



Faculty of Electronics and Computer Engineering

A large, faded version of the UTeM logo is centered in the background behind the title text.

**PERFORMANCE EVALUATION ON QUANTIZED WEIGHT FOR
CONVOLUTIONAL NEURAL NETWORK BASED OBJECT
DETECTION**

Mohd Hasbullah bin Putra

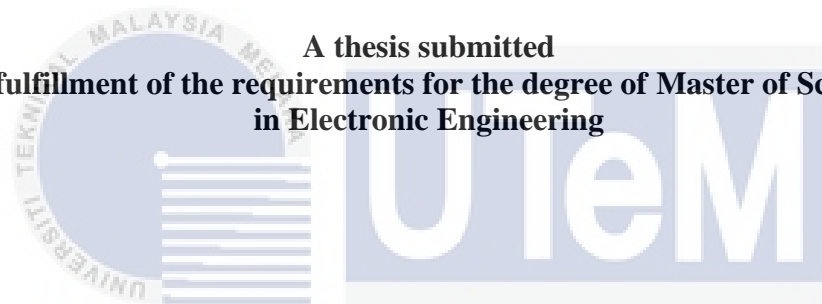
Master of Science in Electronic Engineering

2021

**PERFORMANCE EVALUATION ON QUANTIZED WEIGHT FOR
CONVOLUTIONAL NEURAL NETWORK BASED OBJECT DETECTION**

MOHD HASBULLAH BIN PUTRA

**A thesis submitted
in fulfillment of the requirements for the degree of Master of Science
in Electronic Engineering**



اونيورسيتي تيكنيكل مليسيا ملاك

UNIVERSITY OF ELECTRONICS AND COMPUTER ENGINEERING

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

2021

DECLARATION

I declare that this thesis entitled “Performance Evaluation on Quantized Weight for Convolutional Neural Network based Object Detection” is the result of my own research except as cited in the references. The thesis has not been accepted for any degree and is not concurrently submitted in candidature of any other degree.



Signature :

Name : MOHD HASBULLAH BIN PUTRA

Date :

اونيورسيتي تيكنيكل مليسيا ملاك

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

APPROVAL

I hereby declare that I have read this thesis and in my opinion this thesis is sufficient in terms of scope and quality for the award of Master of Science in Electronic Engineering.

| | | |
|--|-----------------|---|
|  | Signature | |
| | Supervisor Name | :PROFESSOR DR. ZULKALNAIN BIN MOHD YUSSOF |
| | Date | اونيومر سیتی تیکنیکل ملیہ ملاک |

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

DEDICATION

To my beloved parents

“Life is like a bicycle, to keep the balance we need to keep going”



Abstract

A Convolutional Neural Network (CNN) based object detection is an emerging topic in the image processing field and has become the state-of-the-art in computer vision and machine learning. The traditional system in object detection uses a handcrafted feature extractor which is less robust in applications. By applying the CNN approach in the field, the accuracy of object detections can increase significantly. However, the use of deep CNN architecture model leads to high computation. In this research, a real-time CNN based object detection system is presented. The system is designed based on the modified You Only Look Once (modified-YOLO) architecture which is constructed with only 7 CNN layers. The grid cell parameter value of the system is varied to evaluate its effectiveness and ability in detecting small size objects upon deployment. The experimental results demonstrate that even with 7 convolutional layers, modified-YOLO can provide good detection accuracy and real-time operation achieving the best miss rate (*MR*) of 22.7% *MR*. Although the scores show an increase in the *MR*, the visual qualitative evaluation using randomly captured images indicate that the 7 layers modified-YOLO architecture with 11x11 grid cells can correctly and easily detect small objects. This makes the modified YOLO architecture which has been reduced in terms of complexity a suitable candidate for use in real-time operation. In order to further reduce the complexity of the CNN system, the trained floating-point weights are quantized. Three types of scalar quantization are used to quantize the CNN weights namely symmetric uniform quantizer, asymmetric uniform quantizer and non-uniform quantizer designed using *k*-means algorithm. The quantization reduces the storage and computation requirements. The quantitative results using the *MR* standard metric indicate that the non-uniform quantizer provides the best results compared to the other quantization methods. Using 6-bit precision non-uniformly quantized weights yields detection performance comparable to the CNN network using floating-point weights. Additionally, based on the qualitative results, the CNN network with 4-bit non-uniform quantization weights is able to detect the person objects correctly.

