

Master's Thesis



Czech
Technical
University
in Prague

F3

Faculty of Electrical Engineering
Department of Cybernetics

Deep Neural Networks in Embedded Systems

Bc. Mykhaylo Zelenskyy

Supervisor: Ing. Lukáš Hrubý
April 2019

Acknowledgements

Declaration

I declare that this work is all my own work and I have cited all sources I have used in the bibliography. Prague, dubna 24, 2019

Prohlašuji, že jsem předloženou práci vypracoval samostatně, a že jsem uvedl veškerou použitou literaturu. V Praze,
24. April 2019

Abstract

Keywords:

Supervisor: Ing. Lukáš Hrubý

Abstrakt

Klíčová slova:

Překlad názvu: Hluboké neuronové sítě
ve vestavěných systémech

Contents

1 Introduction	1
1.1 Nvidia Jetson	1
2 Related works	3
3 Neural networks	5
3.1 Architecture	6
3.1.1 Convolutional layer.....	6
3.1.2 Non-linearity layer	6
3.1.3 Pooling layer	7
3.1.4 Fully connected layer	7
3.2 Training of CNN	7
3.2.1 Backpropagation	8
3.2.2 Dropout	9
3.2.3 Batch normalization	9
4 Used networks	11
4.1 ResNet	11
4.2 RetinaNet	12
4.3 SqueezeDet	12
4.4 YOLO	13
4.4.1 YOLO v1	13
4.4.2 YOLO v2	14
4.4.3 YOLO v3	14
5 Experiments	19
5.1 Datasets	19
5.1.1 KITTI dataset	19
5.1.2 GoodVision dataset	19
5.2 Performance evaluation	20
5.3 Inference time evaluation	21
5.3.1 Mixed precision calculation..	22
A Darknet-19	25
B Darknet-53	27
C Sample from videos used for evaluation	29
D Bibliography	31

Figures

1.1 Nvidia Jetson comparison https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/	1
1.2 Tensor core operation https://devblogs.nvidia.com/programming-tensor-cores-cuda-9/	2
3.1 Neuron cell[1]	5
3.2 Convolutional neural network[2] .	6
3.3 Convolutional layer[3]	6
3.4 Convolution computation[3]	7
3.5 Comparation of ReLu and sigmoid non-linearities	8
3.6 Pooling layer[3]	9
4.1 Residual block.....	11
4.2 ResNet architecture	12
4.3 Feature Pyramid Network	16
4.4 Fire module in SqueezeNet	16
4.5 SqueezeDet comparision with other state-of-art solutions	17
4.6 YOLO v1 architecture.....	17
4.7 YOLO v2 improvement.....	17
4.8 YOLO v3 comparation	18
5.1 Image samples from KITTI dataset	20
5.2 Image samples from GoodVision dataset	21
5.3 Intersection over Union https://www.pyimagesearch.com/2016/11/07/intersection-over-union-iou-for-object-detection/	22
A.1 Darknet-19	25
B.1 Darknet-53	27
C.1 Samples from used videos	30

Tables

5.2 Characteristics of video for inference time evaluation	22
5.1 CNN evaluation	23
5.3 PC1 technical specification	23
5.4 PC2 technical specification	24
5.5 Inference time evaluation	24
5.6 Inference time using mixed precision calculation.....	24

Chapter 1

Introduction

1.1 Nvidia Jetson

Nvidia Jetson is a series of embedded modules producing by Nvidia designed for accelerating machine learning applications. First board, Jetson TK1, was presented in 2014. Jetson TK2 was announced in 2017 and was designed for low power systems like smaller camera drones. The next module called Jetson Xavier introduced in 2018 brings up to $20\times$ acceleration compared to predecessor devices with power efficiency being improved $10\times$. The newest Nvidia Nano was announced in 2019 and is focused on hobbyist robotics thanks to its low price. The last three modules are compared in 1.1. As we can

	Jetson Nano™	Jetson TX2			Jetson AGX Xavier™
		TX2 16GB	TX2 4GB	TX2i	
GPU	NVIDIA Maxwell™ architecture with 128 NVIDIA CUDA™ cores	NVIDIA Pascal™ architecture with 256 NVIDIA CUDA cores			NVIDIA Volta™ architecture with 512 NVIDIA CUDA cores and 64 Tensor cores
CPU	Quad-core ARM® Cortex®-A57 MPCore processor	Dual-core Denver 2 64-bit CPU and quad-core ARM A57 complex			8-core ARM v8.2 64-bit CPU, 8 MB L2 + 4 MB L3
Memory	4 GB 64-bit LPDDR4	8 GB 128-bit LPDDR4	4 GB 128-bit LPDDR4	8 GB 128-bit LPDDR4	16 GB 128-bit LPDDR4x
Storage	16 GB eMMC 5.1	32 GB eMMC 5.1	16 GB eMMC 5.1	32 GB eMMC 5.1	32 GB eMMC 5.1
Video Encode	4K @ 30 [H.264/H.265]	2x 4K @ 30 (HEVC)		8x 4K @ 60 (HEVC)	
Video Decode	4K @ 60 [H.264/H.265]	2x 4K @ 30 12-bit support		12x 4K @ 30 12-bit support	
Connectivity	Wi-Fi requires external chip	Wi-Fi onboard	Wi-Fi requires external chip	Wi-Fi requires external chip	
		Gigabit Ethernet			
Camera	12 lanes [3x4 or 4x2] MIPI CSI-2, D-PHY 1.1 [1.5 Gbps]	12 lanes MIPI CSI-2, D-PHY 1.2 [30 Gbps]		16 lanes MIPI CSI-2, 8 SLVS-EC D-PHY [40 Gbps], C-PHY [109 Gbps]	
Size	69.6 mm x 45 mm	87 mm x 50 mm		100 mm x 87 mm	
Mechanical	260-pin edge connector	400-pin connector		699-pin connector	

Figure 1.1: Nvidia Jetson comparison
<https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/>

see, unlike other boards Jetson Xavier is build around NVIDIA Volta™ GPU with tensor cores. These tensor cores accelerate large matrix operations and can perform mixed-precision matrix multiplication and accumulate calculation in a single operation. Each core provides a $4\times 4\times 4$ matrix processing array which performs the operation $D = A \cdot B + C$ shown in 1.2, where A, B, C, D are 4×4 matrices. This operation is crucial for most of machine learning applications, especially in deep learning, because, as we will discuss in further

1. Introduction

chapters, output of each neuron in neural networks are calculated in similar way. Another advantage of Jetson Xavier is using of Nvidia Deep Learning

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix}_{\text{FP16 or FP32}} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix}_{\text{FP16}} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}_{\text{FP16 or FP32}}$$

Figure 1.2: Tensor core operation

<https://devblogs.nvidia.com/programming-tensor-cores-cuda-9/>

Accelerator engines¹. These engines improve energy efficiency and free up the GPU to run more complex networks and implemented dynamic tasks and has up to 5 trillion operations per second (TOPS) INT8 or 2.5 TFLOPS FP16 performance with a power consumption of only 0.5-1.5W. They also supports accelerating of CNN convolution, deconvolution, activation functions, min/max/mean pooling, local response normalization, and fully-connected layers.

