
Car Driving Without Cameras

By

JONES MABEA AGWATA



University of
BRISTOL

Department of Computer Science

A dissertation submitted to the University of Bristol in accordance with the requirements of the degree of MASTER OF SCIENCE in the Faculty of Engineering.

SEPTEMBER 2018 | CSMSC-18



0000055012

EXECUTIVE SUMMARY

The need for robust object detection in 3D point clouds has greatly increased with the ongoing push for autonomous vehicles (AVs). Most of these systems use Light Detection and Ranging (LiDAR), cameras or a combination of both in order to perform object detection. LiDAR presents objects as point clouds in a 3D space thus offering critical shape information of objects in view. However, this representation is sparse. As a result, LiDAR-based detection performs poorly as compared to multimodal methods that have helped overcome this. Nonetheless, multimodal methods are often complex to set up and synchronise with the cost of components running into thousands of pounds. In light of this, reducing the number of sensors while still maintaining or even improving the accuracy has become a topic of interest by industry players who are developing AVs. With the cost of LiDAR declining following recent improvements in solid-state LiDAR technology, dropping cameras in favour of LiDAR has become more plausible. As compared to cameras that suffer drawbacks such as visibility issues in extreme weather, LiDAR is extremely robust and accurate in various conditions.

This project will be 65% type I (Software Development) and 35% type II (Theoretical).

DEDICATION AND ACKNOWLEDGEMENTS

First and foremost, I'd like to thank my parents for their support throughout my studies. I'm extremely grateful to my supervisors, Dr. Luis Gonzalez and Dr. Raul Santos-Rodriguez who have offered indispensable insight that has been truly beneficial to completing this project. In addition, I'd like to thank Dr. Sion Hannuna for providing me with necessary resources.

TABLE OF CONTENTS

	Page
List of Tables	vii
List of Figures	ix
1 Introduction	1
1.1 Aims and Objectives	2
1.2 Deliverables	2
1.3 Report structure	3
2 Background	5
2.1 Components of an AV	6
2.2 Industrial Approaches	9
2.2.1 Hardware Considerations	10
2.2.2 Camera	11
2.2.3 Camera and Radar	12
2.2.4 LiDAR, Camera and Radar	12
2.2.5 LiDAR	12
2.3 Object Detection	12
2.3.1 Object Detection using Classic Computer Vision Methods	12
2.3.2 Object Detection using Deep Learning	13
2.4 Conclusion	16
3 Implementation	17
3.1 Datasets	17
3.1.1 Context Classification	17
3.1.2 LiDAR Performance Evaluation	18
3.2 Can We Automatically Detect The Context in Driving Scenes?	18
3.2.1 Image Context Detection	19
3.2.2 PointCloud Context Detection	20
3.2.3 Implementation Issues	21

TABLE OF CONTENTS

3.3	Do Some Object Detection Methods work better in Different Contexts?	22
3.3.1	VoxelNet	22
3.3.2	AVOD	24
3.3.3	Evaluation	25
3.3.4	Is VoxelNet Performance Valid for Different Datasets?	25
3.4	Conclusion	27
4	Results and Analysis	29
4.1	Can We Automatically Detect Scene Contexts?	29
4.1.1	Metrics	29
4.1.2	Analysis	29
4.2	Do Object Detection Models Perform Better in Different Contexts?	30
4.3	Is VoxelNet Performance Valid For Different Datasets?	34
5	Discussion and Future Work	37
5.1	Future Work	37
6	Conclusion	39
A	Appendix A	41
	Bibliography	43

LIST OF TABLES

TABLE	Page
21 Modes of Operation of Common Architectures	16
31 Velodyne HDL-64E and VLP-16 Specifications	27
41 Comparison of context detection methods	30
42 Penalised Point Cloud classifier	30
41 Baseline	31
41 Car AP	32
42 Pedestrian AP	32
43 Inference Time on Single NVIDIA P100 GPU	32
44 Memory	33
A1 Velodyne LiDAR Family	41

LIST OF FIGURES

FIGURE	Page
21 ALVINN architectures. Left: With camera and laser range finder. Right: camera only. Source: [22]	6
21 LiDAR Point Cloud	7
22 Velodyne LiDAR family. From left: Velodyne HDL-64E , HDL-32E , VLP-16	8
23 Components of Waymo’s Self Driving Car, Source:[36]	9
21 AV Leaderboard. Navigant Research. Source:[28]	11
21 NN vs CNN Source: https://cs231n.github.io/convolutional-networks	13
22 Faster R-CNN overview. Source:[27]	15
31 Urban scene	19
32 Non-Urban scene	19
31 Semantic Histograms	20
31 Ambiguous scene	22
31 VoxelNet Architecture	25
32 AVOD Architecture	25
31 SIL Data Frame CSV	26
32 Veloview Frame of SIL data	26
33 L-CAS Annotation Tool	28
41 Temperature boxplots of evaluation runs.Voxelnet(blue) AVOD(gold)	32
42 Power boxplots of evaluation runs. Voxelnet(blue) AVOD(gold). Last boxplot is of both models.	32
41 *	34
42 *	34
43 Bird Eye View Difference between SIL and KITTI	34
44 Point Cloud Density	35

INTRODUCTION

Accelerated by recent advancements in technology, the prospect of Autonomous Vehicles (AVs) driving in public roads is becoming more and more a reality. As this is an emerging field, there are numerous variations of implementations by different companies. Arguably, a key characteristic of these implementations is a large number of perception sensors including cameras, radars and Light detection ranging sensors(LiDAR) that are necessary for mapping the environment around the vehicle in order to safely navigate. Most of these sensors are quite expensive and fusing their input and processing it requires powerful processing units such as GPUs. This process consumes a lot of power and generates a lot of heat. In an effort to reduce the cost of such systems, companies are exploring different ways to reduce the number of sensors while still achieving a high level of navigational accuracy and safety.

Compared to humans who are able to easily adapt to different situations, AVs need to be trained and tested in different possible situations. Different regions in the world pose varying challenges, such as poor visibility in Scandinavian countries due to heavy snow, and on the other hand heavily populated regions with heavy traffic such as New York pose challenges such as high occlusion of objects. An effective way of testing and evaluating AV systems is using a scenario-based approach. This has been the basis of Software-in-Loop and Hardware-in-Loop tests that test these systems in a virtual real-time environment to understand how they would react. This is done by modelling conditions and evaluating them on the software and hardware. []. As a stepping stone, the scenario-based approach can be used to develop AV systems that are tailor-made for certain environments/regions and later on generalised as the systems become more adaptable to different environments. There are some companies and institutions that adopt this approach and have deployed autonomous prototypes in different areas. For example, Uber currently has some AVs in Steel City, Pittsburgh whereas Waymo has deployed some in Phoenix, Arizona[]. Other companies are still developing their AVs within test tracks[]. Using

this approach such companies are able to collect a large amount of data for testing and evaluating future models. However, this data is not always released to the public therefore limiting the variation and amount of data that is available to other parties working on AVs.

Currently, a lot of cutting-edge research on AV systems relies on a small subset of data. Most research is based on the KITTI [10] dataset that was collected to aid the development of AV systems by providing the necessary data to develop models to achieve different task. Such tasks included, stereo, optical, depth, odometry, sceneflow and object detection. As a result it is highly favoured in the development of novel computer vision algorithm. In equal light, relying on few datasets creates a bottleneck in terms of validation as these algorithms may end up performing sub par when evaluated on different datasets or even worse in a real-time setting.

Following this argument, a viable scenario that could be beneficial to explore is how different sensors and corresponding object detection models perform in urban and non-urban contexts. In doing so, companies developing AVs can fit different sensors in AVs that work in different contexts with the aim of maximising their performance in that context while reducing the cost by removing redundant sensors. The most reliable and accurate combination of sensors in terms of and reliability is camera and LiDAR. However, there has been increased interest in developing LiDAR only systems that are able to exploit highly accurate depth information from LiDAR such as Doxel[1] that is an autonomous robot that uses LiDAR to monitor construction sites. As such it may be possible to develop LiDAR only units that can be used in a specific context.

1.1 Aims and Objectives

Following the motivations in the presented discussion , the performance of LiDAR only and multimodal(LiDAR and Camera) models in different contexts will be investigated with the aim of reducing the number of sensors in AVs.

To achieve this aim, the following objectives will need to be fulfilled:

1. **Detect and characterise the context of images and point clouds.**
2. **Evaluate the performance of single sensor and multimodal models in different contexts.**
3. **Validate performance of the single sensor model on a custom dataset. Validate performance of the single sensor model on a custom dataset.**

1.2 Deliverables

The deliverables are categorized in order of:

1. **Context Detection**

- **Image and LiDAR Context Classifier** - Available as Jupyter interactive notebooks including pre-trained models.

2. Multimodal and single sensor evaluation

- **Custom VoxelNet Model** - Modified VoxelNet model including interactive notebooks for training, testing and validating the model.
- **Custom AVOD Model** - Modified AVOD model including interactive notebooks for training, testing and validating the model.

3. Single sensor validation

- Point Cloud dataset obtained from the University of Bristol Smart Internet Lab working on connected and AVs. Tools to convert the dataset into compatible input for VoxelNet will be provided as well as some annotated sample frames for training and evaluation.

4. Report containing the following discussions:

- A review of related research and implementations tackling object detection in AVs.
- Evaluation and analysis of context classifiers and object detection models in different contexts.
- Validation results of single sensor model using external dataset from the Smart Internet Lab on connected and autonomous vehicles.

1.3 Report structure

This report will consist of five main chapters.

- Chapter 2 discusses the different components of AVs, current implementations in the industry, a background on the research that has been undertaken in the field of object detection.
- Chapter 3 details the project execution and the methods undertaken to achieve the objectives.
- In Chapter 4, the results following evaluation of the methods will be discussed and analysed.
- Finally, a concluding chapter discusses the major findings, whether the objectives were achieved and a justification of their implications.

BACKGROUND

Despite the recent interest in autonomous vehicles(AVs), the idea of autonomous vehicles(AVs) is not new. In fact, there had been attempts to autonomously navigate using cameras as early as 1990. Dan Pomerleau, a leading researcher from Carnegie Mellon, had developed a neural network, ALVINN(Autonomous Land Vehicle in Neural Network)[23], that could steer a van using camera input and avoid obstacles using a laser range finder . His system used an artificial neural networks that was able to categorise road and non-road segments and detect objects[22]. Despite the success, a major limitation was the low computational power of computers at that time that led to long training times for artificial neural networks. Due to this limitation, interest in AVs slowly waned off.

It was not until around 2004 when the AV landscape started seeing major development. At around this period, the US congress had passed regulation to replace 33% of military vehicles with AVs. Despite major military investment, there was not much progress and therefore Defense Advanced Research Projects Agency(DARPA) opted to revolutionise the process by making the project a public competition. In February 2003, DARPA hosted the Grand Challenge with a prize of 1\$ million dollars to develop an AV that could finish a 142 mile course that was set up in Mojave Desert to simulate harsh difficult conditions. Of 160 competitors, there was no winner. In the 2005 Grand Challenge, Stanford University had the fastest car to complete the course and consequently won the challenge. Subsequent challenges were held and more and more cars were able to complete the course aided by advances in AI methods. Following these challenges, more research was undertaken to create reliable and accurate sensors. In fact, Velodyne came about as a result of the Grand Challenge and currently is the leading supplier of LiDAR units.

