

Image & Vision Computing  
Comprehensive Exam

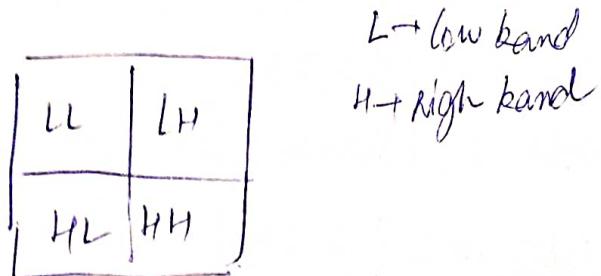
Anubhooti Jain  
D20 CS002.

①

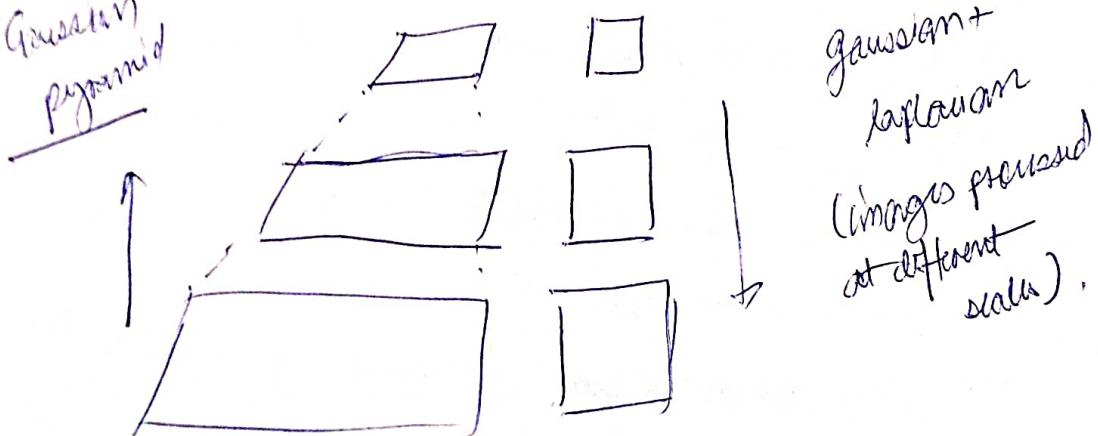
Ques:  
(a) Gaussian Pyramid, DCT, and DWT

Q1(1) Gaussian pyramid is used for processing images at different scales while DCT and DWT are transformations in frequency-band. DCT is direct cosine transformation and DWT is Discrete Wavelet transformation. So, the transformation procedure is completely different.

for DWT



Gaussian  
pyramid



difference ② The information processed

Gaussian → spatial information is more visible and effected as images are upsampled & downsampled.

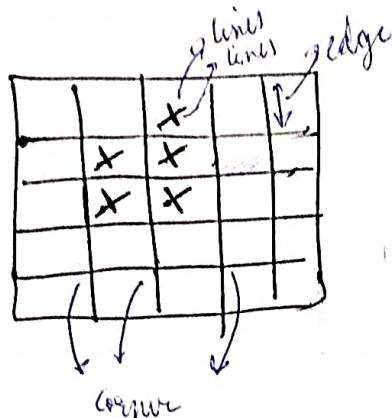
DCT → Information is more visible in frequency domain but changes are visible in spatial domain as well.

DWT → Frequency-based information. We can manipulate low and

⑨

High frequency bands to manipulate the spatial domain.

### Ques 1. (b) Corners, lines, edges



In a simple checkboard-type image corners, edges and lines can be detected as shown.

### for edge detection

- we have filters available for edge detection. It is known that first order derivatives can help detect an edge

- An edge is noted by immediate change

detecting : line b/w background  
edge and circle.



boundary of circle is  
an edge

→ there will be an immediate change  
in the gradient information

→ we can detect it using first order differentiation

\* Known filters & Sobel filters

↳ they are used for edge detection

↳ a better known version is the  
3x3 sobel filter.

For all 3 - corner, edge, line → we have filters available. Sobel filters are well known for edge detection. There is also Hough transform, though, they are better at detecting shapes.

(2)

## Line detection

As said, Hough transform can be used for line detection. This transforms the parameter space  $[x, y \rightarrow \theta, r]$  and the entire image is changed. Intersections of several points indicate existence of lines. we can easily detect edge based on that.

## Corner detection

Harris corner Det. Detection is one well-known algorithm for corner detection. It usually has a reference image (a typical deckboard) with respect to which we can extract corners. It also works on the differential gradient with direction to detect corners.

Ans 2 (ii)  $f = \begin{bmatrix} 130 & 11 & 67 \\ 10 & 10 & 50 \\ 80 & 90 & 100 \end{bmatrix}$

The center value changes with mean and median filter.

