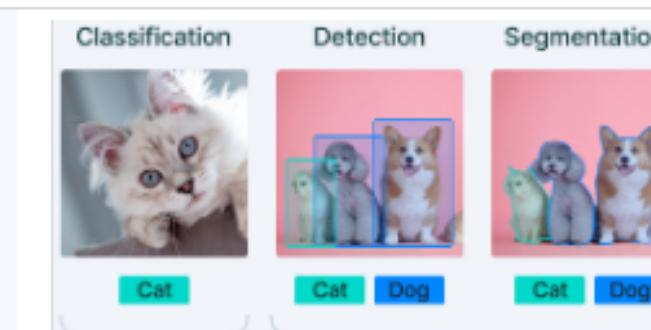


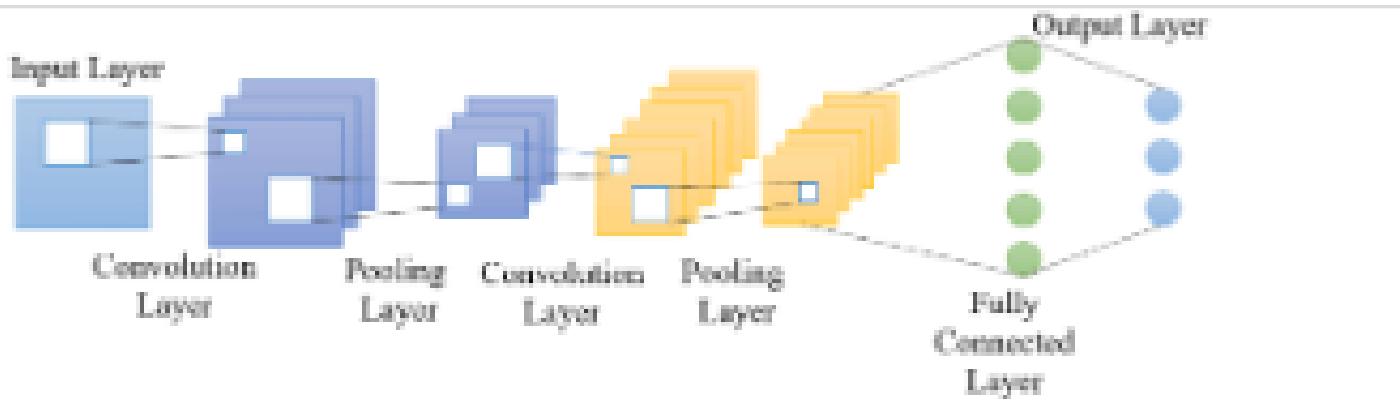
COMPUTER VISION CHEATSHEET

L1 : INTRODUCTION TO COMPUTER VISION (CNN) CHEAT SHEET

- Computer vision is a field of artificial intelligence (AI) that enables computers and systems to derive meaningful information from digital images, videos and other visual inputs.



A. WHAT IS A CONVOLUTIONAL NEURAL NETWORK ?



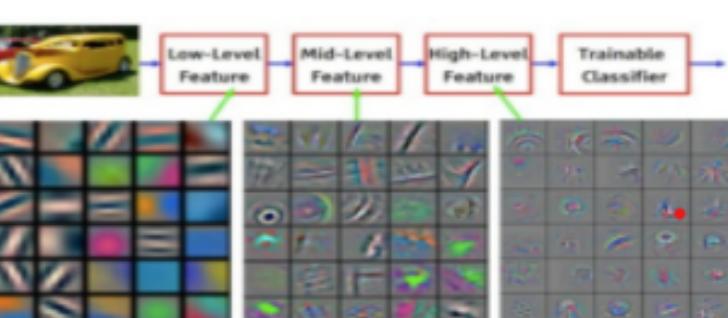
- A Convolutional Neural Network (CNN) is a deep learning model used for image and signal processing, where convolutional and pooling operations are applied to input data to extract features and make predictions.

B. WHY USE CNN OVER MLP?

- Feature extraction:** CNNs are particularly well-suited for computer vision tasks because they can automatically learn and extract meaningful features from images and videos, while traditional machine learning algorithms like MLPs require manual feature engineering.
- Handling large datasets:** CNNs can handle large datasets and can be trained on a large set of image data, while MLPs are not that efficient in handling large datasets.
- Handling image variations:** CNNs can handle variations such as different lighting conditions, angles, and rotations, while MLPs can struggle with these variations.

D. CNN FEATURE LEVELS

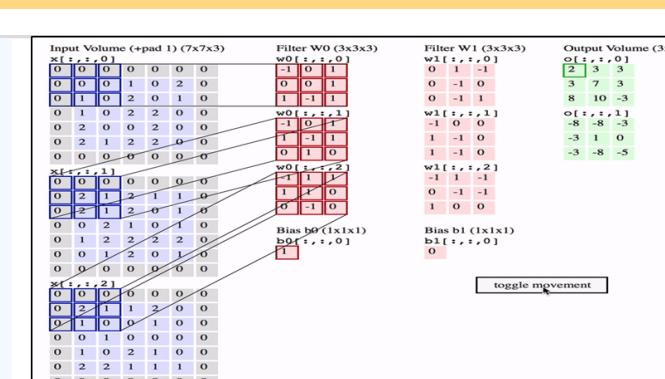
Convolutional Neural Network



- In a Convolutional Neural Network (CNN), features can be classified into different levels based on their complexity and abstraction, such as low-level features (e.g. edges, textures), mid-level features (e.g. parts of objects), and high-level features (e.g. whole objects, scenes).

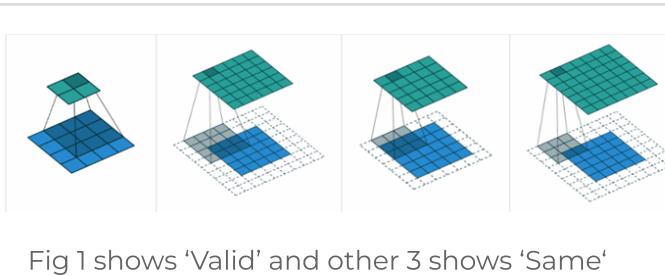
E. WHAT IS A CONVOLUTIONAL AND CONVOLUTION OPERATIONS ?

- A kernel in a CNN is a small matrix used for performing a convolution operation which involves element-wise multiplication followed by addition and activation (such as ReLU)

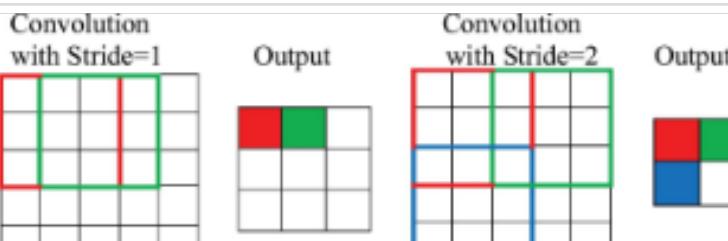


D. WHAT IS PADDING ?

- Valid and same padding in CNNs are methods of handling the spatial dimensions of input and output matrices during the convolution operation, with "valid" reducing the output size and "same" preserving the input size in the output.



E. WHAT IS STRIDE ?



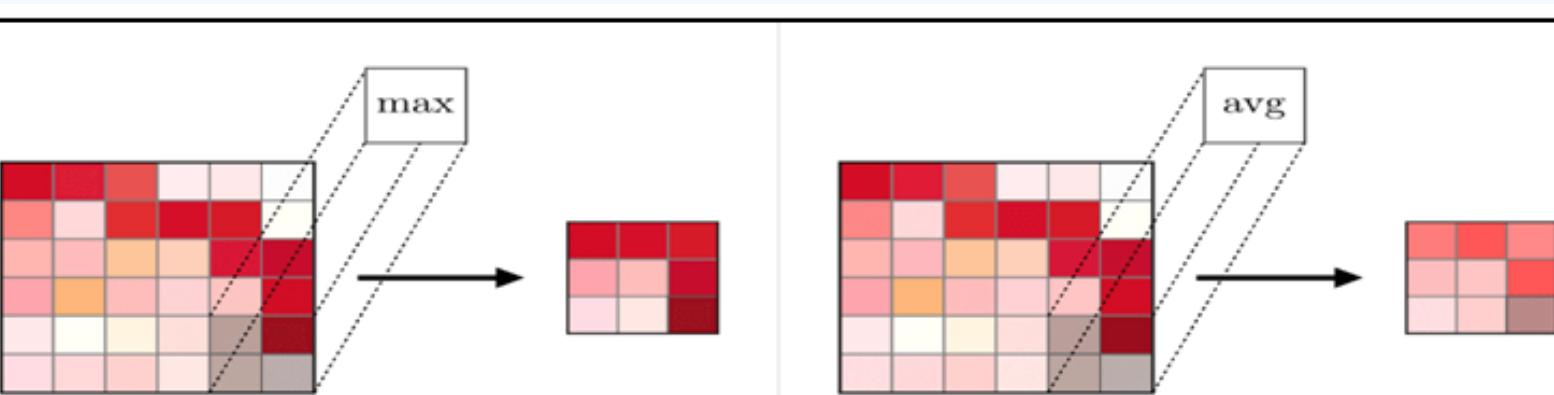
- In Convolutional Neural Networks (CNNs), strides are the number of pixels the convolution filter moves when scanning over the input.

H. SHAPE OF OUTPUT LAYER

$$\text{output} = \frac{\text{input} - \text{kernel_size} + 2 * \text{padding}}{\text{stride}} + 1$$

H. WHAT IS POOLING ?

Pooling layers provide an approach to down-sampling feature maps by summarizing the presence of features in patches of the feature map.



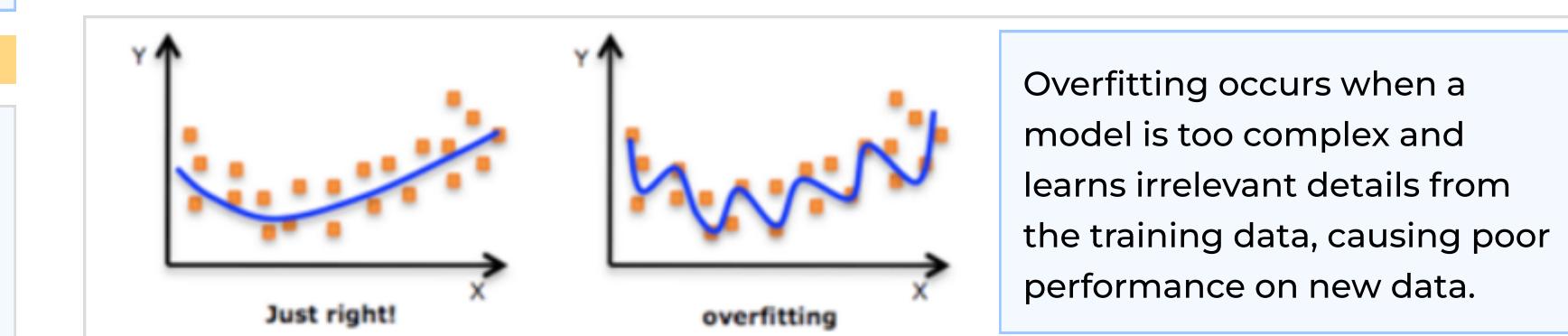
- Max-Pooling** - Each pooling operation selects the maximum value of the current view.
- Average-pooling** - Each pooling operation averages the values of the current view

D. WHAT IS PADDING ?

- For Convolution layer with k filters of size f x f and the input channel of image c: $\text{Total Parameter} = k \times (f \times f \times c_{in}) + k$
- For a Fully Connected layer with n neurons and number of input l: $\text{Total Parameter} = n \times (l+1)$
- By calculating the number of parameters in each layer and summing them up, you can obtain the total number of parameters in the entire CNN.

