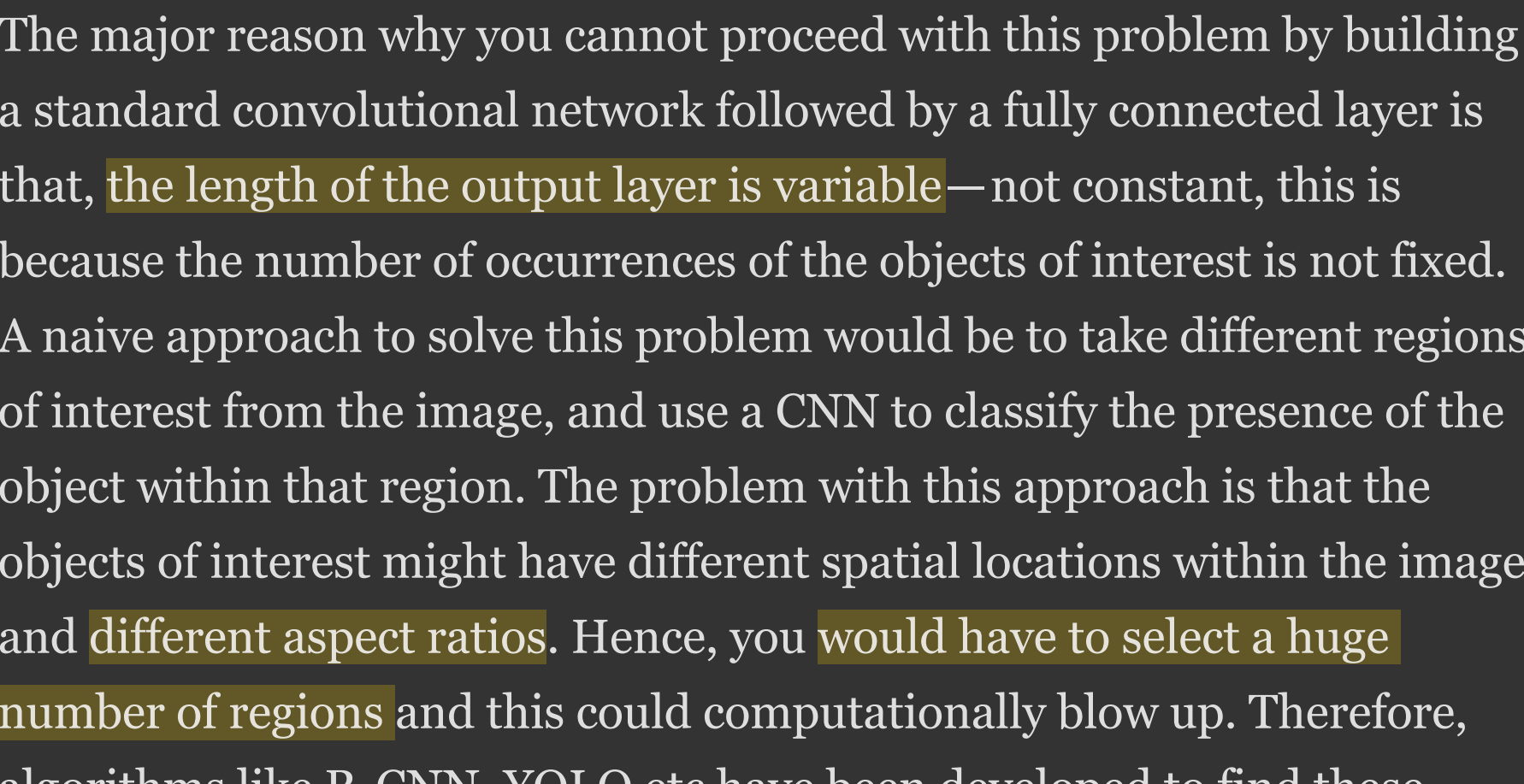


Introduction

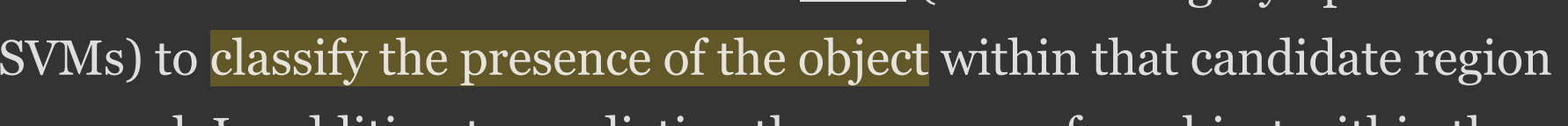
Computer vision is an interdisciplinary field that has been gaining huge amounts of traction in the recent years(since CNN) and self-driving cars have taken centre stage. Another integral part of computer vision is object detection. Object detection aids in pose estimation, vehicle detection, surveillance etc. The difference between object detection algorithms and classification algorithms is that in detection algorithms, we try to draw a bounding box around the object of interest to locate it within the image. Also, you might not necessarily draw just one bounding box in an object detection case, there could be many bounding boxes representing different objects of interest within the image and you would not know how many beforehand.



The major reason why you cannot proceed with this problem by building a standard convolutional network followed by a fully connected layer is that, **the length of the output layer is variable**—not constant, this is because the