# Object Detection and Orientation Estimation for Autonomous Driving

**Jinyi Lu**
Information Networking Institute
jinyil@andrew.cmu.edu

**Xiaoqing Tao**
Information Networking Institute
xtao@andrew.cmu.edu

**Keywords:** [computer vision, object detection, object orientation estimation, autonomous driving]

## 1 Overview

### 1.1 Introduction

Nowadays cameras are available onboard of of almost every new car produced in the last few years. Computer vision provides a very cost effective solution not only to improve safety, but also to one of the holy grails of AI, fully autonomous self-driving cars. In this project we are planning to use deep neural networks to solve the object detection and object orientation estimation problems for autonomous driving.

There are lots of potential challenges that we need to solve for example, due to the weather, road conditons and car location, images from the car cameras will have a high-variety, which requires high robustness for our model. And besides the classical object detection task, we want to further estimate the 3D orientation from the 2D images. Last, but not least our system need to produce a good result within a limited runtime in order to be used in practise.

### 1.2 Dataset

The dataset that we are planning to use is the KITTI Vision Benchmark Suite [1]. It's developed for use in mobile robotics and autonomous driving research. So it contains several novel challenging benchmarks for the tasks of stereo, optical flow, visual odometry/SLAM and 3D object detection. In our project, we mainly focus on the object detection and orientation estimation task. The corresponding benchmark [1] consists of 7481 training images and 7518 test images, comprising a total of 80,256 labeled objects (up to 15 cars and 30 pedestrians are visible per image). All images are color and saved as png.

Dataset statistics

### 1.3 Evaluation

For evaluation, the benchmark is split into three parts: First, we need to evaluate the classical 2D object detection by measuring performance using the well established average precision (AP) metric as described in [2]. Detections are iteratively assigned to ground truth labels starting with the largest overlap, measured by bounding box intersection over union. True positives are required to overlap by more than 50% and multiple detections of the same object are counted as false positives.

---

[1] http://www.cvlibs.net/datasets/kitti/eval_object.php

Second, we assess the performance of jointly detecting objects and estimating their 3D orientation using a novel measure which is called the average orientation similarity (AOS) [1] and is defined as:

$$AOC = \frac{1}{11} \sum_{r \in \{0,0.1,..,1\}} \max_{\tilde{r}:\tilde{r} \geq r} s(\tilde{r}) \qquad (1)$$

Here, $r = \frac{TP}{TP+FN}$ is the PASCAL object detection recall, where detected 2D bounding boxes are correct if they overlap by at least 50% with a ground truth bounding box. The orientation similarity $s \in [0,1]$ at recall r is a normalized ($[0..1]$) variant of the cosine similarity defined as

$$s(r) = \frac{1}{|D(r)|} \sum_{i \in D(r)} \frac{1 + cos\Delta_{theta}^{(i)}}{2} \delta_i \qquad (2)$$

where $D(r)$ denotes the set of all object detections at recall rate $r$ and $\Delta_{theta}^{(i)}$ is the difference in angle between estimated and ground truth orientation of detection $i$. To penalize multiple detections which explain a single object, we set $\delta_i = 1$ if detection $i$ has been assigned to a ground truth bounding box (overlaps by at least 50%) and $\delta_i = 0$ if it has not been assigned.

Finally, we will also evaluate pure classification (16 bins for cars) and regression (continuous orientation) performance on the task of 3D object orientation estimation in terms of orientation similarity.

## 2 Related work

Traditional methods for object detection usually utilize image features,such as SIFT and HOG. We investigated three methods utilizing deep convolutional network in object detection.

### 2.1 R-CNN: Rich feature hierarchies for accurate object detection and semantic segmentation

The first challenge for object detection is how to implement localizing within an image. The new method proposed in this paper[3] combines CNN with region proposals. So it is called R-CNN. The general detection process is to extract about 2000 region proposals for each input image. Then utilize CNN to compute a fixed length feature for each region proposal. Finally utilize linear SVM to classify each region.

The second challenge is how to train a large CNN using limited labeled data. This paper proposed to pre-train CNN on a large auxiliary data set, then continually train on a small data set. This method significantly improved the accuracy of objection detection compared with feature learning models. But training is expensive in space and time, and detection is also slow at test time.

### 2.2 Fast R-CNN

This paper[4] talked about how to train a detection network faster. R-CNN is very slow because it extracts feature for each region proposal and there are many duplicate computations. The process could be faster if we share computation.

This paper proposed to modify R-CNN's architecture by taking an image and multiple regions of interests as input. Region proposal method usually depends on Selective Search. Each region of interest is pooled into a fixed-size feature map and fully connected layers are used to extract features. There two output vectors: softmax probabilities and per-class bounding-box regression offsets.The second one is to reduce mislocalization.

### 2.3 Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks

Fast-CNN achieves near real-time rates using very deep networks, when ignoring the time spent on region proposals. Select Search takes about 2 seconds per image on the CPU. Faster R-CNN[5] adds a fully convolutional network to compute region proposals directly. Then use the same detector on the proposed regions as what Fast R-CNN does.

# 3 Experiment Results

In this section, we will discuss some of the primary experiment results that we have achieved.

Firstly, we directly use the provided pretrained yolo model to make prediction our development set. This yolo model is trained using VOC2012 dataset so it contains about 20 categories including car, people, bicycle, etc. In order to solve our problem, we consider the people it detected as the pedestrian category that we want and detected bicycle as the cyclist that we want.

# 4 Project Plan

We plan to implement 2D object detection using faster R-CNN by the second milestone. After that, we will focus on 3D object detection. We may also try YOLO[6] method and see if there are any improvement. Detailed plan depends on time and compute power.

# References

[1] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.

[2] M Everingham, L Van Gool, CKI Williams, J Winn, and A Zisserman. The pascal visual object classes challenge 2011 (voc 2011) results (2011). In *URL http://www. pascal-network. org/challenges/VOC/voc2011/workshop/index. html*, 2010.

[3] Trevor Darrell Jitendra Malik Ross Girshick, Jeff Donahue. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Conference on Computer Vision and Pattern Recognition(CVPR)*, 2014.

[4] Ross Girshick. Fast r-cnn. In *IEEE International Conference on Computer Vision (ICCV)*, 2015.

[5] Ross Girshick Jian Su Shaoqing Ren, Kaiming He. Faster r-cnn: Towards real-time object detection with region proposal networks. In *IEEE International Conference on Computer Vision (ICCV)*, 2015.

[6] Ross Girshick Ali Farhadi Joseph Redmon, Santosh Divvala. You only look once: Unified, real-time object detection. In *arXiv preprint arXiv:1506.02640*, 2015.