



Department of Computational and Data Sciences



Computer Vision

Deepak Subramani

Assistant Professor

Dept. of Computational and Data Science

Indian Institute of Science Bengaluru



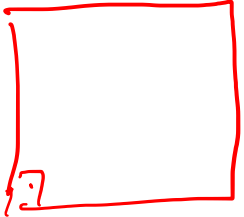
Lecture and Assignment Guide

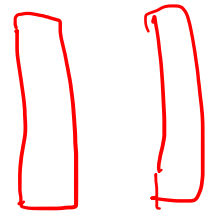
- This Slide Deck has Material for 6 hours of teaching divided into Parts 1-6
- We will go through
 - Week 01
 - Part 01 - Convolutional and Pooling Layers; AST 01
 - Part 02 - Transfer Learning and Modern CV Design Principle; AST 02
 - Week 02
 - Part 01 - Modern Convolutional Building Blocks for Image Classification; AST 03
 - Part 02 - Object Localization
 - Interpreting what convolutions learn (Advanced topic) – AST 03
 - Week 03
 - Part 01 - Object Detection (YOLO), Image Segmentation – Lec 05
 - Part 02 - Practical CVOps
 - AST04 – Object Detection with YOLO
 - Week 04
 - Revision
 - AST05 – Image Segmentation
- Additional Reading material to go in depth of math with references and code references are provided with the marking of “Additional Material” or “Additional Discussion” etc

Cross Entropy

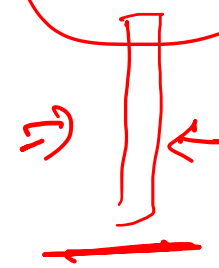
→ Encoder

Encoding

⇒ 
~~224~~ 224 × 224 × 3
 1.5 Lakh
 ≈ 150k

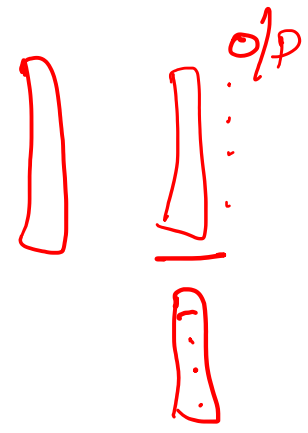


Vector



→ 1000
 → 1024

d-dim



neurons
 # classes

Softmax

Reality

TP	FP
FN	TN

Pre

Representation

"Class imbalance"

