

Places205-VGGNet Models for Scene Recognition

Limin Wang^{1,2} Sheng Guo¹ Weilin Huang^{1,2} Yu Qiao^{1,2}

¹Shenzhen Institutes of Advanced Technology, CAS, China

²The Chinese University of Hong Kong, Hong Kong, China

07wanglimin@gmail.com, {sheng.guo, wl.huang, yu.qiao}@siat.ac.cn

Abstract

VGGNets have turned out to be effective for object recognition in still images. However, it is unable to yield good performance by directly adapting the VGGNet models trained on the ImageNet dataset for scene recognition. This report describes our implementation of training the VGGNets on the large-scale Places205 dataset. Specifically, we train three VGGNet models, namely VGGNet-11, VGGNet-13, and VGGNet-16, by using a Multi-GPU extension of Caffe toolbox with high computational efficiency. We verify the performance of trained Places205-VGGNet models on three datasets: MIT67, SUN397, and Places205. Our trained models achieve the state-of-the-art performance on these datasets and are made public available¹.

1. Introduction

Convolutional networks (ConvNets) [5] have achieved great success for image classification [4]. In the recent ILSVRC competition [7], several successful architectures were proposed for object recognition, such as GoogLeNet [9] and VGGNet [8]. However, directly adapting these models trained on the ImageNet dataset [2] to the task of scene recognition cannot yield good performance. Besides training complicated VGGNets on a large-scale scene dataset is non-trivial, which requires large computational resource and numerous training skills. In this report, we train high-performance VGGNet models for scene recognition on the Places205 dataset [13]. The contribution of this report is twofold:

- Our trained Places205-VGGNet models achieve the state-of-the-art performance on the Places205 dataset [13]. As the training of VGGNet is very time consuming, we release our models to advance the further research on scene recognition.
- We transfer the trained models to other scene datasets, including the MIT67 [6] and SUN397 [12], and ex-

tract ConvNet features off-the-shelf. Our trained Places205-VGGNet models achieve the best performance on these two datasets.

2. Implementation Details

The VGGNets are originally developed for object recognition and detection [8]. They have very deep convolutional architectures with smaller sizes of convolutional kernel (3×3), stride (1×1), and pooling window (2×2). There are four different network structures, ranging from 11 layers to 19 layers. The model capability is increased when the network goes deeper, but imposing a heavier computational cost. Following original implementation of [8], we start with training an 11-layer VGGNet, and then train deeper VGGNets subsequently by using the pre-trained 11-layer model for initialization.

Specifically, we implement ConvNets by using the public Caffe toolbox [3]. As the computational cost and memory consumption of VGGNets are much larger than other architectures (e.g. GoogLeNet), we use Multi-GPU extension of Caffe [11], which is publicly available². Meanwhile, this extension provides more data augmentation techniques, such as *corner cropping strategy* and *multi-scale cropping method*, which have been proved to be effective for action recognition in videos. Therefore we also adopt these two augmentation techniques.

The training of ConvNets is performed with mini-batch gradient descent method, where the batch size is set to 256 and the momentum is 0.9. To reduce the effect of overfitting, the training was regularized by weight decay (the L2 penalty multiplier set to 0.0005) and dropout for the first two fully connected layers (with ratio of 0.5). During training phase, the images are resized to 256×256 . For multi-scale training, we randomly select the width and height of cropped regions from $\{256, 224, 198, 168\}$. These cropped regions are then resized to 224×224 for further processing. We start with training the 11-layer VGGNet, where network weights are randomly initialized with Gaussian distribution

¹<https://github.com/wanglimin/Places205-VGGNet/tree/master> ²https://github.com/yjxiong/caffe/tree/action_recog

Method	top-1 val/test	top-5 val/test
Places205-AlexNet [13]	50.4/50.0	80.9/81.1
Places205-GoogLeNet [1]	-/55.5	-/85.7
Places205-CNDS-8 [10]	54.7/55.7	84.1/85.8
Places205-VGGNet-11	58.6/59.0	87.6/87.6
Places205-VGGNet-13	60.2/60.1	88.1/88.5
Places205-VGGNet-16	60.6/60.3	88.5/88.8

Table 1. Performance comparison of different network architectures on the dataset of Places205.