¹<http://nvdla.org/>



Chapter 2

Related works

Chapter 3

Neural networks

Artificial neural networks are systems inspired by a brain. The basic computation unit in a brain is a neuron (see 3.1), which has input and output. The input is a dendritic tree connected to the outputs of other neurons called axons. Neurons operate in a single direction from the input to the output and their output is binary. Neurons are also basic computation elements of artificial neural networks. Similarly to biological neural networks, it can have several inputs and outputs. Every neuron can be described by function $f(\omega \cdot \mathbf{x} + b)$, where \mathbf{x} is the input, ω denotes weights, b is a bias and f is the activation function. There are several types of artificial neural networks that are commonly used in machine learning. The most popular type used in object detection is convolutional neural network (CNN), which will be described in the following section.

Neuron

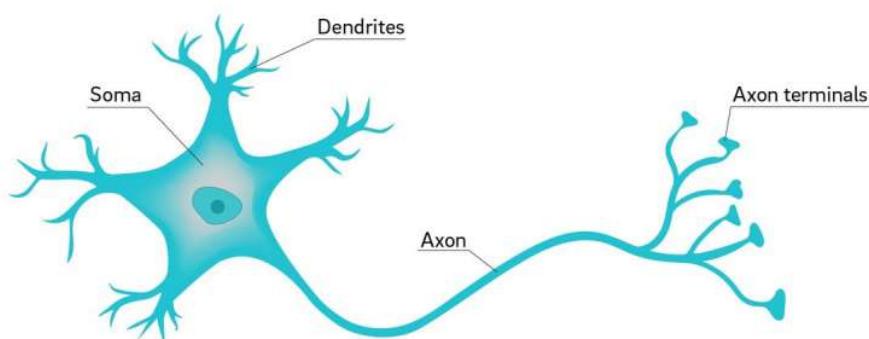


Figure 3.1: Neuron cell[1]

3.1 Architecture

All CNN models has a similar architecture as it is shown in 3.2. The input of such neural network is an image. CNN consists of a series of convolution and pooling operations followed by fully connected layers. These operations are described in the next paragraphs.

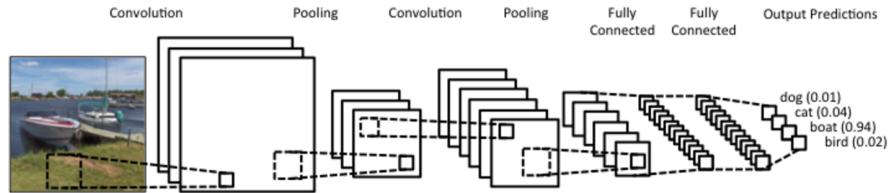


Figure 3.2: Convolutional neural network[2]

3.1.1 Convolutional layer

Convolutional layers consist of neurons placed in a grid of size $N \times M \times C$, where N, M denotes width and height of convolutional filter and C is number of channels in the previous layer (see 3.3). The filter moves from the left to the right with a certain stride until it completes processing width, then it moves down by the same stride to the beginning of the image and repeats the process till the whole image is traversed. The process computes convolution as it is shown in 3.4. Calculated feature map is usually smaller than the input, but it is possible to preserve the same dimensionality by using padding to surround the input with zeros.

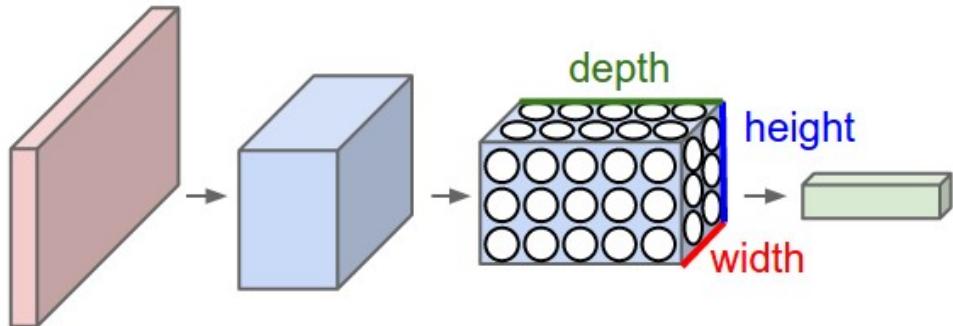


Figure 3.3: Convolutional layer[3]

3.1.2 Non-linearity layer

A non-linearity layer consists of an activation function that takes calculated feature map and creates the activation map as its output. The most common non-linearities used in CNN are sigmoid and ReLu 3.5.

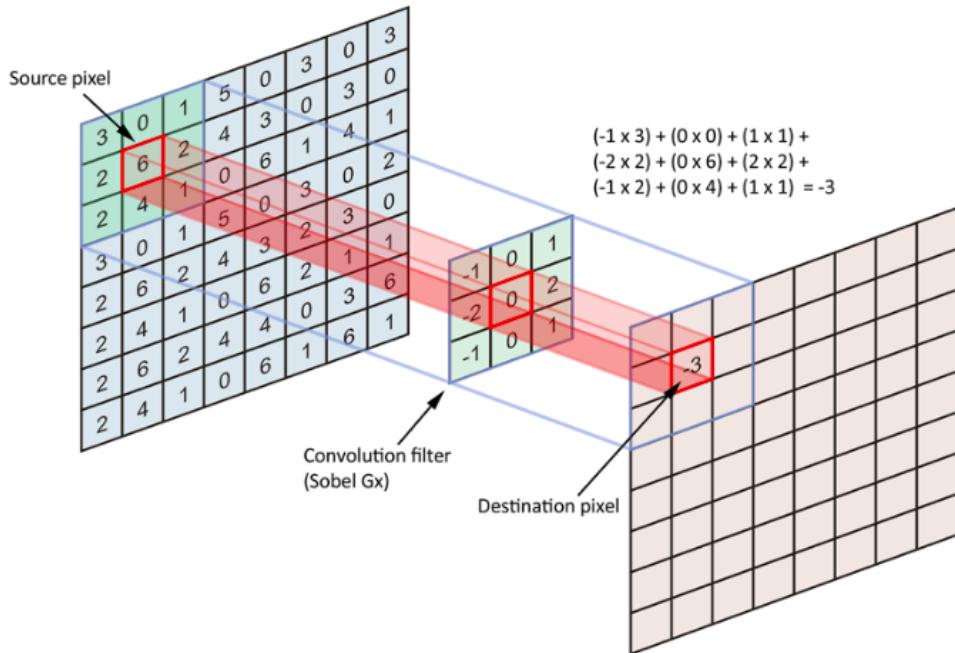


Figure 3.4: Convolution computation[3]

