# Neural Networks at a Fraction with Pruned Quaternions

Sahel Mohammad Iqbal, Subhankar Mishra

National Institute of Science Education and Research, India

{sahelm.iqbal, smishra}@niser.ac.in

## Abstract

Contemporary state-of-the-art neural networks have increasingly large numbers of parameters, which prevents their deployment on devices with limited computational power. Pruning is one technique to remove unnecessary weights and reduce resource requirements for training and inference. In addition, for ML tasks where the input data is multi-dimensional, using higher-dimensional data embeddings such as complex numbers or quaternions has been shown to reduce the parameter count while maintaining accuracy. In this work, we conduct pruning on real and quaternion-valued implementations of different architectures on classification tasks. We find that for some architectures, at very high sparsity levels, quaternion models provide higher accuracies than their real counterparts. For example, at the task of image classification on CIFAR-10 using Conv-4, at 3% of the number of parameters as the original model, the pruned quaternion version outperforms the pruned real by more than 10%. Experiments on various network architectures and datasets show that for deployment in extremely resource-constrained environments, a sparse quaternion network might be a better candidate than a real sparse model of similar architecture.

## 1 Introduction

A key attribute of any neural network architecture is the number of trainable parameters that it has. In general, the greater the number of model parameters, the greater its demands on computational power, time, and energy to train and perform inference. Contemporary state-of-the-art deep neural networks have model parameters that often run into tens or even hundreds of millions [24, 36, 5], imposing great demands on the hardware needed to train these models.

There are several real-world scenarios where we would like to deploy well-performing models to edge devices such as mobile phones. An example would be when dealing with private user data such as images, where sending data to a back-end datacenter for inference would be less than ideal [42]. The best (and biggest) models cannot be run on these resource-constrained computing environments [16]. This leaves us with two options - either use smaller, more specialized architectures (as in MobileNet [16]) or find ways to compress the bigger ones.

There are multiple compression methods by which we can reduce the resource consumption of a model such as parameter pruning and sharing, low-rank factorization and knowledge distillation [3], but in this work, we focus on pruning. Pruning is a method to reduce the number of parameters in a model by removing redundant weights or neurons [20]. Various studies have shown that pruning can drastically reduce model parameter counts while still maintaining accuracy [13, 23, 2]. These pruned models can also be re-trained [7], helping to reduce resource utilization during the training stage as well.

What happens when we prune a model to extreme levels of sparsity, say 90% or more? At this level, we typically see that the accuracy drops off [7, 13], and that the pruned model can no longer match the original model. For this reason, most studies on pruning only prioritize up until the point where the pruned model is no longer on par with the parent [7, 13, 20, 23]. However, we feel that this regime is still attractive because various empirical studies have found that a large-pruned model consistently does better than a small-dense model of equal size [42, 22, 12]. Thus a state-of-the-art model pruned to just 2% of its original size might still provide better accuracy than a miniature model of comparable size, even though the pruned model displays lower accuracy than the original.

Recently, another method of reducing model parameters is undergoing a surge in popularity. Using higher-dimensional data embeddings, such as complex numbers or quaternions, has been successfully shown to reduce model parameters while maintaining accuracy [39, 38, 10, 27]. Quaternions are a 4-dimensional extension to the complex numbers introduced by the mathematician William Rowan Hamilton in 1843 [27], and quaternion neural networks have been built for a variety of ML tasks [43, 29, 10, 28, 32, 4]. Converting a real model to quaternion can lead to a 75% reduction in model parameters (which is explained in more detail in Sec. 3.2), making it a suitable method for model compression.

In the present work, we employ pruning on quaternion networks to see if they have any advantages over their real counterparts at high levels of model sparsities. To the extent of our knowledge, there are no prior studies that explore weight-reduction in neural networks by combining pruning with quaternion-valued neural networks (or any other higher-dimensional data structure). We choose multiple neural networks for image classification on the MNIST [19], CIFAR-10 and CIFAR-100 [18] datasets, build equivalent quaternion representations, and conduct pruning experiments on both real and quaternion implementations. We find that at extreme sparsities (approximately 10% or fewer parameters as the real, unpruned model), the quaternion model outperforms the real. Thus for deploying in a resource-constrained device, a quaternion pruned model might provide the best accuracy out of all available options.

Specifically, the contributions of this work can be summarized as follows. We conduct pruning experiments on real and quaternion-valued implementations of different network architectures. Through this we show that 1) the lottery ticket hypothesis [7] is valid for quaternion models, meaning that pruned quaternion models can be re-trained from scratch to the same accuracy as the unpruned model, and 2) at very high model sparsities, the quaternion equivalent displays higher accuracy than the real network.

# 2 Related Work

## 2.1 Pruning

From the early 1990s, we have known that the majority of weights in a trained neural network can be pruned without sacrificing its accuracy [20, 34, 14]. Earlier works were done on simple architectures, but recently pruning has also been shown to work on much more extensive architectures such as VGG [37] and ResNet [15, 7, 2]. These studies showed that modern state-of-the-art architectures are often heavily over-parameterized and that they only require a far fewer number of parameters to learn the necessary function representations [7].

