

# Weighted Residuals for Very Deep Networks

Falong Shen, Gang Zeng

Peking University

**Abstract.** Deep residual networks have recently shown appealing performance on many challenging computer vision tasks. However, the original residual structure still has some defects making it difficult to converge on very deep networks. In this paper, we introduce a weighted residual network to address the incompatibility between ReLU and element-wise addition and the deep network initialization problem. The weighted residual network is able to learn to combine residuals from different layers effectively and efficiently. The proposed models enjoy a consistent improvement over accuracy and convergence with increasing depths from 100+ layers to 1000+ layers. Besides, the weighted residual networks have little more computation and GPU memory burden than the original residual networks. The networks are optimized by projected stochastic gradient descent. Experiments on CIFAR-10 have shown that our algorithm has a *faster convergence speed* than the original residual networks and reaches a *high accuracy* at 95.3% with a 1192-layer model.

## 1 Introduction

The state-of-the-art model for image classification is built on inception and residual structure [1,2,3]. Lots of works devoted on residual networks are emerging recently [4,5,6,7]. Very deep convolutional networks [8,9], especially with residual units, have shown compelling accuracy and nice convergence behaviors on many challenging computer vision tasks [3,10,11]. Since *vanishing gradients* problem is well handled by batch normalization [12] and highway signal propagation [13], networks with 100+ layers are being developed and trained, even 1000+ layers structure still yields meaningful results when combined with adequate dropout as shown in [6]. He *et al.* [4] also introduced the pre-activation structure to allow the highway signal to be directly propagated through the very deep networks. However they seemed to harness the features with a larger dimension ( $4\times$ ) and adapted multiple  $1\times 1$  convolutional layers to substitute  $3\times 3$  convolutional layers for convergence with 1000+ layers.

A typical convolutional unit is composed of one convolutional layer, one batch normalization layer and one ReLU layer, all of which are performed sequentially [12]. For a residual unit, a central question is how to combine the residual signal and the highway signal, where element-wise addition was proposed in [3]. A natural idea is to perform addition after ReLU activation. However, this leads to a non-negative output from residual branch, which limits the representative ability of the residual unit meaning that it can only enhance the highway signal. He *et al.*

[3] firstly proposed to perform addition between batch normalization and ReLU. In [4], they further proposed to inverse the order of the three layers, performing batch normalization and ReLU before convolutional layers. The question is due to that ReLU activation can only generate positive value which is incompatible with element-wise addition in the residual unit.

As it is non-convex optimization to solve deep networks, an appropriate initialization is important for both faster convergence and a good local minima. The “xavier” [14] and “msra” [15] are popular used for deep networks initialization. However, for networks with depths beyond 100 layers, neither “xavier” nor “msra” works well. The paper of [3] proposed to “warm up” the network with small learning rate and then restore the learning rate to normal value. However, this hand-craft strategy is not that useful for very deep networks, where even a very low learning rate (0.00001) still is not enough to promise convergence and restoring the learning rate has a chance to get rid of the initial convergence [2].

Generally speaking, there are two defects embedded in the training of the original residual networks

- Incompatibility of ReLU and element-wise addition.
- difficultly for networks to converge with depths beyond 1000-layer using “msra” initializer.

The third point resides that a better mode to combine the residuals from different layers are necessary to train very deep networks. For very deep networks, not all layers are that important as 1000-layer networks often perform not much better than 100-layer networks. In fact, lots of layers serve as redundant information and very deep networks tend to over-fit on some tasks.

In this paper, we introduce the weighted residual networks, which learn to combine residuals from different layers effectively and efficiently. All the residual weights are initialized at zeros and optimized with a very small learning rate (0.001), which allows all the residual signals to gradually add to the highway signal. With a group of gradually growing-up residual weights, the 1192-layer residual networks converge even much faster than the 100-layer networks. Finally, the distribution of the learned residual weights is in a symmetry mode ranging in  $[-0.5, 0.5]$ , which implies the incompatibility of ReLU and element-wise addition can be appropriately handled. The networks are optimized by projected stochastic gradient descent with exactly the same training epochs to original residual networks.

We conduct experiments on CIFAR-10 [16] to verify the practicability of the weighted residual networks. Training with the weighted residual networks can converge much faster and reach a higher performance with negligible more computation and GPU memory cost than the original residual networks. The weighted residual networks with depths beyond 1000 layers still converge faster than shallower networks and enjoy a consistent improvement over accuracy with increasing depths from 100+ layers to 1000+ layers without resorting to any hand-craft strategy such as “warm up” [3]. After applying dropout on the residuals, our weighted residual networks reach a very high accuracy (95.3%) on

CIFAR-10 using a 1192-layer model with the same training epochs to the original residual networks (about 164 epochs, 64k iterations).