3.1.3 Pooling layer

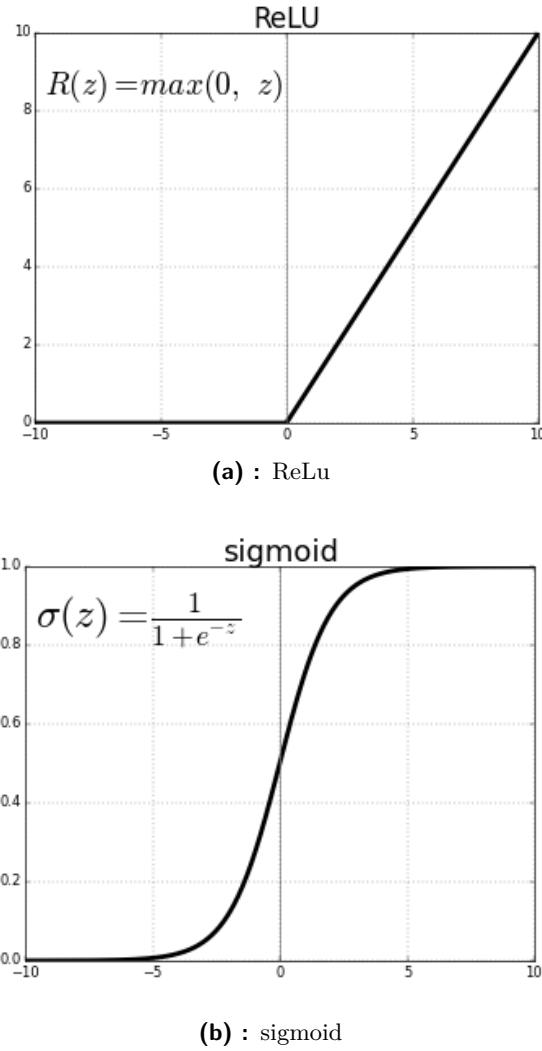
After convolution, pooling layer is used to reduce the dimensionality which enables to reduce number of parameters. Two most common pooling operation are max and min pooling. It simply slides the input with particular stride and choose maximal or minimal value in predefined window (see 3.6). Pooling helps to not overfit CNN and can reduce the training time.

3.1.4 Fully connected layer

In fully connected layers, each neuron is connected to every neuron in the previous layer just like in feedforward neural networks.

3.2 Training of CNN

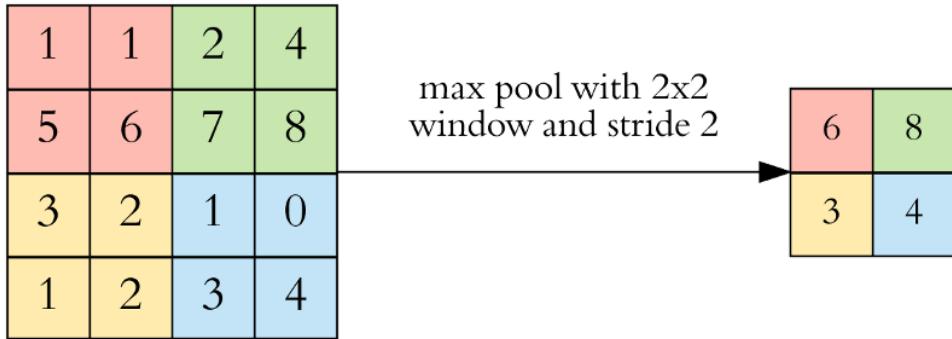
Before we can use any neural network, it must be trained to understand how objects we want it to recognize should look like. The weight of filters are randomized and the filters in lower layers of CNN don't know to look for edges or curves, the filters in higher layers don't know to look for more concrete shapes like wheels, legs, faces. As any supervised learning, CNNs are given a training set of thousands of images with labels to learn features of objects. Learning algorithm is called backpropagation.

**Figure 3.5:** Comparation of ReLu and sigmoid non-linearities

■ 3.2.1 Backpropagation

Backpropagation was firstly introduced in [4] in 1986. This process can be separated into 4 steps: the forward pass, the calculation of loss function, the backward pass and the weight update. During the forward pass, we take a batch of training images and pass it through the network. After first training forward pass output of the network would probably be randomized, because the network isn't able to look for any kind of features, thus isn't able to make any reasonable conclusion about training example. This goes to the next step of backpropagation, the calculation of loss function. Each neural network can have its own loss function depending on what its output is, but the most common loss function used for backpropagation is mean squared error

$$L = \frac{1}{2n} \sum_x \|y(x) - a(x)\|^2, \quad (3.1)$$

**Figure 3.6:** Pooling layer[3]

where n denotes number of training inputs x , $y(x)$ label corresponding to the input x and a is network's output. By training process we want to achieve such a result where the predicted label is the same as the training label. To get this, we want to minimize the amount of loss we have, hence we want to find out which weights most directly contribute to the loss of the networks. This leads to the third step, backward pass, where we calculate partial derivations $\frac{\partial L}{\partial \omega}$, $\frac{\partial L}{\partial b}$. Once we compute the derivation, we can update weights and biases by changing them in the opposite direction of the gradient using learning rate η :

$$\omega = \omega - \eta \frac{\partial L}{\partial \omega} \quad (3.2)$$

$$b = b - \eta \frac{\partial L}{\partial b} \quad (3.3)$$

η is the parameter that defines how big steps learning process will take to update the weights, thus how fast we want the model to converge. However, the learning rate that is too big could result in jumps that are too large and not precise enough to reach the optimal point. Some techniques help to decrease training speed and performance.

3.2.2 Dropout

Overfitting is one of the biggest problems in machine learning. Overfitting means that a model performs well on a training dataset, but fails on test data. To solve this problem in neural networks, we can use so-called dropout. During each training step, an individual neuron can be dropped out of the net with probability $1 - p$ or kept with probability p , so that only a reduced network is trained. The removed neurons are then reinserted into the network with unchanged weights. This method not only decreases overfitting but also improves training speed.

3.2.3 Batch normalization

Assume we have a training set of images with cars that has a particular colour. If we try to use the network that was trained on that dataset, it probably

will not work well on cars with another colour. In that case, we might need to retrain the network by trying to align the distribution of cars in different colours. Batch normalization <https://arxiv.org/pdf/1502.03167.pdf> helps with this problem by reducing the amount of covariance shift in hidden layers. It simply normalizes the output of each layer by subtracting the batch mean and dividing by the batch standard deviation. After this change of activation output, the weights in the next layer are no longer optimal. Therefore, batch normalization adds two trainable parameters to each layer and lets gradient descent do the denormalization by changing only these parameters during each activation. Batch normalization also helps with overfitting, because it adds some noise to each layer's activations. It also allows using higher learning rate, because it makes sure no activation goes high or low.

