
Benchmark of neural network models with reduced number of parameters for image classification

Rares Folea

Department of Computer Science
Politehnica University of Bucharest
Bucharest, Romania
rares.folea@stud.acs.upb.ro

Abstract

This paper analyze the performance of various **deep neural networks**, such as **residual networks**, **VGGs** or **Vision Transformers** against the problem of image classification problems. The analyzed networks are reduced versions of well-known deep neural network architectures that have under 1M trainable parameters. We propose architectural design changes to the networks, as well as looking into pruning techniques on the well-established models. As of 2022, networks under 1M trainable parameters are considered to be composed relatively of a number of **reduced parameters**. We analyze, as well, the influence of various optimizers based on **Stochastic Gradient Descent** algorithms. Most of the evaluations for this project are being performed on the **CIFAR10** dataset, but we present some additional results for the **CIFAR100** dataset.

1 Introduction

Image recognition is a common method for allowing computers to recognize items inside a picture. Image labeling, image search, medical condition recognition, and self-driving cars are just a few of the jobs that use such technologies.[7] The **CIFAR-10** dataset is a set of photographs used to train machine learning and computer vision techniques. This is one of the most used databases for deep neural network studies.

State-of-the-Art models reach phenomenal classification performances on the **CIFAR-10** dataset, with even over 99% accuracy on the task of prediction. However, as outlined in **Table 1**, these models can have tens or hundreds of million parameters. As of Late 2021, the highest percentage of correct classifications is achieved by the **ViT-H/14** model, reaching **99.5%** accuracy. Yet, this model has **632M parameters**. Also other top performing architecture (over 95% accuracy), such as models based on **VGG-19** or **ResNet50** have **20M**, respectively **25M parameters**.

In this paper, we investigate neural network models with a reduced number of parameters, against the task of image classification of on the **CIFAR-10** and **CIFAR-100** dataset.

By *reduced number of parameters*, we want to investigate neural networks with **at most one million number of parameters**. While **1M** parameters can sound completely not a small value for parameters, we will see that state-of-the-art networks can have up to **632M** parameters.

As stated in the **lottery ticket hypothesis**, it is believed that for any given *sufficiently* dense neural network, there exists a random sub-network, within that network, that has *one tenth of parameters*, but would **converge to the same performance in the same amount of training time** [4].



Figure 1: The 10 classes in the CIFAR-10 dataset, with some sample 32x32 images from each class

32 We wanted to analyze the networks are reduced versions of well known deep neural network architectures
 33 that have under 1M trainable parameters, by either proposing architectural design changes to
 34 the networks, as well as looking into pruning techniques on the well-established models.

35 2 Dataset

36 **CIFAR-10** is a computer-vision dataset used for benchmarks for the task of object recognition. The
 37 **CIFAR-10** dataset consists of 60000 32x32 colour images in **10 classes**, with 6000 images per class.
 38 There are *50000 training images and 10000 test images*. [8]

39 Throughout our research, we will benchmark all the neural networks against the **CIFAR-10** dataset.
 40 This ensures an objective comparison between the models and optimizers.

41 The **CIFAR-10** dataset contains *60 thousand 32x32* color images in *10 different classes*. The 10
 42 different classes represent **airplanes, cars, birds, cats, deer, dogs, frogs, horses, ships, and trucks**.
 43 There are **6 thousand** images of each class. All the classes in this dataset are mutually exclusive,
 44 with no overlap between the classes. The total size of the dataset is **163 MB**.

45 For architecture that perform similar, we will use the **CIFAR-100** dataset as a tie-break. Introduced
 46 by the Canadian Institute for Advanced Research, this 100 classes dataset is a subset of the Tiny
 47 Images dataset and consists of 60000 32x32 color images. The 100 classes in the **CIFAR-100** are
 48 grouped into 20 super classes, with 600 images per class. Each image comes with a "fine" label (the
 49 class to which it belongs) and a "coarse" label (the super-class to which it belongs). There are 500
 50 training images and 100 testing images per class [8].

51 3 Related Work

52 3.1 State-of-the-art Models for image classification on CIFAR-10

53 Over the time, the **CIFAR-10** dataset become a benchmark for the task of image classification for
 54 various models.

55 As of Late 2021, the highest percentage of correct classifications is achieved by the **ViT-H/14** model.
 56 Over 10000 training images, the model has achieved a 99.5% accuracy on predicting the class of
 57 the image, meaning less than 50 incorrect prediction over the 10k testing dataset. This has been
 58 initially introduced in 2020 [3], in the paper *"An Image is Worth 16x16 Words: Transformers for*

Model	Number of parameters	Percentage correct
ViT-H/14 (2020)	632M	99.50
LaNet (2019)	44.1M	99.03
DenseNet-BC-190 + Mixup (2017)	25.6M	97.3
VGG-19 with GradInit (2021)	20.03M	94.71

Table 1: The above table presents some important models developed over time for the Image Classification challenge on the CIFAR-10 dataset.

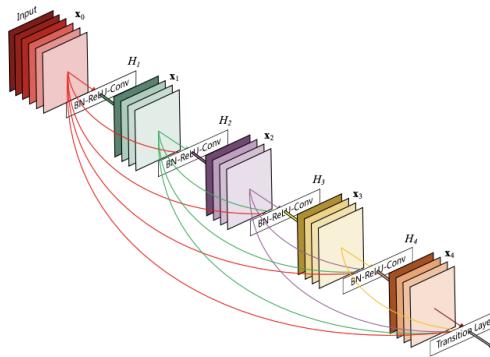


Figure 2: A sample DenseNet architecture, a model that utilises dense connections between layers, through "Dense Blocks", layers where all layers are connected directly with each other.

- 59 *Image Recognition at Scale*". This introduced a Transformer architecture for solving the task of image
60 classification. The model uses attention in conjunction with convolutional networks and to replace
61 certain components of convolutional networks. [3] shows that Convolutional Neural Networks are
62 not required for small image classification tasks, and that a pure transformer applied straight to
63 sequences of picture patches can perform well.
- 64 Another notable model from 2019, **LaNet** [13] was introduced by the "*Sample-Efficient Neural*
65 *Architecture Search by Learning Action Space for Monte Carlo Tree Search*" paper. Such a architecture
66 learns to partition the search space, so that solvers can focus on a smaller region to find better solutions
67 with fewer samples.
- 68 The **DenseNet-BC-190 + Mixup** model [14] introduced in 2017. The model uses at core a **DenseNet**
69 architecture, a type of convolutional neural network that utilises dense connections between layers.
- 70 The paper [14] presents the drawbacks of Large Deep neural networks such as memory and sensitivity
71 to hostile samples. Mixup is a learning principle, which is proposed in this paper as a solution to these
72 problems. What Mixup does essentially is that it trains a neural network using convex combinations
73 of pairs of instances and their labels to train a neural network, allowing it to favor simple linear
74 behavior between training examples.
- 75 Finally, **VGG-19 with GradInit** [15] is another model that performs well on the CIFAR-10 dataset.
- 76 **VGG-19** is a Convolutional Neural Network. There are 19 stands for the number of layers with
77 trainable weights, 16 Convolutional layers and 3 Fully Connected layers.
- 78 **GradInit** [15], an automated and architecture agnostic method for initializing neural networks.
79 GradInit is based on a simple heuristic; the norm of each network layer is adjusted so that a single
80 step of *SGD* or *Adam* with prescribed hyper parameters results in the smallest possible loss value.

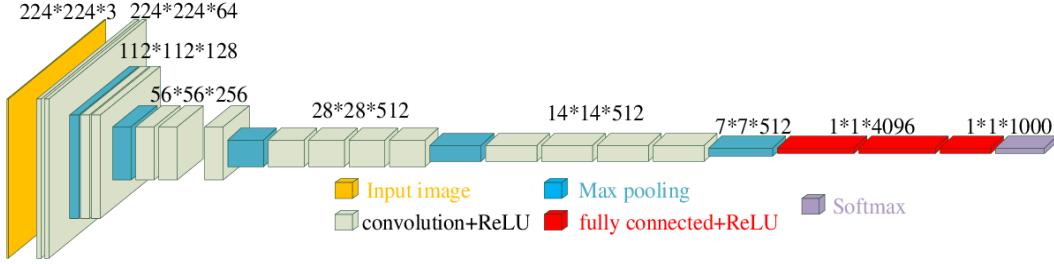


Figure 3: An overall architecture for a VGG-19 network.

81 However, as outlined in **Table 1**, these models have a high numbers of parameter, making the process
 82 of training a significant inconvenient, requiring dedicated hardware and resulting in high energy
 83 commitments.

84 3.2 Architectures of Deep Neural Networks with a reduced number of parameters

85 The objective of this paper is to benchmark the image classification of neural network models with a
 86 reduced number of parameters. In the next subsection, we will analyze the state-of-the-art models for
 87 the classification problem on the **CIFAR-10** dataset. For the architectures that perform really well
 88 (over 90% accuracy) we will also provide some results on the **CIFAR-100** dataset.

