# Multilevel Context Representation for Improving Object Recognition

Andreas Kölsch[1], Muhammad Zeshan Afzal[1,2] and Marcus Liwicki[1,3]

`a_koelsch12@cs.uni-kl.de, afzal@iupr.com, marcus.liwicki@unifr.ch`

[1]MindGarage, University of Kaiserslautern
[2]Insiders Technologies GmbH
[3]University of Fribourg

## Abstract

*In this work, we propose the combined usage of low- and high-level blocks of convolutional neural networks (CNNs) for improving object recognition. While recent research focused on either propagating the context from all layers, e.g. ResNet, (including the very low-level layers) or having multiple loss layers (e.g. GoogLeNet), the importance of the features close to the higher layers is ignored. This paper postulates that the use of context closer to the high-level layers provides the scale and translation invariance and works better than using the top layer only. In particular, we extend AlexNet and GoogLeNet by additional connections in the top $n$ layers. In order to demonstrate the effectiveness of the proposed approach, we evaluated it on the standard ImageNet task. The relative reduction of the classification error is around 1-2% without affecting the computational cost. Furthermore, we show that this approach is orthogonal to typical test data augmentation techniques, as recently introduced by Szegedy et al. (leading to a runtime reduction of 144 during test time).*

## 1. Introduction

While it is quite easy for humans to distinguish between objects in an image, it can be a very challenging task for a computer. Not only can objects appear in different sizes and angles, but also backgrounds and lighting conditions may vary and many other factors can change the way an object is displayed [44]. Furthermore, some object classes have a very low inter-class variance (cf. Fig. 2), while other classes might have a very high intra-class variance (cf. Fig. 3) [13, 10]. Consequently, both the details of an image and the greater context are important for successful classification. Another factor, that makes image classification hard is the presence of several objects in one image (cf. Fig. 4). When
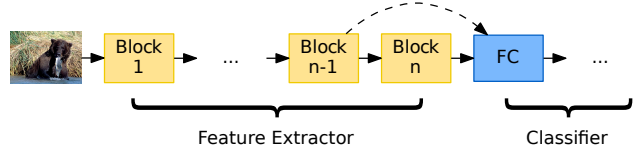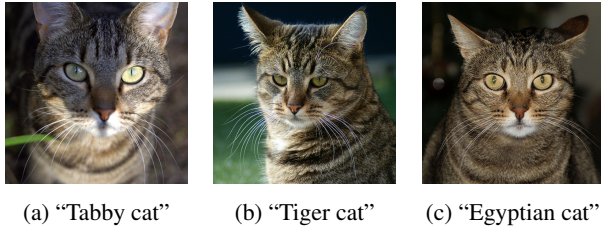


Figure 1: Contribution of this work: Instead of just using the last CNN layer (Block $n$), we use connections to more high-level layers as well in order to include various representations in the recognition.

it comes to realistic applications of object recognition models, the task is exacerbated even further by computational constraints. Some applications, e.g. in autonomous cars, require very fast, yet accurate predictions [40, 23], while other applications, especially those in mobile devices, constrain the amount of available memory and processing power.

In the last years, neural networks have had great success in many tasks related to computer vision. For the task of object recognition, particularly convolutional neural networks (CNNs) have become the state of the art [11]. Thanks to now available computational resources and efficient GPU implementations [17, 5], it is also possible to train very deep CNNs. In 2012, Krizhevsky et al. proposed an eight-layer deep CNN [19], which, due to its groundbreaking performance, laid the foundation of using these networks for image classification. Two years later, 19-layer [38] and 22-layer networks [39] were presented, which performed even better. In 2015, He et al. proposed an even deeper network architecture with 152 layers [14], and most recently the development of ever deeper CNNs culminated with PolyNet with several hundred layers [47] (see Section 2).

While this trend towards deeper networks is evident, the basic architecture of CNNs has never changed fundamentally. They always consist of a number of stacked convolutional layers with a non-linear squashing function [16, 20,

1

(a) "Tabby cat"  (b) "Tiger cat"  (c) "Egyptian cat"

Figure 2: Three similar images from different classes



Figure 3: Even though they look quite different, all three of these images belong to the class "Tiger cat"



Figure 4: An image showing objects of different classes: miniskirt, Hat with a wide brim, sunglasses, plastic bag, purse, person

6], optional local response normalization and pooling for feature extraction [2, 4, 36], followed by a number of fully-connected layers for classification. As illustrated by Zeiler and Fergus, the features extracted by the convolutional layers in the network correspond to more complex patterns from layer to layer [45]. While the features of the lower level layers correspond to dots or edges, the higher level features correspond to patterns, such as faces or entire objects. In the case of AlexNet [19], Zeiler and Fergus describe the features extracted by the two highest level convolutional layers of the network, as highly "class-specific" [45].

