

Theoretical questions

①

R-channel

1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

G-channel

2	2	2	2
2	2	2	2
2	2	2	2
2	2	2	2

B-channel

1	1	1	1
2	2	2	2
3	3	3	3
4	4	4	4

2x3 filter

1	1	1
1	1	1
1	1	1

Result :

45	45
54	54

②

with zero padding:-

Red

Channel

0	0	0	0	0	0
0	1	1	1	1	0
0	1	1	1	1	0
0	1	1	1	1	0
0	1	1	1	1	0
0	0	0	0	0	0

G Channel

0	0	0	0	0	0
0	2	2	2	2	0
0	2	2	2	2	0
0	2	2	2	2	0
0	2	2	2	2	0
0	0	0	0	0	0

B Channel

0	0	0	0	0	0
0	1	1	1	1	0
0	2	2	2	2	0
0	3	3	3	3	0
0	4	4	4	4	0
0	0	0	0	0	0

filter is same as above.

Result

18	27	27	18
30	45	45	30
36	54	54	36
26	37	37	26

③

Using Dilated (a trous) Convolution with dilation rate of 2.

A 3x3 Kernel with a dilation rate of 2 will have the same field of view as a 5x5 Kernel, while only using 9 parameters. Imagine taking a 5x5 Kernel and deleting every second column and row.



Red Channel

0	0	0	0	0	0
0	1	1	1	1	0
0	1	1	1	1	0
0	1	1	1	1	0
0	1	1	1	1	0
0	0	0	0	0	0

Gr-Channel

0	0	0	0	0	0
0	2	2	2	2	0
0	2	2	2	2	0
0	2	2	2	2	0
0	2	2	2	2	0
0	0	0	0	0	0

B Channel

0	0	0	0	0	0
0	1	1	1	1	0
0	2	2	2	2	0
0	3	3	3	3	0
0	4	4	4	4	0
0	0	0	0	0	0

Result:

filter =

1	1	1
-1	1	1
1	1	1

Result =

24	24
20	20

(4)

Template matching interpretation of convolution  $\Rightarrow$  when ever a filter (f) convo with image (I) provide a response in a image form, we always want a higher (brighter) response. we think of each filter applied to the image as a template. when we convolve the image with that filter, we expect the portions of the image that match that particular template to give a higher response.

$\star f \star I = R$  (response) will only higher if f is totally match with I. In this sense we are matching the templates using convolutions.

(5)

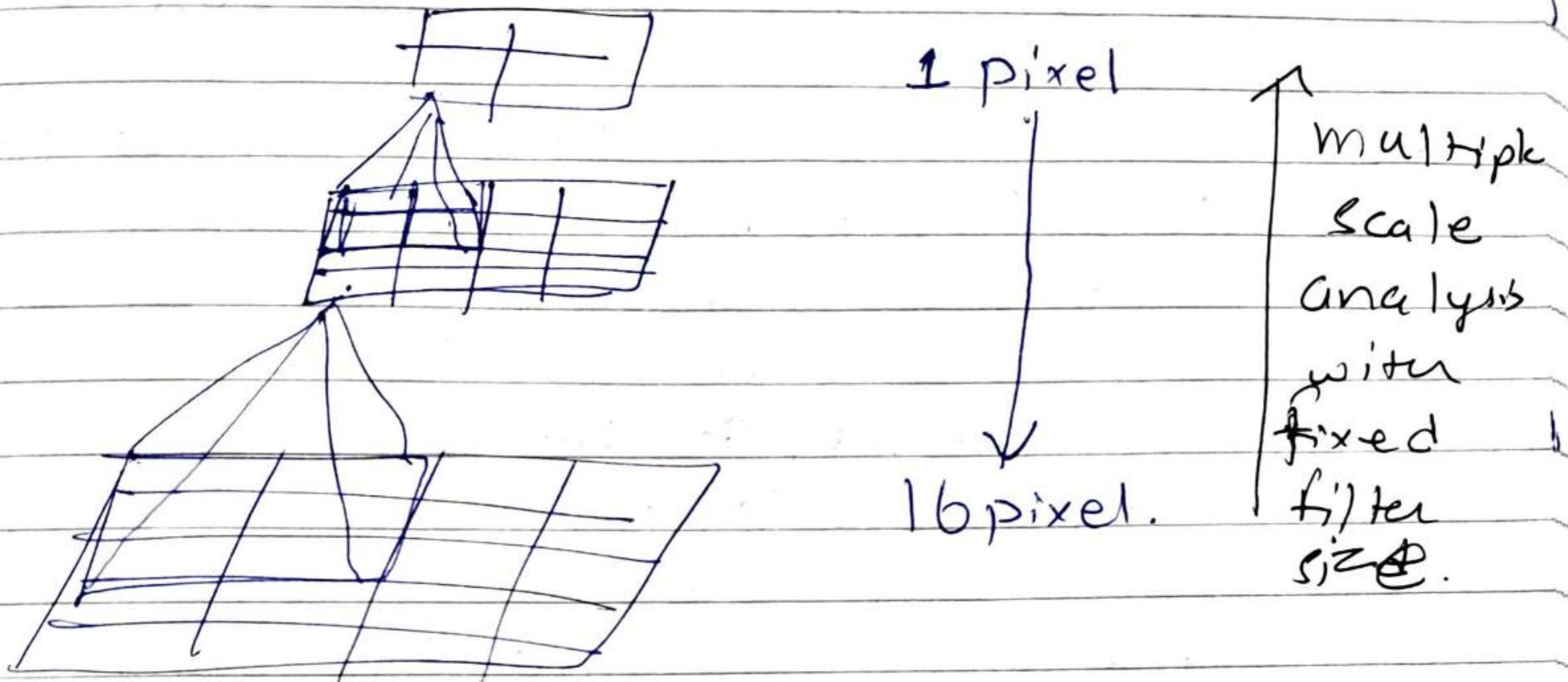
with the fixed size of filter we can do multiple convolution on the same image at different scale. like we have  $8 \times 8$  image and fixed filter of  $2 \times 2$ , after the 1<sup>st</sup> convo convolution we get an output of  $4 \times 4$ . after 2<sup>nd</sup> convolution on with same filter we get  $2 \times 2$ .