With different types of autonomous systems being deployed in vehicles nowadays, the Society of Automotive Engineers(SAE) defined five levels of autonomy to differentiate them:

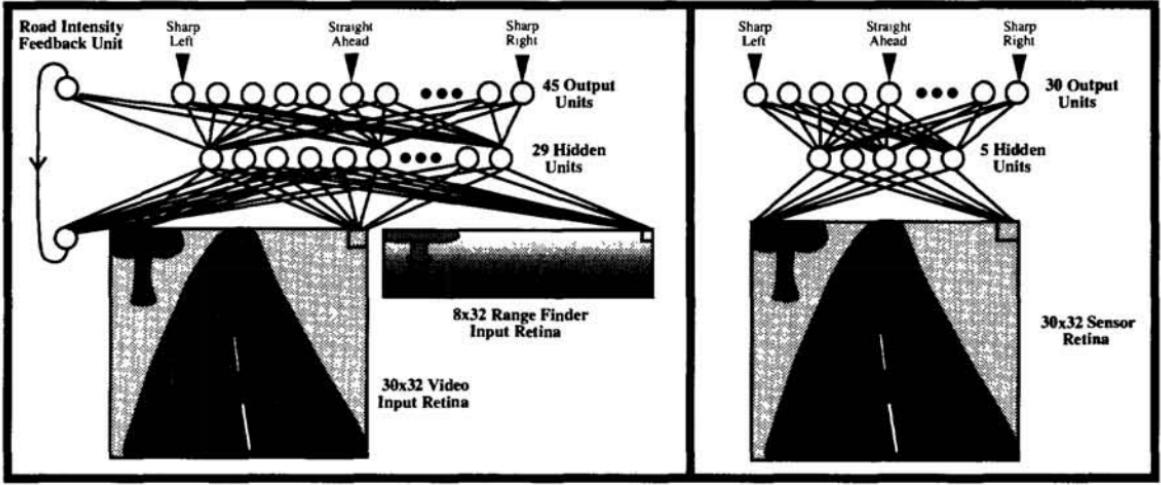


Figure 21: ALVINN architectures. Left: With camera and laser range finder. Right: camera only.
Source: [22]

SAE level	Description
-----------	-------------

- | | |
|---|--|
| 0 | No AV control systems. An example is blind spot indicators. |
| 1 | Basic driver assistance built into vehicle design. An example is the cruise control function. |
| 2 | Basic AV control systems. Driver required to monitor the environment and take back control if need be. A good example is the Tesla Autopilot. |
| 3 | AVs can safely navigate and drive within mapped environment. Driver required to monitor the environment and take back control if need be. |
| 4 | Highly autonomous control capable of handling most conditions but the driver has the option to take control. Renault ISymbioz is an example of a level 4 AV. |
| 5 | Completely autonomous with zero human involvement. |

2.1 Components of an AV

- **LiDAR** - LiDAR provides highly detailed 3D information about the environment around the vehicle and objects in it. LiDAR operates by sending out pulses of lasers and recording the reflections of the pulses from objects. By comparing this with the time taken for the lasers to be reflected(time of flight) and their direction, the distance of these objects can be calculated and mapped in a point cloud. LiDARs have proven to be extremely important for safe navigation due to their high level of accuracy. However due to their low angle resolution they are not able to capture contours of objects. In addition they do not provide texture information and therefore classifying objects can be quite difficult.

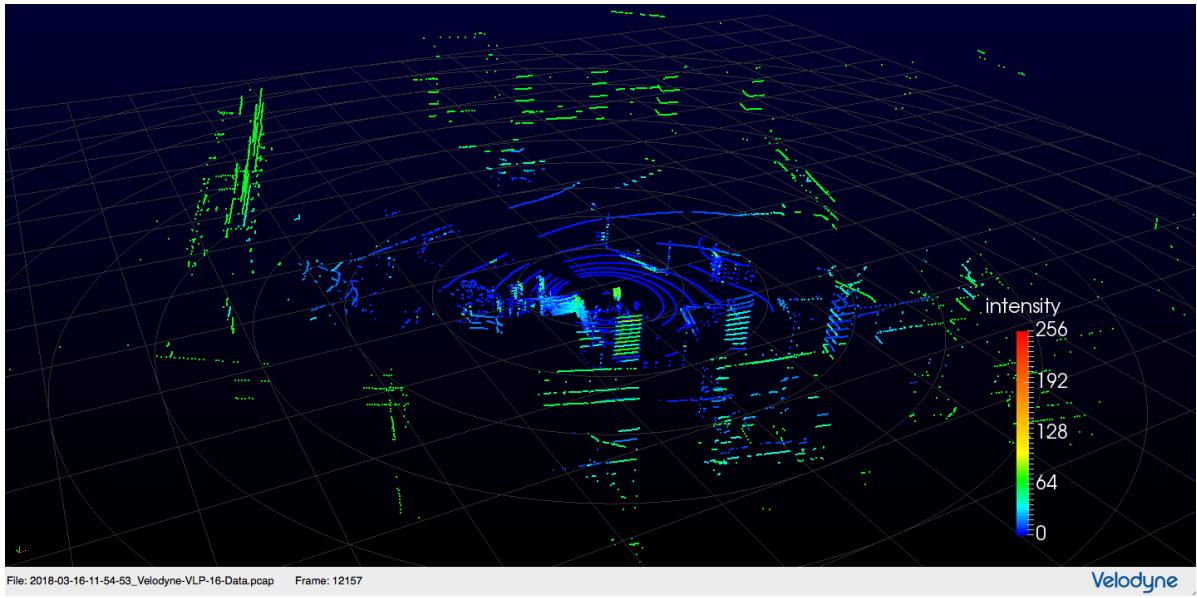


Figure 21: LiDAR Point Cloud

- **Cameras** - Grayscale and color cameras are normally used to detect and classify different objects, markings and signs on the road. Grayscale images tend to offer better edge information than color images. This can be useful for better edge detection. Point clouds can also be obtained by using two cameras to obtain a stereo image that has depth information or by combining a camera and infrared laser to obtain RGB-D images such as in the Kinect sensor.
- **Position Estimators** - Position estimators are a group of sensors used for navigation of the vehicle. These include GPS systems, odometers and gyrometers.
- **Distance Sensors** - Distance sensors such as radars and sonars are important for gauging the distance of objects on the road. Radars are the most commonly used distance sensors and they work by transmitting radio waves and recording the reflected radio waves from objects. As compared to cameras and LiDARs, radars work well in a variety of low visibility scenarios such as poor weather. However, the reflectivity of these radio waves depends on the nature of objects, their size, absorption characteristics and the transmitting power. As such, it is may not be effective for detecting objects with low absorption characteristics such as pedestrians and animals.
- **Processing Unit** - In order to process all the data from sensors on the vehicle, AVs require powerful processing units to do so in real time. Most ML/AI algorithms used for detecting and identifying objects from LiDAR and camera data demand large amounts of processing power. This is achieved through the use of CPUs, GPUs, Field Programmable



Figure 22: Velodyne LiDAR family. From left: Velodyne HDL-64E , HDL-32E , VLP-16

Gate Arrays(FPGA)[4], Application Specific Integrated Circuits(ASICs)[33] or combinations with each other.

These components are used in different operations. In most AVs, the operations can be grouped into three categories.

- **Perception** - This is the first step which involves processing the input from the sensors. In this mode tasks such as object detection and tracking, lane detection, traffic sign detection and recognition are performed.
- **Planning** - This is the next step after detection and recognition tasks are performed . In this stage route and trajectory planning algorithms are run to plan how the vehicle should navigate in the immediate environment as well as a route to a target location. These algorithms are required to handle complex situations to ensure safety of the passengers and other road users.
- **Control** - This stage involves the execution of plans created in the planning stage. This stage is crucial as the actuators involved in steering and movement have to be able to be able to accurately follow the plans. This involves calculation of energy and forces. At this

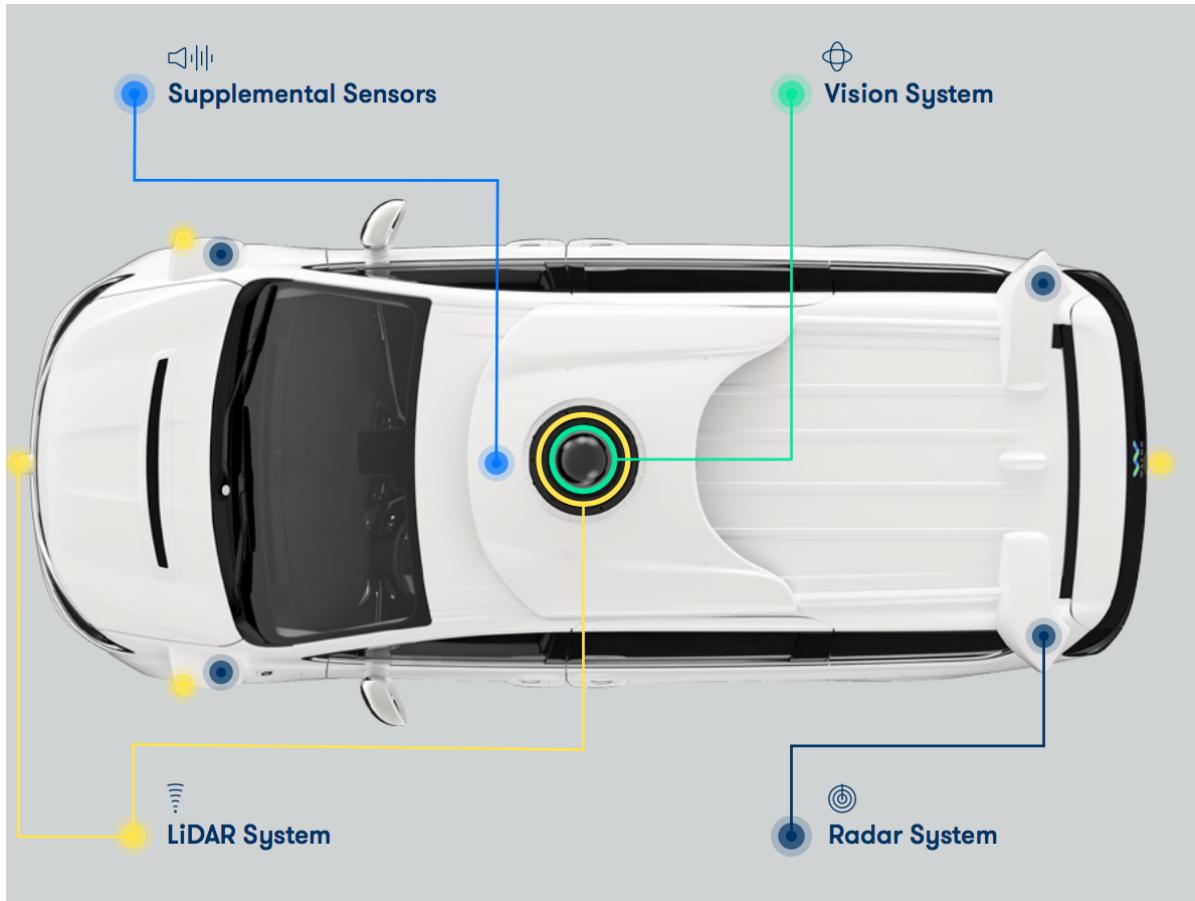


Figure 23: Components of Waymo's Self Driving Car, Source:[36]

stage the trajectories and movement of other road users and objects have to be calculated in order to anticipate and avoid any accidents.

