

Optimized Compression Pipelines For Convolutional Neural Networks

Yiren (Aaron) Zhao
Corpus Christi College



**UNIVERSITY OF
CAMBRIDGE**

*A dissertation submitted to the University of Cambridge
in partial fulfilment of the requirements for the degree of
Master of Philosophy in Advanced Computer Science*

University of Cambridge
Computer Laboratory
William Gates Building
15 JJ Thomson Avenue
Cambridge CB3 0FD
UNITED KINGDOM

Email: yaz21@cam.ac.uk

May 28, 2017

Declaration

I Yiren (Aaron) Zhao of Corpus Christi College, being a candidate for the M.Phil in Advanced Computer Science, hereby declare that this report and the work described in it are my own work, unaided except as may be specified below, and that the report does not contain material that has already been used to any substantial extent for a comparable purpose.

Total word count: 14,235

Signed:

Date:

This dissertation is copyright ©2010 Yiren (Aaron) Zhao.

All trademarks used in this dissertation are hereby acknowledged.

Abstract

This is the abstract. Write a summary of the whole thing. Make sure it fits in one page.

Contents

1	Introduction	1
2	Background	4
2.1	Convolutional Neural Networks	4
2.2	Related Work	6
3	Experiment Setups, Datasets and Trained Networks	9
4	Pruning	12
4.1	Deterministic Pruning	12
4.1.1	Pruning with biases	12
4.1.2	Pruning without biases	16
4.2	Dynamic Network Surgery	17
4.3	Comparison to Existing Works	19
5	Proposed Pruning Strategies	22
5.1	Gradient Profiling	22
5.1.1	Method Description	23
5.1.2	Gradient Profiling and Retraining	24
5.2	Regularization Aided Pruning	26
5.2.1	l_1 and l_2 Norms	27
5.2.2	Summary of Pruning Methods	28
6	Quantization	31
6.1	Weights Sharing	31
6.2	Fixed-point Quantization	33
6.3	Dynamic Fixed-point Quantization	35
7	Proposed Quantization Strategies	38
7.1	Customized Floating-point Quantization	38
7.2	Re-centralized Quantization	41

7.3	Summary of Quantization Methods	46
8	Evaluation	48
8.1	Compression Pipeline	48
8.2	Evaluation of Performance	50
8.3	Evaluation of Compression Techniques	51
9	Summary and Conclusions	52
9.1	Conclusion	52
9.2	Future Works	53

List of Figures

2.1	AlexNet Network Architecture [1].	5
2.2	Fully connected layer followed by an activation function [1]. . .	5
3.1	LeNet5 Network Architecture	10
3.2	CifarNet Network Architecture	11
4.1	Distribution of weights in fc1 layer of pruned and unpruned LeNet5	15
4.2	Mechanism of Dynamic network surgery [2].	17
5.1	<i>Gradient Profiling.</i>	23
5.2	<i>Gradient profiling and Dynamic network surgery</i> without re- training.	25
5.3	<i>Gradient profiling and Dynamic network surgery</i> with 10 epochs retraining.	26
5.4	Effects of $l1$ and $l2$ norms with different hyperparameters. . .	27
6.1	Weights sharing: training (bottom) and inference (top) [3]. . .	32
6.2	Number representation system for fixed-point quantization with n -bit fraction.	34
7.1	Number representation system for customized floating-point quantization with e -bit exponent and m -bit mantissa.	39
7.2	Parameter distribution of the first fully connected layer in <i>LeNet5</i> , and the color coded arithmetic precisions using <i>Cus-</i> <i>tomized floating-point quantization</i> with 1 bit sign, 1 bit ex- ponent and 3 bits fraction. A deeper red color corresponds to a greater arithmetic precision.	41

7.3	Parameter distribution of the first fully connected layer in <i>LeNet5</i> , and the color coded arithmetic precisions using <i>Centralized customized floating-point quantization</i> with 1 bit sign, 1 bit centre, 1 bit exponent and 3 bits fraction. A deeper red color corresponds to a greater arithmetic precision.	42
7.4	Number representation system for <i>Re-centralized customized floating-point quantization</i> with E-bit exponent and M-bit mantissa.	43
7.5	Number representation system for <i>Re-centralized dynamic fixed-point quantization</i> with N-bit fraction.	43
8.1	Overview of <i>Deep Compression</i> [3].	49
8.2	Overview of the proposed compression pipeline.	49

List of Tables

3.1	Number of parameters in LeNet5-431K.	10
3.2	Number of parameters in CifarNet.	11
4.1	Number of parameters of pruned LeNet5-431K.	15
4.2	Number of parameters of pruned <i>Cifarnet</i>	15
4.3	Number of parameters of pruned LeNet5-431K, without pruning biases.	16
4.4	Number of parameters of pruned CifarNet, without pruning biases.	16
4.5	Number of parameters of pruned LeNet5-431K using Dynamic network surgery	18
4.6	Number of parameters of pruned CifarNet using Dynamic network surgery.	19
4.7	LeNet5 pruning summary, CR is the compression rate, ER is the error rate. (<i>Han</i>) is deterministic pruning used by <i>Han et al.</i> (<i>Guo</i>) is the original <i>Dynamic network surgery</i> implemented by <i>Guo et al.</i> (a), (b), (c), (d) are my implementations of various methods. (a) is deterministic pruning with biases, (b) is deterministic pruning without biases, (c) is <i>Dynamic network surgery</i> , (d) is also <i>Dynamic network surgery</i> but with the same error rate as (<i>Guo</i>).	19
4.8	CifarNet pruning summary, CR is the compression rate, ER is the error rate. (a) is deterministic pruning with biases, (b) is deterministic pruning without biases, (c) is <i>Dynamic network surgery</i>	20
5.1	Number of parameters of pruned LeNet5-431K using <i>Gradient profiling</i>	26
5.2	Number of parameters of pruned LeNet5-431K, with $l1$ and $l2$ norms, $\lambda_1 = 1e^{-4}$ and $\lambda_2 = 1e^{-7}$	28

5.3	Number of parameters of pruned CifarNet, with $l1$ and $l2$ norms, $\lambda_1 = 1e^{-5}$ and $\lambda_2 = 1e^{-5}$	28
5.4	LeNet5 Pruning Summary, CR is the compression rate, ER is the error rate. (a) is deterministic pruning with biases, (b) is deterministic pruning without biases, (c) is <i>Dynamic network surgery</i> , (d) is <i>Gradient Profiling</i> and (e) is <i>Regularization Aided Pruning</i>	29
5.5	CifarNet Pruning Summary, CR is the compression rate, ER is the error rate. (a) is deterministic pruning with biases, (b) is deterministic pruning without biases, (c) is <i>Dynamic network surgery</i> and (d) is <i>Regularization Aided Pruning</i>	29
6.1	Fixed-point quantization summary for <i>LeNet5</i> and <i>CifarNet</i> . .	33
6.2	Fixed-point quantization summary for <i>LeNet5</i> and <i>CifarNet</i> . .	35
6.3	Dynamic fixed-point quantization summary for <i>LeNet5</i> and <i>CifarNet</i>	37
7.1	Customized floating-point quantization summary for <i>LeNet5</i> and <i>CifarNet</i> , with 2-bit exponent and various mantissa width. .	40
7.2	Customized floating-point quantization summary for <i>LeNet5</i> , with 3-bit mantissa and various exponent width.	40
7.3	Quantization summary on a pruned <i>LeNet5</i> model. <i>Customized floating-point</i> has 1-bit sign, 1-bit exponent and the rest are mantissa bits. <i>Centralized customized floating-point</i> has 1-bit sign, 1-bit central, 1-bit exponent and the rest are mantissa bits. <i>Centralized dynamic fixed-point</i> has 1-bit sign, 1-bit central and the rest are fraction bits.	44
7.4	Quantization summary on a pruned <i>CifarNet</i> model. <i>Customized floating-point</i> has 1-bit sign, 2-bit exponent and the rest are mantissa bits. <i>Centralized customized floating-point</i> has 1-bit sign, 1-bit central, 2-bit exponent and the rest are mantissa bits. <i>Centralized dynamic fixed-point</i> has 1-bit sign, 1-bit central and the rest are fraction bits.	45
7.5	Quantization summary (bit-width, error rate) on unpruned <i>LeNet5</i> and <i>CifarNet</i> . <i>DFP</i> is <i>Dynamic fixed-point</i> , <i>CFP</i> is <i>Customized floating-point</i> and (<i>Gysel</i>) is <i>Gysel et al.</i> 's implementation of <i>Dynamic fixed-point</i>	46
7.6	Quantization summary (bit-width, error rate) on pruned <i>LeNet5</i> and <i>CifarNet</i> . <i>CFP</i> is <i>Customized floating-point</i> <i>CCFP</i> is <i>Centralized customized floating-point</i> and <i>CDFP</i> is <i>Centralized dynamic fixed-point</i>	47

8.1	LeNet5 Pruning Summary, Pruning and Quantization, CR is the compression rate, ER is the error rate.	50
-----	--	----

Chapter 1

Introduction

Neural Networks have achieved outstanding accuracies in large-scale image classifications and they are becoming a state-of-art technique for solving problems in computer vision [4, 1, 5].

The rise of deep neural network becomes the key for this success. In 1998, *Lecun et al.* proposed a 5-layer network called *LeNet5* for recognising handwritten digits [4]. *Lecun et al.* utilised around 1M parameters to achieve a good accuracy on this particular recognition task. *AlexNet*, designed by *Krizhevsky et al.* for the 2010 ImageNet competition, targeted on a much harder task of categorising 1.2M images into 1000 categories and made use of around 60M parameters [1]. *Coates et al.* scaled up the learning algorithm to utilise over 11B parameters and ran it on 16 machines to recognise unlabelled human faces [6]. It is reasonable to forecast that, in the future, neural networks are becoming deeper and are utilizing a larger number of parameters.

While GPUs stay to be an efficient hardware platform for accelerating neural network learning in a large scale due to their large bandwidths and memory sizes, deploying a network with a large number of parameters on a memory-limited power-sensitive device becomes increasingly difficult. To constraint the problem space, this project focuses on neural network inference rather

than training. There are some recent proposed algorithms, such as *Deep Q Learning* [7] and *Asynchronous Q Learning* [8], that require training of a neural network on a local embedded system. This need of executing network training in local embedded systems is still only essential for a small set of learning algorithms. The more appealing problem currently is how to execute network inference in an efficient manner. To execute network inference, although only forward propagation takes place, mobile systems struggle to achieve a reasonable power efficiency. Two major constraints prevent neural network inference from being applicable in real-life. First, the large size of a neural network makes storing its parameters in embedded systems very challenging. For example, the Alexnet model is over 200MB and VGG-16 model is over 500MB [3]. These large storage overheads post a challenge to the fundamental memory size and memory bandwidth of embedded systems. Second, energy consumption is dominated by memory access [9], and accessing a large number parameters of a neural network can exceed the energy envelope of many power sensitive devices. For instance, the battery of a smartphone struggles with running object classification using AlexNet in real-time for more than an hour [9].

To resolve the above problems, the computer architecture community is now actively researching on novel hardware architectures for neural network inference. Many novel custom hardware architectures have been proposed for neural network inference [10, 11, 12, 13]. Various FPGA-based accelerators have been recently applied on neural network inference [14, 15]. In addition, there is an increasing number of ASIC designs for deep neural network inference [12, 16]. These hardwares normally utilize a large on-chip memory and have custom computing units to calculate matrix dot-products. A custom hardware is definitely beneficial for running neural network inference efficiently, but accessing and storing the large number of parameters in the memory is still a limit for these hardwares.

To build a hardware accelerator, the first step is to consider how to condense the neural network to a reasonable size. A compressed neural network is more amenable to memory-limited systems and also reduces the energy cost of data

movements. Neural network compression is recently catching attentions and many research works have made significant contributions in this field. I would like to systematically summarize state-of-art neural network compression techniques, compare and evaluate various compression techniques. These compression techniques are mainly categorized into three groups: pruning, regularization and quantization. Some of these methods, such as pruning and regularization, encourage sparsity in a neural network; in other words, these methods reduce the number of parameters in a neural network. Other methods, such as quantization, takes a different approach to reduce the size of a neural network: they reduce the number of bits required to represent a single parameter in a neural network. In this report, the aim is to combine several state-of-art compression techniques to build a complete compression pipeline. The compression pipeline is applicable to various neural networks, and outperforms the compression rates of many existing works.

Chapter 2

Background

2.1 Convolutional Neural Networks

Neural networks have been primarily inspired by the model of the actual human brain neural system. Neurons are basic units in a neural network and each neuron provides an output based on its connections (weights) to neurons of the previous layer and an additional bias value. Neural networks are modelled as layers of neurons, in particular, a neural network for image recognition normally has a few convolutional layers and a few fully connected layers. A typical neural network architecture is illustrated in Figure 2.1, in this case, this neural network has five convolutional layers and three fully-connected layers [1].

