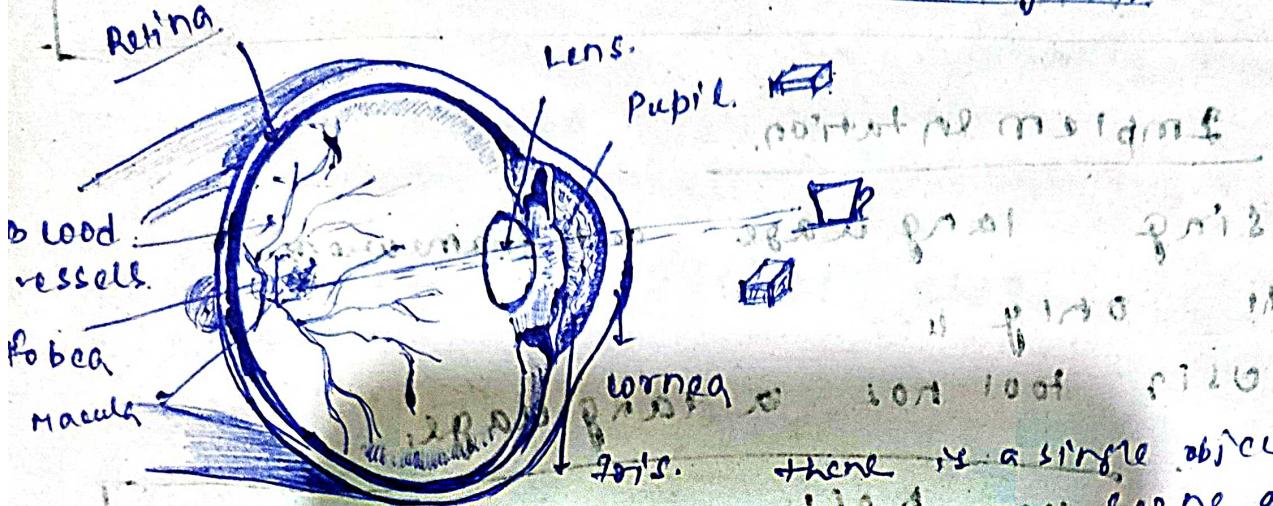


## Understanding attention modules.

### C-BAM and BAM

in paper ~~on~~ on human vision system.



whenever there is a single object in the scene, and lights from all the objects of interest ~~are focused~~ converge upon Macula. of Retina, which is primary functional region of Retina inside eye.

But, whenever a diverse set of objects are present in the scene, and lights from all the objects of interest ~~are focused~~ converge upon Macula. of Retina, which is primary functional region of Retina inside eye our visual perception system kind of use an advanced attention mechanism to focus on single object of interest.

so, we want to use ~~similar~~ kind of mechanism for our computer vision system, so that we need to incorporate Attention mechanism to our vision based deep learning models. And at the same time, the attention mechanism was first popularized by the paper of neural Machine translation, (Attention are you) need

for visual system if we use such attention mechanism the object of interest becomes in focus and surrounding objects/ backgrounds are blurred.

## What are spatial and channel attention modules

The words spatial and channel are most-referred terms in computer vision. These refer to convolution layers.

A convolution layer can be represented by a 3D volume of shape typically  $[c \times h \times w]$ .

where, each slice of the volume refers to a feature map or ~~feature area~~ extracted features are mapped over a depth of  $c$  channels.

$c$  total channels of convolution volume, (often can be considered total feature maps in the layer / total kernels used in the convolution layer)

### What is meant by spatial attention?

① Spatial means, the feature space / domain space encapsulated within feature map or layman word is, the space / domain space which contains the feature map or the feature space or domain space where the convolution features are mapped called feature map.

② Spatial attention presents other attention mechanism ~~mechanism~~ mask on the feature space which defines a single cross-sectional area, which is then domain of interest for classification.



for example in the given feature map  
if  $\rightarrow$  the snake is important object, it is important to attention,

so, spatial attention gives attention to that feature space  $\rightarrow$ , so, spatial features

(which represent the object) attention enhances these features, by generating a spatial mask, by refining spatial spatial attention mask, input to subsequent convolution layers, which improves performance of the model.

What is channel attention.

Usually, ~~filters~~ from ~~convolution~~ layer ~~1~~ have ~~large~~ weights inside the kernels, which make the feature maps work down to the convolution volume and very small close to zero, so why various feature maps containing the convolution volume looks like having similar information, seems like appearing copies of others.

