SnailTyan

────

ResNet　　　　——

||

ResNet　　　——

Tyan
noahsnail.com | CSDN |



https://github.com/SnailTyan/deep-learning-papers-translation

# Deep Residual Learning for Image Recognition

## Abstract

Deeper neural networks are more difficult to train. We present a residual learning framework to ease the training of networks that are substantially deeper than those used previously. We explicitly reformulate the layers as learning residual functions with reference to the layer inputs, instead of learning unreferenced functions. We provide comprehensive empirical evidence showing that these residual networks are easier to optimize, and can gain accuracy from considerably increased depth. On the ImageNet dataset we evaluate residual nets with a depth of up to 152 layers——8× deeper than VGG nets [40] but still having lower complexity. An ensemble of these residual nets achieves 3.57% error on the ImageNet test set. This result won the 1st place on the ILSVRC 2015 classification task. We also present analysis on CIFAR-10 with 100 and 1000 layers.

ImageNet　　　　　　　152　　　　—— VGG[40]　8

ImageNet　　　　3. 57%　　　　　　ILSVRC 2015　　　　　　　　CIFAR-10　　100　1000

The depth of representations is of central importance for many visual recognition tasks. Solely due to our extremely deep representations, we obtain a 28% relative improvement on the COCO object detection dataset. Deep residual nets are foundations of our submissions to ILSVRC & COCO 2015 competitions, where we also won the 1st places on the tasks of ImageNet detection, ImageNet localization, COCO detection, and COCO segmentation.

ImageNet　　　　ImageNet　　　COCO　　COCO　　　　COCO　　　　　　28%　　　　　　ILSVRC　COCO 2015

# 1. Introduction

Deep convolutional neural networks [22, 21] have led to a series of breakthroughs for image classification [21, 49, 39]. Deep networks naturally integrate low/mid/high-level features [49] and classifiers in an end-to-end multi-layer fashion, and the " levels" of features can be enriched by the number of stacked layers (depth). Recent evidence [40, 43] reveals that network depth is of crucial importance, and the leading results [40, 43, 12, 16] on the challenging ImageNet dataset [35] all exploit " very deep" [40] models, with a depth of sixteen [40] to thirty [16]. Many other non-trivial visual recognition tasks [7, 11, 6, 32, 27] have also greatly benefited from very deep models.

## 1.

[22, 21]　　　　　[21, 49, 39]　　　　　　　　　　/ /　　　[49]　　　　　　　　　　　　　　" "
[40, 43]　　　　　　　　　ImageNet　　　　　　" " [40]　　　　　16 [40]　30 [16]　　　　　　　　　　　　[7, 11, 6, 32, 27]

Driven by the significance of depth, a question arises: Is learning better networks as easy as stacking more layers? An obstacle to answering this question was the notorious problem of vanishing/exploding gradients [14, 1, 8], which hamper convergence from the beginning. This problem, however, has been largely addressed by normalized initialization [23, 8, 36, 12] and intermediate normalization layers [16], which enable networks with tens of layers to start converging for stochastic gradient descent (SGD) with backpropagation [22].

[23, 8, 36, 12]　　　　　　　[16]　　　　　　　　　　　　　　　　　　　　　　　　　/　　[14, 1, 8]
　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　SGD

When deeper networks are able to start converging, a degradation problem has been exposed: with the network depth increasing, accuracy gets saturated (which might be unsurprising) and then degrades rapidly. Unexpectedly, such degradation is not caused by overfitting, and adding more layers to a suitably deep model leads to higher training error, as reported in [10, 41] and thoroughly verified by our experiments. Fig. 1 shows a typical example.
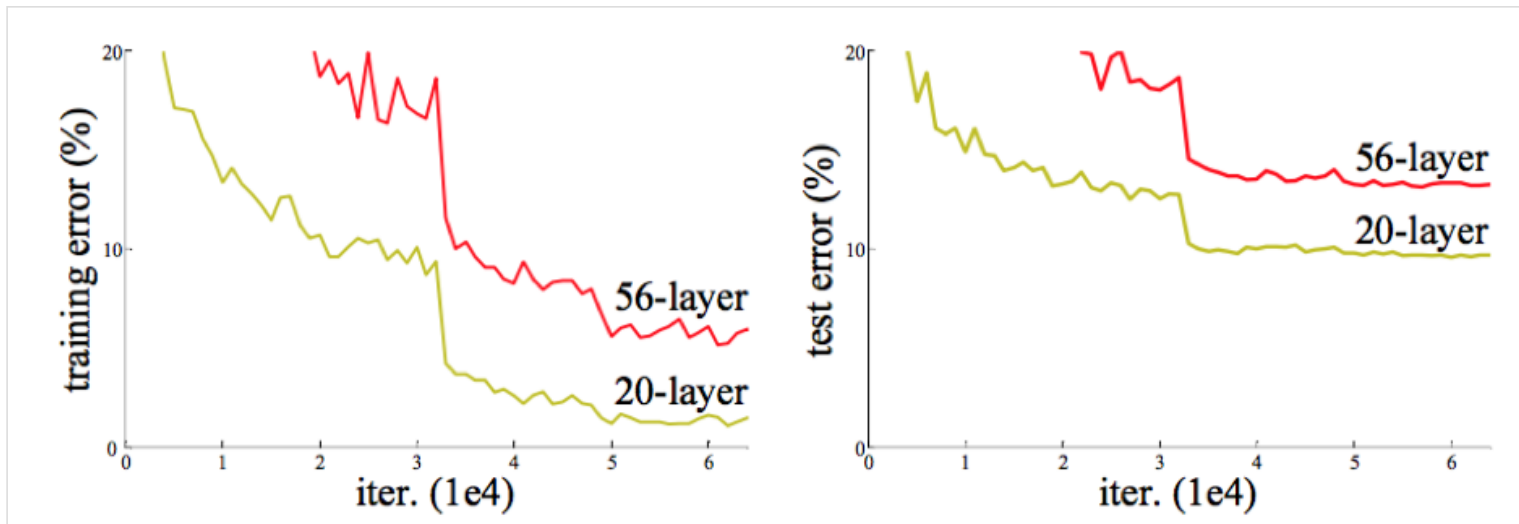
Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer "plain" networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

[10, 41]                                                                                                    1



1 20    56    "    "    CIFAR-10                                                                    ImageNet            4

The degradation (of training accuracy) indicates that not all systems are similarly easy to optimize. Let us consider a shallower architecture and its deeper counterpart that adds more layers onto it. There exists a solution by construction to the deeper model: the added layers are identity mapping, and the other layers are copied from the learned shallower model. The existence of this constructed solution indicates that a deeper model should produce no higher training error than its shallower counterpart. But experiments show that our current solvers on hand are unable to find solutions that are comparably good or better than the constructed solution (or unable to do so in feasible time).

In this paper, we address the degradation problem by introducing a *deep residual learning* framework. Instead of hoping each few stacked layers directly fit a desired underlying mapping, we explicitly let these layers fit a residual mapping. Formally, denoting the desired underlying mapping as $H(x)$, we let the stacked nonlinear layers fit another mapping of $F(x) := H(x) - x$. The original mapping is recast into $F(x) + x$. We hypothesize that it is easier to optimize the residual mapping than to optimize the original, unreferenced mapping. To the extreme, if an identity mapping were optimal, it would be easier to push the residual to zero than to fit an identity mapping by a stack of nonlinear layers.

$$H(x)$$

$$F(x) := H(x) - x \qquad\qquad F(x) + x$$

The formulation of $F(x) + x$ can be realized by feedforward neural networks with " shortcut connections"  (Fig. 2). Shortcut connections [2, 33, 48] are those skipping one or more layers. In our case, the shortcut connections simply perform identity mapping, and their outputs are added to the outputs of the stacked layers (Fig. 2). Identity shortcut connections add neither extra parameter nor computational complexity. The entire network can still be trained end-to-end by SGD with backpropagation, and can be easily implemented using common libraries (e.g., Caffe [19]) without modifying the solvers.



Figure 2. Residual learning: a building block.

$F(x) + x$                    "            "                         2                    [2, 33, 48]
        2
                                                                                           SGD                                        Caffe [19]
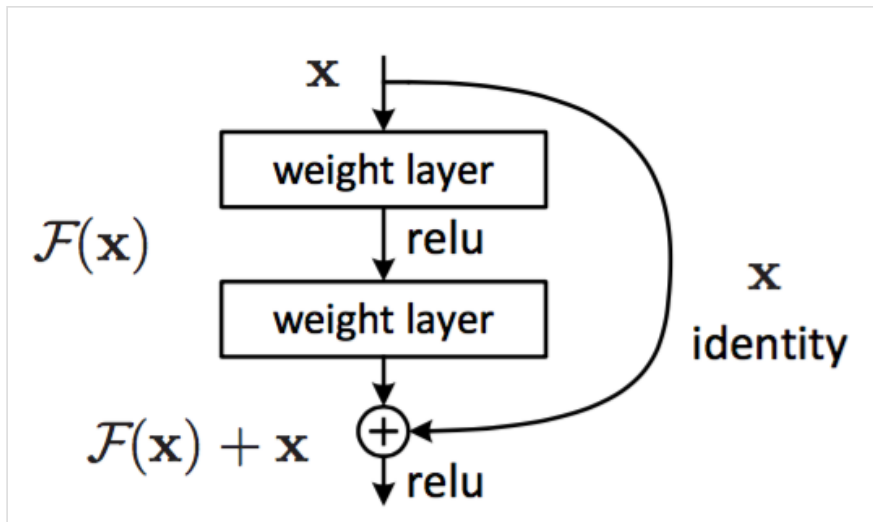
2.