For a single neuron, consider w_i to be a vector of weights, b_i to be a single bias and x_i to be the input vector of a particular neuron i , this neuron generates the following output $y_i = w_i x_i + b_i$. This mathematical model of a single neuron feeding forward is the same as the model of a linear classifier [17]. As illustrated in Figure 2.2, the basic setup of neurons can be seen as a mathematical model of actual neurons in the human brain system. In a neural network, each output y_i of a single neuron has to go through an activation function, this added non-linearity becomes a key for neural network to achieve

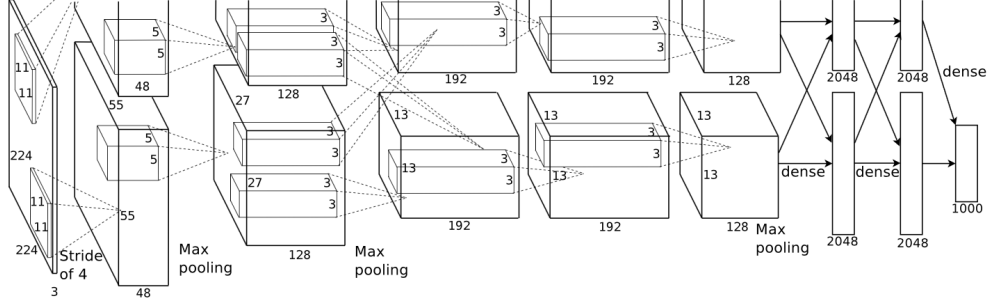


Figure 2.1: AlexNet Network Architecture [1].

good performance. Popular activation functions include *sigmoid* ($f(x) = \frac{1}{1+e^{-x}}$), *tanh* ($f(x) = \tanh(x)$) and *Relu* ($f(x) = \max(0, x)$).

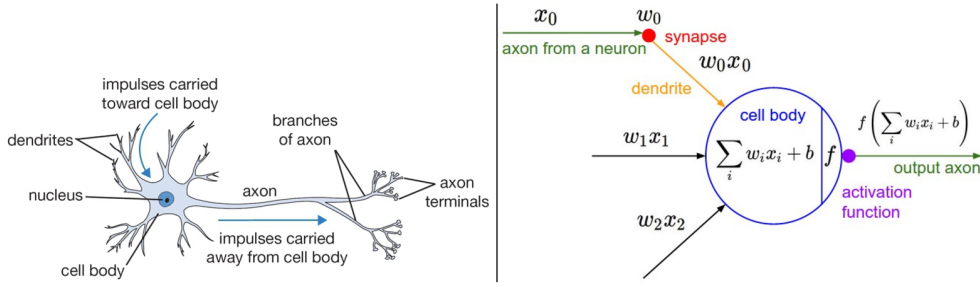


Figure 2.2: Fully connected layer followed by an activation function [1].

As mentioned previously, a number of neurons consists a layer, and a neural network is a layer-wise architecture. A fully connected layer, as its name states, has neurons that are fully connected to all neurons in the previous layer and neurons in a single layer do not share any connections. However, using only fully connected layers lose particular information for image recognition. Convolutional layers requires neurons to be connected to a local region of the previous layer and this helps to capture specific features of a given image. Convolutional layers normally have 3D volumes of neurons [1]. Figure 2.1 shows a typical neural network structure for image recognition. The first few layers are convolutional layers followed by maxpooling layers to extract features from images. As mentioned before, these convolutional layers are all 3D volumes of neurons. The last convolutional layer then flattens

to connect to fully connected layers. These fully connected layers, although are 2D, but normally have a large number of neurons, for instance, the shown figure of AlexNet has 2048 neurons on its first fully connected layer [1].

The current research trend suggests convolutional neural networks will use more layers for achieving a higher accuracy on more complicated problems [4, 1, 6].

2.2 Related Work

A variety of methods have been proposed in the research community to condense neural networks. In general, most of these compressing techniques focus on whether reducing the number of parameters of a neural network or the number of bits used to represent every individual parameter.

Network pruning has been widely used to reduce the number of parameters in convolutional neural networks [4, 18, 19]. Starting from *LeCun et al.*, their idea is to use the second derivative of the loss function to balance the prediction accuracy and model complexity [4]. *Hassibi et al.* proposed to add in the non-diagonal elements of the Hessian matrix into the pruning metric, this is proven to be helpful in reducing excess degrees of freedom in the neural network model [18]. Recently, *Han et al.* utilised pruning and achieved high compression rate on many popular convolutional neural networks [3]. Their pruning technique is tightly coupled with training, after every pruning process, a retraining takes place to bring back the lost accuracies. This class of pruning is also called *Iterative Pruning*. Inspired by *Han et al.*'s work, later on, *Guo et al.* implemented *Dynamic network surgery* that combines pruning with splicing. In the splicing phase of a network surgery, some weights that was previously pruned away would have a chance to recover. *Dynamic network surgery* is able to prune more aggressively and achieved better results than traditional iterative pruning [2]. Apart from network pruning, there are other existing techniques that focus on reducing the number of parameters in a neural network model. These methods focus on encouraging sparsity

of the neural network while training the network. It has long been known that adding $L1$ or $L2$ normalisation brings sparsity into a neural network. However, these regularisers only restrict the magnitudes of weights and thus induces sparsity in a less controlled manner. *Srinivas et al.* proposed a novel regulariser that induce sparsity in a more controlled manner. *Kang et al.* proposed *Shakeout*, which is a new training scheme. The activation functions are customised at each neuron, and their practical results suggest that *Shakeout* induces more sparsity than traditional regularisation methods.

Reducing the number of bits required to represent each individual parameter also directly compresses a given neural network. *Han et al.* used weight sharing to group parameters with similar values [3]. Weights sharing normally considers weights in a layer-wise fashion. For weights in each layer, a clustering algorithm is applied on these weights to group them with various centroid values [3]. These centroids are encoded using a hash function. For later inference operation, only centroids values are stored on-chip and each weight is represented using a hash key that is normally small in terms of bit-width. This idea origins from the *HashNet* proposed by *Chen et al.* [20], but is different from the original work since now the hash function applies on a trained network. Quantization is another popular method for reducing the bit-width of individual parameters. Fixed-point arithmetic has been confirmed as a more energy efficient arithmetic for neural network inference [21]. Quantized networks, also known as low precision networks, utilizes low-precision fixed-point arithmetic to reduce the bit-width of individual parameters in the neural network [22]. Reaching an extreme, neural networks with parameters of only 1-bit width are known as *Binarized Neural Networks* [23]. However, *Binarized Neural Networks* normally require re-implementations to recover its accuracy loss and thus is beyond the scope of this project. For reducing the size of individual parameters, in this project, I will focus on various quantization methods and weight sharing.

Although a large number of compressing techniques have been proposed, these techniques are almost all isolated. Recently, *Han et al.* firstly proposed *Deep Compression* that is a complete compression pipeline utilizing

several different compression techniques. It is possible to build a compression chain that makes use of compression techniques from orthogonal optimization space. *Deep Compression* offers a large compression rate on many modern networks using *deterministic pruning*, *fixed-point quantization*, *weights sharing* and *huffman encoding*. This project aims to evaluate further in this compression optimization space. A larger number of state-of-art compression techniques, including different pruning schemes and quantization schemes are considered for building a better compression pipeline. In addition, this project puts a focus on developing novel pruning and quantization strategies based on some existing compressing techniques.

Chapter 3

Experiment Setups, Datasets and Trained Networks

TensorFlow is the python package used to implement the neural networks [24], it is an open source library for machine learning developed by Google. This package has python APIs that are easy to use, and *Cuda* backend for efficient GPU utilizations. In this project, the lowest level API – *TensorFlow Core* is used, because it provides fine levels of control. Several local GPU machines and Amazon Cloud machines are used to train the networks.

Two datasets, *MNIST* [25] and *CIFAR10* [26], are considered in this project. The first dataset, *MNIST*, consists of handwritten digits. *MNIST* has a training set of 60,000 images and a test set of 10,000 images. These images are a subset of the bigger *NIST* dataset, and all the digits in this dataset have been normalized in size and centered in the middle [25]. The neural network used to recognize *MNIST* is a LeNet5 model [27].

Another bigger dataset considered in this project is the *Cifar10* data set. *Cifar-10* is a subset of the *80 million tiny images dataset* [28]. The *Cifar-10* dataset is divided into five training batches and one test batch, each batch contains 10,000 images [26]. *CifarNet* is the neural network architecture used to recognize images in the *Cifar-10* dataset [29]. It is important to note

that, all these datasets are carefully designed so that all classes are mutually exclusive.

Layer	cov1	cov2	fc1	fc2	total
Params	0.5K	25K	400K	5K	431K

Table 3.1: Number of parameters in LeNet5-431K.

Firstly, I am using a trained *LeNet5* model with 431K parameters that has an error rate of only 0.64%. The original *LeNet5* model proposed by *LeCun et al.* utilizes around 1000K parameters, however, later research works all use a Caffe implementation of *LeNet5* which has only 431K parameters [30, 3, 2]. In this project, the 431K implementation of *LeNet5* is chosen in order to make fair comparisons with other existing research works. The detailed information regarding the number of parameters at each layer of this *LeNet5* network is shown in Table 3.1. The network’s architecture is shown in Figure 3.1. As shown in the figure, this architecture includes two convolutional layers and three fully connected layers.

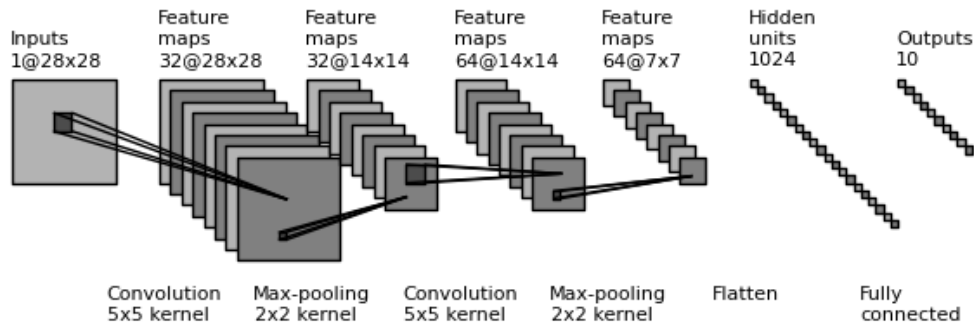


Figure 3.1: LeNet5 Network Architecture

The *Cifar10* dataset is a larger dataset compared to *MNIST*. A *CifarNet* architecture is constructed and trained to an accuracy of 0.82. As expected, since this larger dataset implies harder image recognition task, the trained model uses a larger number of parameters and more layers to achieve a reasonable accuracy. A graphical illustration of the *CifarNet* architecture is shown in Figure 3.2, The *CifarNet* architecture has two convolutional layers

and three fully connected layers. The input images of the dataset now have three channels. The number of parameters at each layer of the *CifarNet* is summarized in Table 3.2.

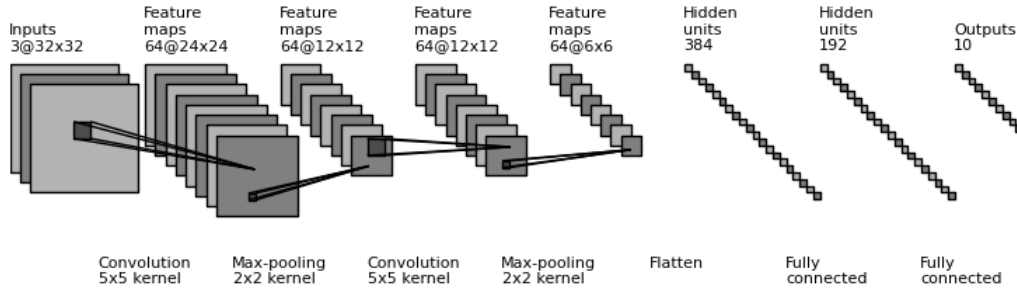


Figure 3.2: CifarNet Network Architecture

Layer	cov1	cov2	fc1	fc2	fc3	total
Params	4.8K	102.4K	885K	74K	2K	1068K

Table 3.2: Number of parameters in CifarNet.

Chapter 4

Pruning

Pruning is an effective method for reducing the redundancies in a network, it inspects all parameters in a layer and moves away connections that have values lower than a certain threshold. Pruning is normally combined with iterative training so that the lost accuracies can be recovered directly. Most existing pruning strategies apply on networks in a deterministic manner [3]. Once a connection is pruned away, it never has a chance to recover. Recently, researchers have discovered more efficient pruning strategies that prune in a non-deterministic fashion [2]. In this section, I would like to implement a range of pruning methods and summarize their performances on selected neural network models.

4.1 Deterministic Pruning

4.1.1 Pruning with biases

Pruning reduces the total number of parameters and finds the minimal neural network topology, but overly pruned network suffers from a significant accuracy drop and even divergence [31]. Consider a set of weights for a neural network of N layers, the weights can be expressed as $\{W_n : 0 \leq n \leq N\}$,

so that W_n is the weights for a given layer n . Similarly, biases are expressed as $\{B_n : 0 \leq n \leq N\}$. To represent a sparse model with pruned weights, a binary matrix $\{M_n : 0 \leq n \leq N\}$ for weights and a binary matrix $\{M_n^b : 0 \leq n \leq N\}$ for biases are used to indicate connections that are kept. For simple deterministic pruning, Equation (4.1) shows how this binary mask is computed using the weights. An arbitrary threshold value t_n is used to determine whether a certain weight variable should be pruned away. I use $h_n(*)$ to denote the discriminative function that produces a binary mask matrix based on the weight matrix.

$$h_n(W_n^{(i,j)}) = \begin{cases} 0, & \text{if } t_n < |W_n^{(i,j)}| \\ 1, & \text{if } t_n \geq |W_n^{(i,j)}| \end{cases} \quad (4.1)$$

The appealing question now is how to determine the threshold value t_n . Instead of picking threshold values in an arbitrary fashion, motivated by the implementation of *Dynamic network surgery* [2], the following metric is used for determining the pruning threshold t_n :

$$t_n = u_n + c\sigma_n \quad (4.2)$$

u_n is the mean value of parameters in layer n , σ_n is the standard deviation of the parameters in layer n . In Equation (4.2), c is a hyperparameter used to define how aggressive this pruning is. This threshold determination works better than defining an arbitrary threshold value because it avoids inspecting the values of parameters at each pruning iteration.