L2: REVISITING CNN DEAL WITH OVERTFITTING CHEAT SHEET

A. WHAT IS OVERTFITTING ?

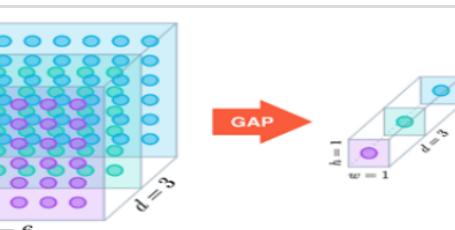


B. METHODS TO AVOID OVERTFITTING

- Regularization: L1/L2, Dropout reduce model complexity, prevent overfitting.
- Early Stopping: Stop training when validation performance declines.
- Data Augmentation: Apply random image transformations to increase dataset.
- Transfer Learning: Use pre-trained models and fine-tune.
- Model Ensemble: Combine multiple models for robust results.
- Reduce Network Complexity: Use smaller filters, fewer layers.

C. GLOBAL AVERAGE POOLING

- Global Average Pooling 2D is a type of pooling operation used in Convolutional Neural Networks (CNNs).



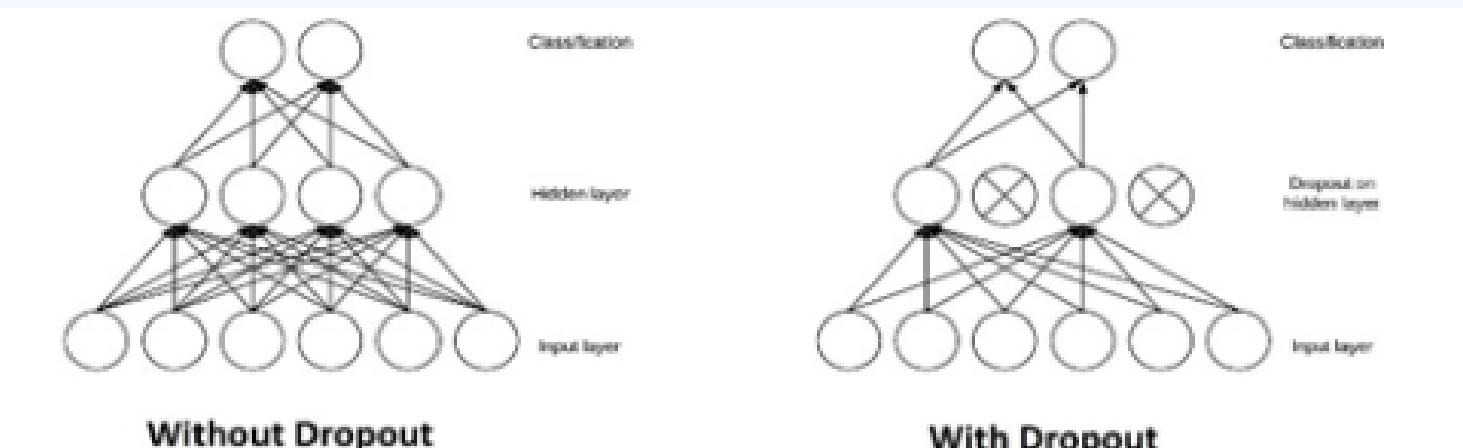
D. BATCH NORMALIZATION

- Batch normalization normalizes neuron activations within batches, improving training stability and speed. It transforms each batch's activations to zero mean and unit standard deviation before passing to the next layer.

$$\hat{x}^{(k)} \leftarrow \frac{x_i^{(k)} - \mu^{(k)}}{\sqrt{\sigma^2(k) + \epsilon}}$$

D. WHAT IS PADDING ?

- Dropout is a deep learning regularization technique that randomly drops out some neurons during training to prevent overfitting, encouraging the network to learn multiple representations of data.



A. WHEN TO USE BATCH NORMALIZATION AND DROPOUT

- Empirically, dropout is most effective placed between Dense layers with 0.5 probability. Lower probability (0.1-0.25) dropout after MaxPooling improves performance.
- Normalizing feature distribution after Conv Layer and passing it to ReLU block clamps negative values to 0, defeating normalization's purpose.

A. WHEN TO USE BATCH NORMALIZATION AND DROPOUT

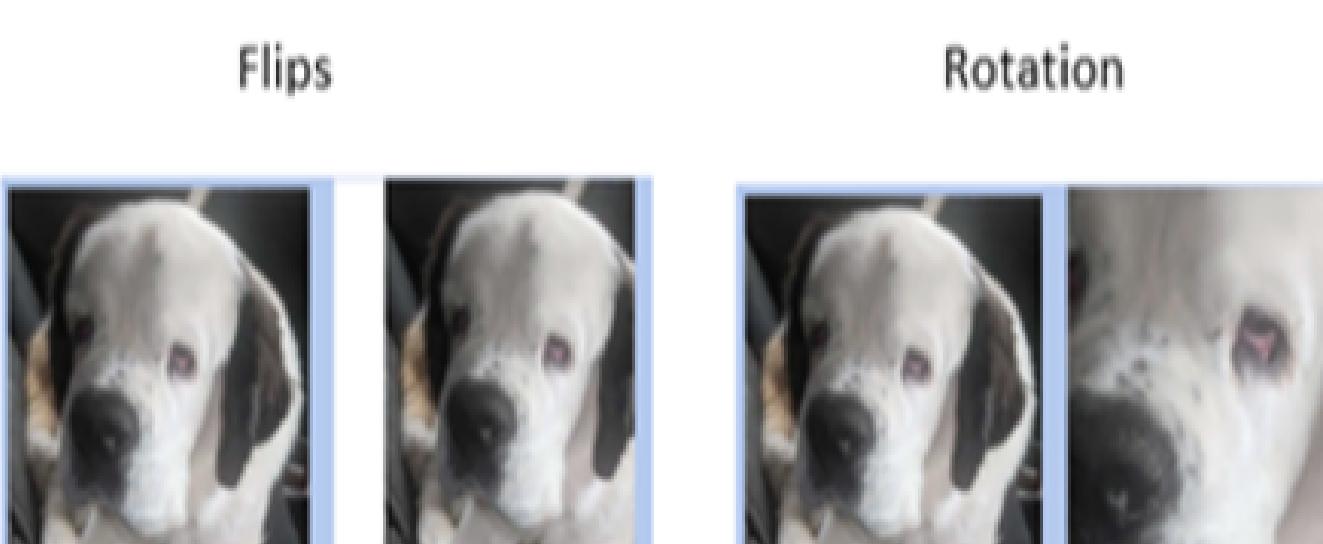
Lasso Regularization(L1)

$$\text{loss} = \sum_{i=0}^n (y_i - X_i \beta)^2 + \sum_{j=0}^m |\beta_j|$$

Ridge Regularization(L2)

$$\text{loss} = \sum_{i=0}^n (y_i - X_i \beta)^2 + \sum_{j=0}^m \beta_j^2$$

D. EXAMPLE OF DATA AUGMENTATION



Translation

Scaling

A. WHEN TO USE BATCH NORMALIZATION AND DROPOUT

- Data augmentation artificially increases training data in CNNs by applying various transformations to input images, preventing overfitting.
- Data augmentation is particularly useful when the size of the training dataset is limited.
- Techniques include rotation, scaling, translation, and flipping.

A. WHEN TO USE BATCH NORMALIZATION AND DROPOUT

- There are many ways to apply augmentation in Tensorflow/Keras, few of them are discussed here:

-- To apply augmentation in TensorFlow/Keras, use Keras Preprocessing Layers such as RandomFlip, RandomRotation, or tf.image methods like stateless_random_flip_up_down and stateless_random_brightness.

-- Use Keras ImageDataGenerator API for a suite of out-of-the-box augmentation

L2: REVISITING CNN DEAL WITH OVERTFITTING CHEAT SHEET

A. FORWARD PROPAGATION IN CNN

- The forward function iterates through every filter in every image channel and it pass over a kernel of specified size over the image and perform the convolution operation between the kernels and the input image.

$$\begin{array}{c} \text{Output Matrix} \\ \begin{matrix} y_{11} & y_{12} & y_{13} \\ y_{21} & y_{22} & y_{23} \\ y_{31} & y_{32} & y_{33} \end{matrix} \end{array} = \begin{array}{c} \text{Image Matrix} \\ \begin{matrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{matrix} \end{array} * \begin{array}{c} \text{Kernel} \\ \begin{matrix} k_{11} & k_{12} \\ k_{21} & k_{22} \end{matrix} \end{array} + \begin{array}{c} \text{Bias} \\ \begin{matrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{matrix} \end{array}$$

$$Y_i = X_i * K_{ij} + B_i, i = 1 \dots \text{depth}, j = 1 \dots n_{\text{filters}}$$

Forward propagation mathematically for the Convolutional Layer is as follows:

B. BACKWARD PROPAGATION IN CNN

Backward Propagation refers to updating the weights and biases of our model with respect to the loss of our model.

The derivative of error with respect to :

- Kernel
- Bias
- Input is calculated during backpropagation

$$\begin{array}{c} I \\ \begin{matrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{matrix} \end{array} \xrightarrow{\text{Conv}} \begin{array}{c} K \\ \begin{matrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{matrix} \end{array} \xrightarrow{\text{ReLU}} \begin{array}{c} Y \\ \begin{matrix} y_{11} & y_{12} \\ y_{21} & y_{22} \end{matrix} \end{array}$$

C. LOSS GRADIENTS

$$\frac{\partial L}{\partial K} = \text{Conv}\left(X, \frac{\partial L}{\partial Y_1}\right)$$

$$\frac{\partial L}{\partial X} = \text{Conv}\left(\text{Padded}\left(\frac{\partial L}{\partial Y_1}\right), 180^\circ \text{Rotated Kernel } K_1\right)$$

$$\frac{\partial L}{\partial B_1} = \sum \left(\frac{\partial L}{\partial Y_1} \right)$$

$$\begin{array}{c} \text{Forward Propagation} \\ \text{ReLU}(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases} \end{array} \quad \begin{array}{c} \text{Backward Propagation} \\ \text{ReLU}'(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases} \end{array}$$

C. LOSS GRADIENTS

Forward Propagation

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

$$\frac{\partial (\text{Softmax})}{\partial x_1} = \frac{e^{x_1} \cdot (e^{x_2} + e^{x_3})}{(e^{x_1} + e^{x_2} + e^{x_3})^2}$$

C. LOSS GRADIENTS

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

$$\begin{array}{c} \text{Forward Propagation} \\ f(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})} \end{array} \quad \begin{array}{c} \text{Backward Propagation} \\ f'(x) = \frac{1}{N} \cdot \left(\frac{1 - y_{true}}{1 - y_{pred}} - \frac{y_{true}}{y_{pred}} \right) \end{array}$$

C. LOSS GRADIENTS

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\tanh'(x) = \frac{1}{N} \cdot \left(\frac{1 - y_{true}}{1 - y_{pred}} - \frac{y_{true}}{y_{pred}} \right)$$

