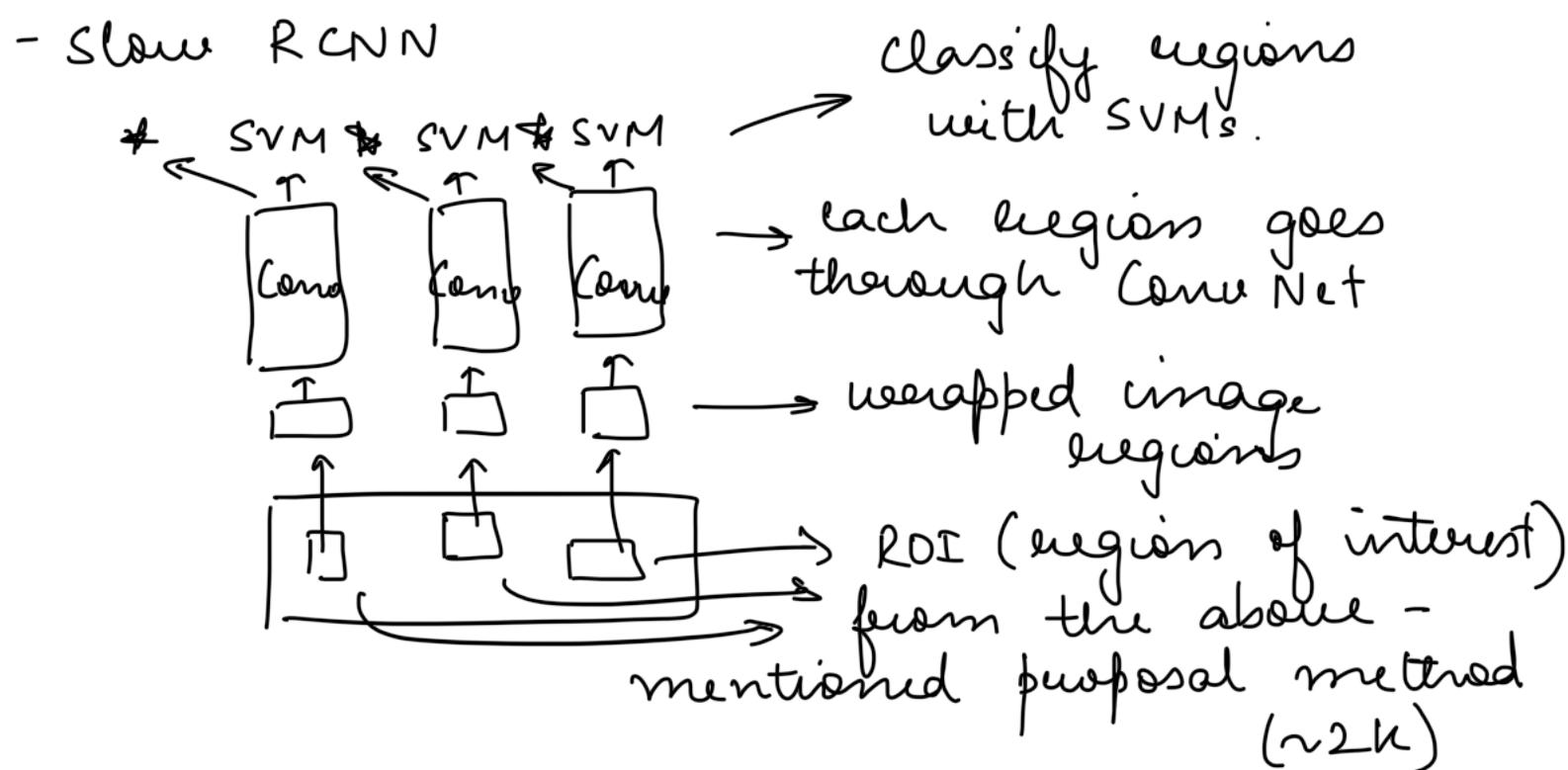


Week-10 DLP

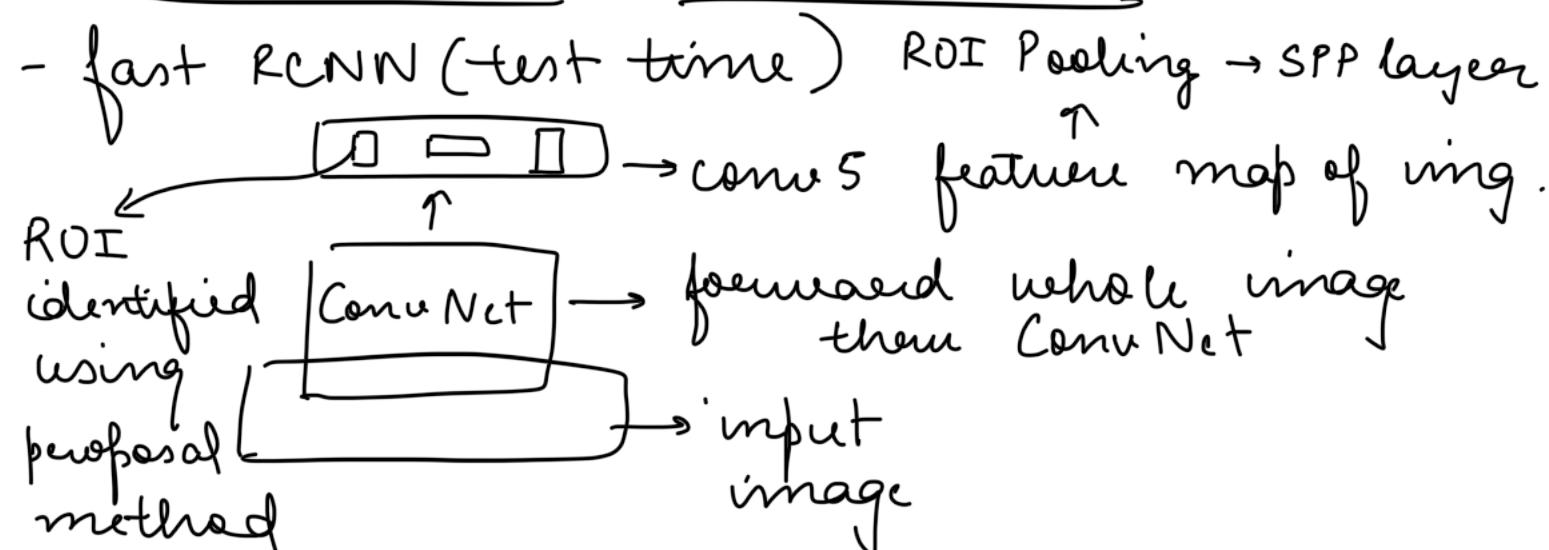
L1 Object detectn - Introduction

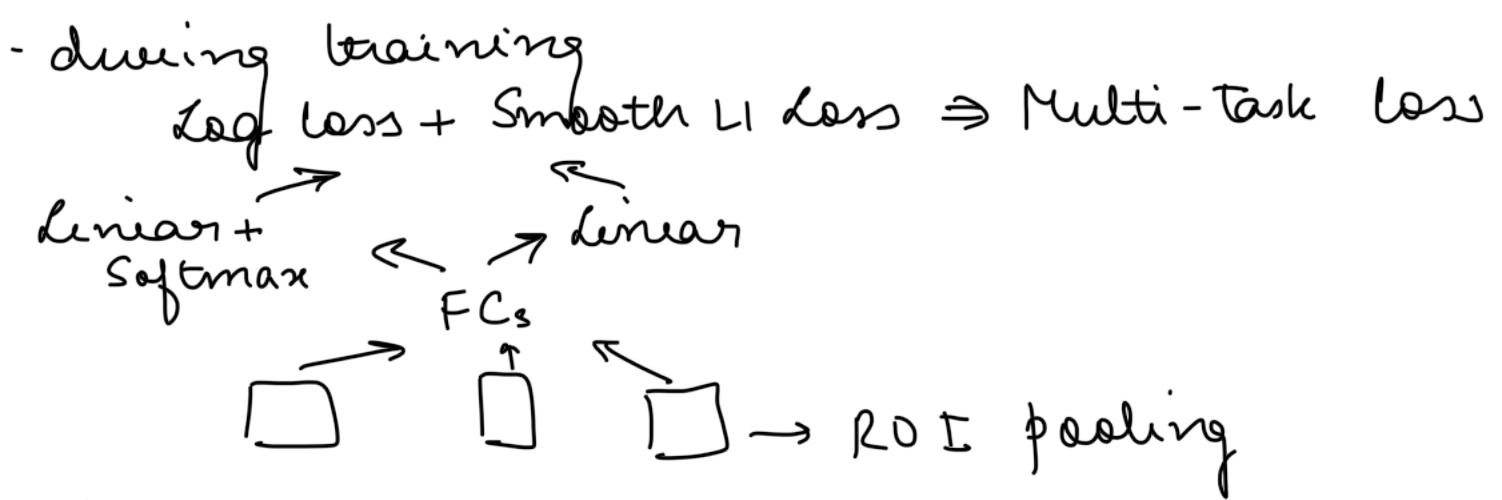
- Object detectn \rightarrow import. task in CV.
 - ↳ The task is to assign a label & a bounding box to all the objects that are present in the image.
- Labels = [1, 2, 3]
either from 1, 2 & 3.
- Basically, it will multiply classes of image classificn. The problem is too many positions and scales to test. If your classifier is fast enough, you can do it.
- R-CNN (Region-based ConvNet) - ConvNets are computationally demanding. But we can't test all positions and scales.
 - ↳ Soln - We will only look at a subset of positions. Choosing them wisely.
- Region proposals \rightarrow find blobs that will most likely to contain objects. The proposal is 'class-agnostic'.
- Process \rightarrow
 1. input image
 2. extract reg. proposals ($\sim 2k$)
 3. compute the CNN features
 4. classify regions