### Mean filter

$$= \frac{1}{9} \times (130 + 11 + 67 + 10 + 10 + 50 + 80 + 90 + 100)$$

$$= \frac{1}{9} (540) = 60.0 \approx 61.$$

Mean  $\rightarrow 61$   
Median  $\rightarrow 61$

### Median filter

- odd size

- 10, 10, 11, 50, (67), 80, 90, 100, 130

$\downarrow$   
 $\therefore 67$

Aus 1 (d)

$$I = \begin{bmatrix} 2 & 1 \\ 3 & 1 \end{bmatrix}_{2 \times 2}$$

(4)

raise by 3 times  $\rightarrow$

(i) bilinear interpolation

$$\begin{array}{cc} 2 & 1 \\ 3 & 1 \end{array} \xrightarrow{\begin{array}{c} 0.2 \\ 0.4 \\ 0.6 \\ 0.8 \end{array}} \begin{bmatrix} 2 & 0.2 & 0.4 & 0.6 & 0.8 & 1 \\ - & - & - & - & - & - \\ - & - & - & - & - & - \\ - & - & - & - & - & - \\ - & - & - & - & - & - \\ - & - & - & - & - & - \end{bmatrix} 6 \times 6$$

$$= \begin{bmatrix} 2 & 1.8 & 1.6 & 1.4 & 1.2 & 1 \\ 2.2 & 1.96 & 1.72 & 1.48 & 1.24 & 1 \\ 2.4 & 2.12 & 1.84 & 1.56 & 1.28 & 1 \\ 2.6 & 2.28 & 1.96 & 1.64 & 1.32 & 1 \\ 2.8 & 2.44 & 2.08 & 1.72 & 1.36 & 1 \\ 3 & 2.6 & 2.2 & 1.8 & 1.4 & 1 \end{bmatrix}$$

$2 \times 0.8 + 5 \times 0.2$

(ii) bicubic Interpolation

(5)

Ans 2 H1:  $f(x) = [60 \underbrace{60}_{1} \underbrace{60}_{1} \underbrace{100}_{1} \underbrace{100}_{1}]$

First and second order derivative.

1<sup>st</sup> order  $\frac{\partial f}{\partial x} = f(x+1) - f(x)$ .

$$\frac{\partial f(n)}{x} = [0, 0, 40, 0, 0]$$

2<sup>nd</sup> order

$$\frac{\partial^2 f}{\partial x^2} = f(x+1) + f(x-1) - 2f(x) \quad [60, 60, 60 \underbrace{100}_{1}, \underbrace{100}_{1}, 100]$$

$$= 60 + 60 - 2 \times 60 = 0$$

$$= 60 + 100 - 2 \times 60 = 160 - 120 = 40$$

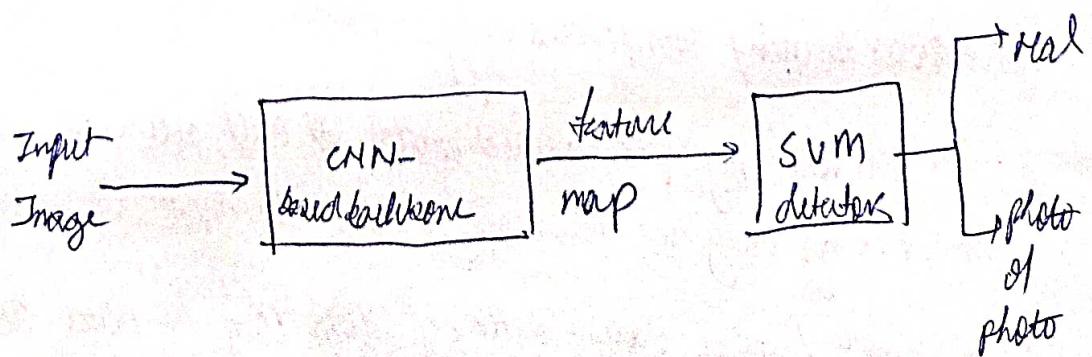
$$= 60 + 100 - 2 \times 100 = 160 - 200 = -40$$

$$= 100 + 100 - 2 \times 100 = 0$$

$$\therefore \frac{\partial^2 f}{\partial x^2} = [0, 40, -40, 0]$$

Ans 2 H1 Defecting photo of photo

Def learning  
problem



(6)

① The deep learning architecture has 2 components

↳ (a) CNN-based backbone — this is used for feature extraction.

(b) SVM detector — this is used to classify between whether the image is real or a photo of a photo.

↳  $I \rightarrow$  features extracted using CNN  
 $O \rightarrow \begin{cases} 0 \text{ for fake/photo of photo} \\ 1 \text{ for real} \end{cases}$

\* Reason : — SVM detectors can learn the boundary between this binary classification task with good accuracy, limited examples, and diverse kernel choice — linear, poly, RBF etc.

for CNN based architecture → we can use a well-known model like VGGNet or ResNet18 or ResNet50.

Reason : — breaking the image into feature maps ~~allows~~ allows to extract important information out of a sparse input like an image. It reduces dimensions of the input without any loss ~~in~~ of crucial information hence making SVM training computationally faster.

