

International Baccalaureate Diploma Program

Computer Science Extended Essay

**ResNet: A Comparative Study of Depth and Width in Neural Network Architectures across  
different datasets**

To what extent does increasing width and depth affect the ResNet's performance (Accuracy,  
Training time, F-1 score, Marginal gains) across various datasets?

3820 Words

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Theoretical Background</b>	<b>4</b>
2.1	ResNet Revisited . . . . .	4
2.1.1	Gradient Vanishing . . . . .	4
2.1.2	Solution provided by ResNet . . . . .	5
2.1.3	ResNet Architecture . . . . .	6
2.1.4	Width in ResNet . . . . .	7
2.1.5	Depth in ResNet . . . . .	8
2.1.6	Wider Or Deeper . . . . .	9
<b>3</b>	<b>Experiment Methodology</b>	<b>10</b>
3.1	Experiment Setup . . . . .	10
3.2	Datasets and Networks used . . . . .	10
3.2.1	Data Augmentation . . . . .	11
3.2.2	Hyperparameters . . . . .	12
3.2.3	Increasing the width and depth . . . . .	12
3.2.4	Measuring the performance of the models . . . . .	13
<b>4</b>	<b>Experiment Result and Analysis</b>	<b>14</b>
4.1	CIFAR-10, FashionMNIST and Analysis . . . . .	14
4.2	GTSRB classification . . . . .	19
<b>5</b>	<b>Conclusion</b>	<b>22</b>
<b>6</b>	<b>Evaluation and future improvements</b>	<b>22</b>

<b>Bibliography</b>	<b>23</b>
<b>Appendix</b>	<b>24</b>

# 1 Introduction

Since the groundbreaking success of the convolutional neural network - AlexNet in 2012, CNNs have experienced exponential growth. Consequently, various new models have emerged every year in the ILSVRC competition to surpass accuracy records. In 2014, Christian Szegedy et al. introduced VGG and GoogleNet, significantly enhancing the model's feature representation capability and learning capacity by creating deeper network architectures. However, deeper models still faced a fundamental challenge - the vanishing gradient problem. This limitation stemmed from inherent design inadequacies in neural networks, impeding effective weight updates in deep neural networks and sometimes rendering networks entirely untrainable. In 2015, this issue is addressed by constructing the deep residual network ResNet using identity mapping by shortcut, successfully implementing networks with up to 1202 layers (He, Zhang, Ren, & Sun, 2016). However, deeper networks also imply longer training times. Since the objective of the ResNet paper was to construct deep networks, it needed to provide more discussion on performance. Subsequently, Wide Residual Network (WRN) was proposed, suggesting that wider and shallower networks outperform deeper networks (Zagoruyko & Komodakis, 2016). However, it did not provide any conclusion to the influence of the datasets to the networks' performance. Thus, two apparent pathways to enhancing neural network models are increasing width or depth. This raises a practical question: Which is more effective, increasing depth or width?

More concretely, it is possible to ask about the gains brought by increasing depth or width to the model's accuracy and to what extent in different scenarios. Can it improve marginal gains and F-1 scores? This paper will address these questions by conducting experimental studies with different widths and depths of ResNet family models on the CIFAR-10, FashionMNIST and GTSRB datasets.

## 2 Theoretical Background

### 2.1 ResNet Revisited

#### 2.1.1 Gradient Vanishing

*Gradient vanishing* is a common issue encountered in deep networks, rooted in the intrinsic mechanism of neural networks. Neural networks update their parameters through backpropagation based on the chain rule in Calculus. This implies that during the gradient descent process, if the derivative of some parts of the activation function remains consistently less than 1, then as the network deepens, the gradient will exponentially decay, leading to the gradient vanishing (Amanatullah, 2023).

To further elaborate, consider a neural network with two layers. Let the input to the network be  $X$ , the weights of the first layer be  $W^{[1]}$ , the bias of the first layer be  $b^{[1]}$ , the activation function of the first layer be  $g^{[1]}$ , the output of the first layer be  $A^{[1]}$ , the weights of the second layer be  $W^{[2]}$ , the bias of the second layer be  $b^{[2]}$ , the activation function of the second layer be  $g^{[2]}$ , and the final output of the network be  $Y$ . Therefore, the mathematical equations can be expressed as follows:

$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ Z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ Y &= g^{[2]}(Z^{[2]}) \end{aligned} \tag{1}$$

The *backpropagation* will compute the gradients of weight parameters  $W^{[1]}$  with respect to the *Loss function* to update the values, which is denoted by  $L$ . The gradients computation is as follows:

$$\frac{\partial L}{\partial W^{[1]}} = \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial Z^{[2]}} \frac{\partial Z^{[2]}}{\partial A^{[1]}} \frac{\partial A^{[1]}}{\partial Z^{[1]}} \frac{\partial Z^{[1]}}{\partial W^{[1]}} \quad (2)$$

To update the weights, the *Gradient descent algorithm* is employed:

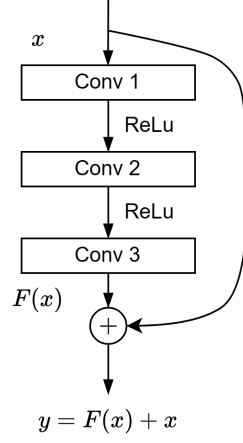
$$W^{[1]} = W^{[1]} - \alpha \frac{\partial L}{\partial W^{[1]}} \quad (3)$$

where  $\alpha$  is the learning rate.

Evidently, when the value of  $\frac{\partial L}{\partial Y} \frac{\partial Y}{\partial Z^{[2]}}$  becomes small enough, at which the fitting capacity of the second layer is strong enough leading to little difference between the predictions and actual values, the derivative would be directly transferred to the value of  $\frac{\partial L}{\partial W^{[1]}}$  and make it small. As a result, based on the Gradient descent equation above,  $W^{[1]}$  would receive little modification. This is just an illustration under the context of a two-layer neural network. As the network depth increases, it can be observed that the weight parameters closer to the input layer are updated more slowly.

### 2.1.2 Solution provided by ResNet

To avoid the gradient vanishing problems faced in deep convolutional neural networks, He and his teams introduced skip connections and designed residual blocks (He et al., 2016). The residual block is the building block of ResNet, and the illustration of the design is shown below:



**Figure 1:** The building block of ResNet architecture, where  $x$  is the input,  $y$  is the output,  $F(x)$  represents three convolutional layers.

To elaborate why this design could make a change, suppose a ResNet with output  $y''$  consisting of two residual blocks, denoted by  $f(X)$  and  $g(X)$  respectively, where  $X$  is the input. According to the building block design above, the output of the network could be expressed as:

$$y'' = f(X) + g(f(X)) \quad (4)$$

Let  $y = f(X)$ ,  $y' = g(f(X))$ . According to the chain rule, the gradient of the loss function  $L$  with respect to the weight parameter  $W$  of the first residual block is expressed as:

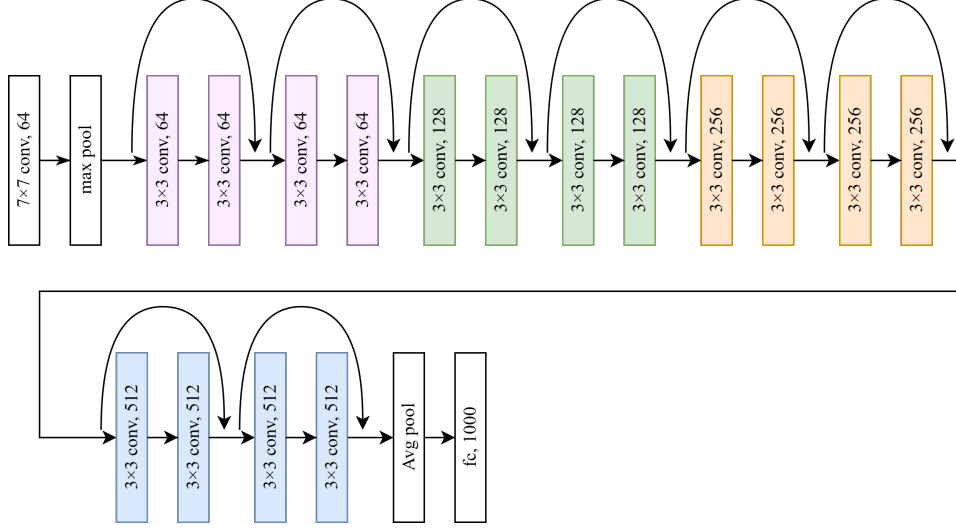
$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y''} \left( \frac{\partial y}{\partial W} + \frac{\partial y'}{\partial W} \right) \quad (5)$$

Even though the value of  $\frac{\partial y'}{\partial W}$  is small enough, there's still a derivative of  $\frac{\partial y}{\partial W}$  existed, therefore, avoid the gradient vanishing problems.

### 2.1.3 ResNet Architecture

Residual networks typically consist of 6 stages, with the first five stages handling convolution operations, while the final stage is a global average pooling layer. This paper only discusses the first

five stages. The first stage consists of a large kernel convolution and a max-pooling layer. The large kernel convolution aims to preserve original pixel information as much as possible, passing the result to the following four stages. These four stages consist of residual blocks with skip connections, each composed of  $n$  residual blocks, as illustrated in Figure 2.



**Figure 2:** An Example of ResNet architecture diagram with 18 layers for 1000-class classification.

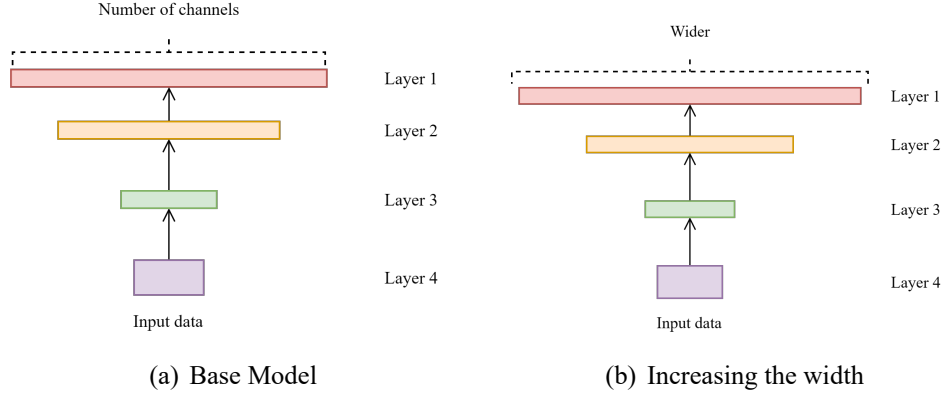
In the depicted diagram, the model employs four distinct colors to delineate the intermediate convolutional modules, housing eight residual blocks, each composed of two  $3 \times 3$  convolutions. These four colors correspond to four stages, each characterized by various convolutional filters, namely  $\{64, 128, 256, 512\}$ . A higher count of filters in a stage enables the model to capture more features. The diagram illustrates a twofold filter increase for each successive stage, forming a pyramidal structure (Paoletti et al., 2019). Such architectural design facilitates the acquisition of more abstract features by the convolutional layers closer to the output, consequently enhancing the model's efficacy in image recognition tasks.