Feature extractor

Precision → Prediction

Triplet loss, Contrastive loss

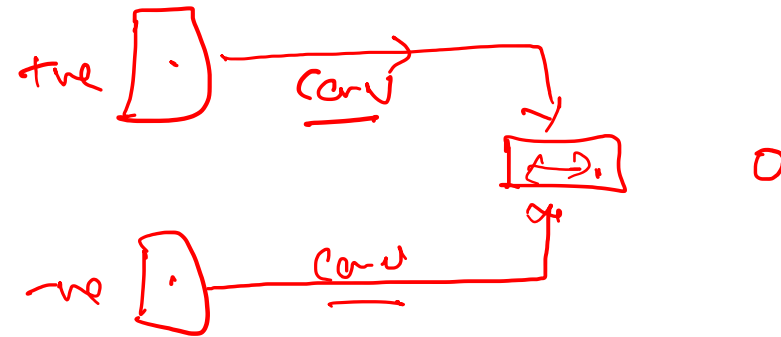
Recall → Reality

(+ve class, +ve class, -ve class)
-ve class, +ve class, -ve class

Focal Loss

A, A, B, C, ..., N

Siamese



Sup Con Loss



Department of Computational and Data Sciences

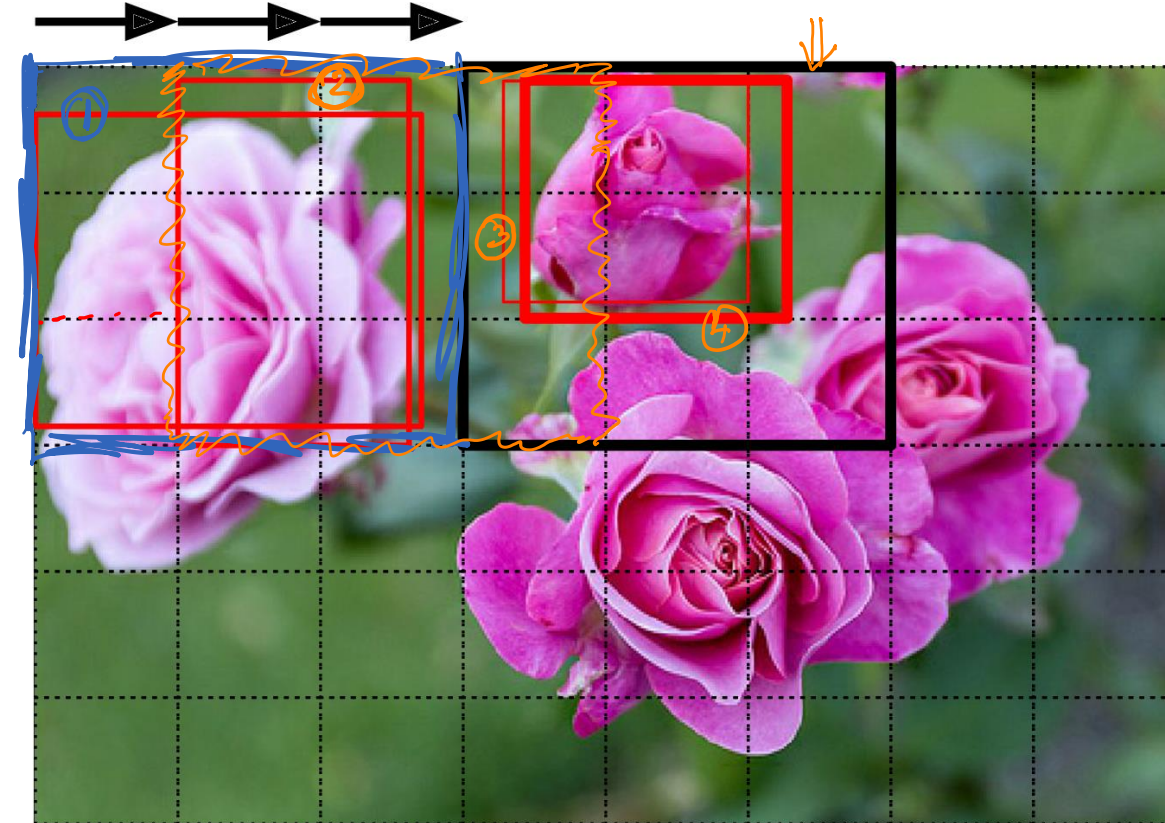


Week 03 – Part 01

Deepak Subramani
Assistant Professor
Dept. of Computational and Data Science
Indian Institute of Science Bengaluru

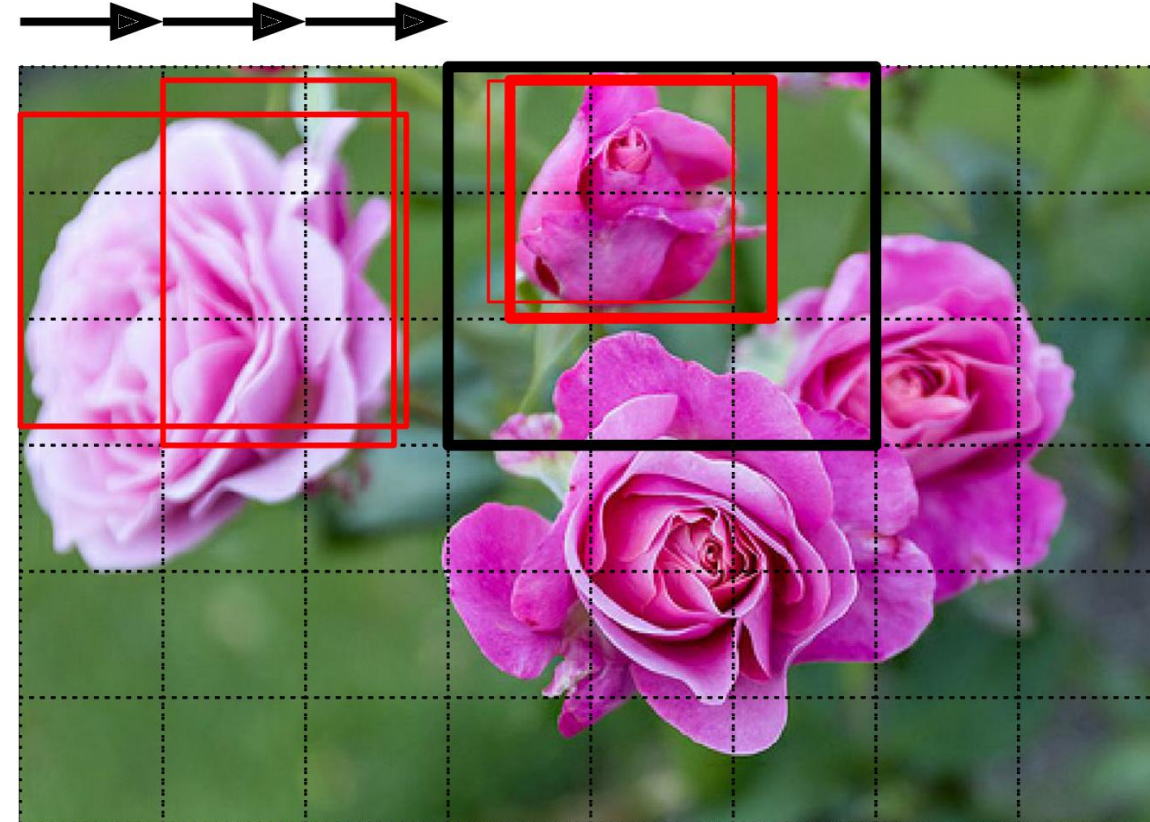
Object Detection

- Object Detection: The task of simultaneously classifying and localizing multiple objects in an image
- Earlier approaches involved sliding a single object detector CNN over an image and use post processing
- Steps
 - Slide bounding box CNN over the image
 - Find bounding boxes
 - Post process