2.2 Industrial Approaches

Despite improvements in traffic safety over the years, over 2.2% of deaths globally are caused by road accidents. Of these accidents, 94% are as a result of human error. A major motivation for the development of AVs is their potential to drastically reduce this by removing the human element. However due to the multidisciplinary nature of AVs there are various factors that need to be considered before public adoption of AV

2.2.0.1 Modelling Risk

Gasser et. al, define safety as the lack of unreasonable risks where risk is the product of the probability of an accident and its severity. From this, it is clear that there are different degrees of

severity each with a various weighting; a fatality is more costly than damage to property. In their study they were able to establish that AVs will have the potential of preventing some accidents, however, they will introduce some risk as a result of their automation as shown in figure ?? that captures this relationship. They further went on to establish the relationship with severity as seen in figure ???. From this relationship they speculated that there is no uniform distribution of risk and in addition, there may be a reduction in severe accidents but an increase in lesser severe accidents. This relationship can be expressed as a ratio that could express the probability of acceptance of a said autonomous system.

$$(2.1) \quad P_{acceptance} = \frac{R_{added\ risk}}{R_{avoided\ risk}}$$

However, this ratio does not necessarily guarantee approval as there are some cases of certain systems such as nuclear plants that the risks outweigh the benefits but due to socio-economic and/or political factors tend to get accepted. In the case for AVs, it is important for

2.2.0.2 Legal and Economic Concerns

2.2.0.3 Acceptance and

2.2.1 Hardware Considerations

In order to be viable, AVs need to meet a few constraints as discussed in [15].

- **Performance**

At the moment there are no clear regulations as to how fast the perception, planning and control pipeline should be. However according to research by Brown et al. [3], humans take around 600 ms to respond and brake when expecting an interruption, however this figure shoots to 850 ms when an unexpected situation arises. In addition, Newell et al [20], established that the fastest human response time is between 100-150ms. These figures can be used as a baseline while developing a pipeline. In this pipeline, there are two factors to consider, namely the frame rate (frequency of the data from sensors) and the processing latency (time taken to process data). As such, an AV system should be able to react within 100ms, that is, faster than the human response time.

- **Storage**

AVs should be able to store maps of different areas with fine granularity for accuracy during localisation. As a result, the storage of these maps can run into tens of Terabytes. Despite recent advances in cloud technology and the emergence of 5G connectivity that is significantly fast. Downloading these maps would take significant time and would also render the car unusable in case of no internet connectivity. As such, the vehicles should have enough storage to store these maps locally.

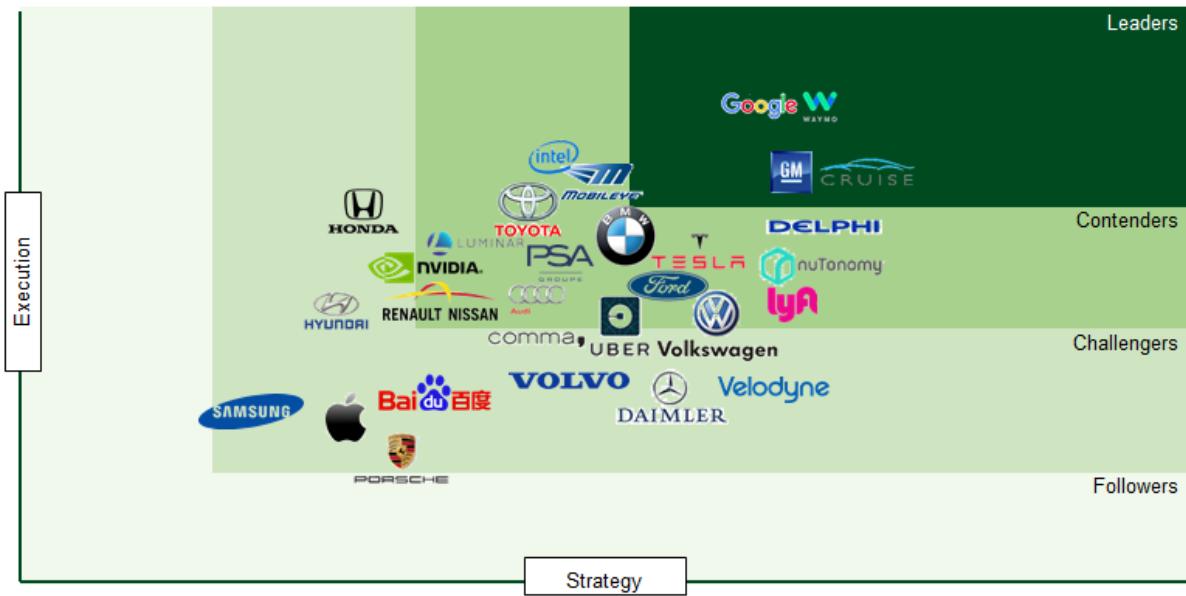


Figure 21: AV Leaderboard. Navigant Research. Source:[28]

- **Power**

Most of the major industry players have moved to electric-vehicles for their AV systems. Depending on the equipment and sensor configurations used in these systems, the power usage can range from 500 watts to 1.5 kilowatts. Given that the AVs have limited battery capacity, heavy power consumption by these systems can lead to poor driving ranges hence making the cars less viable. As such the configuration of these systems has to be carefully considered to ensure a reasonable driving range.

- **Thermal**

Processing components in the AV systems such as CPUs and GPUs require a significant amount of energy to cool. This is necessary to ensure that they operate within their recommended thermal operating range. Failure to do so could result in the failure of these systems. As such, additional cooling systems have to be installed in the vehicle in order to ensure this.

Following the discussion above, the next subsections will review how different industry players are implementing their AV systems with special focus to the issue of perception

2.2.2 Camera

Mobileye[19], is one of the top contenders in the development of camera only AVs. It was recently acquired by Intel and believes that AVs should be able to accurately and safely navigate using cameras only given the fact that humans are able to do so with vision only. Previously, MobilEye

was a supplier of vision systems for Advanced Driver Assistance Systems and had a partnership with Tesla to supply their vision systems prior being acquired by Intel. Given their extensive background in developing these vision systems, the partnership with Intel aims to develop a complete autonomous driving package[12].

2.2.3 Camera and Radar

Another major contender that is using the camera and radar configuration is Tesla. Elon Musk, the founder, believes that LiDAR is not necessary for AV perception following a similar argument as MobilEye. However, this argument was dispelled following a recent crash whereby due to bright sunlight, a Tesla vehicle on 'autopilot' mode was unable to detect a white vehicle thus resulting into an accident leading to loss of life. The invariance of LiDAR to different lighting conditions could have helped in detecting the vehicle.

2.2.4 LiDAR, Camera and Radar

This configuration is used by most of the leading industry players such as Uber, Waymo and GM Cruise. This configuration has proven to be robust and accurate. However, in order to process the amount of fused data, they require large amounts of computational power which in turn may lead to increased power usage.

2.2.5 LiDAR

LiDAR-only systems have not been explored by major companies working on AVs for driving in public roads. However, Doxel, a start up aiming to solve problems in the construction industry by using a robot that autonomously scans and captures the progress of different construction sites. Using these scans, they use a neural network to detect any issues with installations and track changes.

2.3 Object Detection

Object detection has be

2.3.1 Object Detection using Classic Computer Vision Methods

Classic computer vision methods have advanced over the years from simple algorithms such as SIFT[17] and SURF[2] which use descriptors to detect interesting points in images for object detection to much higher level representations such as the Histogram of Gradient(HOG)[8] or Haar Features[35] that were used in classifiers. However, these methods have been unreliable slow and not as accurate. In addition these features are only able to model 2D data. 3D methods such as Intrinsic Shape Signatures(ISS)[37], NARF[34] and Uniform Sampling have since been

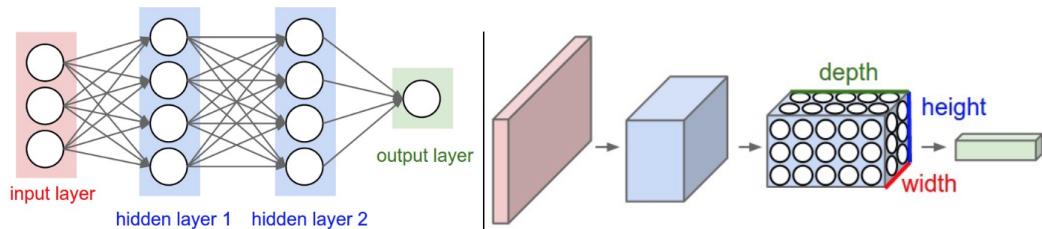
developed but are lacking in terms speed. Furthermore, these feature detectors have to be hand-crafted manually which is a tedious process. As a result, they are rarely used in real-time AV applications due to their speed and accuracy.

2.3.2 Object Detection using Deep Learning

The advancement of neural networks has offset the reliance of classical methods of object detection by allowing for features to be automatically learnt by the network without the need for manual feature extraction. This is attributed to the networks having numerous number of deep layers that are able to capture these features.

Convolutional Neural Networks(CNN)

The idea of CNNs has been there for over two decades now and are inspired by the idea of receptive fields in the primary visual cortex. Currently, they are widely used in image object detection tasks and various CNN architectures have been developed. As compared to regular neural networks which have all neurons connected to each other between hidden layers, neurons in CNNs are connected to a small region in the layer before it(receptive field) and have three dimensions(width,height and depth). By doing so this reduces the number of parameters to be tuned during backpropagation as well as stored.



Left: A regular 3-layer Neural Network. Right: A ConvNet arranges its neurons in three dimensions (width, height, depth), as visualized in one of the layers. Every layer of a ConvNet transforms the 3D input volume to a 3D output volume of neuron activations. In this example, the red input layer holds the image, so its width and height would be the dimensions of the image, and the depth would be 3 (Red, Green, Blue channels).

Figure 21: NN vs CNN Source: <https://cs231n.github.io/convolutional-networks>

Common architectures implement five main layers.

- **Input Layer** Raw image data is fed through this layer.
- **Convolution Layer** This layer calculates the output of neurons connected to a receptive field. The size of the output is determined by the following equation.

$$A_{out} = (A_{in} - K + 2P)/S + 1$$

where A_{in} is the input size, A_{out} is the output size, K is the kernel size, P is the padding size and S is the stride of the kernel.

- **Activation Layer** In this layer the output of the convolution is passed through an activation function such as sigmoid, ReLU or TanH.
- **Pooling Layer** This layer downsamples the output of the activation layer along the width and height dimensions.
- **Fully connected layer** Finally, this layer computes the probability scores for different classes by connecting all neurons from the pooling layer.

Region Proposal Networks(RPN)

R-CNN[] brought forth the idea of region based CNNs with the aim of increasing the accuracy of CNNs. By using selective search, regions of interest(ROI) in input images are proposed. These regions are then resized and then fed through a CNN to generate a bounding box and classification of the image. Finally a classifier, regressor are run on the output. However as this method requires a forward pass of the CNN for each proposal, its speed was quite slow. It also required the CNN, classifier and regressor to be trained separately. Fast R-CNN[11] iterated on R-CNN to further improve the speed and accuracy by using ROI Pooling whereby instead of running each proposal in a CNN, regions that were overlapping were pooled and fed through the CNN thus sharing the computation. Furthermore, this model was end to end and thus combined the classifier and regressor by adding a softmax layer and linear regression layer after the CNN. Faster R-CNN[27] sought to remove the bottleneck imposed by the using selective search to propose regions. By using a sliding window with various anchor boxes on features maps generated by a forward pass of a CNN, bounding boxes on certain regions of the image could be generated. These regions were then passed onto Fast R-CNN to classify the objects.

Depending on the input, these two network architectures have formed the back bone of most 2D and 3D object detection applications.

Monocular Vision Monocular vision is derived from a single camera. This representation is 2D and lacks depth information. Mono3D by Chen et al [5] is a state of the art object detection model that uses monocular images where they were fed through a CNN and the resulting 2D object proposals were extruded into 3D bounding boxes by performing segmentation.

