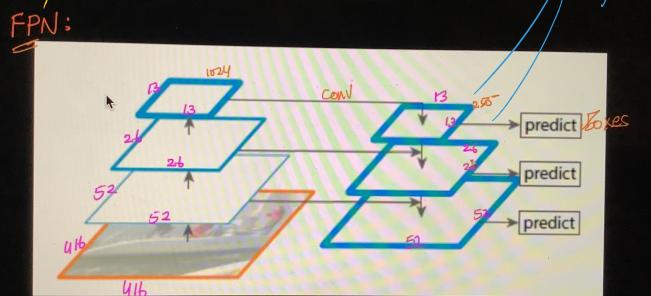
## \* Multi-Scale Prediction

Type	Filters	Size	Output	
Convolutional	32	$3 \times 3$	256 x 256	
Convolutional	64	$3 \times 3 / 2$	128 x 128	
Convolutional	32	$1 \times 1$		
Convolutional	64	$3 \times 3$	128 x 128	
1x	Residual			
Convolutional	128	$3 \times 3 / 2$	64 x 64	
Convolutional	64	$1 \times 1$		
2x	Convolutional	128	$3 \times 3$	64 x 64
Residual				
Convolutional	256	$3 \times 3 / 2$	32 x 32	
Convolutional	128	$1 \times 1$		
8x	Convolutional	256	$3 \times 3$	32 x 32
Residual				
Convolutional	512	$3 \times 3 / 2$	16 x 16	
Convolutional	256	$1 \times 1$		
8x	Convolutional	512	$3 \times 3 / 2$	16 x 16
Residual				
Convolutional	1024	$3 \times 3 / 2$	8 x 8	
Convolutional	512	$1 \times 1$		
4x	Convolutional	1024	$3 \times 3$	8 x 8
Residual				
Avgpool				
Connected	Global		1000	
Softmax				

