

Technische Universität Berlin

Big Data Engineering (DAMS)

Fakultät IV

Ernst-Reuter-Platz 7

10587 Berlin

<https://www.tu.berlin/dams>



Thesis

[Choose
yours: Bach-
elor or Mas-
ter's]

Learned Quantization Schemes for Data-centric ML Pipelines

Anuun Chinbat

Matriculation Number: 0463111

20.01.2025

Supervised by

Prof. Dr. Matthias Boehm

M.Sc. Sebastian Baunsgaard

Hereby I declare that I wrote this thesis myself with the help of no more than the mentioned literature and auxiliary means.

Berlin, 01.01.2024

.....

(*Signature*)

[your name]

Abstract

Machine learning (ML) models are notoriously resource-intensive. Given their widespread application across every-day edge devices, the need to reduce their memory and computing requirements is becoming ever more pressing. Despite the said resource-intensiveness of ML models, at the same time they offer the main ingredient for the remedy to the malady - redundancy - which can be exploited to reduce their memory usage. While the redundancy exploitation of ML models is already a common technique that comes in different forms, starting from weight pruning [25] [27][13] and ending with knowledge distillation [14] [29], quantization presents itself as an especially promising area especially in the sense of learned quantization - the process of making ML models learn their optimal quantization parameters on their own. Hence, the current work employs two techniques that bypass the main issue of learned quantization, that is, the non-differentiability of rounding operations. While the first technique involves custom loss functions that directly take into account quantization goals, the second, more novel approach incorporates a custom scaling factor gradient calculation that takes into account the gradient of the parameters that are being quantized. As a result of these two techniques, a memory usage reduction of up to $\sim 2\times$ is obtained on MNIST, CIFAR10, and Imagenette.

Zusammenfassung

This is a placeholder for the german abstract (Kurzfassung) which should follow the same structure as the abstract. Just testing

Contents

1	Introduction	1
2	Background	3
2.1	Fundamentals of Deep Learning	3
2.1.1	Dense and Convolutional Layers	3
2.1.2	Loss Functions and Regularization	7
2.1.3	Forward-Pass and Back-Propagation	8
2.2	Basics of Quantization	8
2.2.1	Purpose and Definition	8
2.2.2	Core Quantization Approaches	9
2.3	Learned Quantization	12
2.3.1	Trade-offs and Challenges	12
2.3.2	Common Methods	13
3	Learned Quantization Schemes	15
3.1	Nested Quantization Layer	15
3.1.1	Concept and Design	15
3.1.2	Implementation Details	17
3.2	Custom Loss Functions	18
3.2.1	Penalty for Inverse Scale Factor Magnitude	18
3.2.2	Constraint on Bin Count for Quantization	18
3.2.3	Deviation between Quantized and Original Values	18
4	Experiments	19
4.1	Experimental Setup	19
4.2	Hyperparameter Tuning	19
4.2.1	Learning Rate and Regularization Coefficients	19
4.2.2	Quantization Threshold Coefficient	19
4.3	Results and Analysis	19
4.3.1	Overall Results	19
4.3.2	Dataset-Specific Insights	19
4.3.3	Further Implications	19
5	Related Work	21
6	Conclusions	23

List of Acronyms	25
Bibliography	27
Appendix	33

1 Introduction

Such is the life of a modern human being that not a single day passes without a machine learning model toiling away in the background. From unlocking one’s phone with Face ID in the morning to receiving a curated recommendation feed on Netflix in the evening — all is ML — but at what cost?

If we consider GPT-3 as an example, its 175 billion parameters need a whopping 700 gigabytes of storage in total — 4 bytes for each parameter represented in single-precision floating-point format (FP32). This costliness of modern ML models has revitalized interest in the research area of *quantization of Neural Networks* (NNs) which aims to reduce model size by developing methods that directly or indirectly decrease the amount of memory needed to store parameters numbering in the millions or billions. Going back to the GPT-3 example, by directly clamping its FP32 parameters to an 8-bit integer (INT8) representation, we can reduce its storage requirement from 700 to just 175 gigabytes.

But what costs does quantization itself entail? The natural answer to this question would be a reduction in model performance, as a decrease in precision logically implies a decrease in accuracy. However, there is a somewhat counterintuitive phenomenon where quantization improves accuracy by introducing noise, which act like a form of regularization, forcing the model to generalize better [4]. No matter whether quantization results in better or a slightly worse performance, the assumption is that for each model there is an optimal way to quantize it within reasonable degradation ranges. And if the model is able to learn its optimal parameters, it is most likely also able to learn its optimal quantization parameters.

The fact that, indeed, models can learn their optimal quantization parameters has been proven many times in the past. But is there a way to make them do it better - is a question that will always remain and to which this thesis aims to contribute. In that sense, the current work will explore novel ways to tackle the two main problems that the process of learning optimal quantization parameters poses. First, how to overcome the issue of non-differentiability of rounding operations in the back-propagation. Second, how to guide the model to quantize selectively where necessary and adaptively relax the quantization when a certain threshold is reached.

While there is a number of methods dealing with the issue of non-differentiability of rounding operations in quantization, the most popular of them are the Straight-Through Estimator and other similar approaches that employ continuous relaxations of discrete functions to enable gradient-based optimization. Besides these approximation methods,

1 Introduction

other more aggressive techniques — in a sense — avoid non-differentiability altogether and perform quantization based on strictly defined constraints. Binary Connect [4], for example, binarizes weights during both forward pass and backpropagation, with the real values of weights used only during parameter update. However, these aggressive methods, while effective in simpler scenarios, often struggle on more complex datasets. The more effective methods seem to combine both constraints and gradient approximations — just like XNOR-Net enhances BinaryConnect with a gradient computation formula designed for binary weights [33].

Despite the abundance of different methods, there is still room for improvement with regard to the second problem mentioned earlier, namely, not only how to quantize, but also where to do so and when to stop. Therefore, in this thesis, we try to fill this room with the following contributions:

- We introduce a method that directly considers the gradient-to-parameter ratio, which conveys how much the parameter is being adjusted relative to its current value. Based on this ratio, the model learns its quantizer scaling factors, applied at different granularities, alongside the standard trainable parameters.
- We also provide a modular framework that can be easily integrated into a wide range of applications and layers with minimal adjustments, ensuring flexibility and usability.
- We also provide applicable ranges for this ratio for effective quantization for dense and convolutional layers.

2 Background

This chapter addresses the theoretical and contextual background necessary to understand the key concepts and methodologies that form the foundation of the current research. The first section will discuss the basics of deep NNs, upon which the technical setup of this thesis is based. The next section aims to provide a broader context for the term quantization, followed by a final section that explains common techniques of *learned quantization*, as well as the trade-offs and challenges they present.

2.1 Fundamentals of Deep Learning

This section introduces the fundamental concepts of deep learning, beginning with the most basic NN architecture components 2.1.1 and progressing to loss functions with regularization 2.1.2. The concepts of the forward pass and backpropagation will be explained in the last subsection 2.1.3.