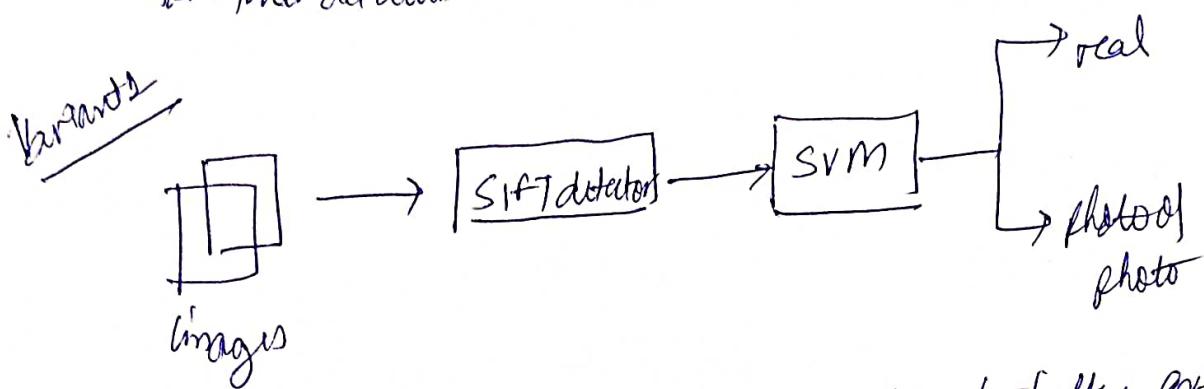
- we can use any pre-trained model. It will allow to reduce training time having only to finetune on our dataset.

Loss : a simple binary cross-entropy loss can be used or better use pre-trained backbone with SVM and its own loss.

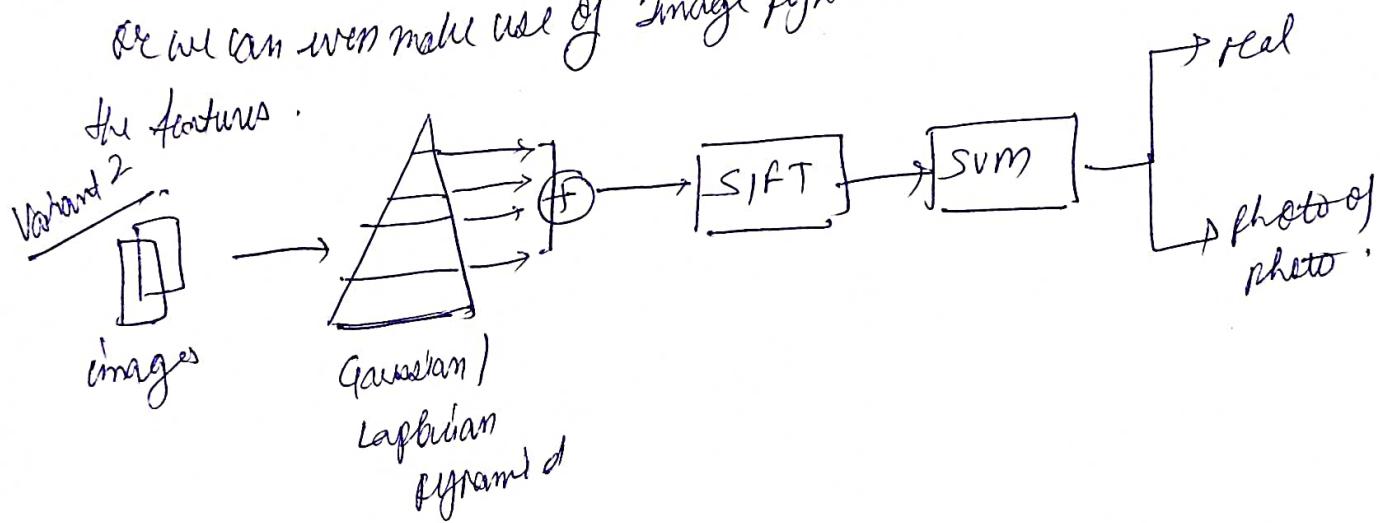
## ⑤ The hand-crafted algorithm

7

↳ we can make use of SIFT features before using SVM detector for final detection.



or we can even make use of image pyramids to better extract the features.



Continuing with Variant ①, we can extract features just like we were doing in deep learning architecture by replacing CNN-based model with SIFT features. SIFT are hand-crafted and uses keypoint detection to extract features. Holm is also ~~an~~ a possible option with gradient & direction and considered in its calculation. The goal is to extract image feature before sending these features to an SVM model for detection training.

## Ques 2: (a) Vanishing Gradient problem

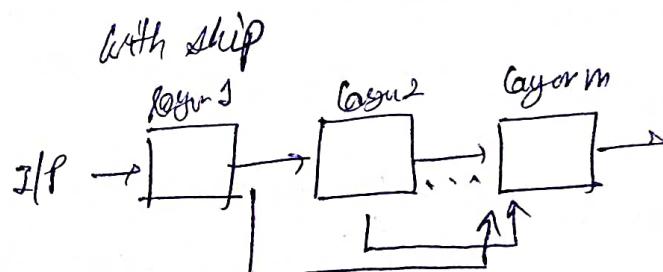
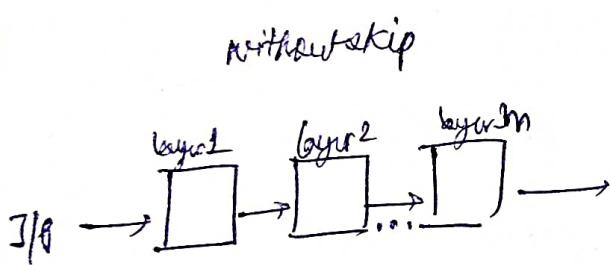
⑧

Vanishing Gradient problem arises as depth of deep neural networks increases. The gradients (smaller values especially) cease to exist before they reach last layers thus losing information and reducing the performance of models.

In context of object detection, we have object detectors which are deep to be able to get image features, thus deep models can easily face the said issue of vanishing gradient.