#### 2.1.4 Width in ResNet

In convolutional neural networks, *width* refers to the capacity of each layer to process input data. One way to measure it is by the number of channels in the convolutional layer, also known as



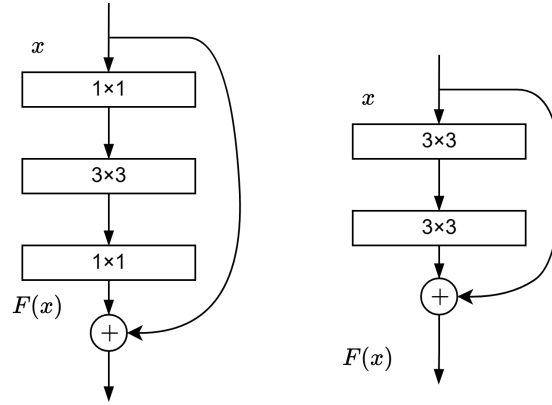
the number of convolutional filters. Figure 3 illustrates the width of a convolutional neural network. In the residual network model shown in Figure 2, the width of the model ranges from  $\{64, 128, 256, 512\}$ . Increasing the width can enhance its learning capacity, allowing it to handle more information from input images. However, it also leads to longer training time and may even result in overfitting, where the model’s performance on the training set significantly surpasses that on the validation set.



**Figure 3:** Demonstration of width in convolutional neural networks.

### 2.1.5 Depth in ResNet

*Depth* refers to the number of layers in a convolutional neural network. The residual neural network depicted in Figure 1 has 17 convolutional layers and one fully connected layer, totaling 18 layers. Increasing the depth of a network is a widely accepted practice in the industry. Deeper models possess more robust capabilities to fit data and, therefore, can perform more complex tasks. In visual models, deeper architectures can extract more detailed and abstract features at different depths, leading to more accurate recognition tasks. However, increasing the depth comes with the cost of higher computational resources and longer training times. Notably, in the ResNet architecture, the authors introduced  $1 \times 1$  convolutions to reduce the model’s parameter count, thereby implementing a bottleneck building block design, as illustrated in Figure 4.



**Figure 4:** The building block of ResNet architecture in practice. Left: a conventional building block. Right: a bottleneck building block

### 2.1.6 Wider Or Deeper

In summary, it is worth noting that increasing the width of a neural network leads to a wider network, while increasing the depth results in a deeper network. Both of these strategies can potentially improve the model's learning capacity to a certain extent. However, the optimal choice between width and depth to maximize performance remains a conundrum. Nonetheless, it is imperative to first comprehend the degree to which each approach enhances the model's learning capabilities before making such determinations. This foundational understanding is indispensable for guiding any further adjustments and optimizations.

It is important to carry out empirical research to ascertain the effectiveness of widening or deepening a network in enhancing its learning capacity in different datasets since there's a great variance in performance while training the networks in different datasets. Besides, it should be noted that there exists a critical threshold beyond which the model becomes overwidened or overdeepened.

To this end, the empirical studies will encompass several components: Firstly, investigating how modifying the width and depth influences the model's performance across various datasets. Lastly, the maximum threshold for practical gains is determined.

### 3 Experiment Methodology

#### 3.1 Experiment Setup

Experiment purpose is to understand the effects of width and depth on the model’s performance in different situations and discover the threshold for practical gains. Reflecting this, the experiment setup would comprise the families of ResNets trained on the universal dataset GTSRB, CIFAR-10 and FashionMNIST, which represents datasets with different dimensions of images.

#### 3.2 Datasets and Networks used

German Traffic Sign Recognition Benchmark (GTSRB) is an image classification dataset consists of 43 classes of traffic signs, splitting into 39209 training images and 12630 test images having the highest resolution among three datasets (Stallkamp, Schlipsing, Salmen, & Igel, 2012). The conventional ResNet architecture will be used in this dataset presented in Table 1.

Layer Name	ResNet-24	ResNet-50	ResNet-101
conv1	$7 \times 7, 64, \text{stride } 2$		
	$3 \times 3, \text{max pool, stride } 2$		
conv2	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$
conv3	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$
conv4	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 23$
conv5	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$

**Table 1:** Conventional ResNet architectures of ResNet-24, ResNet-50 and ResNet-101.

CIFAR-10 is an established computer-vision dataset used for object recognition, which consists of 10 object classes, with 6000 images per class. Each image has dimensions of  $32 \times 32$  pixels (Krizhevsky, 2009). However, the conventional ResNet architectures were built for the ImageNet

dataset, of which each image has dimensions of  $224 \times 224$  pixels, which is completely different from the ones of CIFAR-10. Therefore, a variation of ResNet architectures will be used. In table 2, the adjusted architecture families including ResNet-20, ResNet-56 and ResNet-110 are shown, which would be used in the experiment. The first layer was changed to  $3 \times 3$  convolutions with stride 1, channels of 16. For other three layers below, the numbers of channels was changed to  $\{16, 32, 64\}$  respectively (Ruiz, 2018).

Layer Name	ResNet-20	ResNet-56	ResNet-110
conv1	$3 \times 3, 16, \text{stride } 1$	$3 \times 3, 16, \text{stride } 1$	$3 \times 3, 16, \text{stride } 1$
conv2	$\begin{bmatrix} 3 \times 3, 16 \\ 3 \times 3, 16 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 16 \\ 3 \times 3, 16 \end{bmatrix} \times 9$	$\begin{bmatrix} 3 \times 3, 16 \\ 3 \times 3, 16 \end{bmatrix} \times 18$
conv3	$\begin{bmatrix} 3 \times 3, 32 \\ 3 \times 3, 32 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 32 \\ 3 \times 3, 32 \end{bmatrix} \times 9$	$\begin{bmatrix} 3 \times 3, 32 \\ 3 \times 3, 32 \end{bmatrix} \times 18$
conv4	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 9$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 18$

**Table 2:** Architectures of adjusted ResNet families (ResNet-20, ResNet-56, ResNet-110).

The same model architecture series is also applicable to Fashion-MNIST. This dataset, the smallest among the three, consists of ten thousand grayscale images with dimensions of  $28 \times 28$  pixels, with an input channel count of 1 (Han, Rasul, & Vollgraf, 2017).

### 3.2.1 Data Augmentation

Convolutional neural networks are robust, and augmenting them with increased width and depth may lead to overfitting, where the model’s performance on the training set significantly surpasses that on the validation set. Thus, the experiments will employ data augmentation techniques to mitigate this issue.

Data augmentation involves modifying the original dataset to generate synthetic samples that still represent the same categories as the original samples, thereby preventing the model from simply memorizing the training data. This approach forces the model to generalize better and helps avoid

overfitting.

The experiments will utilize four data augmentation techniques:

- 1) Horizontal flipping
- 2) Vertical flipping
- 3) Rotation of 30 degrees
- 4) Random cropping

These techniques will be applied to the original dataset to generate augmented samples, thereby diversifying the training data and enhancing the model’s ability to generalize to unseen examples.

### **3.2.2 Hyperparameters**

In this experiment, a more stable optimization algorithm, stochastic gradient descent (SGD), will be utilized. For the Fashion-MNIST dataset, the initial learning rate is set to 0.1. The learning rate will decay by a factor of 0.1 at epochs 50, 90, and 110 until completing 150 epochs. The optimizer’s momentum is set to 0.9, and the weight decay is 0.0001. The loss function employed is cross-entropy. For the CIFAR-10 dataset, learning rate decay will occur at epochs 120, 200, and 250 until 300 epochs are completed. Other parameters remain the same as for Fashion-MNIST. For the GTSRB dataset, learning rate decay will occur at epochs 20 until the completion of 40 epochs. The remaining parameters are consistent with the Fashion-MNIST setup.

### **3.2.3 Increasing the width and depth**

The method to increase depth, as illustrated in Tables 1 and 2, involves directly increasing the number of residual blocks. The experiments will utilize the native ResNet models and select models with a depth difference of two-fold. For experiments on the CIFAR-10 and Fashion MNIST

datasets, ResNet-20, ResNet-56, and ResNet-110 will be employed. On the other hand, experiments on the GTSRB dataset will utilize ResNet-24, ResNet-50, and ResNet-101.

In the experiments on width, a width coefficient variable, denoted as  $\omega$ , will be introduced. This variable will be multiplied by the number of channels in each convolutional layer within the residual blocks to increase the model's width, as shown in Table 3. The progression of this ascent will follow a similar pattern to that of increasing depth, i.e., doubling with each step. Therefore, the width coefficients array  $W = \{2, 4\}$  could be derived.

Layer Name	ResNet-20	ResNet-56	ResNet-110
conv1	3×3,16,stride 1	3×3,16,stride 1	3×3,16,stride 1
conv2	$\begin{bmatrix} 3 \times 3, 16\omega \\ 3 \times 3, 16\omega \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 16\omega \\ 3 \times 3, 16\omega \end{bmatrix} \times 9$	$\begin{bmatrix} 3 \times 3, 16\omega \\ 3 \times 3, 16\omega \end{bmatrix} \times 18$
conv3	$\begin{bmatrix} 3 \times 3, 32\omega \\ 3 \times 3, 32\omega \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 32\omega \\ 3 \times 3, 32\omega \end{bmatrix} \times 9$	$\begin{bmatrix} 3 \times 3, 32\omega \\ 3 \times 3, 32\omega \end{bmatrix} \times 18$
conv4	$\begin{bmatrix} 3 \times 3, 64\omega \\ 3 \times 3, 64\omega \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 64\omega \\ 3 \times 3, 64\omega \end{bmatrix} \times 9$	$\begin{bmatrix} 3 \times 3, 64\omega \\ 3 \times 3, 64\omega \end{bmatrix} \times 18$

**Table 3:** Architectures of adjusted ResNet families with widen factor  $\omega$  (ResNet-20, ResNet-56, ResNet-110).

Due to the relatively narrow width of the adjusted ResNet models, there exists considerable room for improvement. In order to further substantiate this assertion, an additional set of experiments will be conducted on the CIFAR-10 and FashionMNIST datasets, introducing a width coefficient of 4.

### 3.2.4 Measuring the performance of the models

The experiments will compare the training time, training accuracy, F-1 score and marginal gains. This comprehensive analysis aims to provide a holistic understanding of the impact of increasing depth and width on model performance.

**Accuracy** Accuracy is a metric that measures how much a model could correctly predicts the outcome.

**Training time** Training time refers to the total time to train the model given the number of epochs.

**F-1 Score** F-1 Score is a more synthesize index that assesses the prediction ability of a model by elaborating on its class-wise performance.

**Marginal gain** Marginal gains a measure of combination of accuracy and training time altogether, which ahs the equation below:

$$\text{Marginal gain} = \frac{\text{Percentage Testing Accuracy}}{\text{Total training time}} \quad (6)$$

## 4 Experiment Result and Analysis

### 4.1 CIFAR-10, FashionMNIST and Analysis

The comparison results of ResNet of different width and depth in CIFAR-10 and FashionMNIST datasets are shown in Table 4.