The contributions of our work presented in this paper have four folds:

- We propose the weighted residual networks, which learn to combine the residuals from each residual unit. The weighted residual networks converge much faster in the training stage and reach a higher accuracy than the original residual networks at little more computation and GPU memory cost.
- The incompatibility of ReLU and element-wise addition can be addressed appropriately by weighted residuals and we clear all the obstacles on the information highway to allow the highway signal to enjoy a unhindered propagation.
- The residuals are gradually added to the highway signal to make the training process more reliable, even networks with depths beyond 1000 layers can converge very fast without the “warm up” strategy.
- We modify the down-sampling step to make the spatial size and feature dimension consistent between highway signal and branched residual signal, without resorting to zero-padding or extra converting matrix.

The weighted residual networks are simple and easy to implement while having surprising practical effectiveness, which makes it particular useful for complicated residual networks in research community and real applications.

## 2 Related works

The residual networks have attracted lots of researchers and many works on it have appeared [4,5,6,7,17]. In the following paragraphs we will review some related works.

The residual networks simplify the highway networks [13] using identity skip connection, which allows information to flow directly and bypass complex layers. The residual networks consist of many residual units. There are two information flows in a residual unit. The highway signal goes through the identity skip connection and the branched residual signal is realized by Conv-BN-ReLU-Conv-BN. The two flows are combined at the end of a residual unit by element-wise addition and then it goes through a ReLU layer for activation. This simple structure is quite powerful and achieved a surprising performance on the imageNet challenge [1] with 150-layer networks [3].

In the original residual networks, the two flows are added up before ReLU activation for a numerical reason that ReLU can only produce *non-negative output*, which means the branched residual signal can only enhance the highway signal. However, intuitively it is not a natural solution as the branched residual signal needs to be “activated”.

He *et al.* [4] proposed to handle the incompatibility between ReLU and element-wise addition by re-arranging these layers to BN-ReLU-Conv-BN-ReLU-Conv and

named it “pre-activation” structure. When applying the “pre-activation” structure, special attention should be taken on the first and the last residual unit of the networks.

To train “residual” networks, it is natural to fit on the “residual” only, which means when the branched residual signal is not presented, the highway signal should still make meaningful results. Under this condition, the branched residual signal can focus on fitting the “residual” in a residual unit. Huang *et al.* [6] proposed a dropout residual network, which randomly drops the branched residual signal in each residual unit. Therefore, when the branched residual signal is presented in a residual unit, it can focus on fitting the “residual”. As this model can be treated as an ensemble of models with different depths, they named it “stochastic depth networks”.

In the convolutional networks, the depth and width are both important for a high performance in image classification [7,3]. The conv1-conv3-conv1 bottleneck structure which used a feature dimension  $4\times$  larger than conv3-conv3 reached a higher performance [4]. Zagoruyko *et al.* [7] used conv3-conv3 with feature dimension  $10\times$  larger and reached the highest performance on CIFAR-10 (4.10%). However, a larger feature dimension costs much more GPU memory and leads a shallower structure. There is a balance between depth and width.

In this paper, we mainly focus on models with depth beyond 100+ layers. We mean to explore how to train a very deep model effectively instead of tuning a more accurate model.

### 3 Weighted Residual Networks

Firstly we will give a brief introduction to the residual networks. The residual networks build the information highway by allowing earlier feature representation to flow unimpededly and directly to the following layers *without* any modification. A residual unit performs the following computation:

$$\mathbf{x}_{i+1} = \text{ReLU}(\mathbf{x}_i + \Delta L_i(\mathbf{x}_i, \theta_i)) \quad (1)$$

Here  $\mathbf{x}_i$  is the input highway signal to the  $i$ -th residual unit.  $\theta_i$  is the filter parameters for the residual unit and it is initialized by “msra”,  $\Delta L_i$  is the residual function, which is realized by a stack of two  $3 \times 3$  convolutional layers. Typically, one convolutional layer should be followed by one batch normalization layer to keep the signal with non-zero variance and one ReLU layer for non-linearity activation. The highway should be clean and unhindered. As it is shown in [4], obstacles on the highway, such as constant scaling and dropout, will make the optimization difficult. A typical residual unit is depicted in Figure 1.

