

Machine Learning Engineer Nanodegree

Capstone

Mahavir Dwivedi

June 6th 2018

Domain Background

After learning a wide spectrum of ML algorithms and tools through the Udacity Machine Learning NanoDegree, I wanted to get my hands dirty on some problem statement where I can implement the concepts of Machine Learning majorly Deep Learning. As I was deciding on my problem statement for my capstone, I came across a lot of advancements going on in the field of Self-Driving Car. This inspired me to pick up a topic which is related to this theme.

There are a lot of aspects involved in making of a successful Self Driven Car, one of the most important aspects is that the car should be able to differentiate between different vehicles running over the road. A self-driven car should be able to recognize the vehicle around itself [2]. As the recent progress in the field of computer vision majorly object-detection, using these techniques I will try to build a deep learning model which can be used for vehicle detection .

As I don't have a Self-driving Car environment I will be using our vehicle detection model in another scenario. I will be setting up a smart Vehicle Detection system which will be installed at an entry/exit gate barrier, where it will be used to maintain the database and to monitor the different type of vehicles entering the premises.

Problem Statement :Following is the block diagram for the problem statement:

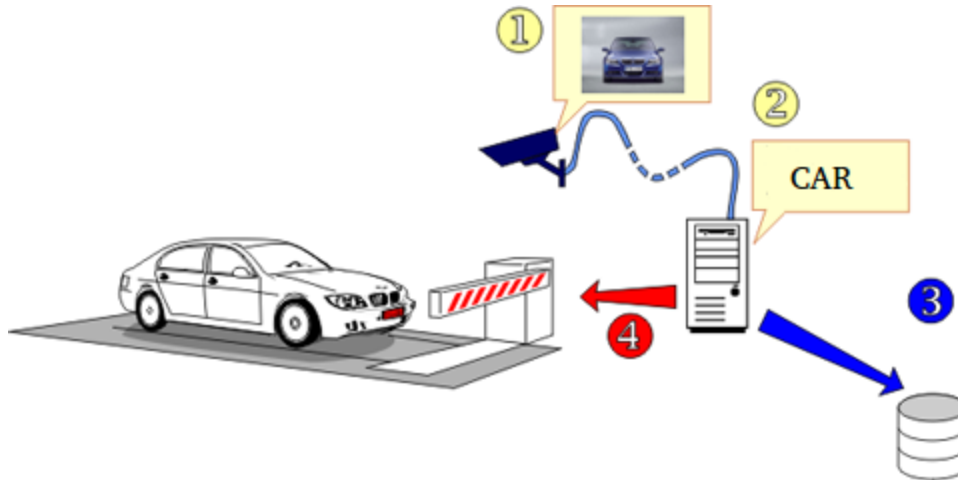


Fig-1 Block Diagram of problem statement

1. The Internet Protocol Camera at the Entry barrier takes the picture
2. The Picture is send to the CPU.
3. The Picture is being processed at the GPU
4. The decision is made to open the gate by the CPU.

Here I am proposing a deep learning model which can detect vehicles in a video stream.Following is the list of vehicles which I propose to detect:

1. Car
2. Bus
3. Motor Bike
4. Bicycle

The recognition will be performed on a video stream. I aim to achieve a processing speed where no vehicle is missed.

As we need to detect the vehicle in a given video frame, I am looking forward to use state of the art object detection models. One such model is Faster-RCNN which has a good processing speed, which is enough for our environment .Faster-RCNN has a good accuracy , as it is better as compared to YOLO for object detection over VOC dataset.[3]

Datasets Exploration :

As I am building a system for detecting vehicles, the training will be performed on following publicly available dataset:

VOC Dataset (<http://host.robots.ox.ac.uk/pascal/VOC/voc2007/>)

Dataset Attributes:

In order to evaluate the classification and detection challenges, the image annotation includes the following attributes for every object in the target set of object classes:—

- *Class*: aeroplane, bird, bicycle, boat, bottle, bus,car, cat, chair, cow, dining table, dog, horse, motorbike,person, potted plant, sheep, sofa, train, tv/monitor.
- *Bounding box*: an axis-aligned bounding box surrounding the extent of the object visible in the image.

Out of the 20 classes:

- *Person*: person
- *Animal*: bird, cat, cow, dog, horse, sheep
- *Vehicle*: aeroplane, bicycle, boat, bus, car, motorbike, train
- *Indoor*: bottle, chair, dining table, potted plant, sofa, tv/monitor

I will be using the model over following 4 vehicle classes: **Vehicle**: Bicycle, Bus, Car, Motorbike