number of occurrences of the objects of interest is not fixed. A naive approach to solve this problem would be to take different regions of interest from the image, and use a CNN to classify the presence of the object within that region. The problem with this approach is that the objects of interest might have different spatial locations within the image and **different aspect ratios**. Hence, you **would have to select a huge number of regions** and this could computationally blow up. Therefore, algorithms like R-CNN, YOLO etc have been developed to find these occurrences and find them fast.

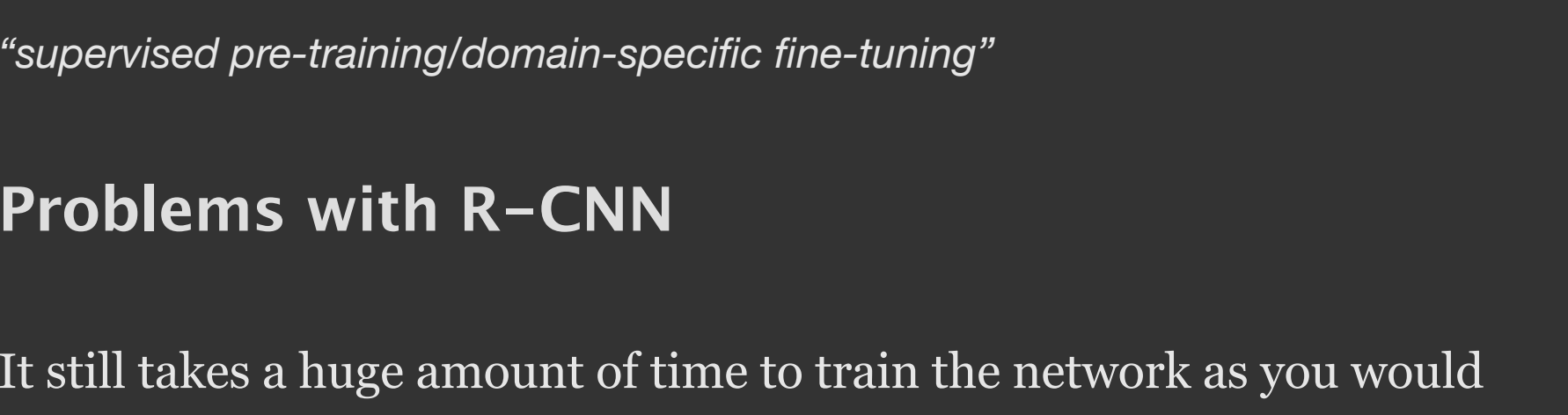
R-CNN

To bypass the problem of selecting a huge number of regions, [Ross Girshick et al.](#) proposed a method where we use selective search to extract just 2000 regions from the image and he called them region proposals. Therefore, now, instead of trying to classify a huge number of regions, you can just work with 2000 regions. These 2000 region proposals are generated using the selective search algorithm which is written below.