Object Detection

- Steps
 - Divide the image into grids
 - Slide an object detector over each 3x3 region
 - Predict one bounding box per 3x3 region
 - Repeat for 4x4 etc
 - In addition to a bounding box, output “objectness” – the probability of the object being present in the bounding box
 - Sigmoid activation trained with binary CE
 - Postprocess to get rid of overlapping bounding boxes
 - Drop BBox with low objectness
 - Find and drop BBox that overlap with the BBox with highest objectness
 - Repeat until no more boxes can be removed



Fully Convolutional Network

- The simple approach in the previous slides are powerful, but slow
- We need to run the CNN many times over the image
- FCN (Fully Convolutional Network) solves this problem
- Key idea: Replace the top dense layers with convolutional layers

- Example – Conv Base C10 followed by D1

$7 \times 7 \times 100 = 4900$

200

- C10 – 100 \times 7 \times 7 [100 maps each of 7×7]
- D1 – 200 neurons – Each neuron has a connection to 100 \times 7 \times 7 from C10 + bias
- Replace D1 with C11 – 200 convolutional filters 7 \times 7+1(V) – Output is 1 \times 1 \times 200
- This C11 is doing the same job as D1 but using convolutions now
 - V – Valid Padding



C11
7 \times 7

FCN: What is the big deal?

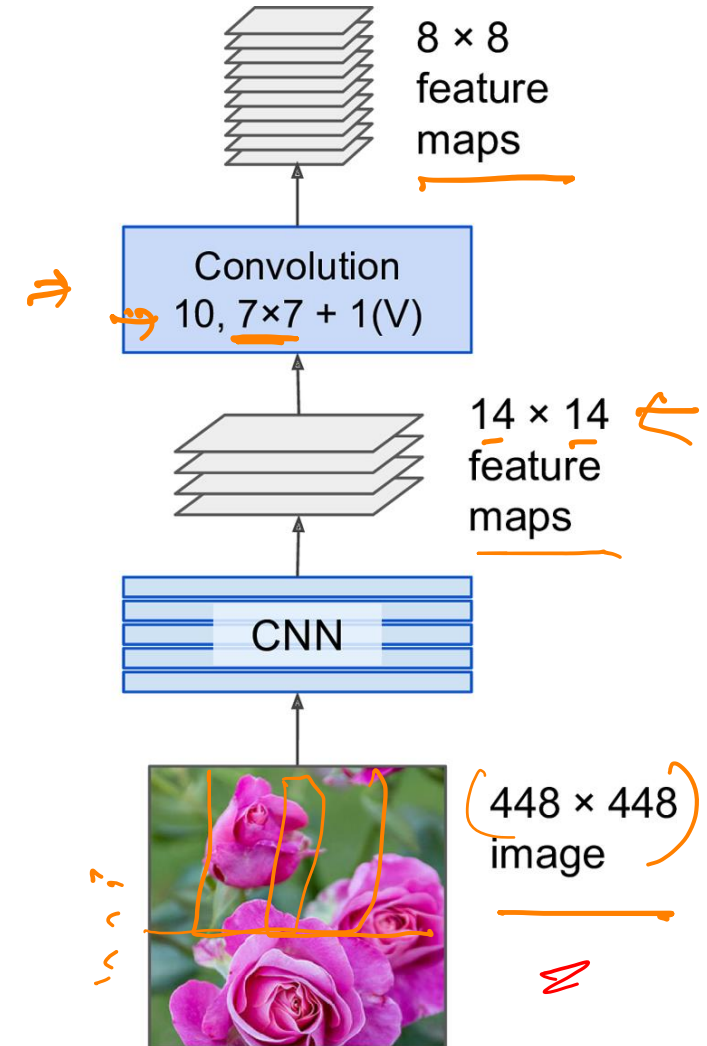
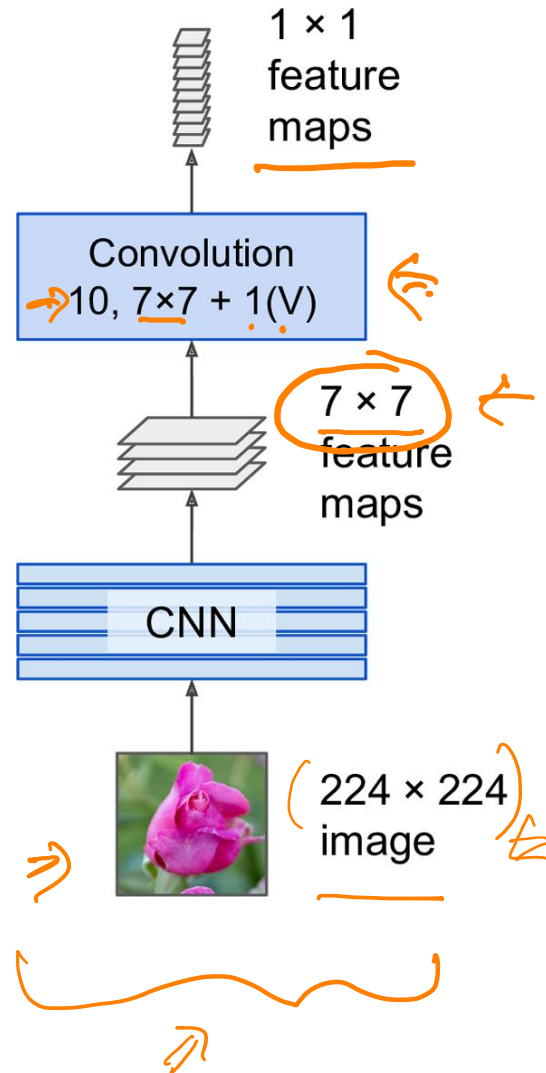
- Dense Layer expects a specific input size
 - In this example: It is $100 \times 7 \times 7$
 - But a Convolution will work even on $100 \times 14 \times 14$ and anything else!
 - It only needs that 100 to be intact
- Let us look at the numbers
 - Trained CNN for flower object detection with input size 224×224
 - 10 outputs on the Dense Layer head – 0-4 (softmax) 1 per class [5 classes], 5 (logistic) objectness score, 6-9 (none) x coordinate, y coordinate, height, width of bounding box
 - Suppose the Conv Base ends at 7×7
 - Replace Dense head with Conv Layer 10 filters of $7 \times 7 + 1(V)$