In a sparse neural network model, weights of each layer n can be easily expressed using an element-wise product between W_n and M_n . The loss function, expressed as $L_n(*)$, is a metric that measures the optimisation target of the neural network training phase. For a given layer of n , the

following equations set the optimization objective:

$$\begin{aligned}
& \min(L_n(W_n \odot M_n + B_n \odot M_n^b)) \\
& M_n = h_n(W_n) \\
& M_n^b = h_n(B_n)
\end{aligned} \tag{4.3}$$

For the math simplicity, the above optimization target is only for a given layer n , a complete model optimization target is different from this single layer target. Nonetheless, Equation (4.3) characterizes a sparse model. The \odot symbol represents hadamard product, which is a binary operation that takes two matrices and produces another matrices that has elements are the product of the elements of the original two matrices. In deterministic pruning, pruned weights cannot recover their values. However, deterministic pruning occurs iteratively – a number of weights is pruned away at each iteration and then retraining occurs to bring back the lost test accuracy. Although pruned weights do not have a chance to recover their values, other existing weights would have a chance to re-learn the correlations between existing connections using Equation (4.3) as a learning target. The weights are updated concurrently by the pruned model with a learning rate of α , and Equation (4.4) shows how weights are updated [2]:

$$W_n \leftarrow W_n - \alpha \frac{\partial}{\partial(W_n \odot M_n + B_n \odot M_n^b)} L(W_n \odot M_n + B_n \odot M_n^b) \tag{4.4}$$

Figure 4.1 shows how pruning affects the weight distribution. In these plots, pruned weights with a value of zero are not plotted. The plots focus on the first fully connected layer of a *LeNet5* model. The unpruned network has weights that are normally distributed and center at zero. The weights spread out to both positive and negative ends, and the values of weights are normally small. In contrast, the second plot shows how pruned weights look like: they gathered at two centers and leaving a blank region around zero.

Table 4.1 shows the detailed pruning information at each layer of the *LeNet5*. Pruning reduces the number of parameters of *LeNet5* architecture without

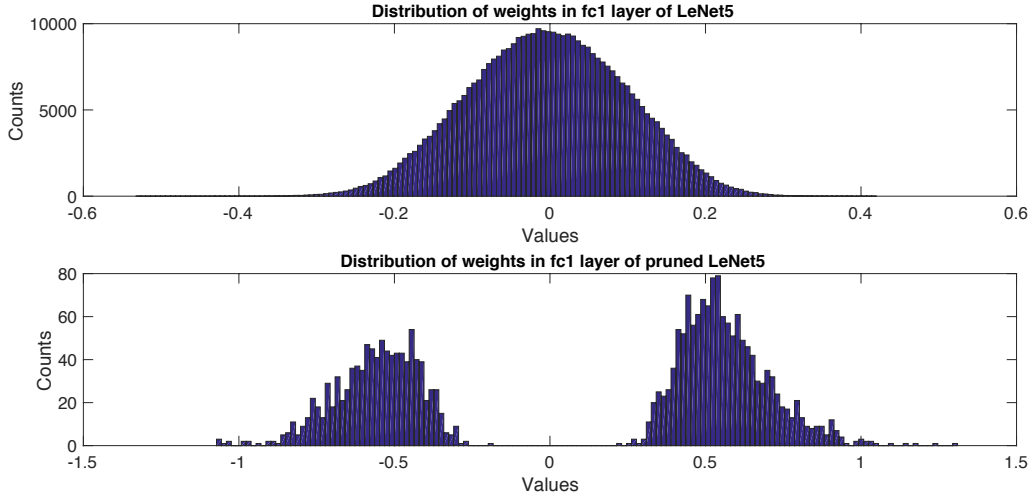


Figure 4.1: Distribution of weights in fc1 layer of pruned and unpruned LeNet5

Layer	cov1	cov2	fc1	fc2
Params	0.5K	25K	400K	5K
Prune %	40%	80%	98%	40%

Table 4.1: Number of parameters of pruned LeNet5-431K.

dropping in its test accuracy. The compressed network is only 3.79% of the size of the original network, which indicates a compression rate of 26.38x.

Table 4.2 shows how the pruning strategy performs on a *CifarNet*. This pruning strategy performs iterative pruning including pruning biases, and the pruned network has the same test accuracy as the original network. As shown in Table 4.2, each layer of the *CifarNet* achieves a lower pruning percentage compared to *LeNet5*. This is a reasonable result, because recognizing images in *MNIST* is an easier task and the *LeNet5* network contains more redundancies than *CifarNet*.

Layer	cov1	cov2	fc1	fc2	fc3
Params	4.8K	102.4K	885K	74K	2K
Prune %	30%	66%	85%	66%	30%

Table 4.2: Number of parameters of pruned *Cifarnet*.

4.1.2 Pruning without biases

The number of biases of a neural network is normally small compared to the number of weights. More importantly, the biases serve as offset values for each individual neuron, these offsets stay invariant to different input images and could have significant impacts on the accuracies of a neural network. It is therefore reasonable to consider pruning only the weights but leaving the biases unpruned.

Similar to the previous section, I compute the layer-wise optimization target of pruning without biases in Equation (4.5).

$$\min(L(W_n \odot M_n + B_n)) \quad (4.5)$$

As the equation states, now the binary mask matrix only applies on weights variables, the biases are kept unchanged. Table 4.3 and Table 4.6 show the pruning results of *LeNet5* and *CifarNet* respectively.

Layer	cov1	cov2	fc1	fc2
Params	0.5K	25K	400K	5K
Prune %	51.92%	79.84%	99.38%	49.90%

Table 4.3: Number of parameters of pruned LeNet5-431K, without pruning biases.

Layer	cov1	cov2	fc1	fc2	fc3
Params	4.8K	102.4K	885K	74K	2K
Prune %	40%	69%	85%	69%	40%

Table 4.4: Number of parameters of pruned CifarNet, without pruning biases.

The results of pruning without biases demonstrate itself to be a more efficient pruning strategy compared to the previous methodology in . Numerically, the compression rate of *LeNet5* when pruned with biases is 26.38x, and now it increases to 37.75x. Similarly, the compression rate of *CifarNet* has increased from 6.10x to 6.19x.

To conclude, pruning without biases give better results and later on in this project, all pruning methods only apply on the weights variables of the neural networks.

4.2 Dynamic Network Surgery

Dynamic network surgery is a pruning method proposed by *Guo et al.* recently [2]. *Guo et al.* proved experimentally that *Dynamic network surgery* outperforms other existing pruning methods by considerable margins [2]. This method combines pruning with splicing. Splicing refers to a procedure that some pruned weights are selected randomly to join the next pruning iteration. The complete procedure of *Dynamic network surgery* is shown in Figure 4.2.

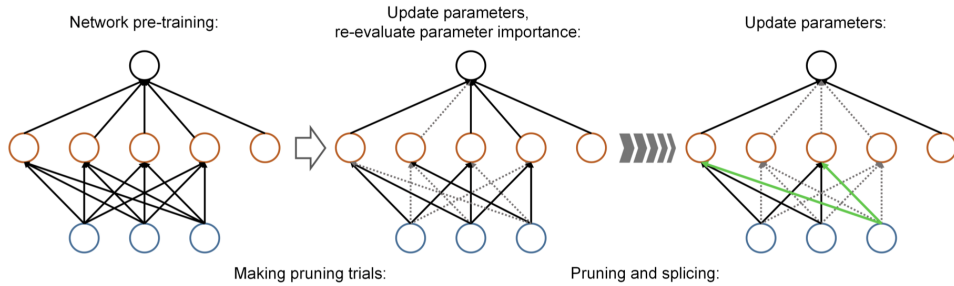


Figure 4.2: Mechanism of Dynamic network surgery [2].

In the previous sections, all pruning strategies are deterministic, meaning that networks might suffer from irretrievable damages once some important weights are lost. In contrast, *Dynamic network surgery* provides a chance for weights to recover. As shown on the rightmost plot in Figure 4.2, the green connections are weights that have been recovered by splicing. The splicing process minimizes the risk of causing an irretrievable damage to the neural network. Consider the parameters defined previously, the major difference

now is on the discriminative function:

$$h_n(W_n^{(i,j)}) = \begin{cases} 0, & \text{if } t_{an} > |W_n^{(i,j)}| \\ M_n, & \text{if } t_{an} \leq |W_n^{(i,j)}| \leq t_{bn} \\ 1, & \text{if } t_{bn} < |W_n^{(i,j)}| \end{cases} \quad (4.6)$$

The discriminative function now have two threshold values, namely t_{an} and t_{bn} for a given layer n . Different from the original discriminative function, for values that stay in between two thresholds, their corresponding binary masks stay unchanged. For weights that have absolute values larger than t_{bn} , their binary masks become one which means these weights are recovered. Weights with absolute values smaller than t_{an} are turned off.

For determining the threshold values, t_{an} and t_{bn} , *Guo et al.* used the following equations:

$$\begin{aligned} t_{an} &= 0.9(u_n + c\sigma_n) \\ t_{bn} &= 1.1(u_n + c\sigma_n) \end{aligned} \quad (4.7)$$

Similar as before, Equation (4.7) made use of the mean and standard deviation at each layer to help determine the pruning thresholds. To test the performance of *Dynamic network surgery*, it is then applied on the two selected neural networks (*LeNet5* and *CifarNet*). The pruning results of both *LeNet5* and *CifarNet* are displayed in Table 4.5 and Table 4.6 respectively.

Layer	cov1	cov2	fc1	fc2
Params	0.5K	25K	400K	5K
Prune %	36.54%	87.90%	99.66%	18.42%

Table 4.5: Number of parameters of pruned LeNet5-431K using Dynamic network surgery

The use of *Dynamic network surgery* significantly improves the results of pruning. The compression rate of *LeNet5* is now 49.05x, and the compression rate of *CifarNet* is 17.66x. The significant increases in compression rates imply that *Dynamic Network Surgery* is a better pruning method. Pruning happens in a nondeterministic fashion in *Dynamic Network Surgery*, if a

Layer	cov1	cov2	fc1	fc2	fc3
Params	4.8K	102.4K	885K	74K	2K
Prune %	53%	87%	95%	82%	26%

Table 4.6: Number of parameters of pruned CifarNet using Dynamic network surgery.

pruned weight finds its importance at later stages of the iterative pruning, it will have a chance to recover. In deterministic pruning, although retraining occurs iteratively, it still causes irretrievable damage to the targeting network.

4.3 Comparison to Existing Works

In this section, I would like to compare the implemented pruning methods to their original implementations.

Model	Layer	Params	%(<i>Han</i> [3])	%(<i>Guo</i> [2])	%(a)	%(b)	%(c)	%(d)
LeNet5	cov1	0.5K	66%	14.2%	60%	48.1%	63.5%	10%
	cov2	25K	12%	3.1%	20%	20.2%	12.1%	6%
	fc1	400K	8%	0.7%	2%	0.6%	0.4%	5.05%
	fc2	5K	19%	4.3%	60%	50.1%	82.6%	5.84%
	total	431K	8%	0.9%	3.8%	2.4%	2.0%	0.89%
	CR	-	12.5x	108x	26x	42x	49x	111x
	ER	-	0.8%	0.91%	0.64%	0.64%	0.64%	0.91%

Table 4.7: LeNet5 pruning summary, CR is the compression rate, ER is the error rate. (*Han*) is deterministic pruning used by *Han et al.*. (*Guo*) is the original *Dynamic network surgery* implemented by *Guo et al.*. (a), (b), (c), (d) are my implementations of various methods. (a) is deterministic pruning with biases, (b) is deterministic pruning without biases, (c) is *Dynamic network surgery*, (d) is also *Dynamic network surgery* but with the same error rate as (*Guo*).

Table 4.7 shows how the original implementations of various pruning methods compared to my implementations. (a) is deterministic pruning including pruning biases, (b) is deterministic pruning without pruning biases. (c) is my implementation of the *Dynamic network surgery* method. (*Han*) is the

deterministic pruning strategy used by *Han et al.* in their *Deep Compression* framework [3]. (*Guo*) is the original implementation of *Dynamic network surgery* [2]. The first few rows show the percentages of parameters that are left at each layer. The "total" row displays the total percentage of parameters that are left in the entire neural network. The row "CR" shows the compression rates achieved using various techniques and "ER" illustrates the final test accuracies of the pruned networks. The compression rates achieved in my implementations are generally better than *Han et al.* both in terms of compression rate and test accuracies. Comparing to the original implementation of *Dynamic network surgery*, (c) shows a smaller compression rate but also lower error rate. Since the targeting network and dataset are the same, this drop in compression rate can be seen as a trade-off between compression rate and error rate. When I choose to tolerate an error rate of 0.91%, as shown in (d), the compression rate reaches 111x, which is very close to the original implementation proposed by *Guo et al.* (108x).

Model	Layer	Params	(a)	(b)	(c)
CifarNet	cov1	4.8K	30%	40%	53%
	cov2	102.4K	66%	69%	87%
	fc1	885K	85%	85%	95%
	fc2	74K	66%	69%	82%
	fc3	2K	30%	40%	26%
	total	1068K	3.8%	2.4%	2.0%
	CR	-	5.41x	5.58x	14.3x
	ER	-	0.8%	0.8%	0.8%

Table 4.8: CifarNet pruning summary, CR is the compression rate, ER is the error rate. (a) is deterministic pruning with biases, (b) is deterministic pruning without biases, (c) is *Dynamic network surgery*.

Table 4.8 shows the pruning results on *CifarNet*, all the pruning methods are implemented by myself. (a) is deterministic pruning including pruning biases, (b) is deterministic pruning without pruning biases. (c) is my implementation of the *Dynamic network surgery* method. *CifarNet* does not have any comparable research works in pruning. This comparison also demonstrates that *Dynamic network surgery* is superior to other pruning strategies.