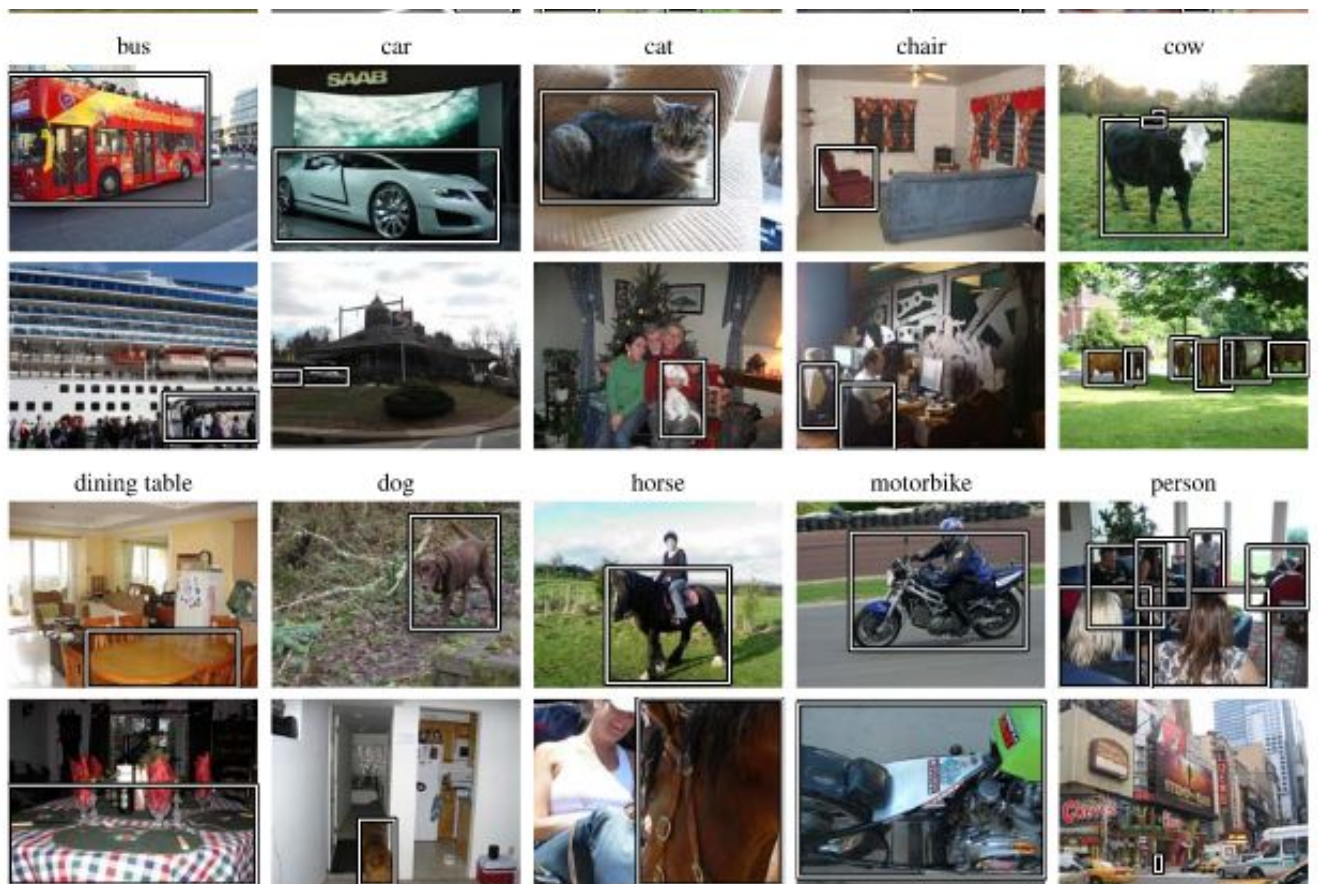


Fig- 2 Example images from the VOC 2007 dataset. These are few of the 20 classes annotated, for each of the classes above in the image two examples are shown. Note the wide range of pose, scale, clutter, occlusion and imaging conditions.

VOC datasets contain significant variability in terms of object size, orientation, pose, illumination, position and occlusion. It is also important to notice that the datasets do not exhibit systematic bias, for example, favouring images with centred objects or good illumination. Similarly, to ensure accurate training and evaluation, image annotations are consistent, accurate and exhaustive for the specified classes.

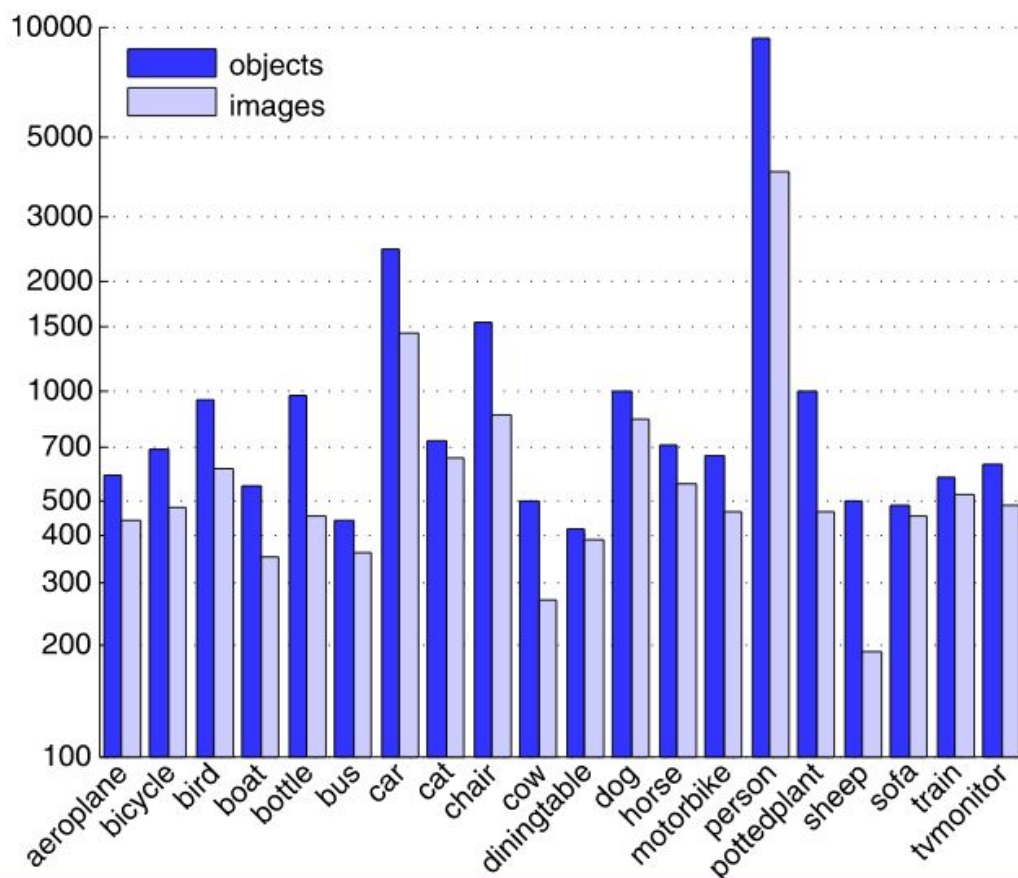


Fig-3 Summary of the entire VOC2007 dataset. Histogram by class of the number of objects and images containing at least one object of the corresponding class.

Evaluation Metrics :

The evaluation matrix I will be using here is mean average precision or “mAP score”. It has become the accepted way to evaluate object detection competitions, such as for the PASCAL VOC, ImageNet, and COCO challenges.

In object detection, evaluation is non trivial, because there are two distinct tasks to measure:

1. Determining whether an object exists in the image (classification)
2. Determining the location of the object (localization, a regression task).

Furthermore, in a typical data set there will be many classes and their distribution is non-uniform. So a simple accuracy-based metric will introduce biases. It is also important to assess the risk of misclassifications. Thus, there is the need to associate a “confidence score” or **model score** with each bounding box detected and to assess the model at various level of confidence.

In order to address these needs, the Average Precision (AP) was introduced. To understand the AP, it is necessary to understand the precision and recall of a classifier. Briefly, in this context, *precision measures the “false positive rate” or the ratio of true object detections to the total number of objects that the classifier predicted*. If you have a precision score of close to 1.0 then there is a high likelihood that whatever the classifier predicts as a positive detection is in fact a correct prediction. *Recall measures*

the “false negative rate” or the ratio of true object detections to the total number of objects in the data set. If you have a recall score close to 1.0 then almost all objects that are in your dataset will be positively detected by the model

To calculate the AP, for a specific class (say a “person”) the precision-recall curve is computed from the model’s detection output, by varying the model score threshold that determines what is counted as a model-predicted positive detection of the class.

The final step to calculating the AP score is to take the average value of the precision across all recall values . This becomes the single value summarizing the shape of the precision-recall curve. To do this unambiguously, the AP score is defined as the mean precision at the set of 11 equally spaced recall values, $Recall_i = [0, 0.1, 0.2, \dots, 1.0]$. Thus,

$$AP = \frac{1}{11} \sum_{Recall_i} Precision(Recall_i)$$

The precision at recall i is taken to be the maximum precision measured at a recall exceeding $Recall_i$.

