

Altering Backward Pass Gradients to Improve Convergence

Bishshoy Das*, Milton Mondal*, Brejesh Lall, Shiv Dutt Joshi, and Sumantra Dutta Roy

Indian Institute of Technology Delhi, Hauz Khas, New Delhi - 110016, India

{bishshoy.das, milton.mondal, brejesh, sdjoshi, sumantra}@ee.iitd.ac.in

Abstract

In standard neural network training, the gradients in the backward pass are determined by the forward pass. As a result, the two stages are coupled. This is how most neural networks are trained hitherto. Gradient modification in the backward pass has seldom been studied in the literature. In this paper we explore decoupled training, where we alter the gradients in the backward pass. We propose a simple yet powerful method called PowerGrad Transform (PGT), that alters the gradients before the weight update in the backward pass and significantly enhances the predictive performance of a convolutional neural network. PGT trains networks to arrive at a better optima at convergence. It is computationally efficient, and adds no additional cost to either memory or compute, but results in improved final accuracies on both the training and test datasets. PowerGrad Transform is easy to integrate into existing training routines, requiring just a few lines of code. With decoupled training, our method improves baseline accuracies for ResNet-50 by 0.73%, for SE-ResNet-50 by 0.66% and by more than 1.0% for the non-normalized ResNet-18 network on the ImageNet classification task.

1. Introduction

Backpropagation is traditionally used for training deep neural networks [8]. Gradients are computed using basic calculus principles to adjust the weights during backpropagation [7]. However, decoupling the backward pass from the forward pass by modifying the gradients to improve training efficiency and final convergent accuracy has hardly been explored. In this paper we explore the landscape of decoupling the forward pass from the backward pass by altering the gradients and subsequently updating the network’s parameters with the modified gradients. There are several ways to achieve gradient modification in the backward pass. We discuss a few techniques in Fig. 1.

Type 0: No modification: In this method, we calculate the gradients using the standard calculus rules and use the

chain rule to calculate the gradients of the rest of the network’s parameters, also known as backpropagation as portrayed in Fig. 1(a). The network is then updated with the gradient descent equation:

$$W_i^{t+1} = W_i^t - \lambda \nabla_{W_i}(L) \quad i = D, D-1, \dots, 1 \quad (1)$$

Type I: Independent gradient modification at multiple points: Here the gradients are first computed using standard procedure and then individually altered as shown in Fig. 1(b). Gradient clipping [11] and Adaptive gradient clipping [1] are examples of such modifications. It can be described as:

$$W_i^{t+1} = W_i^t - \lambda f(\nabla_{W_i}(L)) \quad i = D, D-1, \dots, 1 \quad (2)$$

where the gradients $\nabla_{W_i}(L)$ are transformed using the transformation function ‘ f ’ before the weight update.

Type II: Gradient modification at a point very early in the backward graph: In this type of modification, the gradient is altered at a very early stage in the backward computation graph and then all subsequent gradients are generated using the values obtained with the modified gradients (Fig. 1(c)). Because of the chain rule, network parameters whose gradients are connected to the point of alteration in the computation graph also gets subsequently altered. It can be described as:

$$W_D^{t+1} = W_D^t - \lambda f(\nabla_{W_D}(L)) \quad (3)$$

$$W_i^{t+1} = W_i^t - \lambda \nabla_{W_i}(L)^* \quad i = D-1, \dots, 1 \quad (4)$$

where the gradient $\nabla_{W_D}(L)$ is first transformed using the transformation function ‘ f ’ and then this transformed gradient is propagated through the rest of the backward graph. All other gradient vectors $\nabla_{W_i}(L)^*$ are computed as it is, but because of the early injection of the transformed gradient $\nabla_{W_D}(L)$, all other gradient vectors that are connected to the transformed gradient through the chain rule ($\nabla_{W_i}(L)^*$, $i=D-1, \dots, 1$), gets subsequently altered.

Type I is computationally more expensive than type II as it requires altering the gradients of each and every parameter individually. We propose PowerGrad Transform (PGT), a type II method that modifies the gradients at the softmax layer. The following are the major contributions of this paper:

*Equal contribution.