Solutions to Vanishing gradient problem →

→ skip connections / residual connections



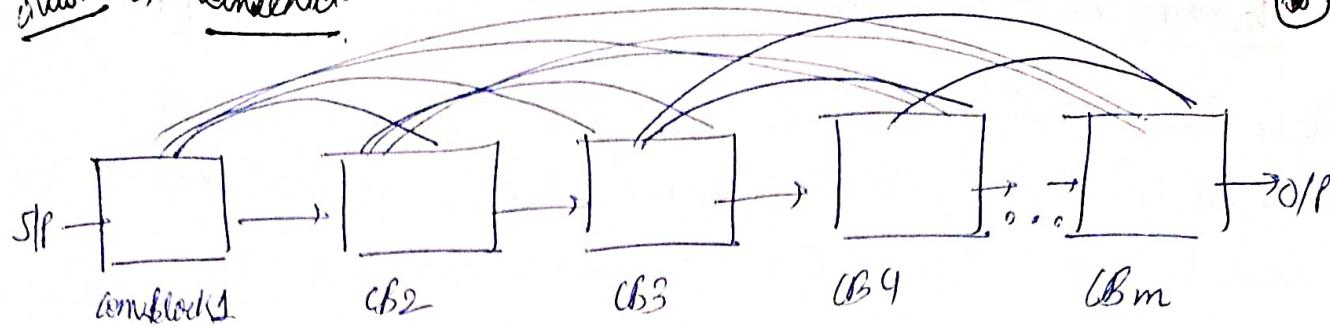
- \* skip connections allow features from previous layers to reach as it is, so as to avoid vanishing gradient problem.
- \* we can have random layers provide their output to further layers skipping one or several layers in between.

models like Resnet and DenseNet uses skip connections. In fact, skip connections were introduced with ResNet model when performance on ~~the~~ models with high depth was suddenly. With ResNet, authors were able to mitigate this issue.

Further, DenseNet went ahead to use all layers being connected to all the other layers. Fully connected Model is also a good example of skip connection.

Ex2 (b) DenseNet

(8) (9)



$CB \rightarrow$  convolution blocks.

With DenseNet, the skip connections were spread to all the blocks in the model architecture. Each block is connected to every other block (as in the figure above). These blocks are also known as dense blocks and is used for Image classification. It can also be utilized to get feature maps by removing the classification layer at the end or pulling maps from later layers.

There are also bottleneck layers in the model. Bottleneck layers have  $1 \times 1$  convolutional layers (that help with dimension reduction) and also uses ReLU activation function with average pooling and batch Normalization.

### Advantages of DenseNet

- ① There is better gradient flow than the models proposed before DenseNet. With skip connections, vanishing gradient issue is mitigated, moreover, by storing it between all the dense blocks (possibly with same size feature maps) it allow for better flow.
- ② With  $1 \times 1$  convolutional layers in the model, the number of parameters are reduced as opposed to other equivalent deep models with similar depth.

(3) Better feature maps  $\rightarrow$  as all the layers are inter-connected, the output from initial layers travel quite far, making the resultant feature maps better than those generated by previous models.

(10)

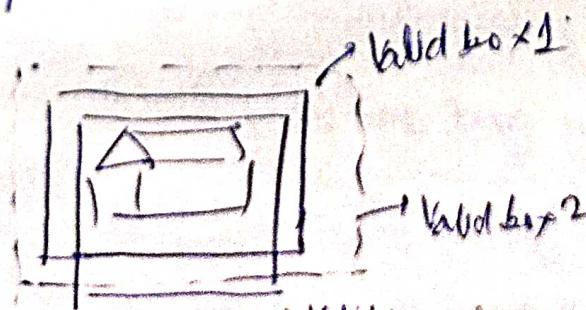
Transition layer: Transition layer has  $1 \times 1$  convolutional layers followed by an average pooling layer. It is used between some blocks so that feature maps can flow easily from one block to another.

### Ques 2 (c) Mean Average Precision (mAP)

mAP is used as an evaluation metric in an object detection algorithm. It allows to compute how accurately a bounding box is predicted by an object detection algorithm. Higher the mAP, better the performance of algorithm.

mAP is dependent of the IoU threshold. IoU  $\rightarrow$  Intersection over Union is a way to compute the overlapping of 2 boxes. How much should the two boxes overlap in order to count as a valid prediction is dependent on the IoU threshold.

For example:



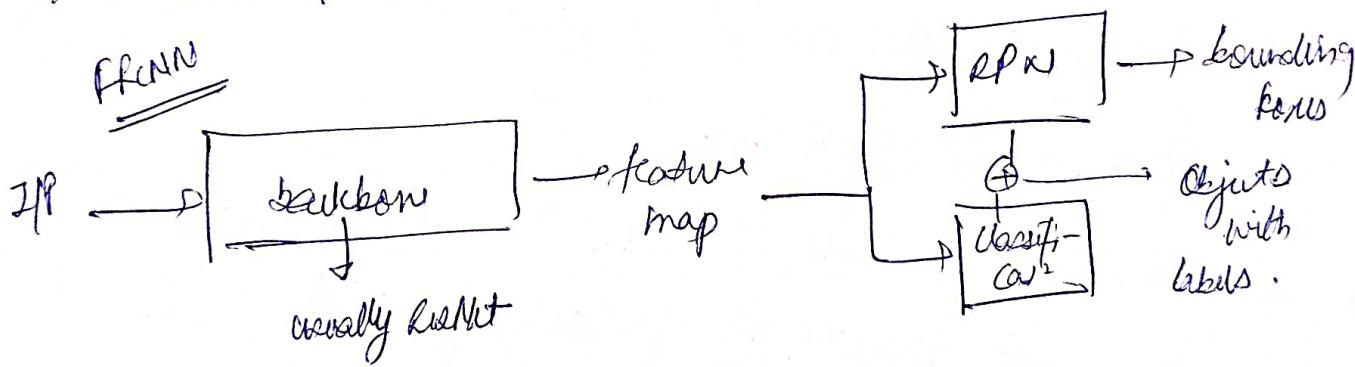
$\rightarrow$  All 3 boxes are valid, however they do not completely overlap, they intersect at different ratios which is computed using IoU.

