

What Deep CNNs Benefit from Global Covariance Pooling: An Optimization Perspective

Qilong Wang¹, Li Zhang¹, Banggu Wu¹, Dongwei Ren¹, Peihua Li², Wangmeng Zuo³, Qinghua Hu^{1,*}

¹ Tianjin Key Lab of Machine Learning, College of Intelligence and Computing, Tianjin University, China

² Dalian University of Technology, China ³ Harbin Institute of Technology, China

Abstract

Recent works have demonstrated that global covariance pooling (GCP) has the ability to improve performance of deep convolutional neural networks (CNNs) on visual classification task. Despite considerable advance, the reasons on effectiveness of GCP on deep CNNs have not been well studied. In this paper, we make an attempt to understand what deep CNNs benefit from GCP in a viewpoint of optimization. Specifically, we explore the effect of GCP on deep CNNs in terms of the Lipschitzness of optimization loss and the predictiveness of gradients, and show that GCP can make the optimization landscape more smooth and the gradients more predictive. Furthermore, we discuss the connection between GCP and second-order optimization for deep CNNs. More importantly, above findings can account for several merits of covariance pooling for training deep CNNs that have not been recognized previously or fully explored, including significant acceleration of network convergence (i.e., the networks trained with GCP can support rapid decay of learning rates, achieving favorable performance while significantly reducing number of training epochs), stronger robustness to distorted examples generated by image corruptions and perturbations, and good generalization ability to different vision tasks, e.g., object detection and instance segmentation. We conduct extensive experiments using various deep CNN models on diversified tasks, and the results provide strong support to our findings.

1. Introduction

Global covariance pooling (GCP) that is used to replace global average pooling (GAP) for aggregating the last convolution activations of deep convolutional neural networks

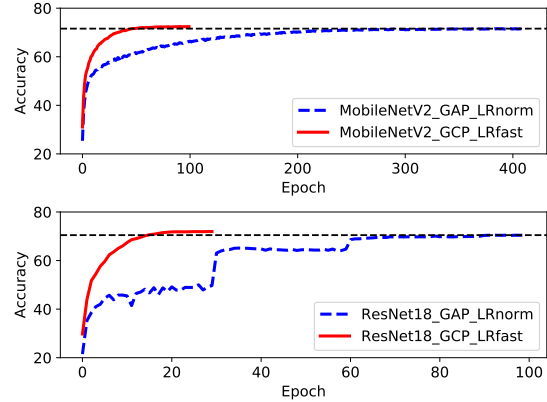


Figure 1. Convergence curves of MobileNetV2 [38] and ResNet18 [16] with global average pooling (GAP) and global covariance pooling (GCP) on ImageNet. Note that the networks with GCP converge much faster while achieving matching or better results. We account for it from an optimization perspective (see Section 3).

(CNNs) has achieved remarkable performance gains on a variety of vision tasks [20, 30, 12, 41, 11, 26, 29, 25]. Existing GCP-based works mainly focus on obtaining better performance using various normalization methods [20, 30, 26, 29] and richer statistics [41, 9, 8, 3] or achieving comparable results with low-dimensional covariance representations [13, 22, 14, 42]. However, the reasons on effectiveness of GCP on deep CNNs have not been well studied. Although some works explain them from the perspectives of statistical modeling [20, 30, 26] or geometry [26], some behaviors of deep CNNs with GCP still lack of reasonable explanations. For example, as illustrated in Figure 1, why GCP can significantly speed up convergence of deep CNNs. Particularly, the networks with GCP can achieve matching or better performance than GAP-based ones, but only use less than 1/4 of training epochs of the latter one.

In this paper, we make an attempt to understand the effectiveness of GCP on deep CNNs from an optimization perspective, thereby explaining behaviors of the networks with GCP in an intuitive way. To this end, we explore

*Qinghua Hu is the corresponding author.

Email: {qlwang, li_zhang, huqinghua}@tju.edu.cn. The work was supported by the National Natural Science Foundation of China (No. 61806140, 61971086, 61925602, U19A2073, 61732011). Qilong Wang was supported by National Postdoctoral Program for Innovative Talents.

the effect of GCP on optimization landscape and gradient computation of deep CNNs, inspired by recent work [39]. Specifically, we first train two widely used CNN models (i.e., MobileNetV2 [38] and ResNet-18 [16]) on large-scale ImageNet [10] using GAP and GCP for aggregating the last convolution activations, respectively. Then, we compare them from optimization perspective, and find that GCP can improve the stability of optimization loss (i.e., Lipschitzness) and the stability of gradients (i.e., predictiveness) over the commonly used GAP. Furthermore, by analyzing back-propagation of GCP and second-order optimization [32, 33, 36] in the context of deep CNNs, we make out that the influence of GCP on optimization shares some similar philosophy with K-FAC [33, 36]. Above findings provide an inspiring view to understand effectiveness of GCP on deep CNNs, which can intuitively and reasonably account for the behaviors of deep CNNs trained with GCP, e.g., GCP makes the optimization landscape more smooth, leading to network convergence to a better local minimum (i.e., better performance as shown in [30, 41, 26]) and much faster convergence as illustrated in Figure 1.

Based on the foregoing findings, we can explain several merits delivered by GCP for training deep CNNs that have not been recognized previously or fully explored. Firstly, since GCP is able to smoothen optimization landscape, deep CNNs with GCP can support rapid decay of learning rates for fast convergence. Meanwhile, previous work [39] shows that improvement of Lipschitzness can accelerate convergence of deep CNNs. To verify this point, we conduct experiments using a variety of deep CNN architectures (i.e., MobileNetV2 [38], ShuffleNet V2 [31] and ResNets [16]) on ImageNet [10]. The results show that, under the setting of rapid decay of learning rates, deep CNNs with GCP achieve comparable performance to GAP-based ones, while using less than a quarter of training epochs. By adjusting schedule of learning rates, deep CNNs with GCP can converge to better local minima using less training epochs.

Secondly, GCP improves the stability of both optimization loss and gradients, so it makes deep CNNs more robust to inputs perturbed by some distortions. Meanwhile, previous work [7] shows that control of the Lipschitz constant is helpful for improving robustness to examples with crafted distortions for confusing classifiers. Therefore, we experiment on recently introduced IMAGENET-C and IMAGENET-P benchmarks [17], where distortions are achieved by common image corruptions and perturbations. The results show that GCP can significantly improve robustness of deep CNNs to image corruptions and perturbations.

Thirdly, GCP usually allows deep CNNs converge to better local minima, and thus deep CNNs pre-trained with GCP can be exploited to provide an effective initialization model for other visual tasks. Therefore, we verify it by transferring the pre-trained CNN models to MS COCO benchmark [28]

for object detection and instance segmentation tasks, and the results indicate that pre-trained CNNs with GCP is superior to GAP-based ones.

The contributions of this paper are concluded as follows. (1) To our best knowledge, we make the first attempt to understand the effectiveness of GCP in the context of deep CNNs from an optimization perspective. Specifically, we show that GCP can improve the Lipschitzness of optimization loss and the predictiveness of gradients. Furthermore, we discuss the connection between GCP and second-order optimization. These findings provide an inspiring view to better understand the behaviors of deep CNNs trained with GCP. (2) More importantly, our findings above can explain several merits of GCP for training deep CNNs that have not been recognized previously or fully explored, including significant acceleration of convergence with rapid decay of learning rates, stronger robustness to distorted examples and good generalization ability to different vision tasks. (3) We conduct extensive experiments using six representative deep CNN architectures on image classification, object detection and instance segmentation, the results of which provide strong support to our findings.

2. Related Work

DeepO₂P [20] and B-CNN [30] are among the first works introducing GCP into deep CNNs. DeepO₂P extends the second-order (covariance) pooling method (O₂P) [5] to deep architecture, while B-CNN captures interactions of localized convolution features by a trainable bilinear pooling. Wang *et al.* [41] propose a global Gaussian distribution embedding network for utilizing the power of probability distribution modeling and deep learning jointly. Li *et al.* [26] present matrix power normalization for GCP and make clear its statistical and geometrical mechanisms. Beyond GCP, some researches [8, 3] propose to use richer statistics for further possible improvement. The aforementioned methods study GCP from perspectives of statistical modeling or Riemannian geometry. Different from them, we interpret the effectiveness of GCP from an optimization perspective, and further explore merits of GCP for training deep CNNs.