Although, they look similar, but these filters are extremely useful for learning different types of features in images like horizontal, vertical edges, diff. patterns, textures in image progression. But, channel attention provides weightage or influence parameter to the channel which is essential for or relevant for learning the classification contributed to perform the task given to the network and boosts the overall model performance (channel attention enhances the importance of channels towards learning both, ~~horizontal~~ and ~~vertical~~ ~~the effectiveness~~)

Authors show that both together worked better as compared to individual attention mechanism.

### BAM

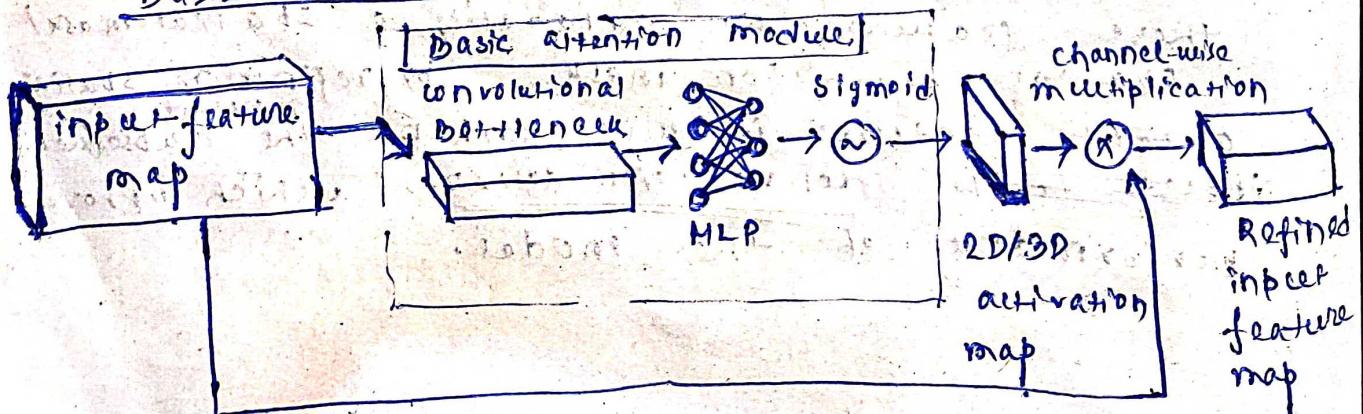
## Understanding Attention Modules: CBAM and BAM

Recently different SOTA networks have leveraged these attention mechanisms that have significantly improved and refined real-time results.

### Attention Modules

Attention modules are used to provide attention on relevant ~~background~~ info in convolution volume, instead of unimportant background info ~~not useful~~ for the task given to the CNN.

### Basic Structure



① It takes  $(C \times H \times W)$  input feature map as input,  $(H \times W)$ . Output: feature map or  $(C \times H \times W)$ , a feature map or attention map. This attention map / attention mask is multiplied with input feature map to get refined output feature map, where relevant information are enhanced as compared to unimportant info.

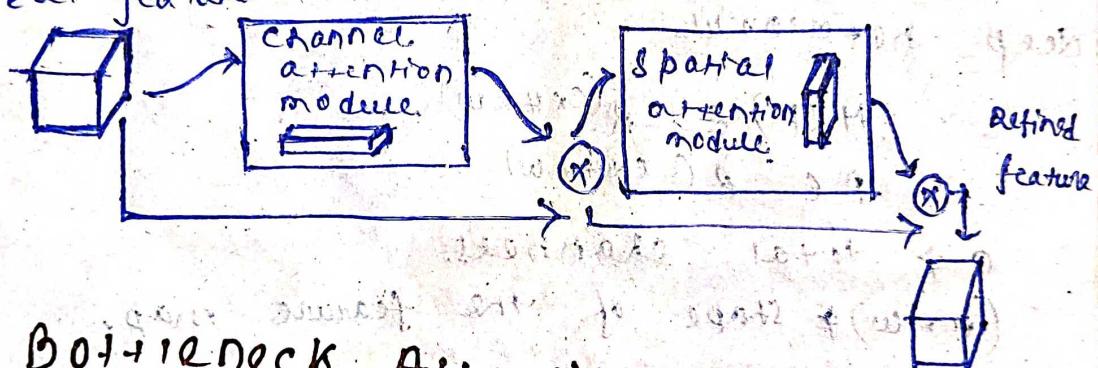
② Generally, attention mechanisms are applied to spatial and channel dimensions.