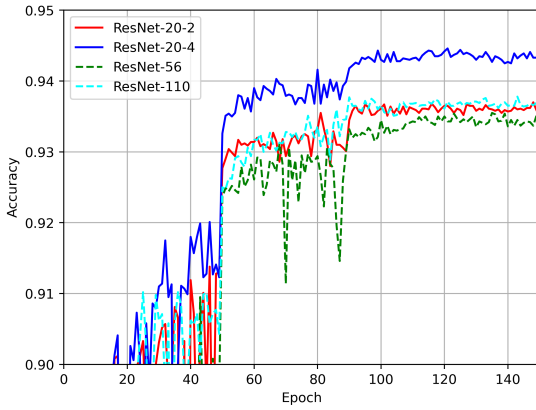
Model	Params	CIFAR-10 classification validation accuracy (%)	FashionMNIST classification validation accuracy (%)	CIFAR-10 Single Epoch Training Time (s)	FashionMNIST Single Epoch Training Time (s)
ResNet-20	0.27M	88.90	93.62	9.4	8.1
ResNet-56	0.85M	89.93	93.54	14.3	14.8
ResNet-110	1.73M	89.5	93.71	23.6	25.3
ResNet-20-2	1.07M	90.77	93.65	9.5	9.1
ResNet-20-4	4.29M	91.83	94.43	24.8	15.48
ResNet-20-8	17.1M	93.14	94.26	77.3	39.4

**Table 4:** The experiment results of ResNet with different width and depth in CIFAR-10 and FashionMNIST datasets. ResNet-20-2 means ResNet-20 plus width coefficient of 2, same as others.

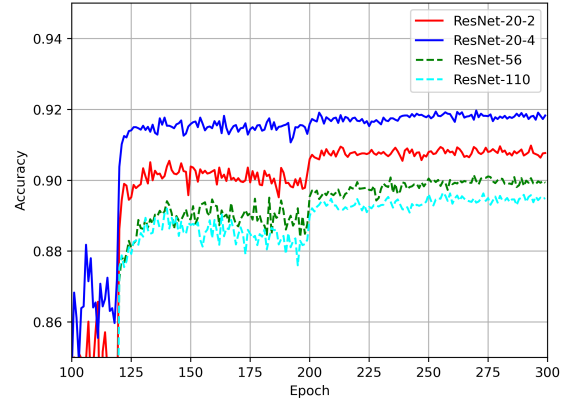
The results in Table 4 show that wider models perform better than deeper models. ResNet-20-8 achieves the highest accuracy of 93.14% in CIFAR-10, which is 3.64% higher than ResNet-110. Even ResNet-20-4 outperforms ResNet-110 by 2.33%. Obviously, compared to increasing depth, increasing the width could gain more in validation accuracy in CIFAR-10 dataset.

In FashionMNIST, wider models still perform better than deeper models in overall. ResNet-20-4 achieves the highest accuracy of 94.43%, which is 0.72% higher than ResNet-110. ResNet-20-2 still outperforms ResNet-56 by 0.11%.

Figures 5(a) and 5(b), respectively, illustrate the learning curves of ResNet-56, ResNet-110, ResNet-20-2, and ResNet-20-44 on the validation sets of the CIFAR-10 and FashionMNIST datasets. On the CIFAR-10 dataset, the accuracy of the width models is higher than that of the depth models. Conversely, on the FashionMNIST dataset, which is smaller in scale and easier to train, the width models still maintain an advantage. However, in the case of ResNet-110, its overall trend indicates a higher accuracy compared to ResNet-20-2.



(a) Validation accuracy in FashionMNIST



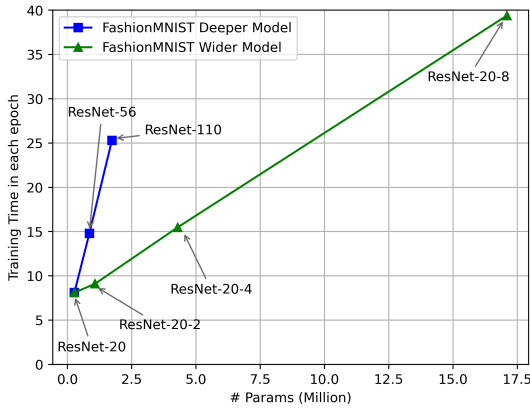
(b) Validation accuracy in CIFAR-10

**Figure 5:** Validation accuracy in FashionMNIST and CIFAR-10 of ResNet-20-2, ResNet-20-4, ResNet-56, ResNet-110.

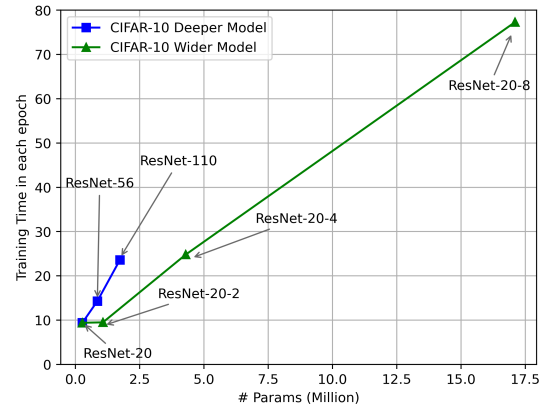
From the perspective of images alone, wide models have more advantages than deep models. While analysis in training time and number of params could provide more insights. According to the training results in Table 4, the parameter quantity of ResNet-20-4 is four times higher than that of the same depth-scaled ResNet-110. Normally, the larger the model's parameter quantity, the longer the training time. This hypothesis has only been proven on the CIFAR-10 dataset, where the single epoch training time of ResNet-20-4 is 4.83% higher than that of ResNet-110, which is not a propor-



tional increase in training time. However, the opposite results were obtained on the FashionMNIST dataset. The single epoch training time of ResNet-20-4 is 63.43% lower than that of ResNet-110, and the single epoch training time of ResNet-20-2 is also 62.64% lower than that of ResNet-56. These results are more intuitively presented in Figure 6. Considering that the size of a single image in both FashionMNIST and CIFAR-10 is relatively small, it can be concluded that wide and shallow models have more advantages both in time and accuracy on relatively small datasets than deep models.



(a) Training time in 1 epoch in FashionMNIST



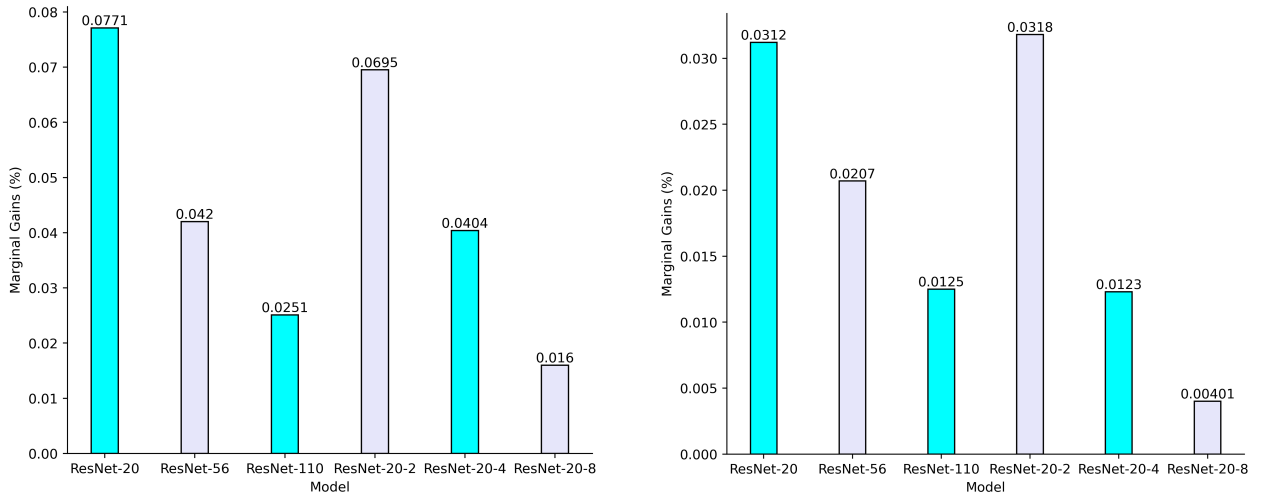
(b) Training time in 1 epoch in CIFAR-10

**Figure 6:** Training time vs number of parameters curves in 1 epoch for wide and deep models.

By contrast, the previous assumption (More parameters, more training time) holds for the ResNet-20-8 model, whose parameter quantity is four times that of ResNet-20-4. The single epoch training time on the CIFAR-10 dataset is also nearly four times that of ResNet-20-4, but the increase in training duration on the FashionMNIST dataset is more minor, only about 2.5 times that of ResNet-20-4. This may be mainly because the FashionMNIST dataset itself consists of grayscale images, i.e., the number of channels in the input image is 1, which makes training models easier. On the other hand, CIFAR-10 consists of RGB color images, with the number of channels in the input image being 3, making training more complex than the former. Nonetheless, regardless of the circumstances, the training time of ResNet-20-8 has a significant increase compared to ResNet-20-4. Even so, ResNet-20-8 did not show an improvement on the FashionMNIST dataset compared to ResNet-20-4, but

rather a decline. While the training time significantly increased compared to ResNet-20-4. This indicates that ResNet-20-8 is overwidened for this dataset, and overwidened models cannot provide any performance gain.

Through further analysis, combining time and accuracy, another indicator could be obtained: Marginal gain, which is how much average percentage increase could be acquired in 1 second. Figure 7(a) and Figure 7(b) respectively show the Marginal Gain of each model in the CIFAR-10 and FashionMNIST datasets.



(a) Marginal gains in FashionMNIST for 150 epoch

(b) Marginal gains in CIFAR-10 for 300 epoch

**Figure 7:** Marginal gains of wide and deep model in FashionMNIST and CIFAR-10 datasets.

The results from the above figure draw a common conclusion: wider networks generally outperform deeper ones in terms of overall performance. When comparing models with the same width scaling and depth scaling factors, it can be seen that on the FashionMNIST dataset, the Marginal gain of ResNet-56 is only 0.042, while that of ResNet-20-2 is 0.0695, which is 65.5% higher than the former. The Marginal gain of ResNet-20-4 is 61.0% higher than that of ResNet-110. However, on the CIFAR-10 dataset, the Marginal gain of ResNet-20-2 is 53.6% higher than that of ResNet-56, but the Marginal Gain of ResNet-20-4 is 1.6% lower than that of ResNet-110. Therefore, it can be seen that overall, increasing the width brings more gains on smaller datasets. But for slightly larger

datasets, such as CIFAR-10 used in the experiment, within a certain range, increasing the width can bring more benefits. However, once this range is exceeded, increasing the depth can bring more benefits compared to increasing the width. Based on the results above, when the width coefficient reaches 8 times or more, the benefits brought by the width will be less than the benefits brought by the depth.

In the CIFAR-10 dataset, the highest Marginal gain is achieved by ResNet-20-2. However, their accuracy is not the highest. This is because the complexity of these two models is not as high as other deeper or wider models, and they are sufficient to handle these datasets, hence they have higher Marginal gains. However, it is clear that these indicators cannot effectively measure the accuracy of the model. In order to more comprehensively explore the accuracy of models with different depths and widths, the experiment also included measurements of other performance indicators in the test set, as shown in Table 5.