2.1.1 Dense and Convolutional Layers

NNs can be considered a mathematical abstraction of the human decision-making process. Consider a scenario where, given an image, you need to say aloud what you see. The two eyes can be regarded as input nodes that receive the initial data, the brain can be seen as a set of *hidden layers* that process this data, and your mouth — the output node that provides the final answer.

A hidden layer, which typically consists of many neurons, is where the magic — or the transformation of data — happens. In its simplest form, within the classic *Multilayer Perceptron* (MLP) model, each hidden layer neuron performs a weighted operation:

$$output = f(w \cdot input + b)$$

where:

- *input* refers to the outputs from the previous layer (or the initial data from input nodes) that are fed into a specific neuron in the hidden layer.
- *w* (weights) is a vector of parameters associated with that specific neuron, defining the importance of each input received by this neuron.
- *b* (bias) is an additional scalar parameter specific to the neuron, which shifts the result of the weighted sum, allowing for more flexibility.

2 Background

- f is the *activation function*, a nonlinear function applied to the weighted sum of inputs and bias in that specific neuron, allowing for more complexity.
- *output* is the result produced by the neuron, which will then be passed on to the next hidden layer (or to the final output layer).

Hidden layers where each neuron is connected to every neuron in the previous layer and every neuron in the next layer are called *dense layers*.

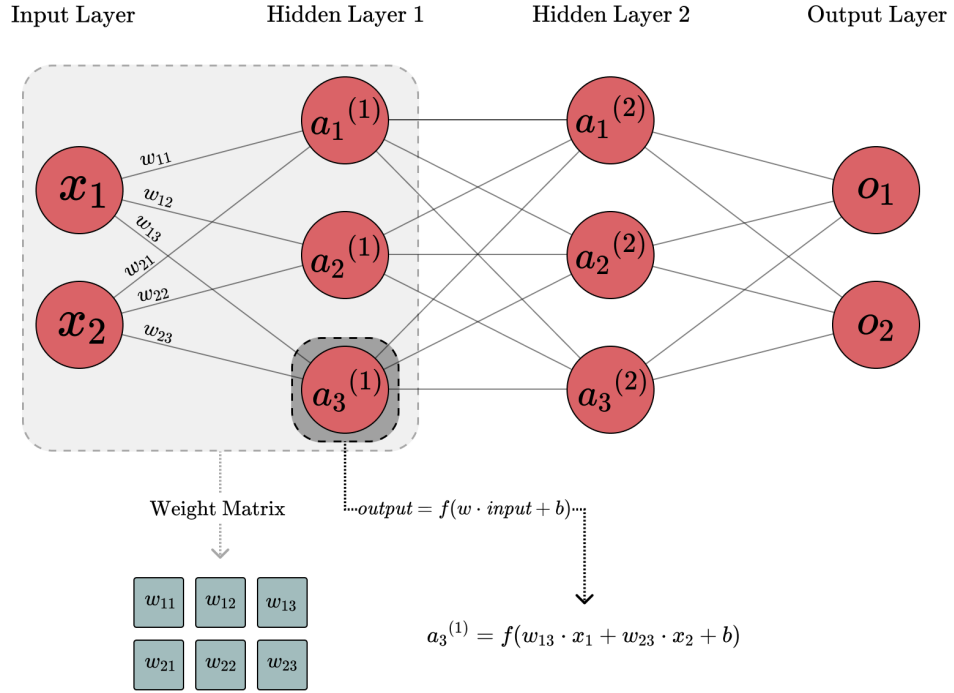


Figure 2.1: An example of a NN with two hidden dense layers, showing the connections between neurons in adjacent layers.

Mathematically, dense layers can be represented as:

$$a = f(W \cdot x + b)$$

where:

- x is the input vector, representing outputs from the previous layer (or initial input data for the first layer).
- W is the *weight matrix*, with each row corresponding to the weight vector w of a specific neuron.

- b represents the bias vector, where each scalar element corresponds to the bias of a specific neuron.
- f is the *activation function* that is applied element-wise.
- a refers to the output vector, representing the activations of all neurons.

This interconnectedness of dense layers introduces the inherent redundancy, or the over-parameterizedness of NNs [9]. It is particularly true in models with a large number of neurons, where W results in a vast number of parameters, which do not contribute to the model accuracy equally [17].

Convolutional layers are another type of hidden layers that involve a *convolution* operation on the input. Intuitively, a standard convolution is a process of sliding a small grid over an input to find patterns. The figure below, for example, shows the application of the Sobel kernel that detects edges on the input image.

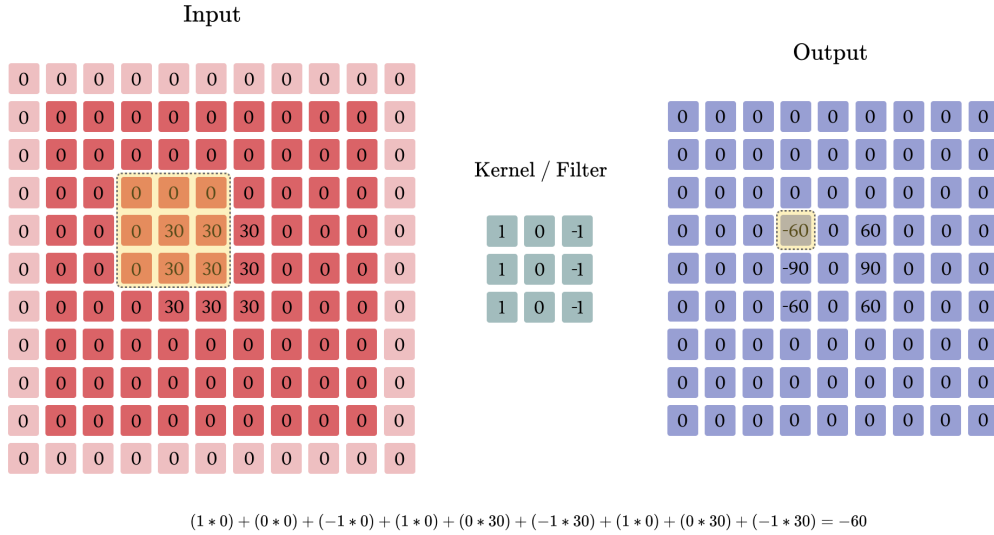


Figure 2.2: A 3×3 kernel (filter) sliding over a padded input matrix to compute the output feature map, demonstrating the interaction between the kernel weights and input values at a specific position.

For multi-channelled inputs, like RGB images, the convolution operation uses a multi-channelled kernel, as shown in Figure 2.3, producing a single-channelled feature map that combines weighted contributions from all input channels. A convolutional layer typically includes multiple such kernels, generating feature maps equal to the number of kernels. After the convolution operation generates the feature maps, a bias term is added to each map, and the activation function is applied element-wise — just like in dense layers.