89 3.2.1 Reduced-ResNet architecture

90 A residual neural network (known as ResNet [5]) is a type of neural network. Skip connections, or
 91 shortcuts, are used by residual neural networks to jump past some layers. The majority of ResNet
 92 models use double or triple layer skips with non-linearities (such as Rectified Linear Unit) and batch
 93 normalization in between.

94 Residual networks are highly configurable and usually the more deeper the ResNet is, the higher ac-
 95 curacy is obtained. Certain ResNet architectures have been wildly studied and proved to provide good
 96 results: **ResNet18**, **ResNet34**, **ResNet50**, **ResNet101**, **ResNet152**. ResNet18 has 18 convolutional
 97 layers, while ResNet152 has 152 convolutional layers.

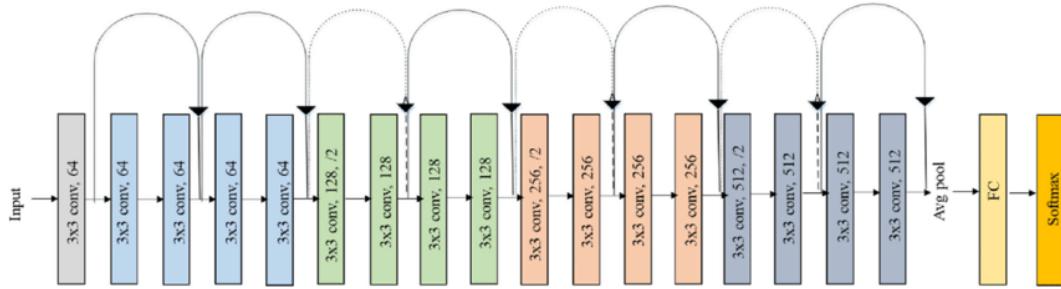


Figure 4: The classical architecture for the Deep Convolutional Network for ResNet18.

98 However, the number of parameters of these networks is quite high. While ResNet18 has around 11
 99 million trainable parameters, ResNet34 has 21M. More deep models such as ResNet50, ResNet101
 100 or ResNet152 have 23M, 42M, respectively 58M parameters.

101 In order to fit the limit of 1M parameters, we have decided to use in benchmarking a different ResNet
 102 architecture, based on ResNet18 but with less output channels for the first 2D convolution. Therefore,

103 instead of outputting the 64 channels, our **ResNet18_Small** (also referred as **ResNet18_S**) neural
 104 network will only output on 16 channels.

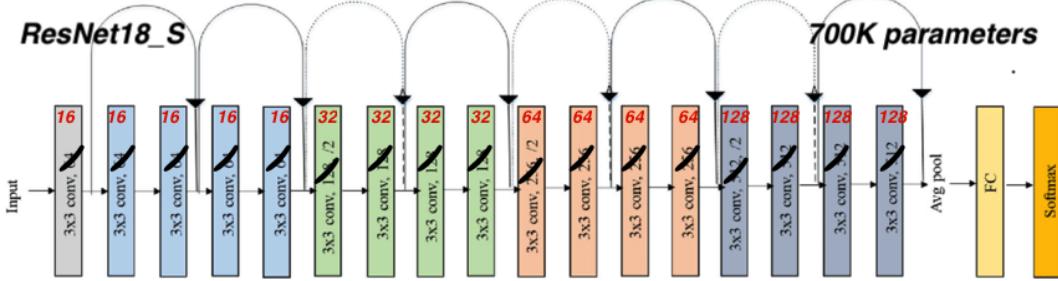


Figure 5: The benchmark-ed ResNet18_S, a slightly modified ResNet18 with less than 1M parameters.

105 The resulting network ResNet18_S has a total of 700K trainable parameters, as opposed to the 11M
 106 trainable parameters of ResNet18, making it a good candidate for our research.

107 3.2.2 Reduced VGG-16 architecture

108 **VGG** are a type of Very Deep Convolutional Networks, initially introduced in [12] are networks of
 109 increasing depth using an architecture with very small (3×3) convolution filters, which shows that
 110 a significant improvement on the prior-art configurations can be achieved by pushing the depth to
 111 16-19 weight layers.

112 In particular, **VGG16** (16 depth-layers) is a convolutional neural network architecture which instead
 113 of having a large number of hyper-parameter, it is focused on having convolution layers of 3×3 filter
 114 with a stride 1 and always used same padding and max-pool layer of a 2×2 filter with the stride 2.
 115 It follows this arrangement of convolution and max pool layers consistently throughout the whole
 116 architecture. In the end it has 2 fully-connected layers followed by a soft-max layer for output. [12].
 117 A classic **VGG16** architecture for **CIFAR-10** has **15245130** trainable parameters.

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			$7 \times 7, 64$, stride 2		
conv2_x	56×56	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1	average pool, 1000-d fc, softmax				
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9

Figure 6: Sample ResNet [5] architectures. For each network (**ResNet18**, **ResNet34**, **ResNet50**, **ResNet101**, **ResNet152**) the main convolutions blocks are similar, but replicated a different number of times, therefore a different in the depth of the convolutional network.



Figure 7: Sample images from the CIFAR-10 dataset.

118 3.2.3 Reduced Visual Transformers(ViT) architecture

119 A Vision Transformer (ViT) is a transformer designed for image recognition and other vision-
 120 processing applications. Usually, transformers have way more parameters than already classical
 121 convolutional deep neural network architectures such as ResNet18 or VGG-16. The SotA ViT-
 122 H/14 transformer that obtained 99.5% accuracy on the CIFAR-10 dataset had over 600M trainable
 123 parameters.

124 3.3 Optimizers

125 Optimizers are algorithms or methods used to change the attributes of your neural network such as
 126 weights and learning rate in order to reduce the losses.

127 We are aiming to benchmark the following suite of algorithms for Stochastic Gradient Descent
 128 Optimization Algorithms for the image classification problems based on **CIFAR10** dataset.

- 129 1. **Stochastic Gradient Descent with Momentum** [9] (SGDM)
- 130 2. **Adam**: A method for stochastic optimization [6] (ADAM)
- 131 3. **AdaBound**: Adaptive Gradient Methods with Dynamic Bound of Learning Rate [10]

132 At each iteration, *stochastic gradient descent* samples a selection of summand functions in order
 133 to reduce the computational cost. In the case of large-scale machine learning problems, this is
 134 particularly effective.

135 **Stochastic gradient descent with momentum** remembers the update of deltas at each iteration, and
 136 determines the next update as a linear combination of the gradient and the previous update. **Adam** [6]
 137 is also an algorithm for first-order gradient-based optimization of stochastic objective functions,
 138 based on adaptive estimates of lower-order moments.

139 [11] presents **AdaBound** as an optimizer that behaves like Adam at the beginning of training, and
 140 gradually transforms to SGD at the end, therefore considered by the author an optimizer that trains as
 141 fast as Adam and as good as SGD.

142 3.4 Pruning techniques

143 Deep learning approaches that are now in use rely on over-parametrized models that are difficult to
 144 implement. Identifying the best strategies for compressing models by minimizing the number of
 145 parameters in them is critical in order to save memory, power, and hardware without sacrificing accu-
 146 racy, deploy lightweight models on devices, and ensure privacy using private on-device computing.
 147 Pruning is used to research the function of fortunate sparse subnetworks and initializations (known as
 148 lottery tickets) as a destructive neural architecture search tool, as well as to analyze the variations in
 149 learning dynamics between over-parametrized and under-parametrized networks [2].

150 Model pruning is a strategy for shrinking the size of a deep learning model by discovering and
 151 removing tiny weights from the model. Model pruning may significantly reduce model size while
 152 also speeding up inference time.

153 Weights are removed on a case-by-case basis in unstructured pruning methods. Weights are removed
 154 in groups using structured pruning methods, such as removing entire channels at a time. Structured
 155 pruning offers better runtime speed (due to the compact computation on fewer channels), but it has a
 156 greater influence on model correctness (due to the fact that it is less selective) [1].
 157 In the global unstructured pruning, instead of specifying individual pruning parameters for individual
 158 layers, we apply to the model all at once [2].
 159 Pruning can be added as an extra step in between training epochs (*iterative pruning*) or done all at
 160 once once model training is complete (*one-shot pruning*).

161 4 Experiment Setup

162  The implementation on which this benchmark of neural network models with reduced number of
 163 parameters the task of image classification has been done is available open-source:

164 <https://github.com/raresraf/benchmark-dnn-small>

165 We used pyTorch, an open source machine learning framework for the python programming language
 166 to implement and test our models. Also, we have used WandB, a central dashboard to keep track of
 167 your hyper parameters, system metrics, and predictions of our models.

168 5 Results

169 This section will present the results of our research, for the ResNet, VGG and ViT architectures.

170 5.1 ResNet spotlight

171 We want to dedicate one entire subsection for analysing some experimental results where we tested a
 172 classic **ResNet18** architecture, as well as a lightweight version, as we have described it as **ResNet18_S**
 173 in **Appendix A**, with different optimizers.