Poll

1. A Fully Convolutional Network is a more efficient network for object detection
 - True, False
2. If a dense layer is replaced with a FCN after training, then the weights can be copied from Dense to the FCN
 - True, False
3. IOU is the typical loss function used for object detection
 - True, False

FCN Example (Cont)

- What happens if we feed 448x448 images to this FCN?
 - Last conv layer is 14x14, and it will produce a 8x8 map
 - What is this 8x8 map? – It is equivalent to sliding the original CNN across the image
- Now the network has to be run only once
- You Only Look Once (YOLO)



11am

YOLO: Introduction

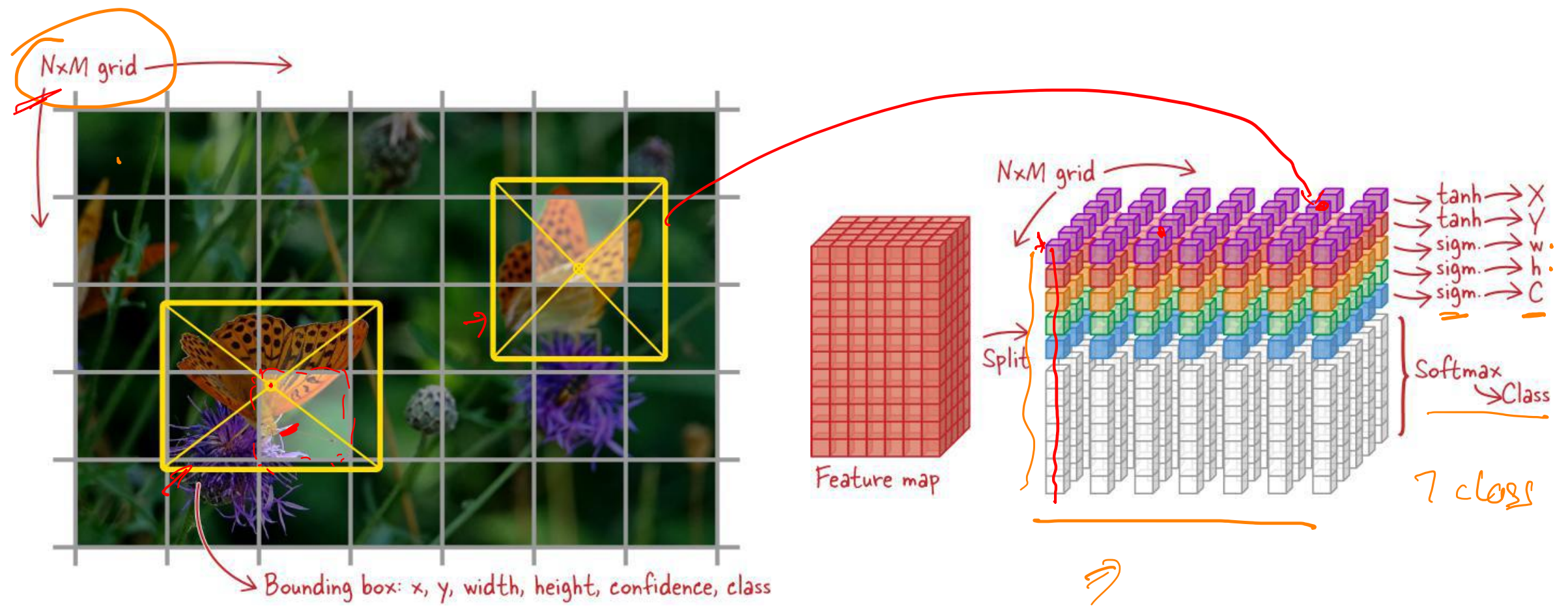
- YOLO (and its variants) – Introduced between 2015-18
- State of the art for real time object detection
- Demo of YOLO V3 - <https://www.youtube.com/watch?v=MPU2HistivI>
- Trained on the PASCAL (Pattern Analysis, Statistical Modeling and Computational Learning) VOC (Visual Object Class) Dataset
- PASCAL VOC challenge is another machine vision challenge for object detection and segmentation
 - Recall: Image net for object classification