2 Background

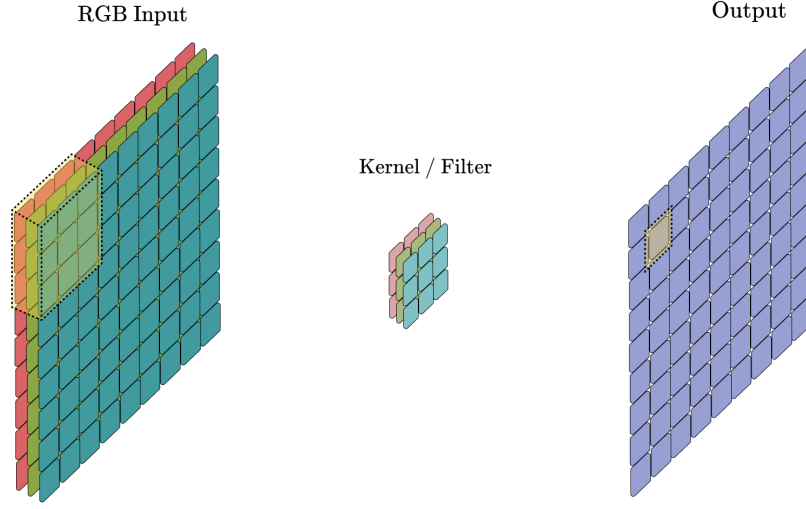


Figure 2.3: A $3 \times 3 \times 3$ kernel (filter) sliding over an RGB input matrix to produce a single-channelled output feature map.

Mathematically, a convolutional layer can be represented as:

$$y_{i,j,k} = \phi \left(\sum_{m=1}^M \sum_{p=1}^P \sum_{q=1}^Q x_{i+p-1,j+q-1,m} \cdot w_{p,q,m,k} + b_k \right)$$

where:

- P, Q are the height and width of the filter, respectively.
- M is the number of input channels.
- $y_{i,j,k}$ denotes the output at position (i, j) for the k -th filter.
- $x_{i+p-1,j+q-1,m}$ is the input at position $(i + p - 1, j + q - 1)$ for the m -th input channel.
- $w_{p,q,m,k}$ represents the weight of the filter at position (p, q) for the m -th input channel and k -th filter.
- b_k is the bias for the k -th filter.
- $\phi(\cdot)$ is the activation function.

In simpler terms, a convolutional layer applies filter weights as it slides over rows (p), columns (q), and channels (m), sums the results, adds bias (b_k), and repeats this for all positions (i, j) and filters (k).

Although convolutional layers often have fewer weight parameters than dense layers in typical architectures, they still contain redundancies [16], presenting an opportunity for quantization. Thus, both dense and convolutional layers will be the focus of this work.

2.1.2 Loss Functions and Regularization

The weights and biases are usually *learnable parameters* that the model adjusts during *training*. The training process of NNs is similar to how our brains learn from mistakes. Given the ground truth, a NN adjusts its learnable parameters using a specific function that compares the ground truth with the output generated by the network, essentially measuring the magnitude of the network’s errors.

This function is called a *loss function*, and depending on the type of question the network aims to answer, it can take many different forms. For example, for the MLP described in Figure 3.1 that generates a binary classification, we would use the *log loss* function. Since the datasets used in this thesis involve multi-class classification, the *sparse categorical cross-entropy* (SCCE) loss function will be used, which measures the difference between the predicted class probabilities and the true labels for each class in the dataset.

Often the loss function alone is not enough for a NN to perform well, as it may lead to overfitting or fail to capture desired generalization properties. This is why a *regularization term* that penalizes unwanted behaviours is added to the loss function.

A typical regularization term is L_2 , which penalizes large weights by adding the sum of the squared weights to the loss. The modified loss function is then expressed as:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{data}} + \lambda \sum_i w_i^2$$

where:

- $\mathcal{L}_{\text{data}}$ is the original loss function (in our case, the SCCE loss function).
- λ is a scalar parameter that controls the strength of the regularization.
- w_i represents each individual weight value in the model.

The current work employs multiple custom regularization terms that encourage specific behaviors in the models while discouraging others. These terms will be discussed in detail in the Experimental Setup section 4.1.

2 Background

2.1.3 Forward-Pass and Back-Propagation

The repetition of the mathematical operations described earlier in the Dense and Convolutional Layers subsection 2.1.1 during model training constitutes the *forward pass*. It is the process where input data is passed through the network layer by layer, with each layer applying its learned weights and biases to produce a final output.

As mentioned in the previous subsection, this output is then compared with the ground truth by the loss function that produces an error. This error is used to update the learnable parameters in W and b during a process called *back-propagation*.

In other words, back-propagation is the method by which the network adjusts its parameters to minimize the error. It calculates the gradient of the loss function with respect to each parameter using the chain rule. W and b are typically updated as follows:

$$W = W - \eta \frac{\partial L}{\partial W}, \quad b = b - \eta \frac{\partial L}{\partial b}$$

where L is the loss function, and η is the learning rate.

For example, consider the weight $w_{1,1}$ represented as the line between x_1 and the hidden layer node $a_1^{(1)}$ in Figure 3.1. The gradient of this weight with respect to the loss is calculated using the chain rule:

$$\frac{\partial L}{\partial w_{1,1}} = \frac{\partial L}{\partial o_1} \cdot \frac{\partial o_1}{\partial a_1^{(1)}} \cdot \frac{\partial a_1^{(1)}}{\partial w_{1,1}}$$

Where:

- $\frac{\partial L}{\partial o_1}$ is the gradient of the loss with respect to o_1 .
- $\frac{\partial o_1}{\partial a_1^{(1)}}$ is the gradient of o_1 with respect to the output of $a_1^{(1)}$.
- $\frac{\partial a_1^{(1)}}{\partial w_{1,1}}$ is the value of x_1 , since $a_1^{(1)}$ is a weighted sum of the inputs.

This shows how each weight contributes to the final error during back-propagation.

2.2 Basics of Quantization

This section aims to answer the *why* question with respect to quantization and further provides a broader understanding of the term regarding its types.

2.2.1 Purpose and Definition

As we become increasingly dependent on deep learning models disguised as everyday tools, the need for these models to function in a resource- and time-efficient manner is

more imperative than ever. The focus on resource efficiency is particularly important, with the research community expressing concerns regarding the environmental effects of large models, the exponential size growth of which continues to significantly outpace that of system hardware [38]. In this regard, studies have examined quantization within the context of Green AI as a method to reduce the carbon footprint of ML models [34].

Aside from the environmental considerations, the mere need to reduce the computational cost and speed of predictive models comes as an apparent business requirement. This requirement is essential when — quite ironically — embedded systems, famous for their compactness, meet ML models, infamous for their complexity. Microcontrollers, for instance, usually are not able to perform floating-point operations, which must therefore be emulated in software, introducing significant overhead. This is why quantization, the process which reduces the memory footprint of a model, is also extensively covered in the realm of embedded systems that inherently prefer integer arithmetic, as well as bitwise operations [33] [40] [1][28].

Another motivation for quantization — although somewhat controversial — is the fact that reducing the bit-width of ML models makes them robust to adversarial attacks in certain cases [26]. This holds significant value in fields, such as autonomous driving, where model vulnerability may result in fatal outcomes. Interestingly enough, the use cases where such robustness is required also demand fast inference, as they rely on real-time predictions. Consider healthcare diagnostics needed for emergency scenarios or military defense mechanisms designed for immediate action.