$$\left\{ \begin{array}{l} \frac{I}{8 \times 8} \times F_{2 \times 2} = O_1_{4 \times 4} \Rightarrow O_1_{4 \times 4} \times F_{2 \times 2} = O_2_{2 \times 2} \end{array} \right\}$$



In this way, as we go higher in the pyramid of Convolutions, our receptive field in the convolved layer corresponds to a large deep receptive field in the original image and so it accounts for differently scaled objects.

⇒ multiple scale analysis.



- ⑥ In each convolution layers, the size of the image reduces. This could lead to a loss of information as ~~ex~~ a result of a decrease in the spatial resolution. to compensate this we are adding the depth (no. of channel) in each convolution layers.

Eg ⇒  $\text{conv}_1 = 3 \times 3 \times \underline{64}$  → number of Channel in  $\text{conv}_1$  (Depth).

⑦  $\left( \text{produce size} = \frac{W - K + 2P + 1}{S} \right)$

$W \rightarrow$  input volume,  $K \rightarrow$  kernel size;  $S \rightarrow$  stride

$P \rightarrow$  amount of zero padding.

Given  $\Rightarrow W = 128$ ,  $P = 0$ ,  $K = 3$ ,  $S = 1$  (16-conv filters).

$$\text{produce size} = \frac{128 - 3 + 0 + 1}{1} = 126.$$

#  $\left\{ \text{Size of resulting tensor} = \underline{\underline{126 \times 126 \times 16}} \right\}$



⑧

Same as previous, this time  $S=2$ .

$$\text{Product size} = \frac{W - K + 2P}{S} + 1$$

$$= \frac{128 - 3 + 2(0)}{2} + 1 = 63$$

{Size of resulting tensor =  $63 \times 63 \times 16$ }

⑨

$1 \times 1$  Convolution  $\Rightarrow$  the advantage of  $1 \times 1$  Convolution filter is we can reduce the number of channel of given input without changing the image (input) size.

eg  $\Rightarrow$   $n - 1 \times 1$  filters through the image, where  $n$  is less than the current number of channels, we would get an output which is the same size as our input, but with ' $n$ '-channels.

⑩

The earlier layers in the convolutions are the Interpretation of Convolution layers: is used to features extraction of input image.

The earlier layers in the convolutions are responsible for learning lower level features extraction in the input image. like as edges, while the deeper layers learn higher-level, more complex features of the images, like eyes, ears or body parts in pictures of animals.

⑪

Red channel

1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

G channel

2	2	2	2
2	2	2	2
2	2	2	2
2	2	2	2

B channel

1	1	1	1
2	2	2	2
3	3	3	3
4	4	4	4

Result after pooling.

1	1
1	1

only companion

2	2
2	2

2	2
4	4



(12)

Purpose of pooling  $\Rightarrow$

The purpose of pooling is to down sample the spatial images of the input image, without having to learn new parameters as in the case of convolutions.

(13)

Data augmentation  $\Rightarrow$  we use data augmentation when we have limited data to use in model or network. it used to increase the dataset size, without actually acquiring new data. eg  $\Rightarrow$  if we have a single image source so here we can create multiple images from original image images by rotating or by moving in a frame frame at different location or shearing the corresponding image. Data augmentation can help to avoid overfitting by increasing the dataset size.