(mean set to 0 and deviation set to 0.01). The learning rate is initially set as 0.01, and decreased to its $\frac{1}{10}$ every 10k iterations. The whole training process stops at 40k iterations. To train the 13-layer and 16-layer VGGNets, we initialize the first four convolutional layers and first two fully connected layers with the pre-trained 11-layer VGGNet.

For testing the VGGNet models, we follow multi-view classification method [4]. Specifically, we randomly crop regions of 224×224 from four corners and center of the image, whose size is 256×256 . After that, these cropped regions are horizontally flipped. Therefore, we obtain 10 views, each of which is fed into ConvNet models for prediction. The final prediction score is the average value of the 10 predictions.

3. Experiments

In this section, we describe our experimental details and results. The training of VGGNets on the Places205 dataset is implemented with a Multi-GPU extension of Caffe [11]. In our experiment, we use 4 GTX Titan-X GPUs and the whole training time of VGGNet-16 is around 2 weeks. To test the performance of our trained Places205-VGGNet models, we conduct experiments on three datasets, namely Places205, MIT67, and SUN397.

First we perform evaluation on the Places205 [13] and the results are summarized in Table 1. We compare with other deep network architectures, like AlexNet [4], GoogLeNet [9], and CNDS-8 [10], and observe that VGGNets obtain much better performance than theirs on this dataset.

To further verify the effectiveness of Places205-VGGNet models on scene recognition, we transfer the learned representations to the MIT67 [6] and SUN397 [12] datasets. Specifically, we extract fc6 features and normalize them with ℓ_2 -norm. Then we employ linear SVMs as classifiers for scene category prediction. The experimental results are shown in Table 2. We compare our Places205-VGGNet models with other public model and our models achieve the best performance on these two challenging datasets.

Model	MIT67	SUN397
ImageNet-VGGNet-16 [8]	67.7	51.7
Places205-AlexNet [13]	68.2	54.3
Places205-CNDS-8 [10]	76.1	60.7
Places205-GoogLeNet [1]	76.3	61.1
Places205-VGGNet-11	82.0	65.3
Places205-VGGNet-13	81.9	66.7
Places205-VGGNet-16	81.2	66.9

Table 2. Performance comparison of transferred representations from different models on the MIT67 and SUN397 datasets.

4. Conclusions

In this report, we describe our implementation of training the VGGNets on the large-scale Places205 dataset with a Multi-GPU extension of Caffe. The trained Places205-VGGNet models achieve the state-of-the-art performance on three scene recognition benchmarks, namely Places205, SUN397, and MIT67. We release our trained Places205-VGGNet models for further research in scene recognition.

References

- [1] <http://places.csail.mit.edu/user/leaderboard.php>. 2
- [2] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and F. Li. ImageNet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009. 1
- [3] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. B. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *CoRR*, abs/1408.5093. 1
- [4] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *NIPS*, pages 1106–1114, 2012. 1, 2
- [5] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 1
- [6] A. Quattoni and A. Torralba. Recognizing indoor scenes. In *CVPR*, pages 413–420, 2009. 1, 2
- [7] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. S. Bernstein, A. C. Berg, and F. Li. ImageNet large scale visual recognition challenge. *CoRR*, abs/1409.0575, 2014. 1
- [8] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 1, 2
- [9] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. abs/1409.4842, 2014. 1, 2
- [10] L. Wang, C. Lee, Z. Tu, and S. Lazebnik. Training deeper convolutional networks with deep supervision. *CoRR*, abs/1505.02496, 2015. 2
- [11] L. Wang, Y. Xiong, Z. Wang, and Y. Qiao. Towards good practices for very deep two-stream ConvNets. *CoRR*, abs/1507.02159, 2015. 1, 2
- [12] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. SUN database: Large-scale scene recognition from abbey to zoo. In *CVPR*, pages 3485–3492, 2010. 1, 2
- [13] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *NIPS*, pages 487–495, 2014. 1, 2