The list of reasons why quantization is useful may go on for a while, but regardless of the motivation, the essence of the term itself — rooted in the early 20th century — remains unchanged: quantization refers to the division of a quantity into a discrete number of small parts [11]. With regard to ML models, it describes the process of dividing higher bit-width numbers into a discrete number of lower bit-width representations without causing significant degradation in performance [9].

Since ML models are generally considered redundant or over-parameterized, there are multiple points where quantization can be applied. Specifically, in this thesis, we apply quantization to the weights and biases of dense layers, as well as the kernels and biases of convolutional layers. Other applications include, but are not limited to, layer activation and layer input quantization (two sides of the same coin), as well as gradient quantization. The bottom line is that wherever there is an opportunity for arithmetic or memory optimization, there is room for quantization.

2.2.2 Core Quantization Approaches

There is a multitude of ways to classify NN quantization methods, a broader overview of which will be covered in the Related Work chapter 5. For now, we will focus on a few basic approaches from the general categories of both *data-driven* and *data-free* methods

2 Background

[39] to provide a basic understanding of the NN quantization process.

The simplest form of data-free quantization, or *post-training quantization* [19], involves converting already trained parameters from FP32 to a lower bit-width format without using the initial training data. A common approach is to apply *uniform quantization* that maps real values to a set number of *bins*. The general formula can be written as:

$$Q(r) = \text{round}\left(\frac{r}{S}\right)$$

Where:

- $Q(\cdot)$ denotes the quantization operation.
- r is the real value of a given model parameter in higher bit-width representation.
- $\text{round}(\cdot)$ is some rounding operation, such as a simple $\text{floor}(\cdot)$.
- S is a scaling factor.

As a result, we essentially end up with a discrete number of values in lower bit precision, instead of an almost continuous range of real numbers as shown in Figure 2.4.

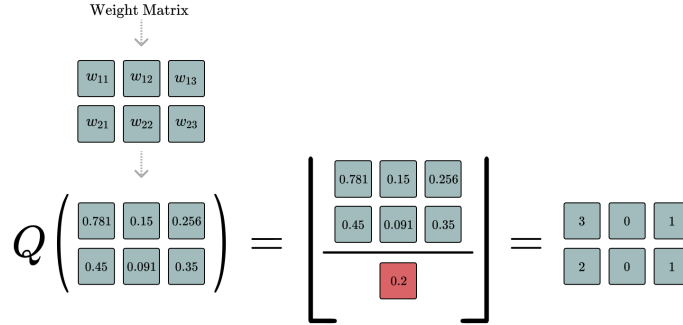


Figure 2.4: An example illustrating the quantization operation on the weight matrix from Figure 3.1, with arbitrary values for demonstration purposes.

Unlike data-free quantization — as the name suggests — data-driven quantization typically involves retraining the model using the initial data. An example of this approach is the Ristretto framework [12], which, similar to data-free methods, first analyzes the trained model to select suitable lower bit-width number formats for its weights. Then, using a portion of the original dataset, the framework determines appropriate formats for layer inputs and outputs. As a next step, based on the validation data, Ristretto adjusts the quantization settings to achieve optimal performance under the given constraints. Finally, the quantized model is fine-tuned using the training data.

A much simpler example of data-driven quantization could be the min-max quantization on input data as shown in Figure 2.5. This method can also be used in a data-free scenario to quantize learned model parameters and is internally used as one of the default techniques in popular ML frameworks like Tensorflow and PyTorch.

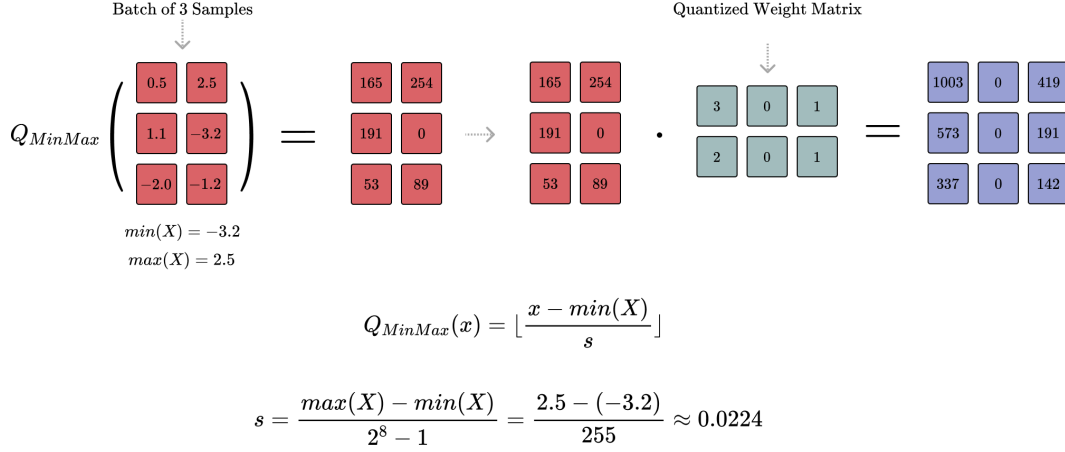


Figure 2.5: An example illustrating min-max quantization of input data to 8 bits, followed by matrix multiplication with the quantized weight matrix from Figure 2.4. Input data has arbitrary values for demonstration purposes.

In the previous subsection, we discussed *where* quantization could be applied in a model, mentioning weights, kernels, and biases as the focus of this thesis. Figures 2.4 and 2.5 show examples of quantization using a scalar scaling factor. However, scaling factors could be applied at varying levels of detail, and this is where the concept of *quantization granularity* comes into play.

Granularity refers to the level of detail at which scaling factors are applied, ranging from a single factor for an entire kernel (coarse granularity) to separate factors for individual spatial locations, channels, or filters (fine granularity). For instance, Figure 2.6 illustrates various possible granularities for the kernels of convolutional layers. Despite this wide range of possibilities, channel-wise quantization is currently the standard for convolutional layers [9], as it helps parallel processing capabilities of accelerators that compute channel outputs independently. For dense layers, row-wise quantization (one scaling factor for weights used by a single output neuron) is more prevalent because it aligns with matrix-vector multiplication, which then can be carried out by specialized linear algebra libraries in an optimized way [22].