We present comprehensive experiments on ImageNet [35] to show the degradation problem and evaluate our method. We show that: 1) Our extremely deep residual nets are easy to optimize, but the counterpart " plain" nets (that simply stack layers) exhibit higher training error when the depth increases; 2) Our deep residual nets can easily enjoy accuracy gains from greatly increased depth, producing results substantially better than previous networks.

ImageNet[35]                                                                    1
2

Similar phenomena are also shown on the CIFAR-10 set [20], suggesting that the optimization difficulties and the effects of our method are not just akin to a particular dataset. We present successfully trained models on this dataset with over 100 layers, and explore models with over 1000 layers.

CIFAR-10          [20]                                                                  100               1000

On the ImageNet classification dataset [35], we obtain excellent results by extremely deep residual nets. Our 152-layer residual net is the deepest network ever presented on ImageNet, while still having lower complexity than VGG nets [40]. Our ensemble has 3. 57% top-5 error on the ImageNet test set, and won the 1st place in the ILSVRC 2015 classification competition. The extremely deep representations also have excellent generalization performance on other recognition tasks, and lead us to further win the 1st places on: ImageNet detection, ImageNet localization, COCO detection, and COCO segmentation in ILSVRC & COCO 2015 competitions. This strong evidence shows that the residual learning principle is generic, and we expect that it is applicable in other vision and non-vision problems.

ImageNet          [35]                          152          ImageNet                    VGG     [40]
ImageNet          3. 57% top-5          ILSVRC 2015                                                      ILSVRC &
COCO 2015         ImageNet     ImageNet     COCO     COCO

## 2. Related Work

**Residual Representations**. In image recognition, VLAD [18] is a representation that encodes by the residual vectors with respect to a dictionary, and Fisher Vector [30] can be formulated as a probabilistic version [18] of VLAD. Both of them are powerful shallow representations for image retrieval and classification [4, 47]. For vector quantization, encoding residual vectors [17] is shown to be more effective than encoding original vectors.

## 2.

VLAD[18]                                       Fisher     [30]     VLAD     [18]                    [4,47]

[17]

In low-level vision and computer graphics, for solving Partial Differential Equations (PDEs), the widely used Multigrid method [3] reformulates the system as subproblems at multiple scales, where each subproblem is responsible for the residual solution between a coarser and a finer scale. An alternative to Multigrid is hierarchical basis preconditioning [44, 45], which relies on variables that represent residual vectors between two scales. It has been shown [3, 44, 45] that these solvers converge much faster than standard solvers that are unaware of the residual nature of the solutions. These methods suggest that a good reformulation or preconditioning can simplify the optimization.

PDE          Multigrid     [3]                                           Multigrid
[44,45]                                         [3,44,45]

**Shortcut Connections**. Practices and theories that lead to shortcut connections [2, 33, 48] have been studied for a long time. An early practice of training multi-layer perceptrons (MLPs) is to add a linear layer connected from the network input to the output [33, 48]. In [43, 24], a few intermediate layers are directly connected to auxiliary classifiers for addressing vanishing/exploding gradients. The papers of [38, 37, 31, 46] propose methods for centering layer responses, gradients, and propagated errors, implemented by shortcut connections. In [43], an " inception" layer is composed of a shortcut branch and a few deeper branches.

[2,33,48]                                    MLP                                    [33,48]    [43,24]
/          [38,37,31,46]                                     [43]          " inception"

Concurrent with our work, " highway networks" [41, 42] present shortcut connections with gating functions [15]. These gates are data-dependent and have parameters, in contrast to our identity shortcuts that are parameter-free. When a gated shortcut is " closed" (approaching zero), the layers in highway networks represent non-residual functions. On the contrary, our formulation always learns residual functions; our identity shortcuts are never closed, and all information is always passed through, with additional residual functions to be learned. In addition, highway networks have not demonstrated accuracy gains with extremely increased depth (e.g., over 100 layers).

" highway networks" [41, 42]              [15]                                                      "   "

100

# 3. Deep Residual Learning

## 3.1. Residual Learning

Let us consider $H(x)$ as an underlying mapping to be fit by a few stacked layers (not necessarily the entire net), with $x$ denoting the inputs to the first of these layers. If one hypothesizes that multiple nonlinear layers can asymptotically approximate complicated functions, then it is equivalent to hypothesize that they can asymptotically approximate the residual functions, i.e., $H(x) - x$ (assuming that the input and output are of the same dimensions). So rather than expect stacked layers to approximate $H(x)$, we explicitly let these layers approximate a residual function $F(x) := H(x) - x$. The original function thus becomes $F(x) + x$. Although both forms should be able to asymptotically approximate the desired functions (as hypothesized), the ease of learning might be different.

## 3.

## 3.1.

$H(x)$                                                $x$
$H(x) - x$(                              )                            $F(x) := H(x) - x$                     $H(x)$                     $F(x) + x$

This reformulation is motivated by the counterintuitive phenomena about the degradation problem (Fig. 1, left). As we discussed in the introduction, if the added layers can be constructed as identity mappings, a deeper model should have training error no greater than its shallower counterpart. The degradation problem suggests that the solvers might have difficulties in approximating identity mappings by multiple nonlinear layers. With the residual learning reformulation, if identity mappings are optimal, the solvers may simply drive the weights of the multiple nonlinear layers toward zero to approach identity mappings.

1

In real cases, it is unlikely that identity mappings are optimal, but our reformulation may help to precondition the problem. If the optimal function is closer to an identity mapping than to a zero mapping, it should be easier for the solver to find the perturbations with reference to an identity mapping, than to learn the function as a new one. We show by experiments (Fig. 7) that the learned residual functions in general have small responses, suggesting that identity mappings provide reasonable preconditioning.
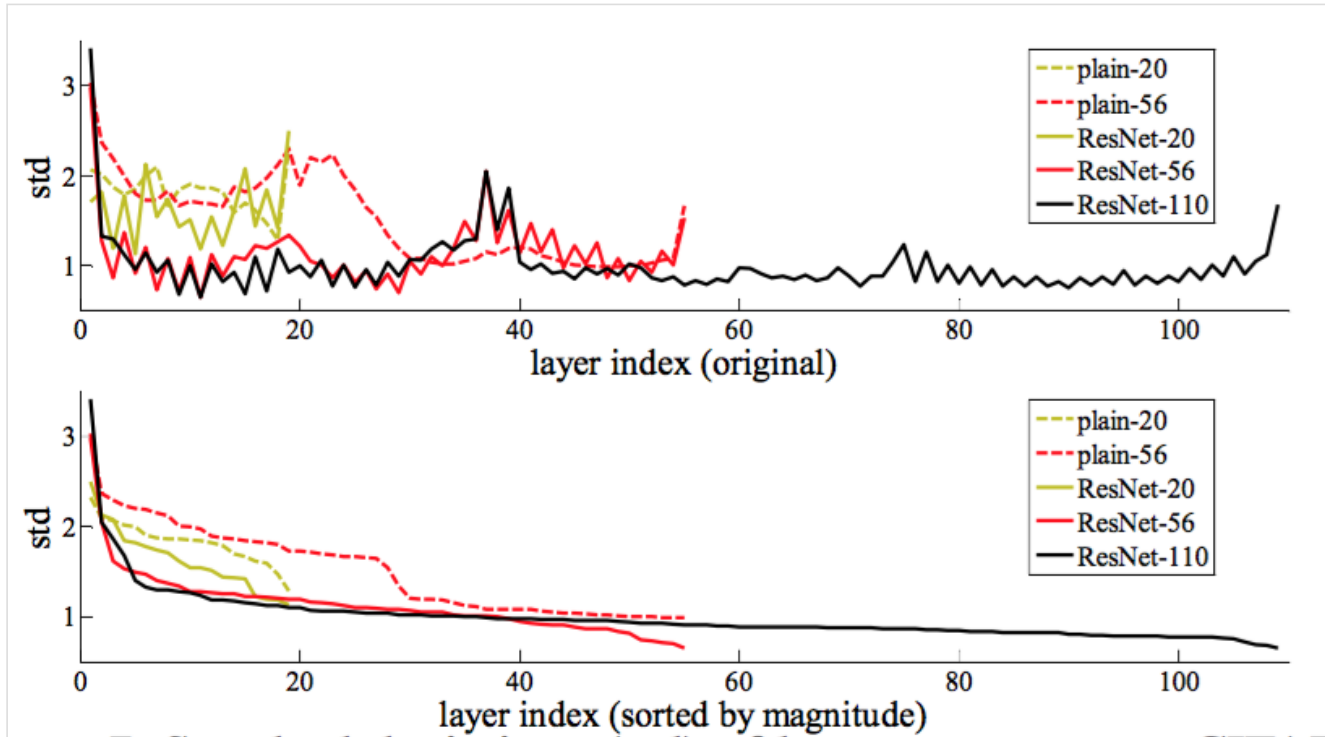


Figure 7. Standard deviations (std) of layer responses on CIFAR-10. The responses are the outputs of each 3×3 layer, after BN and before nonlinearity. Top: the layers are shown in their original order. Bottom: the responses are ranked in descending order.