Stereo Vision Stereo vision involves calibrating two cameras intrinsically and extrinsically by matching similar points in images from both cameras. Using the resulting stereo image, it is possible to estimate the depth of objects in the image. However the accuracy of this depends on a few factors such as the resolution, camera quality, lens distortion and the calibration algorithm. 3DOP [6] is one of the state-of-the-art models that use stereo images in CNNs to detect objects.

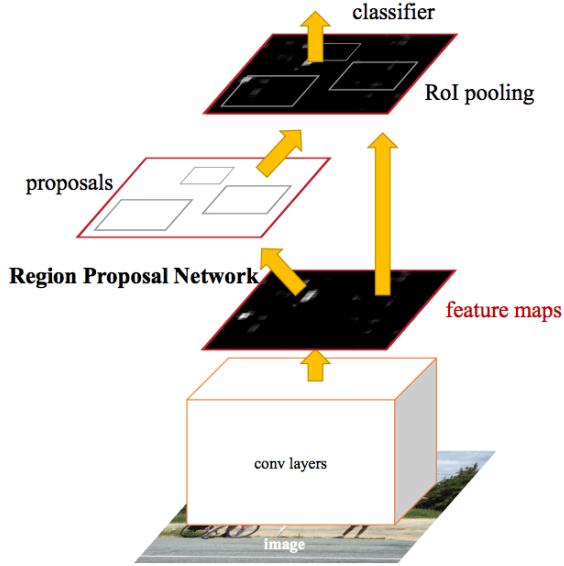


Figure 22: Faster R-CNN overview. Source:[27]

LiDAR+Camera By fusing input from a monocular image and point clouds, it is possible to extract depth information from the images. Chen et al further reiterated their monocular approach into MV3D[7] a multiview 3D object detection network. In their work, they were able to fuse LiDAR and camera input to create a bird’s eye view of the surrounding environment and further perform object detection of features.

LiDAR Point cloud object detection is particularly challenging as established 3D object detection methods for images cannot be applied to point clouds. Point clouds are a set of geometric points in a euclidean space. They are unordered in nature and this presents a particular problem to deep learning methods as they need to be invariant to permutations of the input set. Depending on how they manipulate the point clouds, current point cloud based object detection models can be split into three groups:

1. Direct Manipulation

Pointnet[25] developed by Qi et al was able to overcome this challenge on small sets of point clouds. By using point clouds as direct input into Recurrent Neural Networks[18] They were able to create a global point cloud signature that could classify and segment 3D objects in the point cloud. They further improved this model into Pointnet++[26] which recursively applied PointNet on nested partitions of the point clouds using a hierarchical neural network. In doing so they were able to capture the local structures of the point clouds as a result of their metric nature and thus achieve better performance than PointNet. However this approach has high memory and computational requirements. For both of these publications, the source code was released for public use. Consequently, they have

formed a fundamental foundation for further development of point cloud object detection methods such as F-PointNet[24].

2. Voxel Based

By dividing the pointclouds into 3D voxel grids, points within these voxels can be sampled to perform feature extraction. VoxelNet [38] is a state of the art architecture that uses this method to extract these features through Voxel Feature Extraction layer. The extracted features are then fed through a RPN to generate bounding boxes and classify objects within them.

3. 2D-Based

This involves projecting the pointcloud into a 2D image plane coordinate system which can then be trained on existing image based CNNs such as VGG-16. VeloFCN [14] was one of the first methods to implement this approach using a fully convolutional network. However, VeloFCN lacked bird eye and front view representations.

Network Model	Mode	Code	Inference Time
VoxelNet	<i>LiDAR</i>	Unofficial	100ms
AVOD	<i>LiDAR+Image</i>	Yes	100ms
MV3D	<i>LiDAR+Image</i>	Unofficial	360ms
MONO3D	<i>Mono Image</i>	Yes	4.2s
3DOP	<i>Stereo Image</i>	Yes	3s

Table 21: Modes of Operation of Common Architectures

2.4 Conclusion

This chapter was opened with a brief history and motivation for AVs. Different components and their advantages were discussed

In addition a discussion on current approaches that are being explored where we were able to Shift from classic computer vision to deep learning methods. Where Finally a discussion on From these discussions, I came up with the following hypothesis.

CHAPTER



IMPLEMENTATION

Following the prior discussion, this section seeks to address three challenges affecting object detection models for AVs. This will be necessary to achieve the objectives outlined. Firstly, the necessary datasets to be used for the different evaluation tasks will be outlined. The second task will be to determine whether it is possible to detect the context of images/point clouds from their features. Upon successful classification of images/pointclouds into their respective contexts, the performance of multimodal and LiDAR in different contexts will be assessed. Finally, the performance of the LiDAR only model will be validated using an external dataset to determine if indeed LiDAR based object detection models are extrapolable to different datasets.

3.1 Datasets

3.1.1 Context Classification

1. **Cityscapes** - Cityscapes is a large-scale dataset used to develop pixel and instance level semantic segmentation models. It is composed of scenes from 50 different cities captured using stereo cameras. In 5000 scenes, the images have been annotated on a pixel-level. 20000 scenes have also been coarsely annotated. In addition to this, a benchmark suite to evaluate these models is available.
2. **KITTI Object Detection**- KITTI dataset was collected with the aim of encouraging the development of computer vision and robotic algorithms for AVs. The data was captured by a car fitted with a number of sensors including a high resolution greyscale and colour cameras, Velodyne HDL-64 LiDAR and a GPS/IMU inertial navigation system. The vehicle was driven around a city, rural areas and highways. This dataset is available as part of a benchmark suite for various AV object detection models. The dataset consists of a total

of 14,999 images and corresponding point clouds with 7,481 as labelled training files and 7,518 being unlabelled testing files. The testing set are used for evaluating results through the official test server. However, the test server is not available for general public use.

3.1.2 LiDAR Performance Evaluation

1. **University of Bristol Smart Internet Lab(SIL)** - This dataset was captured from a Velodyne VLP-16 LiDAR on stationary vehicle positioned at the Millennium Square in Bristol as part of the connected and autonomous vehicles project. Available as LiDAR network capture files(.pcap).
2. **KITTI Object Detection**

3.2 Can We Automatically Detect The Context in Driving Scenes?

Determining the context from images and point clouds where no prior information about their location nor context has been provided is quite difficult. This is often the case in a lab setting where the data is obtained from public database. This was indeed the case with images and point clouds from the KITTI dataset that did contain such information and since the frames were discontinuous as they shuffled into training and testing sets, it was difficult to infer the context using the sequence of prior frames. Notably, this was also evident in the Cityscapes dataset. From this observation, it was necessary to create a context detector for images and furthermore point clouds that can be used in a lab setting.

To begin with, I visually classified 3749 images from the KITTI dataset to be used as training and test data for the context detector. Of these 1029 were non-urban and the remaining 2720 urban. The context of the image scenes were determined using the following criteria defined by the UK Department for Environment, Food and Rural Affairs [9] Urban regions are characterised by:

1. High density of road users often including pedestrians, cyclists and vehicles.
2. Presence of numerous buildings that tend to be large and multi-story.
3. Presence of multiple transport facilities such as trams.

On the other hand, non-urban regions are characterised by:

1. Low/medium density of road users mostly vehicles.
2. Few or no buildings, mostly isolated dwellings.
3. Long stretches of motorways surrounded by vegetation.

Following this, the classified data was split on a 80:20 ratio to be used as the training and test set respectively.



Figure 31: Urban scene



Figure 32: Non-Urban scene

3.2.1 Image Context Detection

Images provide a rich representation of a scene that can be exploited through image segmentation. Various image segmentation models have been developed to break down images into coherent regions. Semantic segmentation involves classifying these regions into various classes on a pixel level. Classic computer vision techniques for semantic segmentation utilised Textron Forests[32] and Random Forest classifiers [31], however, deep learning techniques have overtaken them in terms of performance. Developed by Google, Deeplab-V3 is one of the best performing open-source models on the Cityscapes semantic segmentation benchmark table.

Implementation Details

Extracting Semantic Histograms Using the Deeplab-V3 network, I performed semantic segmentation on the classified images. For each image, a semantic histogram containing the number of pixels in each semantic class as seen in figure 31 was calculated.

Training - By matching the semantic histograms with the context label of the image, I was able to train a Support Vector Machine(SVM) using the scikit-learn library to classify the image context. Linear and radial basis functions were tested on the SVM to evaluate their performance.

Evaluation - The trained SVM classifier was then used to classify keypoint descriptors that were extracted from the test point clouds. Given that a single point clouds can contain various features, majority voting was used to establish the dominant feature class.

Penalising misclassifications - In a realistic setting, I speculated that it would be more costly to classify urban as non-urban than non-urban. This is because non-urban settings are less dynamic than urban settings. In order to extend this as part of the AV system, it was important to create a penalising factor to the classifier as defined below:

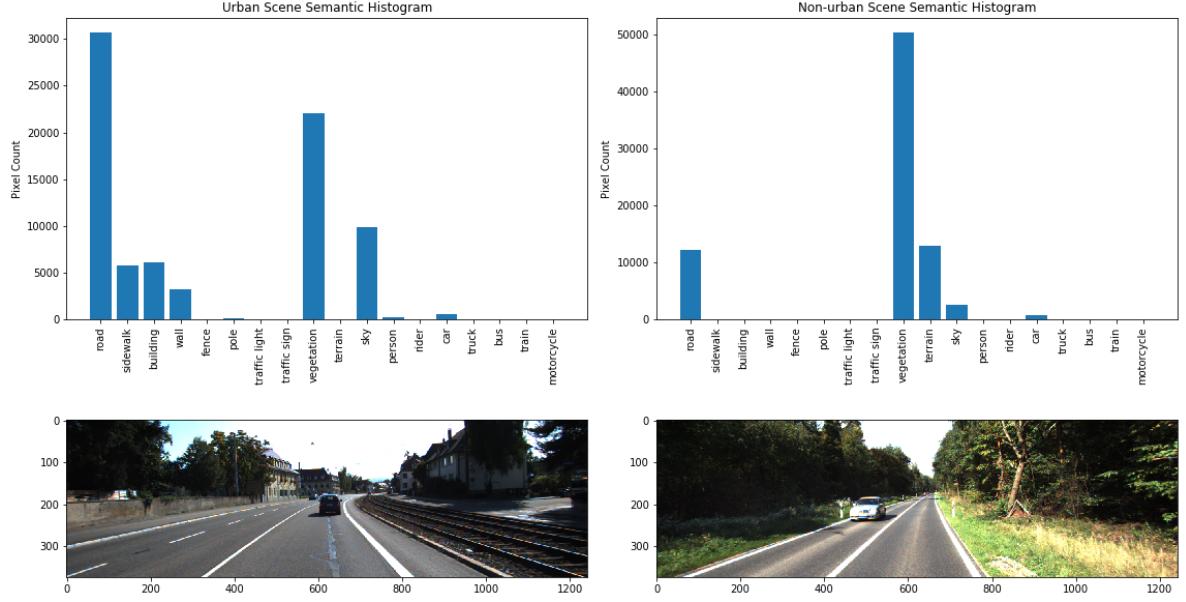


Figure 31: Semantic Histograms

$$R_{urban} = \lambda_{11}P(urban|x) + \lambda_{12}P(non-urban|x)$$

$$R_{non_urban} = \lambda_{21}P(urban|x) + \lambda_{22}P(non-urban|x)$$

$$Context = \begin{cases} \text{Urban,} & \text{if } R_{urban} < R_{non-urban} \\ \text{Non-urban,} & \text{otherwise} \end{cases}$$

where λ is a 2x2 matrix with λ_{11} and λ_{22} set to 0 and λ_{12} set to 1.6 and λ_{21} set to 1.