(11)

The threshold can be set to as high as 0.95 or at an coverage of 0.5. mAP changes based on this threshold and is therefore usually used as an average over several threshold values when reported in a research.

## Ans 2 (d). RNN Vs faster RNN

Both RNN and faster RNN are object detectors. They are used for object detection algorithm with both considering object detection as a regression problem. They are classified as 2-stage detectors. The difference between the two is their architecture. Faster RNN proposed a Region Proposal Network (RPN) to predict the bounding boxes. It uses deep Neural Networks based RPN to predict the boxes which makes it faster and more accurate than the RNN.



RPN is the novel module in faster RNN. To go on the other hand visualizes the object detection problem completely differently. It divides the entire image into a grid and using grid as anchors, it predicts the bounding boxes for objects. It is a 1 stage detector. That is, it is better in terms that there is end-to-end training making the model faster & robust.

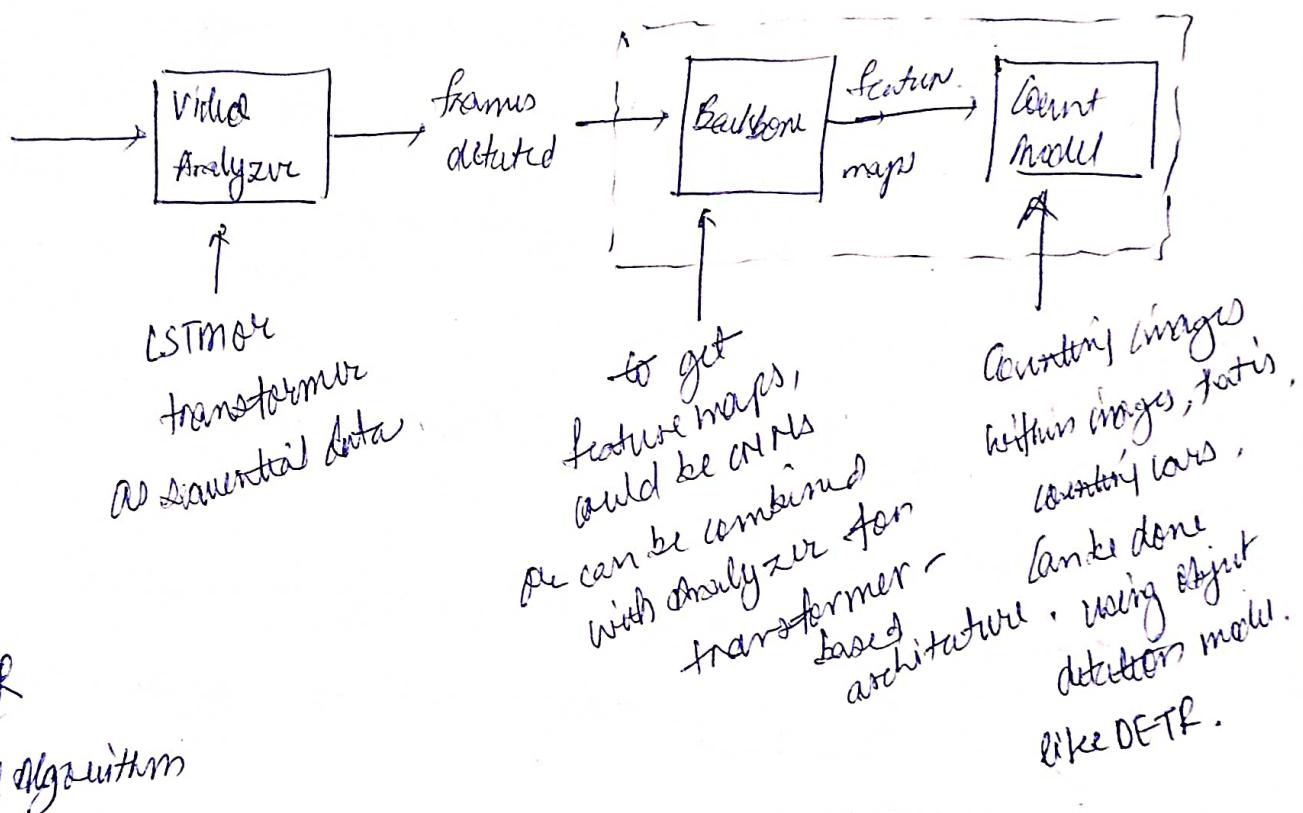
## Ques 1(e) Automated Video Surveillance

(12)

- ↳ Counting cars at toll plaza
- ↳ Motion detection.