In most studies, authors generally attempt to prune after the training process, at the end of which the weights would be ordered based on their contribution to the output. One common example of a heuristic for such ordering is the weight magnitude, where the contribution of individual weights to the output are judged based on their absolute values [13]. However, when training pruned networks from scratch, often they could not match the original accuracy, and the pruned models did much worse [13, 23]. The 'Lottery Ticket Hypothesis' paper [7] showed that these pruned networks could indeed be trained from the beginning, but we had to be pair the model with the initial weights for the unpruned network. This work showed that the benefits from pruning could be realized during the training process, which could potentially reduce the resource requirements of training by a huge amount.
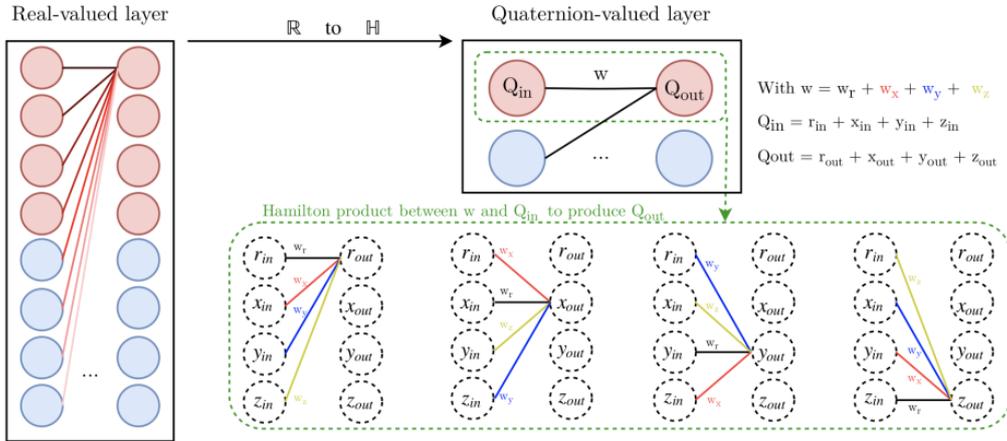


Figure 1: Depiction of how using different linear combinations of coefficients in the Hamilton product results in a reduction in the number of weights in a neural network. Image copyright Tituan Parcollet [29], reproduced with permission.

Pruning is generally of two types. In structured pruning, weights are pruned in groups by removing whole neurons or entire channels [35, 23]. Structured pruning leads to a reduction in the size of the model and improves model inference speeds as the dimensions of weight matrices are reduced [25]. Unstructured pruning, by contrast, is where individual weights of neurons are removed instead of entire neurons or channels [20, 13]. While unstructured pruning reduces the number of parameters, this does not immediately manifest itself through an improvement in inference speeds [30]. This is because unlike structured pruning, here weight matrices remain

the same dimensions but are instead simply made sparse, which current hardware technology is not capable of optimizing [30]. However, empirical studies have shown that unstructured pruning often yields much better results than structured [35]. In addition, with both pruning methods, the resource demands for training are reduced because we now have a far smaller set of weights that we need to optimize. Research to optimize sparse operations on current hardware has shown promising results [6], and hardware that can accelerate sparse-matrix multiplications are being built [33]. This suggests that all the predicted theoretical performance gains from pruning could soon be realized in the near future.

An appealing property of pruned networks, and the one that justifies this work, is that in general, a large-sparse (pruned) model performs better than a small-dense (unpruned) one with an equal number of weights [2]. Multiple studies have shown that for a variety of model architectures on different tasks, sparse models consistently outperform their dense counterparts [42, 22, 17, 12]. This implies that given a resource constraint on the size of the models that a user can run, their best bet at achieving the greatest possible accuracy would be to use a large-pruned model rather than a small-dense one. This is applicable even when the pruned model cannot match the accuracy of the original, unpruned model.

## 2.2 Quaternions

In recent years, there has been a marked increase in works that address the question of whether it is more optimal to use multi-dimensional data embeddings in applications where the input data is multi-dimensional (refer [27] for a comprehensive review). For example, Trabelsi et al. [38] demonstrated that at the task of music transcription, where the input signal is two-dimensional (consisting of magnitude and phase of the signal), complex-valued neural networks outperformed comparable real models. Complex numbers, however, are insufficient to represent higher-dimensional inputs such as the three channels of a color image, which is why some studies extended this idea to four-dimensional quaternions. Zhu et al. [43] compared quaternion and real-valued convolutional neural networks (henceforth referred to as $Q$ and $R$ respectively) with similar architectures and the same number of parameters on the CIFAR-10 dataset. They found that $Q$ achieved faster convergence on the training loss as well as higher classification accuracy on the test set compared to $R$. Gaudet and Maida [10] made a similar comparison with image classification on the CIFAR-10 and CIFAR-100 datasets and image segmentation on the KITTI Road Segmentation dataset [9], but this time with $Q$ having a quarter of the number of parameters as $R$. They reported that on both tasks, quaternion models gave higher accuracy than real and complex networks while having a lower parameter count. Similar advantages for quaternion neural networks over real networks were also found by Parcollet et al. [29] for speech recognition.

# 3 Theory of Quaternions

## 3.1 Quaternion Algebra

Quaternions are a four-dimnensional extension to the complex numbers, and a general quaternion $q$ may be written as

$$q = r + x\mathbf{i} + y\mathbf{j} + z\mathbf{k} \tag{1}$$

where $r, x, y, z \in \mathbb{R}$ and $\mathbf{i}$, $\mathbf{j}$ and $\mathbf{k}$ are complex entities which follow the relations

$$\mathbf{i}^2 = \mathbf{j}^2 = \mathbf{k}^2 = \mathbf{ijk} = -1 \tag{2}$$

Given two quaternions $q_1$ and $q_2$, their product (known as the Hamilton product) is given by

$$\begin{aligned}
q_1 \otimes q_2 = \ & (r_1 r_2 - x_1 x_2 - y_1 y_2 - z_1 z_2) \\
& (r_1 x_2 + x_1 r_2 + y_1 z_2 - z_1 y_2)\ \mathbf{i} \\
& (r_1 y_2 - x_1 z_2 + y_1 r_2 - z_1 x_2)\ \mathbf{j} \\
& (r_1 z_2 + x_1 y_2 - y_1 x_2 - z_1 r_2)\ \mathbf{k}
\end{aligned} \tag{3}$$

Unlike real and complex multiplication, quaternion multiplication is not commutative. Quaternions can be represented as $4*4$ real matrices such that the matrix multiplication between such representations are consistent with the Hamilton product.

$$q = \begin{bmatrix} r & -x & -y & -z \\ x & r & -z & y \\ y & z & r & -x \\ z & -y & x & r \end{bmatrix} \tag{4}$$

This is the representation that is used to calculate quaternion products in a quaternion-valued neural network.