The original residual networks stated above have two defects

**Incompatibility of ReLU and element-wise addition.** The highway signal and the residual signal which is produced by the residual function are combined by the element-wise addition. However, the element-wise addition is operated

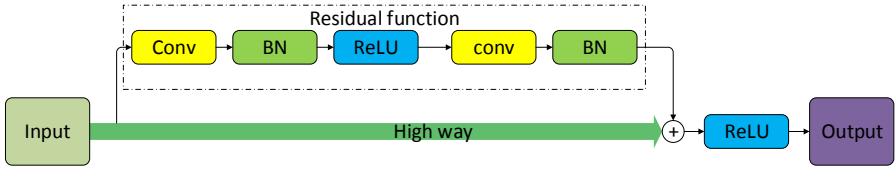


Fig. 1: Diagrammatic sketch of a residual unit. The residual function comprises of two  $3 \times 3$  convolution layers. Each Convolutional layer (Conv) is followed by a batch normalization layer (BN) and a ReLU layer (ReLU). The weights of convolutional layers are initialized by “msra”. The highway signal and the residual signal are combined by element-wise addition.

between the BN layer and the ReLU layer after the second Conv layer. This is mainly due to the ReLU activation function, which produces *non-negative output*. The output of ReLU operation is not compatible with element-wise addition as it can only enhance the highway signal, which limits the representability of the residual function, which is meant to take values in  $(-\infty, +\infty)$ . One can of course resort to designing other activation function which can take values in a larger range or a symmetry mode around zero.

**Initialization of very deep networks.** Very deep networks with depths beyond 1000 layers, even equipped with residual structure, batch normalization and ReLU, still do not converge in the training stage as shown in Figure 5. The paper of [3] proposed to “warm up” the network training with a little learning rate for several epochs and then restore it to the normal learning rate in order to facilitate the initial convergence. However, for deeper networks, even very little learning rate may not work well [2].

In very deep networks, the residuals from each block are added together and make the training hard to converge. One may want to zero all the residuals to start the training. However, *the weights of the convolutional layers* in residual functions should be initialized by “msra” which has little probability to produce all-zero weights.

### 3.1 Weighted Residuals

To address the incompatibility of ReLU and element-wise addition and to get a better initialization for very deep networks, we introduce the weighted residual networks. Formally in a weighted residual networks unit, the computation of the signal is

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \lambda_i \Delta L_i(\mathbf{x}_i, \theta_i), \lambda_i \in (-1, 1), \quad (2)$$

where  $\theta_i$  is the filter parameters and it is initialized by “msra”,  $\lambda_i$  is the weight scalar for the residual and it is initialized by zero with a very small learning rate. The ReLU activation is removed from the highway and  $\Delta L_i$  is realized by two Conv-BN-ReLUs.

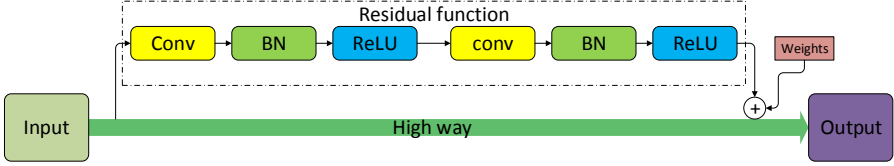


Fig. 2: Diagrammatic sketch of a weighted residual unit. We move the **ReLU** from highway to the branchway, which allows the highway signal to flow unobstructedly through the very deep networks. The residual signal is weighted by a scalar which is initialized by zero in the training stage. In our experiments, the overall convergence is promised when all the residuals are gradually added to the highway signal. The weight takes values in  $(-1, 1)$  to overcome the limitation of the **ReLU** activation function.

For any deep blocks, the feature representation  $\mathbf{x}_{i+k}$  in the  $(i+k)$ -th layer can be expressed as a summation of the input layer representation  $\mathbf{x}_i$  and a series of weighted residual functions,

$$\mathbf{x}_{i+k} = \mathbf{x}_i + \sum_{j=1}^k \lambda_{i+j} \Delta L_{i+j}(\mathbf{x}_{i+j}, \theta_{i+j}), \lambda_{i+j} \in (-1, 1). \quad (3)$$

In the back-propagation stage, the gradient of any layer does not vanish when filter parameter  $\theta_{i+j}$  is arbitrarily small. Note that the pre-activation structure proposed in [4] also has a similar property by converting the order of **Conv-BN-ReLU** to **BN-ReLU-Conv**.