174 **ResNet18** obtained **94.13 %** accuracy while testing on the CIFAR-10 dataset when mixed with
 175 AdaBound optimizer. Interesting, the second best result was obtained by the **ResNet18_S** network
 176 with AdaBound optimized, outclassing the **ResNet18** architecture performance when combined with
 177 SGDM or ADAM optimizers. Also, it is worth noticing that the **ResNet18_S** network with ADAM
 178 optimizer performed almost as good as the **ResNet18** SGDM.

Architecture	Optimizer	Test accuracy
ResNet18	<i>SGDM</i>	89.4
	<i>ADAM</i>	89.97
	<i>AdaBound</i>	94.13
ResNet18_S	<i>SGDM</i>	86.85
	<i>ADAM</i>	89.27
	<i>AdaBound</i>	91.72

Table 2: The test accuracy, on the **CIFAR-10** dataset, for multiple configurations of the ResNet networks architectures and used optimizers.

179 For the ResNet18 and ResNet18_S architectures confusion matrix results on the CIFAR-10 dataset
 180 please refer to Appendix D.

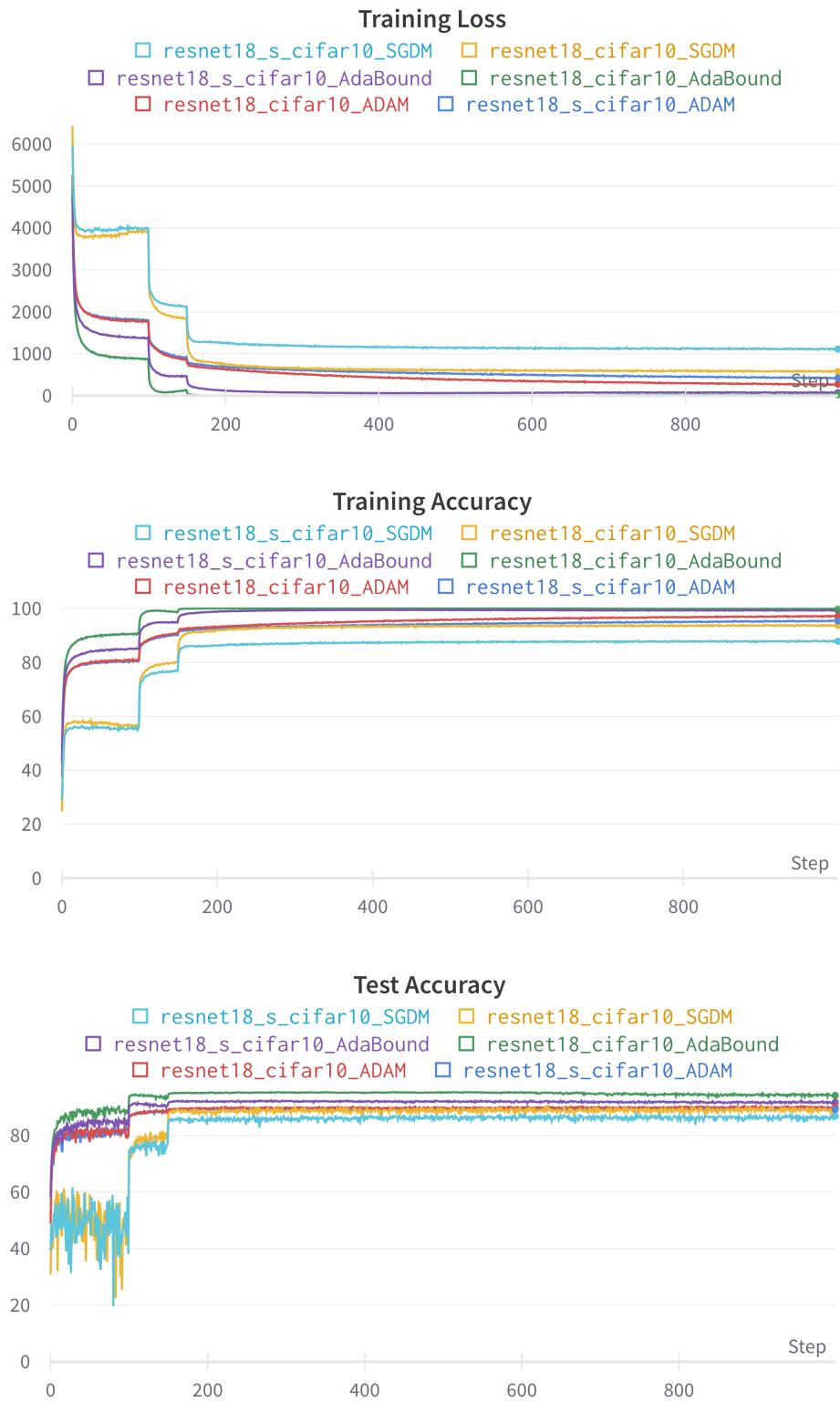


Figure 8: Overall view on the training loss error, training accuracy and testing accuracy obtained on the CIFAR-10 dataset obtained when using different optimizers, such as **SGDM**, **ADAM** and **AdaBound**, during 1000-epoch training of ResNet18/ResNet18_S networks.

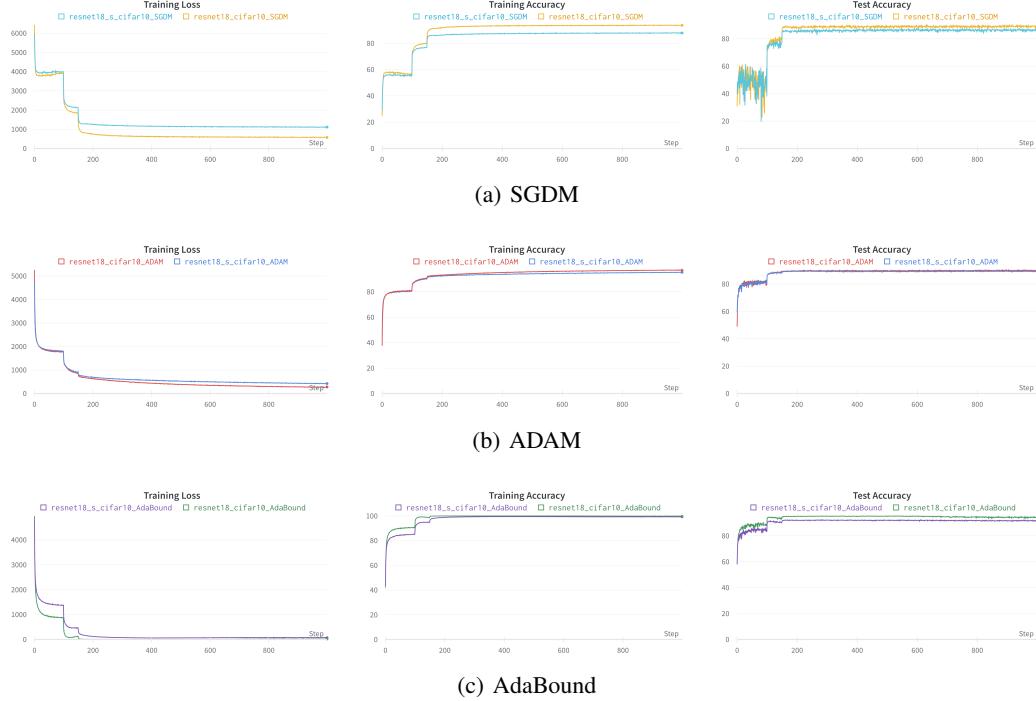


Figure 9: The training loss error, training accuracy and testing accuracy obtained on the CIFAR-10 dataset obtained when using the optimizer: (a) **SGDM** (b) **ADAM** (c) **AdaBound**, during 1000-epoch training of ResNet18/ResNet18_S networks.

181 We ran a similar scenario for the CIFAR-100 dataset, training the same networks on a 100-epoch
 182 scenario with various optimizers. The results showed that best accuracy (58%) has been obtained
 183 while using a ResNet18 with AdaBound optimizer, followed with a 4% decrease in quality by the
 184 ResNet18_S architecture, with the same optimizer. Interesting enough, the latter performed better
 185 than the traditional ResNet18 architecture while using ADAM optimizer, with more than 4.2%.

Architecture	Optimizer	Test accuracy
ResNet18	<i>SGDM</i>	1
	<i>ADAM</i>	39.04
	<i>AdaBound</i>	58
ResNet18_S	<i>SGDM</i>	1
	<i>ADAM</i>	43.31
	<i>AdaBound</i>	54.05

Table 3: The test accuracy, on the **CIFAR-100** dataset, for multiple configurations of the ResNet networks architectures and used optimizers.