Up until now, I have been discussing only the classification task. For the localization component (was the object’s **location** correctly predicted?) I must consider the amount of overlap between the part of the image segmented as true by the model vs. that part of the image where the object is actually located.

Localization and Intersection over Union

In order to evaluate the model on the task of object localization, we must first determine how well the model predicted the location of the object. Usually, this is done by drawing

a bounding box around the object of interest, but in some cases it is an N-sided polygon or even pixel by pixel segmentation. For all of these cases, the localization task is typically evaluated on the Intersection over Union threshold (IoU). Many good explanations of IoU exist but the basic idea is that *it summarizes how well the ground truth object overlaps the object boundary predicted by the model.*

Putting it all together :

Now that I have defined Average Precision (AP) and seen how the IoU threshold affects it, the mean Average Precision or mAP score is calculated by taking the mean AP over all classes and/or over all IoU thresholds, depending on the competition. For example:

- PASCAL VOC2007 challenge only 1 IoU threshold was considered: 0.5 so the mAP was averaged over all 20 object classes.
- For the COCO 2017 challenge, the mAP was averaged over all 80 object categories and all 10 IoU thresholds.

Averaging over the 10 IoU thresholds rather than only considering one generous threshold of $\text{IoU} \geq 0.5$ tends to reward models that are better at precise localization.

Model object detections are determined to be true or false depending upon the IoU threshold. This IoU threshold(s) for each competition vary, but in the COCO challenge, for example, 10 different IoU thresholds are considered, from 0.5 to 0.95 in steps of 0.05. For a specific object (say, 'person') this is what the precision-recall curves may look like when calculated at the different IoU thresholds of the COCO challenge:

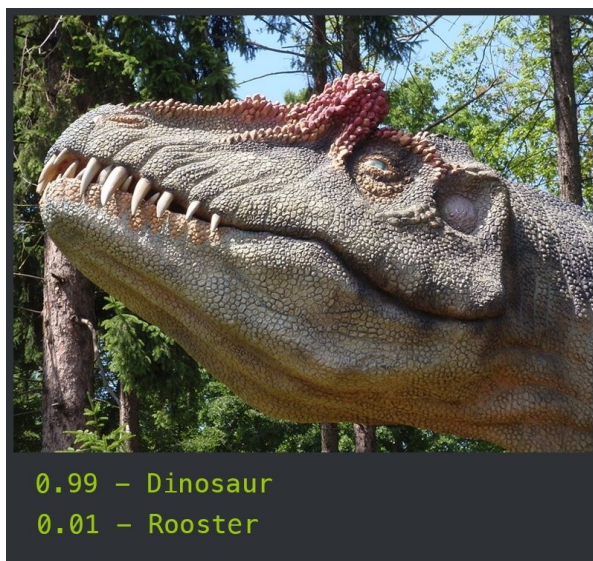
Algorithms And Techniques:

Given the problem statement of vehicle detection in a video stream, I went through the various Computer vision models to arrive at the solution.

As the task is to figure out all the vehicles present in a given frame at a given time, I didn't use a classification model, rather I built my pipeline based on an object detection model.

Classification:

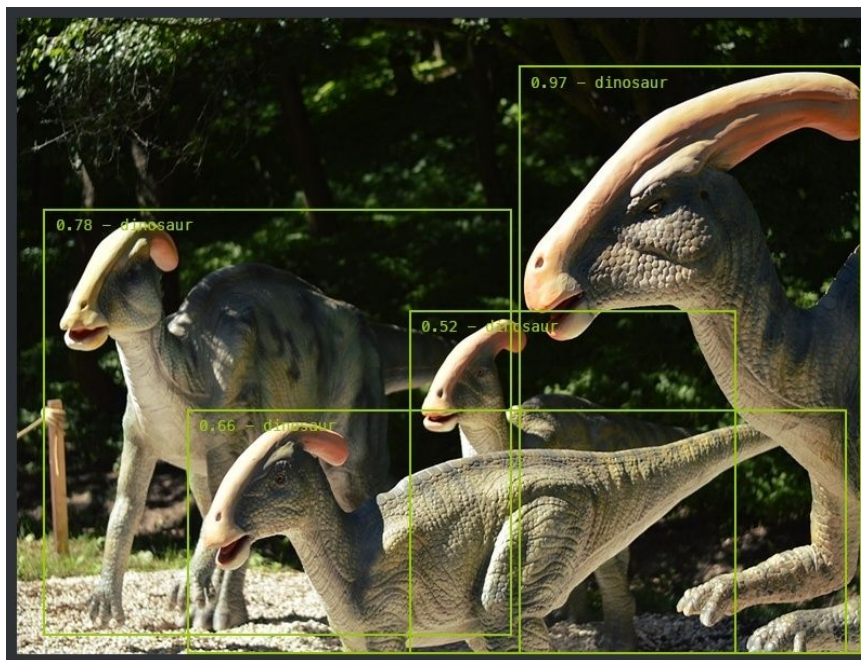
Probably the most well-known problem in computer vision. It consists of classifying an image into one of many different categories. One of the most popular datasets used in academia is ImageNet, composed of millions of classified images, (partially) utilized in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) annual competition. In recent years classification models have surpassed human performance and it has been considered practically solved. While there are plenty of challenges to image classification, there are also plenty of write-ups on how it's usually solved and which are the remaining challenges.



A classification model gives us the probability of all the classes being present in that picture. And then we output that the class with highest probability in that picture as our target class.

Object detection

Iterating over the problem of localization plus classification we end up with the need for detecting and classifying multiple objects at the same time. Object detection is the problem of finding and classifying a variable number of objects on an image. The important difference is the “variable” part. In contrast with problems like classification, the output of object detection is variable in length, since the number of objects detected may change from image to image. Whereas in an object detection model the model will give us all the objects with their respective classes and their positions respectively.



Problems and challenges with object detection :

Let's start getting deeper into which are the main issues of object detection.

- **Variable number of objects**

We already mentioned the part about a variable number of objects, but we omitted why it's a problem at all. When training machine learning models, you usually need to represent data into fixed-sized vectors. Since the number of objects in the image is not known beforehand, we would not know the correct number of outputs. Because of this, some post-processing is required, which adds complexity to the model. Historically, the variable number of outputs has been tackled using a sliding window based approach, generating the fixed-sized features of that window for all the different positions of it. After getting all predictions, some are discarded and some are merged to get the final result.

- **Sizing**

Another big challenge is the different conceivable sizes of objects. When doing simple classification, you expect and want to classify objects that cover most of the image. On the other hand, some of the objects you may want to find could be as small as a dozen pixels (or a small percentage of the original image). Traditionally this has been solved with using sliding windows of different sizes, which is simple but very inefficient.