COMPUTER VISION CHEATSHEET

L9 : OBJECT SEGMENTATION CHEAT SHEET

A. WHAT IS OBJECT SEGMENTATION ?

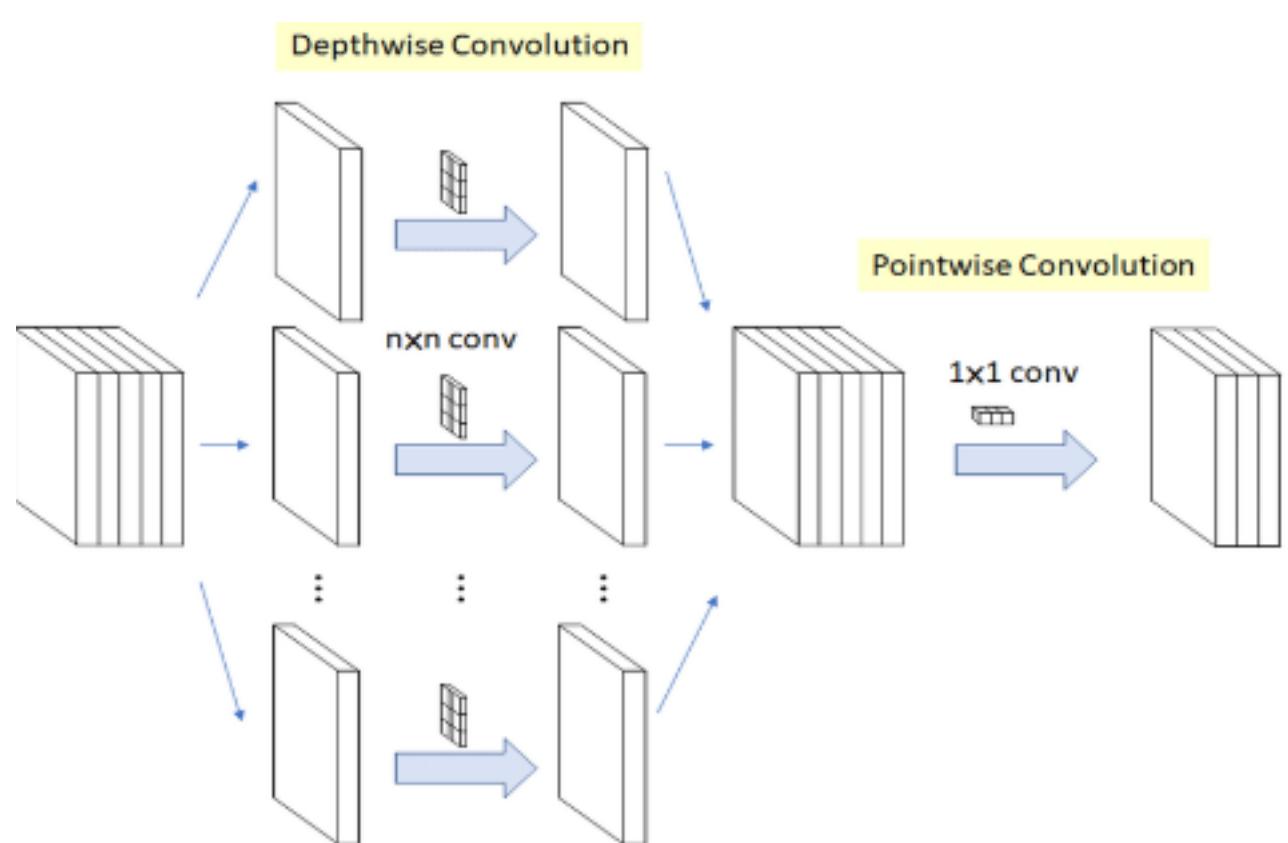
- MobileNet is a compact and efficient CNN architecture ideal for mobile and embedded devices.
- It uses Depthwise-Separable convolution and pointwise convolutions, yielding good results while being computationally inexpensive.
- Despite having 10x fewer parameters than ResNet50, it still offers great results.



A. WHAT IS OBJECT SEGMENTATION ?

- Depthwise Separable Convolution splits the computation into two steps:

 1. Depthwise convolution : It applies a single convolutional filter per each input channel.
 2. Pointwise convolution: It is used to create a linear combination of the output of the depthwise convolution.



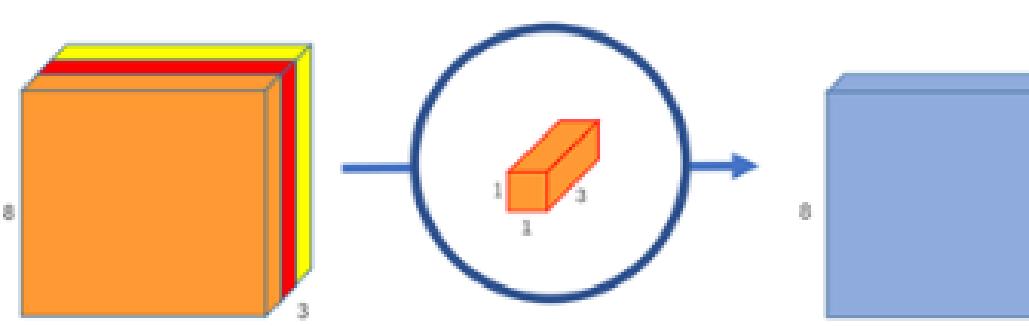
A. WHAT IS OBJECT SEGMENTATION ?

- In Depthwise Convolution, we apply a single convolutional filter for each input channel



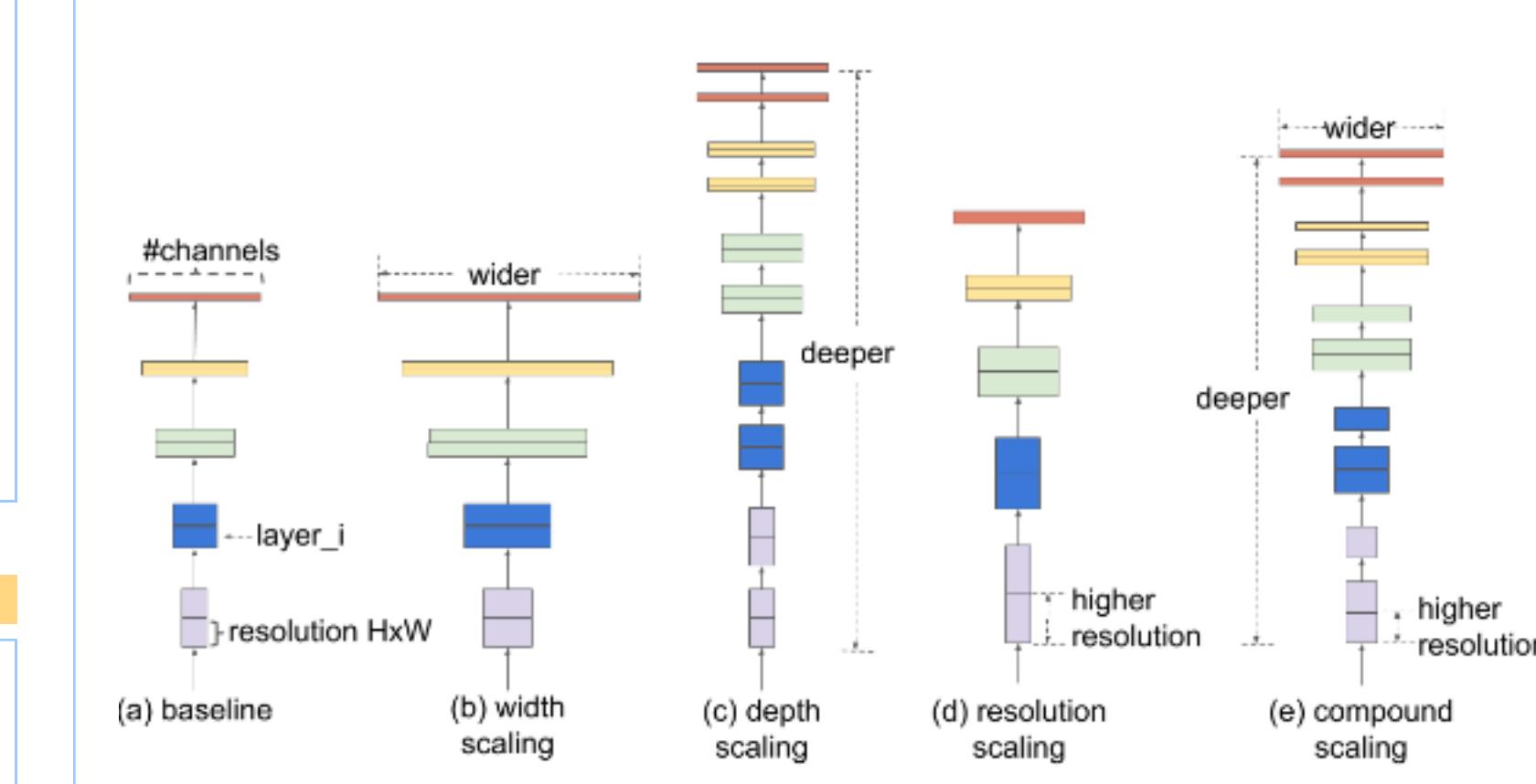
A. WHAT IS OBJECT SEGMENTATION ?

- In Pointwise Convolution, we use a 1x1 kernel, which iterates across each and every point. This kernel has a depth equal to the number of channels in the input picture.

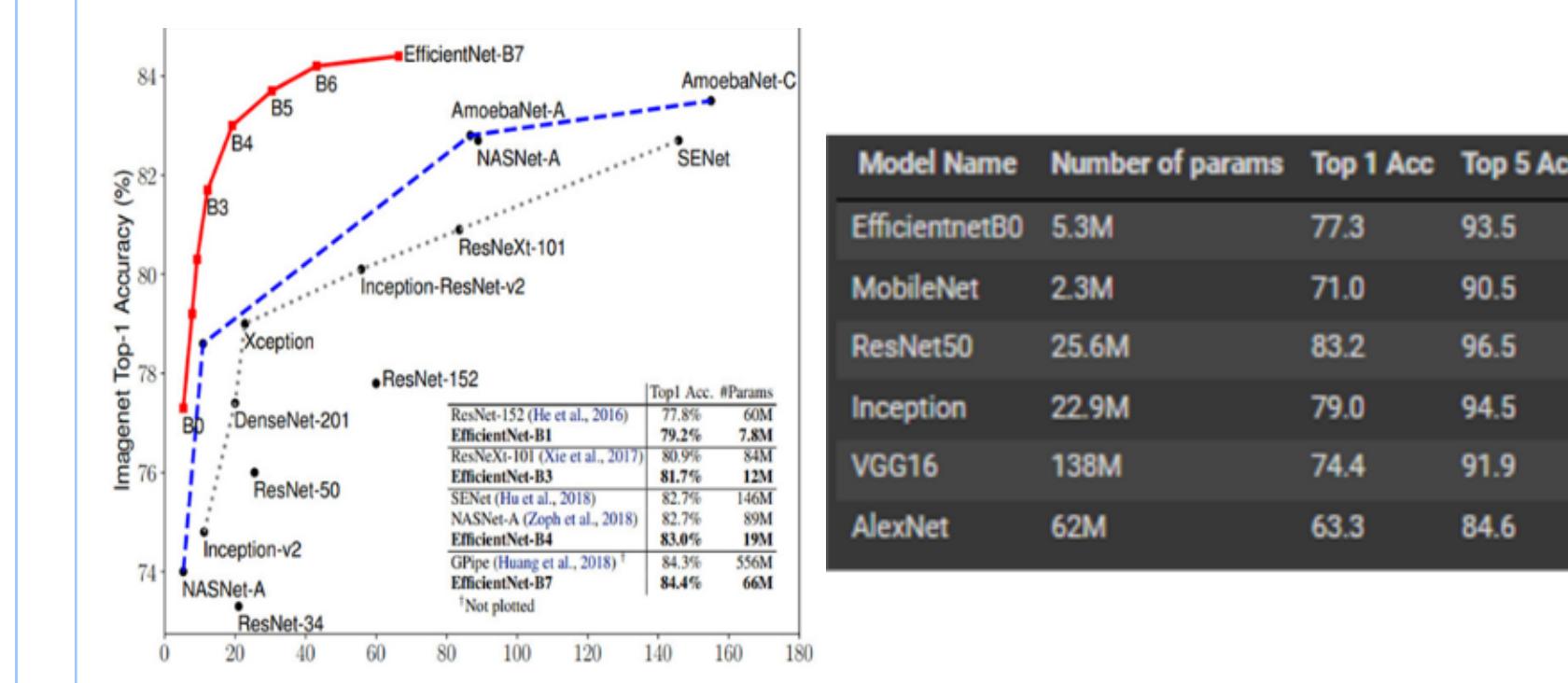


A. WHAT IS OBJECT SEGMENTATION ?

- EfficientNet is a CNN architecture that beats all previous ones with fewer parameters.
- It's a family of neural networks with scaling emphasis using the compound coefficient technique.
- It offers eight models of increasing size, from B0 to B7.