In this section, by comparing various pruning strategies, two important observations can be summarized. First, biases of a neural network should not be pruned since they are only offset values of each neuron. Second, *Dynamic network surgery* is proven to have the best performance, indicating that deterministic pruning can cause damages to neural networks.

Chapter 5

Proposed Pruning Strategies

In this chapter, I will mainly describe two novel pruning strategies. The first proposed pruning technique is called *Gradient profiling*, this method benefits pruning with limited retraining resources. The second method is *Regularization aided pruning*. By putting extra regularization terms into a cost function, later in this chapter, I demonstrate that regularizers encourage sparsity of a neural network and thus induce better pruning results.

5.1 Gradient Profiling

All pruning methods implemented in the previous chapter only focus on pruning parameters based on their values. These implementations, however, have a common problem: weights with small values are ignored, although ignoring them could potentially hurt test accuracies. It is possible for a weight with small value to stay at an important location so it might have a large impact on accuracy. In this section, I would like to propose a new method called *Gradient profiling* to inspect weights that are small in values but important for test accuracies.

5.1.1 Method Description

The *Gradient profiling* method is nearly identical to *Dynamic network surgery*, apart from that it inspects the gradient changes for one epoch. The reason for inspecting all gradients of one epoch is to find enough observations on the entire training dataset. The description of the method is illustrated in Figure 5.1.

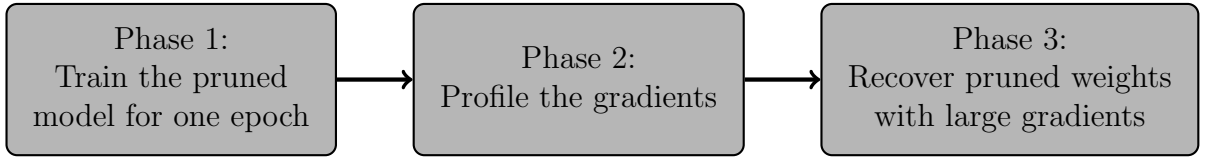


Figure 5.1: *Gradient Profiling*.

The method is broken down into three phases, in the first phase, a pruned neural network is trained only for one epoch. The second phase collects the gradients from that one epoch training. The third phase is to recover pruned weights based on profiled gradients data. The amount of pruned weights that are recovered is controlled by a arbitrarily defined hyperparamter. In this section, the amount of recovered weights is equal to 10% of the total weights remained. The hypothesis is that the gradients of a network serve as indications for weights importance. If an important weight is pruned away, the hypothesis is that the neural network would keep passing this weight a large gradient. Later in the gradient profiling mechanism, pruned weights with large gradients are recognized as important weights and recovered.

From a mathematical point of view, let $\{G_n : 0 \leq n \leq N\}$ to represent the profiled gradients from phase 2 in Figure 5.1. A discriminative function $h_{gn}(\cdot)$ is used to determine a new mask $\{GM_n : 0 \leq n \leq N\}$ for each layer.

$$GM_n = h_{gn}(G_n^{(i,j)}) = \begin{cases} 0, & \text{if } t_{gn} < |W_n^{(i,j)}| \\ 1, & \text{if } t_{gn} \geq |W_n^{(i,j)}| \end{cases} \quad (5.1)$$

t_{gn} is a threshold value determined arbitrarily, as mentioned before, this value

is set to include the top 10% of the total weights remained in the targeting layer. The loss function therefore becomes:

$$\min(L_n(W_n \odot M_n + W_n \odot GM_n + B_n)) \quad (5.2)$$

The training strategy follows the same change:

$$W_n \leftarrow W_n - \alpha \frac{\partial}{\partial(W_n \odot M_n + W_n \odot GM_n + B_n)} L(W_n \odot M_n + W_n \odot GM_n + B_n) \quad (5.3)$$

5.1.2 Gradient Profiling and Retraining

To fully understand the effect of *Gradient profiling*, I would like to compare it to *Dynamic network surgery* on various level of retraining. It is common to combine pruning methods with retraining, however, some researchers argue that retraining takes a significant amount of time and they want to avoid it [32]. So, in this section, I consider the following three retraining strategies:

1. No retrain after pruning.
2. Retrain for 10 epochs after pruning.
3. Retrain for 300 epochs after pruning.

Figure 5.2 is the results when both *Gradient profiling* and *Dynamic network surgery* are applied on *LeNet5* without any retraining. The horizontal axis displays the amount of weights that are left compared to the original network. The vertical axis shows the test accuracies when the network is compressed to various sizes. As expected, because of the fact that some important weights are recovered by *Gradient profiling*, it shows a greater test accuracy than *Dynamic network surgery* when the size of the neural network is compressed below 0.4 of its original size. When the size of the neural network is relatively large (size above 0.5), the loss of test accuracies are not apparent and therefore both methods show similar performances.

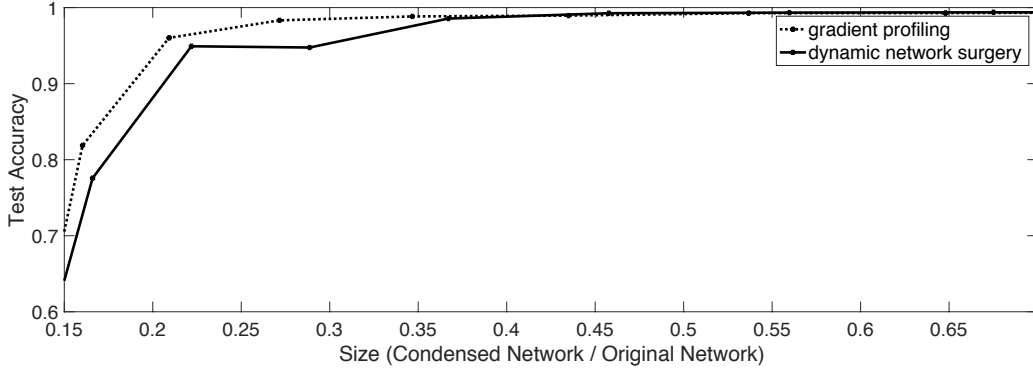


Figure 5.2: *Gradient profiling* and *Dynamic network surgery* without retraining.

It is obvious that retraining recovers the lost test accuracies by forcing the remaining parameters to learn to adapt with each other. However, the amount of time spent on retraining is a major drawback of this method. It is interesting to observe how pruning strategies are combined with *limited retraining*. I define *limited retraining* as a concept that only a very small number of retraining epochs are allowed after pruning. In this case, I pick to retrain only 10 epochs after pruning. Two different pruning methods, *Gradient profiling* and *Dynamic Network Surgery*, are considered and compared. Figure 5.3 shows the results when 10 epochs of retraining is applied. As expected, the test accuracy of *Gradient Profiling* drops at a slower rate compared to *Dynamic network surgery*. It is clear that between the size of 0.2 and 0.4, the test accuracies of *Gradient Profiling* is significantly higher.

Finally, I combine pruning with retraining of 300 epochs, and tested them on *LeNet5*. Table 5.1 shows the pruning results when using *Gradient profiling*. The compression rate is 48x, which is slightly lower than the compression rate of *Dynamic network surgery* (49x). This proves that, if retrained to convergence, *Gradient profiling* is not superior to *Dynamic network surgery*. The large amount of retraining time give weights ability to learn to adapt with each other. The large weights now have enough time to learn, and previously important small weights might lose their importance because large valued weights have already had enough time to adjust their values.

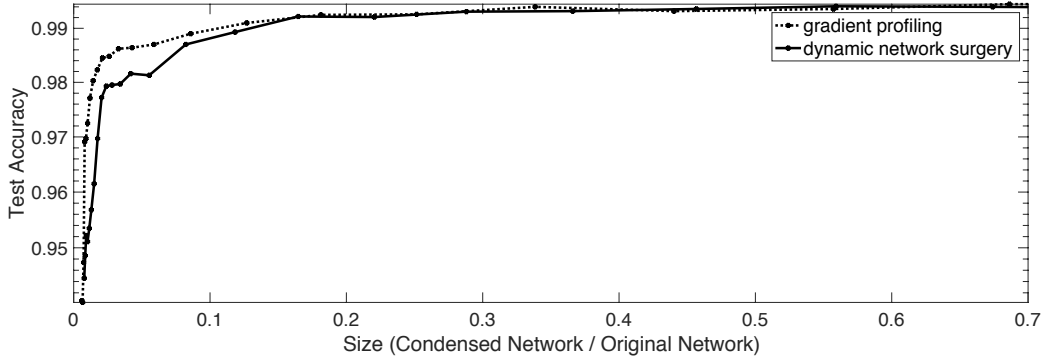


Figure 5.3: *Gradient profiling* and *Dynamic network surgery* with 10 epochs retraining.

Layer	cov1	cov2	fc1	fc2
Params	0.5K	25K	400K	5K
Prune %	44.20%	91.08%	99.31%	28.87%

Table 5.1: Number of parameters of pruned LeNet5-431K using *Gradient profiling*

To summarize, *Gradient profiling* can be a very efficient pruning strategy when retraining resource is limited. Its enhancement on compression rate becomes limited if long retraining time is allowed. In this project, since the topic is to construct a compression pipeline that achieves the best compression rate, I used *Dynamic network surgery* rather than *Gradient profiling* in later sections. However, the importance and effectiveness of *Gradient Profiling* under limited retraining resources still worths further evaluations.

5.2 Regularization Aided Pruning

Regularization methods are popular for preventing the neural networks from overfitting. Some regularization methods, such as l_1/l_2 norms and *Shakeout*, achieve regularization by encouraging sparsities in the model. The use of regularizers leads to a sparse model, and this might benefit the pruning process. In this section, I would like to investigate how regularization could potentially lead to efficient pruning.

5.2.1 l_1 and l_2 Norms

l_1 and l_2 norms, also known as least absolute errors (LAE) and least squares respectively, are popular regularization terms that could be appended to a neural network's cost function [33]. Mathematically, the l_p norm of a given vector v :

$$||v||_p = \left(\sum_{i=1}^n |v_i|^p \right)^{\frac{1}{p}} \quad (5.4)$$

l_1 and l_2 norms have been proven to be robust to outliers in data and also encourages sparsity via training [33]. The intuitive explanation is to view these additional norms as penalties to large weights, so they prevent particular weights from dominating network. For the combined l_1 and l_2 norms, I add the following term into the cost function of a neural network:

$$\lambda_1|w| + \lambda_2|w|^2 \quad (5.5)$$

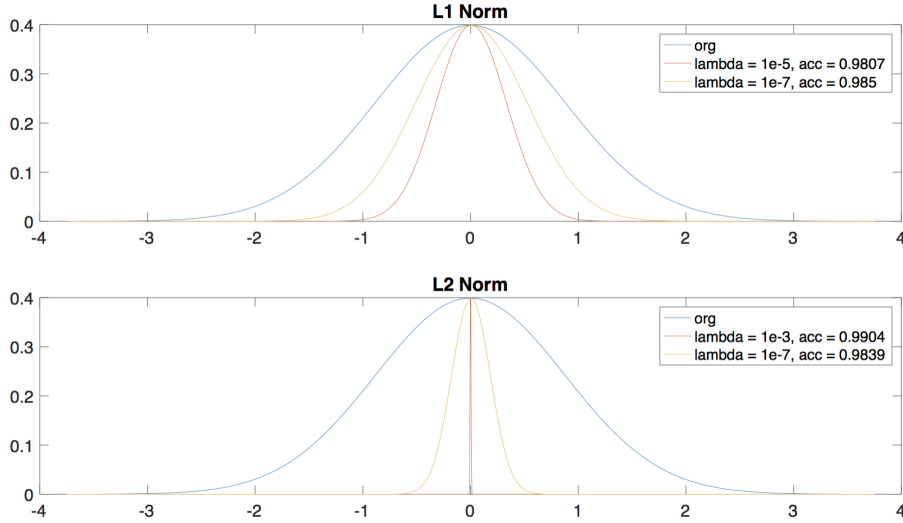


Figure 5.4: Effects of l_1 and l_2 norms with different hyperparameters.

λ_1 and λ_2 are two hyperparameters that are normally chosen arbitrarily. To view the effect of various norms on the weights distribution of a given neural network, Figure 5.4 shows how the weights distribution of a *LeNet5* model

varies when applied with different norms. With larger λ_1 and λ_2 values, the weights distributions are more concentrated to zeros.

Notice in Figure 5.4, l_1 and l_2 norms are shown separately, meaning that λ_2 and λ_1 in Equation (5.5) are set to zeros respectively when chosen to demonstrate the effect of l_1 norm and l_2 norm respectively. For choosing l_1 to be $\lambda_1 = 1e^{-4}$ and $\lambda_2 = 1e^{-7}$, a LeNet5 model is trained to the error rate of 0.64%, which is the same error rate that has been achieved in the precious section. The detailed layer-wise pruning information is shown in Table 5.2.

Layer	cov1	cov2	fc1	fc2
Params	0.5K	25K	400K	5K
Prune %	22.11%	92.47%	99.70%	32.91%

Table 5.2: Number of parameters of pruned LeNet5-431K, with l_1 and l_2 norms, $\lambda_1 = 1e^{-4}$ and $\lambda_2 = 1e^{-7}$.

For *CifarNet*, $\lambda_1 = 1e^{-5}$ and $\lambda_2 = 1e^{-5}$ are selected after exploring a range of values. The layer-wise pruning results of *CifarNet* are shown in Table 5.3.