In Figure 3 we visualize the distribution of the learned residual weights in a 1192-layer model. The residual weight values range around  $(-0.5, 0.5)$  in a symmetry mode, which means the branched residual signal has equal probability to enhance/weaken the highway signal, which means the incompatibility between **ReLU** and element-wise addition is appropriately addressed by the learned residual weights.

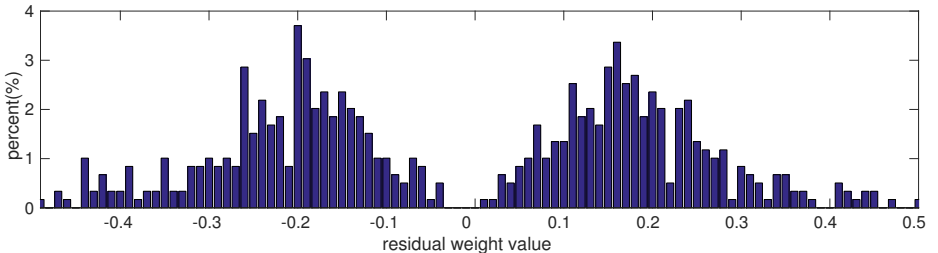


Fig. 3: Distribution of the learned residual weights in a 1192-layer model.

### 3.2 Modification to the structure.

At the beginning of a new block in the original residual networks, the high-way signal is down-sampled by a stride-2 convolution layer while the branched residual signal also need to be halved by a stride-2 convolution layer. When performing the element-wise addition, zeros-padding or convert matrix is necessary to make a matched feature dimension between the two signals.

In our networks as it is shown in Figure 4, we directly halve the feature size at the beginning and the following layers are performed as stated in the previous sections.

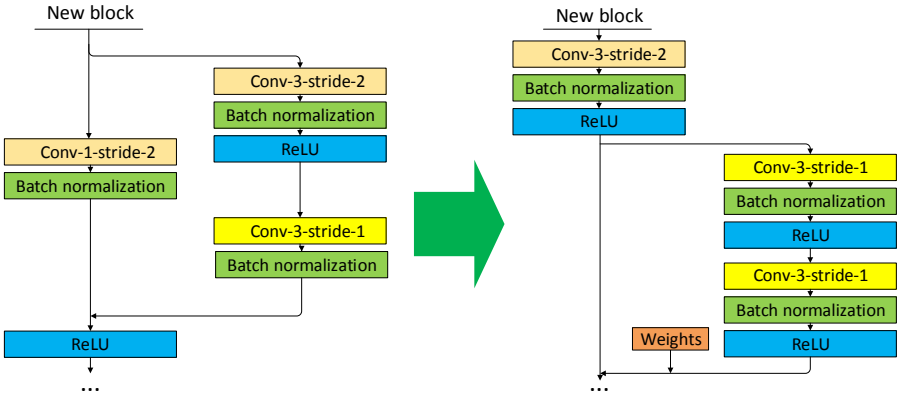


Fig. 4: At the beginning of each new block, the feature map is halved by a conv-3-stride-2 layer.

### 3.3 Optimization.

Given training images and its corresponding ground truth labels  $\{I_i, y_i\}$ , the loss function is the summation of the negative likelihood and the regularized term

$$-\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \log p_{y_i}(I_i, \theta, \lambda) + \frac{1}{2} \|\theta\|_2^2 + \frac{1}{2} \|\lambda\|_2^2, \quad (4)$$

*s.t. each element of  $\lambda \in (-1, 1)$ .*

where  $\theta$  is the network parameters which is initialized by “msra”,  $\lambda$  is the weight vector for the residuals and is initialized by all-zeros. We apply projected SGD to this typical constraint optimization problem. In the  $(t + 1)$ -th iteration, the updated  $\lambda_i^{t+1}$  is projected to the convex set  $S$

$$\lambda_i^{t+1} = P_S(\lambda_i^t + \Delta \lambda_i^t) \quad (5)$$

where the convex set  $S = (-1, 1)$  and  $\Delta\lambda_i^t$  is the gradient of the loss function in Equation 4 with regard to  $\lambda_i^t$ , which is effectively computed by back-propagation [18] in deep networks.

### 3.4 Implementation details.

**Dataset.** CIFAR-10 [16] is a dataset of color images all coming with the same size of  $32 \times 32$ , which consists of 50k training images and 10k testing image in 10 classes. We train our deep model on the *train* set and evaluate the finally trained models on the *test* set. We follow the same residual architecture as proposed in [3].

