

Optimizing Deep Learning Model Efficiency Through Comparative Analysis of Compression Techniques: Knowledge Distillation, Quantization, and Pruning

*

Mikael Kieu¹ and Divyam Sood²

Msc. Applied Data Science

University of Gothenburg

¹Email: guskiemi@student.gu.se

²Email: divyamsood.se@gmail.com

Abstract—This study explores techniques to reduce the sizes of large deep learning models while preserving accuracy and efficiency. We evaluate three key model compression methods: Knowledge Distillation, Weight Pruning, and Quantization. After introducing each technique and its implications, we conduct an empirical analysis on a baseline CNN architecture (Resnet50) with a 525-class bird classification dataset. Our findings highlight that Knowledge Distillation effectively transfers knowledge with minimal accuracy loss, making it ideal for resource-constrained settings. Weight Pruning significantly reduces model size while maintaining high accuracy, and Quantization reduces model size without significant performance degradation. These insights underscore the importance of tailoring model compression to specific requirements for more efficient and scalable neural network deployment.

Index Terms—model compression, quantization, pruning, knowledge distillation, ResNet-50, knowledge transfer, Kullback-Leibler divergence, teacher model, student model, sparsity, fine-tuning.

I. INTRODUCTION

As artificial intelligence and deep learning evolve, it has become increasingly important to develop models that are both efficient and lightweight. Model efficiency is driven by the desire to deploy practical solutions in resource-constrained environments, such as edge devices and cloud-based services [8].

Research and development efforts have been dedicated to achieving model efficiency, resulting in three noteworthy techniques: quantization, pruning, and knowledge distillation. It is important to note that each of these methods focuses on a different aspect of model optimization, but they all share the general goal of reducing computational demands and memory consumption while maintaining the model's predictive capability.

Quantization involves representing the model's weights and activations using fewer bits, which reduces the amount of storage and computation required during inference. Pruning simplifies the model architecture by removing parameters or

neurons that do not have a significant impact on the inference. In knowledge distillation, a larger, complex model is leveraged to train a smaller model, transferring the rich information acquired by the larger model to the smaller model.

To understand the effects of these techniques, this study focuses on implementing these compression techniques on a ResNet-50 architecture [3]. Specifically, the compressed models will be used to classify bird species obtained from [2], with the goal of achieving an equivalent level of performance while substantially reducing the size and parameter count of the model. The objective of this study is to demonstrate how these techniques can be used to deploy highly efficient models without compromising their predictive accuracy. Consequently, it contributes to the broader goal of making AI more accessible and practical for real-world applications, even under resource constraints.

Following is the outline: in Section II, we describe the relevant literature and theory behind the main techniques used in the project. In Section III, the methods used in this study are described, beginning with a discussion of preprocessing and then an explanation of the main methods used to compress the model obtained. A discussion of the results is provided in Section IV, in which the model's performance is analyzed in relation to its memory size and parameter count. Additionally, the model speed will be assessed by examining both the training and inference speeds. Lastly, Section IV concludes the project and discusses potential future work.

II. BACKGROUND

A. Quantization

Quantization, the process of representing network weights and activations with a reduced number of bits, strikes a delicate balance between model accuracy and resource utilization [9]. By quantizing the model's parameters, one can drastically decrease memory consumption and accelerate inference, albeit at the cost of introducing quantization-induced errors. While

floating-point numbers can represent a wide range of values with high precision, quantized integers can only represent a limited range of values with low precision. This error occurs as a result of approximating real numbers with a limited number of discrete values. A variety of techniques are available to minimize the impact of this error [5].

The quantization process can take many forms, including weight quantization (reducing the precision of model parameters), activation quantization (reducing the precision of intermediate activations), and hybrid quantization (combining different quantization methods for weights and activations) in order to optimize memory and computation.

Quantization typically has two forms of training:

- 1) Quantization-Aware Training (QAT) is a technique where the model is trained with quantization effects incorporated into the training process. This allows the model to learn to be robust to quantization and can mitigate the accuracy loss that typically accompanies quantization.
- 2) Post-Training Quantization (PTQ) applies quantization after the model has been trained with full precision. This method is simpler but may result in more substantial accuracy loss compared to QAT.