Layer	cov1	cov2	fc1	fc2	fc3
Params	4.8K	102.4K	885K	74K	2K
Prune %	45%	88%	96%	76%	39.6%

Table 5.3: Number of parameters of pruned *CifarNet*, with l_1 and l_2 norms, $\lambda_1 = 1e^{-5}$ and $\lambda_2 = 1e^{-5}$.

From the results in both Table 5.2, Table 5.3, *Regularization aided pruning* demonstrates its performance by showing greater compression rates on both networks. It can be concluded that the use of regularizers encourage sparsity in the network and thus provide better pruning results.

5.2.2 Summary of Pruning Methods

In this section, I would like to summarize all pruning methods implemented in this project. Some pruning methods haven been compared to their original

implementations in Section 4.3, this section only focuses on comparing my implementations of various pruning strategies on selected datasets.

Model	Layer	Params	(a)	(b)	(c)	(d)	(e)
<i>LeNet5</i>	cov1	0.5K	60%	48.1%	63.5%	65.8%	78.9%
	cov2	25K	20%	20.2%	12.1%	8.9%	7.5%
	fc1	400K	2%	0.6%	0.4%	0.7%	0.3%
	fc2	5K	60%	50.1%	82.6%	72.2%	67.1%
	total	431K	3.8%	2.4%	2.0%	2.0%	1.6%
	CR	-	26x	42x	49x	48x	63x
	ER	-	0.64%	0.64%	0.64%	0.64%	0.64%

Table 5.4: LeNet5 Pruning Summary, CR is the compression rate, ER is the error rate. (a) is deterministic pruning with biases, (b) is deterministic pruning without biases, (c) is *Dynamic network surgery*, (d) is *Gradient Profiling* and (e) is *Regularization Aided Pruning*.

Model	Layer	Params	(a)	(b)	(c)	(d)
CifarNet	cov1	4.8K	30%	40%	53%	45%
	cov2	102.4K	66%	69%	87%	88%
	fc1	885K	85%	85%	95%	96%
	fc2	74K	66%	69%	82%	76%
	fc3	2K	30%	40%	26%	40%
	total	1068K	3.8%	2.4%	2.0%	1.4%
	CR	-	5.41x	5.58x	14.3x	15.4x
	ER	-	0.8%	0.8%	0.8%	0.8%

Table 5.5: *CifarNet* Pruning Summary, CR is the compression rate, ER is the error rate. (a) is deterministic pruning with biases, (b) is deterministic pruning without biases, (c) is *Dynamic network surgery* and (d) is *Regularization Aided Pruning*.

Similar to the previous setup, Table 5.4 shows the pruning results. The percentiles showing for each layer represents the amount of parameters left in that layer. To prune *LeNet5*, it is important to achieve a small percentage on the first fully-connected layer (fc1), since it contains a large number of parameters. Method (a) is deterministic pruning with biases, (b) is deterministic pruning without biases, (c) is *Dynamic network surgery*, (d) is *Gradient profiling* and (e) is *Regularization aided pruning*. The proposed

pruning strategy, *Regularization aided pruning*, achieves the best compression rate: it shows a 1.3x increase in compression rate compared to *Dynamic network surgery*.

The following important observations can be summarized from comparing a range of pruning methods:

1. Pruning without biases is more efficient.
2. Nondeterministic pruning achieves better compression rates since pruned weights now have chance to recover.
3. Gradients can serve as indications for identifying important weights.
4. The use of regularizers helps pruning by encouraging sparsity in the network, and thus *Regularization aided pruning* shows the best pruning results.

Chapter 6

Quantization

One objective of this project is to discover a more efficient number representation system. An efficient number representation system helps to reduce the number of bits required to represent individual parameters, and thus offers a compression on the top of pruning. In this section, I aim to investigate various existing quantization methodologies to further compress the neural networks. Traditionally, a neural network is trained on GPUs in 32-bit floating-point representation, however, this representation becomes problematic on low power devices. First, using 32 bits to store a single parameter is proven to be redundant [34]. Second, floating-point arithmetic operations are generally more power-consuming than fixed-point operations. Methods such as *Weights sharing* would be discussed in this section [3]. Low-precision arithmetics, including *Fixed-point arithmetic*, *Dynamic fixed-point arithmetic* will also be compared and evaluated in this section.

6.1 Weights Sharing

Similar to *HashNet* [20] and *Deep Compression* [3], weights sharing can reduce the number of bits required to represent a parameter by employing a hash function. Designing the hash function firstly requires grouping of

weights. In this case, I performed *K-means* clustering [35] to group weights into n clusters, and therefore n centroid values are produced.

The inference and retraining of the grouped weights are illustrated in Figure 6.1. Consider a small four-by-four matrix, and weights are grouped into 4 groups. Only the cluster indexes need to be stored for network inference. For this particular example, to represent 4 clusters, each weight can be represented using only a 2-bit cluster index. This directly compressed the size of a neural network, and *Han et al.* has summarized an equation for this compression rate [3]. Given k clusters, $\log_2(k)$ bits are required for encoding the index. Consider a neural network with N parameters, and each parameter is b bits wide, the following equation summarizes the compression rate CR :

$$CR = \frac{Nb}{N \log_2(k) + kb} \quad (6.1)$$

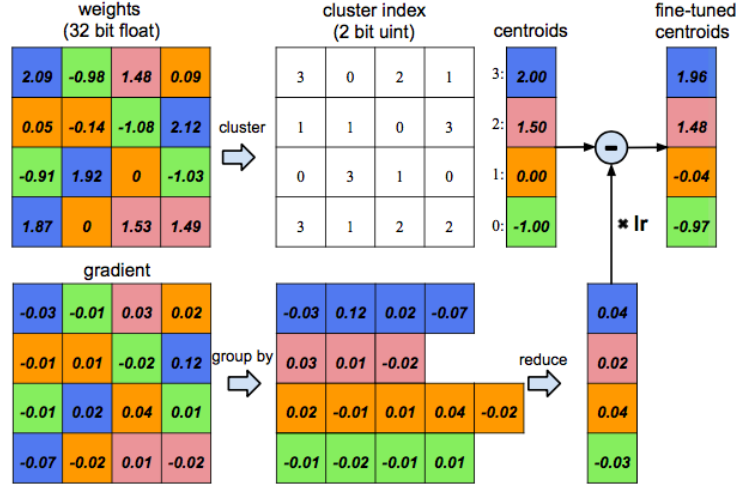


Figure 6.1: Weights sharing: training (bottom) and inference (top) [3].

This *Weights sharing* technique is applied on both *LeNet5* and *CifarNet*. Table 6.1 shows the quantization results of *Weights sharing* on the targeting neural networks. *Weights sharing* is also combined with retraining to bring back the lost test accuracies. Table 6.1 uses two rows for each neural network, one row represent the test accuracy before training and one row shows post

training accuracies. First, as expected, test accuracies are significantly larger when the number of clusters is large. Second, retraining is demonstrated to be an efficient method of bringing back the lost accuracies. *Weights sharing* worked very well on *LeNet5*, with 8 clusters at each layer, *LeNet5* achieves no accuracy loss after retraining. It is important to note that, for 8 clusters, only 3 bits are required to represent the cluster index of each parameter.

Number of clusters	64	32	16	8	4	2
<i>LeNet5</i> , 32-bit floating point accuracy: 99.36%						
Before Retrain	99.36%	99.36%	99.36%	99.29%	98.90%	94.67%
After Retrain	99.36%	99.36%	99.36%	99.36%	99.15%	98.66%
<i>CifarNet</i> , 32-bit floating point accuracy: 82.00%						
Before Retrain	79.94%	79.81%	72.78%	69.00%	25.25%	9.94%
After Retrain	82.21%	82.01%	82.02%	79.05%	65.70%	17.2%

Table 6.1: Fixed-point quantization summary for *LeNet5* and *CifarNet*.

CifarNet is a larger network compared to *LeNet5*. As a result of its larger size, the network requires a larger number of clusters at each layer. According to Table 6.1, at a cluster count of 16, the network achieves no accuracy loss. This means, each weight parameter can be represented using 4 bits, and the neural network is able to achieve the same accuracy as the uncompressed one.

Although weights sharing demonstrated good compression rates on both networks, the arithmetic operations are still in floating-point. Since floating-point arithmetic operations are power consuming, turning into small fixed-point numbers can reduce both energy consumption and circuitry area. In later sections, the focus stays on exploring fixed-point based quantization methods.

6.2 Fixed-point Quantization

One simple quantization strategy when using fixed-point arithmetic is to truncate the least significant bits (LSBs) but leave the most significant bits

(MSBs) unchanged. This method is straightforward and easy to implement.

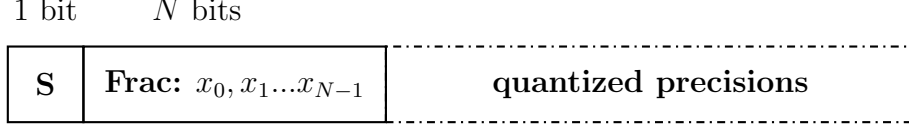


Figure 6.2: Number representation system for fixed-point quantization with n -bit fraction.

Figure 6.2 shows the number representations of a quantized number. In this example, the number is truncated to N -bit fractional bits. The number representation system has 1 bit for sign and n bits for fractions. The number of fractional bits can be therefore changed to track numbers to various level of precisions. Equation (6.2) shows how the number representation of Figure 6.2 can be converted to decimal representation. In this case, S stands for the sign bit, N is the total number of fractional bits. x_i represents the fractional bit values. The representation used in this section is a traditional signed fixed-point representation with N -bit fractions [36].

$$D = (-1)^S \left(\sum_{i=0}^{N-1} 2^{-i} x_i \right) \quad (6.2)$$

Table 6.2 shows how fixed-point quantization performs on *LeNet5* and *Cifar-Net* when the bit-width changes. It is important to note that, the bit-width in this case refers to total bit-width. Total bit-width equals to $N + 1$, where N is the number of fractional bits and 1 bit for integer.

The results in Table 6.2 suggests that 16 bits for *LeNet5* achieves best quantization results without any accuracy loss. 32-bit is redundant since the test accuracy remains the same as 16-bit when retraining is applied. Similar to the previous setup, retraining occurs after quantization to bring back test accuracies. Retraining takes a very important role in fixed-point quantization. For instance, when *LeNet5* is quantized to 4-bit (3-bit fraction and 1-bit sign), its test accuracy before retraining is only 12.74%, but increases

Bit width	32-bit	16-bit	8-bit	4-bit	2-bit
<i>LeNet5</i> , 32-bit floating point accuracy: 99.36%					
Before Retrain	98.98%	97.08%	75.63%	12.74%	9.78%
After Retrain	99.36%	99.36%	99.26%	98.88%	9.78%
<i>CifarNet</i> , 32-bit floating point accuracy: 82.00%					
Before Retrain	32.84%	21.58%	12.16%	10.22%	9.97%
After Retrain	82.00%	79.21%	68.76%	65.76%	9.95%

Table 6.2: Fixed-point quantization summary for *LeNet5* and *CifarNet*.

to 98.88% after retraining. As illustrated in Table 6.2, *CifarNet* only recovers to an accuracy of 79.21% when the precision is 16-bit and 82.00% when the precision is 32-bit. Since *CifarNet* is a larger neural network compared to *LeNet5*, it requires more bits to retain its test accuracy.

6.3 Dynamic Fixed-point Quantization

Dynamic fixed-point arithmetic is a variant of fixed-point arithmetic, where a number of parameters are grouped with a fixed scaling factor [37]. Previously, Courbariaux *et al.* has implemented both *LeNet5* and *CifarNet* using *Dynamic fixed-point* arithmetic [38]. They proposed an algorithm (Algorithm 1) for computing the scaling factors. The scaling factors define the dynamic ranges for a group of weights.

To use *Dynamic fixed-point* in a neural network, Algorithm 1 is applied on each layer of parameters to help weights with various dynamic ranges. The number representation system is the same as *fixed-point* arithmetic shown in Figure 6.2. The major different now is each group of weights has a fixed dynamic range. In Equation (6.3), it shows how a *Dynamic fixed-point* number can be converted to a decimal value. DR represents the fixed dynamic range, consider the scaling factor (s_t) computed in Equation (6.3), the following relation holds: $s_t = 2^{DR}$. These dynamic ranges, or can be also viewed as scaling factors, are stored globally for each layer and thus saves the bit-width for the number representation system.

$$D = (-1)^{s} 2^{DR} \left(\sum_{i=0}^{M-1} 2^{-i} x m_i \right) \quad (6.3)$$

Algorithm 1 Scaling Factor Update

```

1: Initialize: matrix  $M$ , scaling factor  $s_t$  and maximum overflow rate  $r_{max}$ 
2: while  $s_t$  do not converge do
3:   if the overflow rate of  $M > r_{max}$  then
4:      $s_{t+1} \leftarrow 2s_t$ 
5:   else if the overflow rate of  $2M \leq r_{max}$  then
6:      $s_{t+1} \leftarrow s_t/2$ 
7:   else
8:      $s_{t+1} \leftarrow s_t$ 
9:   end if
10: end while

```

Dynamic Fixed-point quantization has been applied on both *LeNet5* and *CifarNet*, the results are shown in Table 6.3. For each bit-width count, it includes the sign bit but excludes the dynamic range. If a bit-width of 4 is given, this means it has 1 bit for sign, 3 bits for fractional and a constant dynamic range that is globally defined for a given layer. Bit-width refers to the number of total bits required to represent a number using *Dynamic Fixed-point*, and the dynamic range is stored as a layer-wise global value. As expected, *LeNet5* reaches reasonable test accuracies at a bit-width of 4, and *CifarNet* achieves no accuracy loss after retrain at a bit-width of 8. These results agree with the experimental results achieved by *Courbariaux et al.* [38] and *Gysel et al.* [39]. However, the test accuracy of retrained *LeNet5* never reaches its original test accuracy (99.36%). This indicates that *Dynamic network surgery* hurts test accuracy in a way that even retraining cannot bring back the lost test accuracies. In both *Courbariaux et al.* and *Gysel et al.*'s work, they used a *LeNet5* with an original test accuracy of 99.17% [38, 39], so they did not observe this effect. Because of the high original test accuracy, my implementation of *LeNet5* is more sensitive to test accuracy deterioration.