DOI reference number: 10.18293/SEKE2023-177

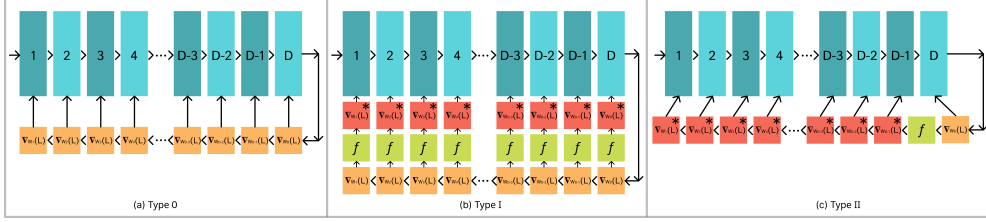


Figure 1: Different ways of altering gradients in the backward pass. Blue blocks denote different layers. Orange blocks indicate the backward graph with unmodified gradients. Green blocks represent transformation functions, while red blocks indicate transformed gradients.

1. We propose PowerGrad Transform (PGT), which decouples the backward and the forward passes of neural network training. PGT alters the gradients in the backward pass before the update step leading to accelerated training and a significant boost in the network’s predictive performance.
2. We perform theoretical analysis of the properties of the PGT and also show that in non-BN networks, PGT can be used to increase the network’s convergence rate and improve the final accuracy.
3. We find that PGT improves the performance for a variety of models (non-BN and BN ResNets, SE-ResNets) using the ImageNet dataset. We empirically conclude that PGT helps a network to improve by locating a more optimum convergence point.

2. Related Works

Gradient Clipping [11] is a gradient modification method that involves clipping/altering the gradients with respect to a predefined threshold value during backward propagation through the network and updating the weights using the clipped gradients [15][13]. By rescaling the gradients, the weight updates are likewise rescaled, significantly reducing the risk of an overflow or underflow [10]. GC can be used for training networks without batch-normalization. At larger batch sizes, the clipping threshold in GC becomes highly sensitive and requires extensive finetuning for various models, batch sizes, and learning rates. Adaptive Gradient Clipping [1] is developed to further enhance backward pass gradients than what is performed by GC. It takes into account the fact that the ratio of the gradient norm to the weight norm can provide an indication of the expected change in a single step of optimization. Label smoothing, introduced by Szegedy et al. [14], utilizes smoothing of the ground truth labels as a method to impose regularization on the logits and the weights of the neural network. AGC performs better than GC in non-normalized networks. However, we show that PGT outperforms both in networks such as ResNets.

Knowledge distillation (KD) [5] is a process in which two networks are trained with hard and soft labels alternatively. Variants of knowledge distillation include self-distillation [16], channel distillation [3]. Even though both PGT and KD require probability manipulation, the key difference is that in the latter the transformation is applied in the forward pass, while PGT is a backward pass modification only. PGT differs from self-knowledge distillation as it neither introduces any additional sub-modules nor creates different ensembles to improve the performance of the model. PGT follows the standard neural network training mechanism with modified gradients.

3. PowerGrad Transform

A neural network with parameters W generates C logits denoted by z for every input vector x . z is given as $z = Wx$. Then a set of probability values p_i are generated from the logits using a softmax function which is defined as $p_i = \frac{e^{z_i}}{\sum_{j=1}^C e^{z_j}}$. p_i and z_i represent the predicted probability values and the logits for the i^{th} class respectively. If the loss function is cross-entropy loss, then the value of the loss is given as $L = -\sum_{i=1}^C q_i \log(p_i)$ where q_i is the ground truth label of the i^{th} class for a particular training example. By standard gradient update rule, we can calculate the gradient of the loss with respect to the logits $\frac{\partial L}{\partial z_i} = p_i - q_i$.