\* We need tracking algorithm → that is detection + prediction. That is, frame is predicted and corrected with respect to previous frames.

### A system



A potential challenge is a lot of frames will have to be processed, with different traffic density at different times of the day. We will need a time-analyzer to detect this change / predict this change so as to change the model based on the time of the day.

### Ques 3, (a) PCNN - Summary

(13)

PCNN is Part-based CNN with 3 modules proposed for fine-grained visual categorization (FGVC) task. The goal is to detect distinctive object parts and learn discriminative features from global images as well as from discriminative object parts.

#### (1) Squeeze and Excitation block (SE) block.

- uses global information
- adaptively recalibrates channel-wise feature responses
- goal is to selectively emphasize informative channels
  - & suppress less useful ones:

#### (2) Part Localization Network (PLN).

- it is a deep fully convolutional network to locate distinctive object parts.
- trained under unsupervised learning by developing the excited feature maps due to SE blocks, with each learned part detector.
- Number of part ~~det~~ detector is upto us. In paper = 4.

#### (3) Part Classification Network (PCN)

- has 2 streams
  - Stream 1: classifying each individual object part into image-level categories.
  - Stream 2: concatenates part-level local features and global image feature from an image together into a joint feature with softmax function for final classification.

(14)

Ans: Multi-duplicate Focal Loss (MFL) is proposed

→ MFL : loss function is reshaped to focus on training on hard examples and thus down-weighting easy examples.

Things to note :-

- (1) a set of  $C \times S$  part detectors are learned.
- (2) both FCN and PN are trained jointly in an end-to-end fashion.
- (3) Metric based learning motivates the loss function with focusing harder examples from the easy examples.

$$\text{Loss} = \underbrace{\text{Metric learning loss}}_{\text{L}_m} + \underbrace{\text{Post classification loss}}_{\text{L}_{\text{post}}} + \underbrace{\text{Image classification loss}}_{\text{L}_{\text{cls}}}$$

reshaped

Limitations

→ the part bounding boxes in the paper is set at  $96 \times 96$  pixels with original input as  $448 \times 448$  pixels, parts which are smaller in size might not be detected in this case. we need flexibility with the choice of part bounding boxes.

→ Another limitation is the use of VGG Net as backbone. It might reduce the performance of the proposed PCNN in different domain. And with different backbone the initial settings (initialization) will have to be changed.

### Ques 3 (b) Solutions to PCNN limitation

(15)

Limitation → ① Rigid size choice for part bounding boxes.

Solution → use a different architecture for PCNN on the entire PCNN. Architectures like Transformer-based detectors can work well with variable length boxes. Or for localization networks, as no ~~loss~~ bounding boxes are available, we can add a regularization term to the loss function penalizing for ~~shorter~~ longer part boxes so as to force the architecture model to look at smaller part boxes.

$$\text{Eqn ①} \quad L_{\text{PCNN}} = \min_{n=1}^N \sum_{i=1}^{|X|} \underbrace{\text{clss}(m_n(x_i))}_{\text{detection function}} + \lambda_2 \underbrace{\text{div}(\text{mn}(x_i))}_{\text{diversity function}}$$

where,  
 $m_n(x) = \text{sigmoid}(O_n^T U)$

For each cluster we can add a regularized term  $\alpha$ .

$$L'_{\text{PCNN}} = L_{\text{PCNN}} + \alpha \|\text{distribution variance}\|$$

→ Regularizing the variance should help keep the part detector within the size limits but also extend. That is, without explicitly partitioning at 96x96 pixels we ~~can~~ can use variance to limit the size.

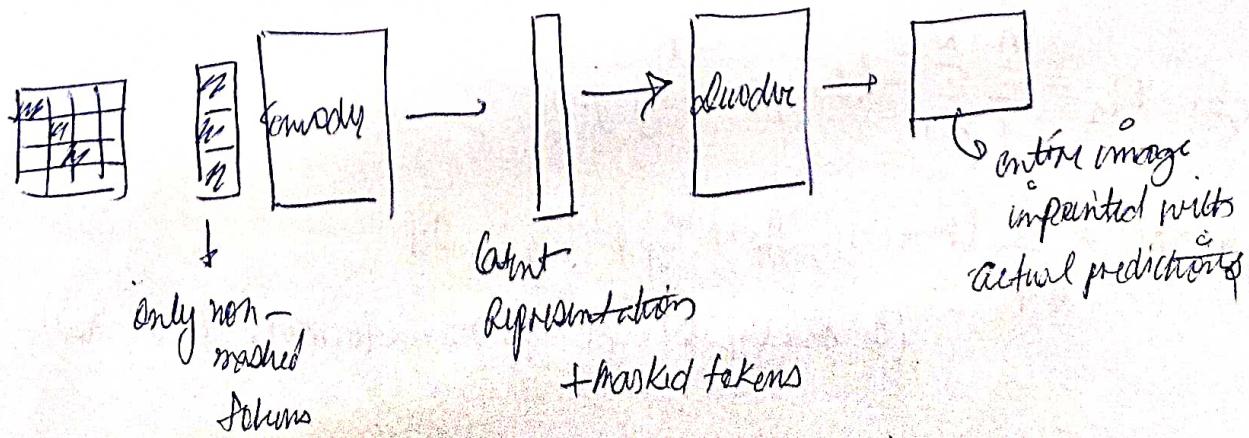
Ans 2) Another sol<sup>n</sup> could be to utilize the Image pyramids to break the part detection into smaller parts without having to explicitly implicitly depending on 1 constant size.