## 3.2 How does weight-reduction happen?

Consider a small section of a fully-connected neural network with four input and four output neurons. In the case of a real-valued implementation, this layer would require $4*4 = 16$ weights. However, if we were to view the above network as consisting of one input quaternion and one output quaternion, then using the Hamilton product, we would only require one quaternion weight, or four real weights, to connect them. Thus provided the number of neurons in all layers are divisible by 4, we can obtain a 75% reduction in the number of parameters of a network by converting it to a quaternion implementation. This is shown graphically in Fig. 1. This weight-reduction is explained in greater detail in [29, 27].

# 4 Methodology

Our experiments were run on classification tasks on the MNIST [19], CIFAR-10 and CIFAR-100 [18] datasets. We chose to demonstrate our experiments at image classification tasks following the lead of some of the most important works in network

| Model | Lenet-300-100 | Conv-2 | Conv-4 | Conv-6 |
|---|---|---|---|---|
| Datasets | MNIST | CIFAR-10 | CIFAR-10 CIFAR-100 | CIFAR-10 CIFAR-100 |
| Conv Layers | | $2*64$, pool | $64*64$, pool $128*128$, pool | $64*64$, pool $128*128$, pool $256*256$, pool |
| FC Layers | 300, 100, 10 | 256, 256, 10/100 | 256, 256, 10/100 | 256, 256, 10/100 |
| All/Conv Weights (Real) | 266.6K | 4.30M/38K | 2.42M/260K | 2.26M/1.14M |
| All/Conv Weights (Quat) | 67.7K | 1.08M/9.9K | 609K/65K | 569K/287K |
| Training epochs/Batch size | 40/60 | 40/60 | 40/60 | 60/60 |
| Optimizer/Learning rate | Adam/1.2e-3 | Adam/2e-4 | Adam/3e-4 | Adam/3e-4 |

Table 1: Model architectures, datasets and hyper-parameters tested in this paper. The number of weights for Conv-2,4,6 are reported for CIFAR-10 classification. Architectures and hyper-parameters have been kept the same as those used in [7] to facilitate a direct comparison.

compression [7, 8] so that our results may be contrasted with state-of-the-art. We used the fully-connected Lenet-300-100 architecture from [21], and the Conv-2, Conv-4, and Conv-6 convolutional models from [7]. Complete details about model architectures and hyper-parameters are given in Tab. 1.

Our goal in this work was to compare pruning on real and quaternion-valued implementations of different architectures. To do this, we first constructed quaternion equivalents of every model with the condition that they both have the same number of real neurons. $Q$ would thus have one-fourth the number of quaternion neurons (since four real neurons constitute a quaternion neuron) and one-fourth the number of weights as $R$ (because of the weight-sharing property of the Hamilton product explained in Sec. 3.2). We used a real output layer for $Q$ because for the MNIST and CIFAR-10 datasets, the number of output labels (10) is not divisible by 4. An alternative option here was to use 10 quaternion neurons for the output layer and then take their norms. We chose not to do this because this would break the equality condition that we just stated, that the number of real neurons across implementations be equal. All other network layers were replaced by equivalent quaternion implementations. The hyper-parameters used for training are also the same as [7] in order to make direct comparisons.

For the MNIST dataset, since the images are grayscale and have only one channel, the input image is flattened and each set of four pixels are fed to a quaternion neuron of Lenet-300-100. On the other hand, for color vision tasks such as classification on the CIFAR-10, it makes more sense to treat the RGB channels of each pixel as belonging to a single quaternion neuron. To do this, we need to add one more channel to the input images, and the equality condition has to be relaxed for the input channel. A few different options exist, such as a channel with all 0s or using an additional layer to learn the fourth channel [10]. For our implementation, we chose to use the grayscale values of the input image as the additional channel.

For pruning experiments, we use iterative pruning, where the model is pruned and

then re-trained for a certain number of epochs over multiple iterations. We prune the networks using a global pruning technique with 20% of the weights in the model pruned in each pruning iteration. We chose global pruning over layer-wise pruning because global-pruning can find smaller lottery tickets for larger networks [7], and we wanted to keep the methodology consistent across all the architectures that we test. We only prune weights and exempt biases, as biases constitute only a tiny proportion of the total parameters in a model. At each level of model sparsity, the pruned model is re-trained from scratch using the initial weights for the same duration as the original model (as done in [7]), and it is the accuracy of these re-trained pruned models that we have reported throughout this paper. Our emphasis is on re-trainable pruned models because our primary concern is with the practical benefits of pruning conferred during training. The models are pruned until their accuracies drop below a threshold (which we chose to be 30%) for two successive pruning iterations. This was done to save computational resources by preventing pruning the model beyond the point where it is of any practical use.