By observing the dynamic range of fixed-point quantization, I can practically

Bit width	32-bit	16-bit	8-bit	4-bit	2-bit
<i>LeNet5</i> , 32-bit floating point accuracy: 99.36%					
Before Retrain	88.40%	88.32%	86.92%	79.10%	71.26%
After Retrain	99.29%	99.30%	99.31%	99.31%	99.00%
<i>CifarNet</i> , 32-bit floating point accuracy: 82.00%					
Before Retrain	58.99%	26.25%	10.48%	9.98%	9.98%
After Retrain	82.00%	82.12%	82.03%	72.18%	30.53%

Table 6.3: Dynamic fixed-point quantization summary for *LeNet5* and *CifarNet*.

find the reason for this loss of accuracy caused by using *Dynamic fixed-point*. Consider the scaling factor determined using Algorithm 1, this scaling factor restricts the range of values that this arithmetic could represent. For instance, a scaling factor of 0.5 is determined for the first fully connected layer by inspecting the pre-quantized *LeNet5* model. However, for a *LeNet5* model that is quantized using fixed-point arithmetic and achieved no accuracy loss, its maximum value in that layer reaches 0.406 but minimum value reaches -0.55 . In this case, if *Dynamic fixed-point* is applied, values that are lower than -0.5 will never be reached. The hypothesis is that the loss of dynamic range can potentially affect retraining since now the quantized weights is restricted to certain numerical ranges. This hypothesis is not theoretically proven, but is hinted by comparing both the test accuracies and numerical ranges of two different arithmetics. A theoretical proof of this hypothesis is beyond the scope of this project, but could be evaluated in future works.

Chapter 7

Proposed Quantization Strategies

In the previous chapter, a number of existing quantization methods have been implemented and evaluated. In this chapter, I would like to propose some new quantization methods for compressing neural networks. First, a *Customized Floating-point* arithmetic is proposed and it overcomes certain limitations of *Dynamic Fixed-point* arithmetic and also showed better performance compared to *Dynamic Fixed-point*. Second, to quantize a pruned model, I applied some bi-centralized arithmetics to fully explore the weights distribution by re-centralizing the arithmetics.

7.1 Customized Floating-point Quantization

In this section, I propose a new number representation system that stands in the middle between *Dynamic fixed-point* arithmetic and *Floating-point* arithmetic.

In *Dynamic fixed-point*, each number is represented as: $(-1)^S 2^{-DR} \sum_{i=0}^{M-1} 2^i x_i$. M denotes the mantissa width, s is the sign bit, DR is the fractional length and x are the mantissa bits [38, 39]. At various layers of the neural network,

Dynamic fixed-point assigns fixed-point numbers with different layer-wise *DR* values. This *DR* value stays unchanged for the grouped weights in one layer. In contrast, *Floating-point arithmetic* has a changing exponent for each individual weight, and therefore the decimal point is "floating". A normal floating-point number can be expressed as $(-1)^S 2^{exp-offset} \sum_{i=0}^{M-1} 2^i x_i$. In this case, an *offset* value is included to help the arithmetic to cover a large range of values [36].

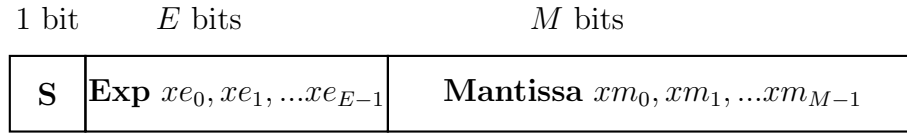


Figure 7.1: Number representation system for customized floating-point quantization with e-bit exponent and m-bit mantissa.

Different from both arithmetics, the customized floating-point arithmetic only needs to track values smaller than 1, so the *offset* value in normal floating-point is eliminated. Consider the number representation system in Figure 7.1, the following equation summarizes its conversion to a decimal value:

$$D = (-1)^S 2^{-\sum_{i=0}^{E-1} 2^{-i} xe_i} \left(\sum_{i=0}^{M-1} 2^{-i} xm_i \right) \quad (7.1)$$

In Equation (7.1), D represents expected decimal value, S is the sign bit, M is the mantissa width and E is the exponent width. xm_i are the mantissa bits and xe_i are the exponent bits.

As shown in Equation (7.1), each individual range now can take its own dynamic range, so this is different from *Dynamic fixed-point* where only fixed dynamic range is allowed for a group of parameters. More importantly the dynamic range now is restricted to be lower than 1, this makes sure the arithmetic tracks to lowest possible precisions given limited bit-width.

Table 7.1 shows the quantization results. Using only 6 bits, *Customized*

floating-point is able to generate a retrained *LeNet5* that has no accuracy loss. For *CifarNet*, a bit-width of 8 is required for it to recover to the original accuracy. In Table 7.1, the bit-width includes both a sign bit and a 2-bit exponent. The exponent is chosen to be 2 bits for both *LeNet5* and *CifarNet*. For instance, given a bit width of 6 bits, it consumes 1 bit for sign, 2 bits for exponent and only 3 bits are left for mantissa.

Bit width	32-bit	16-bit	8-bit	6-bit	4-bit
<i>LeNet5</i> , 32-bit floating point accuracy: 99.36%					
Before Retrain	99.36%	99.36%	99.31%	98.4%	9.79%
After Retrain	99.36%	99.36%	99.36%	99.36%	11.35%
<i>CifarNet</i> , 32-bit floating point accuracy: 82.00%					
Before Retrain	34.54%	25.20%	19.15%	15.17%	10.35%
After Retrain	82.02%	82.00%	82.10%	79.82%	70.73%

Table 7.1: Customized floating-point quantization summary for *LeNet5* and *CifarNet*, with 2-bit exponent and various mantissa width.

The results in Table 7.1 demonstrate a good compression rate, however, using *Customized Floating-point* arithmetic, the width of exponent has to be determined beforehand. Table 7.2 shows how varying the exponent width can affect test accuracies.

Exponent bit width	4-bit	3-bit	2-bit	1-bit	0-bit
<i>LeNet5</i> , 32-bit floating point accuracy: 99.36%					
Before Retrain	98.41%	98.41%	98.40%	97.91%	9.79%
After Retrain	99.36%	99.36%	99.36%	99.36%	11.35%

Table 7.2: Customized floating-point quantization summary for *LeNet5*, with 3-bit mantissa and various exponent width.

The exponent bit width has a locally optimal value of 1 bit when the mantissa width is set to 3. Taking the sign bit into account, the arithmetic would utilize 5 bits for recovering to the original accuracy. In addition, unlike the *Dynamic Fixed-point* method used in Section 6.3, this method is able to span all numerical ranges and thus the test accuracy is able to recover fully.

7.2 Re-centralized Quantization

All previous quantization methods focused on an unpruned *LeNet5* model, in this section, I would like to propose a new quantization method specifically designed for quantizing pruned neural networks. I first show the motivation of this proposed method by comparing to a previous implementation. I then describe the mechanism and show the quantization results.

Motivation

Previously, all quantization methods focused on unpruned networks and proved good compression rates. The high precision region of these applied arithmetics focuses on the center of the weights distribution.

Color coded arithmetic precisions

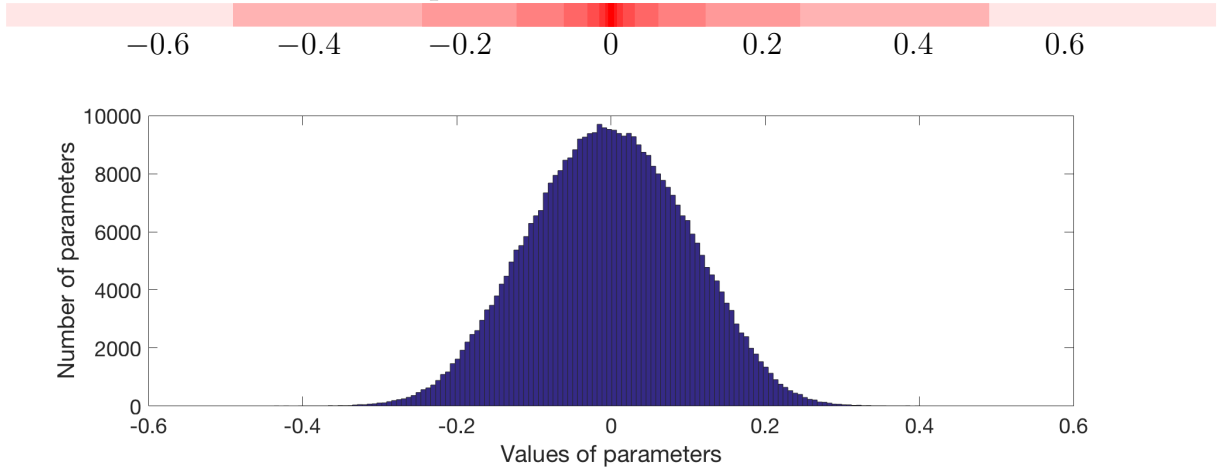


Figure 7.2: Parameter distribution of the first fully connected layer in *LeNet5*, and the color coded arithmetic precisions using *Customized floating-point quantization* with 1 bit sign, 1 bit exponent and 3 bits fraction. A deeper red color corresponds to a greater arithmetic precision.

Given an example in Figure 7.2, the chosen arithmetic has the ability to represent numbers more precisely in the region where most weights stay on. The color coded bar is deepest at values near zero, and the histogram shows its peak at zero as well. This partially explains why quantizations with

dynamic ranges work well, because now greater precisions can be explored near zeros. However, for the pruned model, this advantage of using dynamic ranges in arithmetics disappears.

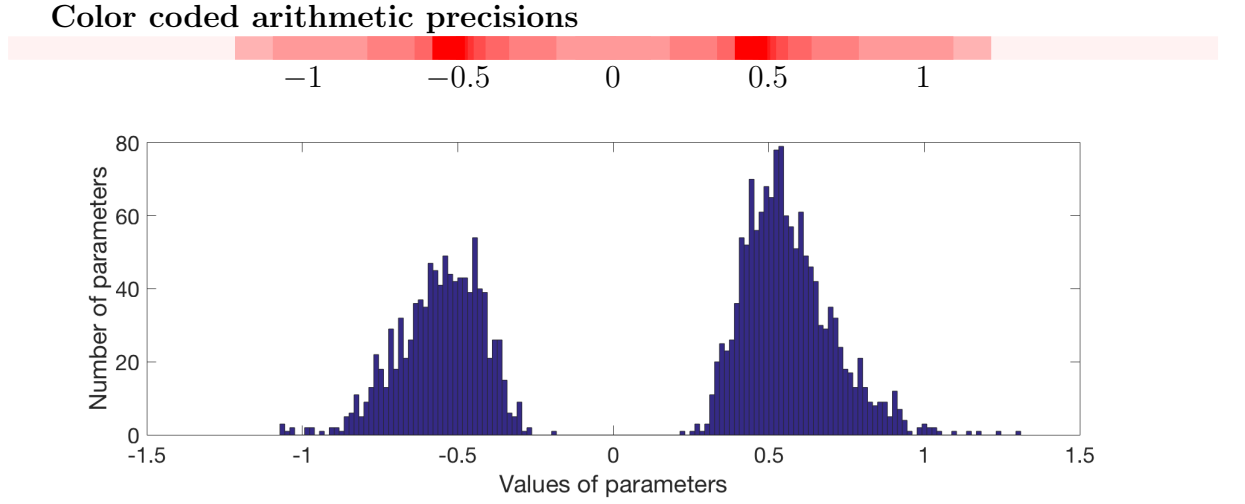


Figure 7.3: Parameter distribution of the first fully connected layer in *LeNet5*, and the color coded arithmetic precisions using *Centralized customized floating-point quantization* with 1 bit sign, 1 bit centre, 1 bit exponent and 3 bits fraction. A deeper red color corresponds to a greater arithmetic precision.

As shown in Figure 7.3, the weights distribution has a binomial shape, if the same arithmetic used before applies on this weights distribution, the middle empty region will have the richest number representations. The waste in number representation suggests that a suitable arithmetic for pruned models should have the ability to represent numbers that are away from zeros in high precisions. The color coded arithmetic precision bar suggests the preferred precision intensities. The arithmetic should have high precision on two center points that are away from zero, and have low precision representation for weights that are near zero because these weights have been pruned away already.

Method Description

The methods introduced in this section is called *Re-centralization*. The shape of weights distribution of a pruned model is binomial, however, its central values of two peaks can be different for every layer. In this method, the weights parameters are inspected before quantization occurs and two central values are determined layer-wise. Two central values correspond to the two peaks of the binomial distribution and are stored globally for parameters in the same layer. For each layer, two central values are recorded and I use 1 bit (C bit) to help a single parameter to determine which central value it is associated with. This re-centralization technique is applied on both *Dynamic fixed-point* and *Customized floating-point*.