R-CNN

To know more about the selective search algorithm, follow this [link](#). These 2000 candidate region proposals are warped into a square and fed into a convolutional neural network that produces a 4096-dimensional feature vector as output. The CNN acts as a feature extractor and the output dense layer consists of the features extracted from the image and the extracted features are fed into an SVM (a set of category-specified SVMs) to **classify the presence of the object** within that candidate region proposal. In addition to predicting the presence of an object within the region proposals, the algorithm also predicts four values which are offset values to **increase the precision of the bounding box**. For example, given a region proposal, the algorithm would have predicted the presence of a person but the face of that person within that region proposal could’ve been cut in half. Therefore, the offset values help in adjusting the bounding box of the region proposal.



R-CNN

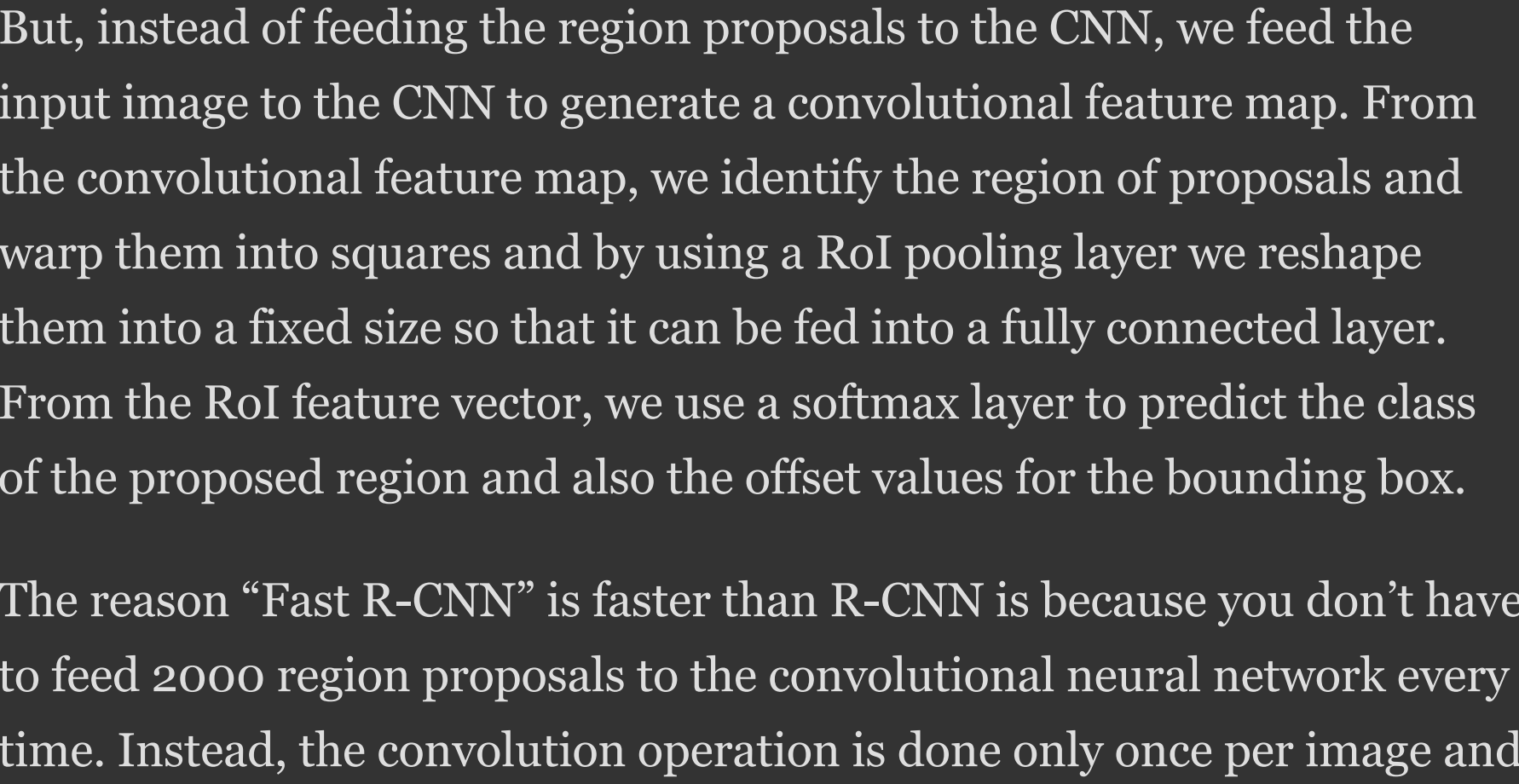
“supervised pre-training/domain-specific fine-tuning”