- **Modeling**

A third challenge is solving two problems at the same time. How do we combine the two different types of requirements: location and classification into, ideally, a single model?

Before diving into deep learning and how to tackle these challenges, let's do a quick run-up of the classical methods:

Classical approach :

Although there have been many different types of methods throughout the years, we want to focus on the two most popular ones (which are still widely used).

The first one is the Viola-Jones framework proposed in 2001 by Paul Viola and Michael Jones in the paper Robust Real-time Object Detection. The approach is fast and relatively simple, so much that it's the algorithm implemented in point-and-shoot cameras which allows real-time face detection with little processing power. We won't go into details on how it works and how to train it, but at the high level, it works by generating different (possibly thousands) simple binary classifiers using Haar features. These classifiers are assessed with a multi-scale sliding window in cascade and dropped early in case of a negative classification.

Another traditional and similar method is using Histogram of Oriented Gradients (HOG) features and Support Vector Machine (SVM) for classification. It still requires a multi-scale sliding window, and even though it's superior to Viola-Jones, it's much slower

Deep learning approach

It's not news that deep learning has been a real game changer in machine learning, especially in computer vision. In a similar way that deep learning models have crushed other classical models on the task of image classification, deep learning models are now state of the art in object detection as well.

Following is an overview on how the deep learning approach has evolved in the last couple of years:

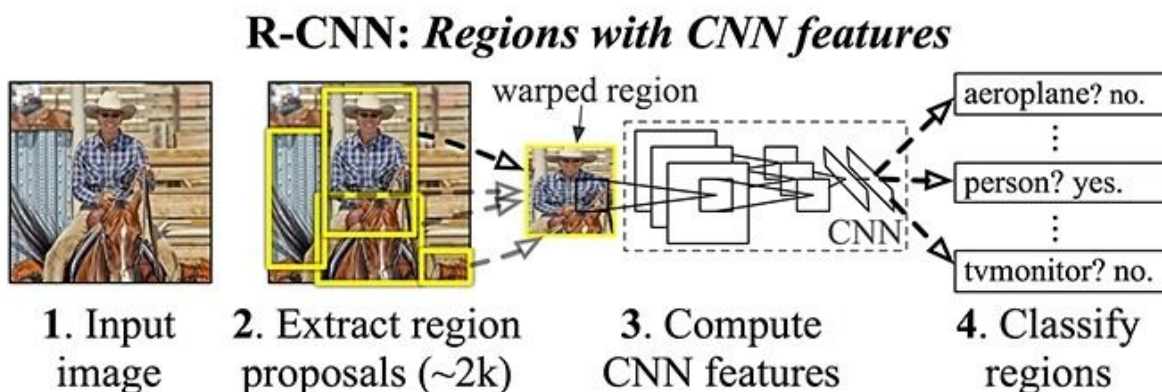
- **OverFeat**

One of the first advances in using deep learning for object detection was OverFeat from NYU published in 2013. They proposed a multi-scale sliding window algorithm using Convolutional Neural Networks (CNNs).

- **R-CNN**

Quickly after OverFeat, Regions with CNN features or R-CNN from Ross Girshick, et al. at the UC Berkeley was published which boasted an almost 50% improvement on the object detection challenge. What they proposed was a three stage approach:

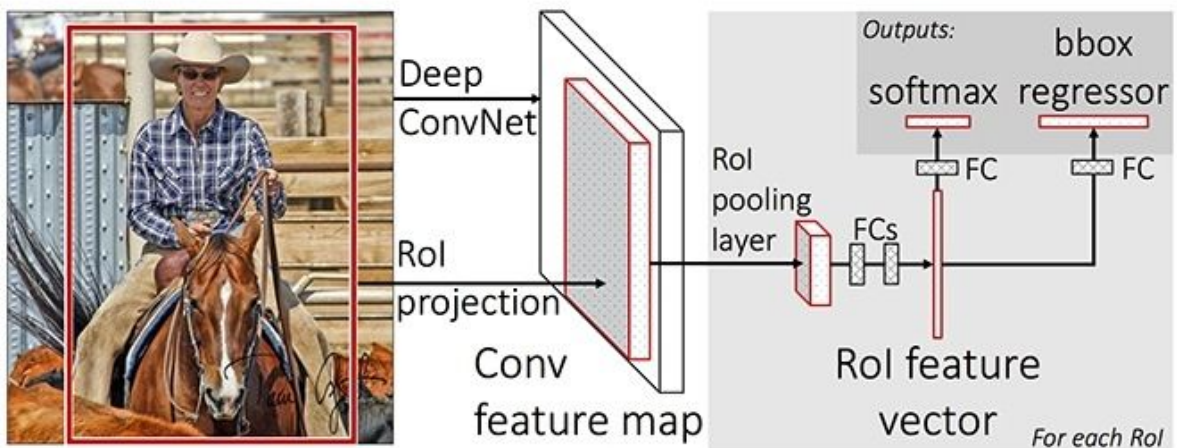
- Extract possible objects using a region proposal method (the most popular one being Selective Search).
- Extract features from each region using a CNN.
- Classify each region with SVMs.



While it achieved great results, the training had lots of problems. To train it you first had to generate proposals for the training dataset, apply the CNN feature extraction to every single one (which usually takes over 200GB for the Pascal 2012 train dataset) and then finally train the SVM classifiers.

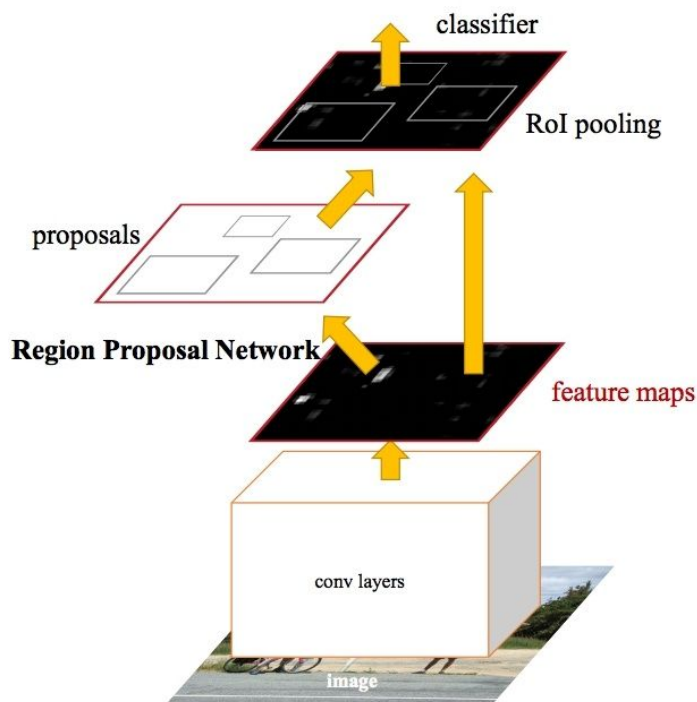
- **Fast R-CNN**

This approach quickly evolved into a purer deep learning one, when a year later Ross Girshick (now at Microsoft Research) published Fast R-CNN. Similar to R-CNN, it used Selective Search to generate object proposals, but instead of extracting all of them independently and using SVM classifiers, it applied the CNN on the complete image and then used both Region of Interest (RoI) Pooling on the feature map with a final feed forward network for classification and regression. Not only was this approach faster, but having the RoI Pooling layer and the fully connected layers allowed the model to be end-to-end differentiable and easier to train. The biggest downside was that the model still relied on Selective Search (or any other region proposal algorithm), which became the bottleneck when using it for inference.