7

7        CIFAR-10        std        3×3        BN

## 3.2. Identity Mapping by Shortcuts

We adopt residual learning to every few stacked layers. A building block is shown in Fig. 2. Formally, in this paper we consider a building block defined as:

$$y = F(x, \{W_i\}) + x \tag{1}$$

Here $x$ and $y$ are the input and output vectors of the layers considered. The function $F(x, \{W_i\})$ represents the residual mapping to be learned. For the example in Fig. 2 that has two layers, $F = W_2 \sigma(W_1x)$ in which $\sigma$ denotes ReLU [29] and the biases are omitted for simplifying notations. The operation $F + x$ is performed by a shortcut connection and element-wise addition. We adopt the second nonlinearity after the addition (i.e., $\sigma(y)$, see Fig. 2).

## 3.2.

2

$$y = F(x, \{W_i\}) + x \tag{1}$$

$x$   $y$                          $F(x, \{W_i\})$                        2                $F = W_2 \sigma(W_1x)$   $\sigma$      ReLU[29]                        $F + x$
                                                                    $\sigma(y)$        2

The shortcut connections in Eqn.(1) introduce neither extra parameter nor computation complexity. This is not only attractive in practice but also important in our comparisons between plain and residual networks. We can fairly compare plain/residual networks that simultaneously have the same number of parameters, depth, width, and computational cost (except for the negligible element-wise addition).

(1)

/

The dimensions of $x$ and $F$ must be equal in Eqn.(1). If this is not the case (e.g., when changing the input/output channels), we can perform a linear projection $W_s$ by the shortcut connections to match the dimensions:

$$y = F(x, \{W_i \}) + W_sx.$$ (2)

We can also use a square matrix $W_s$ in Eqn.(1). But we will show by experiments that the identity mapping is sufficient for addressing the degradation problem and is economical, and thus $W_s$ is only used when matching dimensions.

(1)   $x$   $F$                                   /                               $W_s$

$$y = F(x, \{W_i \}) + W_sx.$$ (2)

           (1)      $W_s$                                        $W_s$

The form of the residual function $F$ is flexible. Experiments in this paper involve a function $F$ that has two or three layers (Fig. 5), while more layers are possible. But if $F$ has only a single layer, Eqn.(1) is similar to a linear layer: $y = W_1x + x$, for which we have not observed advantages.
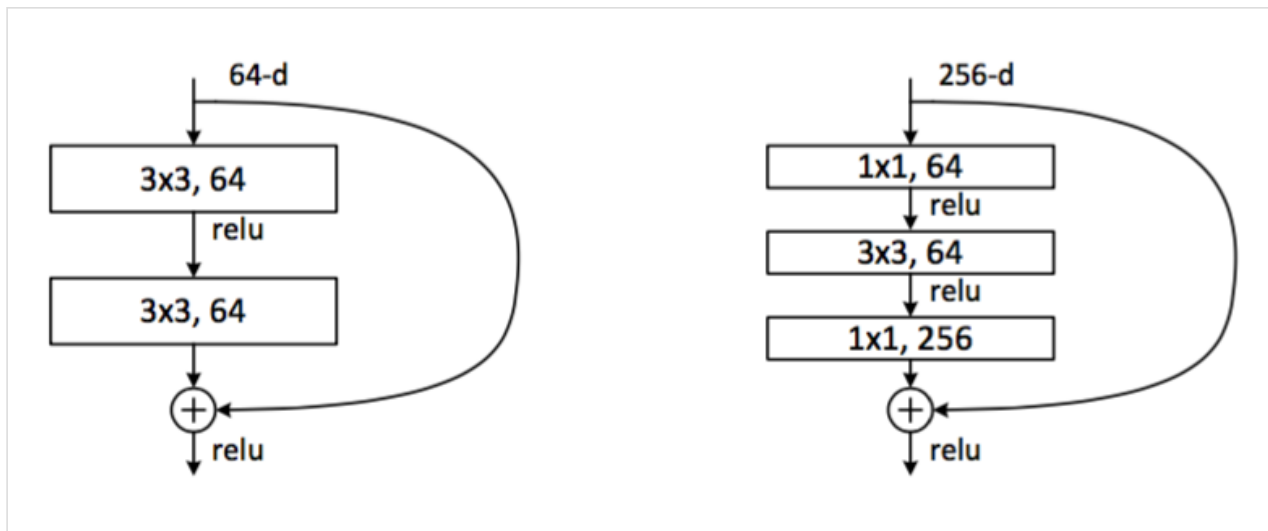


Figure 5. A deeper residual function $F$ for ImageNet. Left: a building block (on 56× 56 feature maps) as in Fig. 3 for ResNet-34. Right: a " bottleneck" building block for ResNet-50/101/152.

      $F$                                     5     $F$                          $F$         (1)           $y = W_1x + x$

5　ImageNet　　　　　　$F$　　　ResNet-34　　　　　56× 56　　　　　3　　　ResNet-50/101/152　"  bottleneck"

We also note that although the above notations are about fully-connected layers for simplicity, they are applicable to convolutional layers. The function $F(x,\{W_i\})$ can represent multiple convolutional layers. The element-wise addition is performed on two feature maps, channel by channel.

$$F(x  \{W_i\})$$

## 3.3. Network Architectures

We have tested various plain/residual nets, and have observed consistent phenomena. To provide instances for discussion, we describe two models for ImageNet as follows.

## 3.3.

　　　　　　　／　　　　　　　　　　　　　　　　　　　　　　　　ImageNet

**Plain Network**. Our plain baselines (Fig. 3, middle) are mainly inspired by the philosophy of VGG nets 40. The convolutional layers mostly have 3× 3 filters and follow two simple design rules: (i) for the same output feature map size, the layers have the same number of filters; and (ii) if the feature map size is halved, the number of filters is doubled so as to preserve the time complexity per layer. We perform downsampling directly by convolutional layers that have a stride of 2. The network ends with a global average pooling layer and a 1000-way fully-connected layer with softmax. The total number of weighted layers is 34 in Fig. 3 (middle).