3.2.2 PointCloud Context Detection

As compared to point clouds from RGB-D or stereo cameras, point clouds from LiDAR lack rich features. This is due to the sparse nature of the point clouds and as a result, the resulting representation tends to have a large number of 'holes'. Despite the availability of point clouds segmentation models, there are no datasets with annotations of outdoor scenes such as Cityscapes. As such using semantic histograms to characterise scene contexts in this case would not be suitable. A different approach was therefore used by first detecting the key points in the cloud, extracting and classifying features from them and later on matching these features in other point clouds.

1. Keypoint extraction

Similar to classic image processing, key points are salient points that are not affected by

any form of transformation or distortion. As such, they are distinctive and repeatable in different points-of-view. They can either be global or local on a given surface with the former being preferred to extract 3D shape information and the latter for 3D object detection. Currently, only a few 3D keypoint detectors have been proposed. That is Intrinsic Shape Signatures(ISS)[37], NARF[34] and Uniform Sampling of point clouds in voxels.

2. Feature description

Keypoints extracted act as features that offer a robust representation of a local or global surface. In order to provide a description of features, a descriptor is used to encode the relevant surrounding of the keypoints. These include include Fast Point Feature Histograms(PFH)[29], Signatures of Histograms of Orientations(SHOT)[30].

3. Feature Matching

Once different feature descriptors have been obtained. They can used for object detection and segmentation using nearest-neighbour functions. For high dimension data as in the case of point-clouds, randomised kd-trees and FLANN are usually used to achieve decent performance.

3.2.2.1 Implementation Details

The above methods were implemented in Python using PyDriver[21] framework and Point Cloud Library. Point Cloud Library is an open-source project for 2D/3D image and point cloud processing. The library contains various software implementations of keypoint extractors, descriptors and feature matching functions. However, the implementations are written in C++ and the available Python wrappers are only for Python 2.7. PyDriver is a framework contains a suite of object detection, classification and tracking functions. In addition, it contains Cythonized wrappers to the Point Cloud Library functions necessary for keypoint detection, extraction and feature matching and can run on Python 3.x. Code is available at <https://github.com/lpltk/pydriver>

2000 images from the classified context database were used as training samples, and the remaining 1749 as test samples.

Training - Keypoints within the ground truth bounding boxes were extracted using ISS with a salient and NonMaxSuppression(NMS) radius of 0.25. SHOT feature descriptor with a radius of 2 was applied on the keypoints to create feature descriptors. Each feature was assigned with the context class of the point cloud. An AdaBoost classifier was then trained on the feature types.

3.2.3 Implementation Issues

Despite being able to visually classify a large number of images, some images were difficult to classify as they exhibited characteristics from both contexts. Furthermore, due to the fact that

the images was not contiguous , I was unable to infer their context by examining preceding image frames. As such, this affected the performance of both classifiers.



Figure 31: Ambiguous scene

3.3 Do Some Object Detection Methods work better in Different Contexts?

Having been able to successfully classify and create a dataset of urban and non-urban images and corresponding point clouds, the next step was evaluating how a LiDAR based and image+LiDAR based object detection model would perform in different contexts. Multimodal methods have been argued to perform better in 3D object detection tasks whereby objects are partially occluded or relatively small. Such characteristics are common in urban areas as compared to non-urban areas which have fewer and more distinct objects. In light of this argument, I assessed whether LiDAR-only models can perform better or at par with multimodal methods in different contexts.

For this task, VoxelNet(LiDAR only) and AVOD(LiDAR+image) proved to be the most suitable candidates. As seen in table 21 both models had the least inference time that would similar inference times, and both were open-source with the only exception being that the VoxelNet implementation was unofficial¹. Both of these models were written in Python and run on TensorFlow Deep Learning Framework.

3.3.1 VoxelNet

VoxelNet is an end to end point cloud based 3D object detection network that was recently published by researchers at Apple Inc. The network is one of the top performing LiDAR only models on the KITTI benchmark. However, the code was never released and only unofficial implementations are currently available online. Code is available at <https://github.com/qianguih/voxelnet>

3.3.1.1 Architecture Overview

To better understand how it works it is important to understand the architecture. VoxelNet is composed of three fundamental blocks; a feature learning layer, CNNs and finally a RPN.

¹Had similar evaluation results as the unreleased original on KITTI benchmark.

3.3. DO SOME OBJECT DETECTION METHODS WORK BETTER IN DIFFERENT CONTEXTS?

1. Feature Learning Layer

In this layer, the point clouds are divided into equal 3D voxels. Points within these voxels are then grouped. For the car detection, 35 points are sampled and for pedestrian and cyclist detections, 35 points are sampled from within the voxels. These sampled points are then fed through multiple FCNs that learn and aggregate features from the points. The result is a 4D tensor that represents these features per voxel.

2. Convolutional Middle Layers

The 4D tensor containing these features is then passed through multiple convolutional middle layers. In doing so the different layers learn different features in an expanding receptive fields. This allows the network to learn complex shape information at different scales.

3. Region Proposal Network

Finally, the feature maps from the convolutional middle layers are fed through a RPN to classify and generate oriented 3D bounding boxes for the objects.

Loss Function

$$L = \alpha \frac{1}{N_{\text{pos}}} \sum_i L_{\text{cls}}(p_i^{\text{pos}}, 1) + \beta \frac{1}{N_{\text{neg}}} \sum_j L_{\text{cls}}(p_j^{\text{neg}}, 0) + \frac{1}{N_{\text{pos}}} \sum_i L_{\text{reg}}(\mathbf{u}_i, \mathbf{u}_i^*)$$

where:

- L_{cls} is the binary cross entropy loss that is calculated for both positive and negative classes.
- α is a weighting parameter for the positive anchors normalised binary cross entropy loss.
- β is a weighting parameter for the negative anchors normalised binary cross entropy loss.
- L_{reg} is the SmoothL1 regression loss.

3.3.1.2 Implementation Details

Training - To encourage generality, the model was trained on 3732 point clouds from the original KITTI training set. VoxelNet was trained both car and pedestrian classes on two NVIDIA P100 GPUs with a batch size of 4 for 160 epochs with $\alpha = 1$ and $\beta = 10$. Adam optimizer was used with an initial learning rate of 0.01 for the epochs ≤ 80 , 0.001 for $80 \leq \text{epoch} \leq 120$ and 0.0001 for $\text{epochs} \geq 120$.

Changes to Loss Function - In an effort to improve performance of pedestrian and cyclist detection, I implemented the focal loss function as defined in [16]. This stemmed from the fact that there were few frames with pedestrians as compared to the car class and as a result the class imbalance was high in that negative anchors \gg positive anchors.

3.3.2 AVOD

AVOD is a state of the art multimodal network that offers improved performance for small object classification by fusing LiDAR point clouds and images to perform detection. Code is available at <https://github.com/kujason/avod>

3.3.2.1 Architecture Overview

AVOD's architecture is comprised of two distinct stages:

1. **Feature Extractors** To create a BEV map, LiDAR point clouds are projected into the XY plane and discretised to create BEV 2D grid. In each cell a height feature is defined by the maximum height of the points in the cell and it's reflectance defines the intensity feature in five slices between [0,2.5] of the Z axis. For each cell, the normalised point cloud density is calculated using $\min(1.0, \frac{\log(N+1)}{\log 16})$ where N is the number of points.

For the image and BEV map similat feature extractors consists of an encoder and decoder are applied. The encoder is a CNN that creates a downsampled feature map resulting in higher representation. The decoder is a Frustum Point Net that upsamples the feature map to original size resulting in higher resolution. The output from the image and point cloud feature is fused by passing it through a 3x3 convolution layer.

2. **Multimodal Fusion Region Proposal Network**

Feature maps from the point cloud and image feature extractors are fused in a RPN to generate proposals. The top proposals are then passed through a second RPN that performs object classification and bounding box regression to generate oriented 3D bounding boxes for the objects.

Implementation Details

Training - Similar to VoxelNet's training method, the model was trained on 3732 point clouds and images from the original KITTI training set. The AVOD-FPN model was trained for both car and pedestrian classes for 120000 steps on two NVIDIA P100 GPUs. The adam optimiser with an exponential decay learning rate was chosen with an initial learning rate of 0.0001 and a decay rate of 0.8 each 30000 steps.

3.3. DO SOME OBJECT DETECTION METHODS WORK BETTER IN DIFFERENT CONTEXTS?

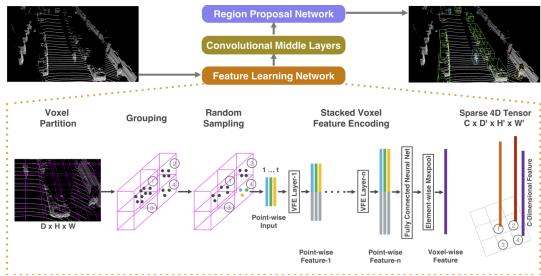


Figure 31: VoxelNet Architecture

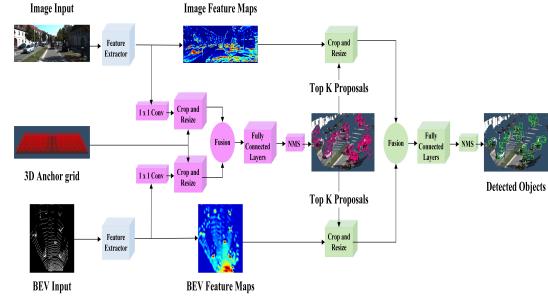


Figure 32: AVOD Architecture

3.3.3 Evaluation

Once both models were trained, they were evaluated on 2058 images/point clouds from the context dataset with equal number of class samples. This was necessary to avoid bias as the urban samples were significantly more than the non-urban samples. For each context, each model was run on a single NVIDIA P100 GPU and the results of 3 runs averaged.

3.3.4 Is VoxelNet Performance Valid for Different Datasets?

As compared to image object detection methods that can perform well regardless of the input image resolution, LiDAR point clouds can vary in different scenes. This can depend on various factors but is largely determined by the LiDAR sensor and its frequency in terms of points per second. Coupled with the sparse nature of point clouds, this implies that LiDAR only models may not be able to generalise to input from different LiDAR sensors with different frequencies. I therefore set out to understand how LiDAR sensors with different frequencies affect the performance of VoxelNet. In this section I aim to validate the VoxelNet model using the SIL dataset that was obtained from a Velodyne VLP-16 sensor with a lower frequency of points than that of KITTIs Velodyne HDL-64E².