- **Faster R-CNN**

Subsequently, Faster R-CNN authored by Shaoqing Ren (also co-authored by Girshick, now at Facebook Research), the third iteration of the R-CNN series. Faster R-CNN added what they called a Region Proposal Network (RPN), in an attempt to get rid of the Selective Search algorithm and make the model completely trainable end-to-end. We won't go into details on what the RPNs does, but in abstract it has the task to output objects based on an "objectness" score. These objects are used by the RoI Pooling and fully connected layers for classification. We will go into much more detail in a subsequent post where we will discuss the architecture in detail.



We used a LPR camera to capture the video stream, once the image is captured it is resized to a particular size. The preprocessed image is then passed to the object detection model where the vehicles of the above mentioned classes are detected. Once vehicles are detected, a bounding box is created on the display screen as well as the same is printed on the screen. Based on the model output further decisions can be taken by the computer.

BenchMark:

I have aimed to achieve couple of benchmarks :

1. The model which I train over VOC dataset should achieve a mAP of 0.699. The above mentioned mAP was achieved by a Faster RCNN model [1]. We aim to achieve a map beyond what was achieved in that paper.
2. As we are building a system for vehicle detection, we shouldn't miss any passing by vehicle. As we know at any of the boom barrier the vehicle stops, hence we can operate at low fps. Given the above specifications we aim to achieve a processing speed of 1 fps. We are setting our benchmark to 1 fps as we want our Vehicle Detection system to run on low end computing machines/portable devices. For achieving a inference speed of 1 fps we don't need a high end gpu machine, given that we implement the preprocessing steps correctly.

Methodology:

- **Data preprocessing:**

1. **Data Normalization:** We have removed the mean of the input color components and set their variance to 1. This has improved the results by 0.8%.
2. **Transform:** Given that, we are using pytorch library, we transform the image into a pytorch tensor.
3. As we are using a RPN based model, We need to preprocess an image for feature extraction.
 - The length of the shorter edge is scaled to a min size set by the user
 - After the scaling, if the length of the longer edge is longer than the set scaled size
 - the image is scaled to fit the longer edge to max image size
 - After resizing the image, the image is subtracted by a mean image value

Implementation :

Model Architecture Faster-rcnn: It all starts with an image, from which we want to obtain:

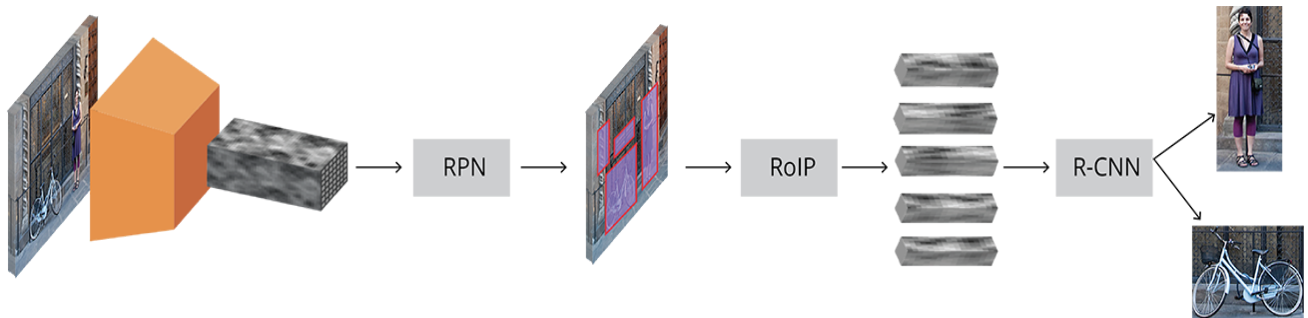


Fig-4 : Faster-RCNN Architecture

- a list of bounding boxes.
- a label assigned to each bounding box.
- a probability for each label and bounding box.

The input images are represented as Height×Width×Depth tensors (multidimensional arrays), which are passed through a pre-trained CNN up until an intermediate layer, ending up with a convolutional feature map. We use this as a feature extractor for the next part.

This technique is very commonly used in the context of Transfer Learning, especially for training a classifier on a small dataset using the weights of a network trained on a bigger dataset.

Next, we have what is called a Region Proposal Network (RPN, for short). Using the features that the CNN computed, it is used to find up to a predefined number of regions (bounding boxes), which may contain objects.

Probably the hardest issue with using Deep Learning (DL) for object detection is generating a variable-length list of bounding boxes. When modeling deep neural networks, the last block is usually a fixed sized tensor output (except when using Recurrent Neural Networks). For example, in image classification, the output is a $(N,)$ shaped tensor, with N being the number of classes, where each scalar in location i contains the probability of that image being labeled i .

The variable-length problem is solved in the RPN by using anchors: fixed sized reference bounding boxes which are placed uniformly throughout the original image.

Instead of having to detect where objects are, we model the problem into two parts. For every anchor, we ask:

Does this anchor contain a relevant object?

After having a list of possible relevant objects and their locations in the original image, it becomes a more straightforward problem to solve. Using the features extracted by the CNN and the bounding boxes with relevant objects, we apply Region of Interest (RoI) Pooling and extract those features which would correspond to the relevant objects into a new tensor.

Finally, comes the R-CNN module, which uses that information to:

- Classify the content in the bounding box (or discard it, using “background” as a label).
- Adjust the bounding box coordinates (so it better fits the object).

Obviously, some major bits of information are missing, but that’s basically the general idea of how Faster R-CNN works.

Base network :

As we mentioned earlier, the first step is using a CNN pre trained for the task of classification (e.g. using ImageNet) and using the output of an intermediate layer. This may sound really simple for people with a deep learning background, but it’s important

to understand how and why it works, as well as visualize what the intermediate layer output looks like.

There is no real consensus on which network architecture is best. The original Faster R-CNN used ZF and VGG pretrained on ImageNet but since then there have been lots of different networks with a varying number of weights. For example, MobileNet, a smaller and efficient network architecture optimized for speed, has approximately 3.3M parameters, while ResNet-152 (yes, 152 layers), once the state of the art in the ImageNet classification competition, has around 60M. Most recently, new architectures like DenseNet are both improving results while lowering the number of parameters.