Pruning experiments were conducted using PyTorch [31]. PyTorch implementations of the various operations and building blocks necessary to build quaternion convolutional neural networks have been borrowed from the hTorch library (MIT License) [41]. This library uses a split activation function [1] where ReLU [11] is applied individually to the four different components of each quaternion. The back-propagation algorithm employed for quaternions is a generalization of those for real and complex networks [26, 28, 38]. All experiments were carried out on a single Nvidia RTX 3090 GPU.
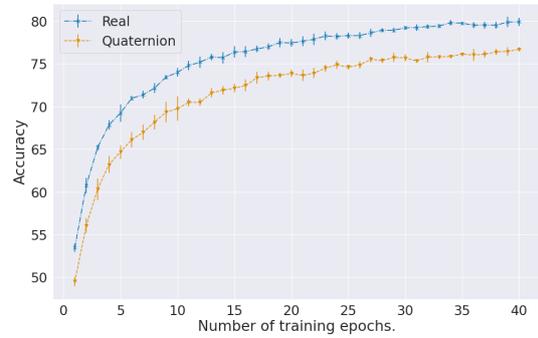
# 5    Experimental Results

A few previous studies using quaternions had found that at certain tasks, quaternions can outperform real networks while having an equal or smaller number of parameters, including for computer vision tasks such as classification which we consider in this paper [43, 10]. In our experiments, where $Q$ has a quarter of the number of parameters as $R$, we found that $Q$ cannot, in general, be trained to the same accuracy as $R$ within a fixed number of iterations. The training results for real and quaternion implementations of the various models are reported in Fig. 2. The relative accuracy difference between real and quaternion implementations are different for different models. We also found that this varied depending on the hyper-parameters used, and so it may be possible that for some set of hyper-parameters, $Q$ could be made as accurate as $R$, reproducing the results given in the references mentioned earlier. We were unable to find this subset in our experimentation.

On re-training pruned models from scratch, we found that just like $R$, pruned $Q$ are capable of matching or exceeding the accuracy of the original, unpruned model (Fig. 3). For example, in Fig. 3a, a pruned model with 51% weights ends training with higher accuracy than the unpruned model. This result held for all six architecture-dataset pairs that we tested. This shows that the lottery ticket hypothesis is valid for quaternion-valued networks as well.
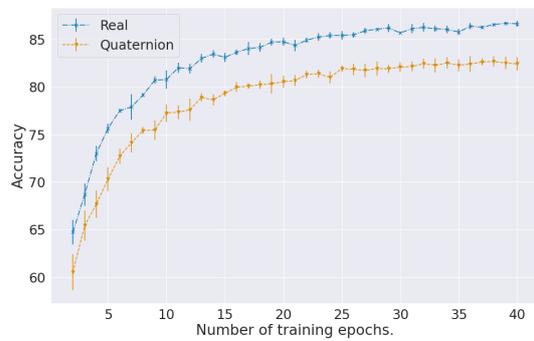
The accuracies of pruned $Q$ and $R$ at various model sparsities are given in Fig. 4. On examining this figure, certain patterns become evident. Consider Conv-4 on CIFAR-10 (Fig. 4c) as an illustrative example. The curve for $Q$ starts at the 25%
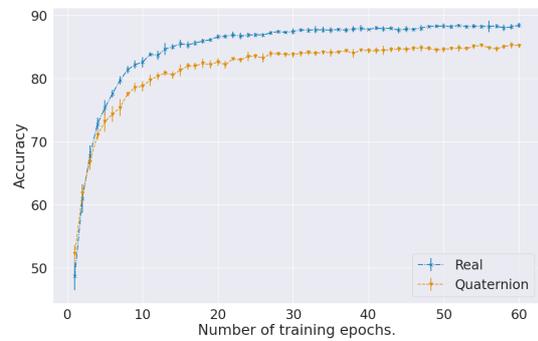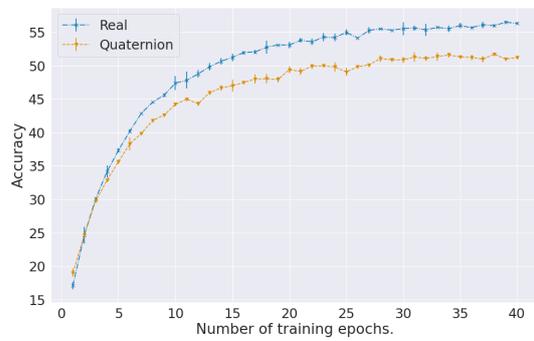
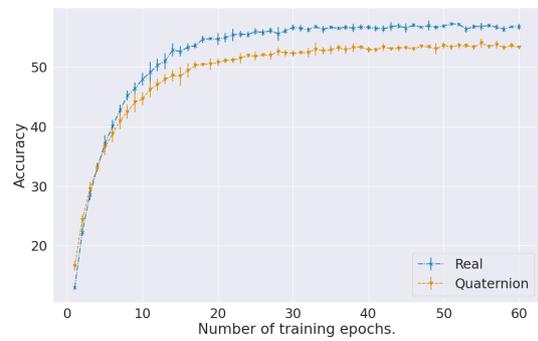(a) Lenet-300-100 on MNIST