PENILAIAN PRESTASI PEMBERAT YANG DIKUANTUMKAN UNTUK PENGESAN OBJEK BERASASKAN RANGKAIAN NEURAL PELINGKARAN

ABSTRAK

Pengesan objek berasaskan Talian Neural Pelingkaran (CNN) adalah satu topik yang berkembang di dalam bidang pemprosesan imej dan telah menjadi kaedah terkini dalam bidang penglihatan komputer dan pembelajaran mesin. Kebanyakan sistem tradisional dalam pengesanan objek menggunakan kaedah pengekstrak ciri yang tetap dan kurang berkesan ketika digunakan. Dengan menggunakan kaedah CNN, ketepatan pengesan objek telah meningkat dengan signifikan. Namun, seni bina CNN yang dalam telah menjurus kepada peningkatan pengiraan kompleks di dalam sistem. Dalam kajian ini, sebuah pengesan objek masa sebenar berasaskan CNN telahpun dibentangkan. Sistem tersebut direka berasaskan pengubahsuaian senibina "You Only Look Once" (modified-YOLO) yang terdiri daripada 7 lapisan CNN. Nilai parameter sel grid di dalam sistem tersebut juga telah dimanipulasi untuk menilai keberkesanan dan kemampuan sistem dalam mengesan objek yang kecil ketika digunakan. Keputusan kajian menunjukkan bahawa walaupun dengan hanya 7 lapisan CNN yang digunakan, "modified-YOLO" mampu memberi keputusan yang tepat dalam masa sebenar dengan pencapaian kadar sasaran (MR) terbaik sebanyak 22.7%MR. Walaupun keputusan menunjukkan peningkatan MR, hasil penilaian kualitatif visual dengan menggunakan imej rawak yang ditangkap menunjukkan bahawa senibina YOLO dengan 7 lapisan CNN dengan parameter sel grid 11x11 mampu untuk mengesan objek kecil dengan tepat dengan mudah. Ini menjadikan pengubahsuaian senibina YOLO yang telah dikurangkan pengiraan kompleks sebagai calon yang sesuai untuk kegunaan masa sebenar. Untuk meneruskan pengurangan pengiraan kompleks dalam sistem, satu simulasi untuk menilai prestasi titik tetap CNN telah dijalankan dengan mengkuantumkan pengubahsuaian senibina YOLO. Tiga kaedah pengkuantuman skalar telah digunakan untuk mengkuantum pemberat CNN iaitu pengkuantuman seragam simetri, pengkuantuman seragam tidak simetri dan pengkuantuman tidak seragam yang menggunakan algoritma k-means. Pengkuantuman memberi kelebihan dengan mengurangkan pengiraan kompleks dan pada masa yang sama mengurangkan penggunaan memori. Keputusan pengkuantuman dianalisis secara kuantitatif menggunakan metrik standard MR menunjukkan pengkuantuman tidak seragam menghasilkan keputusan yang terbaik berbanding kaedah yang lain. Dengan hanya menggunakan kejituan 6-bit, kaedah ini menghasilkan pengesan yang setanding dengan titik apung CNN. Selain itu, berdasarkan keputusan kualitatif, rangkaian CNN yang sama mampu untuk mengesan objek dengan tepat.

ACKNOWLEDGEMENTS

First and foremost, I would like to take this opportunity to express my sincere acknowledgement to my supervisor Professor Dr. Zulkalnain Bin Mohd Yussof from the Faculty of Electronics and Computer Engineering Universiti Teknikal Malaysia Melaka (UTeM) for his essential supervision, support, and encouragement towards the completion of this thesis.

I would also like to express my greatest gratitude to Dr. Sani Irwan bin Salim from Faculty of Electronics and Computer Engineering, co-supervisor of this project for his advice and suggestions in evaluation. Special thanks to Associate Professor Dr. Lim Kim Chuan, project leader for CREST grant for the financial support throughout this project.

Special thanks to my parents who consistently encourage me to finish this research work and become the main reason for me to keep on going. Thanks to all my friends, especially those who are in MLSP lab and Post Graduate lab for all the knowledge sharing and joy throughout the journey.