VGG :

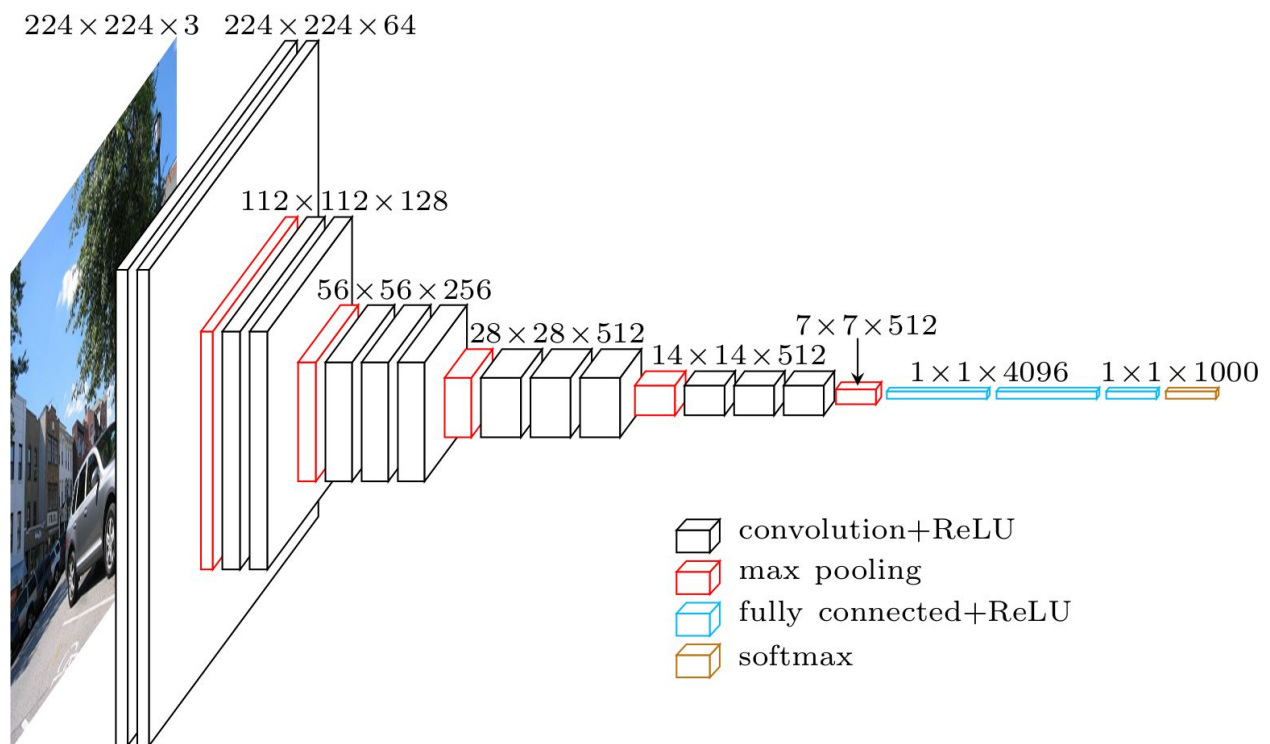


Fig - 6 VGG is being used as the base network

VGG, whose name comes from the team which used it in the ImageNet ILSVRC 2014 competition, was published in the paper “Very Deep Convolutional Networks for Large-Scale Image Recognition” by Karen Simonyan and Andrew Zisserman. By today’s standards it would not be considered very deep, but at the time it more than doubled the number of layers commonly used and kickstarted the “deeper → more capacity → better” wave (when training is possible).

When using VGG for classification, the input is a $224 \times 224 \times 3$ tensor (that means a 224×224 pixel RGB image). This has to remain fixed for classification because the final block of the network uses fully-connected (FC) layers (instead of convolutional), which require a fixed length input. This is usually done by flattening the output of the last convolutional layer, getting a rank 1 tensor, before using the FC layers.

Since we are going to use the output of an intermediate convolutional layer, the size of the input is not our problem. At least, it is not the problem of this module since only convolutional layers are used. Let’s get a bit more into low-level details and define which convolutional layer we are going to use. The paper does not specify which layer to use; but in the official implementation you can see they use the output of conv5/conv5_1 layer.

Each convolutional layer creates abstractions based on the previous information. The first layers usually learn edges, the second finds patterns in edges in order to activate for more complex shapes and so forth. Eventually we end up with a convolutional feature map which has spatial dimensions much smaller than the original image, but greater depth. The width and height of the feature map decrease because of the pooling

applied between convolutional layers and the depth increases based on the number of filters the convolutional layer learns.

In its depth, the convolutional feature map has encoded all the information for the image while maintaining the location of the “things” it has encoded relative to the original image. For example, if there was a red square on the top left of the image and the convolutional layers activate for it, then the information for that red square would still be on the top left of the convolutional feature map.

Challenges Faced:

1. Deciding the size of preprocessed image: We decided the size of our preprocessed image based on the size of the Kernels used in our object detection model(Faster -RCNN).Number of Kernels and their size was decided based on the VOC dataset attributes.
2. As we are using both Opencv and Matplotlib for our processing and display while handling the images. We had to transform the images according to the library we were using.

We used the following function for transforming images so that the given image can be processed by our pytorch tensor: **`image_np.transpose((2, 0, 1))`**

Process on improving upon the algorithms:

1. **Using High End GPU machine**

When we were running our inference model on a 940 MX based GPU machine , my inference speed was 1fps.

As we tried the same inference model on Teslas K80 GPU, the inference speed of the model increased to 6fps.

2. Setting up the confidence Threshold(NMS preset threshold)

We set up the nms confidence threshold initially to 0.7, as a result we were not getting that great results on pictures taken outside the test set. Once we set the threshold value to 0.5 we saw a drastic improvement in our inferences.