- * BBox reg. \Rightarrow apply bounding-box regressors.
- problems with slow R-CNN.
 - \rightarrow log loss, hinge loss, squared loss.
 - \rightarrow training is slow (84 h) + lot of disk space.
 - \rightarrow inference (detchⁿ) is slow. current rate is 47 s / image. It was later fixed by SPP-net.

L2 Fast R CNN & Faster R CNN



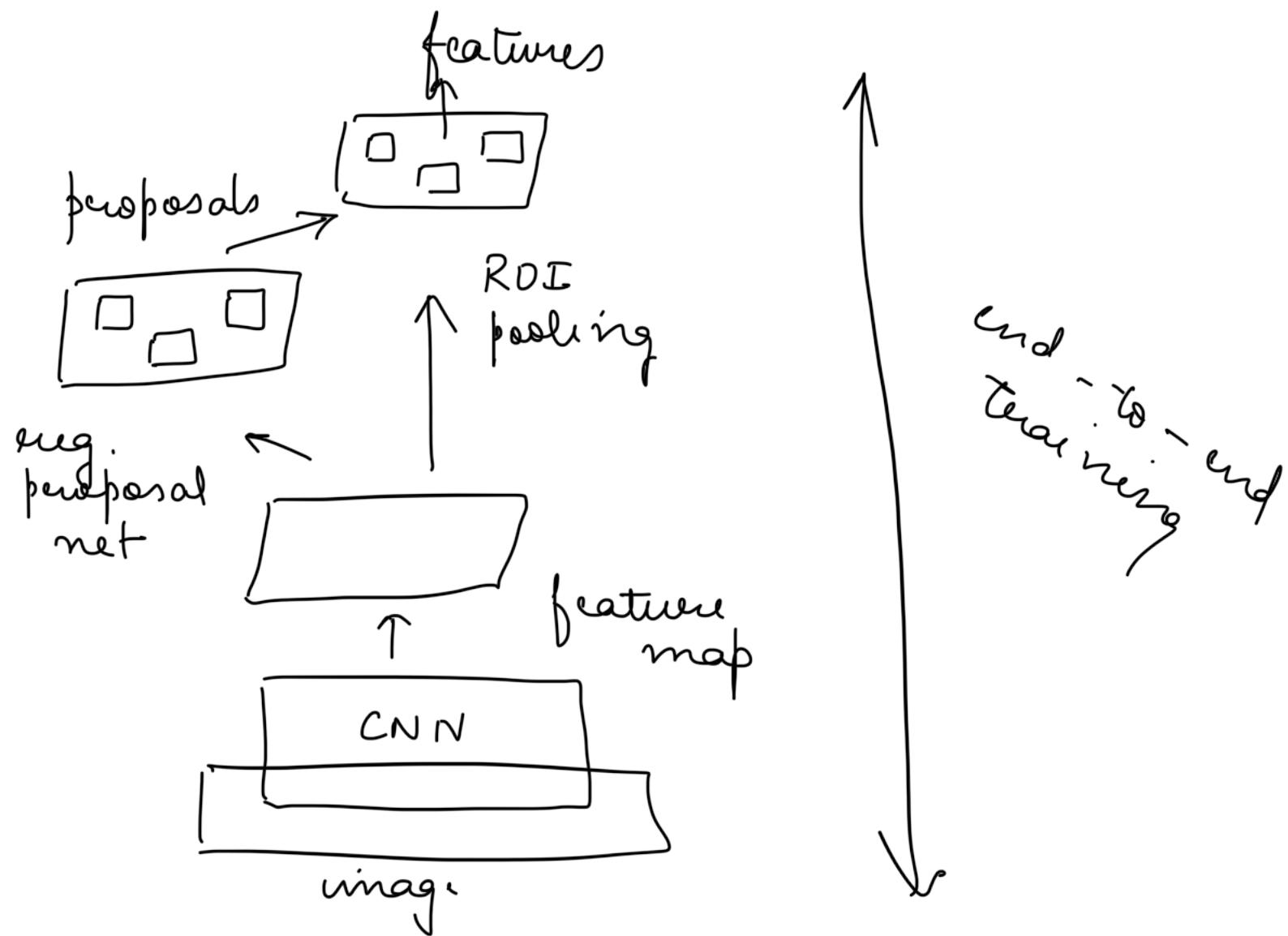


- Cons -

1. out of network regions proposals.
Selective search. 2s / image.

- Faster RCNN

- solely based on CNN.
- no external modules.
- each step is end-to-end.



- Region proposal network
 - slide a small window on the feature map.
 - posⁿ of the window with ref. to image. classify obj./not obj. regress box locations
- (n) scores coordinates
 ↘ ↗
 256 → (4n)
 ↑ →
 feature map n anchors
- it is fully convolutional, it is trained end-to-end.
 - RPN shares convolutional features maps with the detectⁿ network.

L3 YOLO - You Only Look Once

- It is much advantageous over RCNN.
 1. a single neural network for localization and classification
 2. inference happens only once.
 3. looks at the entire image each time leading to less false positives.
- Steps -
 1. image is split into $S \times S$ cells.
 2. for each cell, generate B bounding boxes.
 3. for each bounding box, we have 5 predictions $\rightarrow x, y, w, h, \text{ and confidence.}$

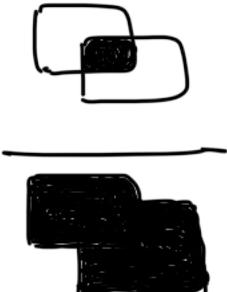