These two attention mechanisms can be applied sequentially or parallelly.

③ Attention mechanism was first experimented in residual architectures.

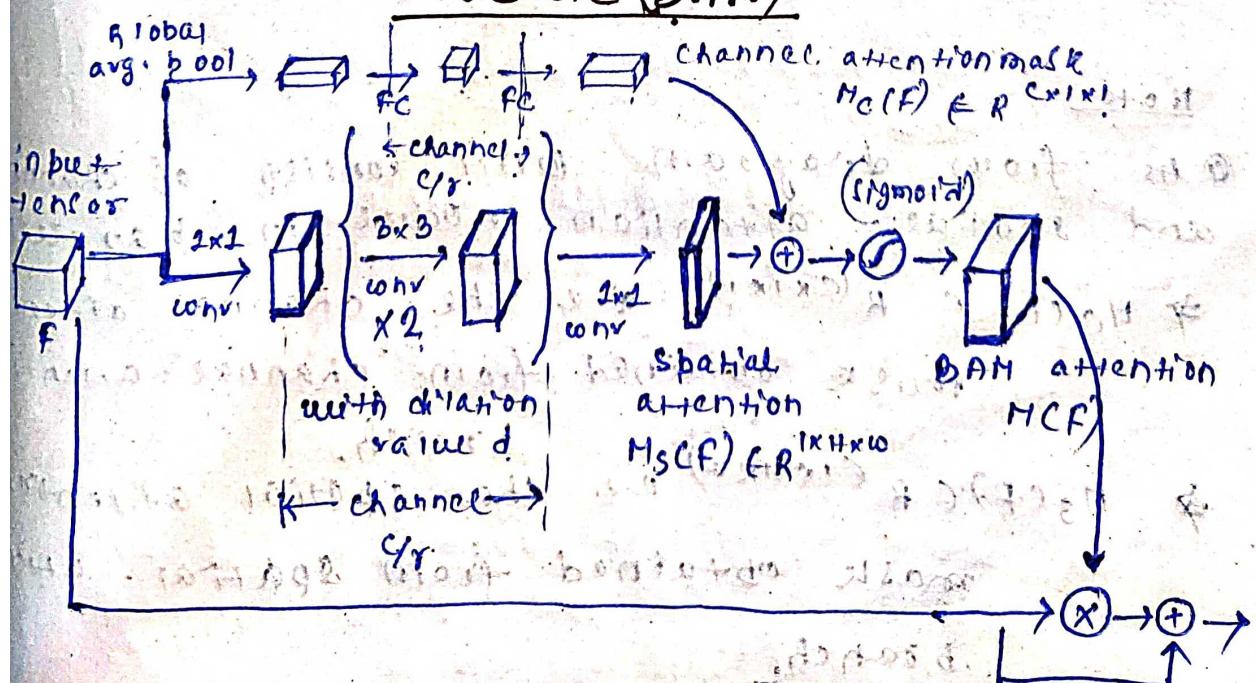
### Overall structure (BAH)

Input features



### Bottom Neck Attention

#### Module (BAH)



[structure of BAH]

$f \Rightarrow$  Input feature map

$f' \Rightarrow$  Output refined/enhanced input feature map.

$$f' = f + f \otimes M(P)$$

$\otimes \Rightarrow$  elementwise multiplication

$M(P) \Rightarrow$  attention mask generated as a function of  $f$  from (BAM)

the step  $f + [f \otimes M(P)]$  acts as the residual step where we add the output of previous layer to the input to next layer so as to reduce the error due to vanishing gradient descent and ensure the ~~propagation~~ flow of information over the

deep networks

$$M(P) \in R^{(C \times H \times W)}$$

$$f \in R^{(C \times H \times W)}$$

$C \Rightarrow$  total channels

$(H \times W) \Rightarrow$  shape of the feature map

### (BAM) diagram

#### Note:

① As from diagram BAM consists of channel and spatial attention module in parallel.

$\Rightarrow M_C(f) \in R^{(C \times H \times W)}$  i.e. the channel attention mask obtained from channel attention branch.

$\Rightarrow M_S(f) \in R^{(C \times H \times W)}$  i.e. the spatial attention mask obtained from spatial attention branch.