TABLE OF CONTENTS

| | PAGE |
|---|-------------|
| DECLARATION | |
| APPROVAL | |
| DEDICATION | |
| ABSTRACT | i |
| ABSTRAK | ii |
| ACKNOWLEDGEMENTS | iii |
| TABLE OF CONTENTS | iv |
| LIST OF TABLES | vi |
| LIST OF FIGURES | vii |
| LIST OF SYMBOLS | x |
| LIST OF PUBLICATIONS | xi |
| | |
| CHAPTER | |
| 1. INTRODUCTION | 1 |
| 1.1 Research background | 1 |
| 1.2 Problem statement | 3 |
| 1.3 Objectives | 5 |
| 1.4 Scopes | 5 |
| 1.5 Research work contribution | 5 |
| 1.6 Thesis outlines | 7 |
| | |
| 2. LITERATURE REVIEW | 8 |
| 2.1 Introduction | 8 |
| 2.2 Artificial Neural Network | 8 |
| 2.3 Basic CNN architecture | 12 |
| 2.4 Region with CNN features (R-CNN) | 13 |
| 2.5 Fast Region-Based Convolutional Neural Network (Fast R-CNN) | 15 |
| 2.6 Faster R-CNN | 16 |
| 2.7 You Only Look Once (YOLO) | 18 |
| 2.8 YOLO 9000 | 19 |
| 2.9 Popular and commonly used neural network architectures | 20 |
| 2.9.1 AlexNet | 21 |
| 2.9.2 Visual Geometry Group (VGG) | 21 |
| 2.9.3 GoogleLeNet | 22 |
| 2.9.4 ResNet | 24 |
| 2.9.5 ReNet | 25 |
| 2.9.6 SqueezeNet | 26 |
| 2.10 Neural network training | 28 |
| 2.10.1 Supervised learning | 29 |
| 2.10.2 Unsupervised learning | 29 |

| | | |
|-----------|--|------------|
| 2.10.3 | Reinforcement learning | 30 |
| 2.11 | Batch lormalization | 31 |
| 2.12 | Quantization | 32 |
| 2.13 | <i>k</i> -means clustering algorithm | 34 |
| 2.14 | Low rank decomposition method | 35 |
| 2.15 | Pruning | 36 |
| 2.16 | Summary | 38 |
| 3. | METHODOLOGY | 40 |
| 3.1 | Introduction | 40 |
| 3.2 | Project flow | 40 |
| 3.3 | Modified-YOLO network architecture | 42 |
| 3.3.1 | Convolutional layer | 44 |
| 3.3.2 | Pooling layer | 46 |
| 3.3.3 | Fully connected layer | 47 |
| 3.3.4 | Activation layer | 49 |
| 3.4 | The datasets | 51 |
| 3.4.1 | CALTECH datasets | 52 |
| 3.4.2 | INRIA datasets | 53 |
| 3.5 | Annotation | 54 |
| 3.6 | Training | 55 |
| 3.7 | Inference | 58 |
| 3.8 | Loss function | 61 |
| 3.9 | Evaluation tools | 63 |
| 3.10 | Quantization | 65 |
| 3.10.1 | Symetric uniform quantization | 66 |
| 3.10.2 | Asymmetric uniform quantization | 68 |
| 3.10.3 | Non-uniform quantization | 69 |
| 3.11 | Weight compression | 70 |
| 3.12 | Summary | 72 |
| 4. | RESULT AND DISCUSSION | 73 |
| 4.1 | Introduction | 73 |
| 4.2 | Floating Point CNN based person detection | 73 |
| 4.3 | CNN based floating point person and car detection | 80 |
| 4.4 | CNN based person detection using quantized weights | 86 |
| 4.4.1 | Bias and weights distribution | 86 |
| 4.4.2 | Symetric uniform quantization | 88 |
| 4.4.3 | Asymmetric uniform quantization | 90 |
| 4.4.4 | Non-uniform quantization | 92 |
| 4.5 | Weight Compression | 94 |
| 4.6 | Summary | 97 |
| 5. | CONCLUSION AND RECOMMENDATIONS OF FUTURE WORK | 98 |
| 5.1 | Conclusion | 98 |
| 5.2 | Future works and improvements | 99 |
| | REFERENCES | 101 |