Chapter 4

Used networks

4.1 ResNet

Residual networks described in [5] are classification networks with an image as the input and object class and confidence score as the output. In this paper, they introduced shortcut connections that are widely used in modern neural networks. One of the biggest problems with training deep neural networks is vanishing and exploding gradient. During backpropagation, a lot of small or large numbers are multiplied to compute gradients. When the network is deep, multiplying of small numbers will become zero (vanished) and multiplying of large numbers will explode. Normally we expect deeper neural network will have more accurate predictions, but the opposite is true, and this degradation problem is caused by the vanishing gradient. This problem can be solved by adding shortcut connection which adds the input to the output after few weight layers, hence the output is $H(x) = F(x) + x$ (see 4.1). There are two types of residual connections. The identity shortcuts can be directly used when both input and output have the same dimension, or extra zero padding can be used when dimensions change. In both cases, no extra parameters are needed. Comparing plain and residual network with 34 layers (see 4.2) Top-1 error drops from 28.54% to 25.03%. On the other hand, if we compare smaller network with 18 layers, Top-1 error changes from 27.94% to 27.88%, which means shortcut connections perform better in deeper networks. ResNets with different number of layers are often used as

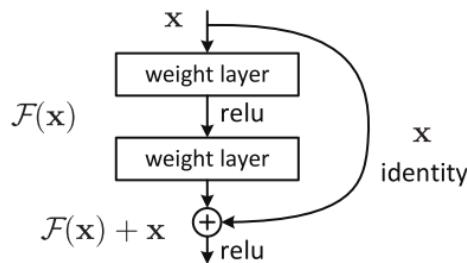


Figure 4.1: Residual block

classification networks (backbone) in detectors such as RetinaNet.

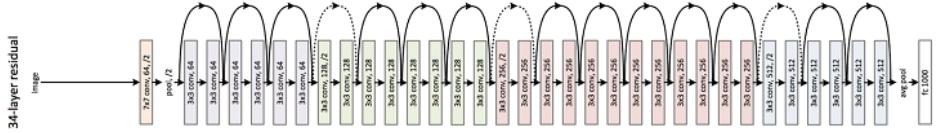


Figure 4.2: ResNet architecture

4.2 RetinaNet

<https://medium.com/@14prakash/the-intuition-behind-retinanet-eb636755607d>

RetinaNet was proposed by Facebook AI Research and its features are described in [6] and [?]. They proposed using anchor boxes instead of predicting bounding boxes. Sizes of the anchor boxes are predefined and used in further predictions. Thus, the network does not predict the final size of the object, but instead, it only adjusts the size of the nearest anchor to the size of the object. Also, they suggested a solution for object detection in different scales. Originally a pyramid of the same image at different scales was used to detect the object. However, this solution is time-consuming and has a high memory demand. Instead, a pyramid of features can be used. Although it is not such efficient for accurate object detection as image pyramids, it provides result faster and with less memory consumption. In [6] authors propose Feature Pyramid Network (FPN) which is fast like the described pyramid of features, but more accurate. Its architecture is seen in 4.3. The other solution, focal loss, solves class imbalance. Instead of normal cross entropy calculated by

$$C(p, y) = - \sum_i y_i \ln p_i \quad (4.1)$$

scaled entropy is used using following equation:

$$C(p, y) = - \sum_i y_i (1 - p_i)^\lambda \ln p_i. \quad (4.2)$$

Here we can see focusing parameter $\lambda \geq 0$ which smoothly adjusts the rate at which easy examples are down weighted and thus training is focused on hard negatives. In this thesis we used Keras implementation of RetinaNet¹ implemented in *TensorFlow* with ResNet50 as a backbone.

4.3 SqueezeDet

SqueezeDet (<https://arxiv.org/pdf/1612.01051.pdf>) is a single stage detection pipeline inspired by YOLO. The main difference between two architectures is that SqueezeDet uses SqueezeNet (<https://arxiv.org/pdf/1602.07360.pdf>) for feature extraction. The building brick of SqueezeNet is called fire module 4.4. Each fire module contains a squeeze layer and an expand layer. Squeeze

¹<https://github.com/fizyr/keras-retinanet>

layers replace 3×3 filters by 1×1 filters to reduce computation complexity 9 times. Following expand layers contain number of 1×1 and 3×3 kernels. Squeeze layers reduce depth of calculated feature map, which means the following 3×3 filters in expand layers have to do fewer computation. Thanks to its architecture, SqueezeDet can be faster and smaller comparing to other state-of-art solutions (see 4.5), and so can be efficiently used on embedded system. The loss function of SqueezeDet is defined as

$$\begin{aligned} & \frac{\lambda_{bbox}}{N_{obj}} \sum_{i=1}^W \sum_{j=1}^H \sum_{k=1}^K I_{ijk} \left[\left(\beta x_{ijk} - \beta x_{ijk}^G \right)^2 + \left(\beta y_{ijk} - \beta y_{ijk}^G \right)^2 \right. \\ & \quad \left. + \left(\beta w_{ijk} - \beta w_{ijk}^G \right)^2 + \left(\beta h_{ijk} - \beta h_{ijk}^G \right)^2 \right] \\ & + \sum_{i=1}^W \sum_{j=1}^H \sum_{k=1}^K \frac{\lambda_{conf}^+}{N_{obj}} I_{ijk} \left(\gamma_{ijk} - \gamma_{ijk}^G \right)^2 + \frac{\lambda_{conf}^-}{WHK - N_{obj}} \bar{I}_{ijk} \gamma_{ijk}^2 \\ & \quad + \frac{1}{N_{obj}} \sum_{i=1}^W \sum_{j=1}^H \sum_{k=1}^K \sum_{c=1}^C I_{ijk} l_c^G \log(p_c) \quad (4.3) \end{aligned}$$

, where first part is the bounding box regression and $(\beta x_{ijk}, \beta y_{ijk}, \beta w_{ijk}, \beta h_{ijk})$ corresponds to the relative coordinate of anchor- k located at grid center- (i, j) . Second part denotes confidence score regression with output γ_{ijk} . The last part is cross-entropy loss for classification. We used a tensorflow implementation of SqueezeDet available to download from GitHub¹.