The PowerGrad Transform technique is now described. We introduce a hyperparameter α , which takes a value between $[0, 1]$ and regulates the degree of gradient modification. The PowerGrad Transform method modifies the predicted probability values in the backward pass as follows:

$$p'_i = \frac{p_i^\alpha}{\sum_{j=1}^C p_j^\alpha} \quad i = 1, \dots, C \quad 0 \leq \alpha \leq 1 \quad (5)$$

The above transformation changes the gradient of the loss with respect to the logits as follows:

$$\widehat{\frac{\partial L}{\partial z_i}} = p'_i - q_i \quad (6)$$

The rest of the backward pass proceeds as usual.

3.1. Properties of the PowerGrad transformation

We use the same setup as described in section 3. To explore the properties of PGT, we start by investigating the effect of the transform on the softmax probabilities.

Lemma 1. For any arbitrary probability distribution P with probability values given by p_i for $i = 1, \dots, C$, the corresponding transformed probability values p'_i given by [Eq. 5] has a threshold $\left(\sum_{j=1}^C p_j^\alpha\right)^{\frac{1}{\alpha-1}}$ and

$$\begin{aligned} p'_i &\geq p_i, \text{ if } p_i \leq \left(\sum_{j=1}^C p_j^\alpha\right)^{\frac{1}{\alpha-1}} \\ p'_i &< p_i, \text{ if } p_i > \left(\sum_{j=1}^C p_j^\alpha\right)^{\frac{1}{\alpha-1}} \end{aligned} \quad (7)$$

We call this threshold, the *stationary threshold*. The stationary threshold is that value of p_i that does not change after the transformation. However, when p_i is greater than the *stationary threshold*, $p'_i < p_i$.

Proposition 1. At $\alpha = 0$, the stationary threshold equals $1/C$ and all values of the transformed distribution p'_i reduces to the uniform distribution for $i = 1, \dots, C$.

Proof. From Eq. (7), we see that the stationary threshold at $\alpha = 0$ is $1/C$. Also, following from the definition of the transformed probabilities (Eq. 5) we conclude that if $\alpha = 0$, then all values of p'_i are $1/C$. Therefore the transformed distribution at $\alpha = 0$ is a uniform distribution.

Theorem 1. For any arbitrary probability distribution P with probability values p_i for $i = 1, \dots, C$, the stationary threshold of the transformed distribution P' with probability values $p'_i = \frac{p_i^\alpha}{\sum_{j=1}^C p_j^\alpha}$, $0 \leq \alpha \leq 1$ is a monotonically non-decreasing function with respect to α .

Proof. To prove monotonicity, we first compute the gradient of the stationary threshold with respect to the variable in concern, α .

$$\begin{aligned} &\frac{\partial}{\partial \alpha} \left(\sum_{j=1}^C p_j^\alpha\right)^{\frac{1}{\alpha-1}} \\ &= \left(\sum_{j=1}^C p_j^\alpha\right)^{\frac{1}{\alpha-1}} \left(\frac{\sum_{j=1}^C p_j^\alpha \log(p_j)}{(\alpha-1)\sum_{j=1}^C p_j^\alpha} - \frac{\log\left(\sum_{j=1}^C p_j^\alpha\right)}{(\alpha-1)^2}\right) \\ &= \frac{1}{\alpha(\alpha-1)^2} \left(\sum_{j=1}^C p_j^\alpha\right)^{\frac{1}{\alpha-1}} \times \\ &\quad \left(\frac{(\alpha-1)\sum_{j=1}^C p_j^\alpha \log(p_j)}{\sum_{j=1}^C p_j^\alpha} - \alpha \log\left(\sum_{j=1}^C p_j^\alpha\right)\right) \end{aligned} \quad (8)$$

If a_1, \dots, a_n and b_1, \dots, b_n are non-negative numbers, then using the log sum inequality,

$$\sum_{j=1}^n a_j \log\left(\frac{a_j}{b_j}\right) \geq \left(\sum_{j=1}^n a_j\right) \log\left(\frac{\sum_{j=1}^n a_j}{\sum_{j=1}^n b_j}\right).$$

Setting $a_j = p_j^\alpha$ and $b_j = 1$, we get the following lower bound

$$\sum_{j=1}^C p_j^\alpha \log(p_j^\alpha) \geq \left(\sum_{j=1}^C p_j^\alpha\right) \log\left(\frac{1}{C} \sum_{j=1}^C p_j^\alpha\right) \quad (9)$$