A. WHAT IS OBJECT SEGMENTATION ?



A. WHAT IS OBJECT SEGMENTATION ?

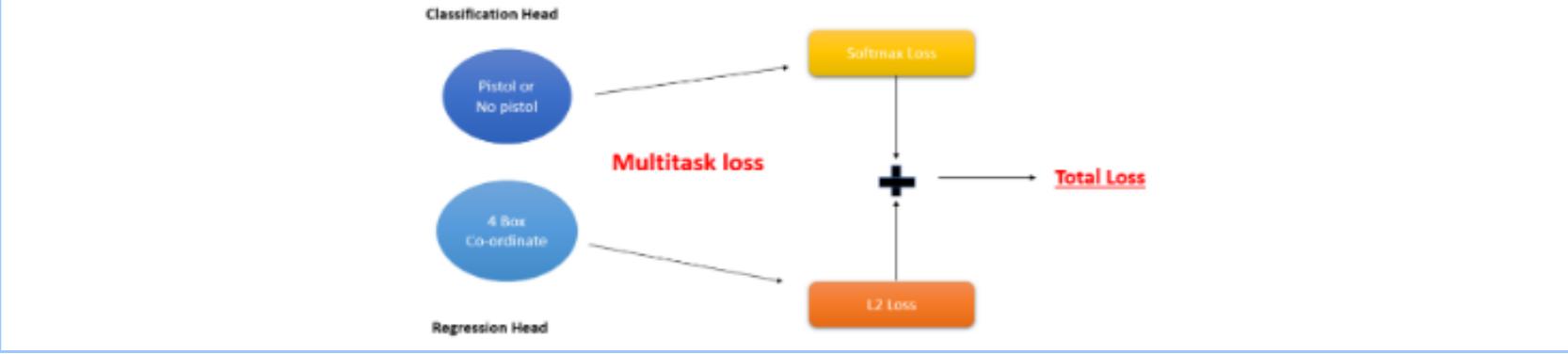
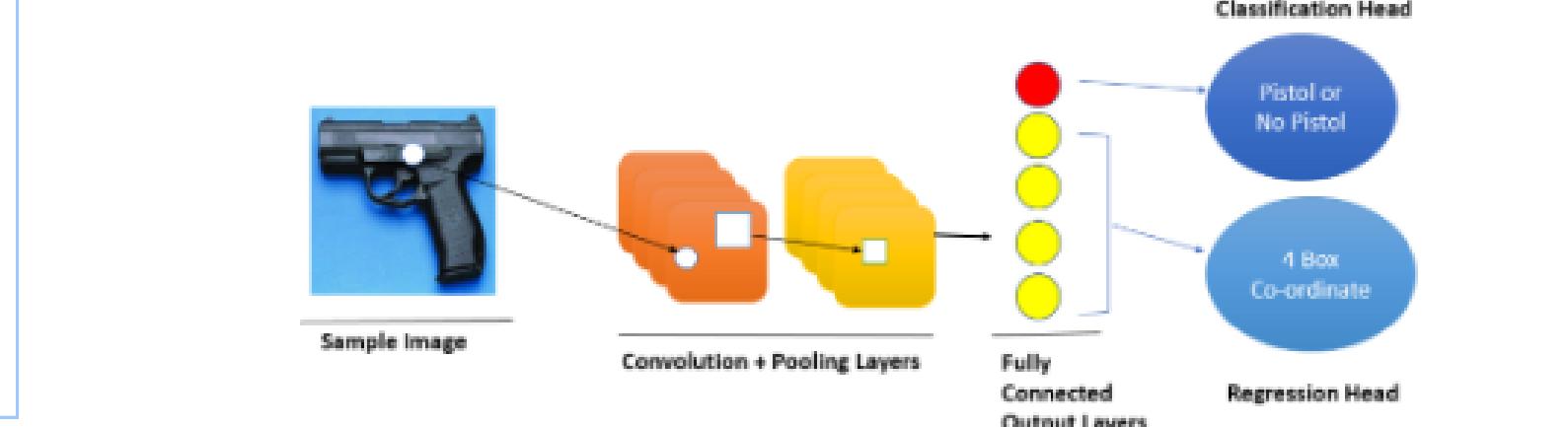
- Size of the new dataset
- Similarity of new dataset with the original dataset
- New Dataset is small and similar : Fine-tuning might lead to overfitting.
- New Dataset is large and similar : Fine-tune the pre-trained network
- New Dataset is large and different : Train a Conv. Network from scratch

L7 : OBJECT DETECTION WITH TWO STAGE METHODS CHEAT SHEET

A. WHAT IS OBJECT SEGMENTATION ?

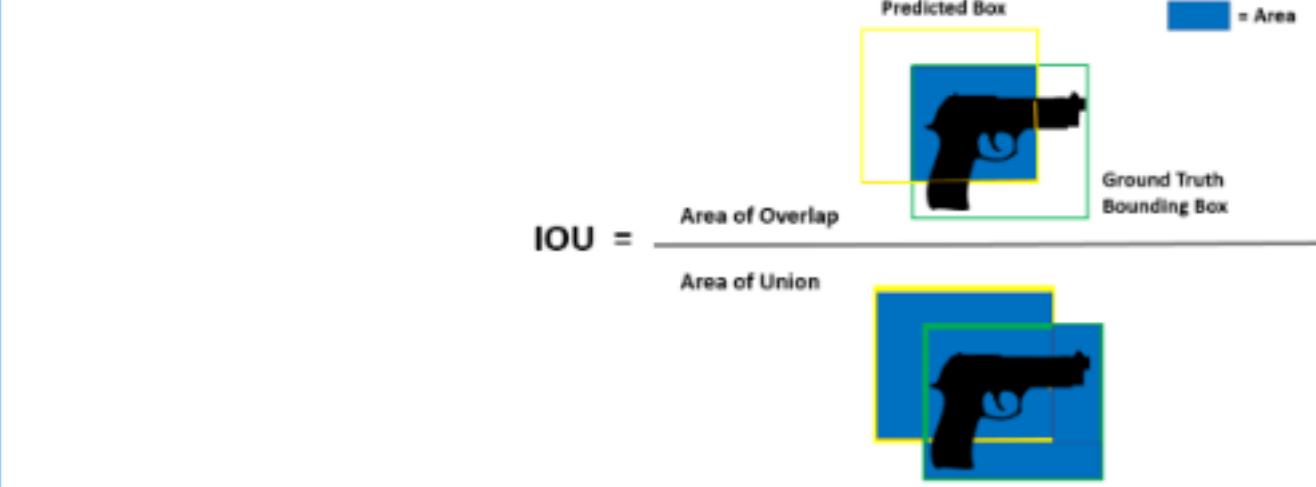
- Object detection is a computer vision task to locate and identify objects in an image or video.
- It answers the following question:
- Where is the object of interest present in the Image?
- How many objects are present in Image?

A. WHAT IS OBJECT SEGMENTATION ?



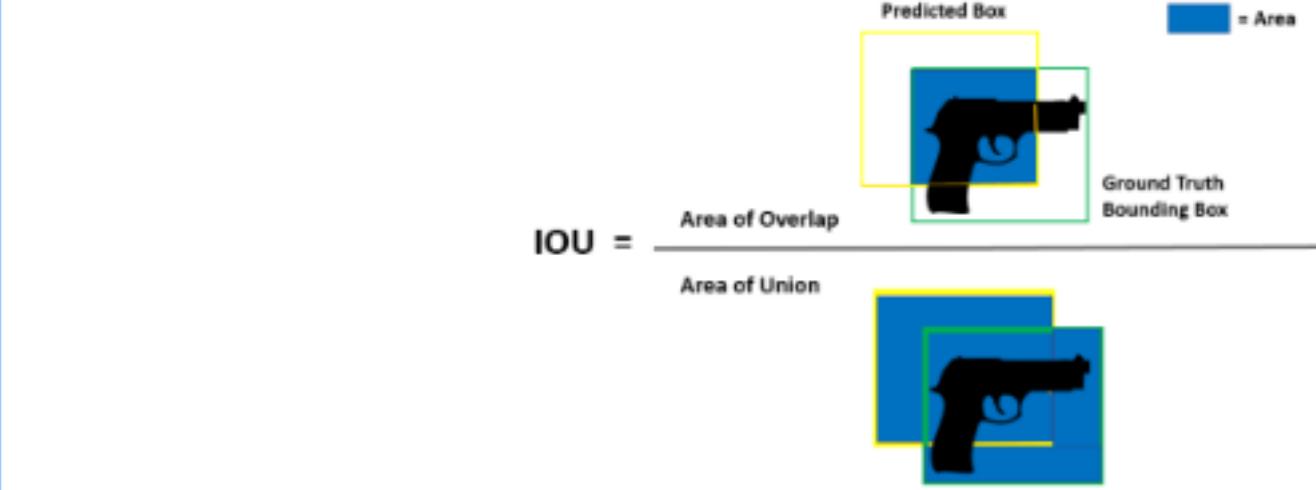
A. WHAT IS OBJECT SEGMENTATION ?

- IoU stands for Intersection over Union, which is a commonly used evaluation metric for object detection tasks.



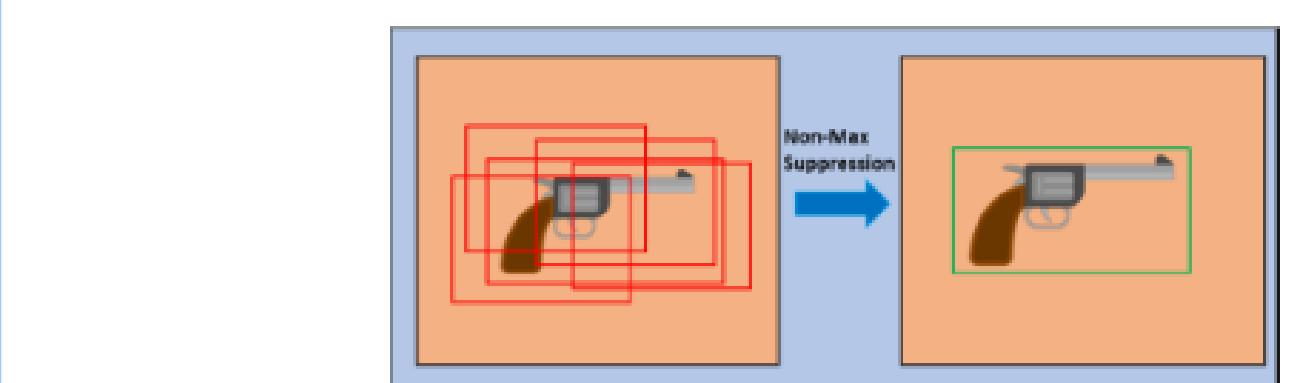
A. WHAT IS OBJECT SEGMENTATION ?

- IoU stands for Intersection over Union, which is a commonly used evaluation metric for object detection tasks.