2 Background

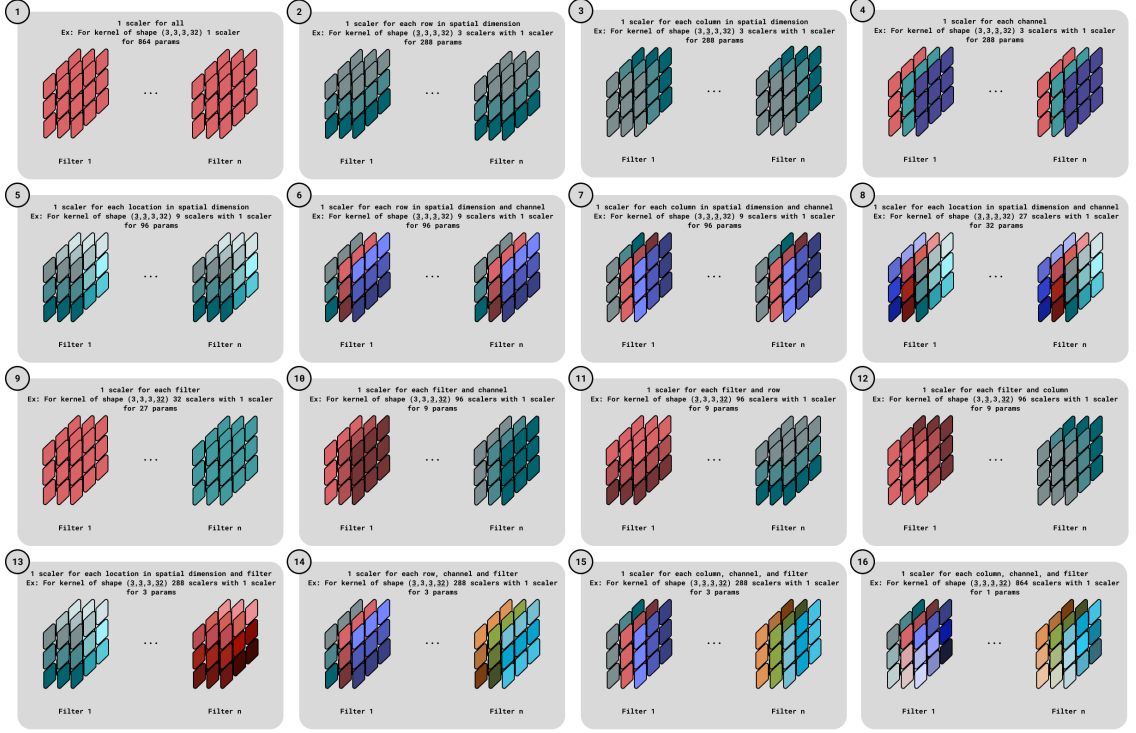


Figure 2.6: A demonstration of the varying application of scaling factors, ranging from a single scalar applied to the entire kernel (1) to separate scalars assigned to spatial dimensions (e.g., 1, 2), channels (4), filters (9), and other granular configurations.

2.3 Learned Quantization

Now that the fundamentals of quantization have been covered, this section introduces key concepts commonly encountered in learned quantization, including its challenges, trade-offs, and the popular techniques used to overcome them.

2.3.1 Trade-offs and Challenges

The inherent — or rather the generally accepted — characteristic of quantization is that it negatively influences performance. This is referred to as the trade-off between quantization and generalization, which reflects the question of how much accuracy — or whatever performance metric is being used — we are willing to sacrifice to gain a reduction in computational cost, memory usage, or inference time. However, the truth is that we usually cannot afford sacrificing anything. This, coupled with the lack of a guarantee that pre-defined *quantizers* can yield optimal results [40] [7], has paved the way for the burgeoning field of learned quantization, which aims to *learn* how to quantize the model

in a manner that mitigates performance loss.

Learned quantization is, however, a double-edged sword in the sense that, despite producing compact results, the cost to achieve them is higher [31]. The obvious reason is the additional computational overhead introduced by learnable quantizers. Thus, it is important to strike a balance between learning optimal quantization and keeping the training process manageable — which explains the prevailing emphasis on simplicity in most learned quantization research.

The main issue in achieving this simplicity is posed by the fact that discretization, in its essence, is non-differentiable — meaning it is challenging to integrate any kind of discretizing operations into gradient-based optimization methods, upon which ML models rely. Using the chain rule back-propagation example from subsection 2.1.3, let's consider a simple flooring operation introduced into the process to better understand the problem.

Suppose we want to quantize activations and apply $\text{floor}(\cdot)$ to hidden layer outputs:

$$a_{1,q}^{(1)} = \text{floor}(a_1^{(1)})$$

As a result, the chain rule becomes:

$$\frac{\partial L}{\partial w_{1,1}} = \frac{\partial L}{\partial a_{1,q}^{(1)}} \cdot \frac{\partial a_{1,q}^{(1)}}{\partial a_1^{(1)}} \cdot \frac{\partial a_1^{(1)}}{\partial w_{1,1}}$$

where $\frac{\partial a_{1,q}^{(1)}}{\partial a_1^{(1)}}$ presents a challenge. Since $a_{1,q}^{(1)} = \text{floor}(a_1^{(1)})$, the derivative is:

$$\frac{\partial a_{1,q}^{(1)}}{\partial a_1^{(1)}} = \begin{cases} 0 & \text{if } a_1^{(1)} \notin \mathbb{Z}, \\ \text{undefined} & \text{if } a_1^{(1)} \in \mathbb{Z}. \end{cases}$$

This means that for most values of $a_1^{(1)}$, which are non-integer, the gradient becomes 0, resulting in $w_{1,1}$ not receiving any updates. For integer values of $a_1^{(1)}$, the backpropagation process fails altogether.

In essence, circumventing the issue of non-differentiability is the fundamental problem that learned quantization aims to solve, all while managing the aforementioned trade-offs to produce a model that is compact, not overly complex to train, and highly performant.

2.3.2 Common Methods

The most common method to address the issue of non-differentiability is to approximate the gradient of quantization operators using the Straight-Through Estimator (STE) [2] [8]. This workaround applies the quantization operation as is during the forward-pass, but replaces the gradient of the piece-wise discontinuous function with that of a continuous

2 Background

identity function. Returning to the example from the previous subsection, if we apply the STE to the problematic gradient, instead of:

$$\frac{\partial a_{1,q}^{(1)}}{\partial a_1^{(1)}} = \begin{cases} 0 & \text{if } a_1^{(1)} \notin \mathbb{Z}, \\ \text{undefined} & \text{if } a_1^{(1)} \in \mathbb{Z}. \end{cases}$$

we approximate it as:

$$\frac{\partial a_{1,q}^{(1)}}{\partial a_1^{(1)}} \approx 1$$

This enables gradient flow, allowing model parameters to receive updates without being hindered by the non-differentiability of the quantization step.

The STE — and other estimators [3] — are the cornerstone of quantization-aware training (QAT) [18], a subfield that falls under the broader umbrella of learned quantization. QAT, however, does not necessarily focus on learning quantization parameters. Instead, it focuses on helping the model adapt to the loss caused by the quantization process during the forward pass, whether or not trainable quantizers are involved.

More often than not QAT and trainable quantization parameters are combined. An example is quantization-interval-learning (QIL) [20], which uses three trainable quantization parameters (the center of the interval, the distance to the center, and a parameter that controls the scaling itself) with piecewise differentiable operations and relies on STE for gradient updates of quantized weights and activations. Similarly, the learned step size quantization (LSQ) approach [7] defines a single trainable quantization parameter (step size) with an explicit gradient formula and also uses STE for standard parameter updates. LQ-Nets [40] depend on STE too, but — unlike the two previous techniques — directly incorporate a quantization error minimization algorithm to calculate binary encodings and the associated bases of model parameters given a predefined bit size.