Problems with R-CNN

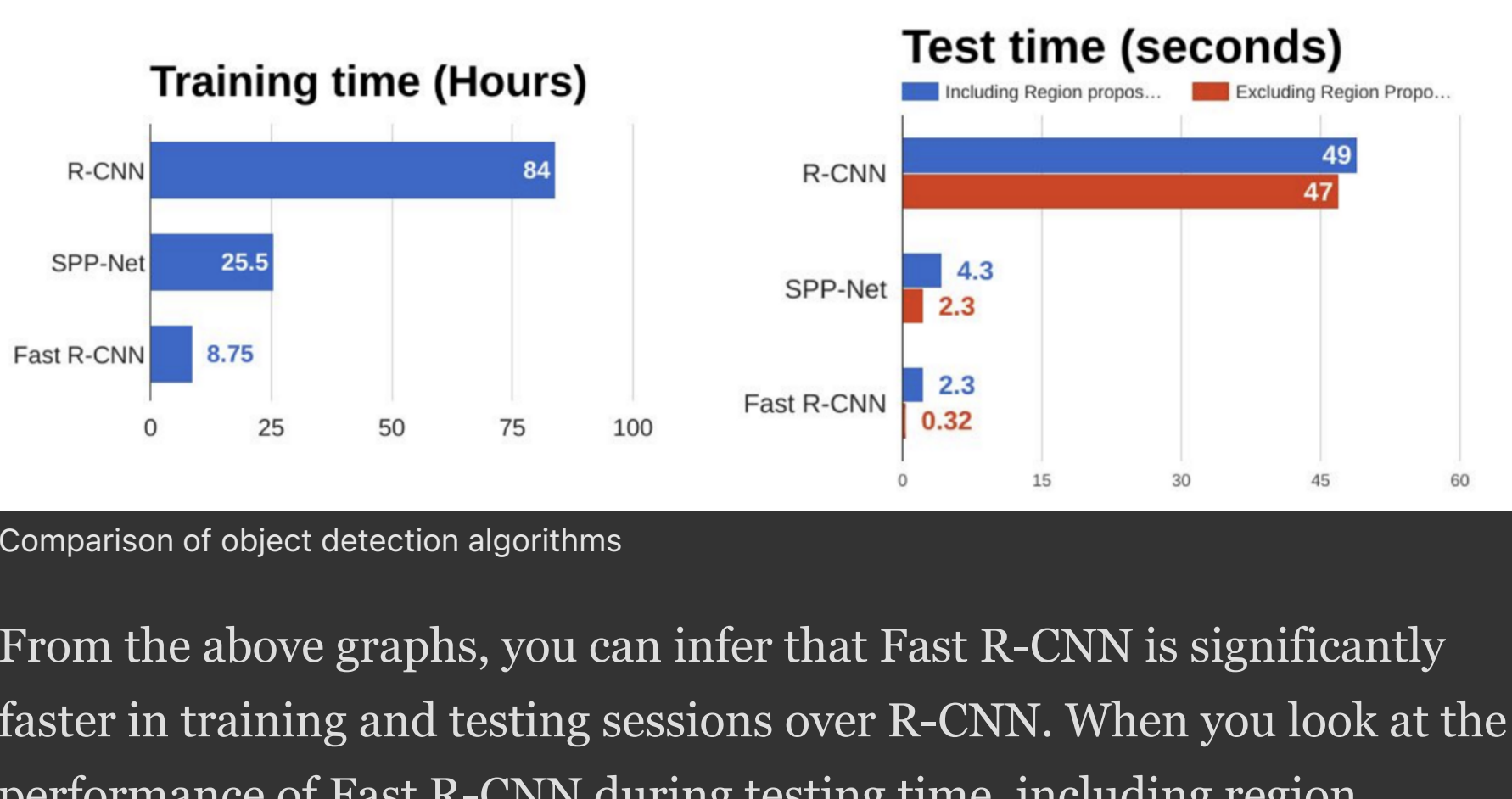
It still takes a huge amount of time to train the network as you would have to classify 2000 region proposals per image. It cannot be implemented real time as it takes around **47 seconds for each test image**.