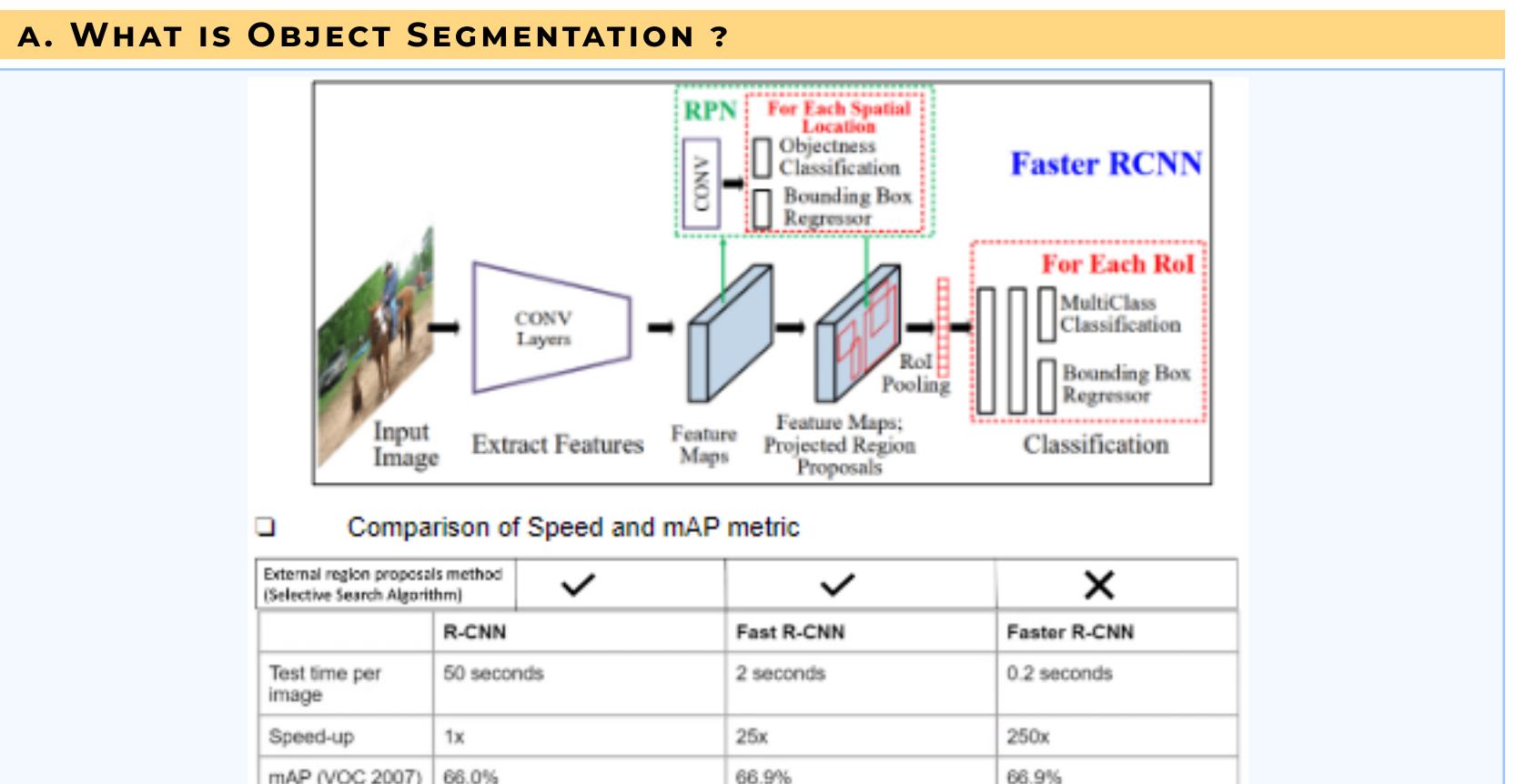
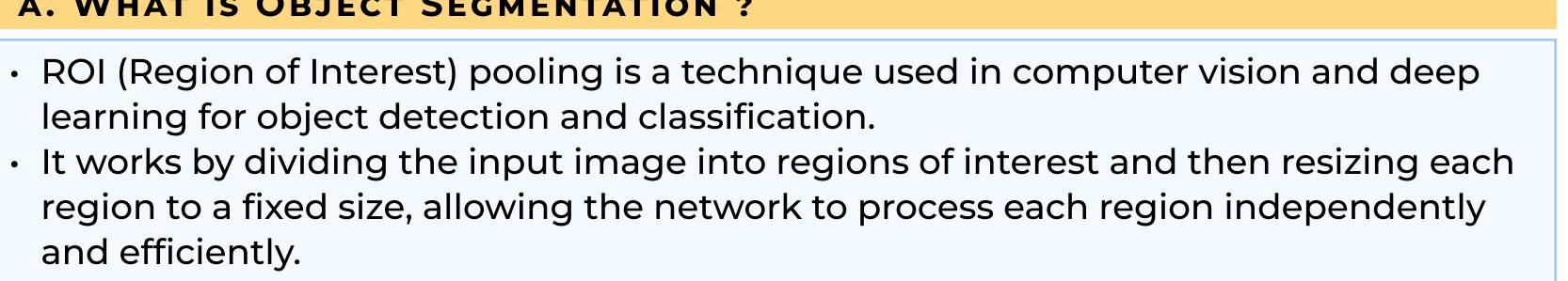
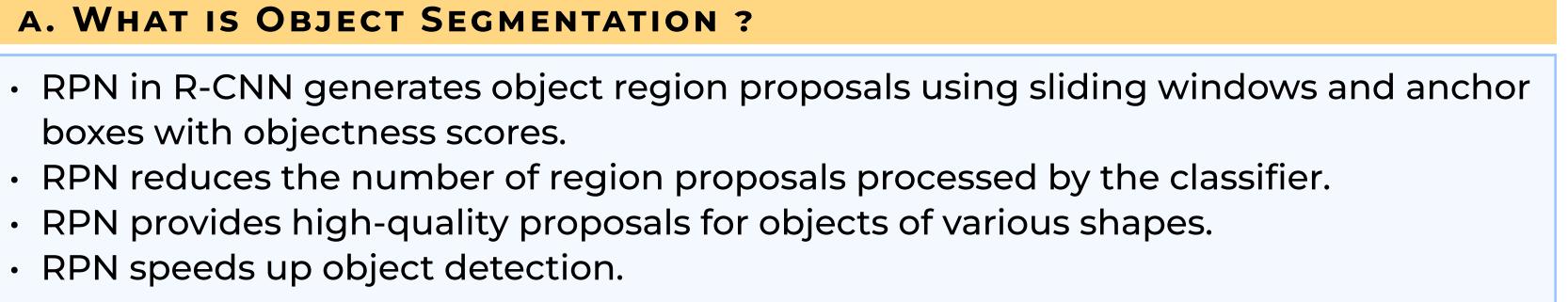
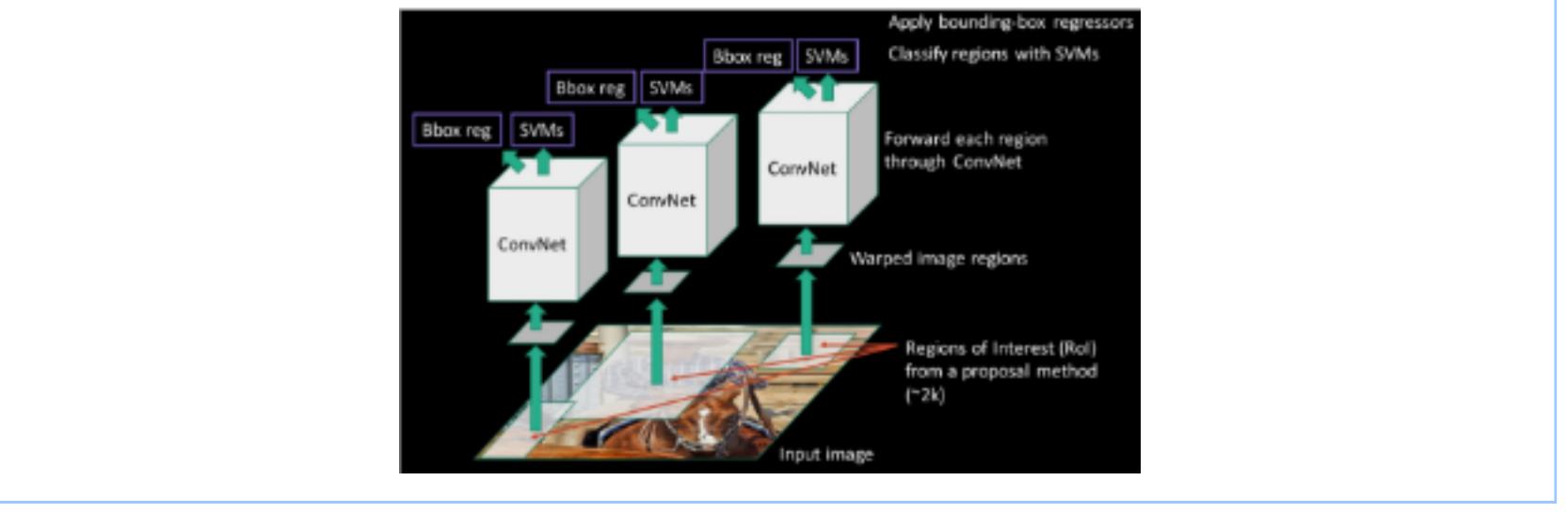
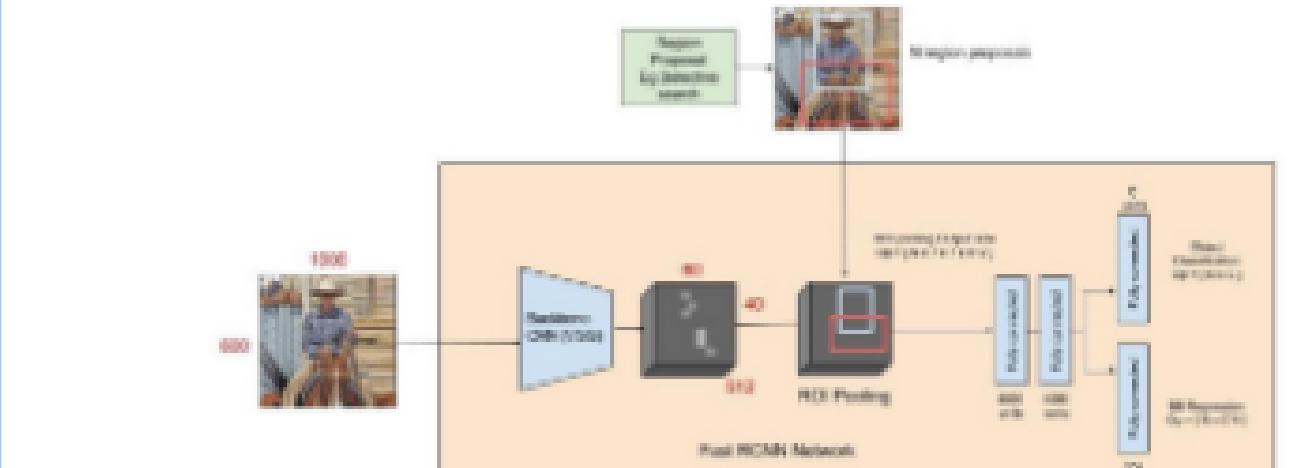
A. WHAT IS OBJECT SEGMENTATION ?

- NMS (Non-Maxima Suppression) is a technique used to suppress redundant bounding boxes (BBOXES) in object detection algorithms in computer vision.
- It works by keeping only the bounding box with the highest confidence score and removing other overlapping boxes, reducing the number of false positive detections.



A. WHAT IS OBJECT SEGMENTATION ?

- Instead of extracting features from each region proposal independently, Fast RCNN uses a shared convolutional layer to extract features from the entire image, reducing the number of computations needed.



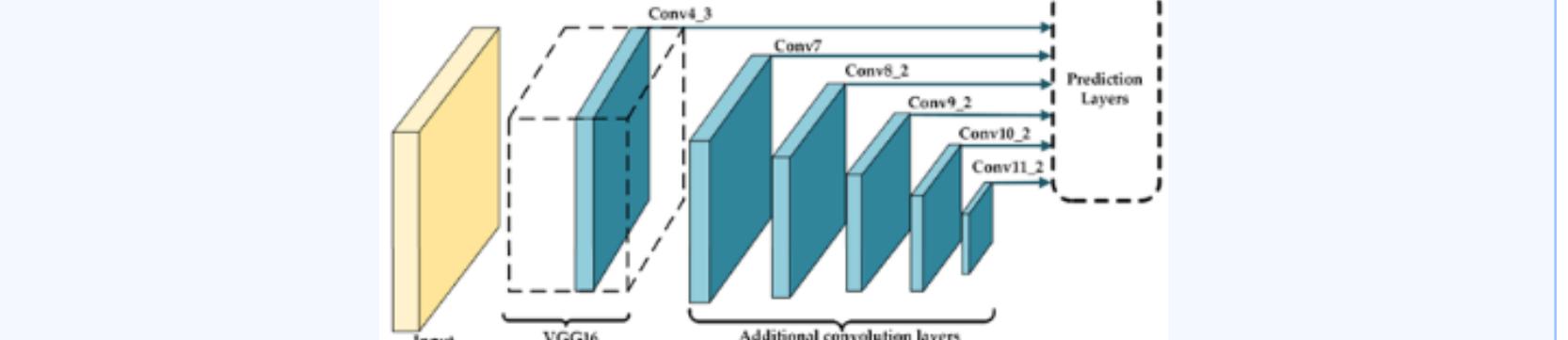
L9 : OBJECT SEGMENTATION CHEAT SHEET

A. WHAT IS OBJECT SEGMENTATION ?

- Computational cost: Slow inference due to two-stage architecture.
- Poor recall: Can miss small or occluded objects.
- Inefficient proposal generation: Selective search is slow and doesn't scale well.
- Limited versatility: Designed for limited object categories.
- Difficult to fine-tune: Challenging for specific tasks with limited training data.
- Limited scalability: Scaling to larger object categories is challenging, leading to slower inference and more memory usage.