B. Knowledge distillation

The concept of knowledge distillation refers to a process where the knowledge acquired by a complex "teacher" model is transferred to a simpler model in order to address the inherent trade-off between model performance and model efficiency [4]. As its name implies, knowledge distillation is the process of distilling the expertise and generalization capabilities of a deep and complex model into a more compact and computationally efficient version.

The knowledge distillation process uses soft targets rather than focusing exclusively on the final predictions of the teacher model. Thus, the teacher's predictions are treated as probability distributions over classes rather than hard labels. Due to this softness, the student model can capture not only what the teacher predicts, but also the uncertainty and nuances associated with those predictions.

An important aspect to knowledge distillation is the loss function, which is typically determined by the Kullback-Leibler divergence or another similarity metric. It is used to determine the degree of alignment between the student's predictions and the teacher's soft targets. The student model strives to minimize this loss during training, thereby effectively learning from the teacher.

There is a temperature parameter (often abbreviated as "T") that determines how soft the teacher's predictions are. When T is higher, the teacher's predictions are softer and more informative, while when T is lower, they are more deterministic. A balance can be struck between fidelity to the teacher and generalization by tuning this parameter. In general, knowledge distillation results in improved performance for the student model as compared to training it from scratch. This is because the student learns not only from the labeled data but

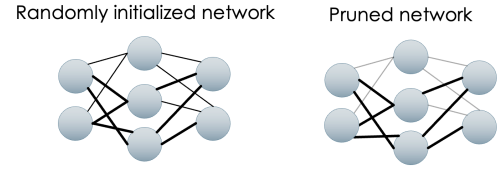


Fig. 1. Illustration of pruning a randomly initialized network based on the lowest absolute values. The width of the connecting lines represent absolute value size.

also from the rich insights and generalization abilities of the teacher.

C. Pruning in DNN

In neural networks, pruning is a crucial technique that is used to systematically eliminate network parameters which contribute minimally to the network's output. By reducing redundancy within the architecture, it is primarily intended to enhance the efficiency of the network. The application of pruning in a judicious manner results in a neural model that retains its predictive accuracy while drastically reducing computational demands during inference [6].

1) *Lottery ticket hypothesis*: The lottery ticket hypothesis introduces an interesting concept for pruning a model [1]. It states that within any randomly initialized feed-forward neural network, there exists a specific sub-network, often referred as the "winning ticket." This winning ticket possesses distinct properties that set it apart from the larger network. After training, the winning ticket will generalize better or in other words it will have higher test accuracy. Moreover, training of the "winning ticket" will require fewer iterations than training the original network. Lastly, the winning ticket typically has a significantly reduced number of parameters, usually less than 10% to 20% of the original network's parameters. Parameters are remarkably reduced in this hypothesis, which is one of its key features.

2) *Finding winning ticket*: Although the lottery ticket hypothesis asserts that winning tickets exist within neural networks, it does not provide an explicit recipe for how they are discovered. A systematic method based on pruning has been proposed by authors in [4] and it follows:

- 1) The process begins with the original neural network, initialized with random weights.
- 2) This network is trained for a fixed number of training iterations, allowing the weights to evolve and adapt to the given dataset.
- 3) Pruning is then employed to systematically eliminate weights, usually based on criteria like the lowest absolute values (L1 norm), and this is typically done on a layer-by-layer basis.
- 4) Throughout this process, a record of the pruning mask is maintained. The pruning mask indicates which weights were pruned, while the pruned weights themselves are discarded.

- 5) After pruning, researchers revert to the original network's weights and apply the pruning mask to retain only the relevant connections.
- 6) Steps 1 to 5 are repeated for multiple iterations, a process known as iterative pruning. These iterations ultimately lead to the discovery of a network with a significantly reduced number of non-pruned weights.

III. METHODS

A. Dataset Selection and Characteristics

The model compression study presented here uses a dataset obtained from [4] that contains 525 bird species. It was intentionally chosen due to its inherent complexity, featuring a large number of classes, making it an ideal benchmark.

In this dataset, 84,635 training images, 2,625 test images (5 per species), and an equivalent number of validation images are included. Each image is uniformly formatted in JPG, maintaining a consistent dimension of 224 x 224 pixels, with RGB color channels. Moreover, each image focuses on a single bird species, and is thoughtfully framed to ensure that the bird occupies at least 50% of the image area.