4.4 YOLO

4.4.1 YOLO v1

A new approach for object detection, YOLO architecture, was presented in [7]. A single neural network is used to predict both bounding boxes and class probabilities, hence an image is evaluated only once. The described system divides the input into a $S \times S$ grid, and if the center of an object falls into a grid cell, this cell is responsible for detecting that object. Each cell also predicts B bounding boxes and confidence score for them. Confidence is defined as

$$score = Pr(Object) \cdot IoU_{pred}^{truth}, \quad (4.4)$$

where $Pr(Object)$ is a probability of an object being inside that bounding box and IoU_{pred}^{truth} denotes intersection over union between ground truth and prediction. Each bounding box consists of $(x, y, w, h, score)$, where (x, y) represents the center of the box and (w, h) denotes its width and height. Each grid cells also predicts conditional probability $C = Pr(Class_i|Object)$. The model consists of 24 convolutional layers followed by 2 fully connected layers as it is shows in 4.6. Training process optimizes loss function

¹<https://github.com/BichenWuUCB/squeezeDet>

$$\begin{aligned}
& \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbf{1}_{ij}^{obj} \left[(x_i - \hat{x}_i)^2 + (y_i + \hat{y}_i)^2 \right] \\
& + \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbf{1}_{ij}^{obj} \left[(\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} + \sqrt{\hat{h}_i})^2 \right] \\
& + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbf{1}_{ij}^{obj} (C_i - \hat{C}_i)^2 \\
& + \lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbf{1}_{ij}^{noobj} (C_i - \hat{C}_i)^2 \\
& + \sum_{i=0}^{S^2} \mathbf{1}_{ij}^{obj} \sum_{c \in classes} (p_i(c) - \hat{p}_i(c))^2 \quad (4.5)
\end{aligned}$$

, where

$$\mathbf{1}_{ij}^{obj} = \begin{cases} 1, & \text{if there is an object.} \\ 0, & \text{otherwise,} \end{cases} \quad (4.6)$$

$\mathbf{1}_{ij}^{noobj}$ is inverse function to $\mathbf{1}_{ij}^{obj}$, λ_{coord} and λ_{noobj} are constant to increase the loss from bounding box coordinate prediction and decrease the loss from confidence prediction for boxes that does not contain objects. While YOLO v1 was faster than most of the existing approaches for object detection, it had relatively low 57.9% mAP on the VOC 2012 test set compared to the existing state of the art.

■ 4.4.2 YOLO v2

A new version of YOLO was introduced in [8]. Authors of this state of the art detector refer to it as a better, faster and stronger version of YOLO. For better performance, they added batch normalization and used images with a bigger resolution to train the network. They also removed fully connected layers and used anchor boxes to predict bounding boxes, which lead to a small decrease in mAP from 69.5% to 69.2%, but it also increased a recall from 81% to 88%. We can see how applied changes improved network performance in 4.7. They also proposed a new classification network called Darknet-19 (see Appendix A) to make YOLO even faster. We can see that Darknet-19 has many 1×1 convolutions to reduce the number of parameters.

■ 4.4.3 YOLO v3

The newest version of YOLO was presented in [9]. Similar to YOLOv2 it predicts bounding boxes using dimension clusters as anchor boxes. The network predicts four coordinates for each bounding box and for training they use a sum of squared error loss. Objectness score for each bounding box is predicted using logistic regression, which should be one if the bounding box prior overlaps a ground truth object by more than any other bounding

box prior. They also use 3 different scales for prediction, which is similar to feature pyramid networks. Deeper extractor called Darknet-53 (see Appendix B) with shortcut connections is used for feature extraction.

Comparing to other state of art solutions, YOLO v3 has similar performance, but it is much faster as it is seen in 4.8. Unlike RetinaNet and SqueezeDet, YOLO uses another neural network framework, Darknet¹, written in C and CUDA.