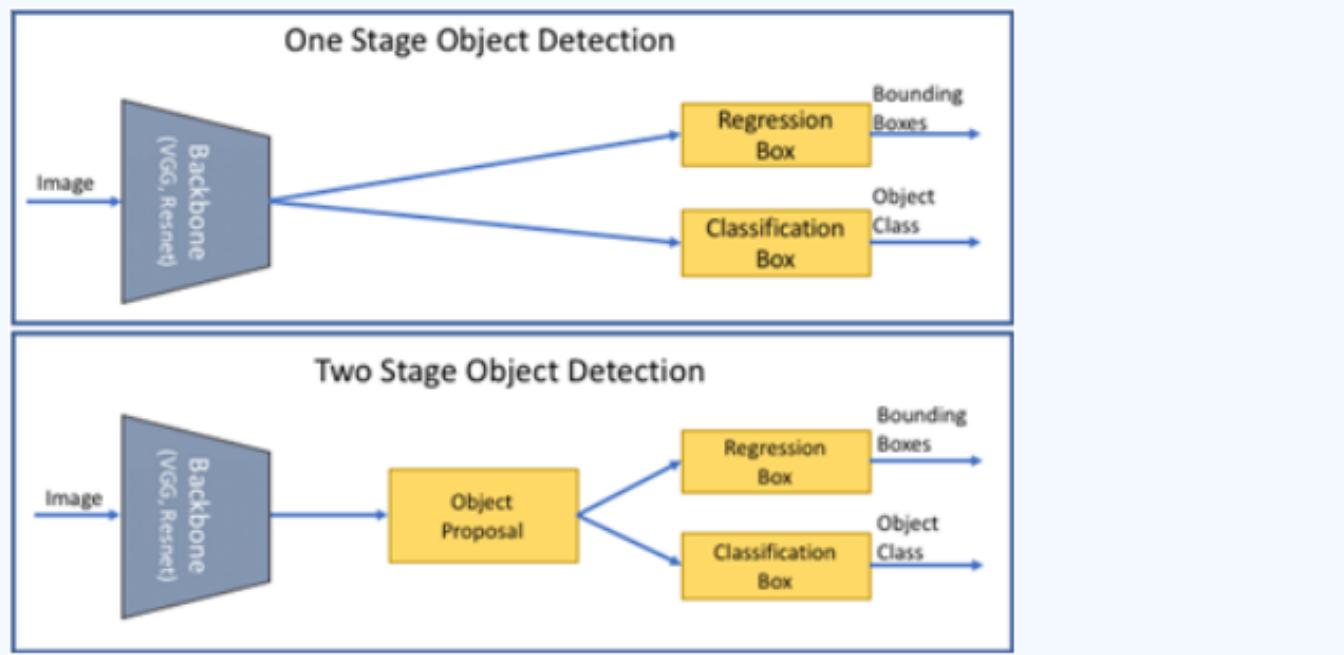
A. WHAT IS OBJECT SEGMENTATION ?

- Object detection method that predicts class labels and bounding boxes in a single stage.
- Offers an efficient alternative to two-stage detectors like R-CNN.
- Designed for real-time processing of large numbers of object categories.



COMPUTER VISION CHEATSHEET

A. WHAT IS OBJECT SEGMENTATION ?

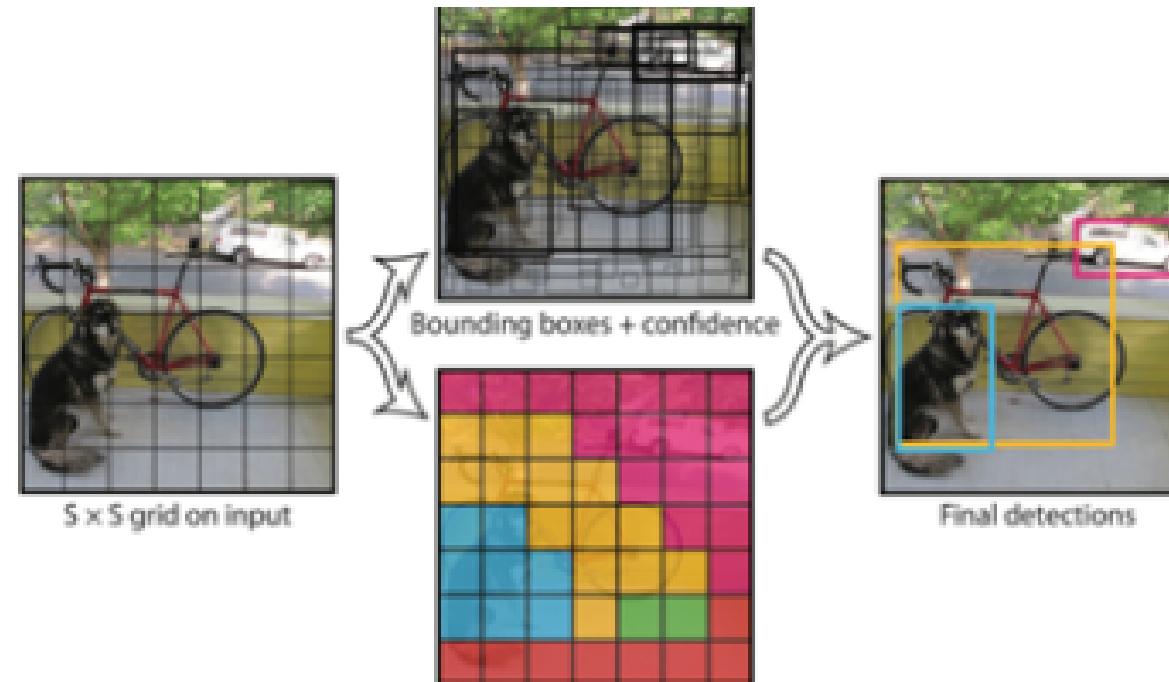


A. WHAT IS OBJECT SEGMENTATION ?

- You Only Look Once (YOLO) is a one-stage object detection method that predicts class labels and bounding boxes for objects within an image.
- It's designed to handle a large number of object categories in real-time and is known for its fast inference time and efficient use of computational resources.

A. WHAT IS OBJECT SEGMENTATION ?

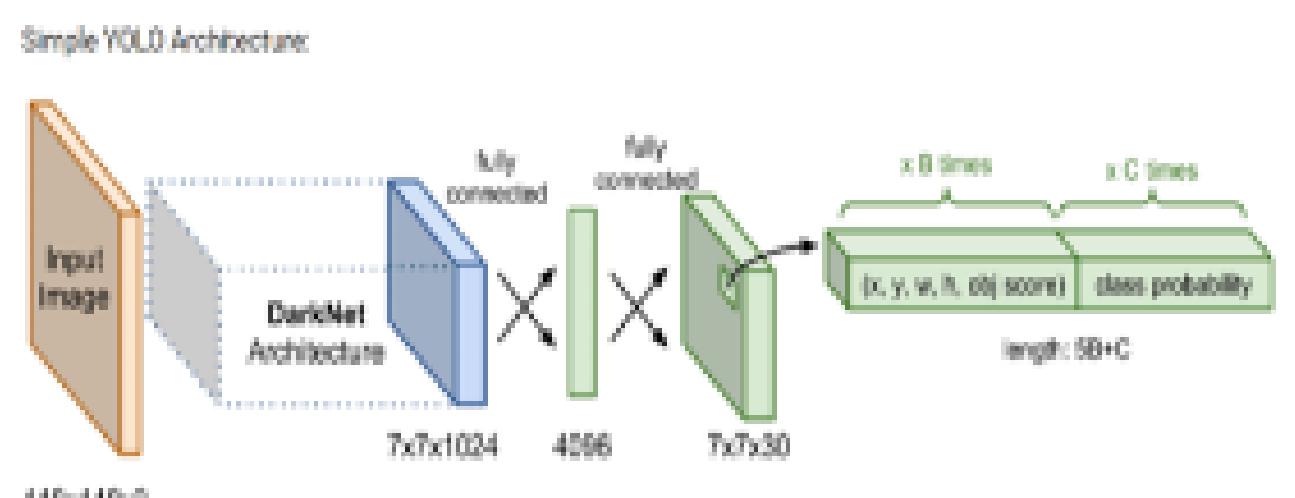
- The bounding box contains four values: x, y, w, h , where (x, y) represents the center of the box. (W, h) defines the width and height of the box.
- Confidence indicates the probability of containing objects in this prediction box, which is the IoU value between the prediction box and the actual box.
- The class probability indicates the class probability of the object, and the YOLOv3 uses a two-class method



A. WHAT IS OBJECT SEGMENTATION ?

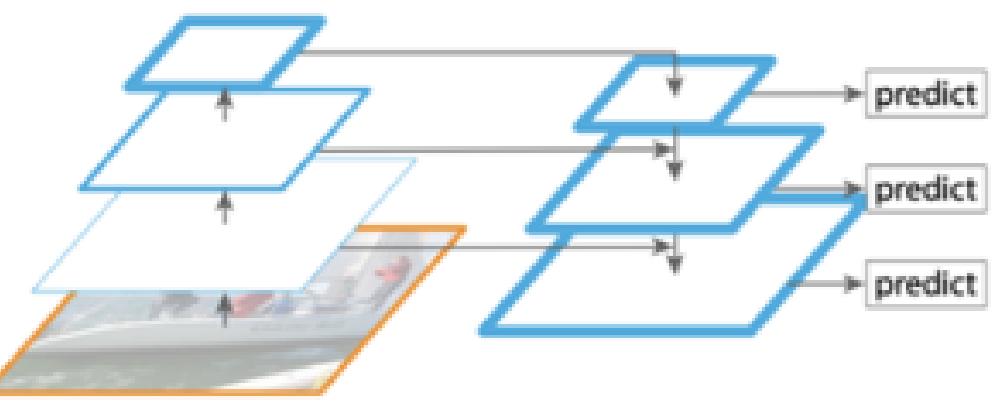
- Both one stage detection methods, like SSD and YOLO evaluate almost 10^4 to 10^5 candidate locations per image.
- But only a few locations contain objects (i.e. Foreground) and rest are just background objects.
- This leads to class imbalance problem.
- Small objects and close - by objects may be missed by YOLO like algorithms
- YOLO can detect a limited number of objects per image, making it less suitable for applications where a large number of objects need to be detected.

A. WHAT IS OBJECT SEGMENTATION ?



A. WHAT IS OBJECT SEGMENTATION ?

- RetinaNet is superior to single-stage object detection methods because it addresses the class imbalance problem, resulting in better accuracy.
- Its two-stage approach utilizes class-specific confidence scores to eliminate false positive detections.
- Additionally, RetinaNet only generates one detection per object, which leads to faster inference times than YOLO and SSD, which use anchor boxes to detect objects.