④ And sigmoid is used to obtain the HCF) attention mask from combination of  $H_C(F)$  and  $H_S(CF)$ , thus

$$HCF = \sigma(H_C(F) + H_S(CF))$$

as,  $\left\{ \begin{array}{l} H_C(F) \in R(C \times C) \\ H_S(CF) \in R(C \times H \times W) \end{array} \right\}$  but,  $HCF \in R(C \times H \times W)$

⑤ before addition to obtain  $H_C(F)$ ,  $H_S(CF)$  are resized to  $(C \times H \times W)$  (by method of broadcasting will be discussed later)

⑥ Note here, to obtain  $HCF$ , to combine  $H_C(F)$  and  $H_S(CF)$ . for sigmoid, we choose elementwise addition over multiplication due to smooth gradient flow. also, the time of backprop.

what does that mean,

if it's addition,

$$\frac{\partial L}{\partial H_C(F)} = \frac{\partial L}{\partial HCF} \cdot \frac{\partial HCF}{\partial H_C(F)}$$

$$= \frac{\partial L}{\partial HCF} \cdot 1 \cdot \cancel{\frac{\partial HCF}{\partial H_C(F)}}$$

but, if multiplication,

$$\frac{\partial L}{\partial H_C(F)} = \frac{\partial L}{\partial HCF} \cdot H_S(CF) \cdot 1$$

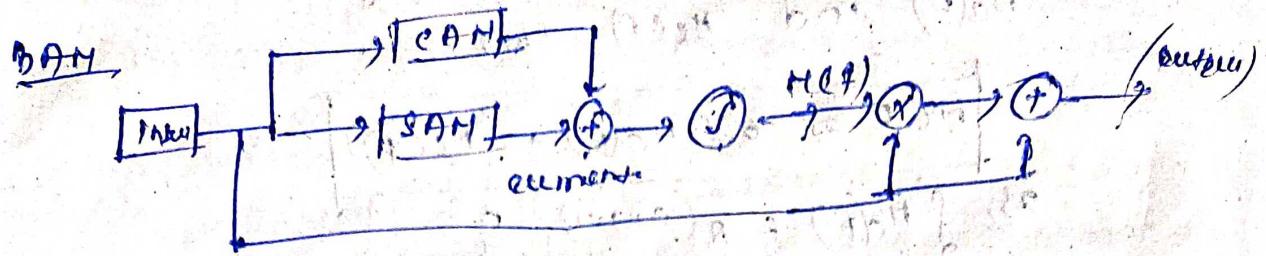
so, this gradient term also depends on

$H_S(CF)$ , this is not robust &

$H_C(F)$  term, only, i.e., it's not smooth flow

and grade deriv in  $H_S(CF)$  branch (maybe wind, Adenitic flow) or turbulent flow in backprop

## Channel Attention Module



### Channel Attention Module

channel attention modules takes the input feature map  $F \in R^{C \times H \times W}$  and generates channel attention mask  $H_C(F) \in R^{C \times 1 \times 1}$

steps to generate channel attention map are ..

1. Do global average pooling of feature map  $F$  and get a channel vector  $F_C \in R^{C \times 1 \times 1}$

2. Pass  $F_C$  to small MLP of one hidden layer followed by one output layer and hidden layer has dimension  $(C/r)$ .  $r$  is the reduction ratio for that.

Channel attention branch.

so, if  $F$  has 1024 channels and reduction ratio is 16, then

total neurons in hidden layer is

$$(1024/16) = 64$$

(BN)

3. Add batch norm layer at the from of this MLP.

$$4. H_C(F) = BN(\text{MLP}(\text{AvgPool}(F)))$$

$$= BN(W_2(\text{ReLU}(\text{AvgPool}(F) + b_0) + b_1))$$

where  $(w_0 \in R^{C \times r \times c}), (b_0 \in R^{C/r})$ ,

$$\underline{W_1 \in R^{C \times C/r}}, \underline{b_1 \in R^C}$$

we used reduction ratio as ~~16~~ because

MLP acts like an auto encoder and the manifold learning is done

and hidden layer acts as latent space

if hidden size =  $C$  then it encodes

acts as identity learning and will not learn 10D dimensional manifold so, it avoided because overfitting from the perspective of the ~~softmax~~ channel activation branch.