The selective search algorithm is a fixed algorithm. Therefore, no learning is happening at that stage. This could lead to the generation of bad candidate region proposals.

Fast R-CNN



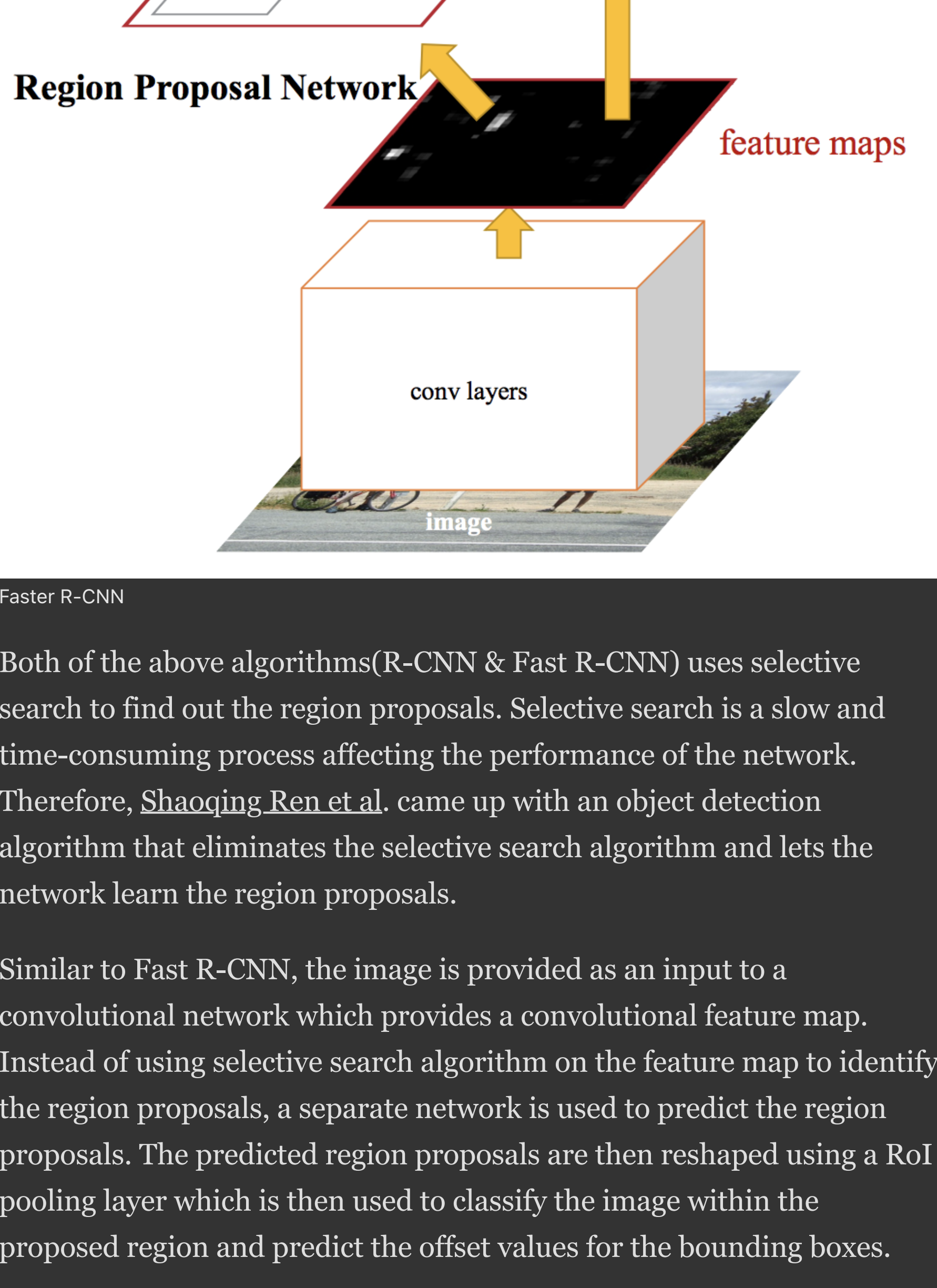
Fast R-CNN



Comparison of object detection algorithms

From the above graphs, you can infer that Fast R-CNN is significantly faster in training and testing sessions over R-CNN. When you look at the performance of Fast R-CNN during testing time, including region proposals slows down the algorithm significantly when compared to not using region proposals. Therefore, region proposals become bottlenecks in Fast R-CNN algorithm affecting its performance.