Substituting (9) in (8), we get:

$$\begin{aligned} \frac{\partial}{\partial \alpha} \left(\sum_{j=1}^C p_j^\alpha\right)^{\frac{1}{\alpha-1}} &\geq \frac{1}{\alpha(\alpha-1)^2} \left(\sum_{j=1}^C p_j^\alpha\right)^{\frac{1}{\alpha-1}} \times \\ &\quad \left((1-\alpha)\log(C) - \log\left(\sum_{j=1}^C p_j^\alpha\right)\right) \end{aligned} \quad (10)$$

p^α is concave, and so by Jensen's inequality we get the following upper bound for the second term:

$$\left(\frac{1}{C} \sum_{j=1}^C p_j^\alpha\right)^\alpha \geq \frac{1}{C} \sum_{j=1}^C p_j^\alpha \quad (11)$$

$$\Rightarrow \log\left(\sum_{j=1}^C p_j^\alpha\right) \leq (1-\alpha)\log(C) \quad (12)$$

Substituting (12) in (10),

$$\frac{\partial}{\partial \alpha} \left(\sum_{j=1}^C p_j^\alpha\right)^{\frac{1}{\alpha-1}} \geq 0 \quad (13)$$

We conclude that the stationary threshold is a monotonic non-decreasing function with respect to α . Also the derivative of PGT function with respect to the true probabilities is non-negative which in turn means that the transformation is an order-preserving map. All values greater than the threshold move towards the threshold after transformation and all values below the threshold also move towards the threshold, and the threshold itself moves monotonically towards $1/C$ as α is decreased from 1 to 0. This concludes that the transformation smooths the original distribution.

4. Experiments

We perform experiments on different variants ResNets using the ImageNet-1K dataset [2]. All models are trained on four V100 GPUs with a batch size of 1024. We utilize a common set of hyperparameters for all experiments, which are as follows: 100 epoch budget, 5 epochs linear

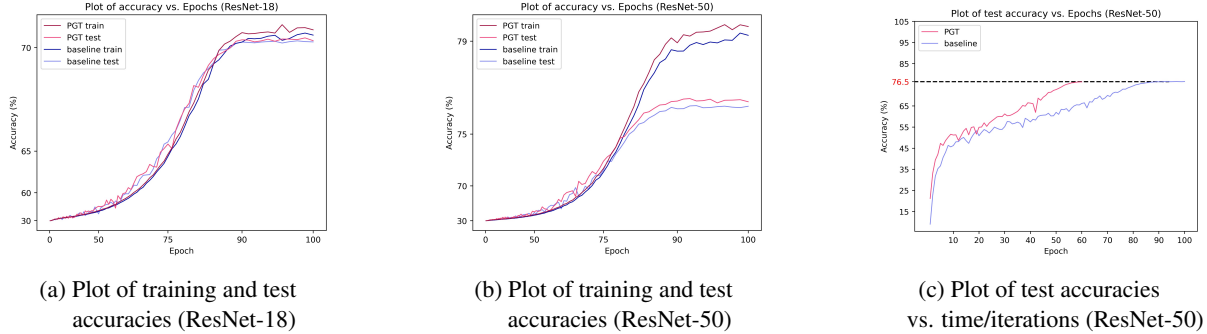


Figure 2: Log-log plots of training and test accuracies and comparison with baseline of batch-normalized variants: (a) ResNet-18 ($\alpha = 0.25$), (b) ResNet-50 ($\alpha = 0.3$). (c) Training speed comparison between PGT (60 epochs) and baseline (100 epochs). They both converge to the same test accuracy (76.5%) on ImageNet-1K. PGT’s accelerated training saves 40% of the epoch budget.