Since deep CNNs are black-boxes themselves, understanding effect of GCP on deep CNNs still is a challenging issue. Many recent works make efforts [43, 44, 39, 2, 45] towards understanding deep CNNs or analyzing the effect of fundamental components, e.g., individual units, batch normalization (BN) [19] and optimization algorithm. Specifically, Zeiler *et al.* [43] and Zhou *et al.* [45] visualize feature maps by deconvolution and image regions of maximal activation units, respectively. Zhang *et al.* [44] and Bjorck *et al.* [2] design a series of controlled experiments to analyze generalization ability of deep CNNs and understand effect of BN on deep CNNs, respectively. In particular, a recent work [39] investigates the effect of BN by exploring opti-

mization landscape of VGG-like networks [40] trained on CIFAR10 [23], while providing theoretical analysis on Lipschitzness improvement using a fully-connected layer. Motivated by [39], we make an attempt to understand *the effect of GCP on deep CNNs* from an optimization perspective.

3. An Optimization Perspective for GCP

In this section, we first revisit global covariance pooling (GCP) for deep CNNs. Then, we analyze smoothing effect of GCP on deep CNNs, and finally discuss its connection to second-order optimization.

3.1. Revisiting GCP

Let $\mathcal{X} \in \mathbb{R}^{W \times H \times D}$ be the output of the last convolution layer of deep CNNs, where W , H and D indicate width, height and dimension of feature map, respectively. To summarize \mathcal{X} as global representations for final prediction, most existing CNNs employ GAP, i.e., $\sum_{i=1}^N \mathbf{X}_i$, where the feature tensor \mathcal{X} is reshaped to a feature matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$ and $N = W \times H$. Many recent works [30, 26] demonstrate the superiority of GCP over GAP. To perform GCP, the sample covariance matrix of \mathbf{X} can be calculated as

$$\Sigma = \mathbf{X}^T \mathbf{J} \mathbf{X}, \quad \mathbf{J} = \frac{1}{N} (\mathbf{I} - \frac{1}{N} \mathbf{1} \mathbf{1}^T), \quad (1)$$

where \mathbf{I} is the $N \times N$ identity matrix, and $\mathbf{1}$ is a N -dimensional vector of all elements being one.

Normalization plays an important role in GCP, and different normalization methods have been studied, including matrix logarithm normalization [20], element-wise power normalization followed by ℓ_2 normalization [30], matrix square-root normalization [41, 29, 25] and matrix power normalization [26]. Among them, matrix square-root (i.e., the power of $1/2$) normalization is preferred considering its promising performance on both large-scale and small-scale visual classification tasks. Therefore, this paper uses GCP with matrix square-root normalization, i.e.,

$$\Sigma^{\frac{1}{2}} = \mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}^T, \quad (2)$$

where \mathbf{U} and $\mathbf{\Lambda}$ are the matrix of eigenvectors and the diagonal matrix of eigenvalues of Σ , respectively.

3.2. Smoothing Effect of GCP

To understand the mechanism of GCP, we investigate effect of GCP on optimization landscape of deep CNNs. Inspired by [39], we explore effect of GCP on stability of optimization loss (i.e., Lipschitzness) and stability of gradients (i.e., predictiveness). Specifically, for examining stability of optimization loss, we measure how loss \mathcal{L} changes along direction of current gradient at each training step. Given the input \mathbf{X} , variation of optimization loss is calculated as

$$\Delta_l = \mathcal{L}(\mathbf{X} + \eta_l \nabla \mathcal{L}(\mathbf{X})), \eta_l \in [a, b], \quad (3)$$

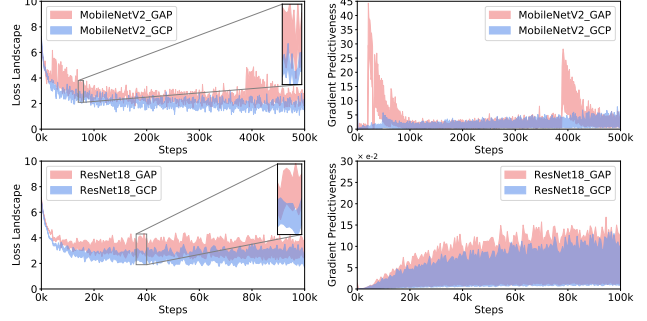


Figure 2. Comparison of (top) MobileNetV2 and (bottom) ResNet-18 trained with GCP or GAP in terms of (left) loss Lipschitzness and (right) gradient predictiveness. Detailed descriptions and discussions can be found in Section 3.2.

where $\nabla \mathcal{L}(\mathbf{X})$ indicates the gradient of loss with respect to the input \mathbf{X} , and η_l is step size of gradient descent. Clearly, smaller variation (range) of Δ_l indicates that optimization landscape is smoother and more easily controlled during training process [39]. As shown in [39], range of Δ_l in Eqn. (3) reflects the Lipschitzness of optimization loss.

To examine the stability of gradients, we measure how the gradient of loss \mathcal{L} changes by computing the Euclidean distance between the gradient of loss and gradients along the original gradient direction on an interval of step sizes. Thus, the gradient predictiveness can be formulated as

$$\Delta_g = \|\nabla \mathcal{L}(\mathbf{X}) - \nabla \mathcal{L}(\mathbf{X} + \eta_g \nabla \mathcal{L}(\mathbf{X}))\|_2, \eta_g \in [a, b], \quad (4)$$

where η_g is step size. Similar to the stability of optimization loss, smaller range of Δ_g implies that gradient is more insensitive to step size, having better gradient predictiveness. The loss Lipschitzness and gradient predictiveness greatly affect optimization of deep CNNs, i.e., smoothness of landscape and robustness to hyper-parameters.

According to the discussions above, we train the networks with GCP or GAP while comparing their loss Lipschitzness and gradient predictiveness to analyze smoothing effect of GCP. Without loss of generality, we employ the widely used MobileNetV2 [38] and ResNet-18 [16] as backbone models, and train them on large-scale ImageNet [10]. For training MobileNetV2 and ResNet-18 with GCP, following [26], we reduce dimension of the last convolution activations to 256, and train the networks using stochastic gradient descent (SGD) with the same hyper-parameters in [38] and [16]. Besides, we do not use *dropout* operation to avoid randomness for MobileNetV2, and discard down-sampling operation in *conv5_x* for ResNet-18. The left endpoints in the ranges of Δ_l and Δ_g are set to the initial learning rates and right endpoints are as large as possible while ensuring stable training of GAP. As such, $\eta_l (\eta_g) \in [0.045, 1.5]$ and $\eta_l (\eta_g) \in [0.1, 75]$ are for MobileNetV2 and ResNet-18, respectively.

The behaviors of MobileNetV2 with 500K steps (~ 37 epochs) and ResNet-18 with 100K steps (~ 20 epochs) trained with GAP and GCP in terms of loss Lipschitzness and gradient predictiveness are shown in top and bottom of Figure 2, respectively. For both MobileNetV2 and ResNet-18, we observe that the networks with GCP have smaller variations of the optimization loss (i.e., Δ_l) than GAP-based ones, while optimization losses of the networks with GCP are consistently lower than those trained with GAP. These results demonstrate that GCP can improve Lipschitzness of optimization loss while converging faster under the same setting with GAP-based ones. Meanwhile, for changes of the gradient, Δ_g of the networks with GCP is more stable than GAP-based ones, suggesting the networks with GCP have better gradient predictiveness. In addition, the jumps of both loss landscape and gradient predictiveness for GAP with MobileNetV2 indicate that variations of loss and gradient are considerably large, suggesting varying step sizes are likely to drive the loss uncontrollably higher. In contrast, the variations of one with GCP are fairly small and consistent, suggesting GCP helps for stable training. *In a nutshell, GCP has the ability to smoothen optimization landscape of deep CNNs and improve gradient predictiveness.*

3.3. Connection to Second-order Optimization

Furthermore, we analyze back-propagation (BP) of GCP to explore its effect on optimization of deep CNNs. Let \mathbf{X} be output of the last convolution layer, the gradient of loss \mathcal{L} with respect to \mathbf{X} for GAP layer can be computed as $\frac{\partial \mathcal{L}}{\partial \mathbf{X}} = \mathbf{C}^T \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GAP}}}$, where $\mathbf{Z}_{\text{GAP}} = \sum_{i=1}^N \mathbf{X}_i$ and \mathbf{C} is a constant matrix. To update weights \mathbf{W} of the last convolution layer, gradient descent is performed as follows:

$$\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t - \eta \nabla_{\mathbf{W}} \mathcal{L}_t, \quad (5)$$

where $\nabla_{\mathbf{W}} \mathcal{L}_t = \frac{\partial \mathcal{L}}{\partial \mathbf{X}} \frac{\partial \mathbf{X}}{\partial \mathbf{W}_t} = \mathbf{C}^T \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GAP}}} \frac{\partial \mathbf{X}}{\partial \mathbf{W}_t}$, and t indicates t -th iteration.