B. Baseline Model - ResNet-50

In the initial phases of our experimentation, we explored various model architectures, including MobileNetV2, ResNet-50, SqueezeNet, and ResNet-152, with the goal of establishing a suitable baseline model. We systematically replaced the original classification head with a custom one capable of producing 525 outputs, corresponding to the number of unique bird species in the dataset. Three distinct training strategies were employed:

- Freezing all layers except the final one, utilizing pre-trained weights from ImageNet.
- Unfreezing all layers and training with pre-trained weights from ImageNet.
- Unfreezing all layers and training with randomly initialized weights.

Ultimately, our findings indicated that the second approach, involving the unfreezing of all layers and training with pre-trained ImageNet weights, yielded the highest accuracy, reaching an impressive 98%. The choice of the ResNet-50 architecture was based on its well-documented success in image classification tasks, thanks to its deep architecture and outstanding performance. Additionally, standard data augmentation techniques were applied during training, and the model was optimized using cross-entropy loss function.

C. Quantization

The baseline model was a ResNet-50 model previously trained on the bird classification dataset. To prepare the model for quantization, it was transitioned into evaluation mode and configured with a suitable quantization setting optimized for the 'x86' hardware platform.

The calibration phase followed, wherein a subset of training data was designated as a calibration dataset. Approximately 10 batches of this dataset were utilized to collect activation

TABLE I
COMPARISON OF BASELINE AND QUANTIZED RESNET-50 ARCHITECTURE

	Inference time	Accuracy
Baseline	0.1672 s	98.63%
Quantized	0.069 s	98.43%

statistics for quantization parameter determination. Activation fusion was then applied to the model, specifically targeting convolutional layers with ReLU activation functions to consolidate operations and enhance computational efficiency. Subsequently, the prepared model underwent conversion to a quantized form using PyTorch's quantization.convert function. This transformation involved converting both model weights and activations into lower-precision formats, typically 8-bit integers, effectively reducing memory requirements. Following quantization, the performance of the quantized model was evaluated on the test dataset, assessing inference speed and classification accuracy, the results can be found in table III-C.

D. Knowledge distillation

This approach aimed to distill knowledge from a powerful ResNet-50 teacher model, originally trained on a comprehensive bird classification dataset, into a more compact SqueezeNet student model. A pretrained ResNet-50 teacher model was loaded and acted as the source of "knowledge". In parallel, a SqueezeNet student model was built, specifically tailored for the classification of birds. Designed to be more lightweight with fewer parameters.

In the training phase, the student model was initially untrained and yielded an accuracy score of 70% on the validation set after 30 epoch, the results are shown in figure 2. During training, the student model's objective was twofold: mimic the teacher's predictions while maintaining competitive classification performance. The loss function employed was a combination of the Kullback-Leibler Divergence (KLDivLoss) between the softmax logits of the teacher and student models, scaled by temperature, and the standard cross-entropy loss. By combining these factors, the student model inherited the teacher's expertise and retained its own discriminatory abilities. The training process was carried out for 30 epochs, allowing the student model to progressively distill the knowledge from the teacher, as depicted in figure 3 .

As a result of knowledge distillation, a more lightweight and computationally efficient SqueezeNet student model was developed. Despite its smaller size, the model was still able to absorb the valuable insights from the ResNet-50 teacher model.

As shown in table II. the baseline model has an accuracy of 97% and SqueezeNet without distillation reached an accuracy score of 60% at 30 epochs. The accuracy of ResNet rose to 84% when ResNet50 was used as a teacher model, thus showing a significant improvement in accuracy along with a substantial reduction in size.

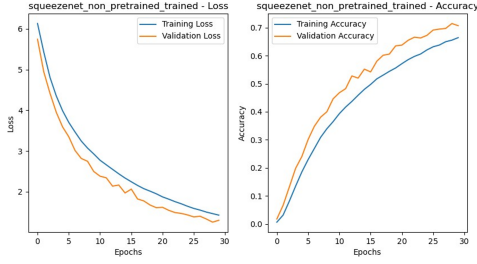


Fig. 2. Comparison of accuracy and loss for training and validation data on Squeezenet network

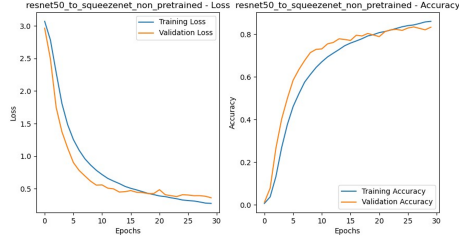


Fig. 3. Comparison of accuracy and loss for training and validation data on a distilled Squeezenet network