Train test graphs :

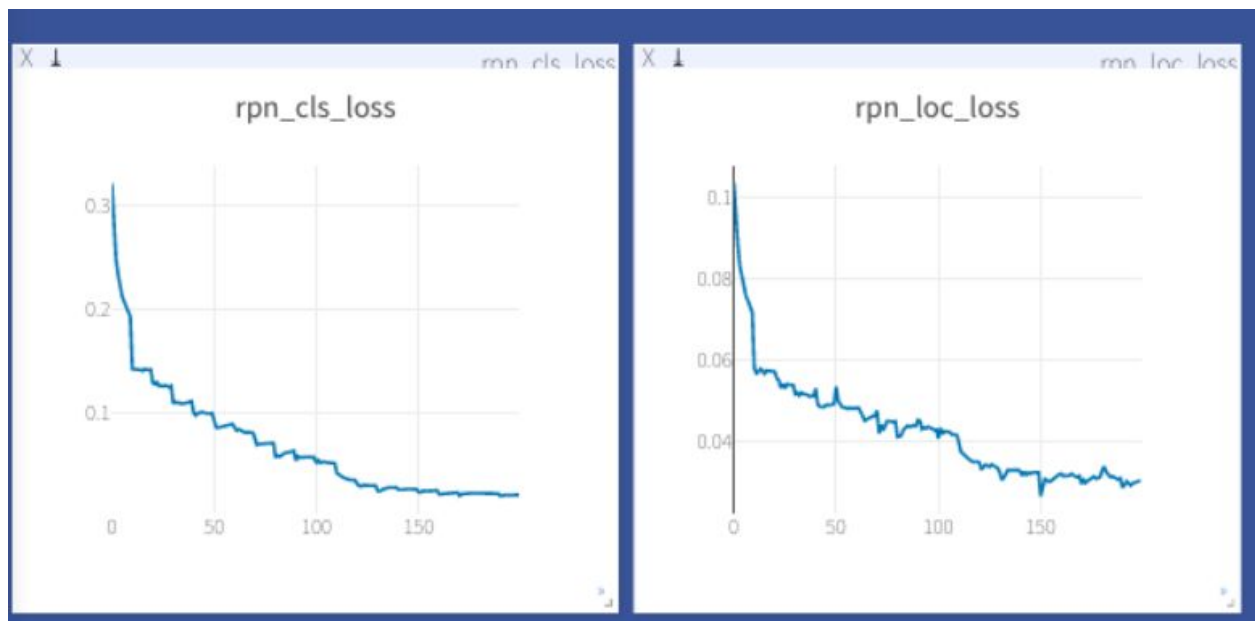


Fig 7: Rpn Losses plotted in above figure.

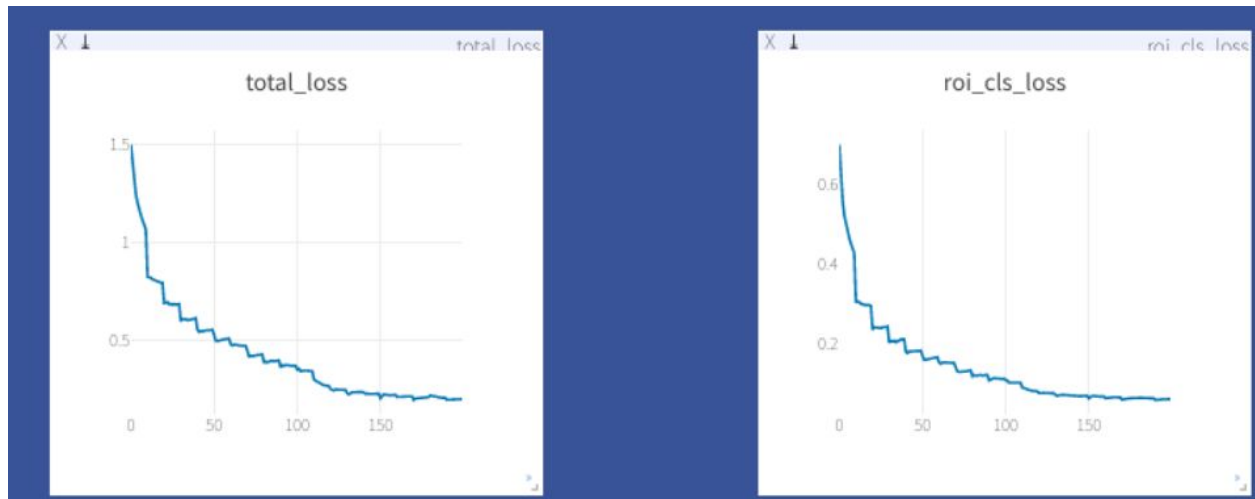


Fig-8 Total Loss being plotted using visdom server

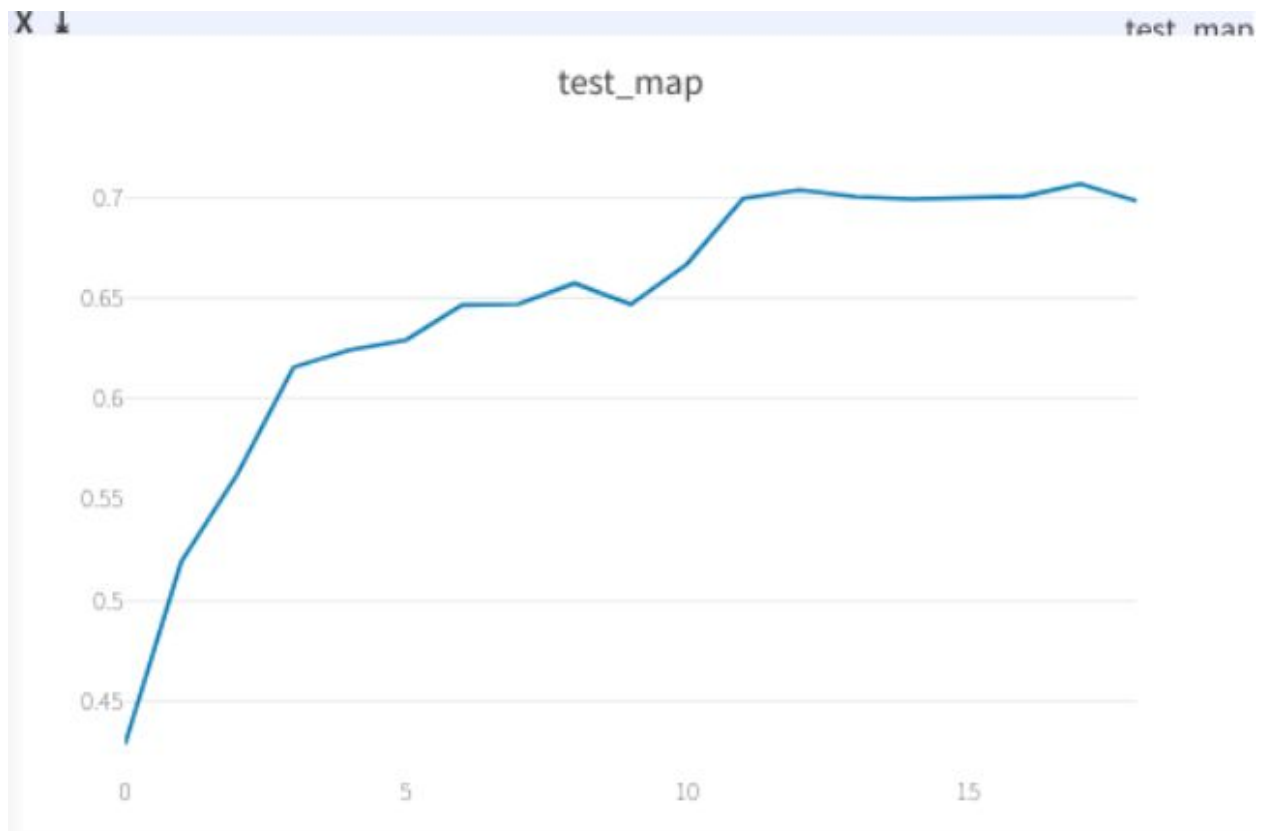


Fig -9 : plot for test_map

Train test speed :