Model	CIFAR-10 classification accuracy (%)	CIFAR-10 classification F-1 Score	FashionMNIST classification accuracy (%)	FashionMNIST classification F-1 Score
ResNet20	88.90	88.17	93.61	93.60
ResNet56	89.93	89.11	93.51	93.50
ResNet110	89.50	90.58	93.70	93.68
ResNet20-2	90.77	90.63	93.59	94.27
ResNet20-4	91.83	91.32	94.33	94.07
ResNet20-8	93.14	92.88	94.17	94.16

**Table 5:** The F-1 Score of ResNet with different width and depth in CIFAR-10 and FashionMNIST dataset.

Table 5 shows that the F-1 score remains consistent with the results based on accuracy metrics. For most models, the F1 score is slightly lower than the accuracy, except for ResNet-110 on the CIFAR-10 dataset, where the F-1 score is 1% higher than the accuracy. This discrepancy may be attributed to oscillations during the training of ResNet-110. Nevertheless, regardless of this observation, for relatively smaller datasets (with small image sizes and fewer categories), depth models exhibit robust learning capacity. However, choosing wider and shallower models proves to be the optimal solution.

## 4.2 GTSRB classification

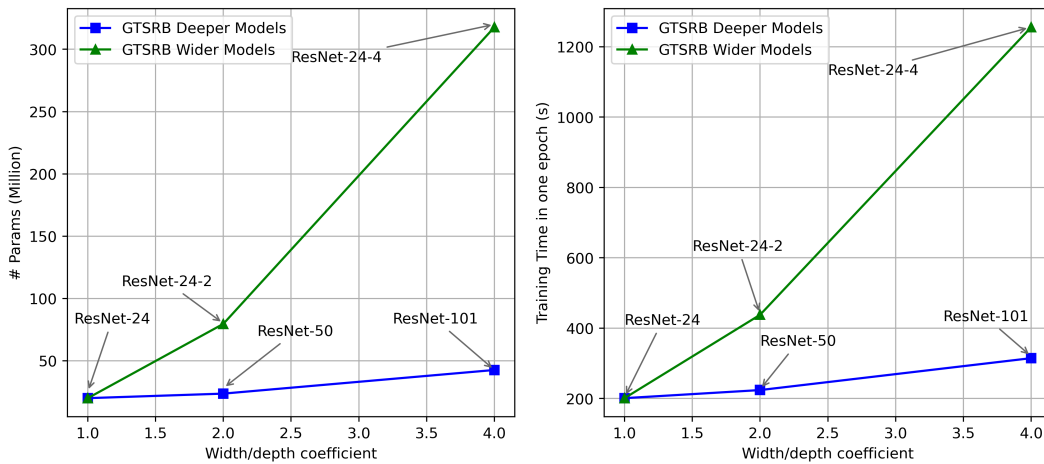
The comparison results of ResNet of different width and depth in GTSRB dataset are shown in Table 6.

Model	Params	GTSRB classification validation accuracy (%)	GTSRB Single Epoch Training Time (s)
ResNet-24	19.96M	92.01	201
ResNet-50	23.60M	94.44	224
ResNet-101	42.58M	94.84	314.7
ResNet-24-2	79.58M	93.04	438
ResNet-24-4	317.75M	92.55	1255.6

**Table 6:** The experiment results of ResNet with different width and depth in GTSRB dataset.

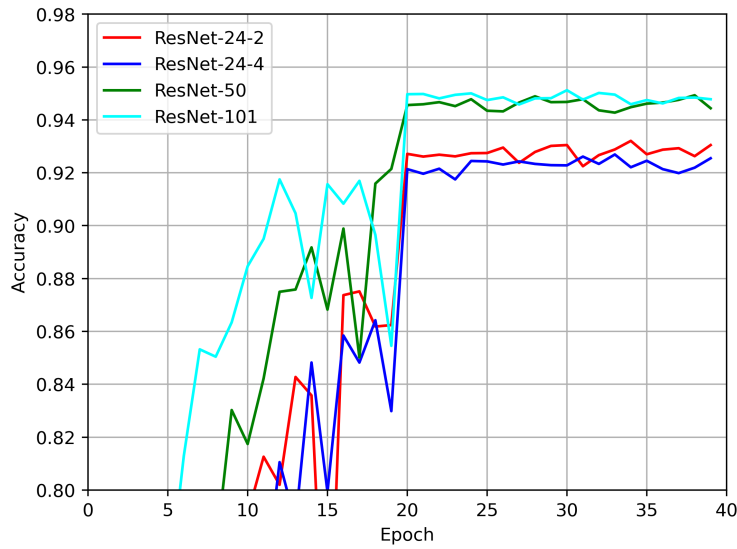
The results in Table 6 show that deeper models perform better in GTSRB dataset. ResNet-110 achieves the highest accuracy of 94.84%, which is 2.47% higher than ResNet-24-4. Even ResNet-50 outperforms ResNet-24-4 by 2.04%.

Additionally, Table 7 presents the relevant data for each model. It is evident that models based on the conventional architecture of ResNet have a substantial parameter count. Increasing the width leads to an exponential growth in both the model's training time and parameter count. Conversely, increasing the depth results in a relatively smoother growth in both parameter count and training time for the model. This variation is more visually represented in Figure 8.



**Figure 8:** The impact of increasing either width or depth of ResNet on number of parameters and training time in one epoch.

Meanwhile, there is another intriguing phenomenon observed when increasing the width of ResNet-24 to four times its original width. Surprisingly, the accuracy of the model decreases by 0.45% compared to ResNet-24-2. Additionally, the training time of ResNet-24-4 is approximately three times longer than that of ResNet-24-2. It is evident that after doubling the width, the benefits of improving model performance become negative. Therefore, increasing the width is not an optimal choice for enhancing performance on the GTSRB dataset.



**Figure 9:** Validation accuracy in GTSRB of ResNet-24-2, ResNet-24-4, ResNet-50, ResNet-101

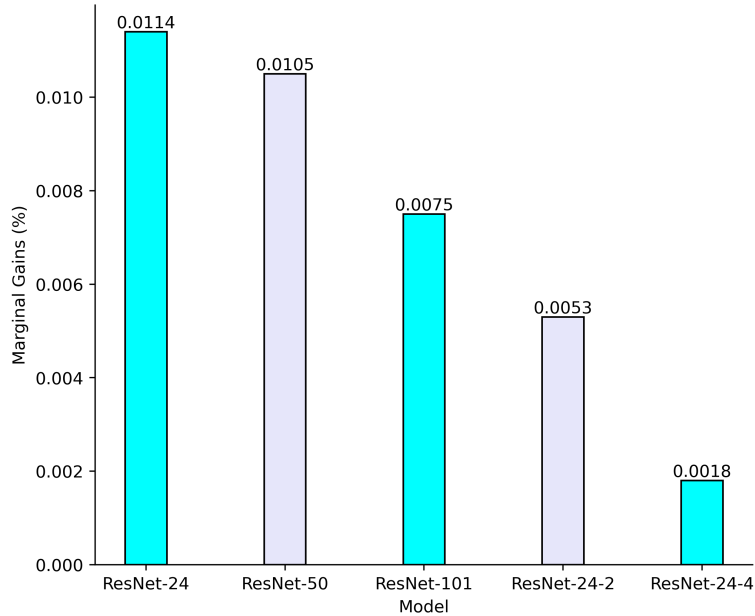
Figure 9 displays the learning curves of ResNet-24-2, ResNet-24-4, ResNet-50, and ResNet-101, further confirming the earlier conclusion: increasing the depth of the model is the optimal solution for the GTSRB dataset. Moreover, in Figure 9, ResNet-101 demonstrates superior learning capabilities right from the start, surpassing the other three models in accuracy. Although after the first learning rate decay, the accuracy of the model remains within the same range of 94% to 95% as ResNet-50, the overall trend of the learning curve for ResNet-101 is higher than that of ResNet-50.

Table 7 displays the F-1 Score of different models, which align closely with their accuracy metrics, further affirming the aforementioned observations.

Model	GTSRB classification validation accuracy (%)	F-1 Score (%)
ResNet-24	92.01	91.80
ResNet-50	94.44	94.21
ResNet-101	94.84	94.90
ResNet-24-2	93.04	92.46
ResNet-24-4	92.55	92.30

**Table 7:** The F-1 scores of ResNet-24-2, ResNet-24-4, ResNet-24, ResNet-50 and ResNet-101 in GTSRB dataset.

From Figure 10, it can be clearly seen that for models with the same width and depth coefficients, deep models have a higher marginal gain, which is consistent with the conclusions drawn previously. Through comprehensive analysis, it becomes evident that the model performance on this dataset contrasts sharply with CIFAR-10 and FashionMNIST. This discrepancy stems from the scale difference between the datasets, with GTSRB representing a large-scale dataset while CIFAR-10 and FashionMNIST are small-scale datasets. Therefore, it can be concluded that for ResNet models, increasing depth leads to better performance on large-scale datasets.



**Figure 10:** Marginal gains of different width and depth model in GTSRB dataset.

## 5 Conclusion

This paper, based on the design of residual neural networks, delves into the impact of increasing depth and width on model performance metrics. It explores the question of whether to prioritize depth or width augmentation when dealing with different datasets. Experimental results indicate that for smaller datasets like CIFAR-10 and FashionMNIST, opting for wider and shallower networks yields optimal results. Conversely, for larger datasets such as GTSRB, selecting deeper networks proves to be a superior solution. These findings are expected to provide valuable insights for future research and model training in the field of computer vision.

## 6 Evaluation and future improvements

The experimental outcomes align with the projected expectations. Hyperparameter tuning has facilitated the attainment of superior training accuracy across diverse datasets. For instance, the discrepancy between the experimental accuracy on the CIFAR-10 dataset and the reported results in the literature is a mere 2%. The juxtaposition of wide and deep models across an array of datasets corroborates the initial hypotheses. Nonetheless, a limitation surfaces in the form of the exclusive use of GTSRB for experiments involving larger datasets, which does not conclusively substantiate the findings. To fortify the conclusions, the inclusion of additional datasets, such as ImageNet or Flowers-101, in the experimental design is recommended. Both aforementioned datasets comprise large-sized images, thereby fulfilling the experimental prerequisites. Moreover, the study confines the evaluation of the performance of wide and deep models to classification tasks. Given that visual tasks extend beyond image recognition to encompass object detection, edge detection, among others, the incorporation of these contexts in future experiments would undoubtedly augment the robustness of the study's conclusions.