¹<https://pjreddie.com/darknet/>

4. Used networks

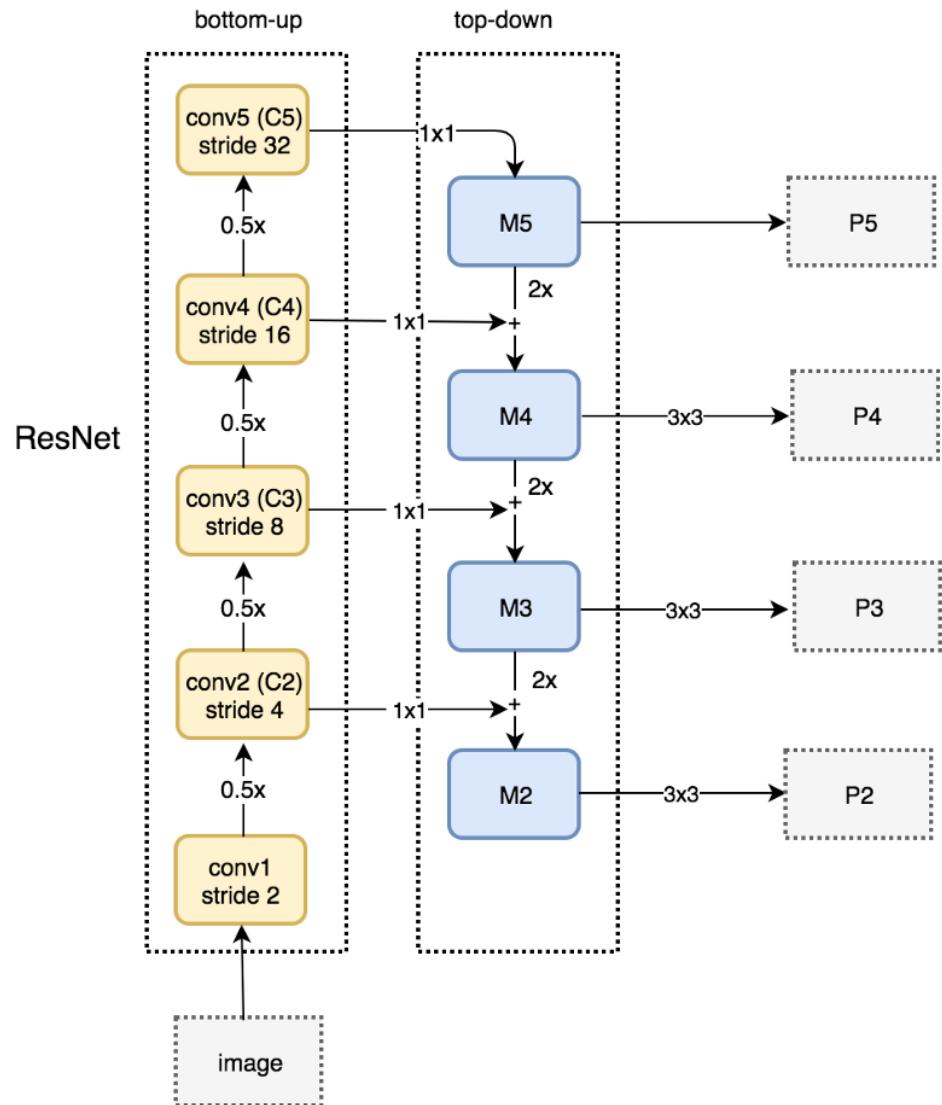


Figure 4.3: Feature Pyramid Network

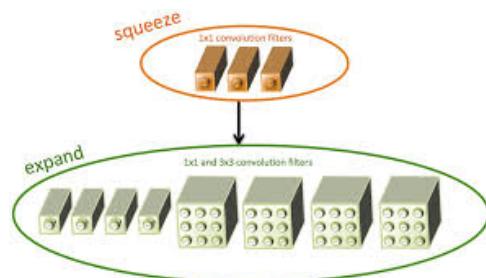
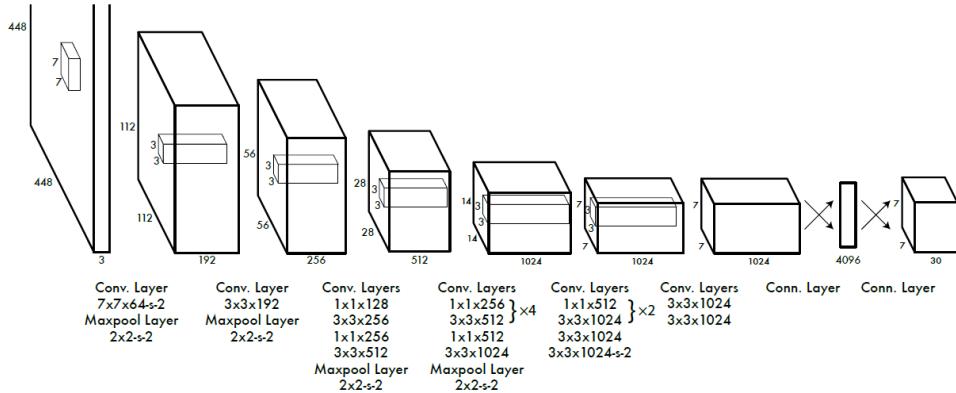


Figure 4.4: Fire module in SqueezeNet

model	Model Size (MB)	FLOPs $\times 10^9$	Activation Memory Footprint (MB)	Average GPU Power (W)	Inference Speed (FPS)	Energy Efficiency (J/frame)	mAP*
SqueezeDet	7.9	9.7	117.0	80.9	57.2	1.4	76.7
SqueezeDet: scale-up	7.9	22.5	263.3	89.9	31.3	2.9	72.4
SqueezeDet: scale-down	7.9	5.3	65.8	77.8	92.5	0.84	73.2
SqueezeDet: 16 anchors	9.4	11.0	117.4	82.9	51.4	1.6	66.9
SqueezeDet+	26.8	77.2	252.7	128.3	32.1	4.0	80.4
VGG16+ConvDet	57.4	288.4	540.4	153.9	16.6	9.3	79.1
ResNet50+ConvDet	35.1	61.3	369.0	95.4	22.5	4.2	76.1
Faster-RCNN + VGG16 [1]	485	-	-	200.1	1.7	117.7	-
Faster-RCNN + AlexNet [1]	240	-	-	143.1	2.9	49.3	-
YOLO**	753	-	-	187.3	25.8	7.3	-

Figure 4.5: SqueezeDet comparison with other state-of-art solutions**Figure 4.6:** YOLO v1 architecture

	YOLO	YOLOv2						
batch norm?	✓	✓	✓	✓	✓	✓	✓	✓
hi-res classifier?		✓	✓	✓	✓	✓	✓	✓
convolutional?		✓	✓	✓	✓	✓	✓	✓
anchor boxes?		✓	✓					
new network?			✓	✓	✓	✓	✓	✓
dimension priors?				✓	✓	✓	✓	✓
location prediction?					✓	✓	✓	✓
passthrough?						✓	✓	✓
multi-scale?						✓	✓	✓
hi-res detector?							✓	
VOC2007 mAP	63.4	65.8	69.5	69.2	69.6	74.4	75.4	76.8
								78.6

Figure 4.7: YOLO v2 improvement

4. Used networks

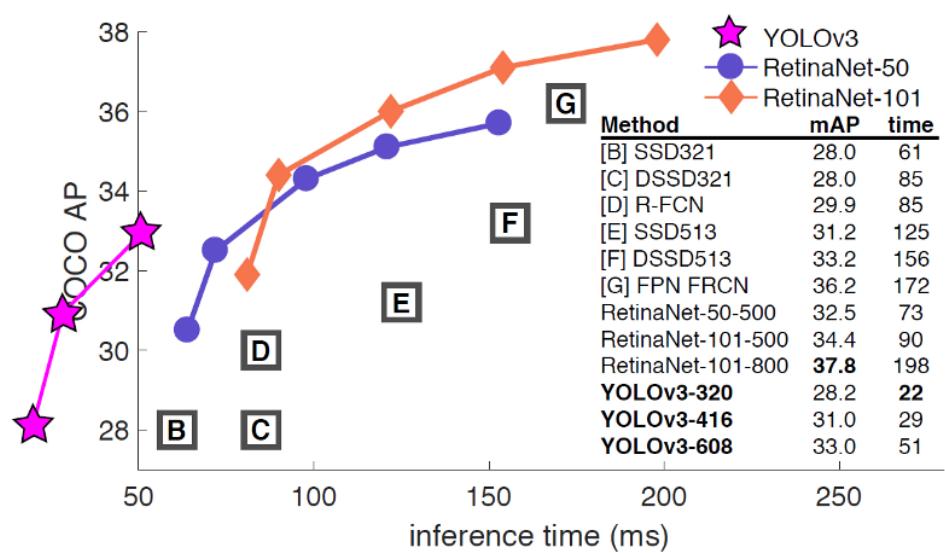


Figure 4.8: YOLO v3 comparation

Chapter 5

Experiments

For experiments, we chose SqueezeDet, YOLOv3, and RetinaNet. We trained all three networks on KITTI (<http://www.cvlibs.net/datasets/kitti/>) and internal GoodVision dataset specifically prepared for our use case. Thanks to the ability of RetinaNet and YOLOv3 to adapt to the size of the input image, we were able to test different image size ratios with no need to retrain these networks. YOLOv3 was evaluated with two different input image resolutions: 608×608 and 418×418 ; for evaluation of RetinaNet, we used three different largest side sizes: 1024, 608 and 418. Unfortunately, SqueezeDet does not adapt its layers to an input image, so we used a network with input image resolution 1242×375 defined by the pre-trained model.

5.1 Datasets

5.1.1 KITTI dataset

KITTI dataset consists of 7481 training images and 7518 test images with a total of 80256 labeled objects. Original datasets has 9 different type of objects: “Car”, “Van”, “Truck”, “Pedestrian”, “Person_sitting”, “Cyclist”, “Tram”, “Misc” or “DontCare”. We used only 5 of them for training: “Car”, “Van”, “Truck”, “Pedestrian”, “Person_sitting”, though we merged “Person_sitting” and “Pedestrian” classes into single “Person” object type. Regardless of lack of common classes like “Bus”, “Bicycle”, “Motorcycle”, it fits our needs, because all data were gathered by driving around cities, in country areas, and on highways, and so our networks would be trained on real-life traffic data with no redundant information 5.1.