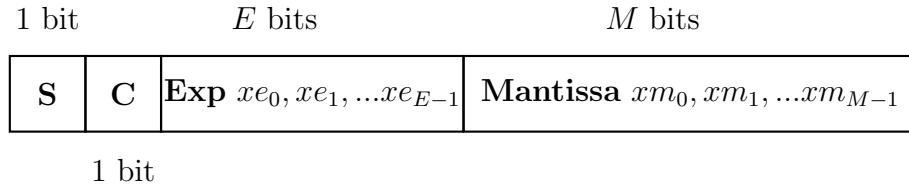


Figure 7.4: Number representation system for *Re-centralized customized floating-point quantization* with E-bit exponent and M-bit mantissa.

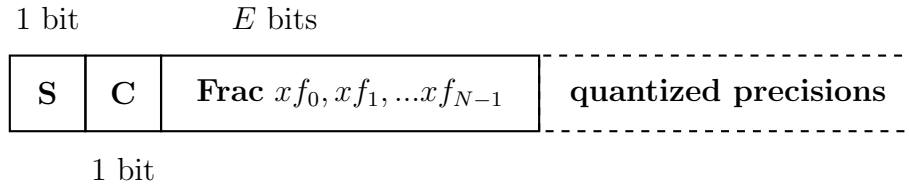


Figure 7.5: Number representation system for *Re-centralized dynamic fixed-point quantization* with N-bit fraction.

A *Re-centralized customized floating-point* number representation system is displayed in Figure 7.4, the extra bit C is also shown to determine which centric value this parameter belongs to. Similarly, a *Re-centralized dynamic fixed-point* arithmetic is shown in Figure 7.5. To convert from this *Re-centralized customized floating-point* arithmetic to a normal decimal value, consider two layer-wise central values to be C_{pos} and C_{neg} . The bit stored to

identify the central value is C , we have:

$$D = (-1)^S 2^{-\sum_{i=0}^{E-1} 2^{-i} x e_i} \left(\sum_{i=0}^{M-1} 2^{-i} x m_i \right) + (C C_{pos} + (1 - C) C_{neg}) \quad (7.2)$$

Similarly, *Re-centralized customized floating-point* arithmetic converts to the decimal value following the equation below:

$$D = (-1)^S 2^{-DR} \left(\sum_{i=0}^{N-1} 2^{-i} x f_i \right) + (C C_{pos} + (1 - C) C_{neg}) \quad (7.3)$$

Results

To demonstrate the performance gain from re-centralization, pruned models of both *LeNet5* and *CifarNet* are considered. I show a comparison between *Customized floating-point quantization*, *Centralized dynamic fixed-point* and *Centralized customized floating-point* to observe the effect of re-centralizations. The results of quantization on a pruned *LeNet5* are shown in Table 7.3, and Table 7.4 shows the quantization results on pruned *CifarNet*.

Bit width	32-bit	16-bit	8-bit	6-bit	4-bit
<i>LeNet5, Customized floating-point</i>					
Before Retrain	99.36%	99.36%	99.31%	98.40%	73.94%
After Retrain	99.36%	99.36%	99.36%	99.36%	99.36%
<i>LeNet5, Centralized dynamic fixed-point</i>					
Before Retrain	99.04%	99.04%	99.06%	98.93%	91.86%
After Retrain	99.36%	99.36%	99.36%	99.36%	99.36%
<i>LeNet5, Centralized customized floating-point</i>					
Before Retrain	98.73%	98.72%	98.67%	98.04%	21.81%
After Retrain	99.36%	99.36%	99.36%	99.36%	99.36%

Table 7.3: Quantization summary on a pruned *LeNet5* model. *Customized floating-point* has 1-bit sign, 1-bit exponent and the rest are mantissa bits. *Centralized customized floating-point* has 1-bit sign, 1-bit central, 1-bit exponent and the rest are mantissa bits. *Centralized dynamic fixed-point* has 1-bit sign, 1-bit central and the rest are fraction bits.

As shown in Table 7.3, *Customized floating-point* arithmetic has similar post-train quantization results to both *Centralized customized floating-point* and *Centralized dynamic fixed-point* on *LeNet5*. However, the pre-trained results of using *Centralized dynamic fixed-point* is actually best one. The limited mantissa width causes *Centralized customized floating-point* to have a bad pre-train accuracy. The pre-train accuracy of this simple network is tightly linked to the number of mantissa/fractional bits available. The underlying issue is that *LeNet5* itself is too simple, it seems like a large precision near the central values would not introduce critical precision improvements when no retraining takes place.

Bit width	32-bit	16-bit	8-bit	6-bit
<i>CifarNet, Customized floating-point</i>				
Before Retrain	34.30%	29.19%	30.58%	20.56%
After Retrain	82.02%	82.01%	81.85%	81.12%
<i>CifarNet, Centralized dynamic fixed-point</i>				
Before Retrain	81.18%	81.17%	80.06%	73.23%
After Retrain	82.02%	82.06%	82.05%	82.10%
<i>CifarNet, Centralized customized floating-point</i>				
Before Retrain	59.65%	59.61%	58.23%	44.94%
After Retrain	82.01%	82.11%	82.00%	81.79%

Table 7.4: Quantization summary on a pruned *CifarNet* model. *Customized floating-point* has 1-bit sign, 2-bit exponent and the rest are mantissa bits. *Centralized customized floating-point* has 1-bit sign, 1-bit central, 2-bit exponent and the rest are mantissa bits. *Centralized dynamic fixed-point* has 1-bit sign, 1-bit central and the rest are fraction bits.

Table 7.4 shows the quantization results when various quantization methods are applied on the pruned *CifarNet*, *Centralized Customized Floating-point* and *Centralized dynamic fixed-point* shows a significant improvement in pre-train accuracy compared to *Customized floating-point*. Conceptually, centralized arithmetics re-focus the arithmetic centers to fit the binomial weights distributions, and this should increase the precisions near central values. Using *Centralized customized floating-point*, it achieves the required accuracy at a precision of only 6 bits.

7.3 Summary of Quantization Methods

To summarize, I would like to compare the implemented quantization methods on the targeting networks.

Consider unpruned *LeNet5* and *CifarNet* as targeting networks, Table 7.5 shows the bit-widths of various networks after quantizations. As mentioned previously, *Dynamic fixed-point* (*DFP*) restricts the range of weights and thus some weights can never have a chance to recover. As shown in the table, the error rate of using *Dynamic fixed-point* is slightly higher on *LeNet5* due to this effect. In contrast, *CifarNet* shows comparable quantization results for both *Dynamic fixed-point* and *Customized floating-point* (*CFP*). *Gysel et al.* obtained their best quantization results using *Dynamic fixed-point* on the unpruned *LeNet5* and *CifarNet*, I added their results to Table 7.5.

For quantizing unpruned models, *Customized floating-point* shows best results on the two selected networks: it showed a good compression rates on both networks without any loss of test accuracies.

Model	<i>Fixed-point</i>	<i>DFP</i>	<i>CFP</i>	(<i>Gysel</i>)
<i>LeNet5</i>	16 bits(6.4%)	4 bits(6.9%)	4 bits(6.4%)	4 bits(9.0%)
<i>CifarNet</i>	32 bits(18%)	8 bits(18%)	8 bits(18%)	8 bits(18.3%)

Table 7.5: Quantization summary (bit-width, error rate) on unpruned *LeNet5* and *CifarNet*. *DFP* is *Dynamic fixed-point*, *CFP* is *Customized floating-point* and (*Gysel*) is *Gysel et al.*’s implemenetation of *Dynamic fixed-point*.

Quantization on pruned networks works differently from the unpruned ones. As illustrated before, pruned network have central values that are non-zeros. Re-centralization significantly helps number representation systems to focus on the region with a large number of parameters using only one extra central bit. Two re-centralization methods, named *Centralized customized floating-point* (*CCFP*) and *Centralized dynamic fixed-point* (*CDFP*) respectively, are applied on the pruned networks. Table 7.6 shows the quantization results, both arithmetics achieve same quantization results on the pruned *LeNet5* model, but *Centered dynamic fixed-point* shows a better performance on the

pruned *CifarNet* model by achieving zero accuracy loss using only 6 bits.

Model	<i>CFP</i>	<i>CCFP</i>	<i>CDFP</i>
<i>LeNet5</i>	4 bits(6.4%)	4 bits(6.4%)	4 bits(6.4%)
<i>CifarNet</i>	8 bits(18.2%)	8 bits(18%)	6 bits(18%)

Table 7.6: Quantization summary (bit-width, error rate) on pruned *LeNet5* and *CifarNet*. *CFP* is *Customized floating-point* *CCFP* is *Centralized customized floating-point* and *CDFP* is *Centralized dynamic fixed-point*.

Chapter 8

Evaluation

8.1 Compression Pipeline

Han et al. proposed a three-stage compression pipeline called *Deep Compression*. This work is most comparable to my compression pipeline. Other research works almost focus only on whether pruning or quantization in an isolated manner. In this section, I would like to compare my compression pipeline to *Deep Compression*

As shown in Figure 8.1, *Han et al.* focused on reducing the number of parameters first using *Deterministic pruning* iteratively. This iterative pruning process is combined with retraining to bring back the lost accuracies. They then utilized both *Fixed-point quantization* and *Weights sharing* to minimize the number representation of each individual parameter. Weights are quantized to limited precisions and encoded using a code book to reduce the amount of representations. Finally, they tried to compress the number representations further using an encoding scheme [3].

Deep Compression is proven successful on many popular neural networks. To give a fair comparison between my compression pipeline and *Deep Compression*, the encoding stage is not considered for two reasons. First, encoding schemes normally encode symbols by their occurrences, such an encoding

scheme can be attached to the end of any compression pipelines. Second, the use of encoding indicates a need for encoder and decoder in hardware, this might affect energy efficiency of the hardware platform.

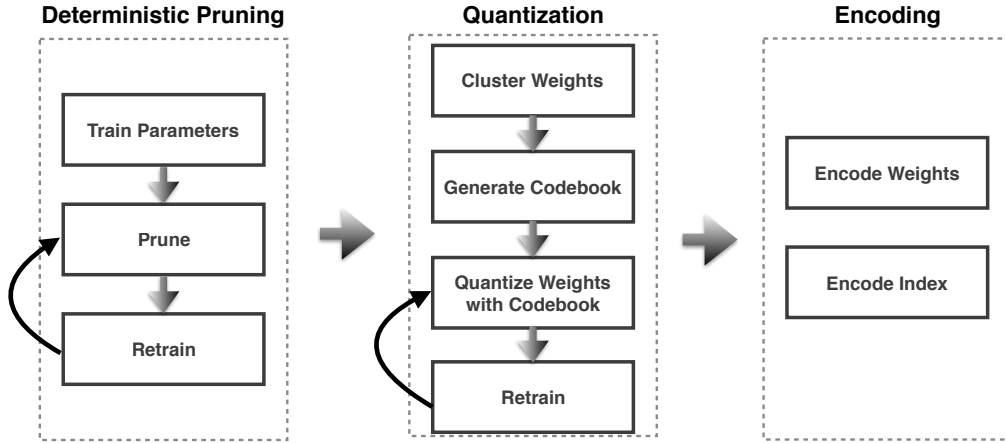


Figure 8.1: Overview of *Deep Compression* [3].

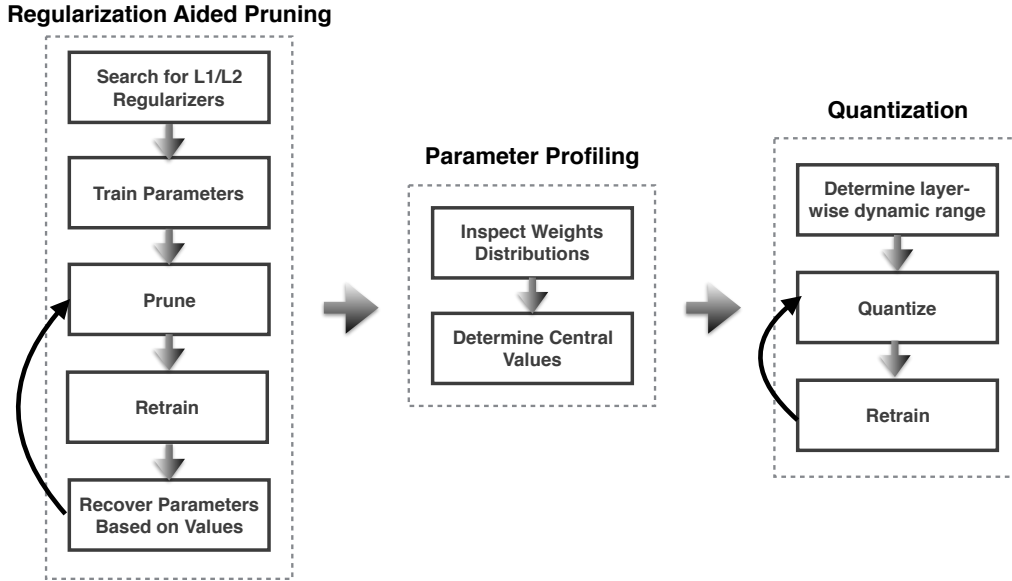


Figure 8.2: Overview of the proposed compression pipeline.

In my implementation, the focus stays on pruning and quantization, an overview of the complete compression pipeline is shown in Figure 8.2. *Regu-*

larization aided pruning is performed to reduce the number of parameters in a neural network. The best hyperparameter combinations for the regularizers come from an exhaustive search. The rest of the process is identical to *Dynamic network surgery*, where pruned weights have a chance to recover. The second phase is to profile the parameters, this profiling serves quantization since some important information is extracted from the pre-quantized model and the central values are determined. The quantization method used is *Centralized dynamic fixed-point quantization*. The dynamic ranges are stored in a layer-wise fashion, each quantization is combined with retraining to recover the test accuracies

8.2 Evaluation of Performance

To fully understand the performance difference between my proposed compression pipeline and *Deep Compression*, Table 8.1 summarizes the performance differences on *LeNet5*.