### Spatial Attention Map (SAM)

Spatial attention map is used to provide attention in spatial domain,

so, takes input  $F_E (C \times H \times W)$

or  $F_R$

$$H_S(F) \in (1 \times H \times W)$$

1. Input feature map  $F_E$  is passed through chains of  $(1 \times 1)$ ,  $(3 \times 3)$  convolution layers.  $(3 \times 3)$  convolution layer consists of dilated convolution having a dilation rate ( $d=4$ ). Dilation is used to increase effective receptive field (for dilated convolution note), these  $(3 \times 3)$  convolution layers containing  $(C/8)$  ~~modulated~~ kernels.

2. after that  $(1 \times 1)$  convolution layer, it generates a spatial attention mask.

Now,  $R_F (C \times H \times W)$  passed through  $(1 \times 1)$  convolution layer to generate spatial attention mask,

$(C/8)$  reduced channel helps as ~~latent space~~ to manifold latent space and similar task features extract and reduce total number of parameters

3. Add batch normalization at the front of convolution layers stack  $f_2$

$$MS(f_2) = BN(f_2^{(1 \times 1)}(f_2^{(3 \times 3)} / \text{dilate} (f_2^{(1 \times 1)}(F)))$$

then we combine

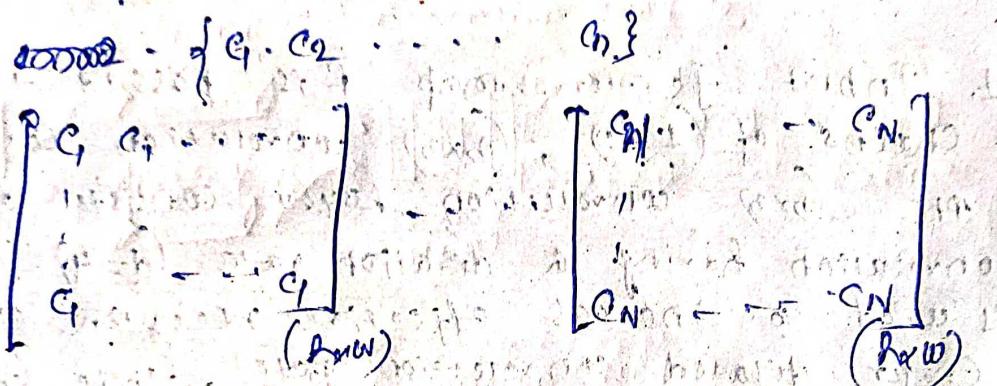
$$MCP_2 \rightarrow f(Mo + R)$$

$$R = F + F(R) MCP$$

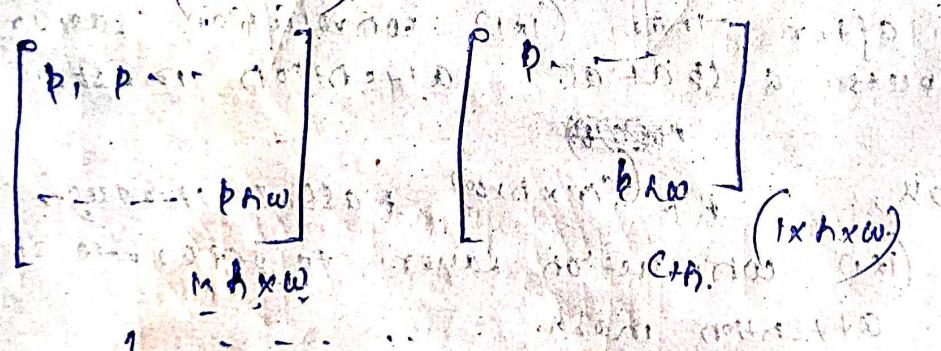
as,  $R$  has structure followed by the  $Mo$  and resembles the Raster family, as it was implemented by Rosner.

Now  $MCP \rightarrow MCP(R)$

i)  $Mo(F) \& R(C_{xw}) \Rightarrow$  broadcasted to,  $(R_{xw})$



ii) Similarly,  $Mo(F) \& R(C_{xw})$  broadcasted to  $(C_{xw})$



so, ~~they~~ either they are added

for  $M's$   $\rightarrow$  all the channels have  $(Cxw)$  spatial attention along the focus or enhance on the relevant features on feature map.