## Bibliography

- Amanatullah. (2023). *Vanishing gradient problem in deep learning: Understanding, intuition, and solutions*. <https://medium.com/@amanatulla1606/vanishing-gradient-problem-in-deep-learning-understanding-intuition-and-solutions-da90ef4ecb54>.
- Han, X., Rasul, K., & Vollgraf, R. (2017). Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. <https://arxiv.org/pdf/1708.07747v2.pdf>.
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770-778. doi: 10.1109/cvpr.2016.90
- Krizhevsky, A. (2009). *Cifar-10 and cifar-100 datasets*. <https://www.cs.toronto.edu/~kriz/cifar.html>.
- Paoletti, M. E., Haut, J. M., Fernandez-Beltran, R., Plaza, J., Plaza, A. J., & Pla, F. (2019). Deep pyramidal residual networks for spectral-spatial hyperspectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 57, 740-754. doi: 10.1109/tgrs.2018.2860125
- Ruiz, P. (2018). *Resnets for cifar-10*. <https://towardsdatascience.com/resnets-for-cifar-10-e63e900524e0>.
- Stallkamp, J., Schlipsing, M., Salmen, J., & Igel, C. (2012). Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, 32, 323-332. doi: 10.1016/j.neunet.2012.02.016
- Zagoruyko, S., & Komodakis, N. (2016). Wide residual networks. *Proceedings of the British Machine Vision Conference 2016*. doi: 10.5244/c.30.87



**Appendix**

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import torch.nn.init as init
5 from torch.autograd import Variable
6 import torch.optim as optim
7 from torchvision import datasets
8 from torchvision import transforms

```

```

1 device = (torch.device('cuda') if torch.cuda.is_available() else
torch.device('cpu'))
2 print(f"Training on device {device}.")

```

```

1 __all__ = ['ResNet', 'resnet20', 'resnet32', 'resnet44', 'resnet56',
'resnet110', 'resnet1202']
2
3 def _weights_init(m):
4     classname = m.__class__.__name__
5     #print(classname)
6     if isinstance(m, nn.Linear) or isinstance(m, nn.Conv2d):
7         init.kaiming_normal_(m.weight)
8
9 class LambdaLayer(nn.Module):
10     def __init__(self, lambd):
11         super(LambdaLayer, self).__init__()
12         self.lambd = lambd
13
14     def forward(self, x):
15         return self.lambd(x)
16
17
18 class BasicBlock(nn.Module):
19     expansion = 1
20
21     def __init__(self, in_planes, planes, stride=1, option='A'):
22         super(BasicBlock, self).__init__()
23         self.conv1 = nn.Conv2d(in_planes, planes, kernel_size=3,
stride=stride, padding=1, bias=False)
24         self.bn1 = nn.BatchNorm2d(planes)
25         self.conv2 = nn.Conv2d(planes, planes, kernel_size=3, stride=1,
padding=1, bias=False)
26         self.bn2 = nn.BatchNorm2d(planes)
27
28         self.shortcut = nn.Sequential()
29         if stride != 1 or in_planes != planes:
30             if option == 'A':
31                 self.shortcut = LambdaLayer(lambda x:
32                                             F.pad(x[:, :, ::2, ::2], (0, 0,
0, 0, planes//4, planes//4), "constant", 0))
33             elif option == 'B':
34                 self.shortcut = nn.Sequential(

```

```

35         nn.Conv2d(in_planes, self.expansion * planes,
kernel_size=1, stride=stride, bias=False),
36         nn.BatchNorm2d(self.expansion * planes)
37     )
38
39     def forward(self, x):
40         out = F.relu(self.bn1(self.conv1(x)))
41         out = self.bn2(self.conv2(out))
42         out += self.shortcut(x)
43         out = F.relu(out)
44         return out

```

```

1  class ResNet20_w(nn.Module):
2      def __init__(self, block, num_blocks, width_f, num_classes=10):
3          super(ResNet20_w, self).__init__()
4          widths = [16 * width_f, 32 * width_f, 64 * width_f]
5          self.in_planes = widths[0]
6          self.conv1 = nn.Conv2d(3, widths[0], kernel_size=3, stride=1,
padding=1, bias=False)
7          self.bn1 = nn.BatchNorm2d(widths[0])
8          self.layer1 = self._make_layer(block, widths[0], num_blocks[0],
stride=1)
9          self.layer2 = self._make_layer(block, widths[1], num_blocks[1],
stride=2)
10         self.layer3 = self._make_layer(block, widths[2], num_blocks[2],
stride=2)
11         self.linear = nn.Linear(widths[2], num_classes)
12
13         self.apply(_weights_init)
14
15     def _make_layer(self, block, planes, num_blocks, stride):
16         strides = [stride] + [1]*(num_blocks-1)
17         layers = []
18         for stride in strides:
19             layers.append(block(self.in_planes, planes, stride))
20             self.in_planes = planes * block.expansion
21
22         return nn.Sequential(*layers)
23
24     def forward(self, x):
25         out = F.relu(self.bn1(self.conv1(x)))
26         out = self.layer1(out)
27         out = self.layer2(out)
28         out = self.layer3(out)
29         out = F.avg_pool2d(out, out.size()[3])
30         out = out.view(out.size(0), -1)
31         out = self.linear(out)
32         return out

```

```

1 def resnet20_2():
2     return ResNet20_w(BasicBlock, [3, 3, 3], 2)
3
4 def resnet20_4():
5     return ResNet20_w(BasicBlock, [3, 3, 3], 4)
6
7 def resnet20_8():
8     return ResNet20_w(BasicBlock, [3, 3, 3], 8)

```

```

1 class ResNet(nn.Module):
2     def __init__(self, block, num_blocks, num_classes=10):
3         super(ResNet, self).__init__()
4         self.in_planes = 16
5
6         self.conv1 = nn.Conv2d(3, 16, kernel_size=3, stride=1, padding=1,
7 bias=False)
8         self.bn1 = nn.BatchNorm2d(16)
9         self.layer1 = self._make_layer(block, 16, num_blocks[0], stride=1)
10        self.layer2 = self._make_layer(block, 32, num_blocks[1], stride=2)
11        self.layer3 = self._make_layer(block, 64, num_blocks[2], stride=2)
12        self.linear = nn.Linear(64, num_classes)
13
14        self.apply(_weights_init)
15
16        def _make_layer(self, block, planes, num_blocks, stride):
17            strides = [stride] + [1]*(num_blocks-1)
18            layers = []
19            for stride in strides:
20                layers.append(block(self.in_planes, planes, stride))
21                self.in_planes = planes * block.expansion
22
23            return nn.Sequential(*layers)
24
25        def forward(self, x):
26            out = F.relu(self.bn1(self.conv1(x)))
27            out = self.layer1(out)
28            out = self.layer2(out)
29            out = self.layer3(out)
30            out = F.avg_pool2d(out, out.size()[3])
31            out = out.view(out.size(0), -1)
32            out = self.linear(out)
33            return out

```

```

1 def resnet20():
2     return ResNet(BasicBlock, [3, 3, 3])
3
4
5 def resnet32():
6     return ResNet(BasicBlock, [5, 5, 5])
7
8
9 def resnet44():
10    return ResNet(BasicBlock, [7, 7, 7])
11

```

```

12
13 def resnet56():
14     return ResNet(BasicBlock, [9, 9, 9])
15
16
17 def resnet110():
18     return ResNet(BasicBlock, [18, 18, 18])
19
20
21 def resnet1202():
22     return ResNet(BasicBlock, [200, 200, 200])

```

```

1 import pickle
2
3 class ModelLogger:
4     def __init__(self, path, model_name, params_num):
5         self.loss_train = []
6         self.acc_train = []
7         self.loss_val = []
8         self.acc_val = []
9         self.params_num = params_num
10        self.model_name = model_name
11        self.path = path
12
13        basename = os.path.basename(self.path)
14        name, ext = os.path.splitext(basename)
15        dir_name = os.path.dirname(self.path)
16        self.data_path = dir_name + "/" + name + "_data.pkl"
17
18        self.startingTime = ""
19        self.start_time = time.time()
20        self.firstLog = True
21        self.duration = 0.0
22
23    def Log(self, string):
24        if self.firstLog == True:
25            startingTime = datetime.datetime.now()
26            self.firstLog = False
27            with open(self.path, 'w') as file:
28                file.write("Training {} starting at {}
29\n".format(self.model_name, startingTime))
30                file.write(f"#Params: {self.params_num}\n")
31                file.write("=====\n")
32                file.write(string + "\n")
33        else:
34            with open(self.path, 'a') as file:
35                file.write(string + "\n")
36
37    def writeLog(self):
38        end_time = time.time()
39        with open(self.data_path, 'wb') as file:
40            pickle.dump([self.loss_train, self.acc_train, self.loss_val,
41self.acc_val, end_time-self.start_time], file)
42
43    def ReadLog(self):

```

```

42         with open(self.data_path, 'rb') as file:
43             loaded_array = pickle.load(file)
44             self.loss_train, self.acc_train, self.loss_val, self.acc_val,
self.duration = loaded_array

```

```

1  import datetime
2  import time
3  import os
4
5  def training_loop(n_epochs, optimizer, scheduler, model, loss_fn,
train_loader, valid_loader, dir, model_name):
6      # Write into log
7      logger = ModelLogger(dir, model_name, sum([p.numel() for p in
model.parameters()])))
8
9      for epoch in range(1, n_epochs + 1):
10         start_time = time.time()
11
12         # Training set
13         model.train()
14         loss_train, acc_train = train(optimizer, model, loss_fn,
train_loader)
15         logger.loss_train.append(loss_train)
16         logger.acc_train.append(acc_train)
17
18         # Valid set
19         model.eval()
20         loss_val, acc_val = valid(model, loss_fn, valid_loader)
21         logger.loss_val.append(loss_val)
22         logger.acc_val.append(acc_val)
23
24         if not scheduler is None:
25             scheduler.step()
26
27         end_time = time.time()
28         epoch_time = end_time - start_time
29
30         # log data to file
31         logging_str = f"Epoch {epoch}/{n_epochs} \nTraining - Loss:
{loss_train}, Acc: {acc_train} \nValidation - Loss: {loss_val}, Acc:
{acc_val}\n" + f"Time taken for epoch: {epoch_time:.2f} seconds\n"
32         logger.Log(logging_str)
33
34         # Jupyter ouput
35         if (epoch == 1 or epoch % 10 == 0):
36             print('{} Epoch {}'.format(datetime.datetime.now(), epoch))
37             print("Training Loss: {} | Acc: {}".format(loss_train /
len(train_loader), acc_train))
38             print("Validation Loss: {} | Acc: {} ".format(loss_val /
len(valid_loader), acc_val))
39             print()
40
41         # Write to file
42         logger.WriteLog()
43

```

```
44     return logger
```

```
1  def valid(model, loss_fn, valid_loader):
2      loss_valid = 0.0
3      correct = 0
4      total = 0
5      #bar = ProgressBar().start()
6      with torch.no_grad():
7          for imgs, labels in valid_loader:
8              imgs = imgs.to(device=device)
9              labels = labels.to(device=device)
10             outputs = model(imgs)
11             # calculate the percentage
12             _, predicted = torch.max(outputs, dim=1)
13             total += labels.shape[0]
14             correct += int((predicted == labels).sum())
15
16             loss = loss_fn(outputs, labels)
17             loss_valid += loss.item()
18         #bar.finish()
19     return loss_valid, correct / total
20
21 def train(optimizer, model, loss_fn, train_loader):
22     loss_train = 0.0
23     correct = 0
24     total = 0
25     #bar = ProgressBar().start()
26     for imgs, labels in train_loader:
27         imgs = imgs.to(device=device)
28         labels = labels.to(device=device)
29         outputs = model(imgs)
30         # calculate the percentage
31         _, predicted = torch.max(outputs, dim=1)
32         total += labels.shape[0]
33         correct += int((predicted == labels).sum())
34
35         optimizer.zero_grad()
36
37         loss = loss_fn(outputs, labels)
38
39         # l2_lambda = 0.001
40         # l2_norm = sum(p.pow(2.0).sum() for p in model.parameters())
41         # loss = loss + l2_lambda * l2_norm
42
43         loss.backward()
44
45         optimizer.step()
46
47         loss_train += loss.item()
48     #bar.finish()
49     return loss_train, correct / total
```