Model	Layer	Params	%(Han, P)	%(Han, P+Q)	%P	%P+Q
LeNet5	cov1	0.5K	66%	78.5%	78.9%	12.2%
	cov2	25K	12%	6.0%	7.5%	1.2%
	fc1	400K	8%	2.7%	0.3%	0.047%
	fc2	5K	19%	6.9%	67.1%	10.5%
	total	431K	8%(12x)	3.05%(33x)	1.6%(63x)	0.25%(403x)
	ER	-	0.8%	0.8%	0.64%	0.64%

Table 8.1: LeNet5 Pruning Summary, Pruning and Quantization, CR is the compression rate, ER is the error rate.

As shown in Table 8.1, putting a *LeNet5* into the proposed compression pipeline, the neural network obtains a compression rate of 403x. This compression rate shows an 12x increase compared to *Deep Compression*. In terms of pruning, *Regularization aided pruning* provides a better result, showing an increase of 51x compared to *Deterministic pruning*. Later on, *Deep Compression* achieved a 5-bit quantization using weights sharing and fixed-point

quantization. The proposed compression pipeline uses *Centralized dynamic fixed-point arithmetic* and quantized the entire model to 5 bits.

It is also important to note that the use of *Weights sharing* in *Deep Compression* gives significant energy overhead. Each weight now is encoded using a codebook. When neural network inference occurs, each weight fetches the correct value using the stored index as an address. This extra weights fetching process is the cost of using *Weights sharing*. *Yang et al.* suggested one of the major energy consumption in network inference is data movement [9]. Since *Weights sharing* added this redundant data movement, its energy consumption should be evaluated more carefully to evaluate the trade-off between energy saving and compression rates. In contrast, using *Centralized dynamic fixed-point arithmetic*, the arithmetic operators can be designed specifically for this kind of arithmetic and thus avoids the weights fetching stage.

8.3 Evaluation of Compression Techniques

In my compression pipeline, two main phases are pruning and quantization. The experimental results of various pruning techniques are summarized in Section 5.2.2. Non-deterministic pruning methods are generally better than deterministic pruning. The best performance pruning strategy is *Regularization Aided Pruning*, however, it might take a significantly larger amount of time depending on how experienced the designer is in finding the correct regularization parameters. In contrast, *Dynamic network surgery* also proves to have good compression rates and avoids the efforts of searching for the regularization hyperparameters.

In terms of quantization, the suitable methods for pruned and unpruned models are different. For unpruned models, *Customized floating-point* proves its performance is comparable to the popular *Dynamic fixed-point* arithmetic, moreover, it does not restrict any dynamic ranges in the number representation system. Consider pruned models, an idea of re-centralizing the arithmetic representations is suggested.

Chapter 9

Summary and Conclusions

9.1 Conclusion

In this project, I have investigated a number of pruning and quantization methods and built a complete compression pipeline based on the methods that have the best performances. A number of existing pruning/quantization have been explored in this project, and they are listed below:

1. Deterministic pruning
2. Dynamic network surgery
3. Fixed-point quantization
4. Dynamic fixed-point quantization

Inspired by these existing compression techniques, some novel pruning/quantization methods are also proposed:

1. Gradient profiling pruning
2. Regularization aided pruning
3. Customized floating-point quantization
4. Centralized quantization

The large range of exploration of compression techniques provided some important empirical results. For pruning, as summarized previously in Section 4.3 and Section 5.2.2, it can be concluded that pruning without biases and pruning in a non-deterministic manner give significantly better compression results. Regularizations also help improving the compression rates. The proposed pruning strategy, *Regularization aided pruning*, becomes the best performance pruning method. In terms of quantization, the results suggest *Customized floating-point quantization* achieves best quantization results on unpruned models due to its ability of tracking numbers to high precisions that are near zeros. However, re-centralization of arithmetics gives significant improvement on pruned model. The practical results proved that the preferred arithmetics are different for dense and sparse models because the weights distributions are different.

At the end, I selected *Regularization aided pruning* and *Centralized quantization* for building a complete compression pipeline. The proposed compression pipeline achieves a better compression results than existing compression techniques.

9.2 Future Works

Possible future works can be broken down into two parts: pruning and quantization. One possible future work is to focus on reducing the retraining time in pruning. Most pruning methods require at least one complete retraining after the pruning process [31, 40]. Shorten the retraining time can significantly reduce the computation time, but it might cause a bad compression result. *Gradient profiling pruning* shows promising results with limited retraining resources. It obtains information of weights importance during retraining and used this information to prune weights in a more selective manner. However, its performance should be further verified on other networks. *Regularization aided pruning* is proven useful in achieving high compression rates. The amount of time spent on choosing feasible hyperparameters for

regularizers are entirely empirical at this stage. One possible future work is to automate this process. It has been empirically observed that the feasible regularization loss (loss caused by regularizers) is around ten percents of the value of the cost function. If such a numerical ratio can be confirmed on a large range of networks, it is easy to automate the hyperparameter definitions by inspecting this numerical ratio. Another potential future work is to explore more efficient regularizers. Regularizers help pruning because they encourage sparsity during the retraining phase. It is therefore logical to consider using regularizers that are better at encouraging zeros for achieving a better compression. *Shakeout* is a new regularizer proposed by *Kang et al.* [41]. *Kang et al.* proved *Shakeout* encourages more sparsity than traditional l_1 and l_2 norms [41]. It is possible to use *Shakeout* as a regularizer for *Regularization Aided Pruning* in future developments.

With respect to quantization, I have demonstrated that *Customized floating-point* can be more efficient than the popular *Dynamic fixed-point* method when the neural network model is dense. It is noted that *Dynamic fixed-point* arithmetic fails to recover test accuracies if the trained network is sensitive to reduced dynamic ranges. A theoretical proof and more throughout experimentations should carry on in this direction to prove whether this is a general phenomenon. Some novel arithmetics, such as logarithmic arithmetics, are catching attentions recently [42] because of their low-cost hardware multiplications. Part of the future works is to propose energy consumption models for the proposed arithmetic operators in this project. These energy estimations could be helpful for exploring possible hardware accelerator architectures.

In this project, pruning and quantization are treated as individual stages and retraining takes place in each of these stages. One interesting future work is to combine pruning, quantization and retraining all in one stage. Previously, a number of research works have tried to combine pruning with retraining into one process and thus eliminates the need of pruning iteratively [18, 43]. The strategy is to add the weights masks into the cost function so that the network encourages more values to be pruned away. It is interesting to see whether there is a chance to fit quantization into this method, for instance,

the bit-width might be considered as a metric in the cost function and a lower bit-width is preferred.

Bibliography

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
- [2] Yiwen Guo, Anbang Yao, and Yurong Chen. Dynamic network surgery for efficient dnns. *CoRR*, abs/1608.04493, 2016.
- [3] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. *CoRR*, abs/1510.00149, 2015.
- [4] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [5] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [6] Adam Coates, Brody Huval, Tao Wang, David J. Wu, Bryan Catanzaro, and Andrew Y. Ng. Deep learning with COTS HPC systems. In *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*, pages 1337–1345, 2013.
- [7] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [8] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy P Lillicrap, Tim Harley, David Silver, and Koray

- Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning*, 2016.
- [9] Tien-Ju Yang, Yu-Hsin Chen, and Vivienne Sze. Designing energy-efficient convolutional neural networks using energy-aware pruning. *CoRR*, abs/1611.05128, 2016.
 - [10] Yunji Chen, Tao Luo, Shaoli Liu, Shijin Zhang, Liqiang He, Jia Wang, Ling Li, Tianshi Chen, Zhiwei Xu, Ninghui Sun, et al. Dadiannao: A machine-learning supercomputer. In *Proceedings of the 47th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 609–622. IEEE Computer Society, 2014.
 - [11] Tianshi Chen, Zidong Du, Ninghui Sun, Jia Wang, Chengyong Wu, Yunji Chen, and Olivier Temam. Diannao: A small-footprint high-throughput accelerator for ubiquitous machine-learning. In *ACM Sigplan Notices*, volume 49, pages 269–284. ACM, 2014.
 - [12] Song Han, Xingyu Liu, Huizi Mao, Jing Pu, Ardavan Pedram, Mark A Horowitz, and William J Dally. Eie: efficient inference engine on compressed deep neural network. In *Proceedings of the 43rd International Symposium on Computer Architecture*, pages 243–254. IEEE Press, 2016.
 - [13] Song Han, Junlong Kang, Huizi Mao, Yiming Hu, Xin Li, Yubin Li, Dongliang Xie, Hong Luo, Song Yao, Yu Wang, et al. Ese: Efficient speech recognition engine with compressed lstm on fpga. *arXiv preprint arXiv:1612.00694*, 2016.
 - [14] Chen Zhang, Peng Li, Guangyu Sun, Yijin Guan, Bingjun Xiao, and Jason Cong. Optimizing fpga-based accelerator design for deep convolutional neural networks. In *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pages 161–170. ACM, 2015.
 - [15] Eriko Nurvitadhi, Ganesh Venkatesh, Jaewoong Sim, Debbie Marr, Randy Huang, Jason Ong Gee Hock, Yeong Tat Liew, Krishnan Srivatsan, Duncan Moss, Suchit Subhaschandra, and Guy Boudoukh. Can fpgas beat gpus in accelerating next-generation deep neural networks? In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, FPGA ’17, pages 5–14, New York, NY, USA, 2017. ACM.
 - [16] Yu-Hsin Chen, Tushar Krishna, Joel S Emer, and Vivienne Sze. Eyerriss: An energy-efficient reconfigurable accelerator for deep convolu-

- tional neural networks. *IEEE Journal of Solid-State Circuits*, 52(1):127–138, 2017.
- [17] Dunja Mladenić, Janez Brank, Marko Grobelnik, and Natasa Milic-Frayling. Feature selection using linear classifier weights: interaction with classification models. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 234–241. ACM, 2004.
 - [18] Babak Hassibi, David G Stork, et al. Second order derivatives for network pruning: Optimal brain surgeon. *Advances in neural information processing systems*, pages 164–164, 1993.
 - [19] Suraj Srinivas and R Venkatesh Babu. Data-free parameter pruning for deep neural networks. *arXiv preprint arXiv:1507.06149*, 2015.
 - [20] Wenlin Chen, James Wilson, Stephen Tyree, Kilian Weinberger, and Yixin Chen. Compressing neural networks with the hashing trick. In *International Conference on Machine Learning*, pages 2285–2294, 2015.
 - [21] Bert Moons and Marian Verhelst. An energy-efficient precision-scalable convnet processor in a 40-nm cmos. *IEEE Journal of Solid-State Circuits*, 2016.
 - [22] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Quantized neural networks: Training neural networks with low precision weights and activations. *arXiv preprint arXiv:1609.07061*, 2016.
 - [23] Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks: Training deep neural networks with weights and activations constrained to ± 1 or ± 1 . *arXiv preprint arXiv:1602.02830*, 2016.
 - [24] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Gregory S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian J. Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Józefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Gordon Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul A. Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda B. Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *CoRR*, abs/1603.04467, 2016.

- [25] Yann LeCun, Corinna Cortes, and Christopher JC Burges. The mnist database of handwritten digits, 1998.
- [26] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. The cifar-10 dataset, 2014.
- [27] Yann LeCun et al. Lenet-5, convolutional neural networks. *URL: <http://yann.lecun.com/exdb/lenet>*, 2015.
- [28] Antonio Torralba, Rob Fergus, and William T Freeman. 80 million tiny images: A large data set for nonparametric object and scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 30(11):1958–1970, 2008.
- [29] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- [30] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [31] Georg Thimm and Emile Fiesler. Pruning of neural networks. Technical report, IDIAP, 1997.
- [32] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient transfer learning. *arXiv preprint arXiv:1611.06440*, 2016.
- [33] Feiping Nie, Heng Huang, Xiao Cai, and Chris H Ding. Efficient and robust feature selection via joint l2, 1-norms minimization. In *Advances in neural information processing systems*, pages 1813–1821, 2010.
- [34] Naveen Mellempudi, Abhisek Kundu, Dipankar Das, Dheevatsa Mudigere, and Bharat Kaul. Mixed low-precision deep learning inference using dynamic fixed point. *arXiv preprint arXiv:1701.08978*, 2017.
- [35] Tapas Kanungo, David M Mount, Nathan S Netanyahu, Christine D Piatko, Ruth Silverman, and Angela Y Wu. An efficient k-means clustering algorithm: Analysis and implementation. *IEEE transactions on pattern analysis and machine intelligence*, 24(7):881–892, 2002.
- [36] Miloš D Ercegovic and Tomas Lang. *Digital arithmetic*. Elsevier, 2004.
- [37] Darrell Williamson. Dynamically scaled fixed point arithmetic. In *Communications, Computers and Signal Processing, 1991., IEEE Pacific Rim Conference on*, pages 315–318. IEEE, 1991.

- [38] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Training deep neural networks with low precision multiplications. *arXiv preprint arXiv:1412.7024*, 2014.
- [39] Philipp Gysel. Ristretto: Hardware-oriented approximation of convolutional neural networks. *CoRR*, abs/1605.06402, 2016.
- [40] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016.
- [41] Guoliang Kang, Jun Li, and Dacheng Tao. Shakeout: a new regularized deep neural network training scheme. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, pages 1751–1757. AAAI Press, 2016.
- [42] Hokchhay Tann, Soheil Hashemi, Iris Bahar, and Sherief Reda. Hardware-software codesign of accurate, multiplier-free deep neural networks. *arXiv preprint arXiv:1705.04288*, 2017.
- [43] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient transfer learning. *CoRR*, abs/1611.06440, 2016.