Our code is built on the open source deep learning framework Caffe [19]. We use a weight decay of 0.0001 and momentum of 0.9 with batch size of 128. The initial learning rate is 0.1 without “warm up” for any model. The initial learning rate for residual weights is set to 0.001. The filter parameters are initialized by “msra” [15]. The residual weights are set to all-zeros. We are not meant to push the state-of-the-art performance on CIFAR-10 so we follow the same training strategy as [3]. All the models are trained for 64k iterations and the learning rate is divided by 10 at 32k and 48k iterations. We also adapt the simple data augmentation as it is shown in [20]: 4 pixels are padded around the training images with zero-values and a translated or mirrored  $32 \times 32$  crop is fed into the networks. We do not have *val* set and the model at the end of training is used to perform on the *test* set. In the test stage, the original  $32 \times 32$  images are evaluated.

The network contains three blocks and the feature map is halved twice. There are totally  $6n + 4$  layers as it is shown in Table 1. We compare  $n = \{1, 3, 9, 18, 48, 100, 198\}$ , which leads to 10, 22, 58, 112, 292, 604 and 1192-layer networks.

Table 1: Our weighted residual network comprises of three blocks similar to [3]. The first block is started with a  $3 \times 3$  convolutional layer with stride 1. The latter two blocks are started with a  $3 \times 3$  convolutional layer with stride 2. There are  $6n + 4$  layers in the whole networks (with the final FC layer).

feature size	$32 \times 32$	$16 \times 16$	$8 \times 8$
filter number	16	32	64
layer number	$2n+1$	$2n+1$	$2n+1$

## 4 Experiments

In this section we present and analyze the experiment results on CIFAR-10 to demonstrate the effectiveness of the weighted residual networks.



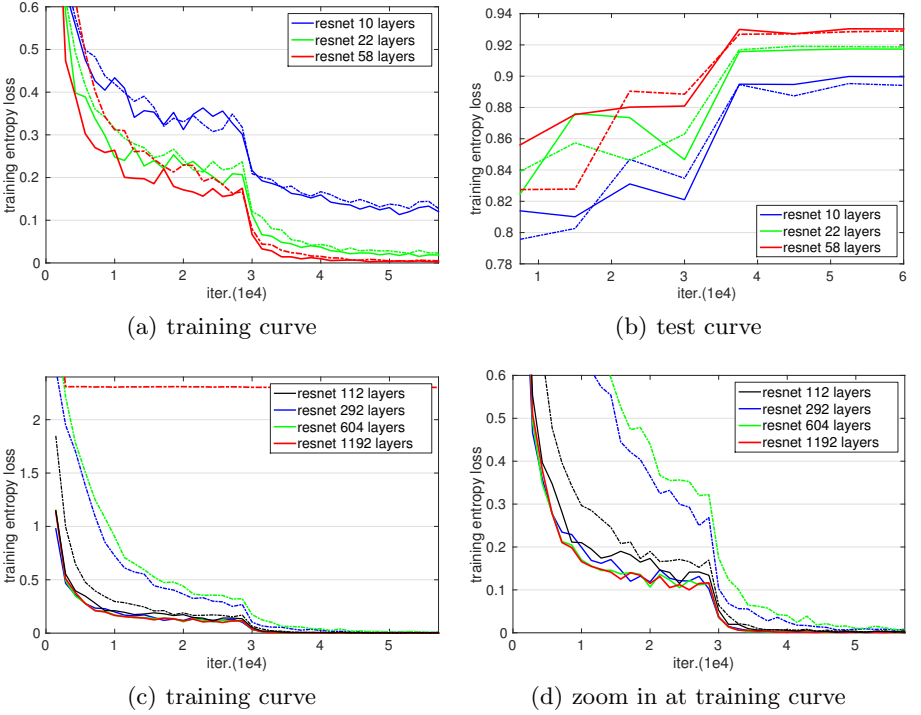


Fig. 5: Comparisons of the original residual networks and the weighted residual networks on the CIFAR-10. The **bold lines** denote the weighted residual networks and the **dashed lines** denote the original residual networks. The **top left** figure is the training entropy loss and the **top right** figure is the corresponding test accuracy on shallow networks. The **bottom left** figure records the training entropy loss for very deep networks and the **bottom right** figure is the zoomed version for more details.