```
1  n_classes = 10
```

```
2
```

```

3 def test(val_loader, model, loss_fn):
4     net.eval()
5     test_loss = 0
6     target_num = torch.zeros((1, n_classes))
7     predict_num = torch.zeros((1, n_classes))
8     acc_num = torch.zeros((1, n_classes))
9
10    with torch.no_grad():
11        for imgs, labels in val_loader:
12            imgs, labels = imgs.to(device), labels.to(device)
13            outputs = model(imgs)
14            loss = loss_fn(outputs, labels)
15
16            test_loss += loss.item()
17            _, predicted = outputs.max(1)
18
19            pre_mask = torch.zeros(outputs.size()).scatter_(1,
predicted.cpu().view(-1, 1), 1.)
20            predict_num += pre_mask.sum(0)
21            tar_mask = torch.zeros(outputs.size()).scatter_(1,
targets.data.cpu().view(-1, 1), 1.)
22            target_num += tar_mask.sum(0)
23            acc_mask = pre_mask * tar_mask
24            acc_num += acc_mask.sum(0)
25
26            recall = acc_num / target_num
27            precision = acc_num / predict_num
28            F1 = 2 * recall * precision / (recall + precision)
29
30            print('Recall {}, Precision {}, F1-score
{}'.format(torch.sum(recall)/10, torch.sum(precision)/10,
torch.sum(F1)/10))

```

## CIFAR-10 Data Processing

```

1 data_path = './data/CIFAR10/'
2 cifar10 = datasets.CIFAR10(data_path, train=True, download=True)
3 cifar10_val = datasets.CIFAR10(data_path, train=False, download=True)

```

```

1 data_path = "./data/CIFAR10"
2 cifar10 = datasets.CIFAR10(data_path, train=True, download=False,
transform=transforms.Compose([
3     transforms.ToTensor(),
4     transforms.Normalize((0.4914, 0.4822, 0.4465), (0.2470, 0.2435,
0.2616)),
5     transforms.RandomHorizontalFlip(p=0.5),
6     transforms.RandomVerticalFlip(p=0.5),
7     transforms.RandomRotation(30),
8     transforms.RandomResizedCrop(32, scale=(0.75, 1.0), ratio=(1.0, 1.0),
antialias=True)
9 ]))
10
11 cifar10_val = datasets.CIFAR10(data_path, train=False, download=False,
transform=transforms.Compose([

```



```

12     transforms.ToTensor(),
13     transforms.Normalize((0.4914, 0.4822, 0.4465), (0.2470, 0.2435,
14     0.2616)))
    ]))

```

## CIFAR-10 Training

```

1 train_loader = torch.utils.data.DataLoader(cifar10, batch_size=128,
    pin_memory=True, num_workers=8, shuffle=True)
2 valid_loader = torch.utils.data.DataLoader(cifar10_val, batch_size=128,
    pin_memory=True, num_workers=8, shuffle=True)
3
4 model = resnet20_2().to(device=device)
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
    weight_decay=1e-4, nesterov=True)
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[120, 200,
    250], gamma=0.1)
7 loss_fn = nn.CrossEntropyLoss()

```

```

1 logger_20_2_cifar = training_loop(300, optimizer, scheduler, model,
    loss_fn, train_loader, valid_loader, "./Log/Exp/CIFAR/ResNet20_2",
    "Resnet20-2")

```

```

1 torch.save(model.state_dict(), "./Models/CIFAR/resnet20_2.pth")

```

```

1 test(valid_loader, ResNet_20_2, loss_fn)

```

```

1 train_loader = torch.utils.data.DataLoader(cifar10, batch_size=128,
    pin_memory=True, num_workers=8, shuffle=True)
2 valid_loader = torch.utils.data.DataLoader(cifar10_val, batch_size=128,
    pin_memory=True, num_workers=8, shuffle=True)
3
4 model = resnet20_4().to(device=device)
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
    weight_decay=1e-4, nesterov=True)
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[120, 200,
    250], gamma=0.1)
7 loss_fn = nn.CrossEntropyLoss()

```

```

1 logger_20_4_cifar = training_loop(300, optimizer, scheduler, model,
    loss_fn, train_loader, valid_loader, "./Log/Exp/CIFAR/ResNet20_4",
    "Resnet20-4")

```

```

1 torch.save(model.state_dict(), "./Models/CIFAR/resnet20_4.pth")

```

```

1 test(valid_loader, model, loss_fn)

```

```

1 train_loader = torch.utils.data.DataLoader(cifar10, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
2 valid_loader = torch.utils.data.DataLoader(cifar10_val, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
3
4 model = resnet20_8().to(device=device)
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
  weight_decay=1e-4, nesterov=True)
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[120, 200,
  250], gamma=0.1)
7 loss_fn = nn.CrossEntropyLoss()

```

```

1 logger_20_8_cifar = training_loop(300, optimizer, scheduler, model,
  loss_fn, train_loader, valid_loader, "./Log/Exp/CIFAR/ResNet20_8",
  "Resnet20-8")

```

```

1 torch.save(model.state_dict(), "./Models/CIFAR/resnet20_8.pth")

```

```

1 test(valid_loader, model, loss_fn)

```

```

1 logger_20_2_cifar.ReadLog()
2 logger_20_4_cifar.ReadLog()
3 logger_20_8_cifar.ReadLog()

```

```

1 train_loader = torch.utils.data.DataLoader(cifar10, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
2 valid_loader = torch.utils.data.DataLoader(cifar10_val, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
3
4 model = resnet20().to(device=device)
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
  weight_decay=1e-4, nesterov=True)
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[120, 200,
  250], gamma=0.1)
7 loss_fn = nn.CrossEntropyLoss()

```

```

1 logger_20_cifar = training_loop(300, optimizer, scheduler, model, loss_fn,
  train_loader, valid_loader, "./Log/Exp/CIFAR/ResNet20", "Resnet20")

```

```

1 torch.save(model.state_dict(), "./Models/CIFAR/resnet20.pth")

```

```

1 test(valid_loader, model, loss_fn)

```

```

1 train_loader = torch.utils.data.DataLoader(cifar10, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
2 valid_loader = torch.utils.data.DataLoader(cifar10_val, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
3
4 model = resnet56().to(device=device)
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
  weight_decay=1e-4, nesterov=True)
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[120, 200,
  250], gamma=0.1)
7 loss_fn = nn.CrossEntropyLoss()

```

```

1 logger_56_cifar = training_loop(300, optimizer, scheduler, model, loss_fn,
  train_loader, valid_loader, "./Log/Exp/CIFAR/ResNet56", "Resnet56")

```

```

1 torch.save(model.state_dict(), "./Models/CIFAR/resnet56.pth")

```

```

1 test(valid_loader, model, loss_fn)

```

```

1 train_loader = torch.utils.data.DataLoader(cifar10, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
2 valid_loader = torch.utils.data.DataLoader(cifar10_val, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
3
4 model = resnet110().to(device=device)
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
  weight_decay=1e-4, nesterov=True)
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[120, 200,
  250], gamma=0.1)
7 loss_fn = nn.CrossEntropyLoss()

```

```

1 logger_110_cifar = training_loop(300, optimizer, scheduler, model, loss_fn,
  train_loader, valid_loader, "./Log/Exp/CIFAR/ResNet110", "Resnet110")

```

```

1 torch.save(model.state_dict(), "./Models/CIFAR/resnet110.pth")

```

```

1 test(valid_loader, model, loss_fn)

```

```

1 logger_20_cifar.ReadLog()
2 logger_56_cifar.ReadLog()
3 logger_110_cifar.ReadLog()

```

## Plotting graph

```

1 import matplotlib.pyplot as plt
2
3 plt.plot(logger_20_2_cifar.acc_val, label='ResNet-20-2', color='r')
4 plt.plot(logger_20_4_cifar.acc_val, label='ResNet-20-4', color='b')
5

```

```

6 plt.plot(logger_56_cifar.acc_val, label='ResNet-56', ls='--', color='g')
7 plt.plot(logger_110_cifar.acc_val, label='ResNet-110', ls='--',
  color='cyan')
8
9 plt.legend()
10 plt.xlim(100, 300)
11 plt.ylim(0.85, 0.95)
12 plt.grid()
13 plt.xlabel('Epoch')
14 plt.ylabel('Accuracy')
15 plt.savefig("./Imgs/CIFAR10_Cmp", dpi=400)
16 plt.show()

```

```

1 CIFAR_time_deep = [9.4,14.3,23.6]
2 CIFAR_Params_deep = [0.27, 0.853, 1.73]
3 CIFAR_time_wdith = [9.4,9.5,24.8,77.3]
4 CIFAR_Params_width = [0.27, 1.07, 4.29, 17.1]
5
6 CIFAR_Model_names = ["ResNet-20", "ResNet-56", "ResNet-110"]
7 CIFAR_Model_width_names = ["ResNet-20", "ResNet-20-2", "ResNet-20-4",
  "ResNet-20-8"]

```

```

1 fig, ax = plt.subplots()
2 plt.plot(CIFAR_Params_deep, CIFAR_time_deep, 's-b', label='CIFAR-10 Deeper
  Model')
3 plt.plot(CIFAR_Params_width, CIFAR_time_wdith, '^-g', label='CIFAR-10 wider
  Model')
4 plt.legend()
5 plt.grid()
6 plt.xlabel("# Params (Million)")
7 plt.ylabel("Training Time in each epoch")
8 plt.ylim(0, 80)
9
10 ax.annotate('ResNet-20-8', xy=(17, 75), xytext=(14.5, 62),
  arrowprops=dict(arrowstyle = '->', color = 'dimgray'))
11 ax.annotate('ResNet-20-4', xy=(4.5, 24), xytext=(7, 32),
  arrowprops=dict(arrowstyle = '->', color = 'dimgray'))
12 ax.annotate('ResNet-20-2', xy=(1.1, 8.8), xytext=(3.2, 15),
  arrowprops=dict(arrowstyle = '->', color = 'dimgray'))
13 ax.annotate('ResNet-20', xy=(0.27, 9.4), xytext=(0.3, 1.5),
  arrowprops=dict(arrowstyle = '->', color = 'dimgray'))
14 ax.annotate('ResNet-56', xy=(0.853, 14.5), xytext=(-0.4, 40),
  arrowprops=dict(arrowstyle = '->', color = 'dimgray'))
15 ax.annotate('ResNet-110', xy=(1.73, 24), xytext=(2.5, 45),
  arrowprops=dict(arrowstyle = '->', color = 'dimgray'))
16
17 plt.savefig('CIFAR-10_wide_vs_Deep', dpi=400)
18 plt.show()

```

```

1 import matplotlib.pyplot as plt
2
3 y = [0.0312, 0.0207, 0.0125, 0.0318, 0.0123, 0.00401]

```