INSERT STUFF ABOUT QUANTIZATION ENCOURAGING REGULARIZATION

In this thesis, we will utilize the STE but introduce a custom gradient calculation for the scale factor gradients based on a specific type of gradient sensitivity in the model parameters. Additionally, we will explore custom loss regularization terms that encourage quantization and systematically compare them across different datasets.

3 Learned Quantization Schemes

This chapter introduces two custom learned quantization schemes — approaches that allow models to learn to quantize themselves with adjustable aggressiveness. The first one, a custom quantization layer featuring tailored logic and a threshold for scale updates, will be discussed in the first section. The second scheme, which focuses on custom regularization terms with a configurable penalty rate, will be covered second.

3.1 Nested Quantization Layer

To separate the quantization logic from the usual structure of NN layers, we define a nested quantization layer that can be used within a standard layer. This approach provides usability, making it easy to extend the logic to other types of layers beyond dense and convolutional ones. The implementation details will follow after we first explain the core logic it incorporates.

3.1.1 Concept and Design

Our nested quantization layer has only one trainable parameter, *scale* (s), which serves as the scaling factor for model parameter quantization. The quantization itself is performed using a simple flooring operation:

$$P_{quantized} = \text{floor}\left(\frac{P}{s}\right)$$

where P denotes the parameter being quantized, such as a weight, bias, or kernel. The scaling factor s has an adjustable shape, allowing it to be applied at different levels of granularity, such as per-row, per-column, per-channel, or even per-element.

During back-propagation, the scaling factor s is updated using a custom gradient formula. The gradient of the loss with respect to s , denoted as $\nabla_s L$, is computed as:

$$\nabla_s L = g_s \cdot m,$$

Let's consider both multiplication terms separately. g_s is the main "decision maker" on whether to increase the scale and, therefore, quantize more. It is based on a hyperparameter threshold λ , which is compared against the ratio r .

$$g_s = \begin{cases} 0, & \text{if } r \geq \lambda, \\ -\tanh(\lambda - r), & \text{if } r < \lambda, \end{cases}$$

3 Learned Quantization Schemes

In turn, r is the ratio between the gradient of the model parameter with respect to the loss and its absolute value:

$$r = \frac{|\nabla_P L|}{\max(\epsilon, |P|)}$$

In essence, it conveys the relative impact of the gradient on the parameter's value. A large ratio indicates that the parameter is not ready for aggressive quantization because small perturbations can lead to significant changes in its optimization. Conversely, parameters with a small r are better candidates for quantization since they are less sensitive.

The decision to replace P with ϵ when $P = 0$ ensures that the corresponding r becomes large, effectively resisting quantization. This behavior is valid regardless of the parameter's sensitivity — if the zero parameter is sensitive, quantization could disrupt future optimization steps, and if it is not sensitive, quantization adds no value since S has a positive non-zero constraint and $\frac{P}{S}$ will remain zero.

The motivation behind using $\tanh(x)$ primarily stems from two key reasons. First, it is bounded (in our use case, to $[-1, 0)$), which prevents excessive gradient magnitudes. Second, unlike sigmoid, it does not require any additional rescaling since it is symmetric around 0. An additional point is that $\tanh(x)$ saturates comparably faster, potentially allowing for more decisive gradients, but it is uncertain how much influence this has.

Now that g_s is covered, let's take a look at m defined as:

$$m = \max(|P_{\text{quantized}}|)$$

where the shape of m corresponds to the shape of s . For example, if s is a row-wise scaler, then m will hold the maximum value from each row of $|P_{\text{quantized}}|$. Similarly if s is a scalar scaler, then m represents the maximum across the entire layer parameter.

A larger m indicates a wider range of quantized values, implying the parameter can tolerate coarser quantization. In contrast, a smaller m means a narrower range, where aggressive quantization could be rather harmful. As a result, by multiplying g_s with m , the adjustment to the scale becomes proportional to the parameter's range. This encourages more aggressive quantization for parameters with larger ranges while being more conservative for smaller ones.

To sum this part up, the intuition is that the gradient adjustment for the scale factor s adapts dynamically based on both the sensitivity of the parameter (g_s) and its range (m). Sensitive parameters are left with a "zero vote," while the less sensitive ones determine how much quantization they can tolerate.

The final touch is that the scale gradients are initially calculated for each parameter value but are then aggregated along the corresponding granularity axes. This reflects the

collective behavior of parameters within the same granularity, where only those deemed quantizable and with a meaningful "say" contribute to the overall adjustment, while sensitive parameters express their resistance with a "zero vote."

3.1.2 Implementation Details

As the name "nested quantization layer" suggests — this layer is implemented in a way that it is initialized from within a model layer itself. Conceptually, the resulting layer with one or more nested quantization layers is illustrated in the figure below.

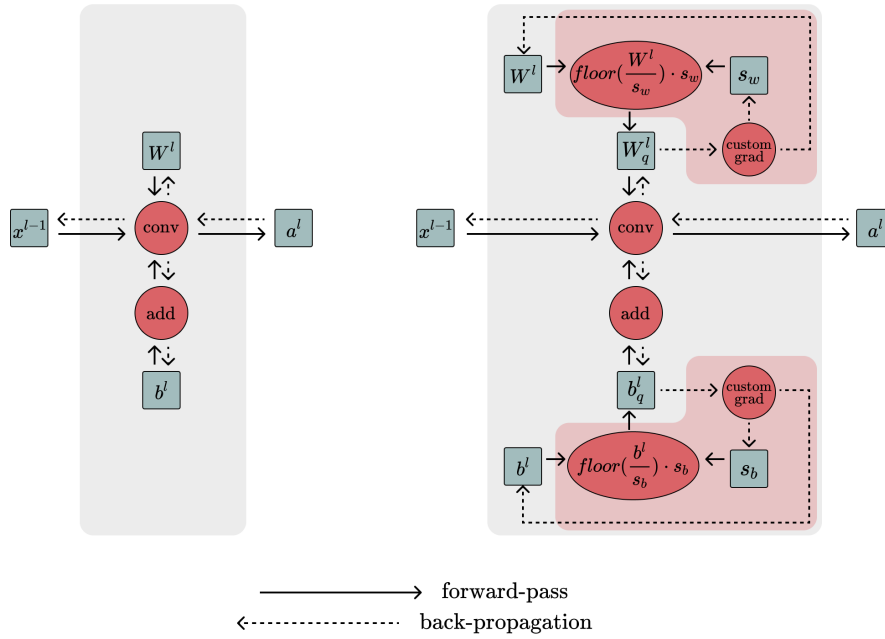


Figure 3.1: A standard convolutional layer (left) and its integration with the nested quantization layer (right) for both weights and bias. Quantization logic is applied to weights and biases during the forward pass, with trainable scaling factors updated using custom gradients in the backward pass.

In this section, we introduce custom layers built upon Tensorflow's `tf.keras.Layer` class, which serves as the base for all Keras layers. Each custom layer also leverages Tensorflow's `tf.custom_gradient` decorator to define its own gradient computation. For clarity, the upcoming subsections start by showing how to define the corresponding standard, non-quantized layer using `tf.keras.Layer` and `tf.custom_gradient`, then move on to the specific quantized implementations.