LIST OF TABLES

| TABLE | TITLE | PAGE |
|-------|--|------|
| 2.1 | SqueezeNet network architecture dimensions | 28 |
| 2.2 | CNN comparisons | 38 |
| 2.3 | Commonly used Neural Network architecture | 39 |
| 4.1 | Comparisons of person detection results on the CALTECH dataset | 79 |
| 4.2 | mAP, frame per second and AP performance | 85 |
| 4.3 | Miss rate at 0.1 FPPI for uniform quantized weights | 90 |
| 4.4 | Miss rate at 0.1 FPPI for non-uniform quantized weights | 92 |
| 4.5 | Miss rate at 0.1 FPPI for k -means quantized weights | 94 |
| 4.6 | Number of total weights (6-bit compression) | 95 |
| 4.7 | Number of total bits and size after quantization | 95 |

LIST OF FIGURES

| FIGURE | TITLE | PAGE |
|--------|---|------|
| 1.1 | Selective search in detecting multiple sizes and location objects | 4 |
| 2.1 | Example of biological neuron | 9 |
| 2.2 | An Artificial Neural Network basic design | 10 |
| 2.3 | (a) Feed forward. (b) Back propagation on a perceptron | 11 |
| 2.4 | The CNN model proposed for handwritten zip code recognition | 12 |
| 2.5 | LeNet-5 CNN architecture for handwritten digits recognition | 13 |
| 2.6 | Region Convolutional Neural Network for object detection | 14 |
| 2.7 | ROI Pooling | 15 |
| 2.8 | Fast R-CNN architecture with the unified training | 16 |
| 2.9 | Anchor boxes Scoring | 17 |
| 2.10 | YOLO Original network architecture | 19 |
| 2.11 | AlexNet CNN architecture | 21 |
| 2.12 | VGG-16 CNN architecture | 22 |
| 2.13 | Inception module with dimensionality reduction from the GoogLeNet architecture | 23 |
| 2.14 | Explanation for residual block from the ResNet architecture | 24 |
| 2.15 | One layer of ReNet architecture. The layer is modeling horizontal and vertical spatial dependence | 26 |
| 2.16 | SqueezeNet with 1×1 kernel | 27 |
| 2.17 | Markov Decision Process | 31 |
| 2.18 | A uniform quantizer | 33 |
| 2.19 | Channel and filter wise pruning | 37 |
| 3.1 | Process flowchart for research work | 41 |
| 3.2 | Output tensor form modified-YOLO | 43 |
| 3.3 | Unified detection of modified-YOLO | 43 |

| | | |
|------|--|----|
| 3.4 | The modified-YOLO network architecture | 44 |
| 3.5 | Generation of feature maps through convolution process with stride | 45 |
| 3.6 | The zero padding around the input array | 46 |
| 3.7 | Down sampling using max-pooling | 47 |
| 3.8 | A perceptron in the fully connected layer of modified YOLO | 48 |
| 3.9 | ReLU activation function | 50 |
| 3.10 | Leaky ReLu activation function | 51 |
| 3.11 | Example of overfitting problem | 52 |
| 3.12 | Examples of CALTECH dataset images (a) images of pedestrian without occlusion (b) Image of the pedestrian with occlusions | 53 |
| 3.13 | Examples of images in INRIA datasets | 54 |
| 3.14 | Bounding box ground truth labeling | 55 |
| 3.15 | Modified YOLO prediction representation | 56 |
| 3.16 | Batch mode training process | 57 |
| 3.17 | Methods to calculate the IOU for training error detection | 58 |
| 3.18 | Flowchart for inference | 59 |
| 3.19 | Image is divided into $S \times S$ grid cells | 60 |
| 3.20 | Each grid cells will predict 2 bounding boxes | 60 |
| 3.21 | The grid cell will decide whether the box contain any object for detection | 61 |
| 3.22 | IOU accuracy and training loss | 63 |
| 3.23 | The example of Hoiem evaluation | 64 |
| 3.24 | Fixed point number format representation | 67 |
| 3.25 | Bins generated for uniform quantization | 67 |
| 3.26 | Bins generated for asymmetric uniform quantization | 69 |
| 3.27 | Centroids generated for non-uniform quantization | 70 |
| 3.28 | Weights and bias of the YOLO network | 70 |
| 4.1 | (a) ground truth annotations (b) M-YOLO_7x7 person detection results compared with the ground truth annotations | 74 |
| 4.2 | (a) ground truth annotations (b) M-YOLO_9x9 person detection results | 75 |
| 4.3 | (a) ground truth annotations (b) M-YOLO_11x11 person detection results | 76 |
| 4.4 | M-YOLO_11x11 tested with a random image | 77 |
| 4.5 | M-YOLO_9x9 tested with a random image | 78 |
| 4.6 | M-YOLO_7x7 tested with a random image | 78 |
| 4.7 | Graph of FPPI against MR | 79 |