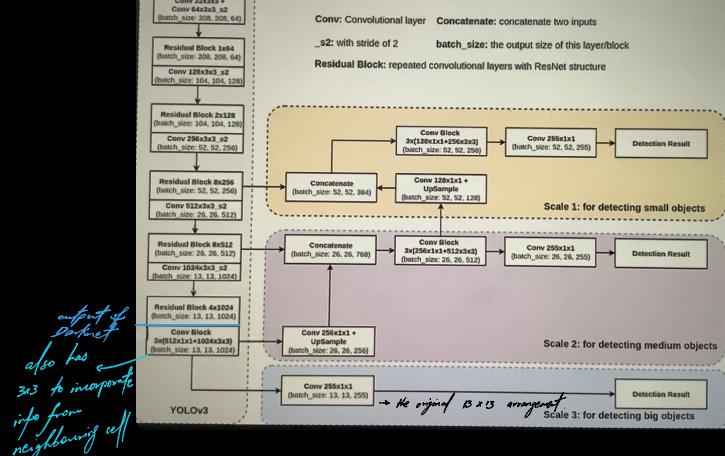
Darknet-53 model

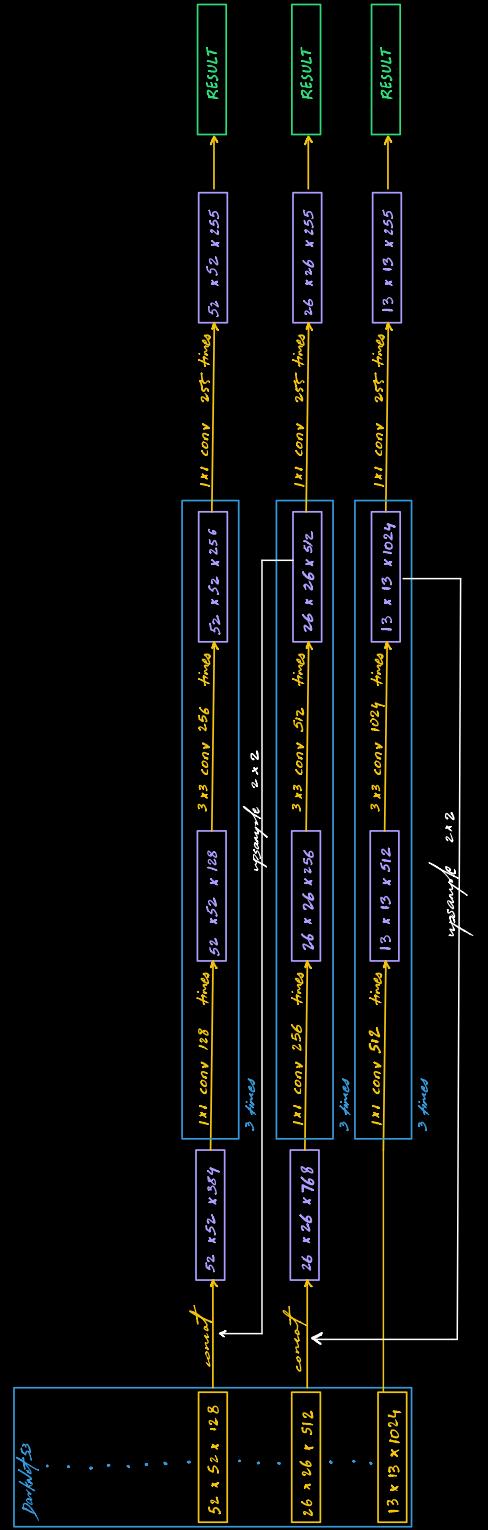


## Feature Pyramid Networks:



## YOLOv3 Network Architecture





\* People know the Dapple  
not well, then they favor  
the remaining party

# code taken from  
Jason brownlee's blog

\* code is pretty obvious, to  
me. so I'm not writing much

# Just for reference

```

32 def make_yolov3_model():
33     input_image = Input(shape=(None, None, 3))
34     # Layer 0 => 4
35     x = _conv_block(input_image, [{"filter": 32, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 0},
36         {"filter": 64, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 1},
37         {"filter": 32, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 2},
38         {"filter": 64, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 3}])
39     # Layer 5 => 8
40     x = _conv_block(x, [{"filter": 128, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 5},  
    → layer_idx 4 is  
    a skip connection.  
    we add it later.
41         {"filter": 64, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 6},
42         {"filter": 128, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 7}])
43     # Layer 9 => 11
44     x = _conv_block(x, [{"filter": 64, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 9},
45         {"filter": 128, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 10}])
46     # Layer 12 => 15
47     x = _conv_block(x, [{"filter": 256, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 12},
48         {"filter": 128, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 13},
49         {"filter": 256, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 14}])
50     # Layer 16 => 36
51     for i in range(7):
52         x = _conv_block(x, [{"filter": 128, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 16+i*3},
53             {"filter": 256, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 17+i*3}])
54     skip_36 = x → will be useful for our concat later on
55     # Layer 37 => 40
56     x = _conv_block(x, [{"filter": 512, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 37},
57         {"filter": 256, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 38},
58         {"filter": 512, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 39}])
59     # Layer 41 => 61
60     for i in range(7): → not 8 because already written one above
61         x = _conv_block(x, [{"filter": 256, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 41+i*3},
62             {"filter": 512, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 42+i*3}])
63     skip_61 = x → saving for FPN concat
64     # Layer 62 => 65
65     x = _conv_block(x, [{"filter": 1024, "kernel": 3, "stride": 2, "bnorm": True, "leaky": True, "layer_idx": 62},
66         {"filter": 512, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 63},
67         {"filter": 1024, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 64}])
68     # Layer 66 => 74
69     for i in range(3):
70         x = _conv_block(x, [{"filter": 512, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 66+i*3},
71             {"filter": 1024, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 67+i*3}])
72     # Layer 75 => 79 point
73     x = _conv_block(x, [{"filter": 512, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 75},
74         {"filter": 1024, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 76},
75         {"filter": 512, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 77},
76         {"filter": 1024, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 78},
77         {"filter": 512, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 79}], skip=False)
78     # Layer 80 => 82
79     yolo_82 = _conv_block(x, [{"filter": 1024, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 80},
80         {"filter": 255, "kernel": 1, "stride": 1, "bnorm": False, "leaky": False, "layer_idx": 81}], skip=False)
81     # Layer 83 => 86
82     x = _conv_block(x, [{"filter": 256, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 84}], skip=False)
83     x = UpSampling2D(2)(x)
84     x = concatenate([x, skip_61])
85     # Layer 87 => 91
86     x = _conv_block(x, [{"filter": 256, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 87},
87         {"filter": 512, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 88},
88         {"filter": 256, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 89},
89         {"filter": 512, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 90},
90         {"filter": 256, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 91}], skip=False)
91     # Layer 92 => 94
92     yolo_94 = _conv_block(x, [{"filter": 512, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 92},
93         {"filter": 255, "kernel": 1, "stride": 1, "bnorm": False, "leaky": False, "layer_idx": 93}], skip=False)
94     # Layer 95 => 98
95     x = _conv_block(x, [{"filter": 128, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 96}], skip=False)
96     x = UpSampling2D(2)(x)
97     x = concatenate([x, skip_36])
98     # Layer 99 => 106
99     yolo_106 = _conv_block(x, [{"filter": 128, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 99},
100         {"filter": 256, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 100},
101         {"filter": 128, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 101},
102         {"filter": 256, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 102},
103         {"filter": 128, "kernel": 1, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 103},
104         {"filter": 256, "kernel": 3, "stride": 1, "bnorm": True, "leaky": True, "layer_idx": 104},
105         {"filter": 255, "kernel": 1, "stride": 1, "bnorm": False, "leaky": False, "layer_idx": 105}], skip=False)
106     model = Model(input_image, [yolo_82, yolo_94, yolo_106])
107     return model