However, to the best of the authors' knowledge, it has not yet been tried to make the information from multiple high-level layers directly available to the classification part of the network to test, whether or not the additional features provide complementary information. This is the main idea of this paper. As such, we hope that providing the classification part of the network with more data leads to better classification results while keeping the additional computational costs for both training and prediction minimal.

The contributions of this work are three-fold. Firstly, we perform some initial experiments motivating the idea of adding connections from multiple levels. Secondly, we present two novel CNN architectures that explicitly use the features extracted by multiple convolutional layers, for classification. These networks do not rely on pre-training, but provide standalone solutions that can be trained from scratch. Thirdly, we evaluate both the performance and computational costs of the novel architectures and compare them to existing approaches.

The remaining sections are organized as follows. Section 2 gives an overview of related work. Section 3 shows the validity of the approach in a constrained setup. In
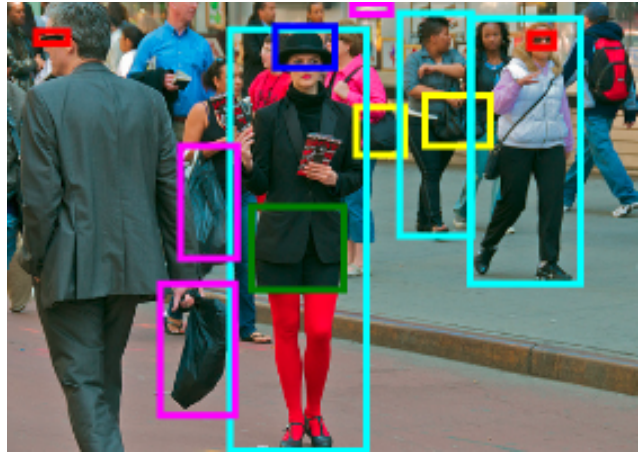
section 4 the novel network architectures and the training methodology are described in detail and large scale experiments are performed. The results are presented in section 5 which is followed by the discussion in section 6 and conclusion in section 7.

## 2. Related Work

Both Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have been used for solving key computer vision problems [21]. Some example areas such as object recognition and detection [19, 39, 14, 47, 42], semantic segmentation [3, 25, 27, 32], describing images with text [18, 43, 41, 7] and also generating the images using text [9, 28, 33, 34, 46] show the effectiveness of the deep learning models for solving challenging tasks. The RNNs did not receive much attention for object recognition and the work in this area is rather sparse due to the following reasons: First by the virtue of their architecture that is suitable for temporal data processing. Secondly, due to the temporal dependencies, the processing of RNNs is inefficient in comparison to CNNs. However, the work of Visin et al. [42] had signified that contextual processing of pixels can help to recognize the objects better. Their proposed architecture was named as ReNet and it was recommended as an alternative approach to CNNs. In another work, Biswas et al. [1] concatenated different variations (degradations) of the same object and processed it temporally. They used each variation of the object as one time step for RNNs. However, their experiments are limited to MNIST [22] and COIL-20 [31] data sets and the validity and the computational viability of their approach remains undiscovered for large data sets e.g., ImageNet [35]. A combination of CNNs and RNNs has also been reported by Linag et al. [24]. Their proposed ar-
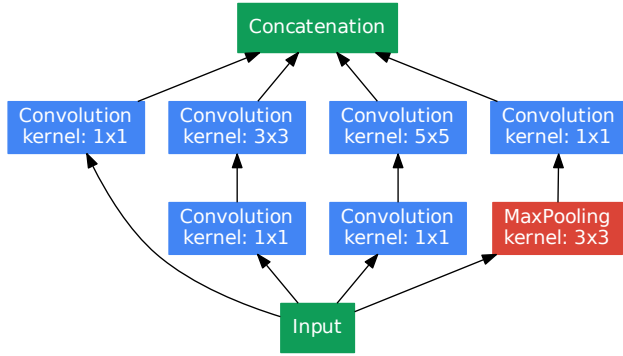
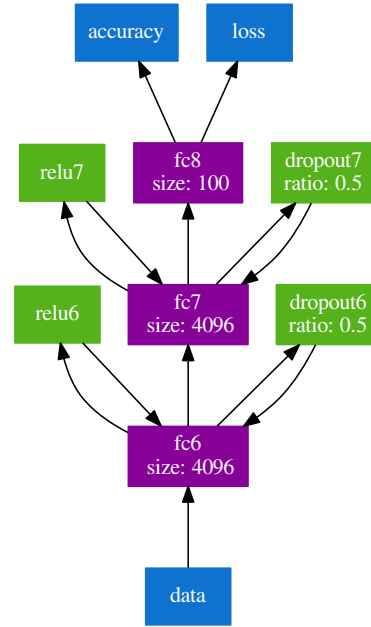Figure 5: Building block of the Inception architecture [39].