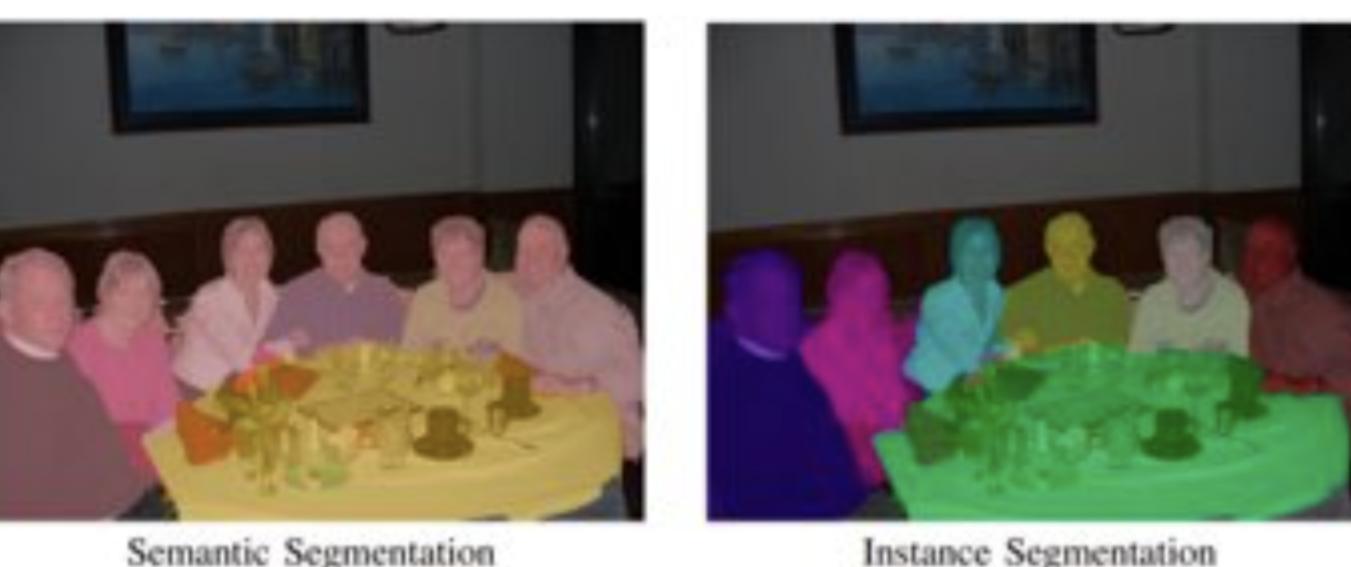


L9 : OBJECT SEGMENTATION CHEATSHEET

A. WHAT IS OBJECT SEGMENTATION ?

Image segmentation is the process of classifying each pixel in the image as belonging to a specific category.

- Semantic segmentation: We treat multiple objects within a single category as one entity
- Instance segmentation: We identify individual objects within these categories



B. TRANPOSED CONVOLUTION

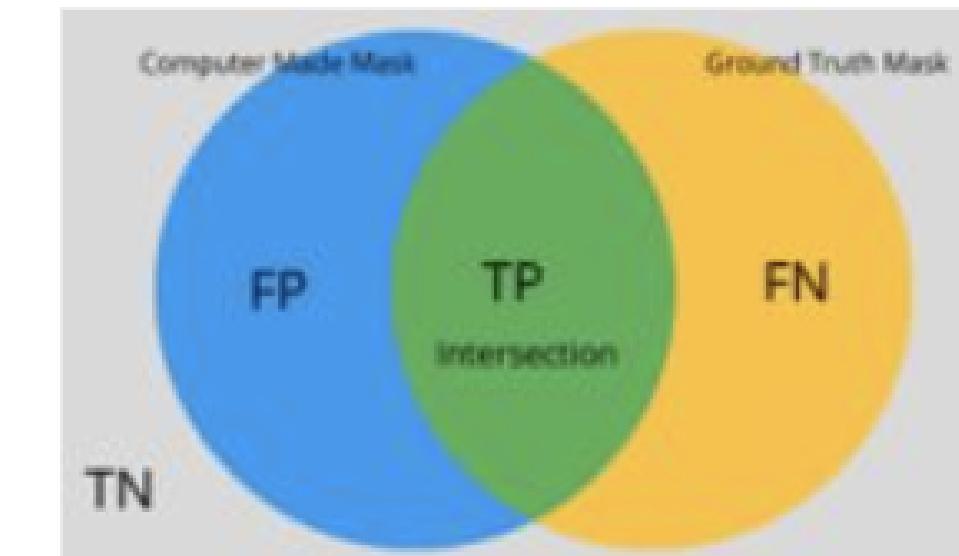
Transposed Convolutions is a method to up-sample the output. It can be considered as an opposite process to any simple CNN.

Output Shape :

$$\text{output size} = (\text{input} - 1) * \text{stride} - 2 * \text{padding} + (\text{kernel size} - 1) + 1$$

C. INTERSECTION OVER UNION

$$\frac{\text{Intersection}}{\text{Union}} = \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}$$

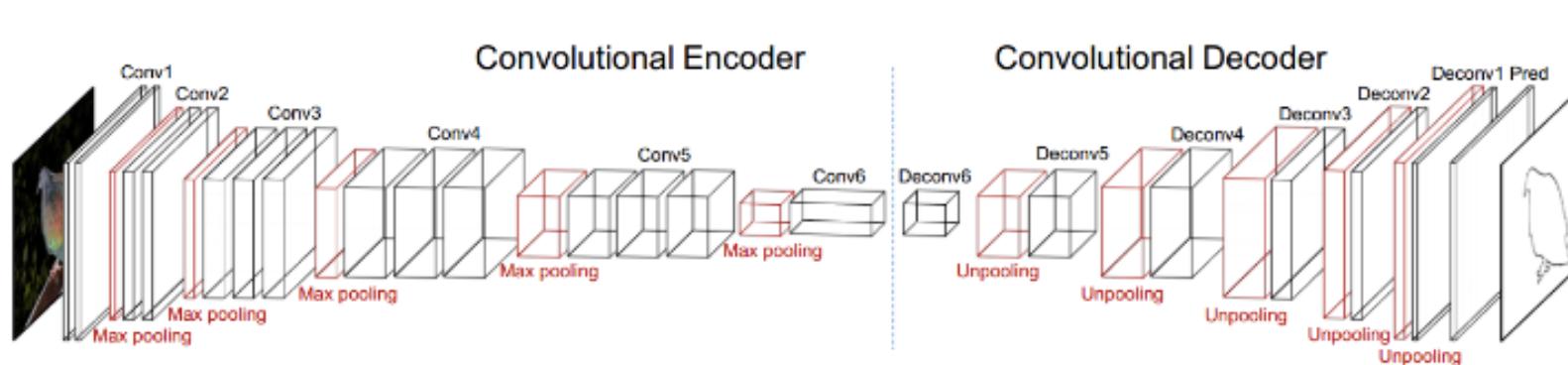


D. DICE COEFFICIENT

$$\text{Dice Coefficient} = \frac{2 \times \text{Intersection}}{\text{Union} + \text{Intersection}} = \frac{2\text{TP}}{2\text{TP} + \text{FN} + \text{FP}}$$

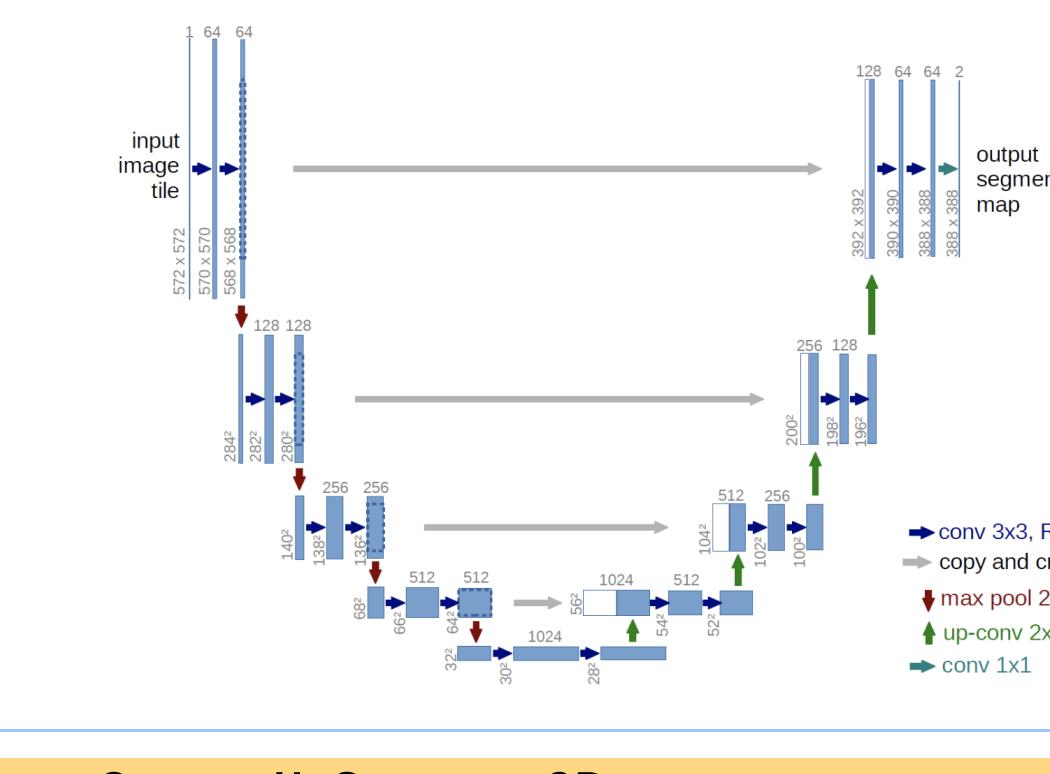
E. ENCODER-DECODER NETWORK

Encoder-Decoder Network is an architecture consisting of two parts, an encoder that compresses input into a lower dimensional representation, and a decoder that reconstructs the output from the encoded



F. U-NET ARCHITECTURE

- U-net Model has a "U" shape architecture with a symmetric Encoder and Decoder.
- It uses Skip Connections between layers of Encoder and Decoder are used to make the information loss as minimal as possible.
- The final output layer produces a per-pixel prediction of the target mask or segmentation.



G. WHAT IS OBJECT UPSAMPLING2D ?

UpSampling2D is a simple scaling up of images by using nearest neighbour or bilinear upsampling.

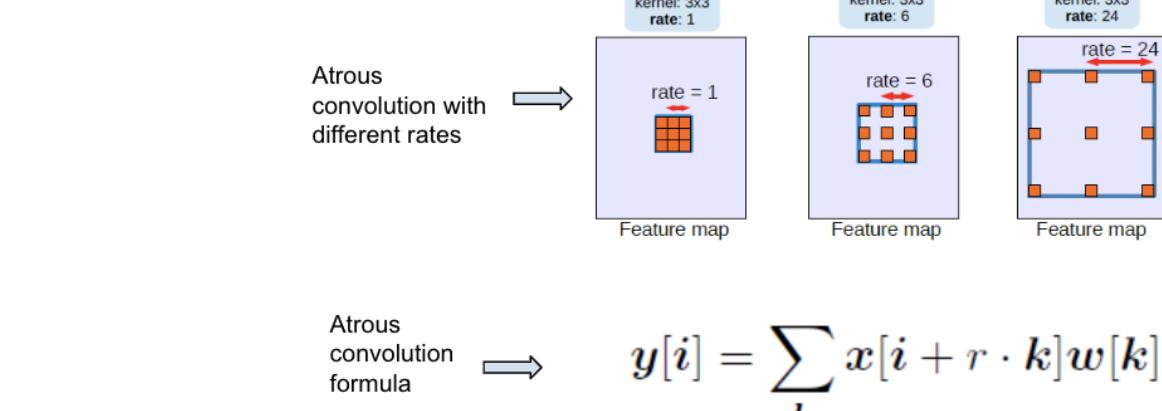
H. MASK R-CNN ARCHITECTURE

- Mask R-CNN is an extension of the Faster R-CNN model for object detection and instance segmentation.
- The RPN generates object proposals by predicting object scores and bounding box coordinates.
- The ROIAlign layer resamples object proposal features to a fixed size to ensure they can be processed by the fully connected layers.

I.

J. WHAT IS ATROUS CONVOLUTION ?

- For each location i on the output y and a filter w , atrous convolution is applied over the input feature map x where the atrous rate r corresponds to the stride with which we sample the input signal.
- It is also called dilated convolution.