(b) Conv-2 on CIFAR-10

(c) Conv-4 on CIFAR-10

(d) Conv-6 on CIFAR-10
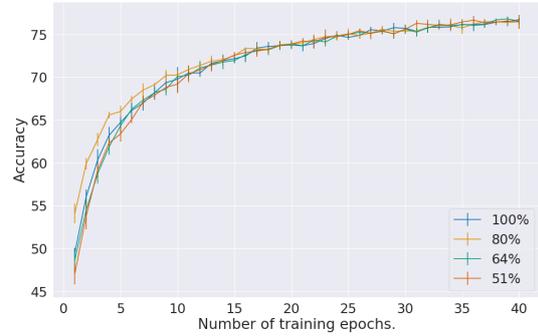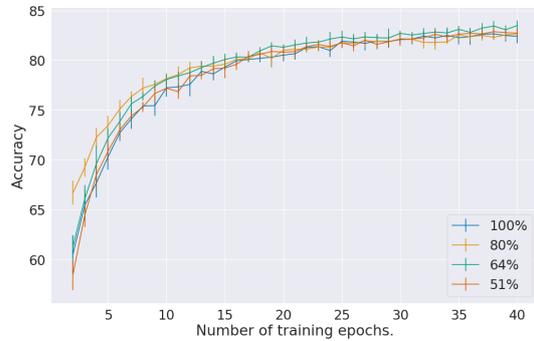
(e) Conv-4 on CIFAR-100

(f) Conv-6 on CIFAR-100

Figure 2: Training results for various architectures for real and quaternion implementations. Results are the mean over 5 trials, and the error bars are the standard deviation.
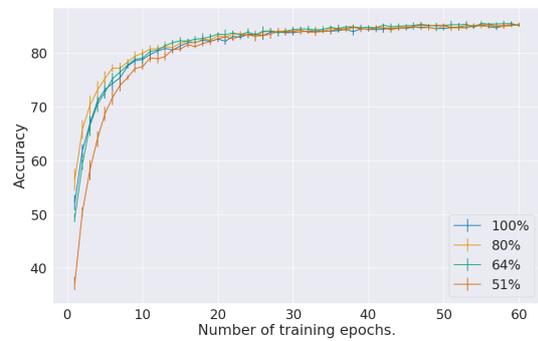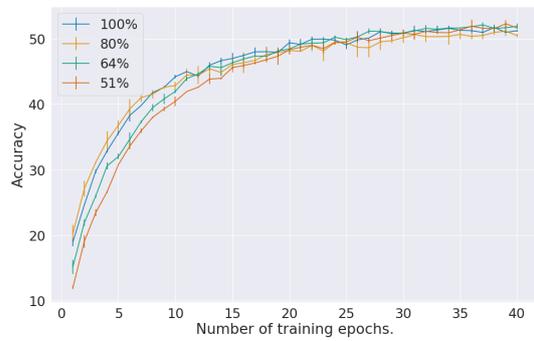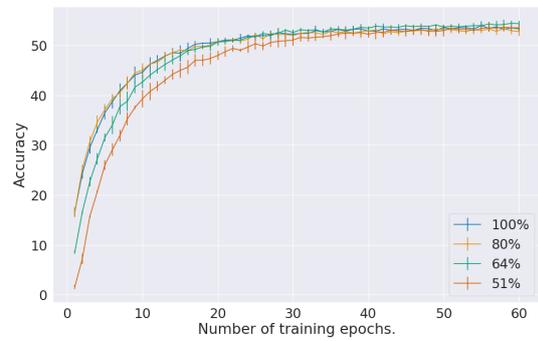
(a) Lenet-300-100 on MNIST

(b) Conv-2 on CIFAR-10
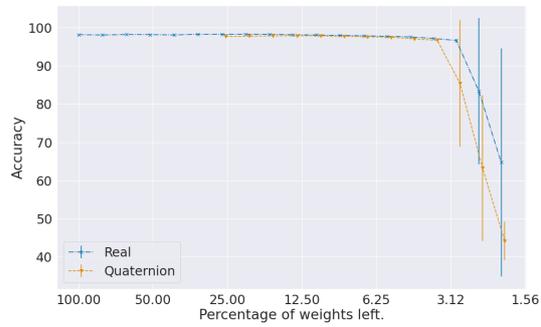
(c) Conv-4 on CIFAR-10

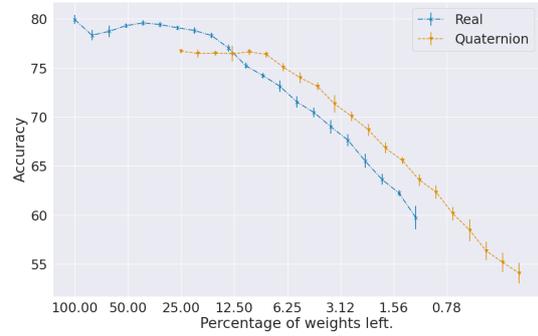(d) Conv-6 on CIFAR-10

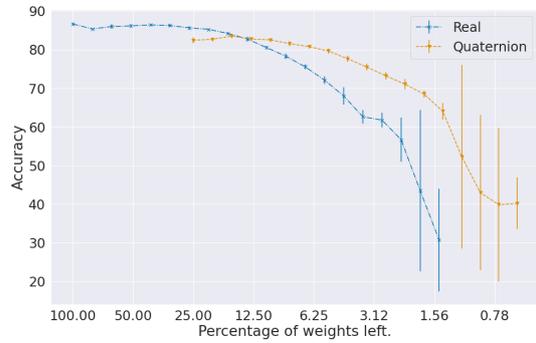(e) Conv-4 on CIFAR-100

(f) Conv-6 on CIFAR-100

Figure 3: Training results for quaternion implementations of models at different sparsities. Pruned models have been re-trained from scratch with the intial weights. Plot labels are the percentage of weights remaining after pruning. Results are the mean over 5 trials, and the error bars are the standard deviation.
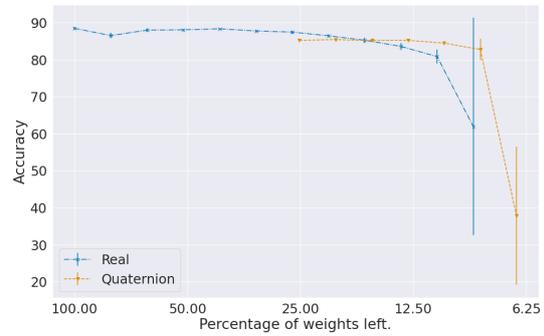
(a) Lenet-300-100 on MNIST