**VGG-19**

**34-layer plain**

**34-layer residual**

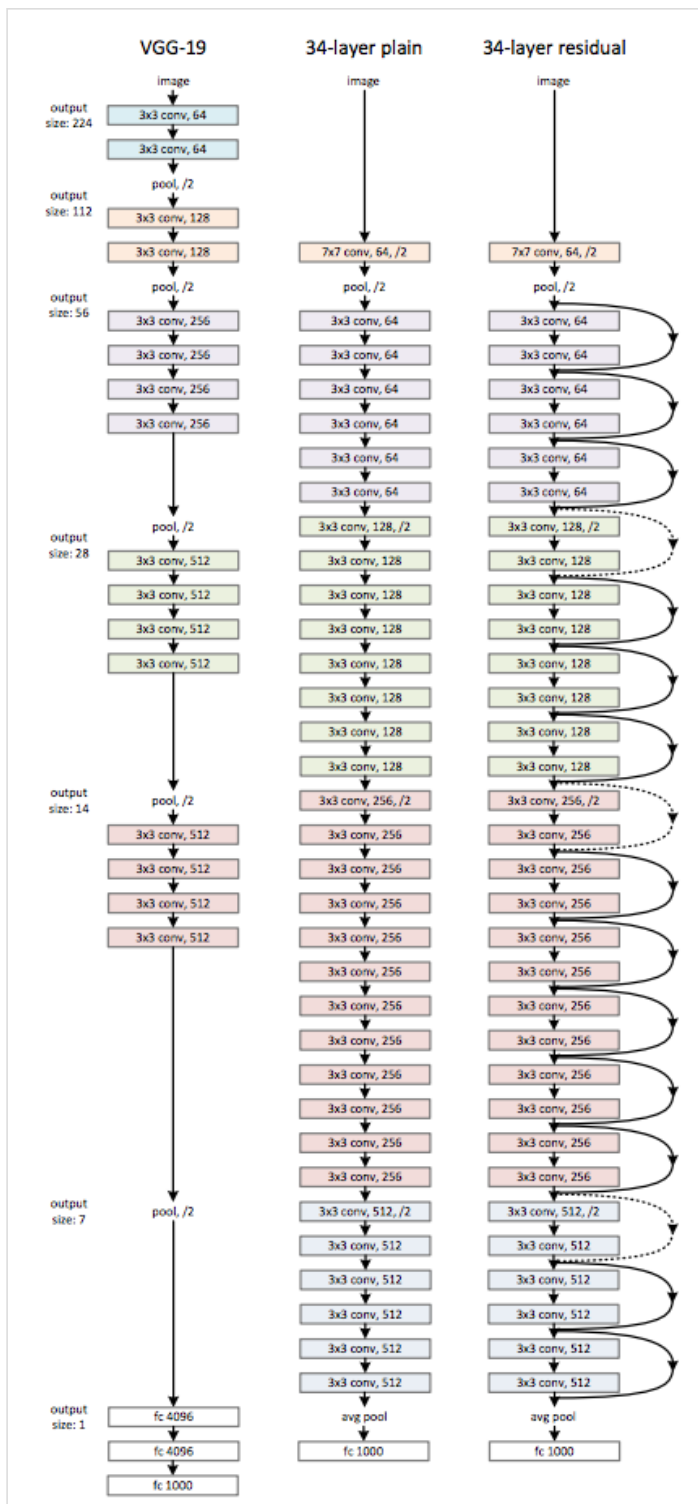| output size: 224 | image | image | image |
| output size: 112 | 3x3 conv, 64 | | |
| | 3x3 conv, 64 | | |
| | pool, /2 | | |
| | 3x3 conv, 128 | 7x7 conv, 64, /2 | 7x7 conv, 64, /2 |
| output size: 56 | 3x3 conv, 128 | | |
| | pool, /2 | pool, /2 | pool, /2 |
| | 3x3 conv, 256 | 3x3 conv, 64 | 3x3 conv, 64 |
| | 3x3 conv, 256 | 3x3 conv, 64 | 3x3 conv, 64 |
| | 3x3 conv, 256 | 3x3 conv, 64 | 3x3 conv, 64 |
| | 3x3 conv, 256 | 3x3 conv, 64 | 3x3 conv, 64 |
| output size: 28 | | 3x3 conv, 64 | 3x3 conv, 64 |
| | | 3x3 conv, 64 | 3x3 conv, 64 |
| | pool, /2 | 3x3 conv, 128, /2 | 3x3 conv, 128, /2 |
| | 3x3 conv, 512 | 3x3 conv, 128 | 3x3 conv, 128 |
| | 3x3 conv, 512 | 3x3 conv, 128 | 3x3 conv, 128 |
| | 3x3 conv, 512 | 3x3 conv, 128 | 3x3 conv, 128 |
| | 3x3 conv, 512 | 3x3 conv, 128 | 3x3 conv, 128 |
| | | 3x3 conv, 128 | 3x3 conv, 128 |
| | | 3x3 conv, 128 | 3x3 conv, 128 |
| | | 3x3 conv, 128 | 3x3 conv, 128 |
| output size: 14 | pool, /2 | 3x3 conv, 256, /2 | 3x3 conv, 256, /2 |
| | 3x3 conv, 512 | 3x3 conv, 256 | 3x3 conv, 256 |
| | 3x3 conv, 512 | 3x3 conv, 256 | 3x3 conv, 256 |
| | 3x3 conv, 512 | 3x3 conv, 256 | 3x3 conv, 256 |
| | 3x3 conv, 512 | 3x3 conv, 256 | 3x3 conv, 256 |
| | | 3x3 conv, 256 | 3x3 conv, 256 |
| | | 3x3 conv, 256 | 3x3 conv, 256 |
| | | 3x3 conv, 256 | 3x3 conv, 256 |
| | | 3x3 conv, 256 | 3x3 conv, 256 |
| | | 3x3 conv, 256 | 3x3 conv, 256 |
| | | 3x3 conv, 256 | 3x3 conv, 256 |
| | | 3x3 conv, 256 | 3x3 conv, 256 |
| | | 3x3 conv, 256 | 3x3 conv, 256 |
| output size: 7 | pool, /2 | 3x3 conv, 512, /2 | 3x3 conv, 512, /2 |
| | | 3x3 conv, 512 | 3x3 conv, 512 |
| | | 3x3 conv, 512 | 3x3 conv, 512 |
| | | 3x3 conv, 512 | 3x3 conv, 512 |
| | | 3x3 conv, 512 | 3x3 conv, 512 |
| | | 3x3 conv, 512 | 3x3 conv, 512 |
| output size: 1 | fc 4096 | avg pool | avg pool |
| | fc 4096 | fc 1000 | fc 1000 |
| | fc 1000 | | |

Figure 3. Example network architectures for ImageNet. Left: the VGG-19 model [40] (19.6 billion FLOPs) as a reference. Middle: a plain network with 34 parameter layers (3.6 billion FLOPs). Right: a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. Table 1 shows more details and other variants.

| layer name | output size | 18-layer | 34-layer | 50-layer | 101-layer | 152-layer |
|---|---|---|---|---|---|---|
| conv1 | 112×112 | 7×7, 64, stride 2 | | | | |
| | | 3×3 max pool, stride 2 | | | | |
| conv2_x | 56×56 | $\begin{bmatrix} 3\times3,\ 64 \\ 3\times3,\ 64 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 64 \\ 3\times3,\ 64 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 64 \\ 3\times3,\ 64 \\ 1\times1,\ 256 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 64 \\ 3\times3,\ 64 \\ 1\times1,\ 256 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 64 \\ 3\times3,\ 64 \\ 1\times1,\ 256 \end{bmatrix}\times3$ |
| conv3_x | 28×28 | $\begin{bmatrix} 3\times3,\ 128 \\ 3\times3,\ 128 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 128 \\ 3\times3,\ 128 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1,\ 128 \\ 3\times3,\ 128 \\ 1\times1,\ 512 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1,\ 128 \\ 3\times3,\ 128 \\ 1\times1,\ 512 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1,\ 128 \\ 3\times3,\ 128 \\ 1\times1,\ 512 \end{bmatrix}\times8$ |
| conv4_x | 14×14 | $\begin{bmatrix} 3\times3,\ 256 \\ 3\times3,\ 256 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 256 \\ 3\times3,\ 256 \end{bmatrix}\times6$ | $\begin{bmatrix} 1\times1,\ 256 \\ 3\times3,\ 256 \\ 1\times1,\ 1024 \end{bmatrix}\times6$ | $\begin{bmatrix} 1\times1,\ 256 \\ 3\times3,\ 256 \\ 1\times1,\ 1024 \end{bmatrix}\times23$ | $\begin{bmatrix} 1\times1,\ 256 \\ 3\times3,\ 256 \\ 1\times1,\ 1024 \end{bmatrix}\times36$ |
| conv5_x | 7×7 | $\begin{bmatrix} 3\times3,\ 512 \\ 3\times3,\ 512 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 512 \\ 3\times3,\ 512 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 512 \\ 3\times3,\ 512 \\ 1\times1,\ 2048 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 512 \\ 3\times3,\ 512 \\ 1\times1,\ 2048 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 512 \\ 3\times3,\ 512 \\ 1\times1,\ 2048 \end{bmatrix}\times3$ |
| | 1×1 | average pool, 1000-d fc, softmax | | | | |
| FLOPs | | $1.8\times10^9$ | $3.6\times10^9$ | $3.8\times10^9$ | $7.6\times10^9$ | $11.3\times10^9$ |

Table 1. Architectures for ImageNet. Building blocks are shown in brackets (see also Fig. 5), with the numbers of blocks stacked. Down-sampling is performed by conv3_1, conv4_1, and conv5_1 with a stride of 2.