warmup phase beginning with a learning rate of 4×10^{-4} and ending with a peak learning rate of 0.4, a momentum of 0.9 and weight decay of 5×10^{-4} , the SGD Nesterov optimizer and mixed precision. In our studies, we employ either a step scheduler (dividing the learning rate by 10 at the 30th, 60th, and 90th epochs) or a cosine decay scheduler [9]. We find $\alpha = 0.25$ and $\alpha = 0.05$ to be good choices for ResNet-18 and ResNet-50, though larger values such as $\alpha = 0.3$ also have good performance as well. We use $\alpha = 0.3$ for Squeeze-and-Excitation variant of ResNet-50 i.e. SE-ResNet-50 [6]. The experimental results are shown in Table 1. ¹

With consistent improvements in training and test accuracies across all cases, we conclude PGT helps networks learn better representations and arrive at better optimas during convergence. Per epoch training and test accuracy plots of ResNet-18 and ResNet-50 (both with and without PGT) are shown in Fig. 2(a,b). Practitioners can also choose to accelerate training and save as much as 40% of the epoch budget [Fig. 2(c)].

4.1. Empirical studies on networks without Batch Normalization

We examine issues that occur in non-normalized networks (networks without BN layers). We use ResNet-18 [4] as the foundation model trained on ImageNet-1K [12]. Deeper networks such as ResNet-34 and ResNet-50 are impossible to train without Batch Normalization courtesy of the increased depth and so we solely focus on ResNet-18. Throughout the training process, we monitor variations in the the per-filter L2-norm of each layer’s weights. In Fig. 3(a), some filters of layer 11 achieve a norm of zero during training. We refer to this event as ‘Zeroing Out’ (Fig. 4), and it occurs when one of the channels (or filters) of a

Table 1: Results and comparison table for networks trained on ImageNet-1K. Best training and test accuracies are highlighted in red and blue respectively. Accuracy differences are highlighted in yellow.

Model	Scheduler	Method	PGT (α)	Train Acc.(%)	Train Diff(%)	Test Acc.(%)	Test Diff(%)
SE-ResNet-50	Cosine	Baseline	-	81.5		77.218	
		PGT	0.3	82.47	+0.97	77.952	+0.734
ResNet-50	Cosine	Baseline	-	79.18		76.56	
		PGT	0.05	79.68	+0.5	77.216	+0.656
ResNet-50	Step	Baseline	-	78.99		75.97	
		PGT	0.05	79.56	+0.57	76.494	+0.524
ResNet-101	Cosine	Baseline	-	82.29		77.896	
		PGT	0.3	83.1	+0.81	78.258	+0.362
SE-ResNet-18	Cosine	Baseline	-	71.42		71.09	
		PGT	0.25	71.6	+0.18	71.436	+0.346
ResNet-18	Step	Baseline	-	69.95		69.704	
		PGT	0.25	70.3	+0.35	69.844	+0.14
ResNet-18	Cosine	Baseline	-	70.38		70.208	
		PGT	0.25	70.53	+0.15	70.298	+0.09

weight tensor gets fully filled with zeros and such filters do not contribute at all to determine the input-output relationship of a dataset, as the feature tensor it produces is also filled with zeros for the corresponding filter. When a filter once zeroes out, it does not recover with further training, as all gradients that it receives in future iterations are all zeros.

Fig. 3(c) is the plot of the final conv layer’s filters (layer 19) and the features output of final global average pooling (GAP) layer respectively. We observe a number of filters and features in the final layer completely zeroing out. Gradient modification methods such as PGT can alleviate the zeroing out phenomena as we observe that the number of zeroed out filters has considerably reduced [Fig. 3(b, d)]. The feature vector after the GAP layer [Fig. 3(e)] directly interfaces with the fully connected layer. Therefore any zeroed out features leads to permanent information loss, as it does not contribute to the learning of decision boundaries in the fully connected layer. Also, zeroed out weights ten-

¹Reproducible code, training recipes, checkpoints and training logs are provided at: <https://github.com/skalian/power-grad-transform>