| | | |
|------|---|----|
| 4.8 | The detection result with image captured from a car dashboard camera | 80 |
| 4.9 | Car and person detection result from INRIA testing dataset | 81 |
| 4.10 | Car and person detection result from a street view captured image | 81 |
| 4.11 | Detection result with 7x7 grid cells | 82 |
| 4.12 | Detection result with 9x9 grid cells | 83 |
| 4.13 | Detection result with 11x11 grid cell | 83 |
| 4.14 | Error analysis chart to show the percentage of localization and background errors of YOLO | 84 |
| 4.15 | Error analysis chart to show the percentage of localization and background errors of M-YOLO_11x11 | 84 |
| 4.16 | Histogram showing the bias distribution in every layer | 87 |
| 4.17 | Histogram showing the weights distribution in every layer | 88 |
| 4.18 | Uniform quantized weights detection results, (a) Floating-point, (b) Q(1,8), (c) Q(1,7), (d) Q(1,6) | 89 |
| 4.19 | FPPI vs MR of uniform quantization weights | 89 |
| 4.20 | Asymmetric uniform quantized weights detection results, (a) Floating-point, (b) 256 bins (8-bit), (c) 128 bins (7-bit), (d) 64 bins (6-bit) | 91 |
| 4.21 | FPPI vs MR of non-uniform quantization weights | 91 |
| 4.22 | Non- uniform quantized weights detection results | 93 |
| 4.23 | FPPI vs MR of non-uniform quantization weights | 93 |
| 4.24 | Total number of compressed bits and the weights file size | 96 |

LIST OF SYMBOLS

| | | |
|-----|---|--------------------------|
| c | - | Class |
| h | - | Height |
| w | - | Width |
| S | - | Number of Grid Cells |
| B | - | Number of Bounding Boxes |



LIST OF PUBLICATIONS

1. Putra, M.H., Yussof, Z.M., Lim, K.C. and Salim, S.I., 2018. Convolutional neural network for person and car detection using yolo framework. *Journal of Telecommunication, Electronic and Computer Engineering (JTEC)*, 10(1-7), pp.67-71.
2. Putra, M.H., Yussof, Z.M., Salim, S.I. and Lim, K.C., 2017. Convolutional Neural Network for Person Detection using YOLO Framework. *Journal of Telecommunication, Electronic and Computer Engineering (JTEC)*, 9(2-13), pp.1-5.



CHAPTER 1

INTRODUCTION

1.1 Research background

Vision-based object detection is a hot research topic among the computer vision research community. In particular, the person and vehicle detection have a direct application in Advanced Driver Assistance System (ADAS), Intelligent Vehicle and Visual Surveillance System. Various methods have been proposed for person, car detection, or general object detection. However, majority of the techniques proposed to focus on achieving high detection accuracy at the expense of high computational complexity. Hence many of these methods are not suitable for real-time applications such as ADAS, which requires both high detection accuracy and high frame per second (FPS) performance.

Before the emergence of the Convolutional Neural Network (CNN), the Deformable Part Model (DPM) using handcrafted features such as Histogram of Oriented Gradient (HOG) has been the state-of-the-art object detector for many years. Inspired by the impressive performance demonstrated in image classification, CNN has been applied to object detection and achieves impressive results. Most notably, Girshick et. al. (2014) proposed the Regions with Convolutional Neural Network (R-CNN) framework for object detection and demonstrated state-of-the-art performance on standard detection benchmarks, for example in PASCAL VOC, with a large margin over the previous arts, which are mostly DPM based.

R-CNN uses a handcrafted selective search algorithm to generate object proposals and CNN classifier for detection tasks. The R-CNN is however computationally expensive

due to the forward pass computation required for each proposal. Girshick (2015) then proposed Fast Region-based Convolutional Network (Fast R-CNN), which reduces computational complexity by sharing convolutional features and pooling object proposals from the last convolutional layer. While Fast R-CNN achieves excellent detection accuracy, its speed is still limited by the bottleneck due to the object proposal generation. Then Girshick et. al. come out with series of proposal to overcome the intensive computation of CNN. Although the series of improvements introduced by the authors have achieved excellent object detection accuracy, it is computationally intensive and not suitable for use in real-time applications.