| layer name | output size | 18-layer | 34-layer | 50-layer | 101-layer | 152-layer |
|---|---|---|---|---|---|---|
| conv1 | 112×112 | 7×7, 64, stride 2 | | | | |
| conv2_x | 56×56 | 3×3 max pool, stride 2 | | | | |
| | | $\begin{bmatrix}3\times3,64\\3\times3,64\end{bmatrix}\times2$ | $\begin{bmatrix}3\times3,64\\3\times3,64\end{bmatrix}\times3$ | $\begin{bmatrix}1\times1,64\\3\times3,64\\1\times1,256\end{bmatrix}\times3$ | $\begin{bmatrix}1\times1,64\\3\times3,64\\1\times1,256\end{bmatrix}\times3$ | $\begin{bmatrix}1\times1,64\\3\times3,64\\1\times1,256\end{bmatrix}\times3$ |
| conv3_x | 28×28 | $\begin{bmatrix}3\times3,128\\3\times3,128\end{bmatrix}\times2$ | $\begin{bmatrix}3\times3,128\\3\times3,128\end{bmatrix}\times4$ | $\begin{bmatrix}1\times1,128\\3\times3,128\\1\times1,512\end{bmatrix}\times4$ | $\begin{bmatrix}1\times1,128\\3\times3,128\\1\times1,512\end{bmatrix}\times4$ | $\begin{bmatrix}1\times1,128\\3\times3,128\\1\times1,512\end{bmatrix}\times8$ |
| conv4_x | 14×14 | $\begin{bmatrix}3\times3,256\\3\times3,256\end{bmatrix}\times2$ | $\begin{bmatrix}3\times3,256\\3\times3,256\end{bmatrix}\times6$ | $\begin{bmatrix}1\times1,256\\3\times3,256\\1\times1,1024\end{bmatrix}\times6$ | $\begin{bmatrix}1\times1,256\\3\times3,256\\1\times1,1024\end{bmatrix}\times23$ | $\begin{bmatrix}1\times1,256\\3\times3,256\\1\times1,1024\end{bmatrix}\times36$ |
| conv5_x | 7×7 | $\begin{bmatrix}3\times3,512\\3\times3,512\end{bmatrix}\times2$ | $\begin{bmatrix}3\times3,512\\3\times3,512\end{bmatrix}\times3$ | $\begin{bmatrix}1\times1,512\\3\times3,512\\1\times1,2048\end{bmatrix}\times3$ | $\begin{bmatrix}1\times1,512\\3\times3,512\\1\times1,2048\end{bmatrix}\times3$ | $\begin{bmatrix}1\times1,512\\3\times3,512\\1\times1,2048\end{bmatrix}\times3$ |
| | 1×1 | average pool, 1000-d fc, softmax | | | | |
| FLOPs | | $1.8\times10^9$ | $3.6\times10^9$ | $3.8\times10^9$ | $7.6\times10^9$ | $11.3\times10^9$ |

1  ImageNet                          5                          2  conv3_1, conv4_1  conv5_1

It is worth noticing that our model has *fewer* filters and *lower* complexity than VGG nets 40. Our 34-layer baseline has 3.6 billion FLOPs (multiply-adds), which is only 18% of VGG-19 (19.6 billion FLOPs).

                VGG    3                          34      36  FLOP( )    VGG-19  196  FLOP    18%

**Residual Network**. Based on the above plain network, we insert shortcut connections (Fig. 3, right) which turn the network into its counterpart residual version. The identity shortcuts (Eqn.(1)) can be directly used when the input and output are of the same dimensions (solid line shortcuts in Fig. 3). When the dimensions increase (dotted line shortcuts in Fig. 3), we consider two options: (A) The shortcut still performs identity mapping, with extra zero entries padded for increasing dimensions. This option introduces no extra parameter; (B) The projection shortcut in Eqn.(2) is used to match dimensions (done by 1×1 convolutions). For both options, when the shortcuts go across feature maps of two sizes, they are performed with a stride of 2.

1              3          3                    A                            3          B      2
        1×1                                          2

## 3.4. Implementation

Our implementation for ImageNet follows the practice in [21, 40]. The image is resized with its shorter side randomly sampled in [256, 480] for scale augmentation [40]. A 224×224 crop is randomly sampled from an image or its horizontal flip, with the per-pixel mean subtracted [21]. The standard color augmentation in [21] is used. We adopt batch normalization (BN) [16] right after each convolution and before activation, following [16]. We initialize the weights as in [12] and train all plain/residual nets from scratch. We use SGD with a mini-batch size of 256. The learning rate starts from 0.1 and is divided by 10 when the error plateaus, and the models are trained for up to $60 × 10^ 4$ iterations. We use a weight decay of 0.0001 and a momentum of 0.9. We do not use dropout [13], following the practice in [16].

## 3.4.

ImageNet                    [21  40]                          [256,480]                    [40]  224×224
[21]        [21]                                          BN  [16]            [12]                                    /                        256
SGD              0.1                        10          $60 × 10^ 4$                          0.0001        0.9        [16]                  [13]

In testing, for comparison studies we adopt the standard 10-crop testing [21]. For best results, we adopt the fully-convolutional form as in [40, 12], and average the scores at multiple scales (images are resized such that the shorter side is in {224, 256, 384, 480, 640}).

10-crop      [21]                          [40, 12]                                                  {224, 256, 384, 480, 640}

# 4. Experiments

## 4.1. ImageNet Classification

We evaluate our method on the ImageNet 2012 classification dataset [35] that consists of 1000 classes. The models are trained on the 1.28 million training images, and evaluated on the 50k validation images. We also obtain a final result on the 100k test images, reported by the test server. We evaluate both top-1 and top-5 error rates.

## 4.

## 4.1. ImageNet

ImageNet 2012            [35]                                      1000                128                        5
10                                  top-1  top-5

**Plain Networks**. We first evaluate 18-layer and 34-layer plain nets. The 34-layer plain net is in Fig. 3 (middle). The 18-layer plain net is of a similar form. See Table 1 for detailed architectures.

18    34                34            3          18                                                  1

The results in Table 2 show that the deeper 34-layer plain net has higher validation error than the shallower 18-layer plain net. To reveal the reasons, in Fig. 4 (left) we compare their training/validation errors during the training procedure. We have observed the degradation problem —— the 34-layer plain net has higher training error throughout the whole training procedure, even though the solution space of the 18-layer plain network is a subspace of that of the 34-layer one.

|           | plain | ResNet |
|-----------|-------|--------|
| 18 layers | 27.94 | 27.88  |
| 34 layers | 28.54 | **25.03** |

Table 2. Top-1 error (%, 10-crop testing) on ImageNet validation. Here the ResNets have no extra parameter compared to their plain counterparts. Fig. 4 shows the training procedures.
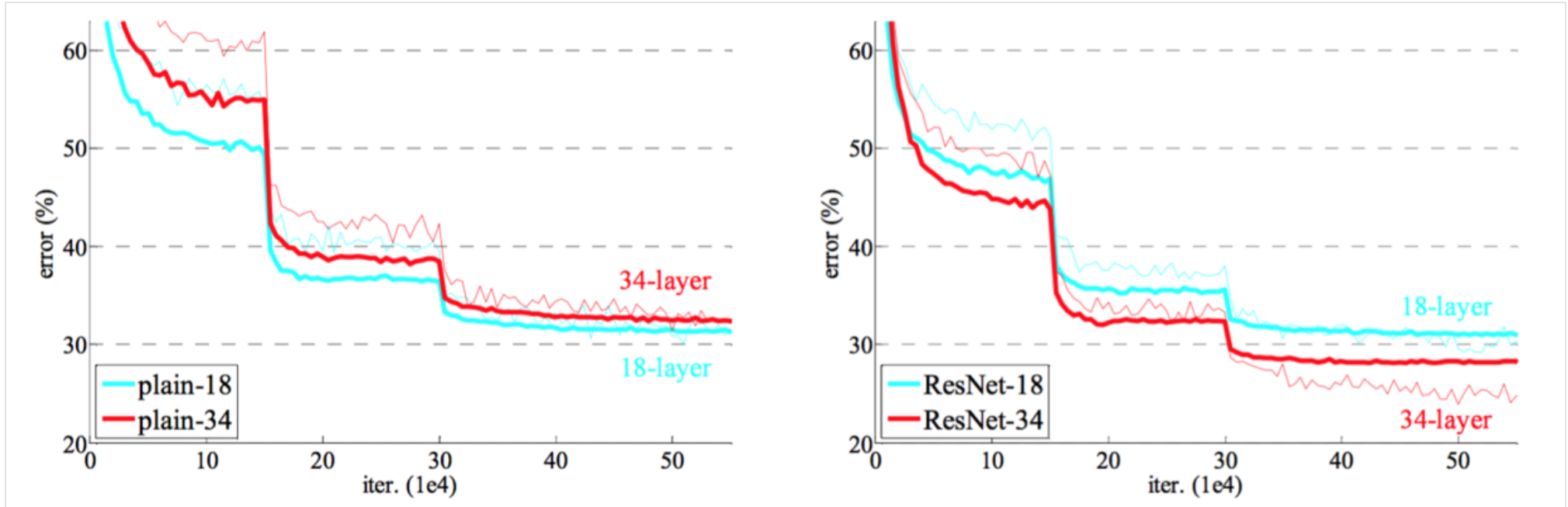


Figure 4. Training on ImageNet. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

2         34          18        4                    /              —— 18
      34                  34

|  | plain | ResNet |
|---|---|---|
| 18 layers | 27.94 | 27.88 |
| 34 layers | 28.54 | **25.03** |

2 ImageNet      Top-1      (%  10          )                    ResNet          4

4　　ImageNet　　　　　　　　　　　　　　　18　34　　　　18　34　　ResNet

We argue that this optimization difficulty is unlikely to be caused by vanishing gradients. These plain networks are trained with BN [16], which ensures forward propagated signals to have non-zero variances. We also verify that the backward propagated gradients exhibit healthy norms with BN. So neither forward nor backward signals vanish. In fact, the 34-layer plain net is still able to achieve competitive accuracy (Table 3), suggesting that the solver works to some extent. We conjecture that the deep plain nets may have exponentially low convergence rates, which impact the reducing of the training error. The reason for such optimization difficulties will be studied in the future.

| model | top-1 err. | top-5 err. |
|---|---|---|
| VGG-16 [40] | 28.07 | 9.33 |
| GoogLeNet [43] | - | 9.15 |
| PReLU-net [12] | 24.27 | 7.38 |
| plain-34 | 28.54 | 10.02 |
| ResNet-34 A | 25.03 | 7.76 |
| ResNet-34 B | 24.52 | 7.46 |
| ResNet-34 C | 24.19 | 7.40 |
| ResNet-50 | 22.85 | 6.71 |
| ResNet-101 | 21.75 | 6.05 |
| ResNet-152 | **21.43** | **5.71** |

Table 3. Error rates (%, 10-crop testing) on ImageNet validation. VGG-16 is based on our test. ResNet-50/101/152 are of option B that only uses projections for increasing dimensions.