Preprocessing

Prior to using the SIL data, I had to preprocess it to the necessary input format. VoxelNet only required four input fields, **XYZ values** and **reflectance** of the points, whereas the dataset was in the form of network capture files(.pcap). To convert them into this format, I exported the capture file into CSV format for each of the frames. This was done in VeloView using the CSV export function. The result was a CSV file with 13 fields as seen in figure 31

²See table 31

CHAPTER 3. IMPLEMENTATION

Points_m_XY2:0	Points_m_XYZ:1	Points_m_XYZ:2	x	y	z	intensity	laser_id	azimuth	distance_m	adjustedtime	timestamp	vertical_angle
0.01759156584739685	5.304635319510943	-1.421434164047241	0.01759156507648251	5.304835470033373	-1.421434195703044	1.0	0.0	19.0	5.492	3433080897.0	3433080897.0	-15.0
0.07194533199071884	20.6107349395752	0.3597639203071594	0.07194532899821898	20.61073481953913	0.3597639062981624	20.0	1.0	20.0	20.614	3433080900.0	3433080900.0	1.0
0.02194887027144432	5.988438129425049	-1.382549166679382	0.02194887062991644	5.98843819471877	-1.38254917997394	1.0	2.0	21.0	6.146	3433080902.0	3433080902.0	-13.0
0.02645202539861202	6.88900899887085	-1.339097499847412	0.02645202482367224	6.889008789131341	-1.339097529562591	1.0	4.0	22.0	7.018	3433080907.0	3433080907.0	-11.0
0.3046736121177673	75.89746856689453	6.64022159576416	0.3046736085365985	75.89747013826725	6.64022172845881	63.0	5.0	23.0	76.188	3433080909.0	3433080909.0	5.0
0.03370168805122375	8.045639038085938	-1.274315118789673	0.03370168894267474	8.045638637731509	-1.274315152217721	1.0	6.0	24.0	8.146	3433080911.0	3433080911.0	-9.0
0.3095385730266571	70.94056701660156	8.710489273071289	0.3095385870564042	70.94056833290217	8.71048945053951	63.0	7.0	25.0	71.474	3433080913.0	3433080913.0	7.0
0.04397721216082573	9.691121101379395	-1.18993227134705	0.04397721326242784	9.691120843322894	-1.1899322690786	1.0	8.0	26.0	9.764	3433080916.0	3433080916.0	-7.0
0.3213215470314026	70.80862426757812	11.2150993347168	0.3213215462835961	70.80862345723331	11.21509966766423	63.0	9.0	26.0	71.69200000000001	3433080918.0	3433080918.0	9.0
0.05691538378596306	12.07733017883301	-1.056676268577576	0.05691538501887421	12.0773041569957	-1.056676225072608	3.0	10.0	27.0	12.124	3433080920.0	3433080920.0	-5.0
0.08184240758419037	16.16960144042969	-0.847423791885376	0.0818424105905633	16.16960230547166	-0.8474238034857464	2.0	12.0	29.0	16.192	3433080925.0	3433080925.0	-3.0

Figure 31: SIL Data Frame CSV

Out of the 13 fields, I extracted the XYZ values and intensity(reflectance) of the points and saved them as Numpy binary files. This resulted in a smaller file sizes(from around 4MB to 70KB) that are also quicker to read than normal CSV files.

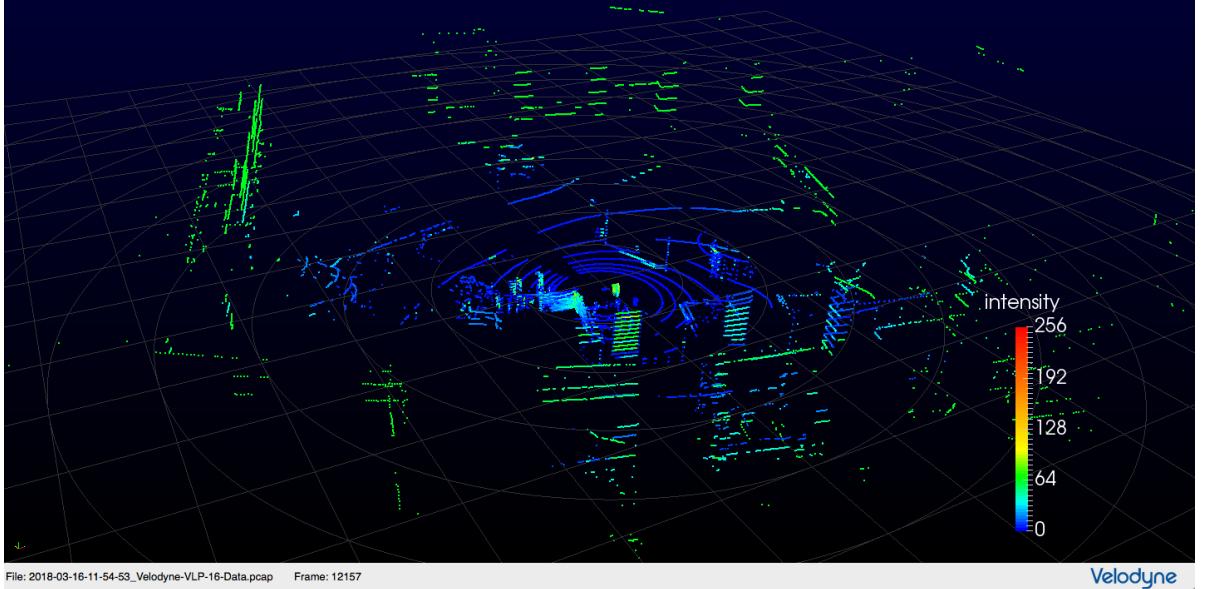


Figure 32: Veloview Frame of SIL data

Finally, I modified VoxelNet's configuration file to include the specifications for the VLP-16 LiDAR sensor. This was essential as the network was initially trained on KITTI point clouds that were obtained from a Velodyne HDL-64E with different specifications as seen in table 31. These specifications are necessary for manipulating the point clouds and therefore directly affect the performance of the network. The modified entries were the angular and vertical resolution, maximum and minimum XYZ values and the height of the VLP-16 LiDAR sensor on the vehicle. Given that there was no image data for visualisation, the image size in the configuration file was set to 0 to prevent unnecessary tensor memory allocation for images that are used for visualisation.

LiDAR	Hor FOV	Ver FOV	Range	Angular Resolution	Points/second	Channels
HDL-64E	360°	26.9°	120m	~0.4°	~2.2 Million	64
VLP-16	360°	± 15°	100m	0.1°	600,000	16

Table 31: Velodyne HDL-64E and VLP-16 Specifications.

Annotating Ground Truth Labels

The SIL dataset did not contain any ground truth labels for the point clouds as compared to the KITTI dataset. As a result, there was no metric to assess the performance of the network. To solve this, objects from the dataset needed to be annotated to obtain some ground truth bounding boxes.

L-CAS Annotation Tool

Developed by the Lincoln Centre for Autonomous Systems Research, this tool provides a semi-automatic labelling function that clusters and highlight regions of interest that may contain objects. The user can then label them depending on the object class such as car, cyclist, pedestrian. The result is then stored in a file containing the objects and their positions in the point cloud. To use this tool, the .pcap capture files were converted into .pcd format using the ROS Velodyne point cloud package³.

3.3.4.1 Validation

Using the annotation tool, I annotated 500 frames to be used for validation. The output was a list of files containing object labels and their bounding boxes that would serve as the ground truth for validation. I then ran VoxelNet on these frames to generate the predicted bounding boxes for different object classes.

3.4 Conclusion

In the beginning of this section I set out to answer three important questions. Following a methodological approach I was able to establish tools to detect and characterise the context of images and point clouds. Building up on this, I was able to identify two neural networks that could potentially perform better in different contexts. Finally, I was able to identify a potential challenge with LiDAR only models and we were able to test whether VoxelNet would still perform given a different dataset.

³See listing A.1

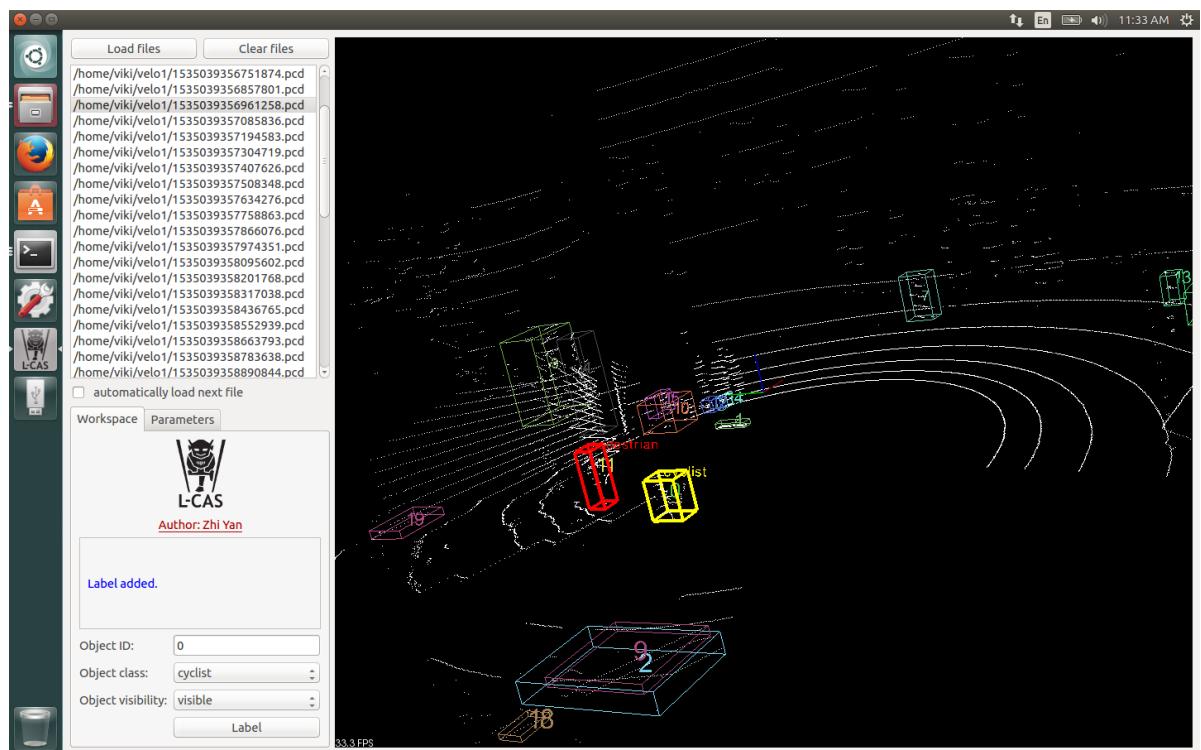


Figure 33: L-CAS Annotation Tool

RESULTS AND ANALYSIS

Following the described methodology in the previous chapter, this chapter seeks to interpret and analyse the results obtained while evaluating

4.1 Can We Automatically Detect Scene Contexts?

4.1.1 Metrics

In evaluating the classifiers, it was crucial to understand the predictive power of the two context classifiers. As such, precision, recall and f1-score were selected as the main performance metrics.

Precision is the percentage of positive predictions that were correct.

$$\text{Precision} = \frac{TP}{TP + FP}$$

Recall is the percentage of positive samples that were correctly predicted.

$$\text{Recall} = \frac{TP}{TP + FN}$$

F1-score is the balanced mean of the precision and recall.

$$2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

4.1.2 Analysis

I established that determining the context from images is much more accurate than from point clouds. As highlighted before, images offer higher representational power and resolution than

sparse point clouds. As seen in figure 41, the f1-score of the image classifier outperformed the point cloud classifier by more than 30% .

Context	Context Detection using PointCloud Feature Matching				Context Detection using Image Segmentation			
	precision	recall	f1-score	support	precision	recall	f1-score	support
<i>urban</i>	0.52	0.45	0.48	206	0.81	0.9	0.85	193
<i>non-urban</i>	0.51	0.58	0.55	206	0.9	0.81	0.85	218
avg / total	0.51	0.51	0.51	412	0.86	0.85	0.85	411

Table 41: Comparison of context detection methods

With regard to the point cloud context classifier, majority voting resulted in around 50:50 precision for both classes. The recall for the non-urban was 7 percent more than urban which was below 50%. By implementing the penalty function in the point cloud classifier, recall for the urban class improved by 17% as seen in table 42.