Ans 3, (a) Input image is 50% masked

Learn patches using a  $3 \times 3$  mask with all values 0

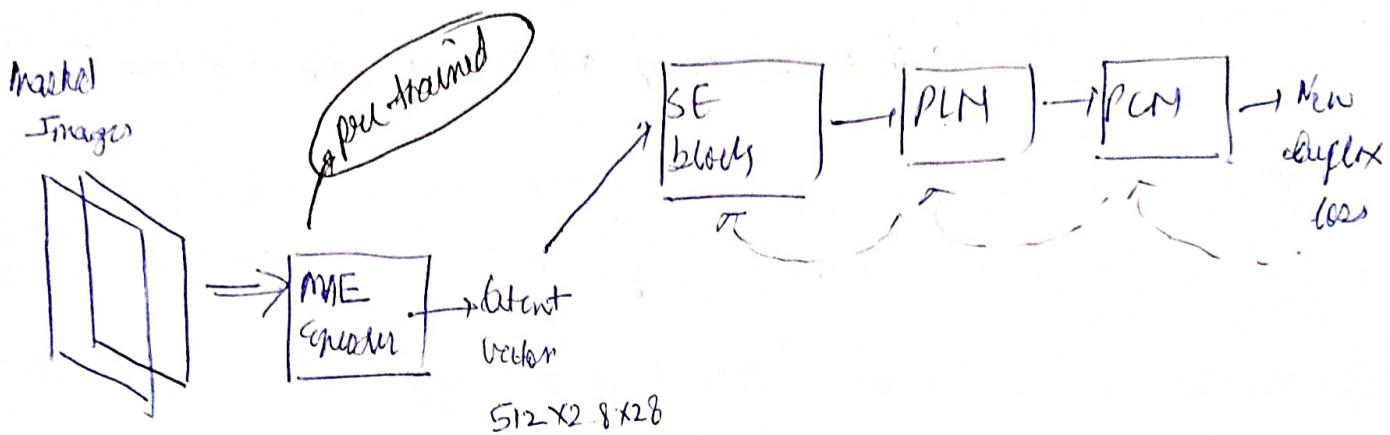
We can use Autoencoders before the PGN networks to handle the Input image that is masked. There exist Masked Autoencoders, that are able to reconstruct an image when it is upto 75% masked. These masked Autoencoders are pre-trained on heavy datasets in self-supervised fashion so we can either fine-tune the masked images or we can use the Autoencoder encoder instead of pre-trained VGG Net to get the features. Due to pre-training, the encoder learns specifically to get the important/essential features of a masked image.

Modified Step:  
Masked Autoencoders (MAEs)



(17)

- MME :- self-supervised pre-trained networks
- trained as vision transformers at backend.
  - can be used instead of VGG



Mathematically, everything remains same as the only change is in the backbone choice to get the feature maps or the input for SE blocks.

### Ques 3. (d) Duplex focal loss:

$$\text{Eq.7: } L_{\text{PCN}} = \min \left[ \underbrace{fL_1 + \lambda_2 fL_2}_{\text{and classification & metric learning}} + L_{\text{cls}} \right]$$

$$\text{Eq.8: } fL_1 = \sum_{i=1}^{1 \times 1} \sum_{n=1}^N - (1 - y_{i,n})^r \log y_{i,n}$$

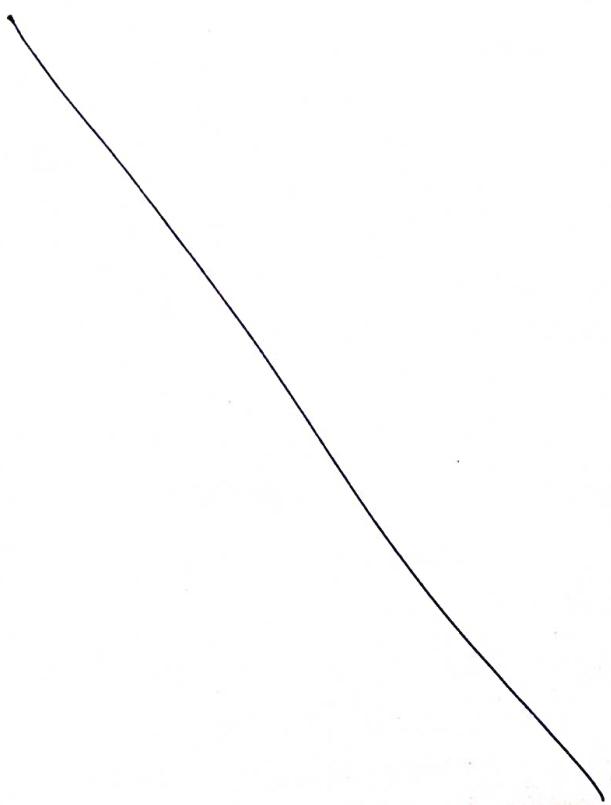
$$\text{Eq.9: } fL_2 = \sum_{q,p} (1 - \max(y_q, y_p))^r \log D^2(x_q, x_p) + (1 - \max(y_q, y_p))^r (1 - \log D^2(x_q, x_p)) \max(0, \log_2 D^2(x_q, x_p))$$

(18)

Replace Eq. 8. with SVM style margin based loss

$$FL_1 = \sum_{i=1}^{|X|} \sum_{n=1}^N -(1-y_{i,n}) \log y_{i,n}$$

Ans 3. b) A handcraft approach will be to get features from handcraft method like SIFT or HOG instead of VGGNet and then pass it through convolution filters through the 3 modules.