YOLO: Basic Idea

- Similar to the algorithm we discussed so far, but with a few engineering modifications – No new concept
- Outputs five bounding boxes (instead of 1) and 20 class probabilities (PASCAL VOC dataset)
- Before training, a K-Means clustering is applied to the training dataset bounding boxes and the centroid of 5 representative bounding boxes are found. This serves as the BBox prior, called as anchor boxes
- The network then predicts the factor by which the anchor box must be resized
- During training, new image dimensions are randomly chosen so that the network learns to deal with different scales

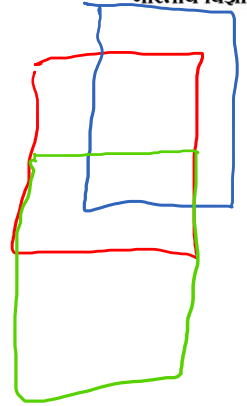


YOLO Visually



Evaluation Metrics

- Object detection is more complex to evaluate than image classification
- We use the following definitions
 - **True Positive (TP):** A correct detection. Detection with $\text{IOU} \geq \text{threshold}$
 - **False Positive (FP):** A wrong detection. Detection with $\text{IOU} < \text{threshold}$
 - **False Negative (FN):** A ground truth not detected
 - **True Negative (TN):** Does not apply. It would represent a corrected misdetection. In the object detection task there are many possible bounding boxes that should not be detected within an image. Thus, TN would be all possible bounding boxes that were correctly not detected (so many possible boxes within an image). That's why it is not used by the metrics.

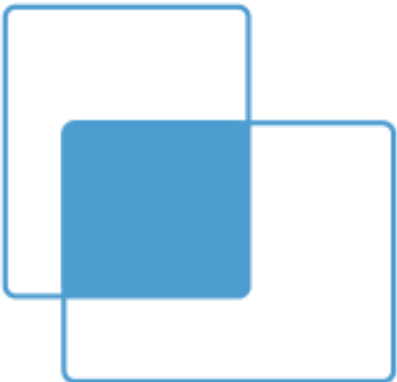


threshold: depending on the metric, it is usually set to 50%, 75% or 95%.

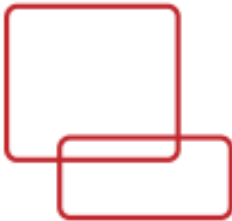


IoU

$$\text{IoU} = \frac{\text{Area of intersection}}{\text{Area of union}}$$



Examples



IoU = 0.1



IoU = 0.3



IoU = 0.6

Accuracy of Object Detection System: mAP

- Metric for Object Detection Accuracy: Mean Average Precision
- Why mean average?
- Precision-Recall Tradeoff
 - $\text{Precision} = \text{TP} / (\text{TP} + \text{FP}) = \text{TP} / (\text{all detections})$
 - $\text{Recall} = \text{TP} / (\text{TP} + \text{FN}) = \text{TP} / (\text{all ground truths})$
 - Higher the precision, lower the recall.
 - But sometimes there is no tradeoff, especially at low ends. This leads to mAP
- Sometimes 90% precision at 10% recall and 96% precision at 20% recall.
Then we should use the 96,20 operating point

Accuracy of Object Detection System: mAP

- To generalize the situation we should look at the maximum precision with atleast some recall (say 20%)
- We keep changing this recall threshold 0-100 and collect the maximum precision possible. Then report the average. This is AP
- If we have more classes, then we have AP for each class. So we take mean AP
- In object detection there is an added complexity – accuracy depends on the IoU also.
- So we report mAP@IoU. We also take $\text{mean}(\text{mAP}@[\text{IoU range}])$ – mean mean average precision!
- Example calculation: <https://github.com/rafaelpadilla/Object-Detection-Metrics>