Penalised Point Cloud Context Classifier				
context	precision	recall	f1-score	support
urban	0.52	0.62	0.54	206
non-urban	0.46	0.33	0.38	206
avg / total	0.47	0.47	0.46	412

Table 42: Penalised Point Cloud classifier

4.2 Do Object Detection Models Perform Better in Different Contexts?

Metrics

KITTI Evaluation Protocol

The models were evaluated using the official KITTI evaluation protocol. Average precision(AP) is the area under a recall-precision curve that was calculated after predictions. Predicted bounding boxes were considered as positive matched if they were above a predefined intersection over union(IoU). IoU is calculated as:

$$IoU = \frac{\text{Overlap area}}{\text{Union area}}$$

Where the overlap area is the region that is shared between the predicted bounding box and the ground truth bounding box and the Union area is the total area of the predicted bounding box and ground truth box. This was calculated for 2D bounding boxes, 3D bounding boxes and Bird-Eye View bounding boxes. The IoU threshold was set at ≥ 0.7 for the car class and ≥ 0.5 for the pedestrian class.

GPU

As mentioned earlier, the resources of an AV are limited and therefore it is encouraged to implement efficient models. In light of this, I monitored the GPU statistics to better understand how different models affect the GPU. During inference, GPU statistics were collected at an interval of 1 second using NVIDIA System Management Interface(nvidia-smi). The following metrics were collected.

- **Power Draw** - Watts
- **Memory Used** - MB
- **Temperature** - °C
- **Clock Speed** - MHz

For uniformity purposes, all the experiments were performed on single NVIDIA P100 GPU with 16GB High Bandwidth Memory.

Baseline

A baseline, was established by running each model was run on the original KITTI validation dataset. The results can be seen in table in 41. In addition, the temperature and power metrics were collected and can be seen in the last boxplot of figures 41 and 42.

Model	Class	Car			Pedestrian		
		Easy	Medium	Hard	Easy	Medium	Hard
AVOD	<i>2D Bounding Box</i>	86.987999	77.10569	68.012459	56.214417	51.969517	46.578247
	<i>BEV Bounding Box</i>	86.00412	74.355942	65.62397	50.94368	49.857227	44.704525
	<i>3D Bounding Box</i>	75.447769	63.742085	54.334251	48.757492	44.84359	43.014023
VoxelNet	<i>2D Bounding Box</i>	65.951279	56.683437	55.972511	66.309807	60.556213	53.479141
	<i>BEV Bounding Box</i>	81.242409	70.672707	65.347595	50.308723	46.726196	41.24551
	<i>3D Bounding Box</i>	40.000149	34.892159	31.376112	34.161884	30.576441	28.982103

Table 41: Baseline

4.2.0.1 Running Each Model Exclusively

To investigate whether the performance of each model varies in different contexts, each model was run exclusively on a single GPU with data from each context dataset.

Analysis

Car - AVOD outperformed VoxelNet on 2D and 3D detections in both contexts. However, VoxelNet performed slightly better in bird-eye-view detections as compared to AVOD.

Model	Context	Urban			Non-Urban		
		Easy	Medium	Hard	Easy	Medium	Hard
AVOD	<i>2D Bounding Box</i>	86.987999	77.10569	68.012459	89.186172	79.754562	78.550514
	<i>BEV Bounding Box</i>	86.00412	74.355942	65.62397	87.11274	76.859818	75.720955
	<i>3D Bounding Box</i>	75.447769	63.742085	54.334251	75.430656	64.199951	62.902302
VoxelNet	<i>2D Bounding Box</i>	69.597939	65.98201	59.284454	77.520409	67.733231	62.281651
	<i>BEV Bounding Box</i>	86.34037	76.187859	68.103294	88.635025	75.361214	69.502548
	<i>3D Bounding Box</i>	73.632263	58.473991	50.744587	68.620865	49.455608	45.809917

Table 41: Car AP

Model	Context	Urban			Non-Urban		
		Easy	Medium	Hard	Easy	Medium	Hard
AVOD	<i>2D Bounding Box</i>	78.626633	77.407684	70.029228	76.096977	70.946465	71.205276
	<i>BEV Bounding Box</i>	80.880447	80.393105	79.296593	77.806831	77.654358	71.792923
	<i>3D Bounding Box</i>	80.649719	80.142876	79.056778	77.538368	71.868484	71.564285
VoxelNet	<i>2D Bounding Box</i>	72.129387	65.245384	65.508232	84.218407	76.52832	76.593803
	<i>BEV Bounding Box</i>	74.766678	73.944069	74.138588	89.830276	80.201164	80.114815
	<i>3D Bounding Box</i>	71.177628	64.247559	64.18145	87.206612	77.964493	78.036652

Table 42: Pedestrian AP

Context	VoxelNet			AVOD		
	min	max	mean	min	max	mean
Urban	0.113	3.813	0.129	0.095	2.457	0.112
Non-urban	0.113	2.224	0.127	0.096	2.506	0.113

Table 43: Inference Time on Single NVIDIA P100 GPU

Pedestrian - As compared to the car class, there was a distinct performance boundary whereby AVOD dominated all detections in the urban context whereas VoxelNet did in the non-urban context.

It is important to bear in mind that for the car class, the IoU threshold for detections is 0.7 as compared to the pedestrian IoU threshold of 0.5. As mentioned earlier while describing the model architectures, due to the fact that AVOD fuses the top proposals from the image and LiDAR point cloud feature maps, the predicted bounding boxes tend to be tighter and thus have a higher IoU. Going by the pedestrian class result, it is plausible to argue that the performance of VoxelNet would have been higher or close to AVOD's given that the IoU threshold was slightly reduced.

GPU - As seen in figures 42 and 41, AVOD generally ran at a lower temperature and consumed less power as compared to VoxelNet. In fact, temperature-wise, VoxelNet ran hotter over 75% of the time. Furthermore, AVOD's temperature and power usage varied less as compared to VoxelNet's which fluctuated along a larger range. With regard to context, by only running the models on a specific context dataset, power usage and temperature decreased for both VoxelNet and AVOD.

4.2. DO OBJECT DETECTION MODELS PERFORM BETTER IN DIFFERENT CONTEXTS?

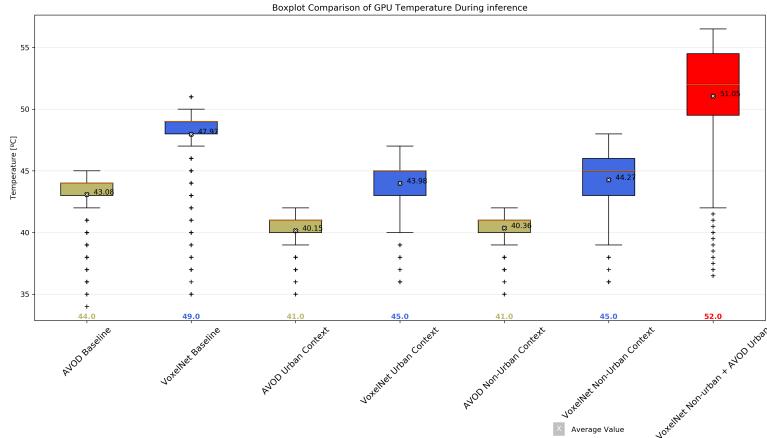


Figure 41: Temperature boxplots of evaluation runs. Voxelnet(blue)
AVOD(gold) AVOD+VoxelNet(Red)

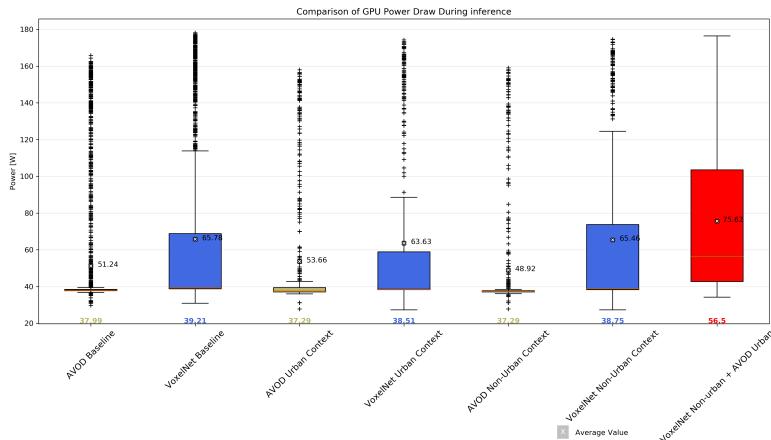


Figure 42: Power boxplots of evaluation runs. Voxelnet(blue),
AVOD(gold), AVOD+VoxelNet(Red)

Running Both Models Simultaneously

Having identified that AVOD is better at detecting cars and pedestrians in urban contexts and VoxelNet in non-urban contexts, I set out to further investigate whether it is possible to have both models running on a single GPU.

Table 44: Memory

	Min	Max
Memory(MB)	4608.50	14320.00

Analysis

Temperature - Running both models simultaneously resulted in significantly higher temperatures as compared running them exclusively. The average temperature was 7°hotter than the maximum mean temperature from running the models exclusively.

Power - Similarly, the mean power usage was 10 watts greater than the maximum mean power from running the models exclusively.

Memory - Running both models resulted in a total memory usage of 14.3 GB out of the total memory of 16GB thereby proving that they are able to run optimally on a GPU with 16GB+ memory.

4.3 Is VoxelNet Performance Valid For Different Datasets?

Metrics

Given that the data was collected from a stationary vehicle in an urban area with mostly pedestrians, VoxelNet was evaluated on the pedestrian class only. However, as there was no corresponding image data for the point clouds, only the 3D bounding boxes and Bird-Eye View bounding boxes with $\text{IoU} \geq 0.5$. were used in calculating the AP.

Analysis

VoxelNet was not able to predict any objects on the SIL dataset. As highlighted before, the VLP-16 LiDAR sensor has a lower frequency than that of the HDL-64E. This therefore affected the point cloud density and is visible in figure 43 . SIL point clouds had an average point cloud density per frame that was lower than that of KITTI's by 42% in 1000 samples as seen in figure 44. This was evaluated by calculating the number of points in each point cloud frame and dividing by the frame size. It is therefore safe to conclude that the point cloud density has a positive correlation to the performance of LiDAR only models.