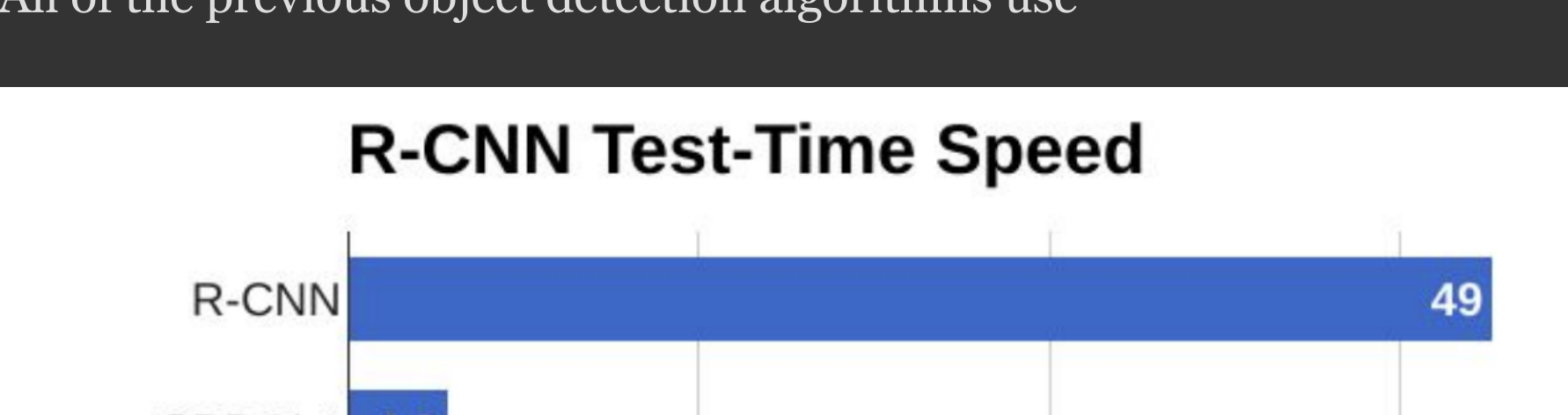
Faster R-CNN



Faster R-CNN

Both of the above algorithms(R-CNN & Fast R-CNN) uses selective search to find out the region proposals. Selective search is a slow and time-consuming process affecting the performance of the network. Therefore, [Shaoqing Ren et al.](#) came up with an object detection algorithm that eliminates the selective search algorithm and lets the network learn the region proposals.

Similar to Fast R-CNN, the image is provided as an input to a convolutional network which provides a convolutional feature map. Instead of using selective search algorithm on the feature map to identify the region proposals, a separate network is used to predict the region proposals. The predicted region proposals are then reshaped using a RoI pooling layer which is then used to classify the image within the proposed region and predict the offset values for the bounding boxes.