(b) Conv-2 on CIFAR-10

(c) Conv-4 on CIFAR-10

(d) Conv-6 on CIFAR-10

(e) Conv-4 on CIFAR-100

(f) Conv-6 on CIFAR-100

Figure 4: Accuracy vs sparsity results for real and quaternion implementations of various architectures. Quaternion curves start at around the 25% mark because unpruned $Q$ have only one-fourth (approx.) the number of parameters as the corresponding unpruned $R$. Models are pruned until their accuracies drop below 30% for two successive pruning iteraions.

mark as by construction, $Q$ only has that many parameters compared to $R$. $Q$ also starts out with lower accuracy than $R$, which can also be seen in Fig. 2c. The regime we are interested in is that of high pruning rates, at around 10% of total weights or lower. Though both models start out at different accuracies, as we get to the 12.5% mark, the curves for $R$ and $Q$ coincide, meaning that at this sparsity level, both perform equally well. On further pruning, the $R$ curve drops below the $Q$ curve, implying that the quaternion outperforms the real model at very high sparsity levels. At about 3.12% sparsity, $Q$ shows close to 75% accuracy while it is around 62% for $R$, a difference of more than 10%. This overall pattern is also repeated for all but one of the model-dataset pairs that we tested, the sole exception being Lenet-300-100 on MNIST (Fig. 4a).
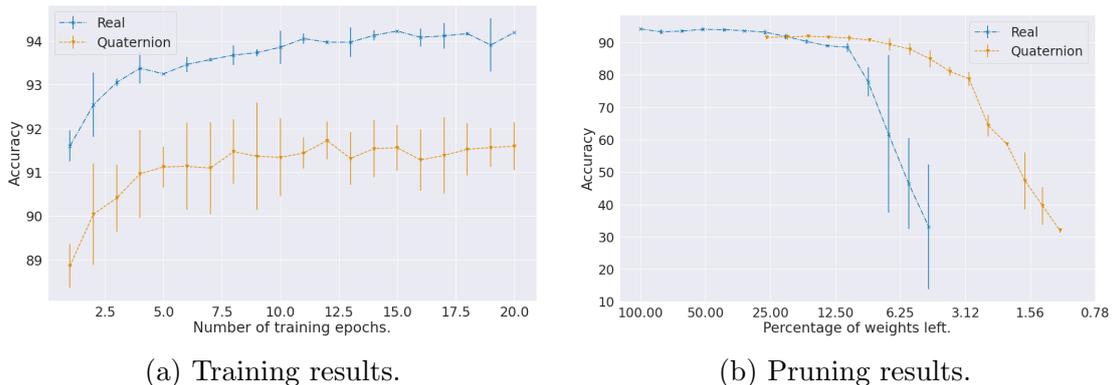


(a) Training results.   (b) Pruning results.

Figure 5: Training and pruning results for Lenet-12 on MNIST.

For Lenet-300-100, however, $R$ outperforms $Q$ at every level of sparsity. This model is different from the other three models that we tested in two ways: 1) it is a fully-connected network (no convolutional layers), and 2) it displays around 98% accuracy on its task, which is considerably higher than any of the other model-dataset pairs. To isolate which of these two properties led to the difference in the sparsity-accuracy trend, we ran pruning experiments on a custom Lenet-12 model (which has a single hidden layer with 12 real neurons), which is also a fully-connected model but with lower accuracy. The results for this model on the MNIST dataset are given in Fig. 5. Here the earlier trend reappears, and $Q$ performs better than $R$ at high sparsity levels. Thus Lenet-300-100 showed divergent behavior not because it is a fully-connected network but because it has very high accuracy at its task. A possible explanation for this may be that this model is so over-parameterized that the real model can be pruned to a great extent without a significant drop in accuracy. This explanation, however, cannot justify why $Q$ underperforms $R$ beyond the 4% sparsity region for Lenet-300-100.

On the whole, this analysis shows that, in general, at extreme model sparsities, quaternion models perform better than their real counterparts. Although quaternion models start out at a disadvantage because of their lower initial accuracy, as the models are reduced in size, their relative performance gap diminishes until at a certain point the real model dips below the quaternion, where it remains for the rest of the pruning process.