Let $\mathbf{Z}_{\text{GCP}} = (\mathbf{X}^T \mathbf{J} \mathbf{X})^{\frac{1}{2}}$. The derivative of loss with respect to \mathbf{X} for GCP layer can be written as

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{X}} &= 2\mathbf{J}\mathbf{X} \left[\mathbf{U} \left(\left(\mathbf{K}^T \circ \left(\mathbf{U}^T 2 \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \right)_{\text{sym}} \mathbf{U} \Lambda^{\frac{1}{2}} \right) \right. \right. \right. \\ &\quad \left. \left. \left. + \left(\frac{1}{2} \Lambda^{-\frac{1}{2}} \mathbf{U}^T \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \mathbf{U} \right)_{\text{diag}} \right) \mathbf{U}^T \right]_{\text{sym}}, \end{aligned} \quad (6)$$

where \mathbf{U} and Λ are the matrix of eigenvectors and the diagonal matrix of eigenvalues of sample covariance of \mathbf{X} , and \mathbf{K} is a mask matrix associated with eigenvalues. Here, \circ denotes matrix Hadamard product; $(\cdot)_{\text{sym}}$ and $(\cdot)_{\text{diag}}$ indicate matrix symmetrization and diagonalization, respectively. More details can refer to [21, 41, 26]. With some assumptions and simplification, Eqn. (6) can be trimmed as

$$\frac{\partial \mathcal{L}}{\partial \mathbf{X}} \approx 2\mathbf{J}\mathbf{X} \left(2\mathbf{K}^T \circ \Lambda^{\frac{1}{2}} + \frac{1}{2} \Lambda^{-\frac{1}{2}} \right) \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}}. \quad (7)$$

Method	Gradient	Remark
GAP	$\eta \mathbf{C}^T \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GAP}}} \frac{\partial \mathbf{X}}{\partial \mathbf{W}_t}$	\mathbf{C} is a constant matrix
GCP	$\approx \mathbf{F}^{-1} \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \frac{\partial \mathbf{X}}{\partial \mathbf{W}_t}$	$\mathbf{F}^{-1} = \eta 2\mathbf{J}\mathbf{X} \left(2\mathbf{K}^T \circ \Lambda^{\frac{1}{2}} + \frac{1}{2} \Lambda^{-\frac{1}{2}} \right)$
K-FAC [36]	$\mathbf{H}^{-1} \mathbf{C}^T \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GAP}}} \frac{\partial \mathbf{X}}{\partial \mathbf{W}_t}$	$\mathbf{H}^{-1} = \eta \left(\frac{\partial \mathcal{L}}{\partial \mathbf{X}} \right)^{-1} \otimes \hat{\mathbf{X}}^{-1}$

Table 1. Comparison of gradients involved in GAP, GCP and GAP with K-FAC [36]. \otimes indicates Kronecker product.

Details of Eqn. (6) and Eqn. (7) can be found in the supplemental file. By substituting Eqn. (7) into Eqn. (5), we can approximatively update the weights of convolution as

$$\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t - \mathbf{F}^{-1} \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \frac{\partial \mathbf{X}}{\partial \mathbf{W}_t}, \quad (8)$$

where $\mathbf{F}^{-1} = \eta 2\mathbf{J}\mathbf{X} \left(2\mathbf{K}^T \circ \Lambda^{\frac{1}{2}} + \frac{1}{2} \Lambda^{-\frac{1}{2}} \right)$.

Previous works [32, 33, 36] show second-order optimization (i.e., $\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t - \mathbf{H}^{-1} \nabla \mathcal{L}_t$) can speed up training of deep neural networks. However, computation of inverse of Hessian matrix (\mathbf{H}^{-1}) is usually very expensive and is sensitive to noise. Therefore, many methods [32, 33, 1] are proposed to approximate \mathbf{H}^{-1} . Recently, K-FAC [33, 36] based on an accurate approximation of the Fisher information matrix has proven to be effective in optimizing deep CNNs, which approximates \mathbf{H}^{-1} using Kronecker product between inverse of input of convolution layer (i.e., $\hat{\mathbf{X}}^{-1}$) and inverse of gradient of the loss with respect to the output (i.e., $\left(\frac{\partial \mathcal{L}}{\partial \mathbf{X}} \right)^{-1}$). The gradients involved in GAP, GCP and GAP with K-FAC are compared in Table 1, where the trimmed BP of GCP (7) shares some similar philosophy with K-FAC. The key difference is that \mathbf{F}^{-1} is computed by the output \mathbf{X} and its eigenvalues, while \mathbf{H}^{-1} is approximated by the input $\hat{\mathbf{X}}$ and the gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{X}}$. The experiments in Section 4.1 show that ResNet-50 with GCP uses less training epochs to achieve result matching that of K-FAC [36], which may indicate that BP of GCP is a potential alternative of pre-conditioner for Hessian matrix.

4. Merits Benefited from GCP

In previous section, we explore the effect of GCP on deep CNNs from an optimization perspective. Specifically, we show smoothing effect of GCP on optimization landscape, and discuss the connection between GCP and second-order optimization. Furthermore, these findings can also account for several merits delivered by GCP for training deep CNNs that have not been recognized previously or fully explored, including significant acceleration of network convergence, stronger robustness to distorted examples generated by image corruptions and perturbations, and good generalization ability to different vision tasks. In this

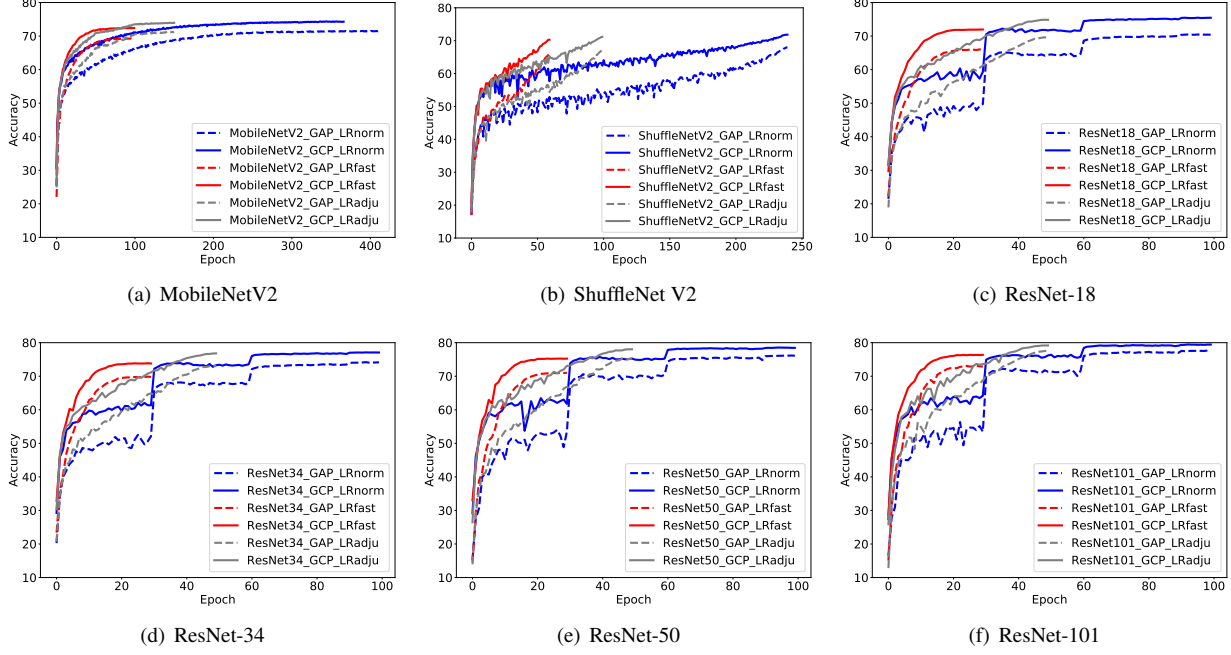


Figure 3. Convergence curves of different deep CNN models trained with GAP and GCP under various settings of learning rates (i.e., LR_{norm} , LR_{fast} and LR_{adju}) on ImageNet.