In order to meet the combined requirement of high object detection accuracy and real-time operation, a different approach of CNN-based object detection named You Only Look Once (YOLO) was proposed by Redmon et. al. (2016). In contrast to region proposal-based object detection algorithms such as Fast R-CNN, YOLO CNN-based algorithm predicts bounding boxes and class probabilities directly from full images in a single evaluation. Instead of having a consecutive pipeline of region proposals and object classification, YOLO implements a different approach by treating the overall input as a single regression problem. This reduces the latency at the expense of accuracy. The full 27 layers YOLO encounters a problem where the network is struggling to detect small sized objects.

However, all of the impressive performances are achieved by utilizing many layers and millions of parameters which leads to high computational and memory storage requirements. For example, the 8-layer AlexNet requires 244 MB of parameters and 1.4 billion floating point operations (GFLOP) to classify a single image, while the deeper CNN VGG-16 requires 552 MB of weight storage and 30.8 GFLOP per image Han et. al. (2015). This prohibits the deployment of deep CNNs on resource-limited embedded devices for example for mobile and portable devices especially involving real-time applications.

Embedded devices typically have constraints in terms of memory, energy and computation resource. Thus, in order to deploy deep neural networks to embedded platforms, some schemes should be applied to reduce the computation and storage requirements while preserving the deep CNN performance.

1.2 Problem statement

The early development of object detection methods used the selective search, which was deployed in most of the DPM based models and some of the early designed CNN models. The selective search method will slide through the whole image pixel by pixel. However, sliding through pixels requires lots of computational power and this process causes bottleneck towards the whole process of computation. Furthermore, the objects to be detected can appear at any location and scale within the input image. Figure 1.1 depicts the example of selective search in detecting the object. In addition, the object boundaries might be less clear compared with other objects that exist in the image. Therefore, all object scales have to be considered in selective search (Uijlings, 2012). Thus, the whole system speed performance will be slow and not suitable to be implemented in a real time detection system. In order to overcome this problem, the selective search method is replaced with another method called the unified detection which is introduced by YOLO and suitable for the real time application. In addition, YOLO network architecture was further modified to enable the model to be deployed in a low power platform.

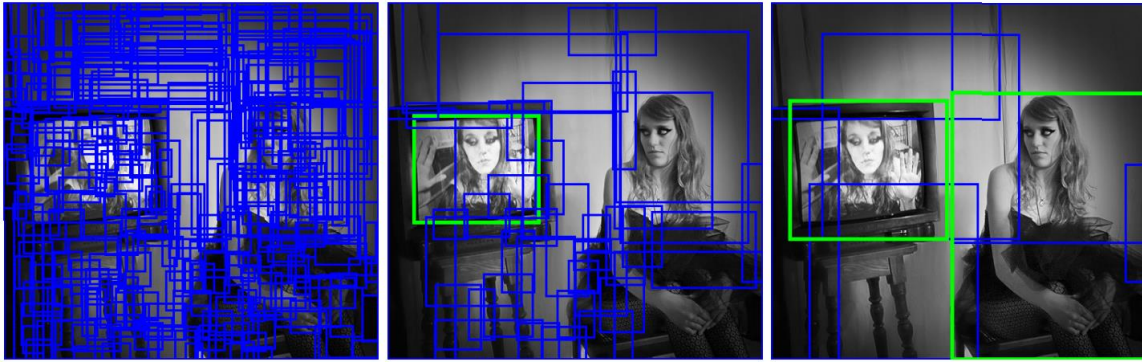


Figure 1.1: Selective search in detecting multiple sizes and location objects

On the other hand, the proposed CNN model YOLO struggles to detect small sized object which appears in the input image. Hypothetically, to remedy the problem, the grid cell parameter in the model is varied. This is to observe whether detection result towards small object can be improved. Theoretically, by increasing the grid cell parameter value should give advantages to detect small objects, but it will also increase the computation complexity of the system. Thus, in this research work, the optimum value of the grid cells parameter will be determined.