E. Pruning

A GitHub repository on pruning obtained from [5] was influential in our work to prune the ResNet-50 model. The script was initially modified to adjust the first convolutional layer to make the model prunable. We computed initial model's accuracy while measuring global sparsity. The implementation provides both layer-wise and global pruning options as part of the pruning process. Using layer-wise pruning, it systematically prunes a specified percentage of weights in each module (Conv2d or Linear) based on their L1 norms, effectively eliminating less important connections. Alternatively, global pruning prunes a specified percentage of weights across all layers at the same time. Following each pruning iteration, fine-tuning was performed in order to recover any accuracy loss due to pruning by adjusting the remaining weights while applying L1 and L2 regularization for sparsity and weight decay. Due to the excessive running time, this iterative process was limited to six iterations. As presented in table III, the finalized pruned model resulted in a global sparsity of 57.58% while achieving a higher accuracy score than the baseline model.

IV. RESULTS

In conclusion, the series of techniques applied to the neural networks in this project has led to significant improvements

TABLE II
COMPARISON OF MODEL SIZE AND PARAMETER

	Parameters	Size
ResNet-50 (Teacher model)	24 583 757	93.98 Mb
SqueezeNet (Student model)	1 004 749	3.83 Mb
StudentNet (Student model)	281 613	1.07 Mb

TABLE III
COMPARISON OF BASELINE AND PRUNED RESNET-50 ARCHITECTURE

Iteration	Global sparsity	Validation accuracy
0	0.00%	97.83%
1	19.56%	98.67%
2	35.19%	98.43%
3	47.72%	98.51%
4	57.78%	99.04%

in both efficiency and accuracy, underscoring their practical importance. The incorporation of knowledge distillation, where complex knowledge from a model like ResNet-50 was distilled into a smaller model like SqueezeNet, resulted in remarkable accuracy enhancements. Notably, SqueezeNet's accuracy surged from an initial 70% to an impressive 84% following knowledge distillation. This demonstrates the efficacy of knowledge transfer and its impact on model performance. Quantization, another instrumental approach employed in this project, considerably reduced the model's size while preserving a high level of accuracy. This not only contributes to a more streamlined model but also highlights the potential for practical deployment where resource constraints are a concern. Furthermore, the technique of pruning was pivotal in optimizing the neural network architecture. By systematically removing redundant layers in the pruned model and iteratively fine-tuning it, we achieved enhanced accuracy beyond the baseline ResNet-50.

For those interested in replicating our results and findings, we have made the project code accessible on our GitHub repository: Github/DML Project Repository, complete with a comprehensive readme. While we've made significant strides within the constraints of time, our future endeavors aim to explore additional methods, such as Low-Rank factorization, and to delve into more nuanced sub-categories within knowledge distillation, pruning, and quantization techniques [7].

REFERENCES

- [1] FRANKLE, J., AND CARBIN, M. The lottery ticket hypothesis: Finding sparse, trainable neural networks, Mar 2019.
- [2] GERRY. Birds 525 species- image classification, Apr 2023.
- [3] HE, K., ZHANG, X., REN, S., AND SUN, J. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016).
- [4] HINTON, G., VINYALS, O., AND DEAN, J. Distilling the knowledge in a neural network, Mar 2015.
- [5] JIANG, C. Efficient quantization techniques for deep neural networks. *2021 International Conference on Signal Processing and Machine Learning (CONF-SPML)* (2021).
- [6] SIETSMA, AND DOW. Neural net pruning-why and how. *IEEE International Conference on Neural Networks* (1988).
- [7] SWAMINATHAN, S., GARG, D., KANNAN, R., AND ANDRES, F. Sparse low rank factorization for deep neural network compression. *Neurocomputing* 398 (2020), 185–196.
- [8] TAYLOR, B., MARCO, V. S., WOLFF, W., ELKHATIB, Y., AND WANG, Z. Adaptive deep learning model selection on embedded systems. *Proceedings of the 19th ACM SIGPLAN/SIGBED International Conference on Languages, Compilers, and Tools for Embedded Systems* (2018).
- [9] YANG, B., LIU, J., ZHOU, L., WANG, Y., AND CHEN, J. Quantization and training of object detection networks with low-precision weights and activations. *Journal of Electronic Imaging* 27, 01 (2018), 1.