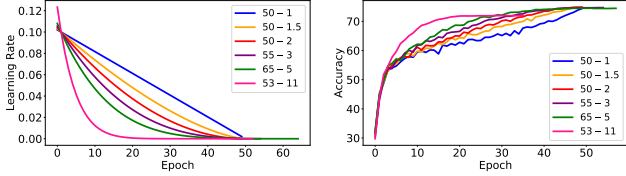


Figure 4. (Left) Curves of various settings ($\rho - e_f$) of lr and (Right) their corresponding convergence curves using ResNet-18 as backbone model on ImageNet. Here e_f and ρ denote the final epoch and power of the polynomial decay in Eqn. (9).

work, all programs are implemented by Pytorch package¹, and run on a workstation with four NVIDIA RTX 2080 Ti GPUs. The detailed descriptions are given as follows.

4.1. Acceleration of Network Convergence

It is well known that training of deep CNNs is a time-consuming process, requiring a mass of computing resources. Therefore, acceleration of network convergence plays a key role in fast training of deep CNNs, especially for large-scale training datasets. Previous study [25] empirically shows the networks trained with GCP converge faster than GAP-based ones, which can be explained by our findings. Specifically, Lipschitzness improvement [39, 35] and connection to second-order optimization [36] brought by GCP can accelerate convergence of networks. Here,

¹https://github.com/ZhangLi-CS/GCP_Optimization

we further show deep CNNs with GCP can support rapid decay of learning rates for significantly improving convergence speed, due to its ability of smoothing optimization landscape. To verify it, we conduct extensive experiments using six representative deep CNN architectures on ImageNet, including MobileNetV2 [38], ShuffleNet V2 [31] and ResNets [16] of 18, 34, 50 and 101 layers.

Specifically, we train all the networks with GAP and GCP by varying the settings of learning rates (lr). Firstly, we use the setting of lr in each original paper, which is indicated by LR_{norm} . As shown in [36], ResNet-50 converges much faster using second-order optimization (K-FAC), in which lr is scheduled by polynomial decay, i.e.,

$$\ell_e = \ell_0 \times \left(1 - \frac{e - e_s}{e_f - e_s}\right)^\rho, \quad (9)$$

where ℓ_0 is the initial lr ; e , e_s and e_f indicate e -th, the start and the final epochs, respectively. The parameter ρ controls decay rate. Inspired by [36], we employ the setting of lr in Eqn. (9) for fast training of ResNets. To determine parameter ρ , we set $\ell_0 = 0.1$ and train ResNet-18 with GCP on ImageNet using various ρ while controlling e_f be less than 65, which is consistent with [36]. Figure 4 illustrates the curve of each lr and its corresponding convergence. Clearly, larger ρ leads to faster convergence but lower accuracy. The setting with ($\rho = 11, e_f = 53$) has the fastest convergence, denoted by LR_{fast} hereafter, which converges within 30 epochs. Among them, the setting with ($\rho = 2, e_f = 50$) makes the best trade-off be-

Backbone Model	Pooling Method	lr	BS	Training Epochs	Matching Epoch	Top-1 Accuracy	Top-5 Accuracy
MobileNetV2	GAP	LR_{norm}	96	400	N/A	71.58	90.30
	GCP	LR_{adj}	192	150	68(\downarrow_{332})	73.97($\uparrow_{2.39}$)	91.54($\uparrow_{1.24}$)
ShuffleNet V2	GAP	LR_{norm}	1,024	240	N/A	67.96	87.84
	GCP	LR_{adj}	1,024	100	78(\downarrow_{162})	71.17($\uparrow_{3.21}$)	89.74($\uparrow_{1.90}$)
ResNet-18	GAP	LR_{norm}	256	100	N/A	70.47	89.62
	GCP	LR_{adj}	256	50	32(\downarrow_{68})	74.86($\uparrow_{4.39}$)	91.81($\uparrow_{2.19}$)
ResNet-34	GAP	LR_{norm}	256	100	N/A	74.19	91.61
	GCP	LR_{adj}	256	50	38(\downarrow_{62})	76.81($\uparrow_{2.62}$)	93.09($\uparrow_{1.48}$)
ResNet-50	GAP	LR_{norm}	256	100	N/A	76.17	92.93
	GCP	LR_{adj}	256	50	40(\downarrow_{60})	78.03($\uparrow_{1.86}$)	93.95($\uparrow_{1.02}$)
	GAP	K-FAC [36]*	4,096	35	N/A	75.10($\downarrow_{1.07}$)	N/A
	GCP	LR_{fast}	256	30	20(\downarrow_{15})	75.31($\downarrow_{0.86}$)	92.15($\downarrow_{0.78}$)
ResNet-101	GAP	LR_{norm}	256	100	N/A	77.67	93.89
	GCP	LR_{adj}	256	50	41(\downarrow_{59})	79.18($\uparrow_{1.51}$)	94.51($\uparrow_{0.62}$)

Table 2. Comparison of different CNNs trained with GAP using LR_{norm} and those trained with GCP using LR_{adj} on ImageNet. ‘Matching Epoch’ indicates that at which the networks with GCP achieve comparable performance with the corresponding networks with GAP undergoing full training epochs. *: The result of K-FAC is duplicated from the original paper [36].

Method		Setting of lr
Backbone	lr	
MobileNet V2	LR_{norm}	0.045×0.98^e
	LR_{fast}	0.06×0.92^e
	LR_{adj}	$\begin{cases} linear(6e^{-2}, 1e^{-3}, 0), \\ linear(1e^{-2}, 1e^{-4}, 50), \\ linear(1e^{-3}, 1e^{-5}, 100) \end{cases}$
ShuffleNetV2	LR_{norm}	$0.5 \times (1 - \frac{step}{t_step})$
	LR_{fast}	
	LR_{adj}	
ResNet	LR_{norm}	$0.1^{(e//30)+1}$
	LR_{fast}	$0.1 \times (1 - \frac{e-1}{52})^{11}$
	LR_{adj}	$0.1 \times (1 - \frac{e-1}{49})^2$

Table 3. Detailed settings of lr for various CNNs.

Method	IMAGENET-C		IMAGENET-P	
	mCE	Relative mCE	mFP	mT5D
MobileNetV2 + GAP	87.1	114.9	79.8	96.5
MobileNetV2 + GCP	81.7($\downarrow_{5.4}$)	110.6($\downarrow_{4.3}$)	64.3($\downarrow_{15.5}$)	87.6($\downarrow_{8.9}$)
ShuffleNet V2 + GAP	92.7	126.7	94.7	108.2
ShuffleNet V2 + GCP	85.2($\downarrow_{7.5}$)	112.6($\downarrow_{14.1}$)	75.2($\downarrow_{19.5}$)	95.5($\downarrow_{12.7}$)
ResNet-18 + GAP	84.7	103.9	72.8	87.0
ResNet-18 + GCP	76.3($\downarrow_{8.4}$)	101.3($\downarrow_{2.6}$)	53.2($\downarrow_{19.6}$)	77.1($\downarrow_{9.9}$)
ResNet-34 + GAP	77.9	98.7	61.7	79.5
ResNet-34 + GCP	72.4($\downarrow_{5.5}$)	96.9($\downarrow_{1.8}$)	47.7($\downarrow_{14.0}$)	72.4($\downarrow_{7.1}$)
ResNet-50 + GAP	76.7	105.0	58.0	78.3
ResNet-50 + GCP	70.7($\downarrow_{6.0}$)	97.9($\downarrow_{7.1}$)	47.5($\downarrow_{10.5}$)	74.6($\downarrow_{3.7}$)
ResNet-101 + GAP	70.3	93.7	52.6	73.9
ResNet-101 + GCP	65.5($\downarrow_{4.8}$)	89.1($\downarrow_{4.6}$)	42.1($\downarrow_{10.5}$)	68.3($\downarrow_{5.6}$)

Table 4. Comparison of GAP and GCP on IMAGENET-C and IMAGENET-P.

tween convergence speed and classification accuracy, denoted by LR_{adj} .