While training and testing our Model on Tesla K80 , following was the train and test speed:

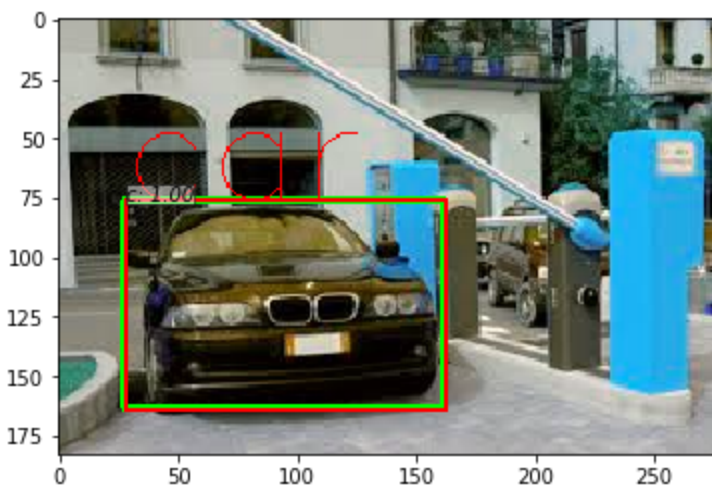
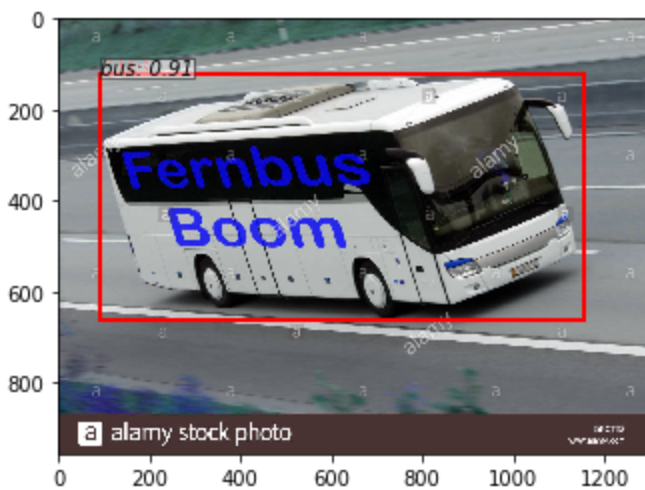
- Train Speed: 4 it/s
- Test Speed: 7 it/s

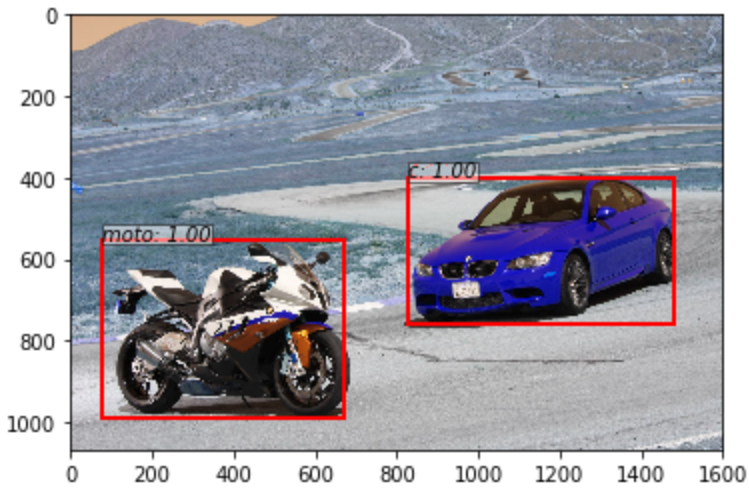
Results:

In the below figures you can see :

1. A bus has been detected
2. A car has been detected
3. A car and a motorbike is detected
4. A car has been detected

Note : In our results c stands for car and motto for motorbike.





The robustness of the model's solution:

- We have trained the model over VOC dataset , as you can see above the model performs well over random images taken from the internet.
- As the model has been trained majorly for four kind of vehicles, its not able to detect other kind of vehicles.

Model evaluation and validation :

	Implementation	mAP
1	origin paper [1]	0.699
2	Our vehicle detection Model	0.7053

Implementation	GPU	Inference	Training
origin paper [1]	K40	5 fps	NA
vehicle-detection-model	TITAN Xp	15-17fps	6fps

Justification:

1. On a system which is meant for Vehicle Detection we couldn't have used a model which predicts other classes too along with vehicles. Hence we came up with a model which is meant just for vehicle detection.

2. On a low end GPU system such as (940 MX) we are able to process videos at a rate of 1 fps.
3. If we use higher end GPU such as Titan Xp we can process videos as 6fps.

Conclusion:

Through this capstone project we came up with a problem statement to detect vehicles. This resulted in a model which can detect vehicles with a state of art accuracy and a processing time where it doesn't miss any vehicle. As we already set up an environment(boom barriers) where the speed of incoming vehicles was slow. To implement the same on a barrier less traffic we will have to introduce high end GPUs. We also need to train our model for more class of vehicles.

- *What was the most interesting aspect of this work ?*

Most of the advancements on deep learning is based on High end GPUs .The most interesting part of this project was to run our model on lower end GPUs such as 940 MX.

- *What was the biggest challenge, and how did you overcome it ?*

As I mentioned in the above question ,the biggest and the most interesting challenge was to get this model working on lower end GPUs. We aimed to achieve a inference speed where we shouldn't miss a vehicle in a live video feed. We removed the redundant frames to achieve a inference speed of 1 fps at a Nvidia 940MX Laptop.

Potential solutions:

As our problem statement is to detect vehicles given in a video stream, we aim to include more class of vehicles.

For this we plan to train the same model over other datasets such as COCO and Udacity self driving car dataset.

We also plan to train it over custom datasets, where we capture images ourselves , annotate them, preprocess them and then train the model. This will make the model more robust to real world images.

References :

[1] Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, <https://arxiv.org/abs/1506.01497>

[2] A closer look at Faster R-CNN for vehicle detection
<https://ieeexplore.ieee.org/document/7535375/>

[3]https://medium.com/@jonathan_hui/object-detection-speed-and-accuracy-comparison-faster-r-cnn-r-fcn-ssd-and-yolo-5425656ae359