Department of Computational and Data Sciences

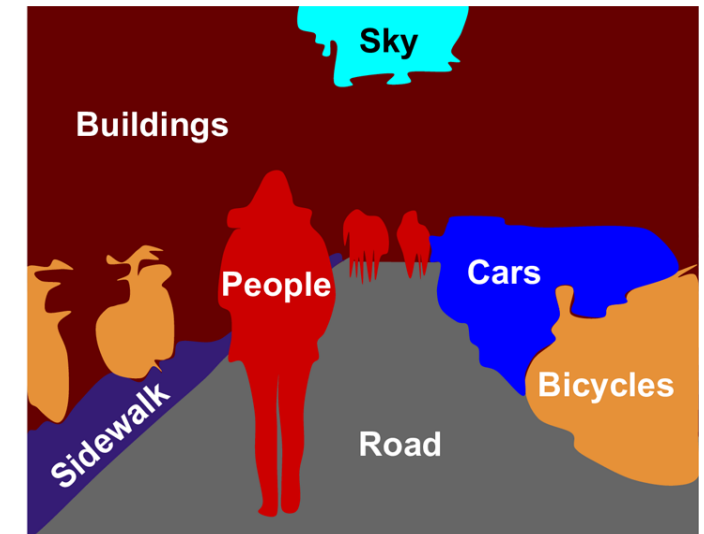
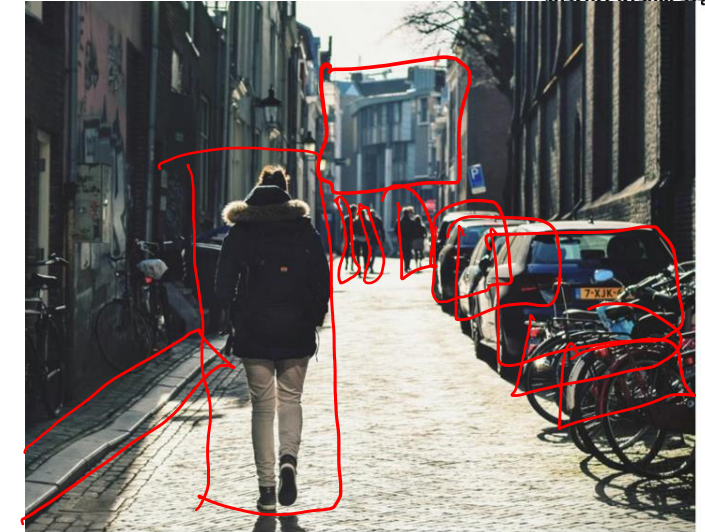


Week 03 Part 02

Deepak Subramani
Assistant Professor
Dept. of Computational and Data Science
Indian Institute of Science Bengaluru

Segmentation

- Semantic Segmentation: Each pixel is classified and given a class label
- Different objects of the same class are not distinguished
 - Two people standing together will be identified as one big lump of people!
- Instance Segmentation: Each object is distinguished from the other
- The most popular architecture is called Deep Lab (current SOTA)
- U-Net is a popular architecture used in Biomedical Image segmentation and recently for Satellite Image Analysis
- Two concepts to be learnt here
 - ✓ Transposed Convolution, with skip connection
 - ✓ Atrous Convolution

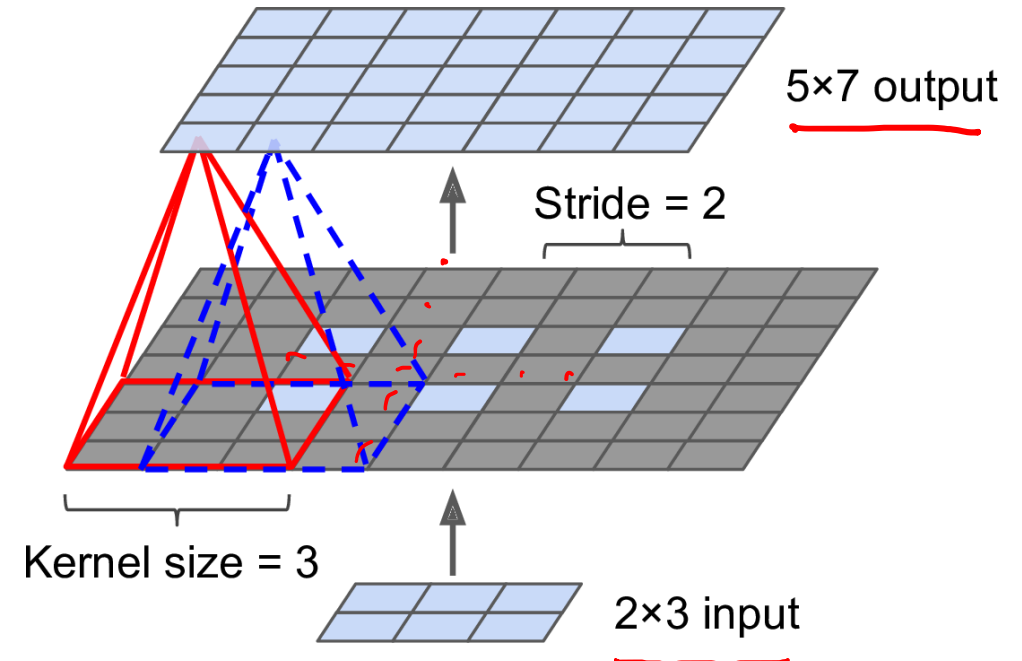


Transposed Convolution: Motivation

- To learn multiscale features efficiently, stride > 2 is necessary in the convolutional base
- But, then how to classify each pixel?
- This needs an upsampling layer
- The idea of Transposed Convolution comes from this need
- The simplest way is a bilinear interpolation, but that works only for 4x or even 8x, but not beyond
- Long et al 2015, introduced this Transpose Convolution, in addition to the FCN that we saw earlier