For MobileNetV2, the original lr is scheduled by exponential decay (i.e., 0.045×0.98^e), and we set LR_{fast} by decreasing base number while increasing the initial lr , i.e., 0.06×0.92^e . LR_{adj} is scheduled by a stage-wise linear decay, i.e., $linear(l_s, l_e, n) = l_s - \frac{l_s - l_e}{50}(e - n)$ for e -th epoch where l_s , l_e and n are initial lr , final lr and start epoch in each stage. The lr of ShuffleNet V2 in the original paper is scheduled by a step-wise linear decay, and number of total steps (i.e., t_step) is $3e5$ (~ 240 epochs). For ShuffleNet V2, LR_{fast} and LR_{adj} are set by reducing training epochs to 60 and 100, respectively. The detailed settings of lr are summarized in Table 3.

The convergence curves of different networks trained with GAP and GCP under various settings of lr are illustrated in Figure 3, from which we have the following observations. (1) Comparing with the networks trained under LR_{norm} , those trained under the settings of LR_{fast} and LR_{adj} suffer from performance degradation. However, performance degradations of the networks with GCP are

Backbone	lr	D	Top-1 Acc.	Top-5 Acc.
MobiNetV2	LR_{norm}	256	74.36	91.90
		128	73.28 ($\downarrow 1.08$)	91.30 ($\downarrow 0.60$)
	LR_{adj}	256	73.97	91.54
		128	72.58 ($\downarrow 1.39$)	90.90 ($\downarrow 0.64$)
ResNet-50	LR_{norm}	256	78.56	94.16
		128	78.21 ($\downarrow 0.35$)	93.93 ($\downarrow 0.23$)
	LR_{adj}	256	78.03	93.95
		128	77.64 ($\downarrow 0.39$)	93.67 ($\downarrow 0.28$)

Table 5. Results of MobileNetV2 and ResNet-50 trained with GCP using different lr and dimensions (D) of input on ImageNet.

less than those based on GAP, especially under LR_{fast} . (2) The networks trained with GCP using LR_{fast} achieve better or matching results using only about 1/4 of the number of epochs for training networks with GAP under LR_{norm} .

Moreover, we compare the networks trained with GAP under LR_{norm} (i.e., the original settings) and those trained with GCP using LR_{adj} . According to the results in Table 2, we make a summary as follows. (1) Comparing to the networks trained with GAP under LR_{norm} , those trained with GCP under LR_{adj} achieve higher accuracies while

using less training epochs. (2) The networks with GCP obtain matching or comparable accuracies to GAP-based ones using much less training epochs, especially for lightweight CNN models, i.e., MobileNetV2 and ShuffleNet V2. For example, MobileNetV2 with GCP achieves matching accuracies to GAP-based one using only 68 epochs, while the latter one needs about 400 epochs. (3) Comparing to the second-order optimization method K-FAC [36], GCP obtains moderate accuracy gain using less training epochs, while achieving comparable accuracies with K-FAC using only 20 epochs. Furthermore, GCP is easier to implement. *The extensive experiments above strongly support our finding: GCP can significantly speed up convergence of deep CNNs with rapid decay of learning rates.*

Additionally, we assess the effect of dimension of covariance representations (COV-Reps) on behavior of convergence using MobileNetV2 (MobiNetV2) and ResNet-50 on ImageNet. If the dimension of input features is D , GCP will output a $D(D + 1)/2$ -dimensional COV-Reps. Here, we set D to 256 (the default setting) and 128, and train the networks under the settings of LR_{norm} and LR_{adj} , respectively. The results are given in Table 5, from which one can see that lower-dimensional COV-Reps still allow faster convergence of deep CNNs, but suffer from larger performance degradation in the case of faster convergence (i.e., LR_{adj}). This indicates that dimension of COV-Reps has a nontrivial effect on the behavior of convergence. Therefore, how to compress COV-Reps while preserving merits of high-dimensional ones is an important issue. Albeit many works are proposed to compress COV-Reps [13, 22, 42], they still have not been verified in large-scale scenarios. A potential solution is to learn compact COV-Reps from high-dimensional ones based on knowledge distillation [18], which will be studied in future.

4.2. Robustness to Distorted Examples

Improvement of loss Lipschitzness and gradient predictiveness brought by GCP make the networks more robust to inputs perturbed by distortions. We also note similar conclusion is stated in [7]. To verify this point, we conduct experiments on recently introduced IMAGENET-C and IMAGENET-P benchmarks [17]. Different from the works [4, 34] that study effect of adversarial distortions as a type of worst-case analysis for network robustness, these two benchmarks are designed to evaluate robustness of deep CNNs to common image corruptions and perturbations, which have a connection with adversarial distortions and play a key role in safety-critical applications.

The IMAGENET-C benchmark performs fifty types of corruptions (e.g., noise, blur, weather and digital) on validation set of ImageNet, and each type of corruption has five levels of severity. The IMAGENET-P benchmark generates a series of perturbation sequences on validation set

of ImageNet by performing more than ten types of perturbations, such as motion and zoom blur, brightness, translation, rotation, scale and tilt perturbations. Following the standard protocol in [17], we train all CNN models on training set (clean images) of ImageNet, and report the results on IMAGENET-C and IMAGENET-P benchmarks. The evaluation metrics include mean Corruption Error (mCE) and Relative mean Corruption Errors (Relative mCE) for IMAGENET-C, mean Flip Rate (mFR) and mean Top-5 Distance (mT5D) for IMAGENET-P. For details of the metrics one can refer to [17]. *Note that lower values indicate better performance for all evaluation metrics.*

For a fair comparison, we employ evaluation code released by the authors. Note that AlexNet [24] is a baseline model, which obtains value of 100 for all evaluation metrics. The results of different deep CNNs with GAP and GCP are given in Table 4, from which we can see that the networks with GCP significantly outperform GAP-based ones, suggesting that GCP can greatly improve the robustness of deep CNNs to common image corruptions and perturbations. Note that VGG-VD19 [40] and VGG-VD19 with BN achieve 88.9, 122.9, 66.9, 78.6 and 81.6, 111.1, 65.1, 80.5 in terms of mCE, Relative mCE, mFR, mT5D, respectively. Albeit BN improves the Lipschitzness, it is not robust to perturbations. In contrast, GCP is robust to both corruptions and perturbations. In [17], many schemes are suggested to improve the robustness to corruptions and perturbations, and our work shows that GCP is a novel and promising solution. Moreover, it is potential to combine GCP with other schemes for further improvement.

4.3. Generalization Ability to Other Tasks

Since GCP usually allows deep CNNs converge to better local minima, the networks with GCP pretrained on large-scale dataset can provide a better initialization model to other vision tasks. That is, they may have good generalization ability. To verify this, we first train the networks with GCP on ImageNet, and then directly apply them to object detection and instance segmentation tasks. Specifically, using ResNet-50 and ResNet-101 as backbone models, we compare performance of the networks trained with GAP and GCP on MS COCO [28] using Faster R-CNN [37] and Mask R-CNN [15] as basic detectors. For training networks with GCP, Li et al. [26] suggest no down-sampling in *conv5_1*. This increases resolution of the last feature maps while resulting in larger computational cost, especially for object detection and instance segmentation where large-size input images are required. To handle this issue, we introduce two strategies: (1) we still use down-sampling as done in the original ResNet, and the method is indicated by GCP_D ; (2) a max-pooling layer with a step size 2 is inserted before *conv5_1* for down-sampling, indicated by GCP_M .