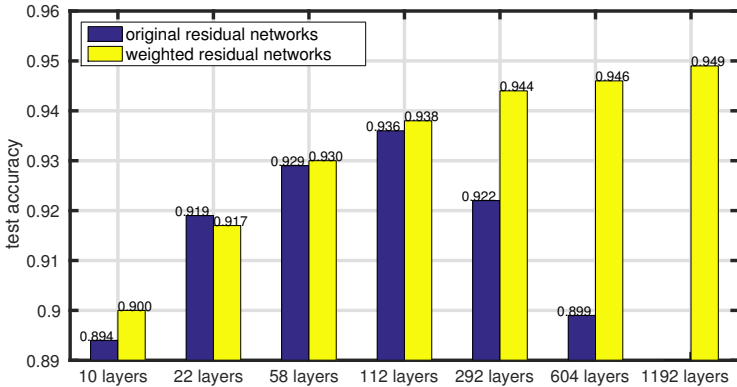


Fig. 6: Test accuracy on CIFAR-10. The original 1192-layer residual network fails to reach a meaningful result in the training stage and we do not report it.

## 4.1 Results

**Convergence.** Firstly we experiment on shallow networks (layer number  $< 100$ ). As it is shown in Figure 5(a) and Figure 5(b), both of the weighted residual networks and the original residual networks have very similar performance of convergence and final accuracy on shallow networks.

Then we conduct experiments on very deep networks (layer number  $> 100$ ). In Figure 5(c), the weighted residual network shows much better performance on convergence in the training stage. In fact, networks with depths beyond 1000 layers still converge faster than the 112-layer networks in Figure 5(d). As contrary, the original residual network does not converge well and the 1192-layer network even does not converge at all as we did not apply the “warm up” strategy. However, even equipped with “warm up”, the original 1192-layer residual network ends with over-fitting and reaches a worse performance than the 112-layer network as it is reported in [3].

**Accuracy.** The overall test accuracy of deep networks on CIFAR-10 is reported in Figure 6. The **blue histograms** denote the performance of the original networks. The accuracy decreases after the layer number is larger than 100. However, for the weighted residual networks, which are denoted as **yellow histograms**, the performance enjoys a consistent improvement with the increasing depths from 10+ layers to 1000+ layers.

The weighted residual networks can always *converge faster* and reach a *higher performance* when there are *more layers* throughout our experiments.

## 4.2 Comparison with the State-of-the-art

In this subsection we compare the weighted residual networks (**WRResNet**) with other recently proposed models. Mainly there are two kinds of models, first of which focus on enlarging the feature dimension and we call them **wide models**, the second of which focus on depths and we call them **deep models**. Note that 1001-layer Pre-activation [4] is both deep (1000+ layers) and wide ( $4\times$  feature dimension) model. The results are presented in Table 2. All these models, except for Highway [13], share similar structures with ResNet [3], including three feature blocks.

Pre-activation [4] adapted a **conv1-conv3-conv1** bottle-neck structure and enlarged the feature dimension by  $4\times$ . Apparently a  $4\times$  wider model enjoys a higher performance but costs more GPU memory. As the GPU memory (12GB for one GTX TITAN X) resource is limited, it is important to tune the model width and depth economically for a very accurate model. WideDim [7] and RiR [5] are two other methods to enlarge the feature dimension for higher accuracy. A clear tendency is that a wider feature is better for higher performance. WideDim adapted a  $10\times$  feature dimension and reached a very high performance (95.8%) on CIFAR-10. Dropout [6] realized stochastic depth networks by applying the dropout operation on the residual signal at exactly the same GPU memory cost. The only defect resides that it needs much more epochs (about  $2\times$ ) to converge at a good performance.

Table 2: Test accuracy(%) on CIFAR-10.

Type	method	depth	feature dim.	epochs	acc.(%)
wide models	Highway	32	-	-	91.2
	Pre-activation	164	64-128-256	200	94.5
	WideDim	28	160-320-640	200	95.8
	RiR	18	96-192-384	82	95.0
deep models	ResNet	1202	16-32-64	164	92.1
	Pre-activation	1001	64-128-256	200	95.4
	Dropout	1202	16-32-64	300	95.1
	<b>WResNet</b>	1192	16-32-64	164	94.9
	<b>WResNet-d</b>	1192	16-32-64	164	95.3

The weighted residual networks make very deep networks training converge faster and reach a good performance while bringing little more computation and GPU memory burden. As time and GPU resource is limited, we have not tuned the model width (feature dim.) or more training epochs and we are meant to explore the effectiveness of the weighted residuals in *training very deep models*. Yet with shorter feature dim., the weighted residual networks still perform much better than the original residual networks and reach a quite meaningful accuracy as shown in Table 2.