4.3. IS VOXELNET PERFORMANCE VALID FOR DIFFERENT DATASETS?

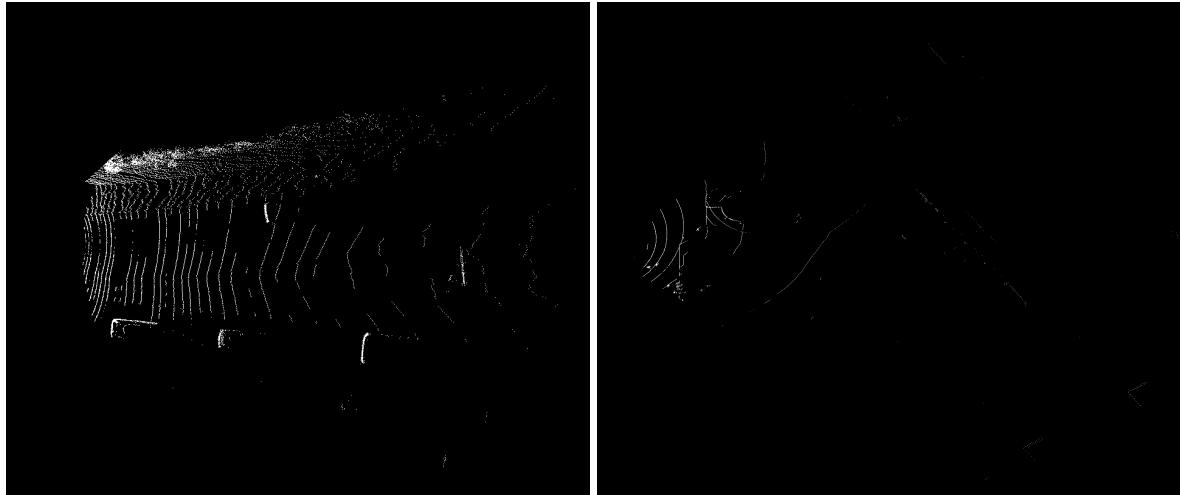


Figure 41: *

KITTI bird eye view

Figure 42: *

SIL bird eye view

Figure 43: Bird Eye View Difference between SIL and KITTI

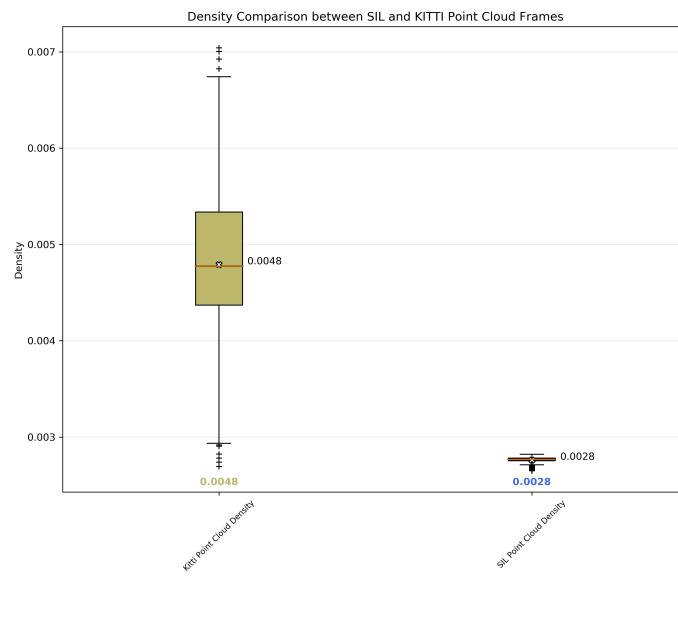


Figure 44: Point Cloud Density

CHAPTER



DISCUSSION AND FUTURE WORK

Following the analysis, we have identified that indeed it is possible to detect and characterise images with relatively high accuracy and point clouds to a much lesser extent. In addition, we have determined that multimodal methods perform better in dynamic environments as compared to LiDAR only models. Finally, we were able to realise the impact of using point clouds from different LiDAR sensors on LiDAR only models.

In this chapter, I wish to further discuss the implications of these findings on the development and production of AVs. To conclude, I will discuss how this work can be further improved.

Despite seeing the higher temperature and power usage, this occurred when both models were running simultaneously. However this would not be the case in an actual implementation as one model will only be used while the other one will be inactive. This can be implemented using a pipeline system whereby both Tensorflow graphs can be loaded on a single GPU, however, depending on the context, the input can be fed into either graph while the other one remains dormant and vice versa. Therefore the usage would be much lesser than actually having running both.

multimodal for small object detection common in urban lidar for the non-urban context
using LiDARs with more points/sec or standardise models to use a specific LiDAR

5.1 Future Work

Context Detection models using rnn End to End Pipeline model

CHAPTER



CONCLUSION

Internationally Agreed Best Practices VLP-128 Solid state lidar

A P P E N D I X



APPENDIX A

LiDAR	Hor FOV	Ver FOV	Range	Angular Resolution	Points/second	Channels
VLS-128	360°	+15°to -25°	300m	0.11°	~9.6 Million	128
HDL-64E	360°	26.9°	120m	~0.4°	~2.2 Million	64
HDL-32E	360°	+10°to -30°	80m-100m	0.1°	~1.39 Million	32
VLP-32C	360°	+15°to -25°	200m	0.1°	~1.2 Million	32
VLP-16	360°	± 15°	100m	0.1°	600,000	16

Table A1: Velodyne LiDAR Family

```

1 #Code to read .pcap files and convert them to .pcd
2
3 #Load existing pcap file and create a sensor publisher
4 $ roslaunch velodyne_pointcloud VLP-16-points.launch pcap:=(path to .pcap file)
5
6 #Receiver to convert sensor messages to .pcd (In separate terminal window)
7 $ rosrun pcl_ros pointcloud_to_pcd input:=/velodyne_points _prefix:=(path to save .pcd
     file)
8

```

Listing A.1: .pcap -> .pcd

BIBLIOGRAPHY

- [1] *Artificial intelligence for construction productivity*, 2018.
- [2] H. BAY, T. TUYTELAARS, AND L. VAN GOOL, *Surf: Speeded up robust features*, in European conference on computer vision, Springer, 2006, pp. 404–417.
- [3] I. BROWN, *Human factors: the journal of the human factors and ergonomics society*, Driver Fatigue, 36 (1994), pp. 298–314.
- [4] S. D. BROWN, R. J. FRANCIS, J. ROSE, AND Z. G. VRANESIC, *Field-programmable gate arrays*, vol. 180, Springer Science & Business Media, 2012.
- [5] X. CHEN, K. KUNDU, Z. ZHANG, H. MA, S. FIDLER, AND R. URTASUN, *Monocular 3d object detection for autonomous driving*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 2147–2156.
- [6] X. CHEN, K. KUNDU, Y. ZHU, H. MA, S. FIDLER, AND R. URTASUN, *3d object proposals using stereo imagery for accurate object class detection*, IEEE transactions on pattern analysis and machine intelligence, 40 (2018), pp. 1259–1272.
- [7] X. CHEN, H. MA, J. WAN, B. LI, AND T. XIA, *Multi-view 3d object detection network for autonomous driving*, in IEEE CVPR, vol. 1, 2017, p. 3.
- [8] N. DALAL AND B. TRIGGS, *Histograms of oriented gradients for human detection*, in Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on, vol. 1, IEEE, 2005, pp. 886–893.
- [9] F. DEPARTMENT FOR ENVIRONMENT AND R. AFFAIRS, *Urban and rural area definitions for policy purposes in england and wales: Methodology (v1.0)*, GOV.UK, (2018).
- [10] A. GEIGER, P. LENZ, AND R. URTASUN, *Are we ready for autonomous driving? the kitti vision benchmark suite*, in Conference on Computer Vision and Pattern Recognition (CVPR), 2012.
- [11] R. GIRSHICK, *Fast r-cnn*, in International Conference on Computer Vision (ICCV), 2015.

BIBLIOGRAPHY

- [12] INTEL, 2018.
<https://www.intel.co.uk/content/www/uk/en/automotive/go-automated-driving.html>.
- [13] J. KU, M. MOZIFIAN, J. LEE, A. HARAKEH, AND S. WASLANDER, *Joint 3d proposal generation and object detection from view aggregation*, arXiv preprint arXiv:1712.02294, (2017).
- [14] B. LI, T. ZHANG, AND T. XIA, *Vehicle detection from 3d lidar using fully convolutional network*, arXiv preprint arXiv:1608.07916, (2016).
- [15] S.-C. LIN, Y. ZHANG, C.-H. HSU, M. SKACH, M. E. HAQUE, L. TANG, AND J. MARS, *The architectural implications of autonomous driving: Constraints and acceleration*, in Proceedings of the Twenty-Third International Conference on Architectural Support for Programming Languages and Operating Systems, ACM, 2018, pp. 751–766.
- [16] T.-Y. LIN, P. GOYAL, R. GIRSHICK, K. HE, AND P. DOLLÁR, *Focal loss for dense object detection*, IEEE transactions on pattern analysis and machine intelligence, (2018).
- [17] D. G. LOWE, *Object recognition from local scale-invariant features*, in Computer vision, 1999. The proceedings of the seventh IEEE international conference on, vol. 2, Ieee, 1999, pp. 1150–1157.
- [18] L. MEDSKER AND L. JAIN, *Recurrent neural networks*, Design and Applications, 5 (2001).
- [19] MOBILEYE, *Enabling autonomous - mobileeye*, 2018.
<https://www.mobileye.com/future-of-mobility/mobileye-enabling-autonomous/>.
- [20] A. NEWELL AND S. K. CARD, *The prospects for psychological science in human-computer interaction*, Human-computer interaction, 1 (1985), pp. 209–242.
- [21] L. PLOTKIN, *Pydriver*, bachelor's thesis (Studienarbeit), Karlsruhe Institute of Technology, Germany, Mar. 2015.
- [22] D. POMERLEAU, *Rapidly adapting artificial neural networks for autonomous navigation*, in Advances in neural information processing systems, 1991, pp. 429–435.
- [23] D. A. POMERLEAU, *Alvinn: An autonomous land vehicle in a neural network*, in Advances in neural information processing systems, 1989, pp. 305–313.
- [24] C. R. QI, W. LIU, C. WU, H. SU, AND L. J. GUIBAS, *Frustum pointnets for 3d object detection from rgbd data*, arXiv preprint arXiv:1711.08488, (2017).
- [25] C. R. QI, H. SU, K. MO, AND L. J. GUIBAS, *net: Deep learning on point sets for 3d classification and segmentation*, Proc. Computer Vision and Pattern Recognition (CVPR), IEEE, 1 (2017), p. 4.

- [26] C. R. QI, L. YI, H. SU, AND L. J. GUIBAS, *Pointnet++: Deep hierarchical feature learning on point sets in a metric space*, in Advances in Neural Information Processing Systems, 2017, pp. 5105–5114.
- [27] S. REN, K. HE, R. GIRSHICK, AND J. SUN, *Faster r-cnn: Towards real-time object detection with region proposal networks*, in Advances in neural information processing systems, 2015, pp. 91–99.
- [28] N. RESEARCH, *Navigant research leaderboard: Automated driving vehicles*, Jan 2018.
<https://www.navigantresearch.com/research/navigant-research-leaderboard-automated-driving-vehicles>.
- [29] R. B. RUSU, N. BLODOW, AND M. BEETZ, *Fast point feature histograms (fpfh) for 3d registration*, in Robotics and Automation, 2009. ICRA'09. IEEE International Conference on, Citeseer, 2009, pp. 3212–3217.
- [30] S. SALTI, F. TOMBARI, AND L. DI STEFANO, *Shot: Unique signatures of histograms for surface and texture description*, Computer Vision and Image Understanding, 125 (2014), pp. 251–264.
- [31] J. SHOTTON, A. FITZGIBBON, M. COOK, T. SHARP, M. FINOCCHIO, R. MOORE, A. KIPMAN, AND A. BLAKE, *Real-time human pose recognition in parts from single depth images*, in Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on, Ieee, 2011, pp. 1297–1304.
- [32] J. SHOTTON, M. JOHNSON, AND R. CIPOLLA, *Semantic texton forests for image categorization and segmentation*, in Computer vision and pattern recognition, 2008. CVPR 2008. IEEE Conference on, IEEE, 2008, pp. 1–8.
- [33] M. J. S. SMITH, *Application-specific integrated circuits*, vol. 7, Addison-Wesley Reading, MA, 1997.
- [34] B. STEDER, R. B. RUSU, K. KONOLIGE, AND W. BURGARD, *Narf: 3d range image features for object recognition*, in Workshop on Defining and Solving Realistic Perception Problems in Personal Robotics at the IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS), vol. 44, 2010.
- [35] P. VIOLA AND M. JONES, *Rapid object detection using a boosted cascade of simple features*, in Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on, vol. 1, IEEE, 2001, pp. I–I.
- [36] WAYMO, *Waymo safety report*, 2018.
<https://waymo.com/safety/>.

BIBLIOGRAPHY

- [37] Y. ZHONG, *Intrinsic shape signatures: A shape descriptor for 3d object recognition*, in Computer Vision Workshops (ICCV Workshops), 2009 IEEE 12th International Conference on, IEEE, 2009, pp. 689–696.
- [38] Y. ZHOU AND O. TUZEL, *Voxelnet: End-to-end learning for point cloud based 3d object detection*, arXiv preprint arXiv:1711.06396, (2017).