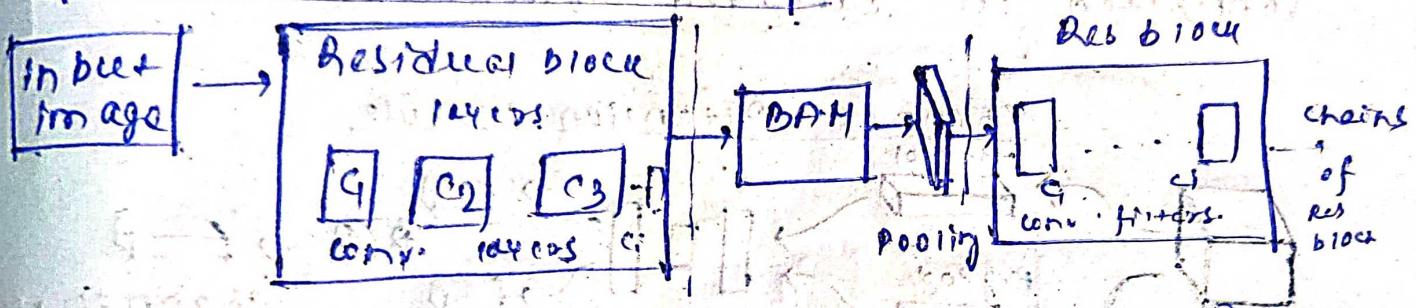
and on other hand represented like

$Mo(C_{xw}) \rightarrow$  focus on the channel which has highest priority for the goal to achieve by NN.

so, when they add as a combined  
use obtain a convolution volume  
(ex fixe) and which

has, feature map, where spatially attention  
or enhanced on relevant features along  
with feature maps having the weights  
enriched based on their relevance.

## Where to keep BAM?



{ Res. block  $\xrightarrow{\text{wts}} \text{BN} \xrightarrow{\text{activation}} \text{pooling}$ }  
but skip connection

After each conv (residual block) BAM is used  
to ~~attenuate denoise~~ denoise low level  
features. At early stages and at then at  
deeper layers it helps to capture or focus  
on exact target which has a high range of  
semanticity

{ semanticity  $\equiv$  meaningfulness }

{ semantic segmentation  $\Rightarrow$  fine classifica- }  
which is more meaningful for that class

semanticity  $\Rightarrow$  the quality that a linguistic  
form has ~~of being~~ so that it can be able to  
convey meanings, in particular by ~~means~~  
reference to the world of physical reality

Inform to  
help me

(1)  
II

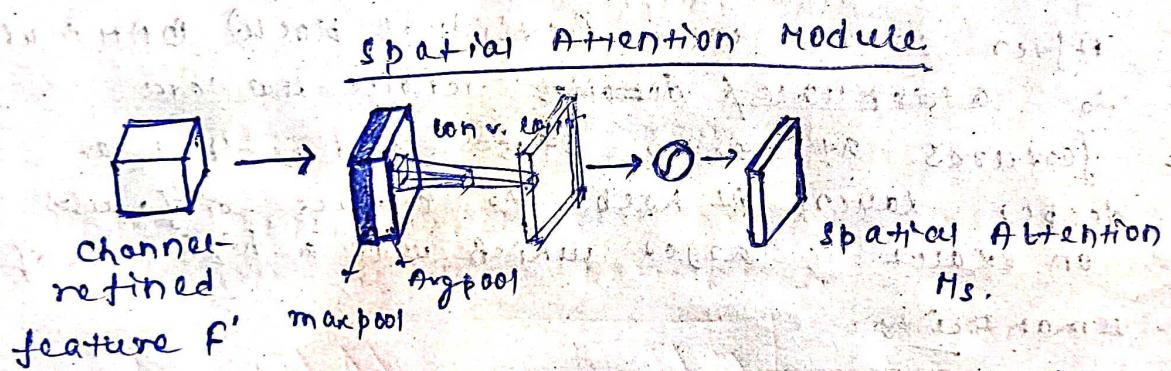
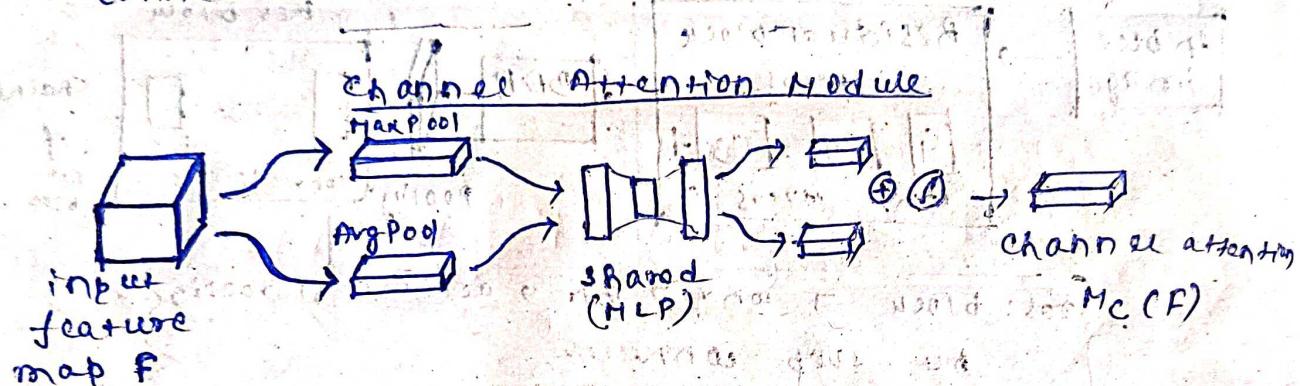
Implement CBAM from AD3 to AP3