## 5.1   Using Early Stopping

In [7], the authors use an early stopping criterion to stop training. The criterion used is that of the minimum validation loss. We repeated our experiments with the same early stopping criterion (with a patience of 10 iterations), but here we saw that $Q$ no longer did better than $R$ at any sparsity level for any of the models that we tested. For example, the pruning results for Conv-2 when run with the early stopping criterion is given in Fig. 6, which should be compared with Fig. 4b.
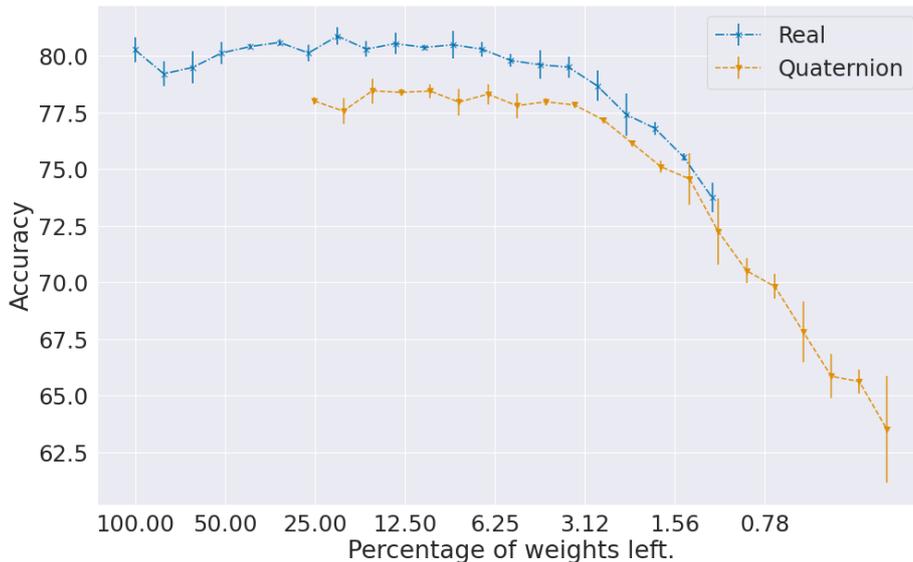


Figure 6: Accuracy vs sparsity results for Conv-2 on CIFAR-10 when using the early stopping criterion.

## 5.2   Larger Models

In addition to the models considered earlier, we also ran similar experiments on Resnet-18 [15] and VGG-16 [37]. These are deeper neural network architectures that require batch normalization layers. The first implementation of a quaternion batch normalization algorithm was given by [10]. This implementation treats a quaternion as a single entity and requires complex matrix operations to calculate the mean and variance of each layer, and is thus extremely computationally intensive (adding a single batchnorm layer to Conv-2 increased the time taken for each trainging iteration by approximately 40 times compared to baseline). Hence we had to opt for another implementation of batch normalization given in [40]. This implementation treats each individual components of quaternions separately, and is hence much faster. This obviously comes at the cost of compromising the relationships between the individual components of quaternions, but since we are already doing this with our use of the split activation function, it may not lead to any further disadvantage.

Neither Resnet nor VGG follow the trend we saw for the smaller models, in that for both of these network architectures $R$ has greater accuracy than $Q$ at all sparsity levels. Whether it is the introduction of the batchnorm layer or the depth of the architectures that is causing this reversal in trend is unclear and needs to be investigated further.

# 6    Conclusions

In this work, we conduct extensive pruning experiments on real and quaternion-valued implementations of different neural network architectures with the objective of checking whether using quaternions provides any advantages in model compression. We first found that pruned quaternion models can be re-trained from scratch to match the original accuracy of the unpruned model, showing that lottery tickets exist for quaternion networks as well. More importantly, our experiments demonstrate that when pruned to high levels of sparsities, quaternion implementations of certain models outperform their complementary real-valued models of equivalent architectures. Hence for ML tasks with multi-dimensional inputs that need to be run on devices with limited computational power, a pruned quaternion model may be a more suitable option than an analogous real network.

## 6.1    Limitations and Future Work

Like all empirical studies, the primary limitation of this work is in the scope of vision tasks, datasets and models tested. Here we tested six architectures on three different datasets at the task of classification. An extension to this work, and one that could further generalize the conclusions reached, is to consider a larger set of models and datasets, and test them on additional vision tasks such as semantic segmentation. Investigating how pruned quaternion implementations of deeper architectures that require batch normalization layers can be made to outperform their real counterparts is also identified as future work.

# References

[1] Paolo Arena et al. "Neural networks for quaternion-valued function approximation". In: *Proceedings of IEEE International Symposium on Circuits and Systems-ISCAS'94*. Vol. 6. IEEE. 1994, pp. 307–310.

[2] Davis Blalock et al. "What is the state of neural network pruning?" In: *arXiv preprint arXiv:2003.03033* (2020).

[3] Yu Cheng et al. "A survey of model compression and acceleration for deep neural networks". In: *arXiv preprint arXiv:1710.09282* (2017).

[4] Danilo Comminiello et al. "Quaternion Convolutional Neural Networks for Detection and Localization of 3D Sound Events". In: *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2019, pp. 8533–8537. DOI: 10.1109/ICASSP.2019.8682711.

[5] Jacob Devlin et al. "Bert: Pre-training of deep bidirectional transformers for language understanding". In: *arXiv preprint arXiv:1810.04805* (2018).

[6] Erich Elsen et al. "Fast sparse convnets". In: *CVPR*. 2020, pp. 14629–14638.

[7] Jonathan Frankle and Michael Carbin. "The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks". In: *ICLR*. 2019.

[8] Jonathan Frankle et al. "Pruning Neural Networks at Initialization: Why Are We Missing the Mark?" In: *International Conference on Learning Representations*. 2021. URL: https://openreview.net/forum?id=Ig-VyQc-MLK.

[9] Jannik Fritsch, Tobias Kuehnl, and Andreas Geiger. "A New Performance Measure and Evaluation Benchmark for Road Detection Algorithms". In: *International Conference on Intelligent Transportation Systems (ITSC)*. 2013.

[10] Chase J. Gaudet and Anthony S. Maida. "Deep Quaternion Networks". In: *Proceedings of the International Joint Conference on Neural Networks* 2018-July (2018). DOI: 10.1109/IJCNN.2018.8489651. arXiv: 1712.04604.