4. also, we do object classification.
- each cell predicts a class probability conditioned on object being present.
e.g. $P(\text{Car} | \text{Object})$
 - $P(\text{Class} | \text{Object}) \times P(\text{Object}) = P(\text{Class})$
- Yolo architecture -
1. input $\rightarrow 448 \times 448 \times 3$
 2. 24 conv. layers
 3. $S = 7, B = 2, C = 20$ (default)
 4. output is $\boxed{S \times S \times (5B + C)}$

L4 Object detection Metrics & Experimental Results

- we use MAP (mean avg. precision)
- | | |
|--|---|
| $\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$ | Actual
+ -
Predicted
+ TP FP
- FN TN |
| $\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$ | |

$$F1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- precision - refers to finding TP out of all positive predictions

- recall - how to find the TP out of all predictions
 - If $\text{IOU} >$ certain threshold (0.5), the prediction is TP. If $\text{IOU} \leq$ threshold, the prediction is a FP.
- $$\text{IOU} = \frac{\text{Area of Overlap}}{\text{Area of Unions}}$$
- 
- AP (Avg. Precision) - area under the precision-recall curve. Avg. of precision values across a range of recall values. This is called AP for that specific class.
 - MAP - mean of the AP values across all object classes.
 - mIoU analysis -
 1. makes far less background errors
 2. but far more localization errors.