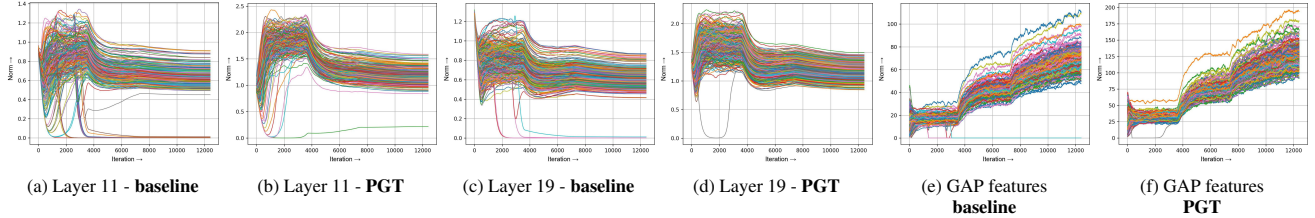


Figure 3: Norm vs iter. plots demonstrating layer characteristics (the zero out phenomena) and the efficacy of PGT over baseline. Each colour represents a different filter or feature vector of a particular layer of a non-BN variant of ResNet-18.

sors lead to zeroed out gradients hence stopping training for all subsequent iterations leading to a collapse in training for the affected layers. With large batch sizes, it is possible that an entire layer zeroes out as shown in the figure below. The final feature tensor [Fig. 3(f)] with PGT enabled, does not contain any zeroed out regions indicating that information loss is mitigated as the features pass on from the feature extracting layers to the fully connected layer.

In Table 2 we find that the baseline performance for high batch sizes (1024) is drastically inferior to baselines for other batch sizes. PGT helps regain some of the lost performance by 0.682% (65.498% vs. 64.816%).

At a batch size of 512, invoking PGT improves the training accuracy baseline by 1.48% and the test accuracy baseline by 0.684%, while at a batch size of 256, the improvement in training and test accuracies are 1.11% and 1.018% respectively. In comparison, the test accuracy improvements obtained by GC and AGC

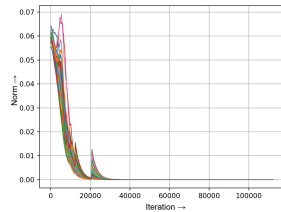


Figure 4: Zeroing out of feature maps in the second layer non-normalized ResNet-18.

at batch size of 256 is much less at 0.27% and 0.5%, respectively. On the training accuracy front, since we get a significant boost (1.48% at batch size of 512 and 1.11% at batch size of 256), it leads us to infer that when PowerGrad Transform is used, the network fits the training dataset more tightly and the convergence optima is significantly superior. When AGC and PGT are combined, we see a tremendous increase in test accuracy of over 2.06% over the baseline.

5. Ablation Study

We conduct ablation studies to investigate the effects of PowerGrad Transform for different values of the hyperparameter (α), where we use the ResNet-50 architecture and combine our proposed method with different schedulers, regularization techniques and different values of α . We report our findings in Table 3. First we examine the effect of

Table 2: Results for non-normalized ResNet-18 on ImageNet-1K. Best training and test accuracies are highlighted in red and blue respectively. Top differences in training and test accuracies are marked in yellow.

Batch Size	Method	PGT (α)	Train Acc.(%)	Train Diff(%)	Test Acc.(%)	Test Diff(%)
1024	Baseline	-	66.27	-	64.816	-
1024	PGT	0.92	66.62	+0.35	65.498	+0.682
512	Baseline	-	68.02	-	66.552	-
512	PGT	0.25	69.5	+1.48	67.236	+0.684
256	Baseline	-	68.86	-	66.796	-
256	GC	-	69.04	+0.18	67.064	+0.268
256	AGC	-	69.06	+0.2	67.298	+0.502
256	PGT	0.25	69.97	+1.11	67.814	+1.018
256	Baseline	-	68.86	-	66.796	-
256	GC+PGT	0.25	68.67	-0.19	67.088	+0.292
256	AGC+PGT	0.25	70.92	+2.06	68.856	+2.06