[11] Xavier Glorot and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks". In: *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings. 2010, pp. 249–256.

[12] Scott Gray, Alec Radford, and Diederik P Kingma. "Gpu kernels for block-sparse weights". In: *arXiv preprint arXiv:1711.09224* 3 (2017).

[13] Song Han et al. "Learning Both Weights and Connections for Efficient Neural Networks". In: *NeurIPS*. Montreal, Canada, 2015, pp. 1135–1143.

[14] Babak Hassibi and David G. Stork. "Second Order Derivatives for Network Pruning: Optimal Brain Surgeon". In: *NeurIPS*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1992, pp. 164–171. ISBN: 1558602747.

[15] Kaiming He et al. "Deep residual learning for image recognition". In: *CVPR*. 2016, pp. 770–778.

[16] Andrew G Howard et al. "Mobilenets: Efficient convolutional neural networks for mobile vision applications". In: *arXiv preprint arXiv:1704.04861* (2017).

[17]  Nal Kalchbrenner et al. "Efficient neural audio synthesis". In: *International Conference on Machine Learning*. PMLR. 2018, pp. 2410–2419.

[18]  Alex Krizhevsky and Geoffrey Hinton. *Learning multiple layers of features from tiny images*. 2009.

[19]  Yann LeCun, Corinna Cortes, and CJ Burges. "MNIST handwritten digit database". In: *ATT Labs [Online]. Available: http://yann.lecun.com/exdb/mnist* 2 (2010).

[20]  Yann LeCun, John S Denker, and Sara A Solla. "Optimal brain damage". In: *NeurIPS*. 1990, pp. 598–605.

[21]  Yann LeCun et al. "Gradient-based learning applied to document recognition". In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.

[22]  Namhoon Lee et al. "A signal propagation perspective for pruning neural networks at initialization". In: *arXiv preprint arXiv:1906.06307* (2019).

[23]  Hao Li et al. "Pruning filters for efficient convnets". In: *ICLR*. 2017.

[24]  Jian Li et al. "DSFD: dual shot face detector". In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2019, pp. 5060–5069.

[25]  Zhuang Liu et al. "Learning efficient convolutional networks through network slimming". In: *ICCV*. 2017.

[26]  Tohru Nitta. "A quaternary version of the back-propagation algorithm". In: *Proceedings of ICNN'95-International Conference on Neural Networks*. Vol. 5. IEEE. 1995, pp. 2753–2756.

[27]  Titouan Parcollet, Mohamed Morchid, and Georges Linarès. "A survey of quaternion neural networks". In: *Artificial Intelligence Review* 53.4 (2020), pp. 2957–2982. ISSN: 15737462. DOI: 10.1007/s10462-019-09752-1. URL: https://doi.org/10.1007/s10462-019-09752-1.

[28]  Titouan Parcollet et al. "Quaternion recurrent neural networks". In: *arXiv preprint arXiv:1806.04418* (2018).

[29]  Titouan Parcollet et al. "Speech recognition with quaternion neural networks". In: *arXiv preprint arXiv:1811.09678* (2018).

[30]  Jongsoo Park et al. "Faster cnns with direct sparse convolutions and guided pruning". In: *ICLR*. 2017.

[31]  Adam Paszke et al. "Pytorch: An imperative style, high-performance deep learning library". In: *NeurIPS*. Vol. 32. 2019, pp. 8026–8037.

[32]  Dario Pavllo et al. "Modeling human motion with quaternion-based neural networks". In: *International Journal of Computer Vision* 128.4 (2020), pp. 855–872.

[33]  Jeff Pool, Abhishek Sawarkar, and Jay Rodge. *Accelerating Inference with Sparsity Using the NVIDIA Ampere Architecture and NVIDIA TensorRT*. https://developer.nvidia.com/blog/accelerating-inference-with-sparsity-using-ampere-and-tensorrt/. Technical Walkthrough. 2021.

[34]  Russell Reed. "Pruning algorithms-a survey". In: *IEEE transactions on Neural Networks* 4.5 (1993), pp. 740–747.

[35] Alex Renda, Jonathan Frankle, and Michael Carbin. "Comparing rewinding and fine-tuning in neural network pruning". In: *ICLR*. 2020.

[36] Karen Simonyan and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition". In: *arXiv preprint arXiv:1409.1556* (2014).

[37] Karen Simonyan and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition". In: *arXiv preprint arXiv:1409.1556* (2014).

[38] Chiheb Trabelsi et al. "Deep complex networks". In: *ICLR*. 2018, pp. 1–19. arXiv: 1705.09792.

[39] Thanh Tran et al. "HABERTOR: An efficient and effective deep hatespeech detector". In: *arXiv preprint arXiv:2010.08865* (2020).

[40] Qilin Yin et al. "Quaternion Convolutional Neural Network for Color Image Classification and Forensics". In: *IEEE Access* 7 (2019), pp. 20293–20301. DOI: 10.1109/ACCESS.2019.2897000.

[41] Giorgio Zannini et al. *ispamm/hTorch*. https://github.com/ispamm/hTorch. Github Repository. 2021.

[42] Michael Zhu and Suyog Gupta. "To prune, or not to prune: exploring the efficacy of pruning for model compression". In: *arXiv preprint arXiv:1710.01878* (2017).

[43] Xuanyu Zhu et al. "Quaternion convolutional neural networks". In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 11212 LNCS (2018), pp. 645–661. ISSN: 16113349. DOI: 10.1007/978-3-030-01237-3_39. arXiv: 1903.00658.