We further apply dropout on the residuals with  $dropout\_ratio = \{0.2, 0.4, 0.6\}$  for three blocks as proposed by [6]. The performance of this model is named as **WResNet-d**. With only about half training epochs of [6], the weighted residual networks with dropout reach a relative very high performance (95.3%).

### 4.3 Analysis

We provide more insights into the weighted residual networks by presenting more details information of results in this subsection.

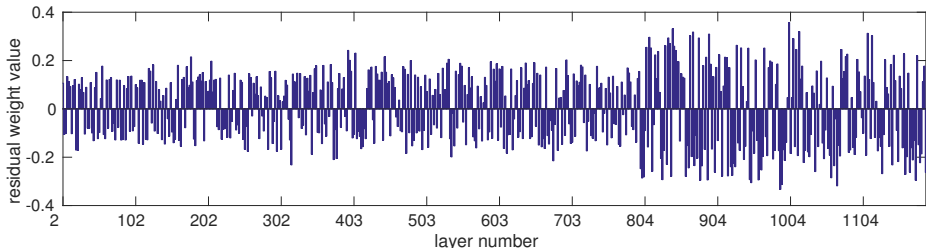
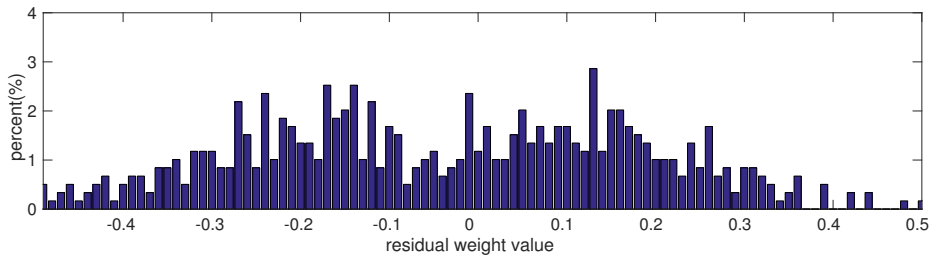
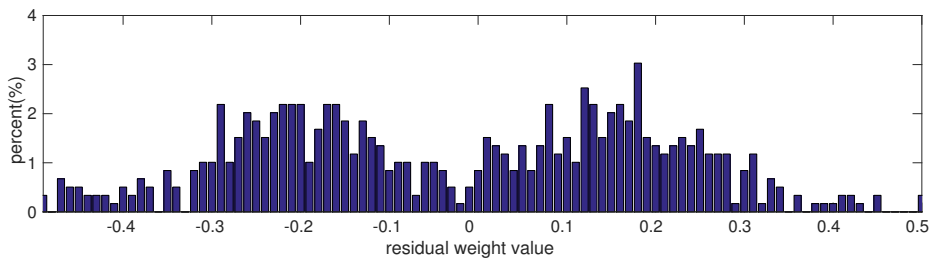


Fig. 7: The residual weight values from a trained 1192-layer model on CIFAR-10.

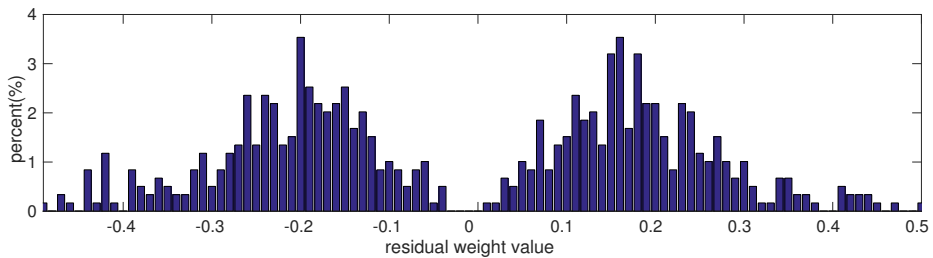
The initial learning rate for the residual weights is set to 0.001 for all models and the residual weights are initialized with zeros. Figure 7 shows the learned