BN [16]

34

3

BN

| model | top-1 err. | top-5 err. |
|---|---|---|
| VGG-16 [40] | 28.07 | 9.33 |
| GoogLeNet [43] | - | 9.15 |
| PReLU-net [12] | 24.27 | 7.38 |
| plain-34 | 28.54 | 10.02 |
| ResNet-34 A | 25.03 | 7.76 |
| ResNet-34 B | 24.52 | 7.46 |
| ResNet-34 C | 24.19 | 7.40 |
| ResNet-50 | 22.85 | 6.71 |
| ResNet-101 | 21.75 | 6.05 |
| ResNet-152 | **21.43** | **5.71** |

3  ImageNet        %  10             VGG16                    ResNet-50/101/152      B

**Residual Networks**. Next we evaluate 18-layer and 34-layer residual nets (ResNets). The baseline architectures are the same as the above plain nets, expect that a shortcut connection is added to each pair of 3×3 filters as in Fig. 3 (right). In the first comparison (Table 2 and Fig. 4 right), we use identity mapping for all shortcuts and zero-padding for increasing dimensions (option A). So they have no extra parameter compared to the plain counterparts.

18    34          ResNets                          3              3×3                        2    4
                       A

We have three major observations from Table 2 and Fig. 4. First, the situation is reversed with residual learning —— the 34-layer ResNet is better than the 18-layer ResNet (by 2.8%). More importantly, the 34-layer ResNet exhibits considerably lower training error and is generalizable to the validation data. This indicates that the degradation problem is well addressed in this setting and we manage to obtain accuracy gains from increased depth.

2    4                                              ——34  ResNet  18  ResNet    2.8            34  ResNet

Second, compared to its plain counterpart, the 34-layer ResNet reduces the top-1 error by 3.5% (Table 2), resulting from the successfully reduced training error (Fig. 4 right vs. left). This comparison verifies the effectiveness of residual learning on extremely deep systems.

34  ResNet        3.5%  top-1

Last, we also note that the 18-layer plain/residual nets are comparably accurate (Table 2), but the 18-layer ResNet converges faster (Fig. 4 right vs. left). When the net is " not overly deep" (18 layers here), the current SGD solver is still able to find good solutions to the plain net. In this case, the ResNet eases the optimization by providing faster convergence at the early stage.

18      /      2    18   ResNet      4      "     "    18      SGD

ResNet

**Identity vs. Projection Shortcuts**. We have shown that parameter-free, identity shortcuts help with training. Next we investigate projection shortcuts (Eqn.(2)). In Table 3 we compare three options: (A) zero-padding shortcuts are used for increasing dimensions, and all shortcuts are parameter-free (the same as Table 2 and Fig. 4 right); (B) projection shortcuts are used for increasing dimensions, and other shortcuts are identity; and (C) all shortcuts are projections.

2      3      (A)

2   4      (B)      C

Table 3 shows that all three options are considerably better than the plain counterpart. B is slightly better than A. We argue that this is because the zero-padded dimensions in A indeed have no residual learning. C is marginally better than B, and we attribute this to the extra parameters introduced by many (thirteen) projection shortcuts. But the small differences among A/B/C indicate that projection shortcuts are not essential for addressing the degradation problem. So we do not use option C in the rest of this paper, to reduce memory/time complexity and model sizes. Identity shortcuts are particularly important for not increasing the complexity of the bottleneck architectures that are introduced below.

3      B   A      A      C   B
A/B/C      C      /

**Deeper Bottleneck Architectures**. Next we describe our deeper nets for ImageNet. Because of concerns on the training time that we can afford, we modify the building block as a bottleneck design. For each residual function $F$ , we use a stack of 3 layers instead of 2 (Fig. 5). The three layers are 1×1, 3×3, and 1×1 convolutions, where the 1×1 layers are responsible for reducing and then increasing (restoring) dimensions, leaving the 3×3 layer a bottleneck with smaller input/output dimensions. Fig. 5 shows an example, where both designs have similar time complexity.

ImageNet               $F$      3      2
5      1×1   3×3   1×1      1×1      3×3      /      5

The parameter-free identity shortcuts are particularly important for the bottleneck architectures. If the identity shortcut in Fig. 5 (right) is replaced with projection, one can show that the time complexity and model size are doubled, as the shortcut is connected to the two high-dimensional ends. So identity shortcuts lead to more efficient models for the bottleneck designs.

5

**50-layer ResNet**: We replace each 2-layer block in the 34-layer net with this 3-layer bottleneck block, resulting in a 50-layer ResNet (Table 1). We use option B for increasing dimensions. This model has 3.8 billion FLOPs.

**50 ResNet**      3      34      2      50   ResNet    1      B      38   FLOP

**101-layer and 152-layer ResNet**: We construct 101-layer and 152-layer ResNets by using more 3-layer blocks (Table 1). Remarkably, although the depth is significantly increased, the 152-layer ResNet (11.3 billion FLOPs) still has lower complexity than VGG-16/19 nets (15.3/19.6 billion FLOPs).

**101 152 ResNet**      3      101    152   ResNets    1      152   ResNet   113   FLOP      VGG-16/19      153/196
FLOP

The 50/101/152-layer ResNets are more accurate than the 34-layer ones by considerable margins (Table 3 and 4). We do not observe the degradation problem and thus enjoy significant accuracy gains from considerably increased depth. The benefits of depth are witnessed for all evaluation metrics (Table 3 and 4).

50/101/152   ResNet   34   ResNet      3   4      3
4

**Comparisons with State-of-the-art Methods**. In Table 4 we compare with the previous best single-model results. Our baseline 34-layer ResNets have achieved very competitive accuracy. Our 152-layer ResNet has a single-model top-5 validation error of 4.49%. This single-model result outperforms all previous ensemble results (Table 5). We combine six models of different depth to form an ensemble (only with two 152-layer ones at the time of submitting). This leads to 3.57% top-5 error on the test set (Table 5). This entry won the 1st place in ILSVRC 2015.

| method | top-1 err. | top-5 err. |
|---|---|---|
| VGG [40] (ILSVRC'14) | - | $8.43^{\dagger}$ |
| GoogLeNet [43] (ILSVRC'14) | - | 7.89 |
| VGG [40] (v5) | 24.4 | 7.1 |
| PReLU-net [12] | 21.59 | 5.71 |
| BN-inception [16] | 21.99 | 5.81 |
| ResNet-34 B | 21.84 | 5.71 |
| ResNet-34 C | 21.53 | 5.60 |
| ResNet-50 | 20.74 | 5.25 |
| ResNet-101 | 19.87 | 4.60 |
| ResNet-152 | **19.38** | **4.49** |

Table 4. Error rates (%) of single-model results on the ImageNet validation set (except reported on the test set).

| method | top-5 err. (test) |
|---|---|
| VGG [40] (ILSVRC'14) | 7.32 |
| GoogLeNet [43] (ILSVRC'14) | 6.66 |
| VGG [40] (v5) | 6.8 |
| PReLU-net [12] | 4.94 |
| BN-inception [16] | 4.82 |
| **ResNet (ILSVRC'15)** | **3.57** |

Table 5. Error rates (%) of ensembles. The top-5 error is on the test set of ImageNet and reported by the test server.

4       34 ResNet       152 ResNet   4.49 top-5
5       152     3.5 top-5    5    2015
ILSVRC

| method | top-1 err. | top-5 err. |
|---|---|---|
| VGG [40] (ILSVRC'14) | - | 8.43[†] |
| GoogLeNet [43] (ILSVRC'14) | - | 7.89 |
| VGG [40] (v5) | 24.4 | 7.1 |
| PReLU-net [12] | 21.59 | 5.71 |
| BN-inception [16] | 21.99 | 5.81 |
| ResNet-34 B | 21.84 | 5.71 |
| ResNet-34 C | 21.53 | 5.60 |
| ResNet-50 | 20.74 | 5.25 |
| ResNet-101 | 19.87 | 4.60 |

| | | |
|---|---|---|
| ResNet-101 | 19.87 | 4.60 |
| **ResNet-152** | **19.38** | **4.49** |

4　　　　　ImageNet　　　　　% (　†　　　　　　　)

| method | top-5 err. (test) |
|---|---|
| VGG [40] (ILSVRC'14) | 7.32 |
| GoogLeNet [43] (ILSVRC'14) | 6.66 |
| VGG [40] (v5) | 6.8 |
| PReLU-net [12] | 4.94 |
| BN-inception [16] | 4.82 |
| **ResNet (ILSVRC'15)** | **3.57** |

5　　　　　　　　(%)  top-5　　　ImageNet

## 4.2. CIFAR-10 and Analysis

We conducted more studies on the CIFAR-10 dataset [20], which consists of 50k training images and 10k testing images in 10 classes. We present experiments trained on the training set and evaluated on the test set. Our focus is on the behaviors of extremely deep networks, but not on pushing the state-of-the-art results, so we intentionally use simple architectures as follows.