(14)

Transfer learning. can be thought of as a feature extraction mechanism in the cases when we have a smaller training dataset at our disposal. by using a model that has been trained on a larger dataset, we use the pre-trained model as an initial layer of our model to extract features. These extracted features can then be fed into a classifier that is customized for our use case.

(15)

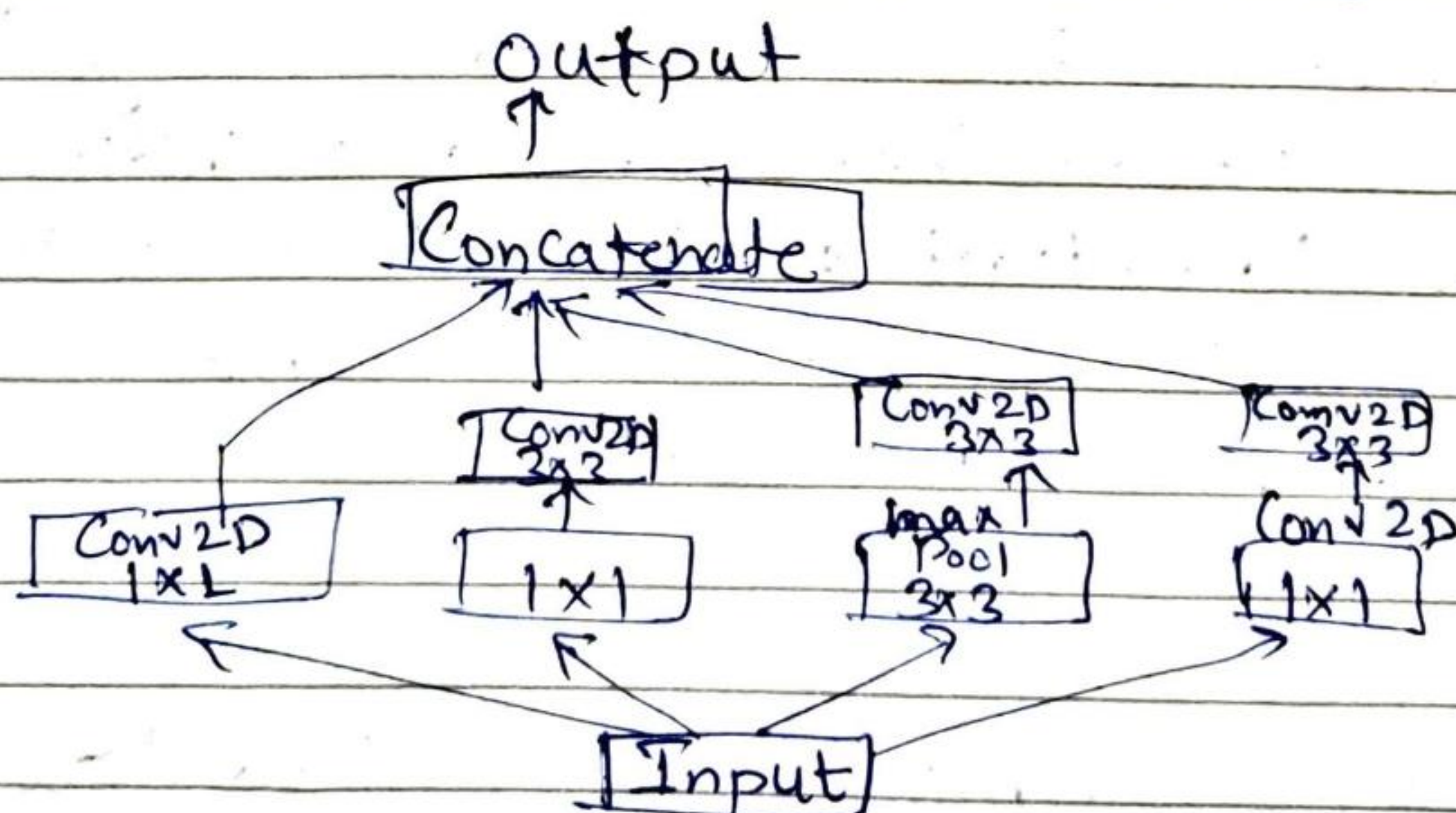
Freezing the coefficients of the pre-trained network.  $\Rightarrow$  we freeze the weights of pre-trained layer that came from transfer learning. because the weights are already trained and learnt. the top few layers of the model would be a custom classifier layers whose weights are initialized with small value or close to zero. So when get a gradient value that, for the classifier layers, would be



Small because the initial weights are small but the weights of the pre-trained layers ~~are not small~~ might be larger since they have already been learnt, which gives large gradient value. In this way we can destroy the weights that have already been learnt. To avoid this scenario the weights of the pre-trained model are frozen.