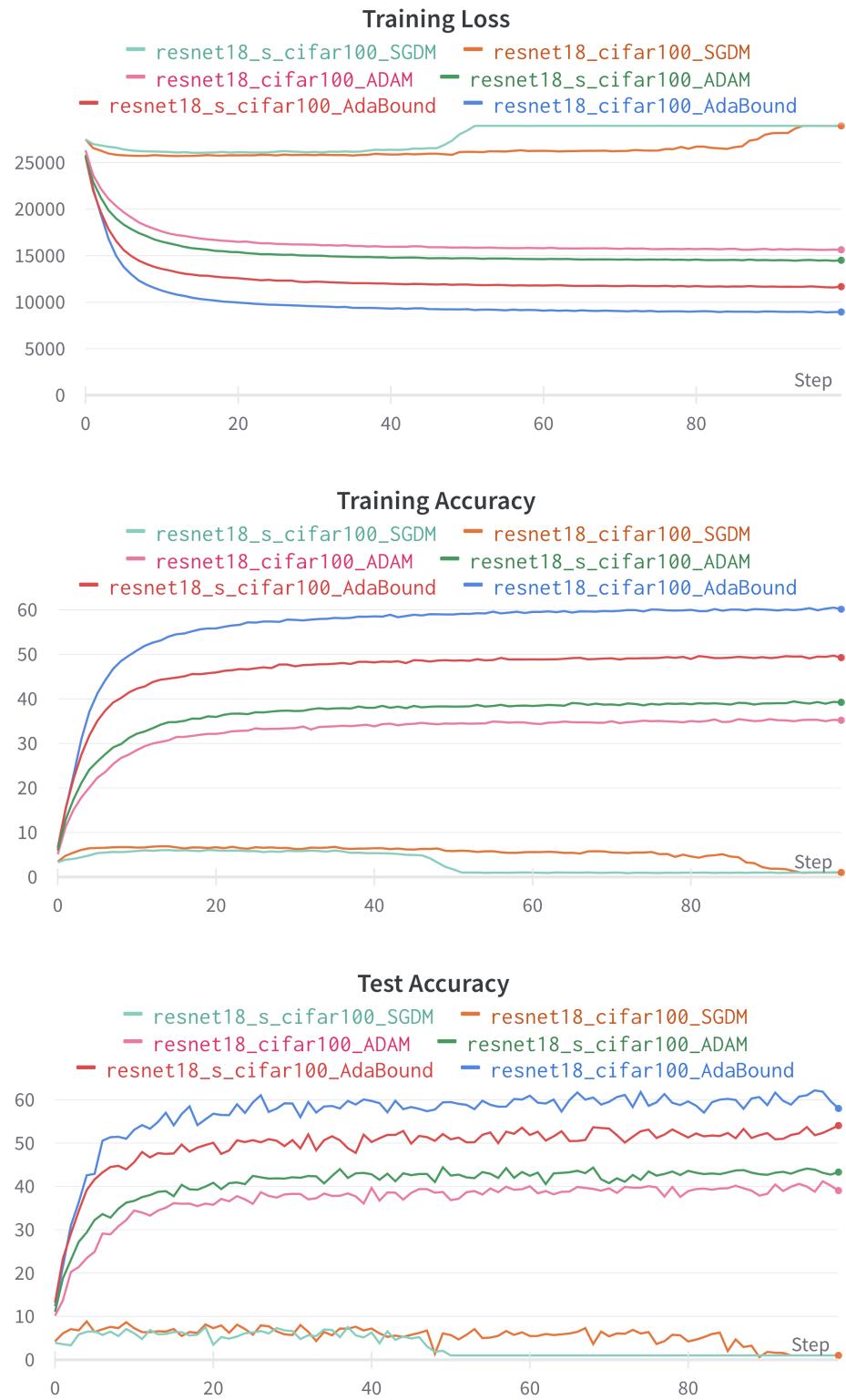


Figure 10: Overall view on the training loss error, training accuracy and testing accuracy obtained on the CIFAR-100 dataset obtained when using different optimizers, such as **SGDM**, **ADAM** and **AdaBound**, during 100-epoch training of ResNet18/ResNet18_S networks.

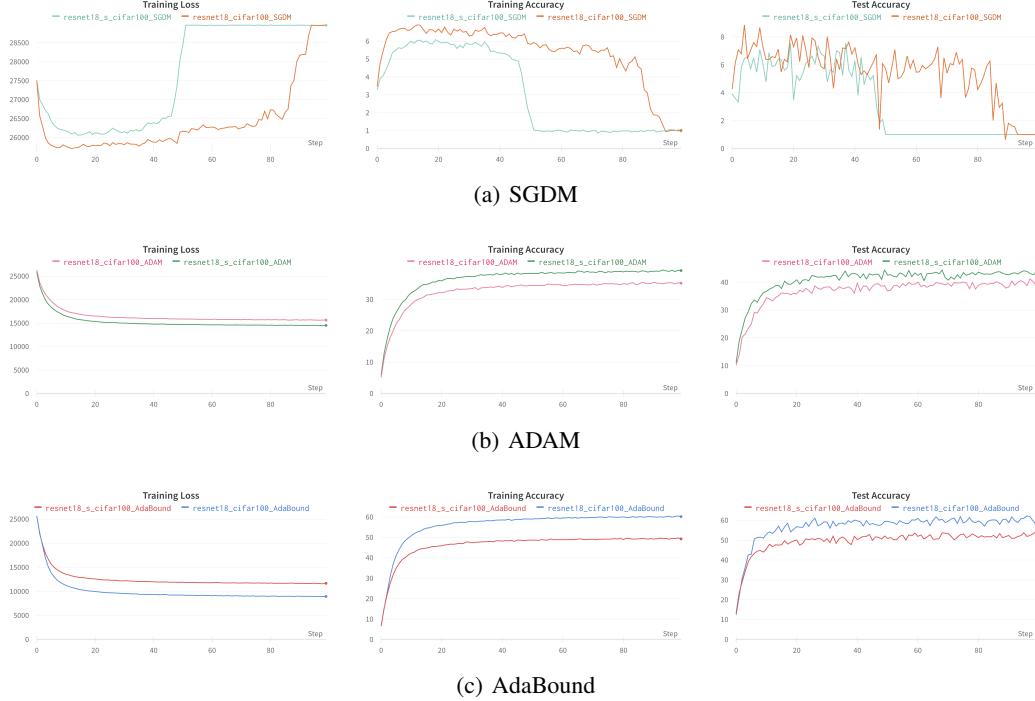


Figure 11: The training loss error, training accuracy and testing accuracy obtained on the CIFAR-100 dataset obtained when using the optimizer: (a) **SGDM** (b) **ADAM** (c) **AdaBound**, during 100-epoch training of ResNet18/ResNet18_S networks.

186 5.2 VGG

187 The results for the VGG architectures were inconclusive. So far, we were not able to induce the
 188 VGG-16 architecture into a training mode that would make meaningful learning progress.

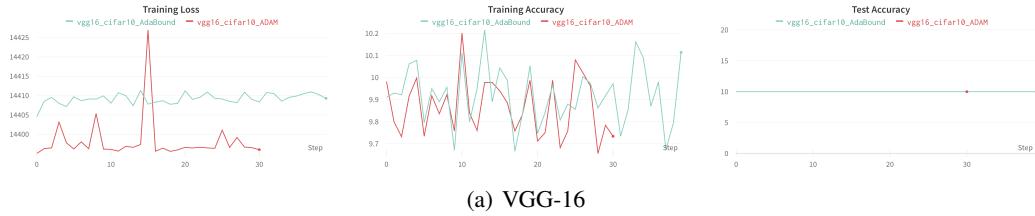


Figure 12: The training loss error, training accuracy and testing accuracy obtained on the CIFAR-10 dataset for the VGG16 architecture during the preliminary training.

189 5.3 Vision transformers

190 The Vision Transformer is a model for image classification that employs a Transformer-like architec-
 191 ture over patches of the image.

Architecture	Optimizer	Test accuracy
ViT	<i>SGDM</i>	10
	<i>ADAM</i>	51.49
	<i>AdaBound</i>	48.38
ViT_S	<i>SGDM</i>	10
	<i>ADAM</i>	50.99
	<i>AdaBound</i>	49.44

Table 4: The test accuracy on the **CIFAR-10** dataset, for multiple configurations of network architectures and used optimizers for the ViT architecture. The SotA ViT model reaches over 99% accuracy on the **CIFAR-10** dataset.