```

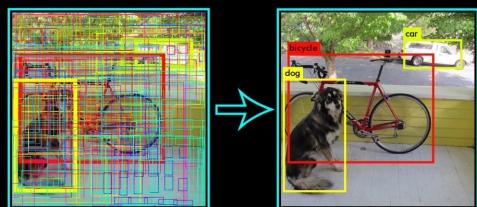
Type	Filters	Size
Convolutional	32	3 x 3
Convolutional	64	3 x 3 / 2
Convolutional	32	1 x 1
Convolutional	64	3 x 3
Residual		
Convolutional	128	3 x 3 / 2
Convolutional	64	1 x 1
Convolutional	128	3 x 3
Residual		
Convolutional	256	3 x 3 / 2
Convolutional	128	1 x 1
Convolutional	256	3 x 3
Residual		
Convolutional	512	3 x 3 / 2
Convolutional	256	1 x 1
Convolutional	512	3 x 3
Residual		
Convolutional	1024	3 x 3 / 2
Convolutional	512	1 x 1
Convolutional	1024	3 x 3
Residual		
Avgpool		Global
Connected		1000
Softmax		

Darknet-53 model

# we still need to fix some bugs  
to make this work  
out of 1089 → only take max( $\rho_i + p_i$ ) > 0.5  
6 bones

\* Non max suppression  $\Rightarrow$  after filtering for  $IoU < 0.5$ , we still have a bunch of boxes left.

# first filter every box with  $P_c \leq 0.6$  - just to filter...



Andrew Ng's video on NMS is better for this one



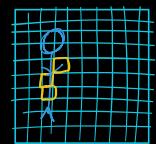
# YOLO v3 uses intelligence & engineering, Google just throws big compute on problems. This is why YOLO v3 is so interesting.

# PJ Reddy: they has all documentation for the C version of YOLO v3, the faster version.

# the documentation is really good, and pretty obvious  
I'm not going to note down any of it.

\* Interested  $\rightarrow$  They have Tiny YOLO v3 for low powered devices.

# conv2DTranspose has more params vs upsampling.  $\therefore$  we choose upsampling



\* each yellow box finds if it got the midpoint of the person, but the person has only one midpoint.

① Pick the box with highest  $P_c$  to highest confidence box call it  $B$

② For every other box, compute its IOU with  $B$

③ Delete the ones that have a high overlap with  $B$  ( $IoU \geq 0.5$ )

in this way we're deleting the non-maximal boxes. Thus the name NMS

depends on how densely packed the objects are in the image.