## 4.2. CIFAR-10

CIFAR-10　　　[20]　　　　　　　　　　10　　　　5　　　　　1

The plain/residual architectures follow the form in Fig. 3 (middle/right). The network inputs are 32× 32 images, with the per-pixel mean subtracted. The first layer is 3× 3 convolutions. Then we use a stack of 6n layers with 3× 3 convolutions on the feature maps of sizes {32, 16, 8} respectively, with 2n layers for each feature map size. The numbers of filters are {16, 32, 64} respectively. The subsampling is performed by convolutions with a stride of 2. The network ends with a global average pooling, a 10-way fully-connected layer, and softmax. There are totally 6n+ 2 stacked weighted layers. The following table summarizes the architecture:

| output map size | 32×32 | 16×16 | 8×8 |
|---|---|---|---|
| # layers | 1+2n | 2n | 2n |
| # filters | 16 | 32 | 64 |

When shortcut connections are used, they are connected to the pairs of 3×3 layers (totally 3n shortcuts). On this dataset we use identity shortcuts in all cases (i.e., option A), so our residual models have exactly the same depth, width, and number of parameters as the plain counterparts.

/        3    /                32×32                    3×3                {32,16,8}              3×3        6n
  2n            {16,32,64}          2                        10          softmax        6n+2

| output map size | 32×32 | 16×16 | 8×8 |
|---|---|---|---|
| # layers | 1+2n | 2n | 2n |
| # filters | 16 | 32 | 64 |

3×3          3n                                                    A

We use a weight decay of 0.0001 and momentum of 0.9, and adopt the weight initialization in [12] and BN [16] but with no dropout. These models are trained with a mini-batch size of 128 on two GPUs. We start with a learning rate of 0.1, divide it by 10 at 32k and 48k iterations, and terminate training at 64k iterations, which is determined on a 45k/5k train/val split. We follow the simple data augmentation in [24] for training: 4 pixels are padded on each side, and a 32×32 crop is randomly sampled from the padded image or its horizontal flip. For testing, we only evaluate the single view of the original 32×32 image.

              0.0001        0.9        [12]  BN[16]                          GPU                128              0.1     32k      48k
        10      64k                        45k/5k        /                  [24]                4
32×32                              32×32

We compare $n = \{3, 5, 7, 9\}$, leading to 20, 32, 44, and 56-layer networks. Fig. 6 (left) shows the behaviors of the plain nets. The deep plain nets suffer from increased depth, and exhibit higher training error when going deeper. This phenomenon is similar to that on ImageNet (Fig. 4, left) and on MNIST (see [41]), suggesting that such an optimization difficulty is a fundamental problem.
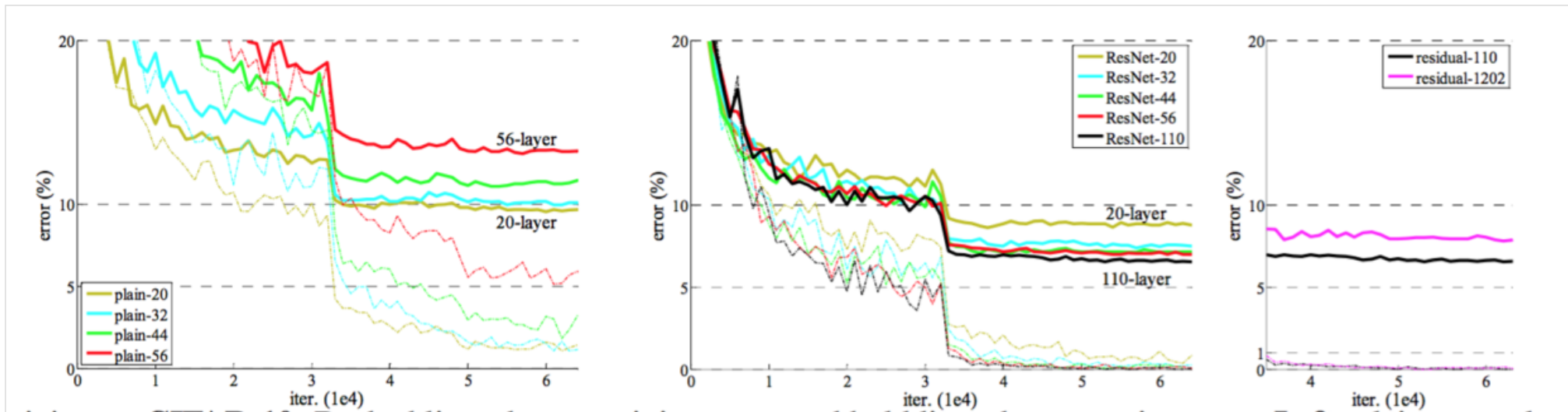
Figure 6. Training on CIFAR-10. Dashed lines denote training error, and bold lines denote testing error. Left: plain networks. The error of plain-110 is higher than 60% and not displayed. Middle: ResNets. Right: ResNets with 110 and 1202 layers.

$n = \{3, 5, 7, 9\}$        20    32    44    56                  6
ImageNet      4          MNIST        [41]



6    CIFAR-10                        110              60%              ResNet      110  ResNet  1202  ResNet

Fig. 6 (middle) shows the behaviors of ResNets. Also similar to the ImageNet cases (Fig. 4, right), our ResNets manage to overcome the optimization difficulty and demonstrate accuracy gains when the depth increases.

6              ResNet        ImageNet            4            ResNet

We further explore $n = 18$ that leads to a 110-layer ResNet. In this case, we find that the initial learning rate of 0.1 is slightly too large to start converging. So we use 0.01 to warm up the training until the training error is below 80% (about 400 iterations), and then go back to 0.1 and continue training. The rest of the learning schedule is as done previously. This 110-layer network converges well (Fig. 6, middle). It has fewer parameters than other deep and thin networks such as FitNet [34] and Highway [41] (Table 6), yet is among the state-of-the-art results (6.43%, Table 6).

| method | # layers | # params | error (%) |
|:---:|:---:|:---:|:---|
| Maxout [9] | | | 9.38 |
| NIN [25] | | | 8.81 |
| DSN [24] | | | 8.22 |
| FitNet [34] | 19 | 2.5M | 8.39 |
| Highway [41, 42] | 19 | 2.3M | 7.54 (7.72±0.16) |
| Highway [41, 42] | 32 | 1.25M | 8.80 |
| ResNet | 20 | 0.27M | 8.75 |
| ResNet | 32 | 0.46M | 7.51 |
| ResNet | 44 | 0.66M | 7.17 |
| ResNet | 56 | 0.85M | 6.97 |
| ResNet | 110 | 1.7M | **6.43** (6.61±0.16) |
| ResNet | 1202 | 19.4M | 7.93 |

Table 6. Classification error on the CIFAR-10 test set. All methods are with data augmentation. For ResNet-110, we run it 5 times and show " best (mean±std)" as in [42].

| method | | | error (%) |
|---|---|---|---|
| Maxout [9] | | | 9.38 |
| NIN [25] | | | 8.81 |
| DSN [24] | | | 8.22 |
| | # layers | # params | |
| FitNet [34] | 19 | 2.5M | 8.39 |
| Highway [41, 42] | 19 | 2.3M | 7.54 (7.72$\pm$0.16) |
| Highway [41, 42] | 32 | 1.25M | 8.80 |
| ResNet | 20 | 0.27M | 8.75 |
| ResNet | 32 | 0.46M | 7.51 |
| ResNet | 44 | 0.66M | 7.17 |
| ResNet | 56 | 0.85M | 6.97 |
| ResNet | 110 | 1.7M | **6.43** (6.61$\pm$0.16) |
| ResNet | 1202 | 19.4M | 7.93 |

6    CIFAR-10                                      ResNet-110        [42]                5        "        (mean±std)"

**Analysis of Layer Responses**. Fig. 7 shows the standard deviations (std) of the layer responses. The responses are the outputs of each 3× 3 layer, after BN and before other nonlinearity (ReLU/addition). For ResNets, this analysis reveals the response strength of the residual functions. Fig. 7 shows that ResNets have generally smaller responses than their plain counterparts. These results support our basic motivation (Sec.3.1) that the residual functions might be generally closer to zero than the non-residual functions. We also notice that the deeper ResNet has smaller magnitudes of responses, as evidenced by the comparisons among ResNet-20, 56, and 110 in Fig. 7. When there are more layers, an individual layer of ResNets tends to modify the signal less.