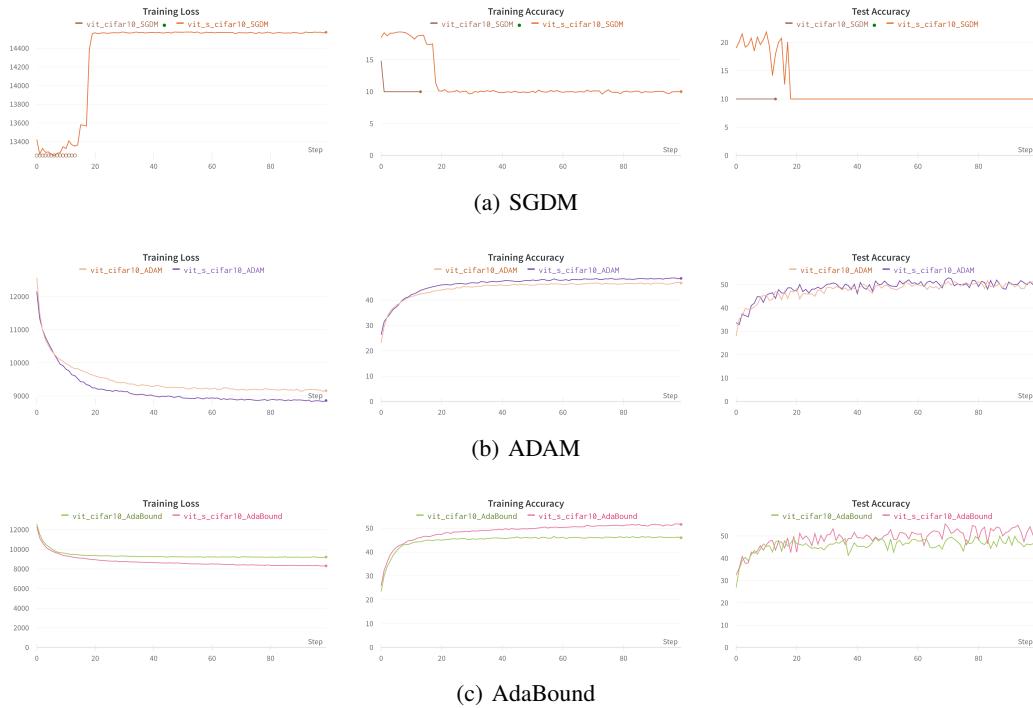


Figure 13: The training loss error, training accuracy and testing accuracy obtained on the CIFAR-10 dataset obtained when using the optimizer: (a) **SGDM** (b) **ADAM** (c) **AdaBound**, during 1000-epoch training of ResNet18/ResNet18_S networks.

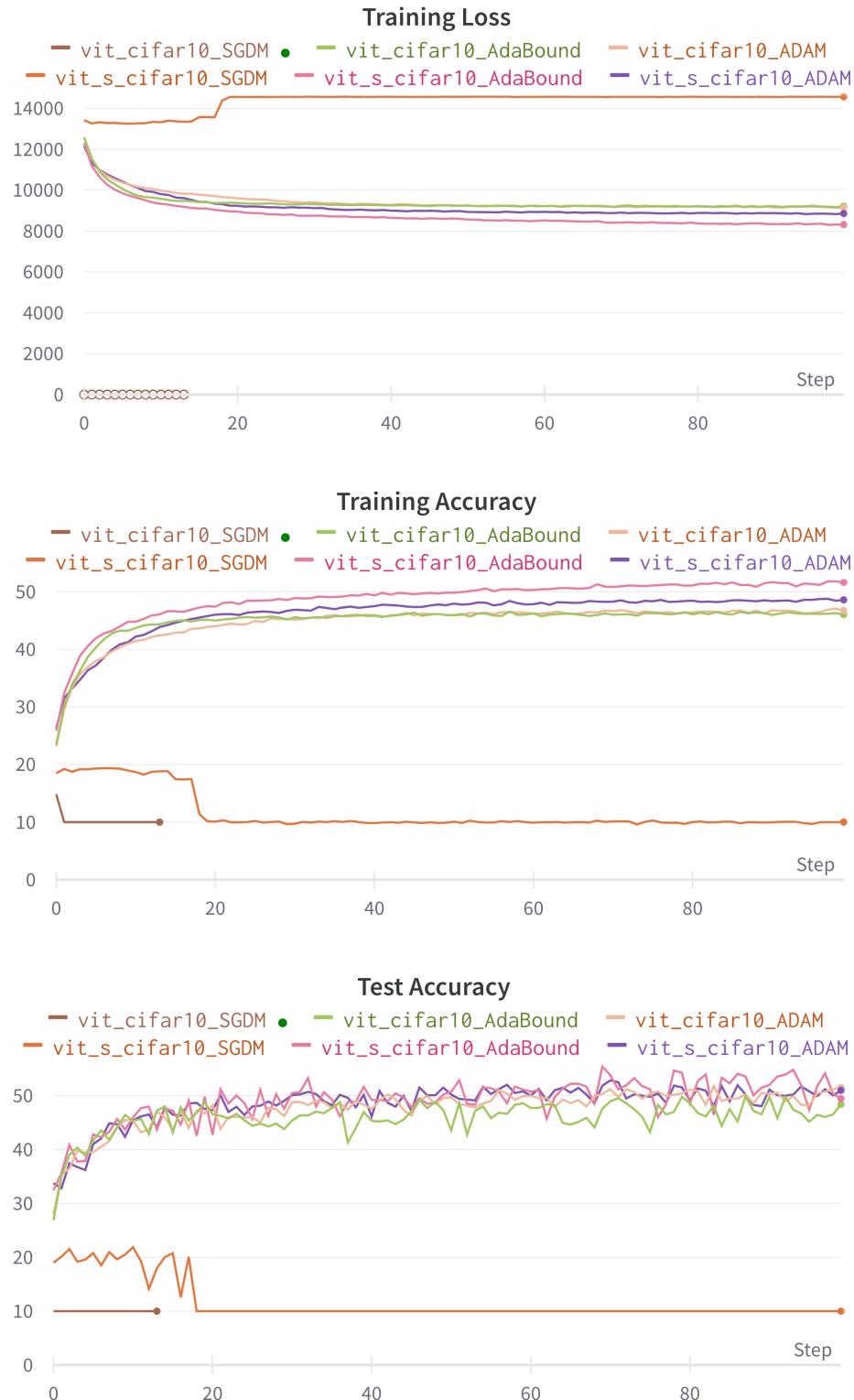


Figure 14: Overall view on the training loss error, training accuracy and testing accuracy obtained on the CIFAR-10 dataset when using different optimizers, such as **SGDM**, **ADAM** and **AdaBound**, during 100-epoch training of ViT/ViT_S networks.

192 We ran a similar scenario for the CIFAR-100 dataset, training the same vision transformers on a
 193 100-epoch scenario with various optimizers.

Architecture	Optimizer	Test accuracy
ViT	<i>SGDM</i>	1
	<i>ADAM</i>	22.97
	<i>AdaBound</i>	25.05
ViT_S	<i>SGDM</i>	3.86
	<i>ADAM</i>	23.56
	<i>AdaBound</i>	21.56

Table 5: The test accuracy, on the **CIFAR-100** dataset, for multiple configurations of the ResNet networks architectures and used optimizers.

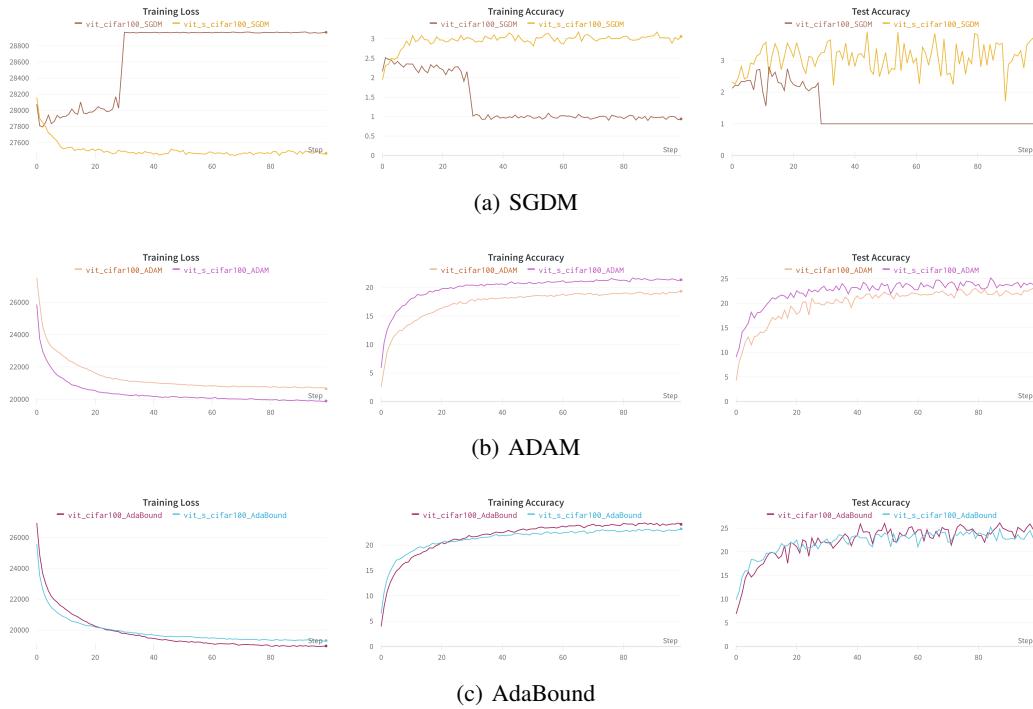


Figure 16: The training loss error, training accuracy and testing accuracy obtained on the CIFAR-100 dataset obtained when using the optimizer: (a) **SGDM** (b) **ADAM** (c) **AdaBound**, during 100-epoch training of ResNet18/ResNet18_S networks.