PAP has 8 TDS we can't do PAP for I transform

shift change

## Convolution Block Attention Module (CBAM)

CBAM consists of a sequence of channel attention module and spatial attention module connected in cascaded structure.



An intermediate feature map  $\mathbf{M}_C F' R$  ( $C \times H \times W$ ) is passed to a channel attention module.

$$\text{Note} \quad [F' = M_C(P) \otimes F]$$

and,  $F'$  is then screen+tracey passed to a spatial attention module.

$$\text{Note} \quad [F'' = M_S(P') \otimes F']$$

$M_C(F) \in R$  ( $C \times H \times W$ ) } with  $M_C$  multiplication, we  
 $M_S(F) \in R$  ( $1 \times H \times W$ ) need

$$M_C(F) |_{\in R} (C \times H \times W)$$

$$= M_S(F) |_{\in R} (C \times H \times W)$$

so,  $M_S(F)$  and  $M_C(F)$  are then broadasted to respective dimensions.

(ii) element-wise multiplication.

### Channel attention map

follows same generation process as DAT.

Here we pass two vector:

- ① Arg. pool of  $F^c$
- ② Max pool of  $F^c$  } instead of direct global arg. pool of DAT

then, Argpool( $F^c$ ) passed to HLP( $C$ )

also Maxpool( $F^c$ ) passed to HLP( $C$ )

so here Argpool and Maxpool both are added to obtain more distinctive ~~feature~~ channel feature.

$$\begin{aligned}H_C(F^c) &= \text{G}(HLP(\text{Argpool}(F^c)) + HLP(\text{Maxpool}(F^c))) \\&= \text{G}(\omega_0(\omega_0(F_{\text{arg}})) + \omega_1(\omega_1(F_{\text{max}})))\end{aligned}$$

$\Rightarrow$  sigmoid function.

$$w_0 \in (C/2 \times C), \quad w_1 \in (C \times C/2)$$

ReLU activation is followed by  $w_0$ ,  
here similar prediction ratio can be  $1:13$   
used to obtain a encoder-like structure to  
obtain a non-linear information of the latent space.  
or,  $w_0$  is layer of  $(C/2)$  nodes.  
and these transformation is non-linear so  
it can understand non-linear features or  
non-linear information which is instilled inside the  
data, thus helps to denoise the  
array of  $(C \times 1 \times 1)$  vector to obtain  
influence or effect of different channels.

### Spatial Attention Map

Take input feature maps  $F$  and generate two intermediate feature map  $[F_{\text{arg}}$  and  $F_{\text{max}}]$