(16) After training our fully connected and classifier layers we can unfreeze few of the top layers of the "Conv-bn-pre-trained" model, and retrain the entire model to fine tune the weights.

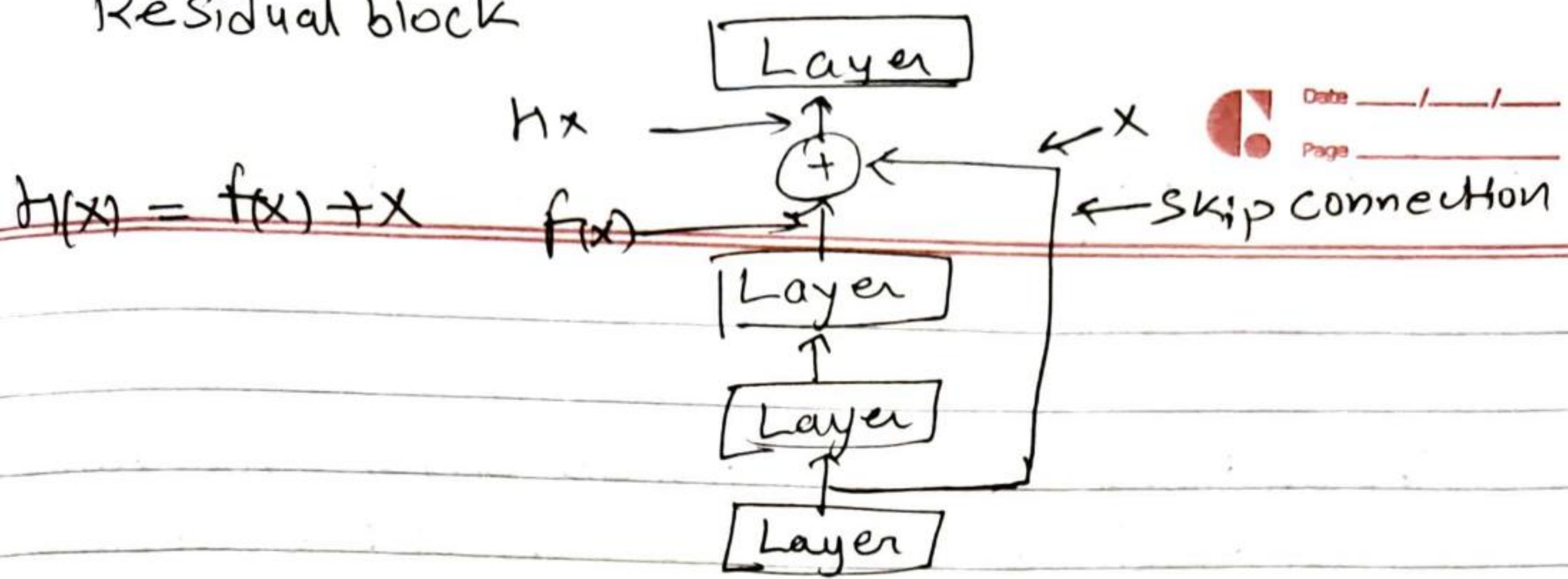
(17) Inception blocks  $\Rightarrow$  we convolve the ~~pre~~ previous layer with multiple filters. The outputs of these filters are then concatenated together. By doing so, we use more receptive fields, the concatenated output of which gives us higher dimensionality in the output. Also, by passing a 1-D convolution before each convolution block, we are able to reduce the number of parameters that need to be estimated.



(18) Residual block include a skip connection where an input layer  $x$  is added to the output of a subsequent layer. These skip connections help with vanishing gradients.



## Residual block



- 19) Once we have trained our convolution model we can create a new model that gives us multiple outputs - one output for each of the layers in our model. Plotting these individual outputs allows us to visualize the activations in each intermediate layer. Doing so allows us to visualize the templates that are learnt by the filters at each individual layer.
- 20) To visualize the filter weights of the trained convolution layers, we can use gradient ascent to find the input that will maximize the response at the filter. Since we think of filters as templates that are being matched to the image, visualizing these filter weights in this way gives us a sense of the pattern or template that each filter looks for in the inputs.
- 21) To visualize the heatmap of class activation, we can follow the below algorithm:
1. Feed an image to the network
  2. Compute gradients of the selected output with respect to each channel of the target layer for which activations are to be computed, say the final convolution layer.
  3. Compute the average gradients of each channel.
  4. Add the activations of each channel weighted by their average gradient magnitude.





Date \_\_\_\_/\_\_\_\_/\_\_\_\_

Page \_\_\_\_\_

5. Superimpose the activations on the original image to get the heatmap.

Visualizing such a heatmap gives us a chance to understand which features of the original image influence the final classification output. Moreover, heatmaps for misclassified images can point towards the problem points that have caused the error.