For a fair comparison, all methods are implemented by

Backbone Model	Method	Detectors	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
ResNet-50	GAP	Faster R-CNN	36.4	58.2	39.2	21.8	40.0	46.2
	GCP _D		36.6($\uparrow 0.2$)	58.4($\uparrow 0.2$)	39.5($\uparrow 0.3$)	21.3($\downarrow 0.5$)	40.8($\uparrow 0.8$)	47.0($\uparrow 0.8$)
	GCP _M		37.1 ($\uparrow 0.7$)	59.1 ($\uparrow 0.9$)	39.9 ($\uparrow 0.7$)	22.0 ($\uparrow 0.2$)	40.9 ($\uparrow 0.9$)	47.6 ($\uparrow 1.4$)
ResNet-101	GAP		38.7	60.6	41.9	22.7	43.2	50.4
	GCP _D		39.5($\uparrow 0.8$)	60.7($\uparrow 0.1$)	43.1($\uparrow 1.2$)	22.9($\uparrow 0.2$)	44.1 ($\uparrow 0.9$)	51.4 ($\uparrow 1.0$)
	GCP _M		39.6 ($\uparrow 0.9$)	61.2 ($\uparrow 0.6$)	43.1 ($\uparrow 1.2$)	23.3 ($\uparrow 0.6$)	43.9($\uparrow 0.7$)	51.3($\uparrow 0.9$)
ResNet-50	GAP	Mask R-CNN	37.2	58.9	40.3	22.2	40.7	48.0
	GCP _D		37.3($\uparrow 0.1$)	58.8($\downarrow 0.1$)	40.4($\uparrow 0.1$)	22.0($\downarrow 0.2$)	41.1($\uparrow 0.4$)	48.2($\uparrow 0.2$)
	GCP _M		37.9 ($\uparrow 0.7$)	59.4 ($\uparrow 0.5$)	41.3 ($\uparrow 1.0$)	22.4 ($\uparrow 0.2$)	41.5 ($\uparrow 0.8$)	49.0 ($\uparrow 1.0$)
ResNet-101	GAP		39.4	60.9	43.3	23.0	43.7	51.4
	GCP _D		40.3($\uparrow 0.9$)	61.5($\uparrow 0.6$)	44.0($\uparrow 0.7$)	24.1 ($\uparrow 1.1$)	44.7($\uparrow 1.0$)	52.5($\uparrow 1.1$)
	GCP _M		40.7 ($\uparrow 1.3$)	62.0 ($\uparrow 1.1$)	44.6 ($\uparrow 1.3$)	23.9($\uparrow 0.9$)	45.2 ($\uparrow 1.5$)	52.9 ($\uparrow 1.5$)

Table 6. Object detection results of various deep CNN models using Faster R-CNN and Mask R-CNN on COCO val2017.

Method	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
R-50 + GAP	34.1	55.5	36.2	16.1	36.7	50.0
R-50 + GCP _D	34.2	55.3	36.4	15.8	37.1	50.1
R-50 + GCP _M	34.7	56.3	36.8	16.4	37.5	50.6
R-101 + GAP	35.9	57.7	38.4	16.8	39.1	53.6
R-101 + GCP _D	36.5	58.2	38.9	17.3	39.9	53.5
R-101 + GCP _M	36.7	58.7	39.1	17.6	39.9	53.7

Table 7. Instance segmentation results of various deep CNN models using Mask R-CNN on COCO val2017.

Backbone	Method	Top-1	Top-5	AP	AP_{50}	AP_{75}
ResNet-18	w/o DS	75.47	92.23	30.0	50.7	31.4
	w/ DS	74.48	91.68	30.3	51.0	32.2
ResNet-50	w/o DS	78.56	94.16	36.6	58.4	39.5
	w/ DS	78.10	94.09	36.8	58.5	39.7
ResNet-101	w/o DS	79.48	94.75	39.5	60.7	43.1
	w/ DS	79.11	94.56	39.6	60.9	43.4

Table 8. Comparison of GCP with or without down-sampling (DS) on ImageNet (columns 3 and 4) and MS COCO (columns 5, 6, and 7). Here, Faster R-CNN is used for object detection.

MMDetection toolkit [6] with the same (default) settings. Specifically, the shorter side of input images are resized to 800, and SGD is used to optimize the networks with a weight decay of $1e-4$, a momentum of 0.9 and a mini-batch size of 8. All detectors are trained within 12 epochs on train2017 of COCO, where the learning rate is initialized to 0.01 and is decreased by a factor of 10 after 8 and 11 epochs, respectively. The results on val2017 are reported for comparison. As listed in Table 6, GCP_M outperforms GCP_D while both of them are superior to GAP. Specifically, for ResNet-50, GCP_M improves GAP by 0.7% in terms of AP using Faster R-CNN and Mask R-CNN as detectors. For ResNet-101, GCP_M outperforms GAP by 0.9% and 1.3% for Faster R-CNN and Mask R-CNN, respectively. For results of instance segmentation in Table 7, GCP_M improves GAP by 0.6% and 0.8% using ResNet-50 (R-50) and ResNet-101 (R-101) as backbone models, respectively. Note that GCP brings more improvement when backbone model and detector are stronger. These results show the pre-trained networks with GCP on large-scale datasets can well generalize to different vision tasks, indicating the networks with GCP can provide better initialization models.

As described above, integration of GCP into ResNets discards down-sampling (DS) operation in *conv5.1* to obtain more sampling features for more promising classification performance. However, it decreases the resolution of *conv5.x* and increases computing cost, especially for large-size input images. Here, we assess its effect on performance of GCP. Specifically, we employ ResNet-18, ResNet-50 and

ResNet-101 as backbone models, and compare the results of GCP with or without DS (i.e., *conv5.1* with a stride of 2). For object detection, we use Faster R-CNN as the basic detector. As shown in Table 8, GCP with DS is inferior to one without DS on ImageNet classification for each model, but achieves better performance on object detection. These results suggest that DS can be introduced for balancing classification accuracy, performance of object detection and model complexity for the networks trained with GCP.

5. Conclusion

In this paper, we made an attempt to analyze the effectiveness of GCP on deep CNNs from an optimization perspective. Specifically, we showed that GCP has the ability to improve Lipschitzness of loss and predictiveness of gradient in context of deep CNNs, and discussed the connection between GCP and second-order optimization. Our findings can account for several merits of GCP for training deep CNNs that have not been recognized previously or fully explored, including significant acceleration of network convergence, stronger robustness to distorted examples, and good generalization ability to different vision tasks. The extensive experimental results provide strong support to our findings. Our work provides an inspiring view to understand GCP and may help researchers explore more merits of GCP in context of deep CNNs. In future, we will investigate the theoretical proofs on smoothing effect of GCP and the rigorous connection to second-order optimization.

References

- [1] Naman Agarwal, Brian Bullins, and Elad Hazan. Second-order stochastic optimization for machine learning in linear time. *J. Mach. Learn. Res.*, 18:116:1–116:40, 2017.
- [2] Johan Bjorck, Carla P. Gomes, Bart Selman, and Kilian Q. Weinberger. Understanding batch normalization. In *NeurIPS*, 2018.
- [3] Sijia Cai, Wangmeng Zuo, and Lei Zhang. Higher-order integration of hierarchical convolutional activations for fine-grained visual categorization. In *ICCV*, 2017.
- [4] Nicholas Carlini, Guy Katz, Clark W. Barrett, and David L. Dill. Ground-truth adversarial examples. *arXiv*, 1709.10207, 2017.
- [5] J. Carreira, R. Caseiro, J. Batista, and C. Sminchisescu. Free-form region description with second-order pooling. *IEEE Trans. Pattern Anal. Mach. Intell.*, 37(6):1177–1189, 2015.
- [6] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv:1906.07155*, 2019.
- [7] Moustapha Cissé, Piotr Bojanowski, Edouard Grave, Yann N. Dauphin, and Nicolas Usunier. Parseval networks: Improving robustness to adversarial examples. In *ICML*, 2017.
- [8] Yin Cui, Feng Zhou, Jiang Wang, Xiao Liu, Yuanqing Lin, and Serge Belongie. Kernel pooling for convolutional neural networks. In *CVPR*, 2017.
- [9] Xiyang Dai, Joe Yue-Hei Ng, and Larry S. Davis. FASON: First and second order information fusion network for texture recognition. In *CVPR*, 2017.
- [10] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [11] Ali Diba, Vivek Sharma, and Luc Van Gool. Deep temporal linear encoding networks. In *CVPR*, 2017.
- [12] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. In *EMNLP*, 2016.
- [13] Yang Gao, Oscar Beijbom, Ning Zhang, and Trevor Darrell. Compact bilinear pooling. In *CVPR*, 2016.
- [14] Mengran Gou, Fei Xiong, Octavia Camps, and Mario Sznaier. Monet: Moments embedding network. In *CVPR*, 2018.
- [15] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask R-CNN. In *ICCV*, 2017.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [17] Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *ICLR*, 2019.
- [18] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. *arXiv:1503.02531*, 2015.
- [19] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- [20] Catalin Ionescu, Orestis Vantzos, and Cristian Sminchisescu. Matrix backpropagation for deep networks with structured layers. In *ICCV*, 2015.
- [21] Catalin Ionescu, Orestis Vantzos, and Cristian Sminchisescu. Training deep networks with structured layers by matrix backpropagation. *arXiv*, 1509.07838, 2015.
- [22] Shu Kong and Charles Fowlkes. Low-rank bilinear pooling for fine-grained classification. In *CVPR*, 2017.
- [23] Alex Krizhevsky. Learning multiple layers of features from tiny images, 2009.
- [24] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [25] Peihua Li, Jiangtao Xie, Qilong Wang, and Zilin Gao. Towards faster training of global covariance pooling networks by iterative matrix square root normalization. In *CVPR*, 2018.
- [26] Peihua Li, Jiangtao Xie, Qilong Wang, and Wangmeng Zuo. Is second-order information helpful for large-scale visual recognition? In *ICCV*, 2017.
- [27] Tsung-Yi Lin, Piotr Dollár, Ross B. Girshick, Kaiming He, Bharath Hariharan, and Serge J. Belongie. Feature pyramid networks for object detection. In *CVPR*, pages 936–944, 2017.
- [28] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, 2014.
- [29] Tsung-Yu Lin and Subhransu Maji. Improved bilinear pooling with CNNs. In *BMVC*, 2017.
- [30] Tsung-Yu Lin, Aruni RoyChowdhury, and Subhransu Maji. Bilinear CNN models for fine-grained visual recognition. In *ICCV*, 2015.
- [31] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet V2: practical guidelines for efficient CNN architecture design. In *ECCV*, 2018.
- [32] James Martens. Deep learning via hessian-free optimization. In *ICML*, 2010.
- [33] James Martens and Roger B. Grosse. Optimizing neural networks with Kronecker-factored approximate curvature. In *ICML*, 2015.
- [34] Jan Hendrik Metzen, Tim Genewein, Volker Fischer, and Bastian Bischoff. On detecting adversarial perturbations. In *ICLR*, 2017.
- [35] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *ICLR*, 2018.
- [36] Kazuki Osawa, Yohei Tsuji, Yuichiro Ueno, Akira Naruse, Rio Yokota, and Satoshi Matsuoka. Large-scale distributed second-order optimization using Kronecker-factored approximate curvature for deep convolutional neural networks. In *CVPR*, 2019.