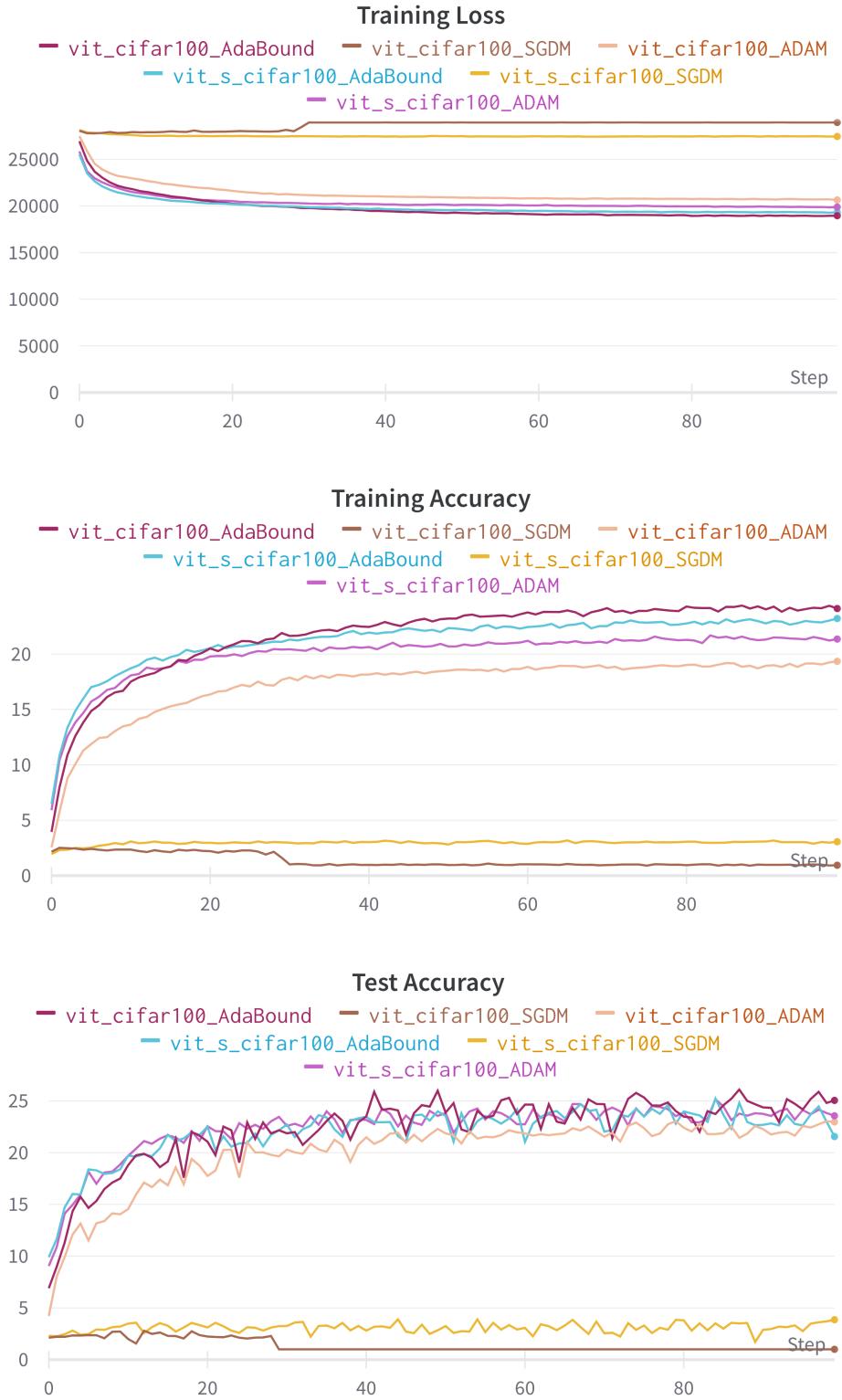


Figure 15: Overall view on the training loss error, training accuracy and testing accuracy obtained on the CIFAR-100 dataset obtained when using different optimizers, such as **SGDM**, **ADAM** and **AdaBound**, during 100-epoch training of the visual transformers networks.

194 **5.4 Experimenting with pruning techniques for the ResNet architecture**

195 Using the trained ResNet18 network on the CIFAR-10 dataset, we applied various pruning techniques,
196 such as random pruning, 11 structured and unstructured pruning and global pruning. The detailed
197 results can be referred in **Appendix E**.

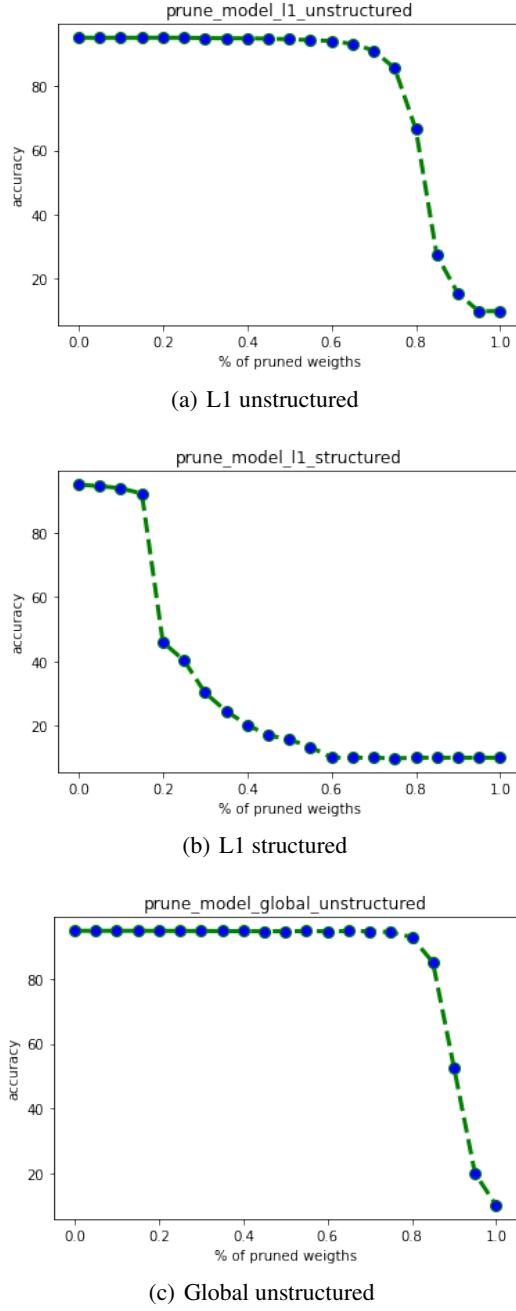


Figure 17: Overall view on the testing accuracy obtained on the CIFAR-10 dataset obtained when using different pruning techniques on the ResNet18 network.

198 **6 Conclusion**

199 Our research has investigated whether we can use traditional deep neural architectures and achieve
200 good performance with a reduced number of trainable parameters. For this, we have used both
201 architectural changes and pruning techniques.

202 The optimizers play a key role in obtaining the desired high-accuracy for the deep neural networks
203 analyzed. A network with tens of millions of trainable parameters may not perform any better than a
204 smaller replica, with fewer trainable parameters but a different optimizer.

205 **ResNet18** obtained **94.13 %** accuracy on the CIFAR-10 dataset when mixed with AdaBound. Interesting,
206 the second best result was obtained by the **ResNet18_S** network with AdaBound optimized,
207 outclassing the **ResNet18** architecture performance when combined with SGDM or ADAM optimizers.
208 With less than 3% accuracy decrease, the ResNet18_S model managed to classify correctly the
209 class of the image in over than 91% with the Adabound optimizer, while only using 700k parameters
210 (down from over 10M, as the original ResNet18 architecture was initially proposed with).

211 We ran a similar scenario for the CIFAR-100 dataset, training the same networks on a 100-epoch
212 scenario with various optimizers. The results showed that best accuracy (58%) has been obtained
213 while using a ResNet18 with AdaBound optimizer, followed with a 4% decrease in quality by the
214 ResNet18_S architecture, with the same optimizer. Interesting enough, the latter performed better
215 than the traditional ResNet18 architecture while using ADAM optimizer, with more than 4.2%.

216 The vision transformers that we have analyzed were relatively small wrt. the total number of
217 parameters as opposed to the SotA vision transformers that obtain benchmark results on the two
218 datasets. As a result, the top performance for these transformers on the CIFAR-10 dataset was just
219 shy over 50%, effectively being outclassed by the traditional ResNet architectures.

220 The pruning techniques proved also to be useful, and by using L1 local and global unstructured
221 pruning, we obtained high accuracy models with less than 50% of the parameters, with only 2% loss
222 in accuracy by only using 20% of the parameters. However, the target of pruning a ResNet18 network
223 to obtain a good accuracy with less than 1M parameters did not succeed.