```

4 x_label = ["ResNet-20", "ResNet-56", "ResNet-110", "ResNet-20-2", "ResNet-
  20-4", "ResNet-20-8"]
5
6 bar_width = 0.3
7
8 fig = plt.figure(figsize=(7, 6), dpi=100)
9 ax = plt.axes()
10 ax.spines["top"].set_visible(False)
11 ax.spines["right"].set_visible(False)
12
13
14 colors = ['cyan', 'lavender']
15
16 bar = plt.bar(x_label, y, width=bar_width, align='center', label='value',
  color=colors, edgecolor='black')
17
18 plt.bar_label(bar, label_type='edge')
19
20 plt.xlabel('Model')
21 plt.ylabel('Marginal Gains (%)')
22
23 plt.savefig("Imgs/widhtAndDepth_CIFAR_Margin.png", bbox_inches='tight',
  dpi=400)
24
25 plt.show()

```

## FashionMNIST Data Processing

```

1 data_path = './data/FashionMinist/'
2 MNIST = datasets.FashionMNIST(data_path, train=True, download=True)
3 MNIST_val = datasets.FashionMNIST(data_path, train=False, download=True)

```

```

1 MNIST = datasets.FashionMNIST(data_path, train=True, download=False,
  transform=transforms.Compose([
2     transforms.ToTensor(),
3     transforms.Normalize((0.1307,), (0.3081,)),
4     transforms.RandomHorizontalFlip(p=0.5),
5     transforms.RandomVerticalFlip(p=0.5),
6     transforms.RandomRotation(30),
7 ]))
8
9 MNIST_val = datasets.FashionMNIST(data_path, train=False, download=False,
  transform=transforms.Compose([
10     transforms.ToTensor(),
11     transforms.Normalize((0.1307,), (0.3081,)),
12 ]))

```

## FashionMNIST training

```
1 train_loader = torch.utils.data.DataLoader(MNIST, batch_size=128,  
    pin_memory=True, num_workers=8, shuffle=True)  
2 valid_loader = torch.utils.data.DataLoader(MNIST_val, batch_size=128,  
    pin_memory=True, num_workers=8, shuffle=True)  
3  
4 model = resnet20_2().to(device=device)  
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,  
    weight_decay=1e-4, nesterov=True)  
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[50, 90,  
    110], gamma=0.1)  
7 loss_fn = nn.CrossEntropyLoss()
```

```
1 logger_20_2_mnist = training_loop(150, optimizer, scheduler, model,  
    loss_fn, train_loader, valid_loader, "./Log/Exp/MNIST/ResNet20_2",  
    "Resnet20-2")
```

```
1 torch.save(model.state_dict(), "./Models/MNIST/resnet20_2.pth")
```

```
1 test(valid_loader, model, loss_fn)
```

```
1 train_loader = torch.utils.data.DataLoader(cifar10, batch_size=128,  
    pin_memory=True, num_workers=8, shuffle=True)  
2 valid_loader = torch.utils.data.DataLoader(cifar10_val, batch_size=128,  
    pin_memory=True, num_workers=8, shuffle=True)  
3  
4 model = resnet20_4().to(device=device)  
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,  
    weight_decay=1e-4, nesterov=True)  
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[50, 90,  
    110], gamma=0.1)  
7 loss_fn = nn.CrossEntropyLoss()
```

```
1 logger_20_4_mnist = training_loop(150, optimizer, scheduler, model,  
    loss_fn, train_loader, valid_loader, "./Log/Exp/MNIST/ResNet20_4",  
    "Resnet20-4")
```

```
1 torch.save(model.state_dict(), "./Models/MNIST/resnet20_4.pth")
```

```
1 test(valid_loader, model, loss_fn)
```

```

1 train_loader = torch.utils.data.DataLoader(cifar10, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
2 valid_loader = torch.utils.data.DataLoader(cifar10_val, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
3
4 model = resnet20_8().to(device=device)
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
  weight_decay=1e-4, nesterov=True)
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[50, 90,
  110], gamma=0.1)
7 loss_fn = nn.CrossEntropyLoss()

```

```

1 logger_20_8_mnist = training_loop(150, optimizer, scheduler, model,
  loss_fn, train_loader, valid_loader, "./Log/Exp/MNIST/ResNet20_8",
  "Resnet20-8")

```

```

1 torch.save(model.state_dict(), "./Models/MNIST/resnet20_8.pth")

```

```

1 test(valid_loader, model, loss_fn)

```

```

1 logger_20_2_mnist.ReadLog()
2 logger_20_4_mnist.ReadLog()
3 logger_20_8_mnist.ReadLog()

```

```

1 train_loader = torch.utils.data.DataLoader(MNIST, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
2 valid_loader = torch.utils.data.DataLoader(MNIST_val, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
3
4 model = resnet20().to(device=device)
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
  weight_decay=1e-4, nesterov=True)
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[50, 90,
  110], gamma=0.1)
7 loss_fn = nn.CrossEntropyLoss()

```

```

1 logger_20_mnist = training_loop(150, optimizer, scheduler, model, loss_fn,
  train_loader, valid_loader, "./Log/Exp/MNIST/ResNet20", "Resnet20")

```

```

1 torch.save(model.state_dict(), 'ResNet20.pth')

```

```

1 test(valid_loader, model, loss_fn)

```

```

1 train_loader = torch.utils.data.DataLoader(MNIST, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
2 valid_loader = torch.utils.data.DataLoader(MNIST_val, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
3
4 model = resnet56().to(device=device)
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
  weight_decay=1e-4, nesterov=True)
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[50, 90,
  110], gamma=0.1)
7 loss_fn = nn.CrossEntropyLoss()

```

```

1 logger_56_mnist = training_loop(150, optimizer, scheduler, model, loss_fn,
  train_loader, valid_loader, "./Log/Exp/MNIST/ResNet56", "Resnet56")

```

```

1 torch.save(model.state_dict(), 'ResNet56.pth')

```

```

1 test(valid_loader, model, loss_fn)

```

```

1 train_loader = torch.utils.data.DataLoader(MNIST, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
2 valid_loader = torch.utils.data.DataLoader(MNIST_val, batch_size=128,
  pin_memory=True, num_workers=8, shuffle=True)
3
4 model = resnet110().to(device=device)
5 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
  weight_decay=1e-4, nesterov=True)
6 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[50, 90,
  110], gamma=0.1)
7 loss_fn = nn.CrossEntropyLoss()

```

```

1 logger_110_mnist = training_loop(150, optimizer, scheduler, model, loss_fn,
  train_loader, valid_loader, "./Log/Exp/MNIST/ResNet110", "Resnet110")

```

```

1 torch.save(model.state_dict(), 'ResNet110.pth')

```

```

1 test(valid_loader, model, loss_fn)

```

```

1 logger_110_mnist.ReadLog()
2 logger_56_mnist.ReadLog()
3 logger_20_mnist.ReadLog()

```

## Plotting graph

```

1 import matplotlib.pyplot as plt
2
3 plt.plot(logger_20_2.acc_val, label='ResNet-20-2', color='r')
4 plt.plot(logger_20_4.acc_val, label='ResNet-20-4', color='b')
5

```



```

6 plt.plot(logger_56_depth.acc_val, label='ResNet-56', ls='--', color='g')
7 plt.plot(logger_110_depth.acc_val, label='ResNet-110', ls='--',
  color='cyan')
8
9 plt.legend()
10 plt.xlim(0, 150)
11 plt.ylim(0.9, 0.95)
12 plt.grid()
13 plt.xlabel('Epoch')
14 plt.ylabel('Accuracy')
15 plt.savefig("./Imgs/MNIST_Cmp", dpi=400)
16 plt.show()

```

```

1 MNIST_time_deep = [8.1,14.8,25.3]
2 MNIST_time_width = [8.1,9.1,15.48,39.4]

```

```

1 fig, ax = plt.subplots()
2 plt.plot(CIFAR_Params_deep, MNIST_time_deep, 's-b', label='FashionMNIST
  Deeper Model')
3 plt.plot(CIFAR_Params_width, MNIST_time_width, '^-g', label='FashionMNIST
  Wider Model')
4
5 ax.annotate('ResNet-20-8', xy=(17, 39), xytext=(14, 32),
  arrowprops=dict(arrowstyle='->', color='dimgray'))
6 ax.annotate('ResNet-20-4', xy=(4.5, 15), xytext=(4, 8),
  arrowprops=dict(arrowstyle='->', color='dimgray'))
7 ax.annotate('ResNet-20-2', xy=(1.1, 8.8), xytext=(2, 4),
  arrowprops=dict(arrowstyle='->', color='dimgray'))
8 ax.annotate('ResNet-20', xy=(0.27, 7.8), xytext=(-0.4, 1.5),
  arrowprops=dict(arrowstyle='->', color='dimgray'))
9 ax.annotate('ResNet-56', xy=(0.853, 15), xytext=(0, 28),
  arrowprops=dict(arrowstyle='->', color='dimgray'))
10 ax.annotate('ResNet-110', xy=(1.8, 25.5), xytext=(3, 26),
  arrowprops=dict(arrowstyle='->', color='dimgray'))
11
12 plt.legend()
13 plt.grid()
14
15 plt.xlabel("# Params (Million)")
16 plt.ylabel("Training Time in each epoch")
17 plt.ylim(0, 40)
18 plt.savefig('FashionMNIST_wide_Vs_Deep', dpi=400)
19 plt.show()

```

```

1 y = [0.0771,0.0420,0.0251,0.0695,0.0404,0.0160]
2 x_label = ["ResNet-20", "ResNet-56", "ResNet-110", "ResNet-20-2", "ResNet-
  20-4", "ResNet-20-8"]
3
4 bar_width = 0.3
5
6 fig = plt.figure(figsize=(7, 6), dpi=100)
7 ax = plt.axes()
8 ax.spines["top"].set_visible(False)

```

```

9 ax.spines["right"].set_visible(False)
10
11
12 colors = ['cyan', 'lavender']
13
14 bar = plt.bar(x_label, y, width=bar_width, align='center', label='value',
15              color=colors, edgecolor='black')
16
17 plt.bar_label(bar, label_type='edge')
18
19 plt.xlabel('Model')
20 plt.ylabel('Marginal Gains (%)')
21
22 plt.savefig("Imgs/widhtAndDepth_MNIST_Margin.png", bbox_inches='tight',
23            dpi=400)
24
25 plt.show()

```

## GTSRB Data Processing

```

1 import zipfile
2
3 with zipfile.ZipFile('./data/GTSRB/archive_2.zip', 'r') as zipf:
4     zipf.extractall('./data/GTSRB')

```

```

1 data_path = "./data/GTSRB/"

```

```

1 meta_df = pd.read_csv('./data/GTSRB/Meta.csv')
2 train_df = pd.read_csv(data_path + 'Train.csv')
3 test_df = pd.read_csv(data_path + 'Test.csv')

```

```

1 train_data_path = os.path.join(data_path, 'Train')
2 test_data_path = os.path.join(data_path, 'Test')
3 meta_data_path = os.path.join(data_path, 'Meta')

```

```

1 transforms_train = transforms.Compose([transforms.RandomRotation(45),
2                                       transforms.ToTensor(),
3                                       transforms.Normalize([0.5, 0.5, 0.5],
4                                       [0.5, 0.5, 0.5]),
5                                       transforms.Resize((225, 225),
6                                       antialias=True),
7                                       transforms.RandomHorizontalFlip(p=0.5),
8                                       transforms.RandomVerticalFlip(p=0.5)])
9
10 transforms_val = transforms.Compose([transforms.ToTensor(),
11                                     transforms.Normalize([0.5, 0.5, 0.5],
12                                     [0.5, 0.5, 0.5]),
13                                     transforms.Resize((225, 225),
14                                     antialias=True)])

```