Ques 4. (a) gray-world assumption

(19)

grey world algorithm assumes that the average reflectance in the scene is achromatic. Hence, the illuminant color distribution is simply estimated by averaging each channel independently. It achieves good visual performance, but performs poorly for extremely deteriorated underwater scenes. It fails to remove the color shift, and looks bluish or suffers from noise and artifacts.

In the paper, author uses Gray-world algorithm to compute the white-balance image but after compensating for the loss of the red channel.

That is,

$$\left[ \begin{array}{c} \text{Red loss} \\ \text{compensation} \end{array} \right] + \left[ \begin{array}{c} \text{Grey world} \\ \text{Algorithm} \end{array} \right] \rightarrow \text{white balanced image.}$$

Observations  
are built based  
on this compensation

The original  
algorithm (with gray  
world assumption)  
is used on the compensated  
image.

Ques 4. (b) problem of saturation of color channel

Saturation of color channel (mostly red channel) can arise due to the ~~the~~ gray world step. The enhancement of red should primarily affect the pixels with small red channel values and should not change pixels that already include a significant

feed component. Basically, if compensation of red or any other channel is performed in those regions that are not highly attenuated it may lead to saturation of that particular channel. (20)

disadv (a) In proposed method, the red channel compensation depends on global characteristics of the image.

$$\text{eq 4. } I_{Rn} = I_{R(n)} + \alpha (\bar{I}_g - \bar{I}_R) \cdot (1 - I_{R(n)}) \cdot I_g \ln$$

\*  $\bar{I}_R$  and  $\bar{I}_g$  in the equation are mean value of  $I_R$  and  $I_g$ , that is, the red and green channels of the image. Entire image is averaged in order to compensate for the red channel.

+ If mean value changes in a local neighborhood, we may be able to better compensate the attenuation.

There is one issue with making red channel compensation globally dependent, the entire image might not be attenuated consistently. Lighting and depth changes in the image, hence, some areas might get over or under compensated. It is thus better to make this compensation locally-dependent rather than globally-dependent. The neighborhood can be defined as a simple order statistic of  $3 \times 3$  or  $5 \times 5$  patch or we can use segmentation mask, though it will add to the computation overload for the image and esp. algorithm

## Ques. (d) Red channel compensation

(21)

↳ locally-dependent.

$$I_{rc}(n) = I_r(n) + \alpha (\bar{I}_g - \bar{I}_{r,n}) (1 - I_{r,n}) I_g(n)$$

↳

Creating a neighborhood of  $K \times K$  patches with a total of  $m$  patches

$$I_{rc, \text{comp}}(n) = I_r(n) + \alpha \left[ \sum_{i,n=1}^{K,m} (\bar{I}_{g,i,n} - \bar{I}_{g,n}) \right] (1 - I_{r,n}) I_g(n)$$

↓

We take the mean over  $K \times K$  neighborhood.

$\bar{I}_{g,n}$ ,  $\bar{I}_{g,i,n} \rightarrow$  "  $n^{\text{th}}$  patch with i value of neighborhood"  
to  $K$

## Pseudo-code

I/P  $\rightarrow$  Image

→ divide the image into  $m$  patches of  $K \times K$  size

→ Compute  $m$  mean values for red & green channel

•  $\bar{I}_{r,n}, \bar{I}_{g,n} \quad \forall n = 1, \dots, m$

→ Compute compensated  $\bar{I}_{rc}(n)$  with this ~~local~~ list of means computed on local features of images.

→ Inverse gray world & get white balancing image.

Ques 4. F) For unequal channel attenuation at different regions of the image. We should apply locally dependent channel compensation strategy. Locally dependent strategy will adjust the local means of channel that are attenuated differently and will therefore compensate differently at the specified region.

(22)

The algorithm described in previous answer can be used here or we can take regions as an input with the underwater image, defining specifically marked regions and using a local mean on these areas.

#### Ques 4. F) Gamma Correction

Gamma correction aims at correcting the global contrast. The correction increases the difference between darker and lighter regions at the cost of a loss of details in the under or over exposed regions. The gamma correction which is helpful when we have too bright or too dull images.

With the proposed method, it helps because in general, underwater images appear to be too bright. Using gamma correction, we can make the image less bright making it more evenly spread.

The issue as stated above with this is the loss of details in under or over-exposed regions. To compensate for this, authors use a sharpened version of the white balanced image. They follow the anchors

masking principle, by blending a blurred version of the image with the image to sharpen using → 23

$$S = I + \beta (I - G * I)$$

↑                      ↓                      ↓  
sharpened      actual      blurred  
image          image        version

Ques 4. (g) Yes, as we move upward in a Gaussian pyramid, we get more and more high frequency components as we keep downsampling the image and applying it with the sharpened filters or the laplacian filters.

Laplacian pyramid is a bandpass pyramid as it allows lowpass frequency. We can upscale an image using the pyramid. It is quite opposite →