- [37] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(6):1137–1149, 2017.
- [38] Mark Sandler, Andrew G. Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, 2018.
- [39] Shibani Santurkar, Dimitris Tsipras, Andrew Ilyas, and Aleksander Madry. How does batch normalization help optimization? In *NeurIPS*, 2018.
- [40] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [41] Qilong Wang, Peihua Li, and Lei Zhang. G²DeNet: Global Gaussian distribution embedding network and its application to visual recognition. In *CVPR*, 2017.
- [42] Kaicheng Yu and Mathieu Salzmann. Statistically-motivated second-order pooling. In *ECCV*, 2018.
- [43] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014.
- [44] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *ICLR*, 2017.
- [45] Bolei Zhou, David Bau, Aude Oliva, and Antonio Torralba. Interpreting deep visual representations via network dissection. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(9):2131–2145, 2019.

Appendix I: Implementation Details for Analyzing Smoothing Effect of GCP

In Section 3.2, we analyze smoothing effect of GCP on deep CNNs in terms of the Lipschitzness of optimization loss and the predictiveness of gradients. Specifically, the Lipschitzness of optimization loss is measured by

$$\Delta_l = \mathcal{L}(\mathbf{X} + \eta_l \nabla \mathcal{L}(\mathbf{X})), \eta_l \in [a, b], \quad (10)$$

and the predictiveness of gradients is measured by

$$\Delta_g = \|\nabla \mathcal{L}(\mathbf{X}) - \nabla \mathcal{L}(\mathbf{X} + \eta_g \nabla \mathcal{L}(\mathbf{X}))\|_2, \eta_g \in [a, b], \quad (11)$$

where \mathbf{X} is the input; $\nabla \mathcal{L}(\mathbf{X})$ indicates the gradient of loss with respect to the input \mathbf{X} ; η_l and η_g indicate step sizes of gradient descent.

To assess effect of GCP on the whole CNN models following [39], we employ output of the first convolution layer as \mathbf{X} to compute Eqns. (10) and (11). Note that the experiments in Section 4.2 demonstrate that the networks with GCP is more robust to input images with perturbations, comparing with those based on GAP. Accordingly, optimization loss of the networks with GCP also is more stable to input images with perturbations. For clear illustration, we calculate the ranges of Δ_l and Δ_g every 1,000 and 500 training steps for MobileNetV2 and ResNet-18, respectively. For calculating the ranges of Δ_l and Δ_g , we uniformly sample 50 points of η_l (and η_g) from [0.045, 1.5]

and [0.1, 75] for MobileNetV2 and ResNet-18, respectively. Then, we plot the ranges of Δ_l and Δ_g determined by the minimum and maximum of the 50 sampled points.

Appendix II: Derivations of Eqn. (6) and Eqn. (7)

As described in Section 3.3, the gradient of the loss with respect to the input \mathbf{X} through GCP layer can be calculated as

$$\frac{\partial \mathcal{L}}{\partial \mathbf{X}} = 2\mathbf{J}\mathbf{X} \left[\mathbf{U} \left(\left(\mathbf{K}^T \circ \left(\mathbf{U}^T 2 \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \right)_{\text{sym}} \mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}} \right) \right) + \left(\frac{1}{2} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{U}^T \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \mathbf{U} \right)_{\text{diag}} \right) \mathbf{U}^T \right]. \quad (12)$$

Here, we give detailed derivations of Eqn. (12) as follows. To perform GCP, we compute the square root of sample covariance matrix of features $\mathbf{X} \in \mathbb{R}^{N \times D}$ as

$$\mathbf{Z}_{\text{GCP}} = \mathbf{\Sigma}^{\frac{1}{2}} = (\mathbf{X}^T \mathbf{J} \mathbf{X})^{\frac{1}{2}} = \mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}^T, \quad (13)$$

where \mathbf{U} and $\mathbf{\Lambda}$ are the matrix of eigenvectors and the diagonal matrix of eigenvalues of sample covariance $\mathbf{\Sigma}$, respectively. As shown in [26], $\frac{\partial \mathcal{L}}{\partial \mathbf{X}}$ can be calculated as

$$\frac{\partial \mathcal{L}}{\partial \mathbf{X}} = \mathbf{J}\mathbf{X} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{\Sigma}} + \left(\frac{\partial \mathcal{L}}{\partial \mathbf{\Sigma}} \right)^T \right), \quad (14)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{\Sigma}} = \mathbf{U} \left(\left(\mathbf{K}^T \circ \left(\mathbf{U}^T \frac{\partial \mathcal{L}}{\partial \mathbf{U}} \right) \right) + \left(\frac{\partial \mathcal{L}}{\partial \mathbf{\Lambda}} \right)_{\text{diag}} \right) \mathbf{U}^T, \quad (15)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{U}} = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} + \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \right)^T \right) \mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}}, \quad (16)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{\Lambda}} = \frac{1}{2} \left(\mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{U}^T \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \mathbf{U} \right)_{\text{diag}}. \quad (17)$$

Let $(\mathbf{A})_{\text{sym}} = \frac{1}{2}(\mathbf{A} + \mathbf{A}^T)$, we can rewrite Eqn. (14) and Eqn. (16) as

$$\frac{\partial \mathcal{L}}{\partial \mathbf{X}} = 2\mathbf{J}\mathbf{X} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{\Sigma}} \right)_{\text{sym}}, \quad (18)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{U}} = 2 \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \right)_{\text{sym}} \mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}}. \quad (19)$$

By substituting Eqns. (15), (17) and (19) into Eqn. (18), we achieve

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \mathbf{X}} &= 2\mathbf{J}\mathbf{X} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{\Sigma}} \right)_{\text{sym}} \\
&= 2\mathbf{J}\mathbf{X} \left(\mathbf{U} \left(\left(\mathbf{K}^T \circ \left(\mathbf{U}^T \frac{\partial \mathcal{L}}{\partial \mathbf{U}} \right) \right) + \left(\frac{\partial \mathcal{L}}{\partial \mathbf{\Lambda}} \right)_{\text{diag}} \right) \mathbf{U}^T \right)_{\text{sym}} \quad (20) \\
&= 2\mathbf{J}\mathbf{X} \left[\mathbf{U} \left(\left(\mathbf{K}^T \circ \left(\mathbf{U}^T 2 \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \right)_{\text{sym}} \mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}} \right) \right) \right. \right. \\
&\quad \left. \left. + \left(\frac{1}{2} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{U}^T \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \mathbf{U} \right)_{\text{diag}} \right) \mathbf{U}^T \right]_{\text{sym}}.
\end{aligned}$$

So far, we obtain Eqn. (12).