PGT on the step scheduler baseline in order to later compare it to the cosine scheduler baseline. **Row-1**) We begin with the step scheduler baseline (75.97%). **Row-2**) PGT improves upon the step scheduler baseline (test set) by a substantial margin with 0.524% (76.494% as opposed to 75.97%). **Row-3**) Introducing the cosine scheduler yields a 0.59% improvement (76.56% vs. 75.97%) over the step scheduler. **Row-4**) After introducing label smoothing, the test accuracy relative to the cosine scheduler baseline increases by only 0.138% (from 76.56% to 76.698%). **Row-5**) However, introducing PGT with $\alpha = 0.3$ alone (without label smoothing) improves the cosine scheduler baseline by 0.326% (76.886% vs. 76.56%). **Row-6**) Combining PGT ($\alpha = 0.3$) with label smoothing improves the performance on the test set further by 0.408% (from 76.56% to 76.968%) and reduces the generalization gap (from 2.54% to 1.5%). However, the impact of combining PGT with label smoothing can vary depending on the value of the hyperparameter (α). **Row-7**) With a PGT hyperparameter value of $\alpha = 0.05$, we notice the greatest performance improvement, 1.246% over the step scheduler test baseline and 0.656% over the cosine scheduler test baseline. **Row-8**) Adding label smoothing to PGT ($\alpha = 0.05$) hurts performance even though it reduces the generalization gap.

Table 3: Ablation study for ResNet-50 on ImageNet-1K.

#Row	Scheduler	Label Smoothing	PGT (α)	Train Acc.(%)	Test Acc.(%)	Gap(%)
1.	Step	\times	\times	78.99	75.97	3.02
2.	Step	\times	0.3	79.56	76.494	3.066
3.	Cosine	\times	\times	79.18	76.56	2.62
4.	Cosine	0.1	\times	78.81	76.698	2.112
5.	Cosine	\times	0.3	79.43	76.886	2.544
6.	Cosine	0.1	0.3	78.47	76.968	1.502
7.	Cosine	\times	0.05	79.68	77.216	2.464
8.	Cosine	0.1	0.05	77.69	76.39	1.3

6. Conclusion

PowerGrad Transform enables a significantly better fit to the dataset as measured by training and test accuracy metrics. With PGT, gradient behavior is enhanced and weights attain better values in normalized networks and degenerate states are avoided in non-BN networks. We provide theoretical analyses of the transformation. With different network topologies and datasets, we are able to show the potential of PGT and explore its impacts from an empirical standpoint. PGT helps the network to improve its learning capabilities by locating a more optimum convergence point and simultaneously speeds up training.

References

- [1] Andrew Brock, Soham De, Samuel L Smith, and Karen Simonyan. High-performance large-scale image recognition without normalization. *arXiv preprint arXiv:2102.06171*, 2021.
- [2] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [3] Shiming Ge, Zhao Luo, Chunhui Zhang, Yingying Hua, and Dacheng Tao. Distilling channels for efficient deep tracking. *IEEE Transactions on Image Processing*, 29:2610–2621, 2019.
- [4] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [5] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [6] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.
- [7] Yann LeCun, D Touresky, G Hinton, and T Sejnowski. A theoretical framework for back-propagation. In *Proceedings of the 1988 connectionist models summer school*, volume 1, pages 21–28, 1988.
- [8] Timothy P Lillicrap, Adam Santoro, Luke Marris, Colin J Akerman, and Geoffrey Hinton. Backpropagation and the brain. *Nature Reviews Neuroscience*, 21(6):335–346, 2020.
- [9] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- [10] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. Understanding the exploding gradient problem. *CoRR, abs/1211.5063*, 2(417):1, 2012.
- [11] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pages 1310–1318. PMLR, 2013.
- [12] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- [13] Samuel Smith, Erich Elsen, and Soham De. On the generalization benefit of noise in stochastic gradient descent. In *International Conference on Machine Learning*, pages 9058–9067. PMLR, 2020.
- [14] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016.
- [15] Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jad-babaie. Why gradient clipping accelerates training: A theoretical justification for adaptivity. *arXiv preprint arXiv:1905.11881*, 2019.
- [16] Ying Zhang, Tao Xiang, Timothy M Hospedales, and Huchuan Lu. Deep mutual learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4320–4328, 2018.