Figure 6: Network architecture used for the experiments presented in section 3. The data layer holds the features extracted beforehand. The fc8 layer is adjusted in size to match the number of classes used in the experiments.

chitecture relies both on the efficiency of CNNS and power of contextual processing of RNNs. As our proposed work aims to extend existing deep CNN architectures, the following paragraphs briefly describe the state-of-the-art deep CNN architectures for object recognition.

One well-known benchmark for several visual tasks, including image classification, is the annual ImageNet Large Scale Visual Recognition Challenge (ILSRVC) [35]. It allows the participating teams to compare the performance of their developed models. The data set used for the competition consists of 1.000 image classes containing 1.2 million labeled images for training, 50.000 labeled validation images and 100.000 unlabeled images which are used to compare the performance of the submitted entries.

In 2012, Krizhevsky et al. have, for the first time, used a deep CNN in the ILSVRC and won the image classification with a significant margin, outperforming all traditional methods [19]. The winning architecture consists of five convolutional layers which are followed by three fully connected layers. The first and second layer use local response normalizations and max-pooling before passing the activations to the next layer. The fifth convolutional layer is again followed by a max-pooling layer which provides the input to the fully-connected layers. To prevent overfitting, the fully-connected layers use dropout [16]. Rectified linear units (ReLU) are used as activation function in all layers [30]. AlexNet achieves a top-5 error rate of $16.422\%$ on the ILSVRC test data set.

Since then, many other teams have developed a variety of different CNN architectures to further reduce the classification error of these networks [14, 38, 39]. First of all, Szegedy et al. presented the 22-layer GoogLeNet in 2014, which achieved a top-5 error rate of $6.656\%$ on the ILSVRC 2014 test data set with an ensemble of 7 networks. The network is using the Inception architecture [39] which is different from other approaches, as it employs not stacked convolutional layers, but stacked building blocks which themselves consist of multiple convolutional layers (cf. Fig. 5).

Therefore, it is a network-in-network approach [26]. As the network is so deep, it employs two auxiliary loss layers during training to allow for efficient backpropagation of the error.

Note that the general idea of adding skip-connections (so-called residual connections) has been recently introduced by He et al. in [14, 15]. However, the motivation and effects are fundamentally different. In ResNet, the skip-connections go through all the layers, letting the layers mainly pass the information through the network with minor additive terms. This mainly overcomes the vanishing gradient problem and allows learning architectures with hundreds of layers. Contrary, the main idea in this paper is to use the extracted features of multiple layers directly for classification. The only work which comes close to our idea is the proposal of Center-Multilayer Features (CMF) as proposed by Seuret et al.[37]. They used stacked convolutional auto-encoders (SCAE) for classifying an image or the center pixel of a patch. This idea is similar to our pre-study (section 3) where we use multiple feature representations from fine-tuned networks. However, their network architecture is shallow and they use features from all the layers.

Apart from novel network architectures, a lot of work has been done to understand how the networks are learning and how they can be improved [12, 29]. Especially visualization techniques have helped to get an understanding of the

Table 1: Accuracy achieved by networks on validation sets after training with inputs from conv5 only and inputs from conv4 (normalized) and conv5.

|  | conv5 | conv4 & conv5 |
|---|---|---|
| 10 classes | 90.60% | 91.60% |
| 20 classes | 83.40% | 84.10% |
| 50 classes | 78.72% | 80.12% |
| 100 classes | 68.36% | 71.10% |

Table 2: Input size of the layers of AlexNet and AlexNet++

|  | AlexNet | AlexNet++ |
|---|---|---|
| conv1 | $3 \times 227 \times 227$ | $3 \times 227 \times 227$ |
| pool1 | $96 \times 55 \times 55$ | $96 \times 55 \times 55$ |
| conv2 | $96 \times 27 \times 27$ | $96 \times 27 \times 27$ |
| pool2 | $256 \times 27 \times 27$ | $256 \times 27 \times 27$ |
| conv3 | $256 \times 13 \times 13$ | $256 \times 13 \times 13$ |
| conv4 | $384 \times 13 \times 13$ | $384 \times 13 \times 13$ |
| conv5 | $384 \times 13 \times 13$ | $384 \times 13 \times 13$ |
| pool5 | $256 \times 13 \times 13$ | $256 \times 13 \times 13$ |
| fc6 | $256 \times 6 \times 6$ | $640 \times 6 \times 6$ |
| fc7 | 4096 | 4096 |
| fc8 | 4096 | 4096 |
| prob | 1000 | 1000 |

convolutional layers [45]. Despite all the research that has been done to develop new CNN architectures, novel architectures typically employ more layers or make use ensembling techniques [48]. However, both approaches typically require more computational resources at training and inference time.