L10: SIGNATURE VERIFICATION USING SIAMESE NETWORK CHEAT SHEET

A. WHAT IS A SIAMESE NETWORK?

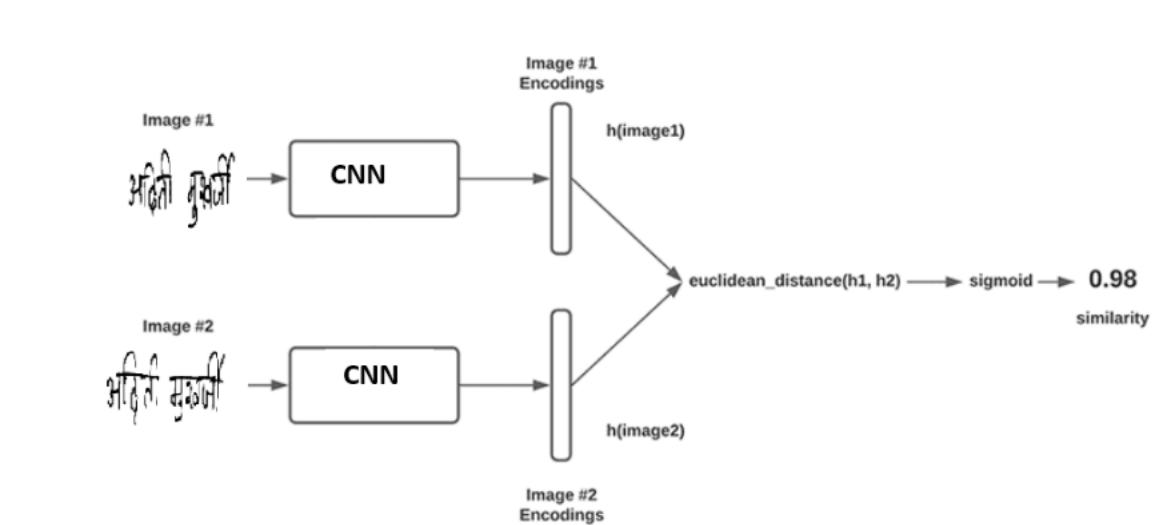
The Siamese Network architecture uses identical sub-networks to extract features from input data for similarity measurement or verification tasks, such as face recognition and signature verification.

B. CONTRASTIVE LOSS

- The contrastive loss function uses two feature representations and a binary label indicating similarity to compute Euclidean distance and apply a penalty.
- Similar pairs are encouraged to be close and dissimilar pairs far. It trains a Siamese Network to learn a discriminative embedding space for accurate similarity measurement, one-shot learning, and verification.



C. ARCHITECTURE OF SIAMESE NETWORK



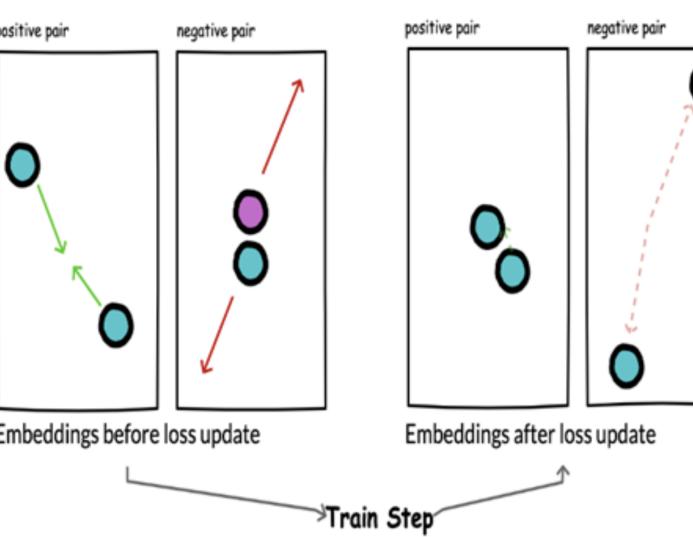
C. CONTRASTIVE LOSS EQUATION

$$Y * D^2 + (1 - Y) * \max(margin - D, 0)^2$$

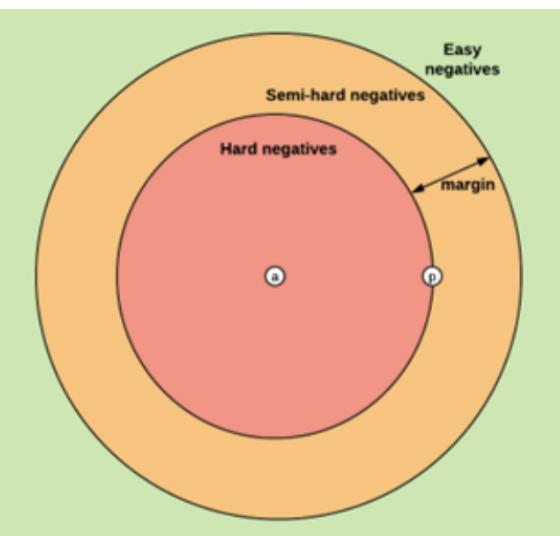
- The label Y is 1 for the same-class image pairs and 0 for different-class pairs.
- The Euclidean distance between sister network embeddings is variable D .
- The max function selects the larger value between 0 and the difference between margin m and distance.

COMPUTER VISION CHEATSHEET

D. EMBEDDINGS BEFORE AND AFTER CONTRASTIVE LOSS



J. TRIPLET MINING



E. TRIPLET LOSS

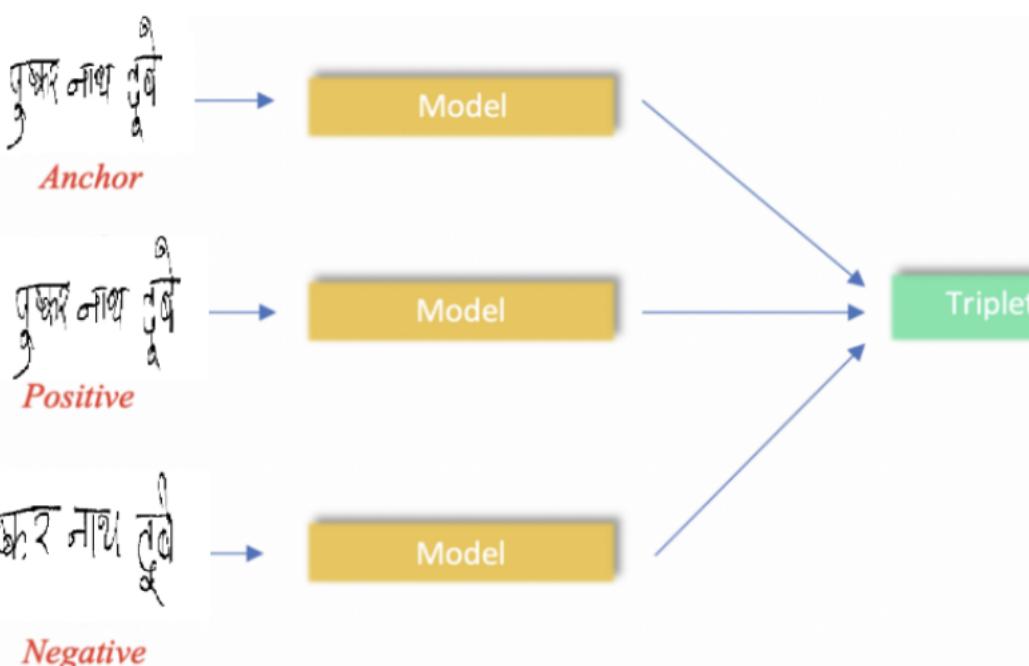
- A triplet includes an anchor, positive sample, and negative sample.
- The anchor and positive samples are similar, while the anchor and negative samples differ.
- The triplet loss function minimizes the difference in distances between anchor-positive and anchor-negative pairs, encouraging embeddings where the anchor-positive distance is less than anchor-negative by a margin.

F. TRIPLET LOSS EQUATION

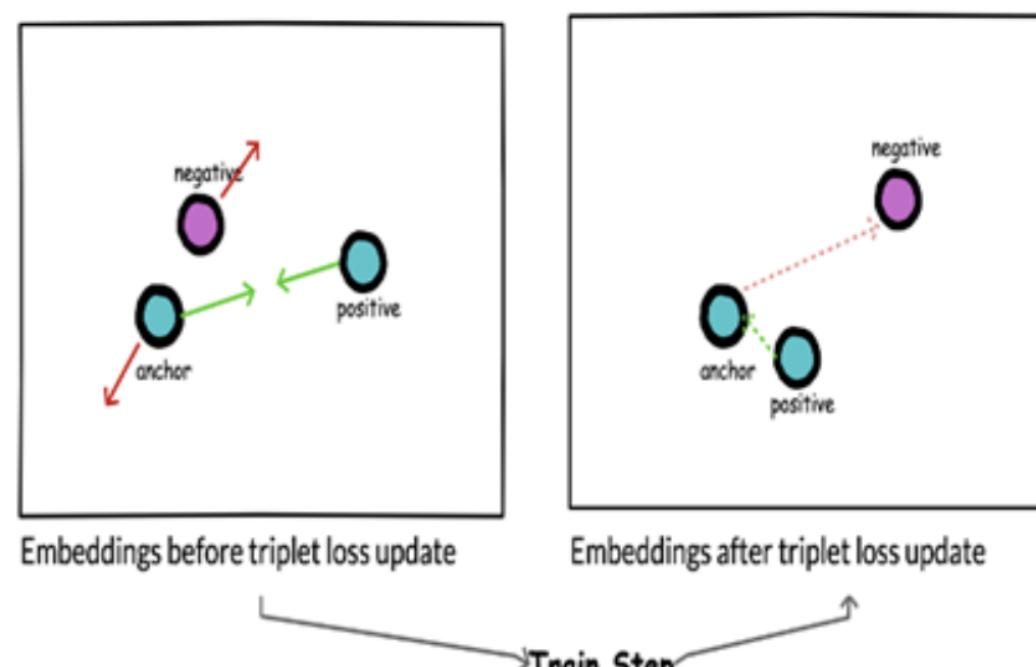
$$L(A, P, N) = \max(0, d(A, P) - d(A, N) + \alpha)$$

- $d(A,P)$ is anchor-to-positive-distance and $d(A,N)$ is anchor-to-negative-distance and α is the margin

G. TRIPLET LOSS



H. EMBEDDINGS BEFORE AND AFTER TRIPLET LOSS



I. TYPES OF TRIPLET MINING

- Easy triplets: $d(A,N) > d(A,P) + \alpha$, loss=0, nothing to learn.
- Hard triplets: $d(A,N) < d(A,P)$, high loss, backpropagation.
- Semi-hard triplets: $d(A,P) < d(A,N) < d(A,P) + \alpha$, positive loss.

C. WHAT IS DISCRIMINATIVE LOSS?

One common example of a discriminative loss function is binary cross-entropy loss, which is defined as:

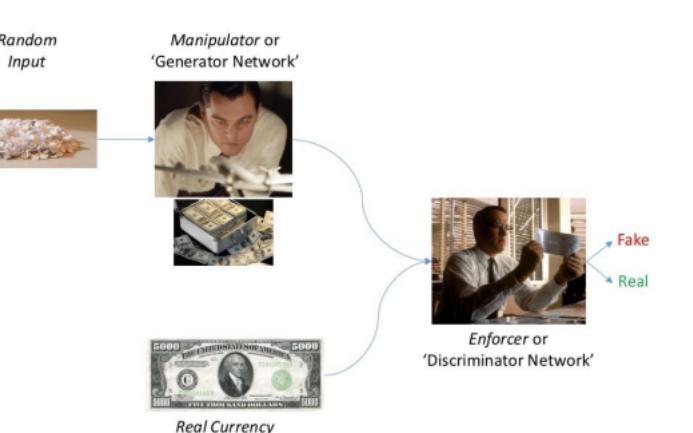
$$\text{Loss} = -[y * \log(D(x)) + (1-y) * \log(1 - D(x))]$$

E. WHAT IS GENERATIVE LOSS?

The generative loss in GANs is a scalar value that represents the discrepancy between the generated samples and the real data. The goal of the generator is to minimize this loss so that the generated samples become as similar as possible to the real data.

$$\text{Loss} = -(1/m) \sum [y \log(p) + (1-y) \log(1-p)]$$

F. SIMPLE WORKING OF GANS



G. APPLICATION OF GANS

- Image Generation
- Image Translation
- Data Augmentation
- Text Generation
- Video Generation

H. GANS APPLICATIONS EXAMPLE - 1

CycleGAN: CycleGAN is a GAN that is used for image-to-image translation, where the goal is to transform an image from one domain to another, such as converting a horse into a zebra, or a photograph into a painting.



I. GANS APPLICATIONS EXAMPLE - 2

SRGAN is a Generative Adversarial Network used for single image super-resolution. It increases the resolution of an image, making it appear clearer and more detailed.



J. GANS APPLICATIONS EXAMPLE - 3

StyleGAN: StyleGAN is a Generative Adversarial Network (GAN) used for synthesizing new images of human faces that are highly realistic and diverse.