Yang You, Igor Gitman, and Boris Ginsburg - Large batch training of convolutional networks This might have info on the ratio thingie

3.2 Custom Loss Functions

3.2.1 Penalty for Inverse Scale Factor Magnitude

3.2.2 Constraint on Bin Count for Quantization

3.2.3 Deviation between Quantized and Original Values

4 Experiments

This chapter provides details about the experiments conducted within the context of this thesis.

You'll need these: there's a formula for precision requirement in [21]

4.1 Experimental Setup

All experiments are carried out on machine XYZ.

4.2 Hyperparameter Tuning

4.2.1 Learning Rate and Regularization Coefficients

4.2.2 Quantization Threshold Coefficient

4.3 Results and Analysis

4.3.1 Overall Results

4.3.2 Dataset-Specific Insights

4.3.3 Further Implications

4 *Experiments*

5 Related Work

A significant amount of scientific work has been done on QAT. This research can be categorized based on different characteristics, which are covered separately in the following paragraphs.

Model architecture.

RNN - [30]

CNN - [33]

CNN - [4] DNN - [10] Transformer bases - [23]

Quantization target parameters.

weights and activations - [24]

weights and activations - [17]

weights - [32]

weights - [30]

binary weights and input activations - [33]

gradients - [41] layer inputs and weights - [39] weights and activations - [40] weights

activations and gradients - [41] **Granularity of quantization.**

Handling of differentiability. STE - [40]

Quantization precision.

binary weights and input activations - [4]

binary weights and activations - [17]

binary weights and input activations - [33]

ternary weights - [30]

higher precision for more important parameters and lower precision for less important ones - [21]

mixed precision - [36] mixed precision - [5]

Integration with pruning & other techniques.

pruning and Huffman Coding - [13]

distillation - [32]

higher precision for more important parameters and lower precision for less important ones or pruning - [21] knowledge distillation [37]

Modifications to loss functions. regularization term WaveQ - [6]

proposes a regularization term into the loss function to push the weight values towards +1 and -1 [35]

introduces a novel loss formulation where each quantization has different importance - [15] **Other interesting approaches.** k-means for parameters [10]

I BERT article has a good related work overview

5 *Related Work*

6 Conclusions

This chapter summarizes the contributions of the thesis and provides an outlook into future work.

Problems:

List of Acronyms

ML	Machine Learning
----	------------------

List of Acronyms

Bibliography

- [1] Dorra Ben Khalifa and Matthieu Martel. “Rigorous Floating-Point to Fixed-Point Quantization of Deep Neural Networks on STM32 Micro-controllers”. In: *10th International Conference on Control, Decision and Information Technologies, CoDIT 2024, Vallette, Malta, July 1-4, 2024*. IEEE, 2024, pp. 1201–1206. DOI: 10.1109/CODIT62066.2024.10708400. URL: <https://doi.org/10.1109/CoDIT62066.2024.10708400>.
- [2] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. “Estimating or Propagating Gradients Through Stochastic Neurons for Conditional Computation”. In: *arXiv preprint arXiv:1308.3432* (2013).
- [3] Jun Chen et al. “Propagating Asymptotic-Estimated Gradients for Low Bitwidth Quantized Neural Networks”. In: *IEEE J. Sel. Top. Signal Process.* 14.4 (2020), pp. 848–859. DOI: 10.1109/JSTSP.2020.2966327. URL: <https://doi.org/10.1109/JSTSP.2020.2966327>.
- [4] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. “BinaryConnect: Training Deep Neural Networks with binary weights during propagations”. In: *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*. Ed. by Corinna Cortes et al. 2015, pp. 3123–3131. URL: <https://proceedings.neurips.cc/paper/2015/hash/3e15cc11f979ed25912dff5b0669f2cd-Abstract.html>.
- [5] Yinpeng Dong et al. “Stochastic Quantization for Learning Accurate Low-Bit Deep Neural Networks”. In: *Int. J. Comput. Vis.* 127.11-12 (2019), pp. 1629–1642. DOI: 10.1007/S11263-019-01168-2. URL: <https://doi.org/10.1007/s11263-019-01168-2>.
- [6] Ahmed T. Elthakeb et al. “Gradient-Based Deep Quantization of Neural Networks through Sinusoidal Adaptive Regularization”. In: *CoRR* abs/2003.00146 (2020). arXiv: 2003.00146. URL: <https://arxiv.org/abs/2003.00146>.
- [7] Steven K. Esser et al. “Learned Step Size quantization”. In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL: <https://openreview.net/forum?id=rkg066VKDS>.
- [8] Angela Fan et al. “Training with Quantization Noise for Extreme Model Compression”. In: *Proceedings of the International Conference on Learning Representations (ICLR)*. Facebook AI Research, LORIA, Inria. 2021.

Bibliography

- [9] Amir Gholami et al. “A Survey of Quantization Methods for Efficient Neural Network Inference”. In: *arXiv preprint arXiv:2103.13630* (2021).
- [10] Yunchao Gong et al. “Compressing Deep Convolutional Networks using Vector Quantization”. In: *CoRR* abs/1412.6115 (2014). arXiv: 1412.6115. URL: <http://arxiv.org/abs/1412.6115>.
- [11] Robert M. Gray and David L. Neuhoff. “Quantization”. In: *IEEE Transactions on Information Theory* 44.6 (1998), pp. 2325–2383.
- [12] Philipp Gysel et al. “Ristretto: A Framework for Empirical Study of Resource-Efficient Inference in Convolutional Neural Networks”. In: *IEEE Trans. Neural Networks Learn. Syst.* 29.11 (2018), pp. 5784–5789. DOI: 10.1109/TNNLS.2018.2808319. URL: <https://doi.org/10.1109/TNNLS.2018.2808319>.
- [13] Song Han, Huizi Mao, and William J. Dally. “Deep Compression: Compressing Deep Neural Network with Pruning, Trained Quantization and Huffman Coding”. In: *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*. Ed. by Yoshua Bengio and Yann LeCun. 2016. URL: <http://arxiv.org/abs/1510.00149>.
- [14] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. “Distilling the Knowledge in a Neural Network”. In: *CoRR* abs/1503.02531 (2015). arXiv: 1503.02531. URL: <http://arxiv.org/abs/1503.02531>.
- [15] Lu Hou, Quanming Yao, and James T. Kwok. “Loss-aware Binarization of Deep Networks”. In: *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL: <https://openreview.net/forum?id=S1oWlN911>.
- [16] Gao Huang et al. “Densely Connected Convolutional Networks”. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 2017, pp. 2261–2269. DOI: 10.1109/CVPR.2017.243. URL: <https://doi.org/10.1109/CVPR.2017.243>.
- [17] Itay Hubara et al. “Quantized Neural Networks: Training Neural Networks with Low Precision Weights and Activations”. In: *CoRR* abs/1609.07061 (2016). arXiv: 1609.07061. URL: <http://arxiv.org/abs/1609.07061>.
- [18] Benoit Jacob et al. “Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2018, pp. 2704–2713.
- [19] Chutian Jiang. “Efficient Quantization Techniques for Deep Neural Networks”. In: *Proceedings of the 2021 International Conference on Signal Processing and Machine Learning*. 2021.