5.1.2 GoodVision dataset

Internal GoodVision dataset was specifically created for object detection and recognition in traffic. It consists of 4917 images with more than 130000 labeled objects. Unlike in KITTI dataset, such classes as “Bicycle”, “Bus”, “Motorcycle” are presented here, which means neural networks trained on this dataset will be more complex and more adequate for our use case. Images



Figure 5.1: Image samples from KITTI dataset

for this dataset were gathered from cameras placed at different heights and angles to roads and highways. Samples from the dataset are presented in 5.2.

5.2 Performance evaluation

For our evaluation, we calculated average precision for every class in the evaluation set. We had to calculate precision and recall, which are defined as

$$\text{Precision} = \frac{TP}{TP + FP}, \quad (5.1)$$

$$\text{Recall} = \frac{TP}{TP + FN}, \quad (5.2)$$

where TP is true positives, FP is false positives, and FN is false negatives. It means precision measures how accurate predictions are, while recall refers to the percentage of total relevant results correctly classified by our network. To determine if the prediction is a true positive or a false positive, intersection over union (IoU) has to be measured. As picture 5.3 shows, IoU is simply the ratio between the area of overlap between the ground truth bounding box and predicted bounding box, and the area encompassed by both ground truth bounding box and predicted bounding box. If IoU is over some predefined threshold, the prediction is considered to be true. Otherwise, it is a false positive. For evaluation, we chose this threshold to be 0.5. Then we calculate precision/recall curve (<https://acutearetesting.org/en/articles/precision-recall-curves-what-are-they-and-how-are-they-used>) and average precision is area under that curve and is calculated for every class independently. Mean



Figure 5.2: Image samples from GoodVision dataset

average precision is simply average value of average precisions across all classes. Results of the evaluation are presented in table 5.1. YOLOv3 performs better than both SqueezeDet and RetinaNet at both resolution. Although YOLOv3 trained on KITTI dataset indicates bigger mAP on KITTI evaluation set, it has worse performance on GoodVision evaluation set, while YOLOv3 trained on GoodVision dataset has almost 10% better mAP. We can notice this trend for all 3 CNNs: while KITTI-trained model has better performance on KITTI test set, it fails on GoodVision set. If we compare models trained on KITTI and GoodVision dataset using AP of classes presented in KITTI dataset only, the difference will be even bigger.

5.3 Inference time evaluation

For computing on the edge, we should ensure that time delay between receiving input data, in this case stream frame, and providing output data to a user is as small as possible. Time of object detection and classification can be the biggest bottleneck in such systems, therefore it is necessary to choose such architecture that provides the best result with the least possible inference

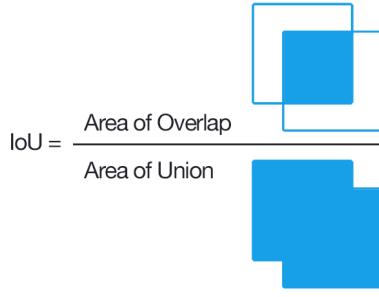


Figure 5.3: Intersection over Union

<https://www.pyimagesearch.com/2016/11/07/intersection-over-union-iou-for-object-detection/>

time. Videos for this evaluation were taken from video databases with free access such as YouTube or Pexels. All videos have different camera view with various object count, video resolution and FPS, as shown in table 5.2. This can affect neural network detector performance, hence we should also compare various input resolutions of neural networks. Video samples presented in Appendix C.1 shows obvious difference between videos scenes. Those scenes represents typical use case for our system.

Video file	FPS	Video Resolution
5.4 4K Camera Road in Thailand.mp4	30	1280x720
Cars Driving On Street.mp4	30	1920x1080
Cars On Highway.mp4	25	1920x1080
Cars On The Road.mp4	50	1280x720
City Traffic.mp4	30	1920x1088
Day Traffic Sample Video Dataset.mp4	30	432x240
Pedestrian and Traffic, Human Activity Recognition Video ,DataSet By UET Peshawar.mp4	30	1280x720
Pexels Videos 1601538.mp4	25	1920x1080
Pexels Videos 2577.mp4	30	1920x1088
Pexels Videos 2670.mp4	25	1920x1088
Pexels Videos 3047.mp4	30	1920x1088
Pexels Videos 948404.mp4	24	3840x2178
moderate_traffic.mp4	30	1280x720

Table 5.2: Characteristics of video for inference time evaluation

We compared the time necessary to process one image frame by CNN using Jetson Xavier and other two different PCs with CPU and GPU specifications presented in tables 5.3 and 5.4.

As we can see from table 5.5, PC1 performs best for all three CNNs. On the other hand, PC2 has worse inference time comparing to Jetson Xavier using YOLOv3 and RetinaNet and almost the same while using squeezeDet. It is worth mentioning that YOLOv3 has similar inference time for both tested resolutions, and squeezeDet has the best inference time among all tested CNNs.

5.3.1 Mixed precision calculation

We have discussed in section 1.1 that Nvidia Jetson has tensor cores to accelerate matrix calculation using mixed precision matrix multiplication

5.3. Inference time evaluation

Model	Training dataset	Testing dataset	Bicycle AP [%]	Bus AP [%]	Motorcycle AP [%]	Car AP [%]	Person AP [%]	Truck AP [%]	Van AP [%]	mAP [%]
YOLOv3 608	GoodVision	GoodVision	62,5	72,4	54,4	89,3	91,2	83,2	80,1	76,2
	KITTI	—	—	—	81,4	89,3	80,6	73,5	81,2	
	GoodVision	—	—	—	58,6	72,3	73,2	67,0	67,8	
YOLOv3 416	GoodVision	GoodVision	63,6	61,2	40,7	83,4	84,0	72,7	65,2	67,3
	KITTI	—	—	—	73,6	80,7	68,3	59,5	70,5	
	GoodVision	—	—	—	52,1	63,5	62,2	43,0	55,2	
RetinaNet 1024	GoodVision	GoodVision	46,3	34,8	15,4	39,5	32,3	53,8	40,6	37,5
	KITTI	—	—	—	35,2	29,9	51,6	31,1	37,0	
	GoodVision	—	—	—	20,6	23,3	42,2	19,4	26,4	
RetinaNet 608	GoodVision	GoodVision	11,0	12,7	0,2	25,5	18,0	28,3	23,4	17,0
	KITTI	—	—	—	23,0	16,6	27,8	20,3	21,9	
	GoodVision	—	—	—	2,3	16,3	24,6	18,1	15,3	
RetinaNet 416	GoodVision	GoodVision	0,0	5,8	0,6	19,5	12,9	15,8	16,7	10,2
	KITTI	—	—	—	19,1	12,3	14,6	15,1	15,3	
	GoodVision	—	—	—	13,4	7,8	10,6	11,2	10,8	
squeezeDet	GoodVision	GoodVision	1,7	8,2	4,1	34,4	39,2	22,9	12,9	17,6
	KITTI	—	—	—	30,6	38,5	22,5	10,8	25,6	
	GoodVision	—	—	—	3,1	2,9	4,7	4,2	3,7	