With some assumptions and simplification, Eqn. (12) can be trimmed as

$$\frac{\partial \mathcal{L}}{\partial \mathbf{X}} \approx 2\mathbf{J}\mathbf{X} \left(2\mathbf{K}^T \circ \mathbf{\Lambda}^{\frac{1}{2}} + \frac{1}{2} \mathbf{\Lambda}^{-\frac{1}{2}} \right) \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}}, \quad (21)$$

where \circ denotes matrix Hadamard product. In the following, we explain how we obtain Eqn. (21). Specifically, we simplify Eqn. (12) by neglecting $(\cdot)_{\text{sym}}$ and $(\cdot)_{\text{diag}}$ operations. Thus, Eqn. (12) can be approximated by

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \mathbf{X}} &\approx 2\mathbf{J}\mathbf{X} \left[\mathbf{U} \left(\mathbf{K}^T \circ \left(\mathbf{U}^T 2 \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}} \right) + \right. \right. \\
&\quad \left. \left. \frac{1}{2} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{U}^T \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \mathbf{U} \right) \mathbf{U}^T \right]. \quad (22)
\end{aligned}$$

Then, we assume that matrix multiplications between diagonal matrix $\mathbf{\Lambda}$ and orthogonal matrix \mathbf{U} (or symmetric matrix $\frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}}$) in Eqn. (22) satisfy the commutative law of multiplication. So Eqn. (22) can be trimmed as

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \mathbf{X}} &\approx 2\mathbf{J}\mathbf{X} \left[\mathbf{U} \left(\mathbf{K}^T \circ \left(\mathbf{U}^T 2 \mathbf{\Lambda}^{\frac{1}{2}} \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \mathbf{U} \right) + \right. \right. \\
&\quad \left. \left. \frac{1}{2} \mathbf{U}^T \mathbf{\Lambda}^{-\frac{1}{2}} \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \mathbf{U} \right) \mathbf{U}^T \right]. \quad (23)
\end{aligned}$$

Finally, we assume that the mask matrix \mathbf{K} only has effect on the diagonal matrix of eigenvalues $\mathbf{\Lambda}$. So we have

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \mathbf{X}} &\approx 2\mathbf{J}\mathbf{X} \left[\mathbf{U} \left(\mathbf{U}^T 2 \mathbf{K}^T \circ \mathbf{\Lambda}^{\frac{1}{2}} \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \mathbf{U} + \right. \right. \\
&\quad \left. \left. \frac{1}{2} \mathbf{U}^T \mathbf{\Lambda}^{-\frac{1}{2}} \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \mathbf{U} \right) \mathbf{U}^T \right] \\
&= 2\mathbf{J}\mathbf{X} \left[\mathbf{U} \left(\mathbf{U}^T \left(2\mathbf{K}^T \circ \mathbf{\Lambda}^{\frac{1}{2}} \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} + \right. \right. \right. \\
&\quad \left. \left. \frac{1}{2} \mathbf{\Lambda}^{-\frac{1}{2}} \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}} \right) \mathbf{U} \right) \mathbf{U}^T \right] \quad (24) \\
&= 2\mathbf{J}\mathbf{X} \left(2\mathbf{K}^T \circ \mathbf{\Lambda}^{\frac{1}{2}} + \frac{1}{2} \mathbf{\Lambda}^{-\frac{1}{2}} \right) \frac{\partial \mathcal{L}}{\partial \mathbf{Z}_{\text{GCP}}}.
\end{aligned}$$

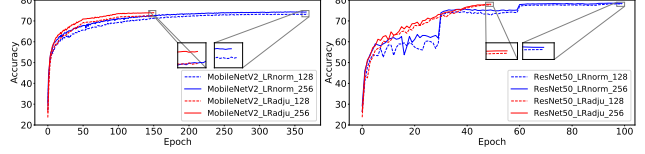


Figure 5. Convergence curves of MobileNetV2 and ResNet-50 trained with GCP under different settings of lr and various dimension (D) of input features on ImageNet.

Note that, in practice, Eqn. (24) is not employed for back-propagation of GCP, but provides a simplified form of Eqn. (12) for discussion on connection with second-order optimization in context of deep CNNs.

Appendix III: Convergence Curves of Networks with GCP under Various Dimensions of Input

In Table 5 of Section 4.1, we gave the results of MobileNetV2 and ResNet-50 with GCP under different settings of lr (i.e., LR_{norm} and LR_{adj}) and various dimension (i.e., $D = 256$ and $D = 128$) of input features on ImageNet. Figure 5 illustrates their corresponding convergence curves, from which we can see that lower-dimensional covariance representations (COV-Reps) share similar behavior with higher-dimensional COV-Reps, but the lower-dimensional COV-Reps suffer from larger performance degradation in the case of faster convergence (i.e., LR_{adj}).

Appendix IV: Implementation Details on Applying Pre-trained Networks with GCP to Other Vision Tasks

To apply the pre-trained networks with GCP to object detection and instance segmentation on MS COCO, we adopt the same strategy with the original GAP-based CNN models [15, 37] and make a modification, i.e., increasing resolution of feature maps in the last stage. The detailed steps are described as follows. All detectors are implemented using MMDetection toolkit [6].

- S.I: Pre-training the networks with GCP on ImageNet [10] without down-sampling in *conv5_1* as suggested in [26];
- S.II: Discarding the GCP layer and the classifier, while introducing Region Proposal Networks (RPN) [37] and Region of Interest (ROI) Pooling [15, 37];
- S.III: Increasing resolution of feature maps in the last stage using GCP_D (i.e., use of down-sampling as done in the original ResNet) and GCP_M (i.e., a max-pooling

layer with a step size 2 is inserted before *conv5_1*) strategies, while introducing feature pyramid networks (FPN) [27];

S_IV: Fine-tuning the whole networks in S_III on MS COCO [28] using the same hyper-parameters with those of the original GAP-based CNN models.

Appendix V: Computational Comparison of GCP and GAP

Here, we compare GCP and GAP in terms of computational cost. The experiments are conducted on large-scale ImageNet using ResNet-18, ResNet-34, ResNet-50 and ResNet-101 as backbone models. The evaluation metrics include network parameters, floating point operations per second (FLOPs), training or inference time per image, and Top-1/Top-5 accuracies. For GCP, size of covariance representations is set to 8k. All models are trained with the same experimental settings and run on a workstation equipped with four Titan Xp GPUs, two Intel(R) Xeon Silver 4112 CPUs @ 2.60GHz, 64G RAM and 480 GB INTEL SSD. From the results in Table 9, we can see that GCP introduces extra $\sim 7\text{M}$ parameters, $\sim 0.4\text{ms}$ training time and $\sim 0.2\text{ms}$ inference time, but increase about 4.6%, 2.6%, 1.5% and 1.8% Top-1 accuracies over GAP-based ResNet-18, ResNet-34, ResNet-50 and ResNet-101, respectively. Besides, GCP achieves matching performance using much lower computational complexity than GAP (e.g., ResNet34+GCP *vs.* ResNet101+GAP and ResNet50+GCP *vs.* ResNet152+GAP). Additionally, GCP with similar computational complexity achieves much better performance than GAP (e.g., ResNet18+GCP *vs.* ResNet34+GAP and ResNet50+GCP *vs.* ResNet101+GAP). Note that we discard down-sampling operation in *conv5_x* for GCP with ResNets, which significantly increases FLOPs. When we use this down-sampling operation, GCP shares similar FLOPs with GAP, leading slight performance decrease.

Table 9. Comparison of GCP and GAP using various ResNets in terms of network parameters, floating point operations per second (FLOPs), training or inference time per image, and classification accuracy.

Methods	Parameter	GFLOPs.	Training time (ms)	Inference time (ms)	Top-1 Err. (%)	Top-5 Err. (%)
ResNet18 + GAP	11.69M	1.81	0.77	0.60	70.47	89.59
ResNet18 + GCP	19.60M	3.11	1.21	0.85	75.07	92.14
ResNet34 + GAP	21.80M	3.66	1.17	0.88	74.19	91.60
ResNet34 + GCP	29.71M	5.56	1.61	1.10	76.80	93.11
ResNet50 + GAP	25.56M	3.86	1.85	1.29	76.02	92.97
ResNet50 + GCP	32.32M	6.19	2.22	1.49	78.56	93.72
ResNet101 + GAP	44.55M	7.57	2.79	1.72	77.67	93.83
ResNet101 + GCP	51.31M	9.90	3.14	1.83	79.47	94.30
ResNet152 + GAP	60.19M	11.28	3.54	2.55	78.13	94.04