- [20] Sangil Jung et al. “Learning to Quantize Deep Networks by Optimizing Quantization Intervals With Task Loss”. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*. Computer Vision Foundation / IEEE, 2019, pp. 4350–4359. DOI: 10.1109/CVPR.2019.00448. URL: http://openaccess.thecvf.com/content%5C_CVPR%5C_2019/html/Jung%5C_Learning%5C_to%5C_Quantize%5C_Deep%5C_Networks%5C_by%5C_Optimizing%5C_Quantization%5C_Intervals%5C_With%5C_CVPR%5C_2019%5C_paper.html.
- [21] Soroosh Khoram and Jing Li. “Adaptive Quantization of Neural Networks”. In: *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL: <https://openreview.net/forum?id=SyOK1Sg0W>.
- [22] Daya Shanker Khudia et al. “FBGEMM: Enabling High-Performance Low-Precision Deep Learning Inference”. In: *CoRR* abs/2101.05615 (2021). arXiv: 2101.05615. URL: <https://arxiv.org/abs/2101.05615>.
- [23] Sehoon Kim et al. “I-BERT: Integer-only BERT Quantization”. In: *Proceedings of the 38th International Conference on Machine Learning*. 2021, pp. 5506–5518.
- [24] Raghuraman Krishnamoorthi. “Quantizing deep convolutional networks for efficient inference: A whitepaper”. In: *arXiv preprint arXiv:1806.08342* (2018).
- [25] Yann LeCun, John S. Denker, and Sara A. Solla. “Optimal Brain Damage”. In: *Advances in Neural Information Processing Systems 2, [NIPS Conference, Denver, Colorado, USA, November 27-30, 1989]*. Ed. by David S. Touretzky. Morgan Kaufmann, 1989, pp. 598–605. URL: <http://papers.nips.cc/paper/250-optimal-brain-damage>.
- [26] Qun Li et al. “Investigating the Impact of Quantization on Adversarial Robustness”. In: *CoRR* abs/2404.05639 (2024). DOI: 10.48550/ARXIV.2404.05639. arXiv: 2404.05639. URL: <https://doi.org/10.48550/arXiv.2404.05639>.
- [27] Pavlo Molchanov et al. “Pruning Convolutional Neural Networks for Resource Efficient Inference”. In: *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL: <https://openreview.net/forum?id=SJGCiw5gl>.
- [28] Pierre-Emmanuel Novac et al. “Quantization and Deployment of Deep Neural Networks on Microcontrollers”. In: *CoRR* abs/2105.13331 (2021). arXiv: 2105.13331. URL: <https://arxiv.org/abs/2105.13331>.
- [29] Kazuki Okado et al. “Channel-wise quantization without accuracy degradation using Δ loss analysis”. In: *ICMLT 2022: 7th International Conference on Machine Learning Technologies, Rome, Italy, March 11 - 13, 2022*. ACM, 2022, pp. 56–61. DOI: 10.1145/3529399.3529409. URL: <https://doi.org/10.1145/3529399.3529409>.

Bibliography

- [30] Joachim Ott et al. “Recurrent Neural Networks With Limited Numerical Precision”. In: *CoRR* abs/1611.07065 (2016). arXiv: 1611.07065. URL: <http://arxiv.org/abs/1611.07065>.
- [31] Eunhyeok Park, Sungjoo Yoo, and Peter Vajda. “Value-Aware Quantization for Training and Inference of Neural Networks”. In: *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part IV*. Ed. by Vittorio Ferrari et al. Vol. 11208. Lecture Notes in Computer Science. Springer, 2018, pp. 608–624. DOI: 10.1007/978-3-030-01225-0_36. URL: https://doi.org/10.1007/978-3-030-01225-0_36.
- [32] Antonio Polino, Razvan Pascanu, and Dan Alistarh. “Model compression via distillation and quantization”. In: *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL: <https://openreview.net/forum?id=S1XolQbRW>.
- [33] Mohammad Rastegari et al. “XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks”. In: *CoRR* abs/1603.05279 (2016). arXiv: 1603.05279. URL: <http://arxiv.org/abs/1603.05279>.
- [34] Álvaro Domingo Reguero, Silverio Martínez-Fernández, and Roberto Verdecchia. “Energy-efficient neural network training through runtime layer freezing, model quantization, and early stopping”. In: *Comput. Stand. Interfaces* 92 (2025), p. 103906. DOI: 10.1016/J.CSI.2024.103906. URL: <https://doi.org/10.1016/j.csi.2024.103906>.
- [35] Wei Tang, Gang Hua, and Liang Wang. “How to Train a Compact Binary Neural Network with High Accuracy?” In: *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*. Ed. by Satinder Singh and Shaul Markovitch. AAAI Press, 2017, pp. 2625–2631. DOI: 10.1609/AAAI.V31I1.10862. URL: <https://doi.org/10.1609/aaai.v31i1.10862>.
- [36] Zhe Wang et al. “RDO-Q: Extremely Fine-Grained Channel-Wise Quantization via Rate-Distortion Optimization”. In: *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XII*. Ed. by Shai Avidan et al. Vol. 13672. Lecture Notes in Computer Science. Springer, 2022, pp. 157–172. DOI: 10.1007/978-3-031-19775-8_10. URL: https://doi.org/10.1007/978-3-031-19775-8_10.
- [37] Yi Wei et al. “Quantization Mimic: Towards Very Tiny CNN for Object Detection”. In: *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VIII*. Ed. by Vittorio Ferrari et al. Vol. 11212. Lecture Notes in Computer Science. Springer, 2018, pp. 274–290. DOI: 10.1007/978-3-030-01237-3_17. URL: https://doi.org/10.1007/978-3-030-01237-3_17.

- [38] Carole-Jean Wu et al. “Sustainable AI: Environmental Implications, Challenges and Opportunities”. In: *CoRR* abs/2111.00364 (2021). arXiv: 2111.00364. URL: <https://arxiv.org/abs/2111.00364>.
- [39] Edouard Yvinec et al. “SPIQ: Data-Free Per-Channel Static Input Quantization”. In: *CoRR* abs/2203.14642 (2022). DOI: 10.48550/ARXIV.2203.14642. arXiv: 2203.14642. URL: <https://doi.org/10.48550/arXiv.2203.14642>.
- [40] Dongqing Zhang et al. “LQ-Nets: Learned Quantization for Highly Accurate and Compact Deep Neural Networks”. In: *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VIII*. Ed. by Vittorio Ferrari et al. Vol. 11212. Lecture Notes in Computer Science. Springer, 2018, pp. 373–390. DOI: 10.1007/978-3-030-01237-3_23. URL: https://doi.org/10.1007/978-3-030-01237-3%5C_23.
- [41] Shuchang Zhou et al. “DoReFa-Net: Training Low Bitwidth Convolutional Neural Networks with Low Bitwidth Gradients”. In: *CoRR* abs/1606.06160 (2016). arXiv: 1606.06160. URL: <http://arxiv.org/abs/1606.06160>.

Bibliography

Appendix

Add additional experimental results that do not need to be directly included in the thesis body.