Table 5.1: CNN evaluation

CPU	Intel(R) Core(TM) i7-8700K CPU @ 3.70GHz
GPU	3584-core 11Gb GeForce GTX 1080 Ti @ 1582MHz

Table 5.3: PC1 technical specification

and accumulation. Darknet framework used for YOLOv3 can activate tensor core directly, because it is written in CUDA. To utilize tensor core using TensorFlow we should use TensorRT platform¹. TensorRT is able to convert TensorFlow CNN graph into supported format and use it for inference with mixed precision calculation. We achieved significant decrease in the YOLOv3 inference time, but in RetinaNet the difference between measured inference time using plain TensorFlow and inference time using TensorRT was small, which is also described in table 5.6. Insignificant improvement in inference FPS is caused by the fact that TensorRT supports only some layers defined in TensorFlow. We used RetinaNet implemented in Keras and some of its layer are custom and not supported in TensorRT. This problem can be solved by reimplementing RetinaNet in TensorFlow or add support for these layers into TensorRT. Although this was a part of the thesis assignment, after some discussions with supervisor, we decided not to improve RetinNet converter, because we have already achieved sufficiently good results with YOLOv3 in both detector precision and its inference time.

¹<https://developer.nvidia.com/tensorrt>

CPU	Intel(R) Core(TM) i5-7300HQ CPU @ 2.50GHz
GPU	640-core 4Gb GeForce GTX 1050 @ 1404MHz

Table 5.4: PC2 technical specification

Model	Machine	Inference time [ms]
YOLOv3 416x416	Jetson Xavier	123
	PC1	35
	PC2	175
YOLOv3 608x608	Jetson Xavier	139
	PC1	39
	PC2	208
squeezeDet	Jetson Xavier	25
	PC1	24
	PC2	8
RetinaNet 1024	Jetson Xavier	210
	PC1	52
	PC2	228
RetinaNet 608	Jetson Xavier	107
	PC1	28
	PC2	100
RetinaNet 416	Jetson Xavier	74
	PC1	19
	PC2	65

Table 5.5: Inference time evaluation

Model	Original inference time [ms]	Mixed-precision inference time [ms]
YOLOv3 416x416	123	49
YOLOv3 608x608	139	91
RetinaNet 1024	210	186
RetinaNet 608	107	95
RetinaNet 416	74	67

Table 5.6: Inference time using mixed precision calculation

Appendix A

Darknet-19

Type	Filters	Size/Stride	Output
Convolutional	32	3×3	224×224
Maxpool		$2 \times 2/2$	112×112
Convolutional	64	3×3	112×112
Maxpool		$2 \times 2/2$	56×56
Convolutional	128	3×3	56×56
Convolutional	64	1×1	56×56
Convolutional	128	3×3	56×56
Maxpool		$2 \times 2/2$	28×28
Convolutional	256	3×3	28×28
Convolutional	128	1×1	28×28
Convolutional	256	3×3	28×28
Maxpool		$2 \times 2/2$	14×14
Convolutional	512	3×3	14×14
Convolutional	256	1×1	14×14
Convolutional	512	3×3	14×14
Convolutional	256	1×1	14×14
Convolutional	512	3×3	14×14
Maxpool		$2 \times 2/2$	7×7
Convolutional	1024	3×3	7×7
Convolutional	512	1×1	7×7
Convolutional	1024	3×3	7×7
Convolutional	512	1×1	7×7
Convolutional	1024	3×3	7×7
Convolutional	1000	1×1	7×7
Avgpool		Global	1000
Softmax			

Figure A.1: Darknet-19

Appendix B

Darknet-53

Type	Filters	Size	Output
Convolutional	32	3×3	256×256
Convolutional	64	$3 \times 3 / 2$	128×128
1x	Convolutional	32	1×1
1x	Convolutional	64	3×3
	Residual		128×128
	Convolutional	128	$3 \times 3 / 2$
			64×64
2x	Convolutional	64	1×1
2x	Convolutional	128	3×3
	Residual		64×64
	Convolutional	256	$3 \times 3 / 2$
			32×32
8x	Convolutional	128	1×1
8x	Convolutional	256	3×3
	Residual		32×32
	Convolutional	512	$3 \times 3 / 2$
			16×16
8x	Convolutional	256	1×1
8x	Convolutional	512	3×3
	Residual		16×16
	Convolutional	1024	$3 \times 3 / 2$
			8×8
4x	Convolutional	512	1×1
4x	Convolutional	1024	3×3
	Residual		8×8
	Avgpool		Global
	Connected		1000
	Softmax		

Figure B.1: Darknet-53

Appendix C

Sample from videos used for evaluation



(a) : 5.4 4K Camera Road in Thailand.mp4



(b) : Cars Driving On Street.mp4



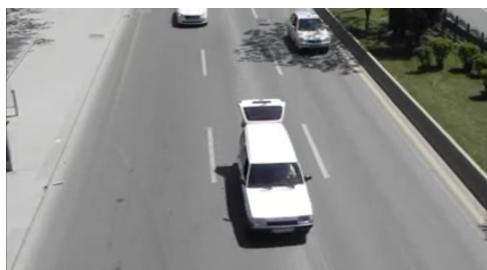
(c) : Cars On Highway.mp4



(d) : Cars On The Road.mp4



(e) : City Traffic.mp4



(f) : Day Traffic Sample Video Dataset.mp4

C. Sample from videos used for evaluation



(g) : Pedestrian and Traffic, Human Activity Recognition Video ,DataSet By UET Peshawar.mp4



(h) : Pexels Videos 1601538.mp4



(i) : Pexels Videos 2577.mp4



(j) : Pexels Videos 2670.mp4



(k) : Pexels Videos 3047.mp4



(l) : Pexels Videos 948404.mp4



(m) : moderate_traffic.mp4

Figure C.1: Samples from used videos

Appendix D

Bibliography

- [1] Science X staff. Why are neuron axons long and spindly? study shows they're optimizing signaling efficiency. <https://medicalxpress.com/news/2018-07-neuron-axons-spindly-theyre-optimizing.html>, Jul 2018.
- [2] Denny Britz. Understanding convolutional neural networks for NLP, Jan 2016.
- [3] Cs231n convolutional neural networks for visual recognition. <http://cs231n.github.io/convolutional-networks>.
- [4] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, 1986.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [6] Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [7] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [8] Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [9] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *CoRR*, abs/1804.02767, 2018.