224 The target of making architectural changes to a ResNet18 network to obtain a good accuracy with
225 less than 1M parameters did succeed, as the ResNet18_S model obtained an accuracy almost as good
226 as the well known ResNet18.

227 **References**

- 228 [1] A developer-friendly guide to model pruning in pytorch, Dec 2020.
- 229 [2] Pruning tutorial, 2020.
- 230 [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
231 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.
232 An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 234 [4] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable
235 neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- 236 [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
237 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
238 pages 770–778, 2016.
- 239 [6] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- 241 [7] Kamil Kłosowski. Image recognition on cifar10 dataset using resnet18 and keras. 2018.

- 242 [8] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
243 2009.
- 244 [9] Yanli Liu, Yuan Gao, and Wotao Yin. An improved analysis of stochastic gradient descent with
245 momentum. *arXiv preprint arXiv:2007.07989*, 2020.
- 246 [10] Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with dynamic
247 bound of learning rate. *arXiv preprint arXiv:1902.09843*, 2019.
- 248 [11] Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with
249 dynamic bound of learning rate. In *Proceedings of the 7th International Conference on Learning
250 Representations*, New Orleans, Louisiana, May 2019.
- 251 [12] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
252 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 253 [13] Linnan Wang, Saining Xie, Teng Li, Rodrigo Fonseca, and Yuandong Tian. Sample-efficient
254 neural architecture search by learning actions for monte carlo tree search. *IEEE Transactions
255 on Pattern Analysis and Machine Intelligence*, 2021.
- 256 [14] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond
257 empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- 258 [15] Chen Zhu, Renkun Ni, Zheng Xu, Kezhi Kong, W Ronny Huang, and Tom Goldstein. Gra-
259 dinit: Learning to initialize neural networks for stable and efficient training. *arXiv preprint
260 arXiv:2102.08098*, 2021.

261 **A ResNet18_S architecture**

262 The architecture of the **ResNet18_S** deep neural network used for the benchmark is as follows:

263	264 Layer (type)	265 Output Shape	Param #
266	Conv2d-1	[-1, 16, 32, 32]	432
267	BatchNorm2d-2	[-1, 16, 32, 32]	32
268	Conv2d-3	[-1, 16, 32, 32]	2,304
269	BatchNorm2d-4	[-1, 16, 32, 32]	32
270	Conv2d-5	[-1, 16, 32, 32]	2,304
271	BatchNorm2d-6	[-1, 16, 32, 32]	32
272	BasicBlock-7	[-1, 16, 32, 32]	0
273	Conv2d-8	[-1, 16, 32, 32]	2,304
274	BatchNorm2d-9	[-1, 16, 32, 32]	32
275	Conv2d-10	[-1, 16, 32, 32]	2,304
276	BatchNorm2d-11	[-1, 16, 32, 32]	32
277	BasicBlock-12	[-1, 16, 32, 32]	0
278	Conv2d-13	[-1, 32, 16, 16]	4,608
279	BatchNorm2d-14	[-1, 32, 16, 16]	64
280	Conv2d-15	[-1, 32, 16, 16]	9,216
281	BatchNorm2d-16	[-1, 32, 16, 16]	64
282	Conv2d-17	[-1, 32, 16, 16]	512
283	BatchNorm2d-18	[-1, 32, 16, 16]	64
284	BasicBlock-19	[-1, 32, 16, 16]	0
285	Conv2d-20	[-1, 32, 16, 16]	9,216
286	BatchNorm2d-21	[-1, 32, 16, 16]	64
287	Conv2d-22	[-1, 32, 16, 16]	9,216
288	BatchNorm2d-23	[-1, 32, 16, 16]	64
289	BasicBlock-24	[-1, 32, 16, 16]	0
290	Conv2d-25	[-1, 64, 8, 8]	18,432
291	BatchNorm2d-26	[-1, 64, 8, 8]	128
292	Conv2d-27	[-1, 64, 8, 8]	36,864
293	BatchNorm2d-28	[-1, 64, 8, 8]	128
294	Conv2d-29	[-1, 64, 8, 8]	2,048
295	BatchNorm2d-30	[-1, 64, 8, 8]	128
296	BasicBlock-31	[-1, 64, 8, 8]	0
297	Conv2d-32	[-1, 64, 8, 8]	36,864
298	BatchNorm2d-33	[-1, 64, 8, 8]	128
299	Conv2d-34	[-1, 64, 8, 8]	36,864
300	BatchNorm2d-35	[-1, 64, 8, 8]	128
301	BasicBlock-36	[-1, 64, 8, 8]	0
302	Conv2d-37	[-1, 128, 4, 4]	73,728
303	BatchNorm2d-38	[-1, 128, 4, 4]	256
304	Conv2d-39	[-1, 128, 4, 4]	147,456
305	BatchNorm2d-40	[-1, 128, 4, 4]	256
306	Conv2d-41	[-1, 128, 4, 4]	8,192
307	BatchNorm2d-42	[-1, 128, 4, 4]	256
308	BasicBlock-43	[-1, 128, 4, 4]	0
309	Conv2d-44	[-1, 128, 4, 4]	147,456
310	BatchNorm2d-45	[-1, 128, 4, 4]	256
311	Conv2d-46	[-1, 128, 4, 4]	147,456
312	BatchNorm2d-47	[-1, 128, 4, 4]	256

```
313      BasicBlock-48           [-1, 128, 4, 4]          0
314      Linear-49              [-1, 10]                1,290
315 =====
316 Total params: 701,466
317 Trainable params: 701,466
318 Non-trainable params: 0
319 -----
320 Input size (MB): 0.01
321 Forward/backward pass size (MB): 2.81
322 Params size (MB): 2.68
323 Estimated Total Size (MB): 5.50
324 -----
```

325 **B VGG16_S architecture**

326 The architecture of the **VGG16_S** deep neural network used for the benchmark is as follows:

327	Layer (type)	Output Shape	Param #
328	<hr/>		
329	Conv2d-1	[-1, 16, 32, 32]	448
330	ReLU-2	[-1, 16, 32, 32]	0
331	Conv2d-3	[-1, 16, 32, 32]	2,320
332	ReLU-4	[-1, 16, 32, 32]	0
333	MaxPool2d-5	[-1, 16, 16, 16]	0
334	Conv2d-6	[-1, 32, 16, 16]	4,640
335	ReLU-7	[-1, 32, 16, 16]	0
336	Conv2d-8	[-1, 32, 16, 16]	9,248
337	ReLU-9	[-1, 32, 16, 16]	0
338	MaxPool2d-10	[-1, 32, 8, 8]	0
339	Conv2d-11	[-1, 64, 8, 8]	18,496
340	ReLU-12	[-1, 64, 8, 8]	0
341	Conv2d-13	[-1, 64, 8, 8]	36,928
342	ReLU-14	[-1, 64, 8, 8]	0
343	Conv2d-15	[-1, 64, 8, 8]	36,928
344	ReLU-16	[-1, 64, 8, 8]	0
345	MaxPool2d-17	[-1, 64, 4, 4]	0
346	Conv2d-18	[-1, 128, 4, 4]	73,856
347	ReLU-19	[-1, 128, 4, 4]	0
348	Conv2d-20	[-1, 128, 4, 4]	147,584
349	ReLU-21	[-1, 128, 4, 4]	0
350	Conv2d-22	[-1, 128, 4, 4]	147,584
351	ReLU-23	[-1, 128, 4, 4]	0
352	MaxPool2d-24	[-1, 128, 2, 2]	0
353	Conv2d-25	[-1, 128, 2, 2]	147,584
354	ReLU-26	[-1, 128, 2, 2]	0
355	Conv2d-27	[-1, 128, 2, 2]	147,584
356	ReLU-28	[-1, 128, 2, 2]	0
357	Conv2d-29	[-1, 128, 2, 2]	147,584
358	ReLU-30	[-1, 128, 2, 2]	0
359	MaxPool2d-31	[-1, 128, 1, 1]	0
360	Dropout-32	[-1, 128]	0
361	Linear-33	[-1, 128]	16,512
362	ReLU-34	[-1, 128]	0
363	Dropout-35	[-1, 128]	0
364	Linear-36	[-1, 128]	16,512
365	ReLU-37	[-1, 128]	0
366	Linear-38	[-1, 10]	1,290
367	<hr/>		
368	Total params: 955,098		
369	Trainable params: 955,098		
370	Non-trainable params: 0		
371	<hr/>		
372	Input size (MB): 0.01		
373	Forward/backward pass size (MB): 1.12		
374	Params size (MB): 3.64		
375	Estimated Total Size (MB): 4.78		