```

1 class TSignsDataset(Dataset):

```

```

2     def __init__(self, df, root_dir, transform=None):
3         self.df = df
4         self.root_dir = root_dir
5         self.transform = transform
6
7     def __len__(self):
8         return len(self.df)
9
10    def __getitem__(self, index):
11        image_path = os.path.join(self.root_dir, self.df.iloc[index, 7]) #the
column of paths in df is 7
12        image = Image.open(image_path)
13        y_class = torch.tensor(self.df.iloc[index, 6]) #the column of
ClsassId in df is 6
14
15        if self.transform:
16            image = self.transform(image)
17
18        return (image, y_class)

```

```

1 train_set = TSignsDataset(train_df, data_path, transform=transforms_train)
2 val_set = TSignsDataset(test_df, data_path, transform=transforms_val)

```

```

1 train_loader = DataLoader(dataset=train_set, batch_size=64, shuffle=True,
pin_memory=True)
2 val_loader = DataLoader(dataset=val_set, batch_size=64, shuffle=False,
pin_memory=True)

```

## ResNet Model Construction

```

1 class block(nn.Module):
2     def __init__(
3         self, in_channels, out_channels, identity_downsample=None,
stride=1):
4         super().__init__()
5         self.expansion = 4
6         self.conv1 =
nn.Conv2d(in_channels, out_channels, kernel_size=1, stride=1, padding=0, bias=False)
7         self.bn1 = nn.BatchNorm2d(out_channels)
8         self.conv2 =
nn.Conv2d(out_channels, out_channels, kernel_size=3, stride=stride, padding=1, bias=False,)
9         self.bn2 = nn.BatchNorm2d(out_channels)
10        self.conv3 = nn.Conv2d(out_channels, out_channels *
self.expansion, kernel_size=1, stride=1, padding=0, bias=False,)
11        self.bn3 = nn.BatchNorm2d(out_channels * self.expansion)
12        self.relu = nn.ReLU()
13        self.identity_downsample = identity_downsample
14        self.stride = stride
15
16    def forward(self, x):
17        identity = x.clone()

```

```

18         x = self.conv1(x)
19         x = self.bn1(x)
20         x = self.relu(x)
21         x = self.conv2(x)
22         x = self.bn2(x)
23         x = self.relu(x)
24         x = self.conv3(x)
25         x = self.bn3(x)
26
27         if self.identity_downsample is not None:
28             identity = self.identity_downsample(identity)
29
30         x += identity
31         x = self.relu(x)
32         return x
33

```

```

1 class ResNet(nn.Module):
2     def __init__(self, block, layers, image_channels, num_classes,
3         width_f):
4         super(ResNet, self).__init__()
5         self.in_channels = 64 * width_f
6         widths = [64 * width_f, 128 * width_f, 256 * width_f, 512 *
7         width_f]
8         self.conv1 = nn.Conv2d(image_channels, widths[0], kernel_size=7,
9         stride=2, padding=3, bias=False)
10        self.bn1 = nn.BatchNorm2d(widths[0])
11        self.relu = nn.ReLU()
12        self.maxpool = nn.MaxPool2d(kernel_size=3, stride=2, padding=1)
13
14        # Essentially the entire ResNet architecture are in these 4 lines
15        below
16        self.layer1 = self._make_layer(block, layers[0],
17        out_channels=widths[0], stride=1)
18        self.layer2 = self._make_layer(block, layers[1],
19        out_channels=widths[1], stride=2)
20        self.layer3 = self._make_layer(block, layers[2],
21        out_channels=widths[2], stride=2)
22        self.layer4 = self._make_layer(block, layers[3],
23        out_channels=widths[3], stride=2)
24
25        self.avgpool = nn.AdaptiveAvgPool2d((1, 1))
26        self.fc = nn.Linear(widths[3] * 4, num_classes)
27
28        def forward(self, x):
29            x = self.conv1(x)
30            x = self.bn1(x)
31            x = self.relu(x)
32            x = self.maxpool(x)
33            x = self.layer1(x)
34            x = self.layer2(x)
35            x = self.layer3(x)
36            x = self.layer4(x)
37
38            x = self.avgpool(x)

```

```

31         x = x.reshape(x.shape[0], -1)
32         x = self.fc(x)
33
34         return x
35
36     def _make_layer(self, block, num_residual_blocks, out_channels,
37 stride):
38         identity_downsample = None
39         layers = []
40
41         if stride != 1 or self.in_channels != out_channels * 4:
42             identity_downsample =
nn.Sequential(nn.Conv2d(self.in_channels, out_channels *
4, kernel_size=1, stride=stride, bias=False)
43 ,nn.BatchNorm2d(out_channels * 4))
44
45         layers.append(block(self.in_channels, out_channels,
46 identity_downsample, stride))
47
48         self.in_channels = out_channels * 4
49
50         for i in range(num_residual_blocks - 1):
51             layers.append(block(self.in_channels, out_channels))
52
53         return nn.Sequential(*layers)

```

```

1 def ResNet24_2(img_channel=3, num_classes=1000):
2     return ResNet(block, [3, 3, 3, 3], img_channel, num_classes, 2)
3
4 def ResNet24_4(img_channel=3, num_classes=1000):
5     return ResNet(block, [3, 3, 3, 3], img_channel, num_classes, 4)

```

```

1 def ResNet24(img_channel=3, num_classes=1000):
2     return ResNet(block, [3, 3, 3, 3], img_channel, num_classes, 1)
3
4 def ResNet50(img_channel=3, num_classes=1000):
5     return ResNet(block, [3, 4, 6, 3], img_channel, num_classes, 1)
6
7 def ResNet101(img_channel=3, num_classes=1000):
8     return ResNet(block, [3, 4, 23, 3], img_channel, num_classes, 1)

```

## GTSRB Training

```

1 model = ResNet24_2(img_channel=3, num_classes=43)
2 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
3 weight_decay=1e-4, nesterov=True)
4 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[20],
5 gamma=0.1)
6 loss_fn = nn.CrossEntropyLoss()

```

```
1 logger_24_2 = training_loop(40, optimizer, scheduler, model, loss_fn,  
    train_loader, val_loader, "./Log/GTSRB/ResNet_24_2", "Resnet24_2")
```

```
1 torch.save(model.state_dict(), "./Models/GTSRB/resnet24_2.pth")
```

```
1 test(val_loader, model, loss_fn)
```

```
1 model = ResNet24_4(img_channel=3, num_classes=43).to(device)  
2 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,  
    weight_decay=1e-4, nesterov=True)  
3 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[20],  
    gamma=0.1)  
4 loss_fn = nn.CrossEntropyLoss()
```

```
1 logger_24_4 = training_loop(40, optimizer, scheduler, model, loss_fn,  
    train_loader, val_loader, "./Log/GTSRB/ResNet_24_4", "Resnet24_4")
```

```
1 torch.save(model.state_dict(), "./Models/GTSRB/resnet24_4.pth")
```

```
1 test(val_loader, model, loss_fn)
```

```
1 model = ResNet24(img_channel=3, num_classes=43).to(device)  
2 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,  
    weight_decay=1e-4, nesterov=True)  
3 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[20],  
    gamma=0.1)  
4 loss_fn = nn.CrossEntropyLoss()
```

```
1 logger_24 = training_loop(40, optimizer, scheduler, model, loss_fn,  
    train_loader, val_loader, "./Log/GTSRB/ResNet_24", "Resnet24")
```

```
1 torch.save(model.state_dict(), "./Models/GTSRB/resnet24.pth")
```

```
1 test(val_loader, model, loss_fn)
```

```
1 model = ResNet50(img_channel=3, num_classes=43).to(device)  
2 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,  
    weight_decay=1e-4, nesterov=True)  
3 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[20],  
    gamma=0.1)  
4 loss_fn = nn.CrossEntropyLoss()
```

```
1 logger_24 = training_loop(40, optimizer, scheduler, model, loss_fn,  
    train_loader, val_loader, "./Log/GTSRB/ResNet_50", "Resnet50")
```

```
1 torch.save(model.state_dict(), "./Models/GTSRB/resnet50.pth")
```

```
1 test(val_loader, model, loss_fn)
```

```
1 model = ResNet101(img_channel=3, num_classes=43).to(device)
2 optimizer = optim.SGD(model.parameters(), lr=1e-1, momentum=0.9,
  weight_decay=1e-4, nesterov=True)
3 scheduler = optim.lr_scheduler.MultiStepLR(optimizer, milestones=[20],
  gamma=0.1)
4 loss_fn = nn.CrossEntropyLoss()
```

```
1 logger_101 = training_loop(40, optimizer, scheduler, model, loss_fn,
  train_loader, val_loader, "./Log/GTSRB/ResNet_101", "Resnet101")
```

```
1 torch.save(model.state_dict(), "./Models/GTSRB/resnet101.pth")
```

```
1 test(val_loader, model, loss_fn)
```

## Plotting

```
1 import matplotlib.pyplot as plt
2
3 plt.plot(logger_24_2.acc_val, label='ResNet-24-2', color='r')
4 plt.plot(logger_24_4.acc_val, label='ResNet-24-4', color='b')
5
6 plt.plot(logger_50.acc_val, label='ResNet-50', ls='--', color='g')
7 plt.plot(logger_101.acc_val, label='ResNet-101', ls='--', color='cyan')
8
9 plt.legend()
10 plt.xlim(0, 40)
11 plt.ylim(0.8, 0.98)
12 plt.grid()
13 plt.xlabel('Epoch')
14 plt.ylabel('Accuracy')
15 plt.savefig("./GTSRB_Cmp", dpi=400)
16 plt.show()
```

```
1 GTSRB_Model_depth_names = ["ResNet-20", "ResNet-56", "ResNet-110"]
2 GTSRB_Model_width_names = ["ResNet-20", "ResNet-20-2", "ResNet-20-4"]
3
4 GTSRB_time_depth = [201, 224, 314.7]
5 GTSRB_Params_depth = [19.96, 23.60, 42.58]
6
7 GTSRB_time_width = [201, 438, 1255.6]
8 GTSRB_Params_width = [19.96, 79.58, 317.75]
```

```
1 factors = [1, 2, 4]
2 plt.figure(figsize=(12, 5))
3
4 plt.subplot(121)
5 plt.plot(factors, GTSRB_Params_depth, 's-b', label='GTSRB Deeper Models')
6 plt.plot(factors, GTSRB_Params_width, '^-g', label='GTSRB Wider Models')
```

```
7 plt.legend()
8 plt.grid()
9 plt.xlabel("width/depth coefficient")
10 plt.ylabel("# Params (Million)")
11 #plt.ylim(0, 80)
12
13 plt.subplot(122)
14 plt.plot(factors, GTSRB_time_depth, 's-b', label='GTSRB Deeper Models')
15 plt.plot(factors, GTSRB_time_wdith, '^-g', label='GTSRB Wider Models')
16 plt.legend()
17 plt.grid()
18 plt.xlabel("width/depth coefficient")
19 plt.ylabel("Training Time in one epoch (s)")
20
21 plt.savefig('GTSRB_wide_Vs_Deep_params_time', dpi=400)
22 plt.show()
```