## 3. Pre-study

To investigate, whether it is worth pursuing the idea of using features from multiple layers, we run a set of small-scale experiments. We use a publicly available model[1] of AlexNet that is pre-trained on ImageNet. This model is used to extract and store the activations generated by the fourth and fifth convolutional layers (conv4 and conv5) during a forward pass of the images. Then, we train a three-layer fully-connected neural network on these activations which, except for the last layer, equals the fully-connected part of AlexNet (cf. Fig. 6). As is common, the last layer is adjusted in size to match the number of classes. We repeat the experiment with several subsets of the ILSVRC data set. The network is trained and evaluated on the activations of 10, 20, 50 and 100 randomly selected classes of the data set.

To get a baseline, we run a set of experiments, in which the network is trained solely on the activations from conv5 that were extracted before. In another set of experiments, we use the same classes, but train the network on the activations of conv4 concatenated with conv5. It turns out, we have to L2-normalize the activations of conv4 to get useful results. With unmodified values, the network fails to learn.

Table 1 shows, that the networks trained on features from conv4 and conv5 yielded a better performance on the validation set than the networks which were trained on the features from conv5 only.

## 4. Large Scale Experiments

With the results from section 3 we step the experiments up and train entire CNNs from scratch on the ILSVRC data set.

### 4.1. Network Architectures

We propose two neural network architectures in this section. The network architectures are strict extensions of the existing networks AlexNet and GoogLeNet. Unlike most CNNs, including AlexNet and GoogLeNet, the proposed networks feed the classification part of the network with information not only from the highest-level convolutional layer, but with information from the two highest-level convolutional layers. We call the enhanced versions of these networks AlexNet++ and GoogLeNet++.

#### 4.1.1 AlexNet++

The original AlexNet is an eight-layer deep CNN in which each layer processes the features extracted by the subjacent layer only. As already said, the proposed AlexNet++ differs from the original AlexNet in the way, that the first fully connected layer (fc6) is not only presented with the activations from conv5, but also with activations from conv4 (cf. Fig. 7). Since after max-pooling, the conv5 layer has only $256 \times 6 \times 6$ values, but conv4 has $384 \times 13 \times 13$ values, we have to balance the number of activations passed to the fc6 layer. Therefore, we add a max-pooling of conv4, before concatenating the values with the conv5 values. Furthermore, to resemble the normalization that was done in the small scale experiments (cf. Section 3), we use *tanh* after the pooling to compress the activations. A detailed comparison of the layer sizes of AlexNet and AlexNet++ is given in table 2.

#### 4.1.2 GoogLeNet++

GoogLeNet is the 22-layer network architecture, that won the ILSVRC2014 [39]. Unlike AlexNet, the architecture consists of stacked building blocks which consist of convolutional layers (cf. Fig. 5). Due to its depth, the original

---

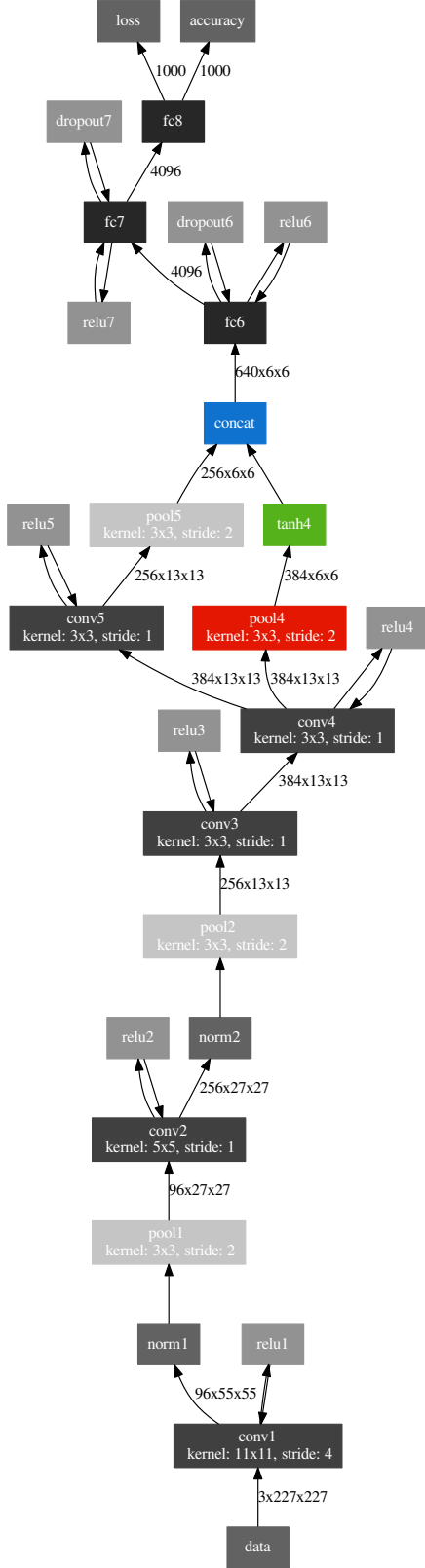[1] https://github.com/BVLC/caffe/tree/master/models/bvlc_alexnet