378 **C ViT_S architecture**

379 The architecture of the **ViT_S** transformer used for the benchmark is as follows:

```

380 ViT(
381     image_size = 32,
382     patch_size = 4,
383     num_classes = 10,
384     dim = 192,
385     depth = 4,
386     heads = 8,
387     mlp_dim = 128,
388     dropout = 0.1,
389     emb_dropout = 0.1
390 )

391 -----
392      Layer (type)          Output Shape       Param #
393 =====
394      Linear-1            [-1, 64, 192]        9,408
395      Dropout-2           [-1, 65, 192]        0
396      LayerNorm-3          [-1, 65, 192]       384
397      Linear-4            [-1, 65, 576]      110,592
398      Linear-5            [-1, 65, 192]      37,056
399      Dropout-6           [-1, 65, 192]        0
400      Attention-7         [-1, 65, 192]        0
401      PreNorm-8           [-1, 65, 192]        0
402      Residual-9          [-1, 65, 192]        0
403      LayerNorm-10         [-1, 65, 192]       384
404      Linear-11           [-1, 65, 128]      24,704
405      GELU-12              [-1, 65, 128]        0
406      Dropout-13           [-1, 65, 128]        0
407      Linear-14           [-1, 65, 192]      24,768
408      Dropout-15           [-1, 65, 192]        0
409      FeedForward-16        [-1, 65, 192]        0
410      PreNorm-17           [-1, 65, 192]        0
411      Residual-18          [-1, 65, 192]        0
412      LayerNorm-19          [-1, 65, 192]       384
413      Linear-20            [-1, 65, 576]      110,592
414      Linear-21            [-1, 65, 192]      37,056
415      Dropout-22           [-1, 65, 192]        0
416      Attention-23          [-1, 65, 192]        0
417      PreNorm-24           [-1, 65, 192]        0
418      Residual-25          [-1, 65, 192]        0
419      LayerNorm-26          [-1, 65, 192]       384
420      Linear-27            [-1, 65, 128]      24,704
421      GELU-28              [-1, 65, 128]        0
422      Dropout-29           [-1, 65, 128]        0
423      Linear-30            [-1, 65, 192]      24,768
424      Dropout-31           [-1, 65, 192]        0
425      FeedForward-32         [-1, 65, 192]        0
426      PreNorm-33           [-1, 65, 192]        0
427      Residual-34          [-1, 65, 192]        0
428      LayerNorm-35          [-1, 65, 192]       384
429      Linear-36            [-1, 65, 576]      110,592

```

430	Linear-37	[-1, 65, 192]	37,056
431	Dropout-38	[-1, 65, 192]	0
432	Attention-39	[-1, 65, 192]	0
433	PreNorm-40	[-1, 65, 192]	0
434	Residual-41	[-1, 65, 192]	0
435	LayerNorm-42	[-1, 65, 192]	384
436	Linear-43	[-1, 65, 128]	24,704
437	GELU-44	[-1, 65, 128]	0
438	Dropout-45	[-1, 65, 128]	0
439	Linear-46	[-1, 65, 192]	24,768
440	Dropout-47	[-1, 65, 192]	0
441	FeedForward-48	[-1, 65, 192]	0
442	PreNorm-49	[-1, 65, 192]	0
443	Residual-50	[-1, 65, 192]	0
444	LayerNorm-51	[-1, 65, 192]	384
445	Linear-52	[-1, 65, 576]	110,592
446	Linear-53	[-1, 65, 192]	37,056
447	Dropout-54	[-1, 65, 192]	0
448	Attention-55	[-1, 65, 192]	0
449	PreNorm-56	[-1, 65, 192]	0
450	Residual-57	[-1, 65, 192]	0
451	LayerNorm-58	[-1, 65, 192]	384
452	Linear-59	[-1, 65, 128]	24,704
453	GELU-60	[-1, 65, 128]	0
454	Dropout-61	[-1, 65, 128]	0
455	Linear-62	[-1, 65, 192]	24,768
456	Dropout-63	[-1, 65, 192]	0
457	FeedForward-64	[-1, 65, 192]	0
458	PreNorm-65	[-1, 65, 192]	0
459	Residual-66	[-1, 65, 192]	0
460	Transformer-67	[-1, 65, 192]	0
461	Identity-68	[-1, 192]	0
462	LayerNorm-69	[-1, 192]	384
463	Linear-70	[-1, 128]	24,704
464	GELU-71	[-1, 128]	0
465	Dropout-72	[-1, 128]	0
466	Linear-73	[-1, 10]	1,290
467	=====		
468	Total params: 827,338		
469	Trainable params: 827,338		
470	Non-trainable params: 0		
471	-----		
472	Input size (MB): 0.01		
473	Forward/backward pass size (MB): 6.76		
474	Params size (MB): 3.16		
475	Estimated Total Size (MB): 9.93		
476	-----		

477 **D ResNet18 and ResNet18_S architectures confusion matrix results on**
478 **CIFAR-10 dataset**

479 **D.1 Resnet18_S with Adabound**

		precision	recall	f1-score	support
482	plane	0.93	0.91	0.92	1000
483	car	0.96	0.96	0.96	1000
484	bird	0.89	0.90	0.89	1000
485	cat	0.81	0.83	0.82	1000
486	deer	0.94	0.91	0.93	1000
487	dog	0.87	0.86	0.87	1000
488	frog	0.95	0.94	0.95	1000
489	horse	0.95	0.94	0.94	1000
490	ship	0.94	0.95	0.94	1000
491	truck	0.95	0.95	0.95	1000
492	accuracy			0.92	10000
494	macro avg	0.92	0.92	0.92	10000
495	weighted avg	0.92	0.92	0.92	10000
496	Test acc	91.720			

498 **D.2 ResNet18 with Adabound**

		precision	recall	f1-score	support
501	plane	0.95	0.96	0.95	1000
502	car	0.97	0.98	0.97	1000
503	bird	0.95	0.92	0.94	1000
504	cat	0.87	0.89	0.88	1000
505	deer	0.97	0.94	0.95	1000
506	dog	0.90	0.92	0.91	1000
507	frog	0.97	0.97	0.97	1000
508	horse	0.98	0.97	0.97	1000
509	ship	0.97	0.96	0.96	1000
510	truck	0.97	0.97	0.97	1000
511	accuracy			0.95	10000
513	macro avg	0.95	0.95	0.95	10000
514	weighted avg	0.95	0.95	0.95	10000

515 **D.3 ResNet18_S with ADAM**

		precision	recall	f1-score	support
518	plane	0.88	0.91	0.90	1000
519	car	0.93	0.96	0.94	1000
520	bird	0.85	0.84	0.85	1000
521	cat	0.81	0.79	0.80	1000
522	deer	0.89	0.89	0.89	1000
523	dog	0.84	0.86	0.85	1000
524	frog	0.90	0.94	0.92	1000
525	horse	0.94	0.91	0.92	1000

```

526      ship      0.94      0.94      0.94      1000
527      truck     0.94      0.93      0.93      1000
528
529      accuracy
530      macro avg   0.89      0.89      0.89      10000
531      weighted avg  0.89      0.89      0.89      10000
532
533 Test acc 89.270

```

534 D.4 ResNet18 with ADAM

		precision	recall	f1-score	support
535					
536					
537	plane	0.89	0.92	0.91	1000
538	car	0.94	0.97	0.95	1000
539	bird	0.88	0.84	0.86	1000
540	cat	0.81	0.79	0.80	1000
541	deer	0.90	0.89	0.90	1000
542	dog	0.84	0.86	0.85	1000
543	frog	0.90	0.94	0.92	1000
544	horse	0.94	0.91	0.92	1000
545	ship	0.94	0.94	0.94	1000
546	truck	0.95	0.94	0.94	1000
547					
548	accuracy			0.90	10000
549	macro avg	0.90	0.90	0.90	10000
550	weighted avg	0.90	0.90	0.90	10000
551					
552	Test acc 89.970				

553 **E Detailed pruning results**

Pruning technique / proportion of pruned connections	0	0.05	0.1	0.15	0.2	0.25	0.3	0.35	0.4	0.45
L1 unstructured	95.0	94.97	95.0	94.98	95.01	94.97	94.89	94.82	94.81	94.71
L1 structured	95.0	94.56	93.87	92.22	45.96	40.3	30.24	24.51	20.25	16.99
Global unstructured	95.0	95.0	95.0	95.0	94.99	94.96	94.94	94.9	94.89	94.86

Table 6: The influence in accuracy based on the proportion of pruned connections in the ResNet18 network, against the CIFAR-10 dataset.

Pruning technique / proportion of pruned connections	0.5	0.55	0.6	0.65	0.7	0.75	0.8	0.85	0.9	0.95	1
L1 unstructured	94.52	94.25	93.98	93.05	91.02	85.48	66.59	27.31	15.48	9.89	10.0
L1 structured	15.67	13.35	10.06	10.01	10.0	9.79	10.0	10.0	10.0	10.0	10.0
Global unstructured	94.87	94.9	94.86	94.91	94.87	94.55	93.01	85.24	52.73	19.85	10.0

Table 7: The influence in accuracy based on the proportion of pruned connections in the ResNet18 network, against the CIFAR-10 dataset. (contd.)