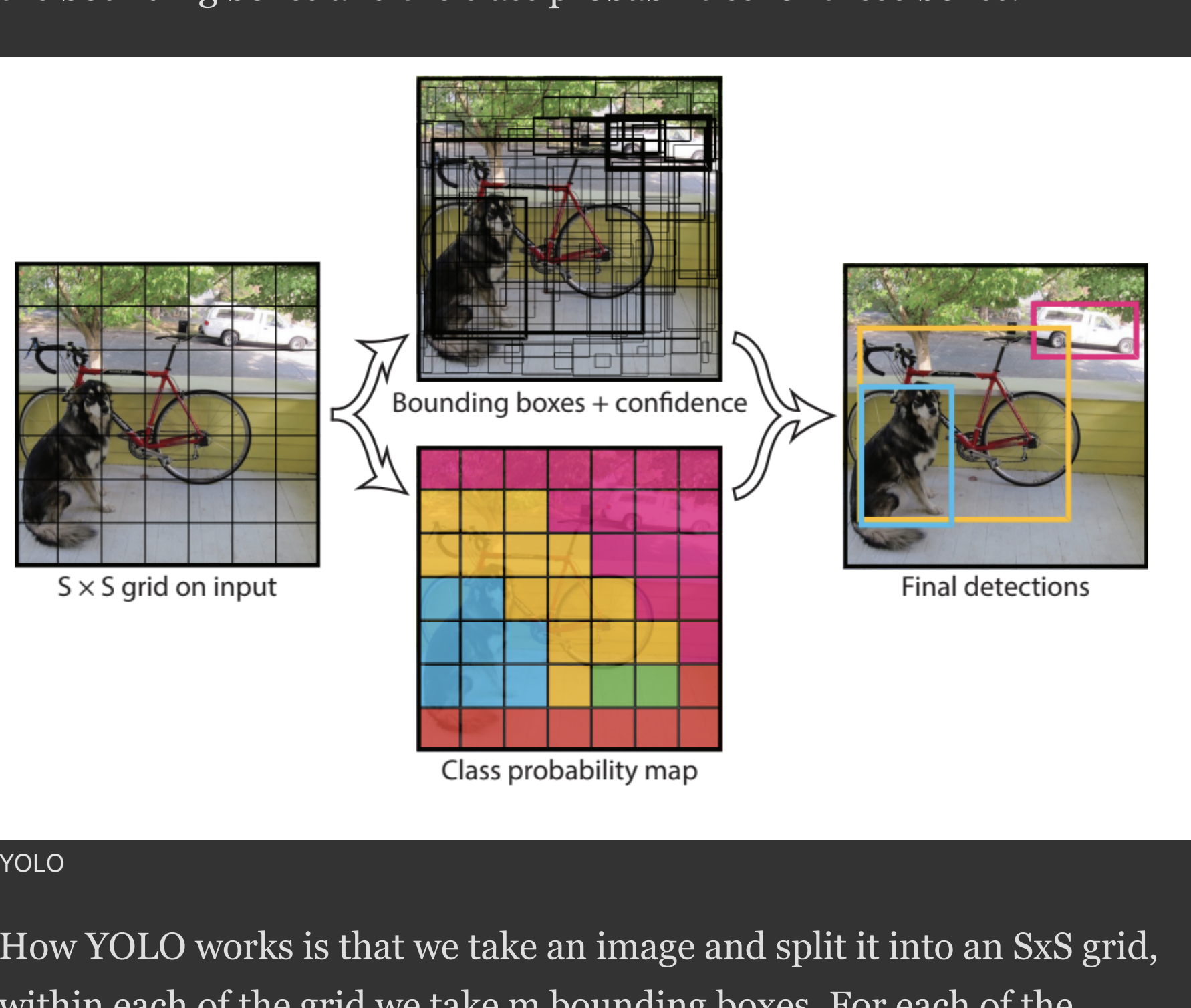


Comparison of test-time speed of object detection algorithms

From the above graph, you can see that Faster R-CNN is much faster than it’s predecessors. Therefore, it can even be used for real-time object detection.

YOLO—You Only Look Once

All of the previous object detection algorithms use



YOLO

How YOLO works is that we take an image and split it into an SxS grid, within each of the grid we take m bounding boxes. For each of the bounding box, the network outputs a class probability and offset values for the bounding box. The bounding boxes having the class probability above a threshold value is selected and used to locate the object within the image.

YOLO is orders of magnitude faster(45 frames per second) than other object detection algorithms. The limitation of YOLO algorithm is that it struggles with small objects within the image, for example it might have difficulties in detecting a flock of birds. This is due to the spatial constraints of the algorithm.

Conclusion

Computer vision conferences have been viewing new radical concepts each year and step by step I guess we are moving towards jaw-dropping performances from AI(if not already!). It only gets better. I hope the concepts were made lucid in this article, thank you :)

References

<https://arxiv.org/pdf/1311.2524.pdf>
<https://arxiv.org/pdf/1504.08083.pdf>
<https://arxiv.org/pdf/1506.01497.pdf>
<https://arxiv.org/pdf/1506.02640v5.pdf>
http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture11.pdf