Transposed Convolution

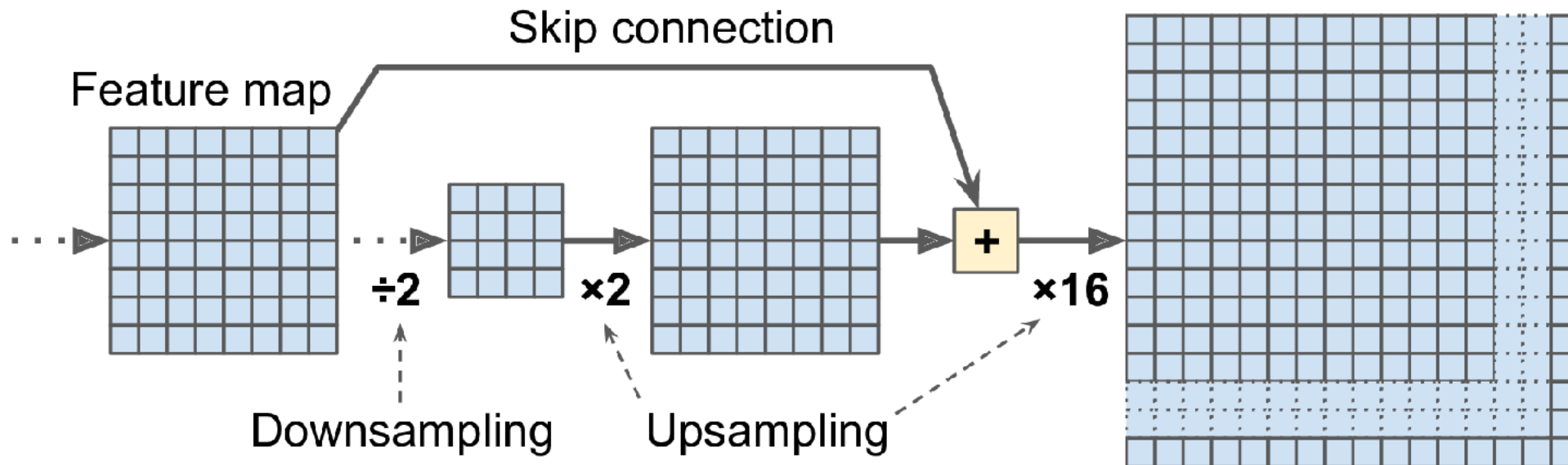
- Think of stretching an image by adding empty rows and columns
- Then on the stretched image do a regular convolution
- Initialize these kernels to do a linear interpolation
- But as the weights are learnable, it does better!





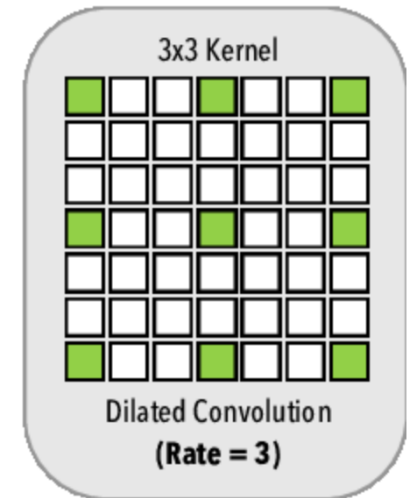
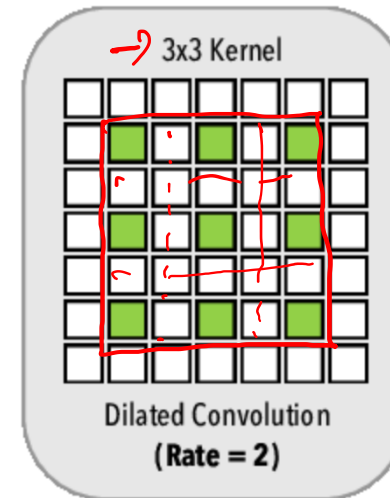
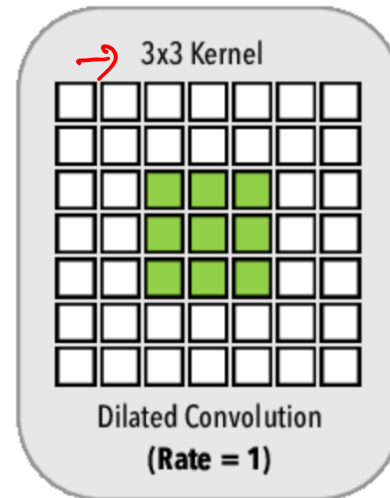
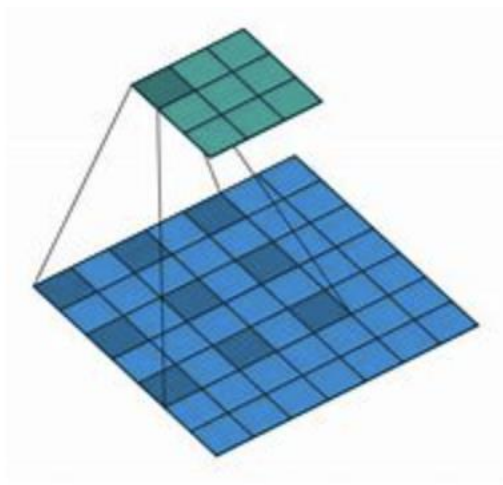
Skip Connection

- The transpose convolution introduces difficulties in gradient flow
- Skip connections are introduced to make the down sampling kernels learn