Figure 7: AlexNet++ architecture with unchanged AlexNet part in gray and the new part colored.
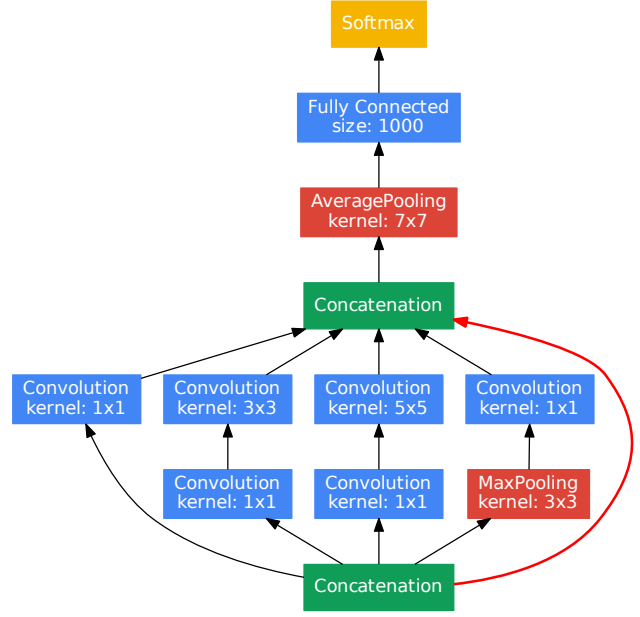


Figure 8: Top part of the GoogLeNet++ architecture. The new connection is marked red. The lower layers of the network are built as proposed by Szegedy et al. [39].

architecture employs three softmax classifiers for efficient error backpropagation during training. At inference time, only the main classifier is used. This classifier gets its input from one fully connected layer, which in turn processes the input from the highest level building block.

The proposed GoogLeNet++ architecture concatenates the activations of the two highest level building blocks, before they are passed to the fully connected layer (cf. Fig. 8). Table 3 gives a detailed comparison of the layer sizes of GoogLeNet and GoogLeNet++.

## 4.2. Training Details

We train a version of AlexNet++ and GoogLeNet++ from scratch on the ILSVRC training data. To allow for a fair comparison, the unchanged versions of AlexNet and GoogLeNet are also trained on the same training data and setup. All networks are trained on a single NVidia Tesla K20X using the Caffe framework [17].

AlexNet and AlexNet++ are trained for 90 epochs using stochastic gradient descent with an initial learning rate of 0.01 which is divided by 10 after 30 and 60 epochs. The momentum is set to 0.9 and the weight decay is 0.0005 as proposed by Krizhevsky et al. [19]. Both networks use a batch size of 256 images to speed up the training. As proposed in the original architecture, we initialize the weights with a gaussian distribution with a mean of zero and a standard deviation of 0.01. The biases are initialized differently from the original approach. While Krizhevsky et al. set the

Table 3: Input size of the layers of GoogLeNet and GoogLeNet++

|  | GoogLeNet | GoogLeNet++ |
|---|---|---|
| conv1 | $3 \times 224 \times 224$ | $3 \times 224 \times 224$ |
| pool1 | $64 \times 112 \times 112$ | $64 \times 112 \times 112$ |
| conv2 | $64 \times 56 \times 56$ | $64 \times 56 \times 56$ |
| pool2 | $192 \times 56 \times 56$ | $192 \times 56 \times 56$ |
| inception3a | $192 \times 28 \times 28$ | $192 \times 28 \times 28$ |
| inception3b | $256 \times 28 \times 28$ | $256 \times 28 \times 28$ |
| pool3 | $480 \times 28 \times 28$ | $480 \times 28 \times 28$ |
| inception4a | $480 \times 14 \times 14$ | $480 \times 14 \times 14$ |
| inception4b | $512 \times 14 \times 14$ | $512 \times 14 \times 14$ |
| inception4c | $512 \times 14 \times 14$ | $512 \times 14 \times 14$ |
| inception4d | $512 \times 14 \times 14$ | $512 \times 14 \times 14$ |
| inception4e | $528 \times 14 \times 14$ | $528 \times 14 \times 14$ |
| pool4 | $832 \times 14 \times 14$ | $832 \times 14 \times 14$ |
| inception5a | $832 \times 7 \times 7$ | $832 \times 7 \times 7$ |
| inception5b | $832 \times 7 \times 7$ | $832 \times 7 \times 7$ |
| pool5 | $1024 \times 7 \times 7$ | $1856 \times 7 \times 7$ |
| fc | $1024 \times 1 \times 1$ | $1856 \times 1 \times 1$ |
| prob | 1000 | 1000 |

biases of the second, fourth, and fifth convolutional layers and the biases of the fully-connected layers to 1, we set them to 0.1, as the network fails to learn with biases set to 1[2].