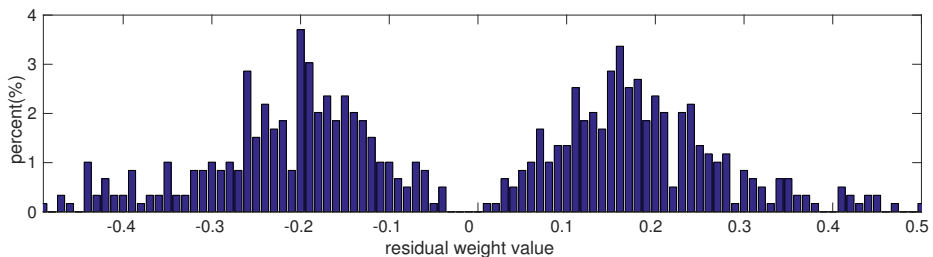
(a) 8k iterations



(b) 16k iterations



(c) 32k iterations



(d) 64k iterations

Fig. 8: Evolution of the distribution of the residual weight values in the training stage of a 1192-layer weighted residual network on CIFAR-10.

residual weight values in each element-wise addition layer in a 1192-layer model. It comprises two parts divided by a visible sharp boundary around the 800-layer and the latter residuals have larger weights. It may imply the residuals from the later layers are more important than earlier layers on the final decisions. We will explore this phenomenon in the future work.

We have also plotted the evolution history of the distribution of the residual weight values as show in Figure 8. At the 8k iteration, the distribution is relative uniform. As more and more training iterations, the distribution begins to concentrate around two peaks. In the 64k iteration, most of the residual weight values are around 0.2 and  $-0.2$  in a symmetry mode indicating that the branched residual signals have equal probability to enhance/weaken the highway signals, which verifies our hypothesis. Therefore the learned residual weights can solve the incompatibility between ReLU activation and element-wise addition appropriately.

## 5 Conclusion

The original residual networks have two defects, 1) Incompatibility between ReLU and element-wise addition. 2) Difficulty for networks to converge with depths beyond 1000-layer using “msra” initializer. In this paper we introduce the weighted residual networks to make very deep residual networks *converge faster* and reach a *higher performance with little more computation and GPU memory burden* than the original residual networks. All the residuals are added to the highway signal gradually by the learned slowly growing-up weights to promise convergence. Experiments on CIFAR-10 have demonstrated the effectiveness of the weighted residual networks for very deep models. It enjoys a consistent improvements over accuracy and convergence with the increasing depths from 100+ layers to 1000+ layers. The weighted residual networks are simple and easy to implement while having surprising practical effectiveness, which makes it particular useful for complicated residual networks in research community and real applications.

## References

1. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *International Journal of Computer Vision* **115** (2015) 211–252
2. Szegedy, C., Ioffe, S., Vanhoucke, V.: Inception-v4, inception-resnet and the impact of residual connections on learning. *arXiv preprint arXiv:1602.07261* (2016)
3. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385* (2015)
4. He, K., Zhang, X., Ren, S., Sun, J.: Identity mappings in deep residual networks. *arXiv preprint arXiv:1603.05027* (2016)
5. Targ, S., Almeida, D., Lyman, K.: Resnet in resnet: Generalizing residual architectures. *arXiv preprint arXiv:1603.08029* (2016)

6. Huang, G., Sun, Y., Liu, Z., Sedra, D., Weinberger, K.: Deep networks with stochastic depth. arXiv preprint arXiv:1603.09382 (2016)
7. Sergey Zagoruyko, N.K.: Deep networks with stochastic depth. arXiv preprint arXiv:1605.07146 (2016)
8. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems. (2012) 1097–1105
9. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
10. Dai, J., He, K., Sun, J.: Instance-aware semantic segmentation via multi-task network cascades. arXiv preprint arXiv:1512.04412 (2015)
11. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: Advances in Neural Information Processing Systems. (2015) 91–99
12. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. arXiv preprint arXiv:1502.03167 (2015)
13. Srivastava, R.K., Greff, K., Schmidhuber, J.: Training very deep networks. In: Advances in Neural Information Processing Systems. (2015) 2368–2376
14. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: International conference on artificial intelligence and statistics. (2010) 249–256
15. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: Proceedings of the IEEE International Conference on Computer Vision. (2015) 1026–1034
16. Krizhevsky, A.: Learning multiple layers of features from tiny images (2009)
17. Shah, A., Kadam, E., Shah, H., Shinde, S.: Deep residual networks with exponential linear unit. arXiv preprint arXiv:1604.04112 (2016)
18. Hecht-Nielsen, R.: Theory of the backpropagation neural network. In: Neural Networks, 1989. IJCNN., International Joint Conference on, IEEE (1989) 593–605
19. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding. In: Proceedings of the ACM International Conference on Multimedia, ACM (2014) 675–678
20. Lee, C.Y., Xie, S., Gallagher, P., Zhang, Z., Tu, Z.: Deeply-supervised nets. arXiv preprint arXiv:1409.5185 (2014)