Atrous Convolution

- The convolution field of view is modified by considering a larger area with zeros added to the filter itself
- Number of learnable parameters is the same as regular convolution, but now the field of view has changed
- This is used in Deep Lab

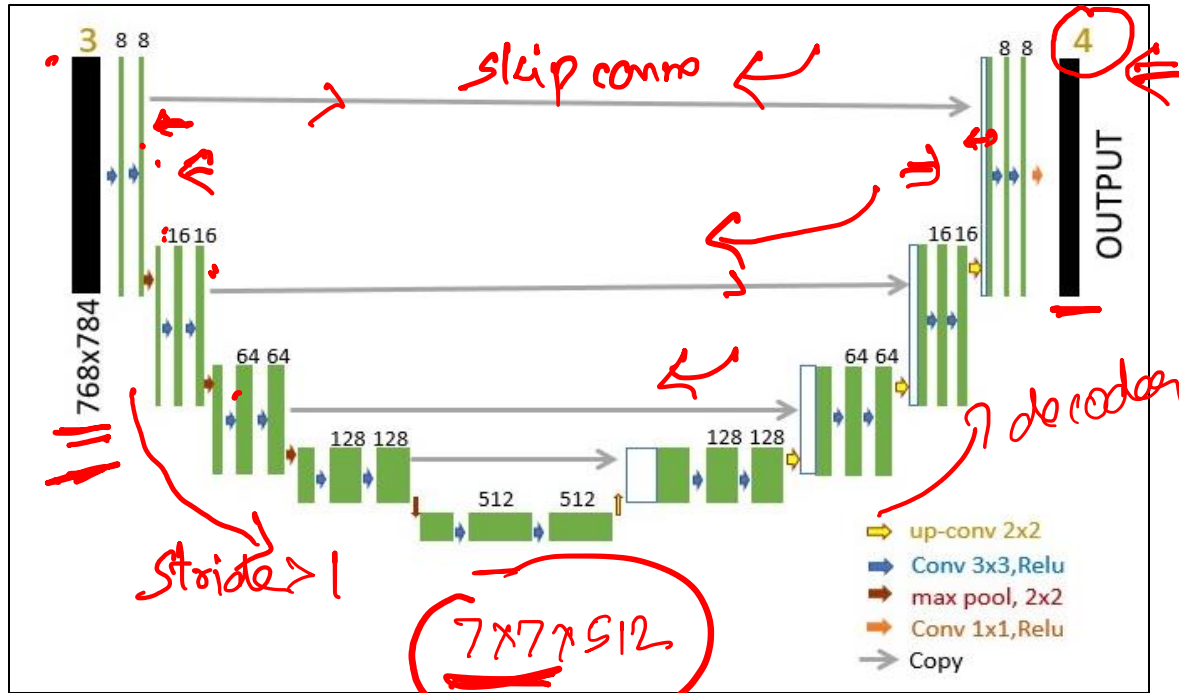


Poll

1. In semantic segmentation, two objects of the same class are distinguished
 - True, False
2. FCN is used for semantic segmentation
 - True, False

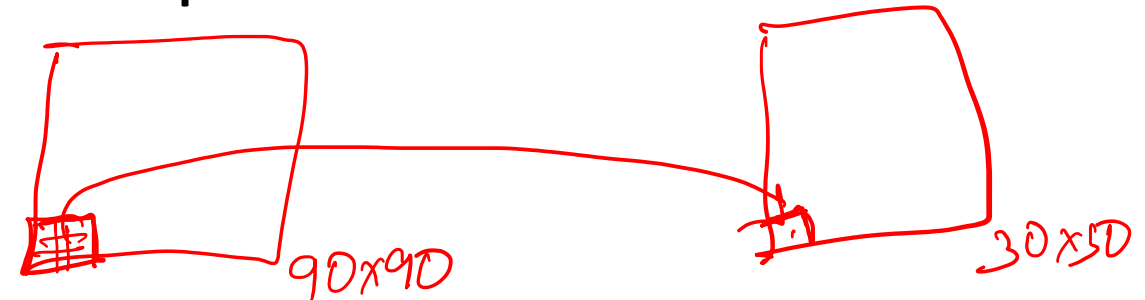


Biomedical Image Segmentation – U nets



U-net architecture by Ronneberger and Olaf. 2015.

- Error/Loss is estimated by comparing target and output of U-net.
- Number of maps in output are equal to number of classes in the problem.



State of the Art as of Dec 2024

- <https://paperswithcode.com/area/computer-vision>
- Image Classification –
 - For speed: Efficient Net (CNN),
 - For accuracy: Vision Transformer }
- Semantic Segmentation –
 - For speed: Unet, DeepLabv3
 - For accuracy: Deformable Convolution (InternImage), Vision Transformer (Segment Anything) SAM
- Instance Segmentation - Mask-R-CNN, RetinaNet
- Object Detection –
 - For speed: YOLOv8 ||
 - For accuracy: Deformable Convolution (InternImage)
 - DETR

Other Tasks in CV

- Pose Estimation
- Object Tracking
- Action Recognition
- Motion Estimation
- Monocular Depth
- Content-aware Image Editing
- Scene Reconstruction (NeRF)