The GoogLeNet and GoogLeNet++ networks are trained for five million iterations with a batch size of 32, i.e. for 133 epochs. We train the networks with stochastic gradient descent with a polynomial update of the learning rate, as proposed by Sergio Guadarrama[3]. This means, the learning rate is updated at every iteration to

$$lr = initial\_lr * \left(1 - \frac{iter}{max\_iter}\right)^{power} \tag{1}$$

In our case, the initial learning rate is set to 0.01 and power is 0.5. The momentum is 0.9, weight decay is set to 0.0002. To initialize the network, we use *normalized initialization* [8].

#### 4.2.1 Data Augmentation

The four networks described above are trained on the ILSVRC training data set with no additional data. To artificially enlarge the data set, we use a rather simple data augmentation technique. Namely, we resize all images to $256 \times 256$ pixels, subtract the mean pixel value in a preprocessing step and then randomly crop patches in the size

(a) ILSVRC image
($2848 \times 2144$ pixels)

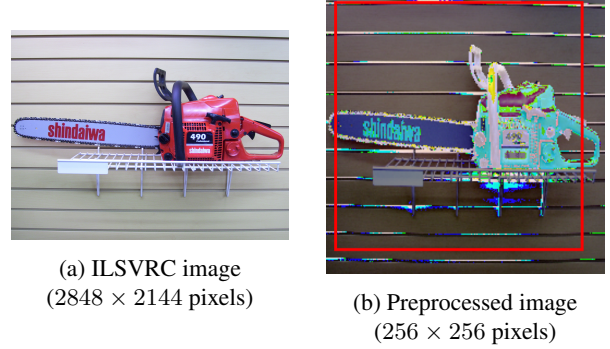(b) Preprocessed image
($256 \times 256$ pixels)

Figure 9: Images are scaled to $256 \times 256$ pixels, the mean pixel value is subtracted and the image is randomly cropped to the network's input size.

of the network input from these, i.e. $227 \times 227$ pixels for AlexNet, $224 \times 224$ pixels for GoogLeNet. Finally, we randomly mirror the images horizontally. An example is given in Fig. 9.

Note, that Krizhevsky et al. and Szegedy et al. use more aggressive data augmentation techniques. However, in the case of GoogLeNet, the given details are very vague and not described in a reproducible way. In the case of the original AlexNet approach, the rectangular images are preprocessed by scaling them such that the shorter edge measures 256 pixels and then a $256 \times 256$ patch is cropped from the center. This scaling preserves the aspect ratio, but does not use the outer regions of an image at all. We decided to not perform more aggressive data augmentation techniques as our idea is orthogonal to the augmentation ideas, i.e., even if more data augmentation is performed, it is still possible to use the Multilevel Context representation.

## 5. Benchmarking and Results

To compare our extended versions of AlexNet and GoogLeNet with the original ones, we use the labeled ILSVRC validation data set. Remember, we not only set out to improve the classification performance of the existing networks, but also, to keep the additional computational costs low. Therefore, this section covers both, the classification accuracy of the models and the time needed for classification and training.

We use two evaluation methods to measure the accuracy of the networks. First, we classify the validation images using a single-crop technique, i.e., we resize them to $256 \times 256$ pixels and use the center crop. Secondly, the networks are tested with the multi-crop technique proposed in [39]. Specifically, the images are resized such that the shorter dimension is 256, 288, 320 and 352 pixels long. Then 3 square crops are taken from each of these images, i.e., the left, center and right part for landscape pictures or the top,

Table 4: Accuracy achieved by networks on validation sets using center crops only and averaging over 144 crops as described by Szegedy et al. [39]

| | Top-1 | | Top-5 | |
| | center crop | 144 crops | center crop | 144 crops |
|---|---|---|---|---|
| AlexNet | 58.17% | 58.54% | 80.82% | 81.18% |
| AlexNet++ | 58.47% | 58.93% | 81.25% | 81.47% |
| GoogLeNet | 68.89% | 68.86% | 88.97% | 89.26% |
| GoogLeNet++ | 69.42% | 69.05% | 89.25% | 89.35% |

Table 5: Time needed for passing one batch through the network. The batch size for AlexNet and AlexNet++ is 256, the batch size for GoogLeNet and GoogLeNet++ is 32.

| | Forward | Backward | $\Sigma$ |
|---|---|---|---|
| AlexNet | 254.87 ms | 654.04 ms | 909.01 ms |
| AlexNet++ | 271.29 ms | 686.76 ms | 958.16 ms |
| GoogLeNet | 112.26 ms | 283.31 ms | 395.71 ms |
| GoogLeNet++ | 112.51 ms | 284.76 ms | 397.42 ms |

center and bottom part for portrait pictures. For each of these square images the 4 corner crops, the center crop and a resized image with the shape of the network input is produced. Finally, all of these images are mirrored as well. In sum, this multi-crop technique produces $4 \times 3 \times 6 \times 2 = 144$ crops. To evaluate the network performance with this approach, we average over the probability vectors of the 144 predictions.