| 7 | std | 3×3 | BN | ReLU/ | ResNets | | 7 | ResNet |
|---|---|---|---|---|---|---|---|---|
| | | 3.1 | | | | ResNet | 7 | ResNet-20 |
| 56 110 | | ResNet | | | | | | |

**Exploring Over 1000 layers**. We explore an aggressively deep model of over 1000 layers. We set $n = 200$ that leads to a 1202-layer network, which is trained as described above. Our method shows no optimization difficulty, and this $10^3$-layer network is able to achieve training error <0.1% (Fig. 6, right). Its test error is still fairly good (7.93%, Table 6).

| 1000 | 1000 | | $n = 200$ | 1202 | | $10^3$ | <0.1 |
|---|---|---|---|---|---|---|---|
| 6 | | 7.93 | 6 | | | | |

But there are still open problems on such aggressively deep models. The testing result of this 1202-layer network is worse than that of our 110-layer network, although both have similar training error. We argue that this is because of overfitting. The 1202-layer network may be unnecessarily large (19.4M) for this small dataset. Strong regularization such as maxout [9] or dropout [13] is applied to obtain the best results ([9, 25, 24, 34]) on this dataset. In this paper, we use no maxout/dropout and just simply impose regularization via deep and thin architectures by design, without distracting from the focus on the difficulties of optimization. But combining with stronger regularization may improve results, which we will study in the future.

| | | 1202 | 110 | | | |
|---|---|---|---|---|---|---|
| 1202 | 19.4M | | maxout[9] | dropout[13] | [9,25,24,34] | maxout/dropout |

## 4.3. Object Detection on PASCAL and MS COCO

Our method has good generalization performance on other recognition tasks. Table 7 and 8 show the object detection baseline results on PASCAL VOC 2007 and 2012 [5] and COCO [26]. We adopt Faster R-CNN [32] as the detection method. Here we are interested in the improvements of replacing VGG-16 [40] with ResNet-101. The detection implementation (see appendix) of using both models is the same, so the gains can only be attributed to better networks. Most remarkably, on the challenging COCO dataset we obtain a 6.0% increase in COCO's standard metric (mAP@ [.5, .95]), which is a 28% relative improvement. This gain is solely due to the learned representations.

| training data | 07+12 | 07++12 |
|---|---|---|
| test data | VOC 07 test | VOC 12 test |
| VGG-16 | 73.2 | 70.4 |
| ResNet-101 | **76.4** | **73.8** |

Table 7. Object detection mAP (%) on the PASCAL VOC 2007/2012 test sets using baseline Faster R-CNN. See also appendix for better results.

| metric | mAP@.5 | mAP@[.5, .95] |
|--------|--------|---------------|
| VGG-16 | 41.5 | 21.2 |
| ResNet-101 | **48.4** | **27.2** |

Table 8. Object detection mAP (%) on the COCO validation set using baseline Faster R-CNN. See also appendix for better results.

### 4.3. PASCAL   MS COCO

7      8          PASCAL VOC 2007   2012[5]      COCO[26]                                          R-CNN[32]
ResNet-101      VGG-16[40]                                                                         COCO          COCO              mAP@
[.5  .95]        6.0              28

| training data | 07+12 | 07++12 |
|---------------|-------|--------|
| test data | VOC 07 test | VOC 12 test |
| VGG-16 | 73.2 | 70.4 |
| ResNet-101 | **76.4** | **73.8** |

7    PASCAL VOC 2007/2012          Faster R-CNN          mAP(%)

| metric | mAP@.5 | mAP@[.5, .95] |
|--------|--------|---------------|
| VGG-16 | 41.5 | 21.2 |
| ResNet-101 | **48.4** | **27.2** |

8    COCO              Faster R-CNN          mAP(%)

Based on deep residual nets, we won the 1st places in several tracks in ILSVRC & COCO 2015 competitions: ImageNet detection, ImageNet localization, COCO detection, and COCO

segmentation. The details are in the appendix.

ILSVRC & COCO 2015                    ImageNet          ImageNet          COCO          COCO

# References

[1] Y.Bengio,P.Simard,andP.Frasconi.Learning long-term dependencies with gradient descent is difficult. IEEE Transactions on Neural Networks, 5(2):157–166, 1994.

[2] C. M. Bishop. Neural networks for pattern recognition. Oxford university press, 1995.

[3] W. L. Briggs, S. F. McCormick, et al. A Multigrid Tutorial. Siam, 2000.

[4] K. Chatfield, V. Lempitsky, A. Vedaldi, and A. Zisserman. The devil is in the details: an evaluation of recent feature encoding methods. In BMVC, 2011.

[5] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The Pascal Visual Object Classes (VOC) Challenge. IJCV, pages 303–338, 2010.

[6] R. Girshick. Fast R-CNN. In ICCV, 2015.

[7] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In CVPR, 2014.

[8] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In AISTATS, 2010.

[9] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio. Maxout networks. arXiv:1302.4389, 2013.

[10] K.Heand J.Sun. Convolutional neural networks at constrained time cost. In CVPR, 2015.

[11] K.He, X.Zhang, S.Ren, and J.Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In ECCV, 2014.

[12] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In ICCV, 2015.

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

[14] S. Hochreiter. Untersuchungen zu dynamischen neuronalen netzen. Diploma thesis, TU Munich, 1991.

[15] S. Hochreiter and J. Schmidhuber. Long short-term memory. Neural computation, 9(8):1735–1780, 1997.

[16] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In ICML, 2015.

[17] H.Jegou, M.Douze, and C.Schmid. Product quantization for nearest neighbor search. TPAMI, 33, 2011.

[18] H. Jegou, F. Perronnin, M. Douze, J. Sanchez, P. Perez, and C. Schmid. Aggregating local image descriptors into compact codes. TPAMI, 2012.

[19] 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:1408.5093, 2014.

[20] A. Krizhevsky. Learning multiple layers of features from tiny images. Tech Report, 2009.

[21] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In NIPS, 2012.

[22] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to hand-written zip code recognition. Neural computation, 1989.

[23] Y.LeCun,L.Bottou,G.B.Orr,and K.-R.Muller. Efficient back prop. In Neural Networks: Tricks of the Trade, pages 9–50. Springer, 1998.

[24] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu. Deeply-supervised nets. arXiv:1409.5185, 2014.

[25] M. Lin, Q. Chen, and S. Yan. Network in network. arXiv:1312.4400,2013.

[26] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollar, and C. L. Zitnick. Microsoft COCO: Common objects in context. In ECCV. 2014.

[27] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In CVPR, 2015.

[28] G. Montufar, R. Pascanu, K. Cho, and Y. Bengio. On the number of linear regions of deep neural networks. In NIPS, 2014.

[29] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In ICML, 2010.

[30] F. Perronnin and C. Dance. Fisher kernels on visual vocabularies for image categorization. In CVPR, 2007.

[31] T. Raiko, H. Valpola, and Y. LeCun. Deep learning made easier by linear transformations in perceptrons. In AISTATS, 2012.

[32] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In NIPS, 2015.

[33] B. D. Ripley. Pattern recognition and neural networks. Cambridge university press, 1996.

[34] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio. Fitnets: Hints for thin deep nets. In ICLR, 2015.

[35] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. arXiv:1409.0575, 2014.

[36] A. M. Saxe, J. L. McClelland, and S. Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. arXiv:1312.6120, 2013.

[37] N.N.Schraudolph. Accelerated gradient descent by factor-centering decomposition. Technical report, 1998.

[38] N. N. Schraudolph. Centering neural network gradient factors. In Neural Networks: Tricks of the Trade, pages 207–226. Springer, 1998.

[39] P.Sermanet, D.Eigen, X.Zhang, M.Mathieu, R.Fergus, and Y.LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. In ICLR, 2014.

[40] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In ICLR, 2015.

[41] R. K. Srivastava, K. Greff, and J. Schmidhuber. Highway networks. arXiv:1505.00387, 2015.

[42] R. K. Srivastava, K. Greff, and J. Schmidhuber. Training very deep networks. 1507.06228, 2015.

[43] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In CVPR, 2015.

[44] R. Szeliski. Fast surface interpolation using hierarchical basis functions. TPAMI, 1990.

[45] R. Szeliski. Locally adapted hierarchical basis preconditioning. In SIGGRAPH, 2006.

[46] T. Vatanen, T. Raiko, H. Valpola, and Y. LeCun. Pushing stochastic gradient towards second-order methods–backpropagation learning with transformations in nonlinearities. In Neural Information Processing, 2013.

[47] A. Vedaldi and B. Fulkerson. VLFeat: An open and portable library of computer vision algorithms, 2008.

[48] W. Venables and B. Ripley. Modern applied statistics with s-plus. 1999.

[49] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional neural networks. In ECCV, 2014.

# Deep Learning

ResNet        ——

YOLO        ——

612

26

1. 1. Deep Residual Learning for Image Recognition