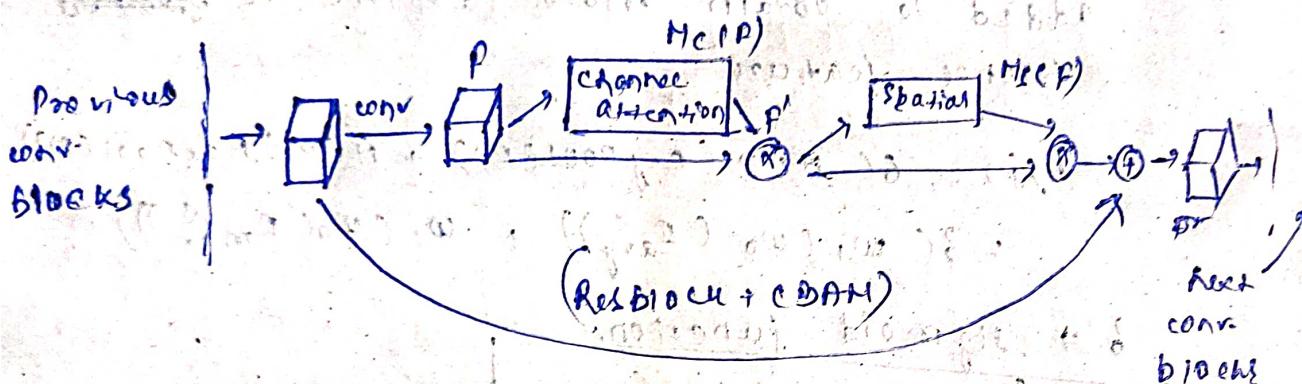
$F_{\text{arg}} \Rightarrow$  average pool }  $\{$   $\text{G}(F^c \text{ for } C \times H \times W)$   
 $F_{\text{max}} \Rightarrow$  max pool }

1. concatenate ArgPool and MaxPool, and pass through a small convolution block of kernel size  $(7 \times 7)$  with dilation = 2

Note: Unlike CBAM, where to increase receptive field used  $d=4$ , whereas here to do so  $(2 \times 2)$  kernel used.

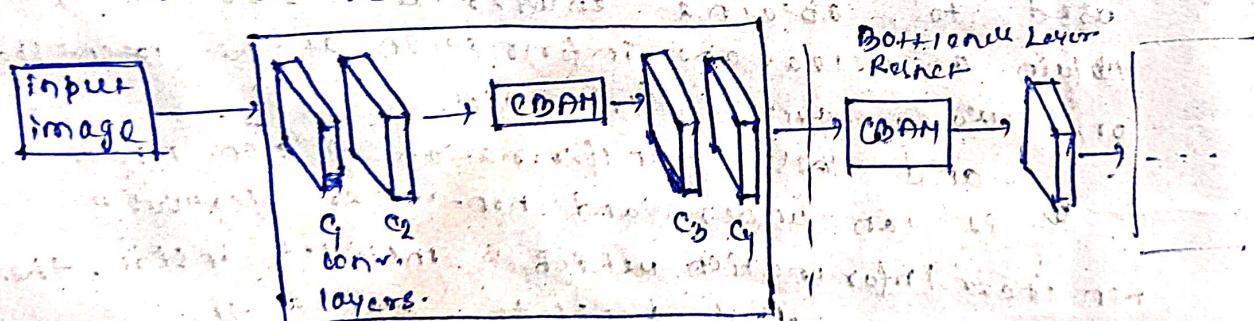
$$\begin{aligned} \text{MSCF} &= \delta(f^{\text{Arg}}(\text{ArgPool}(P); \text{MaxPool}(F))) \\ &= \delta(f^{\text{Arg}}(f^{\text{Arg}}; f^{\text{max}})) \end{aligned}$$

How to use



where to keep CBAM

Unlike BAM, CBAM are placed inside ResNet block as well as the bottlenecks.



what is the difference between BAM and

CBAM

- ① In case of BAM, APC (global average pooling) was used to obtain global statistics over the channels of feature map in that channel attention, whereas in CBAM, maxpool and avg pool used and it was observed that and avg pool combination accounts to generate more salient features from feature map and correlate LRP outputs and exactly global statistic soft.

and max and avg pool as combin petral  
more salient features from feature map and  
concatenate GAP o/p and obtain globalist softcy

- i) In DAN to increase receptive field in  
(SAM), dilated conv used with  $d \geq$ ,  
same in COAH is done by using kernel size  
 $(7 \times 7)$ , normal convolution layer with  
 $d = 1$ .
- ii) In DAN, CAN and SAM are done in  
parallel and then combined by  
element-wise addition followed by  
sigmoid activation, whereas COAH  
CAN and SAM were used sequentially.  
- and at each step since it element-wise  
multiplied with CAN and SAM  
attention mask (Same as DAN)