Table 4 presents the accuracy on the validation data set which is achieved by the four trained networks. As can be seen, the extended networks outperform the original versions with all evaluation techniques. Furthermore, the accuracy achieved by the new models with only one crop is at least comparable to the accuracy achieved by the original models using 144 crops. In the case of GoogLeNet, the Top-1 accuracy of the new network using only one crop is even more than $0.5\%$ higher than the accuracy of the original network using 144 crops (this corresponds to around $2\%$ relative error reduction).

Table 5 shows the time needed for forward and backward passes through the different networks. Obviously, the timings are specific to hardware and software. In our case, all timings are made using caffe 1.0rc3 and a NVidia Tesla K20X. AlexNet++ takes about 5% longer for training, i.e., combined forward and backward pass. However, the accuracy achieved by AlexNet++ using just center crops is comparable and in the case of top-5 accuracy even higher than what AlexNet achieves using 144 crops (see Table 4). Given this result, it is actually reasonable to say AlexNet++ is more than 100 times faster to reach a comparable accuracy.

For GoogLeNet++, the difference in runtime for both forward and backward pass is negligibly small. This is due to the fact, that the GoogLeNet architecture already has a concatenation layer before the fully connected layer and thus, no new layers have to be added to combine the two highest level building blocks (cf. Fig. 8). A combined forward- backward pass with GoogLeNet++ is only $0.4\%$ slower than with the original GoogLeNet. Therefore, training the extended network takes only marginally longer. At inference time, the difference decreases to a mere $0.2\%$ which makes the GoogLeNet++ very well-suited for productive usage.

## 6. Discussion

Our approach postulates that the higher level layers are useful for achieving better performance both in terms of recognition accuracy and computational cost. In order to demonstrate the validity of the claim, we performed a simple experiment using AlexNet. In this experiment, the features from all of the convolutional layers were concatenated. Two observations were made: first, the training took significantly longer (almost 3 times) and secondly, the accuracy was suboptimal (almost 2% less in terms of absolute error). The best performance was achieved by concatenating the features of the top 2 layers.

The networks trained in this work could not reproduce the drastic error reduction that Szegedy et al. have reported [39] by using the $144$ image crops at test time. The main reason for this is that we have used less aggressive data augmentation during training or due to different weight initializations. Unfortunately, Szegedy et al. do not report the

details on how the GoogLeNet model was initialized. Also, they report error reduction by ensembling multiple trained networks. However, this is not in the scope of this paper. The main purpose of this paper is to show that features from lower level convolutional layers that are close to the highest level provide useful and complementary information for classification.

Noteworthy, as shown in this paper, when training our networks with the same data augmentation and hyper-parameters, the extended architectures presented in this work outperform the original ones. This applies to both networks (AlexNet and GoogLeNet) in both evaluation scenarios (no data augmentation during testing or using 144 test samples).

## 7. Conclusion

In this paper, we presented a successful approach for extending existing CNN architectures in order to boost their classification performance. We have demonstrated the effectiveness of this approach by enhancing the network architectures of AlexNet and GoogLeNet and training them from scratch on the ILSVRC data set. We consider networks that are totally different in nature to prove the generality of the proposed approach. Also, it is shown that at almost no additional cost, the relative error rates of the original networks decrease by up to 2%. This fact makes the extended networks a very well suited choice for usage in production environments. The quantitative evaluation signifies that the new approach could be, at inference time, 144 times more efficient than the current approaches while maintaining comparable performance. The proposed approach is not limited to any one of the architectures. We plan to extend the experiments for recurrent and convolutional-recurrent neural networks.

## References

[1] S. Biswas, M. Z. Afzal, and T. Breuel. Using recurrent networks for non-temporal classification tasks. In *Neural Networks (IJCNN), 2014 International Joint Conference on*, pages 135–140. IEEE, 2014. 2

[2] Y.-L. Boureau, J. Ponce, and Y. LeCun. A theoretical analysis of feature pooling in visual recognition. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 111–118, 2010. 2

[3] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062*, 2014. 2

[4] D. C. Ciresan, U. Meier, J. Masci, L. Maria Gambardella, and J. Schmidhuber. Flexible, high performance convolutional neural networks for image classification. In *IJCAI Proceedings-International Joint Conference on Artificial Intelligence*, volume 22, page 1237. Barcelona, Spain, 2011. 2

[5] R. Collobert, K. Kavukcuoglu, and C. Farabet. Torch7: A matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*, 2011. 1

[6] G. E. Dahl, T. N. Sainath, and G. E. Hinton. Improving deep neural networks for lvcsr using rectified linear units and dropout. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 8609–8613. IEEE, 2013. 1