Apart from that, the existing CNN detection systems are computationally expensive and not suitable to be deployed on low powered platform such as FPGAs. Besides, the memory consumption of a CNN system is commonly very large. Therefore, in this research work, the quantization of the CNN weights is introduced, and the performance results of CNN using the quantized weights are compared to the floating-point CNN detection system. This is to ensure that the quantized CNN detection accuracy is still comparable after the quantization process.

1.3 Objectives

The objectives of this research work are as follows:

1. To design an object detection system using convolutional neural network.
2. To evaluate the performance of the designed object detection with other detection systems.
3. To compress the convolutional neural network weights using quantization schemes.
4. To evaluate the performance of the convolutional neural network using quantized weights and compare its performance with the floating point CNN object detection system.

1.4 Scopes

The scopes of this research work are as follows:

1. The input image for both training and testing are resized to 448 x 448 pixels.
2. The detection system detects persons for single object detection and combination of persons with car for multiple object detection.
3. The evaluation and datasets selection is based on ADAS standard.
4. The designed network architecture is trained to detect specific types of objects related to the driving assistance based on ADAS standard.
5. Only scalar quantization (SQ) will be used for weights quantization.

1.5 Research work contribution

In this research, the object detection system has been designed. The network architecture was inspired by Darknet and the network is compatible with the YOLO framework. Typically, a network for CNN will have a large number of layers, which is called Deep Neural Network (DNN). Having a greater number of layers will enhance the detection

capabilities of the CNN but will increase the computational power to run the system. Therefore, the proposed network architecture is designed to have less convolutional layers in order to ensure the network is capable to run in real-time. By reducing the number of layers, the expensive computations are reduced, and the FPS performance of the system can be boosted. Additionally, during the training phase of CNN, weights will require a lot of computation and can consume up to several days of training. To speed up the process, hardware accelerators like CUDA GPU is being utilized to further accelerate the training speed.

To enhance the accuracy of the detection, the grid cell parameter is varied in order to achieve optimum accuracy for the detection system. After the optimum parameter has been achieved, the yielded weights are then compressed using several quantization methods namely uniform quantization, asymmetric uniform quantization, and non-uniform quantization. The goal of weight compression is to reduce the network storage requirement while minimizing performance loss due to quantization.

Several network model compression schemes have been proposed in the literature review including network pruning, quantization, and low-rank matrix decomposition. The well-known fact partly motivates these previous works that the deep CNN network weights have significant redundancies as suggested by Cheng et. al. (2015). In our work, three different scalar quantization (SQ) techniques are used to quantize the network weights that are pre-trained with floating-point precision. No re-training is performed on the quantized network. The three SQ schemes applied are symmetric uniform quantizer, asymmetric uniform quantizer and non-uniform quantizer designed using the k -means clustering algorithm. All the weights in each layer of the YOLO networks are quantized. We then evaluate and compare the performance of the quantized networks using the three SQ schemes with a different number of quantization levels.

1.6 Thesis outlines

The rest of the thesis is organized as follows, Chapter 2 gives a brief review of the prior art research on object detection systems specifically on the convolutional neural network. Chapter 3 describes the proposed methods for research work. In the first part, the proposed network architecture is described alongside the fundamentals of CNN. Then the methods for quantization will be described in the last part of the chapter. Chapter 4 describes the result obtained from the research work. The results are being analyzed and discussed in this chapter. Lastly, Chapter 5 concludes the research work and explains the potential for future works of the research.



CHAPTER 2

LITERATURE REVIEW

2.1 Introduction

This chapter presents the existing research on the fundamentals and the application of CNN. The first part of the chapter explains the basics of multi-layer perceptron which is the fundamental and commonly used in the Artificial Neural Network (ANN). Then the basic CNN model will be explained together with the evolution of the CNN design. Next, this chapter will also describe the commonly used network architecture in CNN. Then, the last part of the chapter will briefly describe the fundamental of quantization used in this research work.

2.2 Artificial Neural Network

Artificial Neural Network (ANN) is a system inspired by the biological neural network connection of the neuron cells in the brain. The brain consists of approximately 100 billion neurons, which each of the neurons is communicating with each other by sending electrochemical signals (Zhang, 2019). The neurons are connected through the junctions which are known as synapses. Each neuron is able to receive thousands of connections from other neurons, constantly receiving incoming signals until the signal reaches the cell body. If the accumulation of the signals surpasses a certain value of threshold, a response signal is sent through the axon. The example of biological neuron cells in the human brain is shown in Figure 2.1.