[7] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2625–2634, 2015. 2

[8] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Aistats*, volume 9, pages 249–256, 2010. 6

[9] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 2

[10] C. Göring, A. Freytag, E. Rodner, and J. Denzler. Fine-grained categorization–short summary of our entry for the imagenet challenge 2012. *arXiv preprint arXiv:1310.4759*, 2013. 1

[11] J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu, X. Wang, and G. Wang. Recent advances in convolutional neural networks. *arXiv preprint arXiv:1512.07108*, 2015. 1

[12] K. He and J. Sun. Convolutional neural networks at constrained time cost. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5353–5360, 2015. 3

[13] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015. 1

[14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 1, 2, 3

[15] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *European Conference on Computer Vision*, pages 630–645. Springer, 2016. 3

[16] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012. 1, 3

[17] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014. 1, 5

[18] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137, 2015. 2

[19] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In

*Advances in neural information processing systems*, pages 1097–1105, 2012. 1, 2, 3, 5

[20] Y. LeCun, Y. Bengio, et al. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10):1995, 1995. 1

[21] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015. 2

[22] Y. LeCun and C. Cortes. MNIST handwritten digit database. 2010. 2

[23] Y. LeCun, K. Kavukcuoglu, and C. Farabet. Convolutional networks and applications in vision. In *Circuits and Systems (ISCAS), Proceedings of 2010 IEEE International Symposium on*, pages 253–256. IEEE, 2010. 1

[24] M. Liang and X. Hu. Recurrent convolutional neural network for object recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3367–3375, 2015. 2

[25] G. Lin, C. Shen, A. van den Hengel, and I. Reid. Efficient piecewise training of deep structured models for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3194–3203, 2016. 2

[26] M. Lin, Q. Chen, and S. Yan. Network in network. *arXiv preprint arXiv:1312.4400*, 2013. 3

[27] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. 2

[28] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 2

[29] D. Mishkin, N. Sergievskiy, and J. Matas. Systematic evaluation of cnn advances on the imagenet. *arXiv preprint arXiv:1606.02228*, 2016. 3

[30] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814, 2010. 3

[31] S. A. Nene, S. K. Nayar, H. Murase, et al. Columbia object image library (coil-20). 1996. 2

[32] H. Noh, S. Hong, and B. Han. Learning deconvolution network for semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1520–1528, 2015. 2

[33] A. Odena, C. Olah, and J. Shlens. Conditional image synthesis with auxiliary classifier gans. *arXiv preprint arXiv:1610.09585*, 2016. 2

[34] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee. Generative adversarial text to image synthesis. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 3, 2016. 2

[35] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 2, 3

[36] D. Scherer, A. Müller, and S. Behnke. Evaluation of pooling operations in convolutional architectures for object recogni-
tion. In *International Conference on Artificial Neural Networks*, pages 92–101. Springer, 2010. 2

[37] M. Seuret, R. Ingold, and M. Liwicki. N-light-n: A highly-adaptable java library for document analysis with convolutional auto-encoders and related architectures. In *Frontiers in Handwriting Recognition (ICFHR), 2016 15th International Conference on*, pages 459–464. IEEE, 2016. 3

[38] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 1, 3

[39] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015. 1, 2, 3, 4, 5, 6, 7

[40] A. Teichman and S. Thrun. Practical object recognition in autonomous driving and beyond. In *Advanced Robotics and its Social Impacts (ARSO), 2011 IEEE Workshop on*, pages 35–38. IEEE, 2011. 1

[41] R. Vedantam, S. Bengio, K. Murphy, D. Parikh, and G. Chechik. Context-aware captions from context-agnostic supervision. *arXiv preprint arXiv:1701.02870*, 2017. 2

[42] F. Visin, K. Kastner, K. Cho, M. Matteucci, A. Courville, and Y. Bengio. Renet: A recurrent neural network based alternative to convolutional networks. *arXiv preprint arXiv:1505.00393*, 2015. 2

[43] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, volume 14, pages 77–81, 2015. 2

[44] M. D. Zeiler. *Hierarchical convolutional deep learning in computer vision*. PhD thesis, NEW YORK UNIVERSITY, 2013. 1

[45] M. D. Zeiler and R. Fergus. *Visualizing and Understanding Convolutional Networks*, pages 818–833. Springer International Publishing, Cham, 2014. 2, 3

[46] H. Zhang, T. Xu, H. Li, S. Zhang, X. Huang, X. Wang, and D. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. *arXiv preprint arXiv:1612.03242*, 2016. 2

[47] X. Zhang, Z. Li, C. C. Loy, and D. Lin. Polynet: A pursuit of structural diversity in very deep networks. *arXiv preprint arXiv:1611.05725*, 2016. 1, 2

[48] Z.-H. Zhou, J. Wu, and W. Tang. Ensembling neural networks: many could be better than all. *Artificial intelligence*, 137(1-2):239–263, 2002. 4