# A Residual Learning Framework for Advanced Visual Recognition

**Dr. Narek I. Veloyan**
**Department of Computer Vision and Pattern Analysis Caspian Institute of Technology, Baku, Azerbaijan**

**Prof. Helena M. Stravik**
**Visual Computing and Artificial Intelligence Laboratory Skagen Technical University, Denmark**

**ABSTRACT**

The advancement of deep learning in computer vision has been largely driven by the development of increasingly deep neural networks. However, a fundamental challenge known as the "degradation problem" has hindered progress, where adding more layers to a suitably deep model leads to higher training error, preventing the network from benefiting from increased depth. This phenomenon is not caused by overfitting but by optimization difficulties inherent in training very deep architectures. To address this, we introduce a deep residual learning framework. Instead of expecting stacked layers to learn an underlying mapping directly, we reformulate them to learn a residual function with respect to the layer inputs. This is achieved by introducing "shortcut connections" that perform identity mapping, adding the input of a block of layers to its output. This reformulation simplifies the optimization process, as it is easier for the network to learn perturbations from an identity mapping than to learn the mapping from scratch.

We provide comprehensive empirical evidence on several benchmark datasets. On the ImageNet 2012 classification dataset, our residual networks (ResNets) are substantially deeper than previous models, with up to 152 layers, yet exhibit lower complexity than shallower networks like VGG. These ResNets easily overcome the degradation problem, showing consistent accuracy gains from increased depth and achieving a 3.57% top-5 error on the ImageNet test set with an ensemble model. We also conducted experiments on the CIFAR-10 dataset, successfully training networks with over 1000 layers, demonstrating that our framework effectively resolves the core optimization issues. Furthermore, the representations learned by our deep residual networks generalize exceptionally well to other computer vision tasks. When used as a backbone for object detection on PASCAL VOC and MS COCO, our models achieve state-of-the-art results, with a 28% relative improvement on the COCO detection metric. Our findings establish deep residual learning as a fundamental and effective technique for training extremely deep neural networks, pushing the boundaries of what is possible in image recognition.

**Keywords:** Deep Learning, Residual Learning, Image Recognition, Convolutional Neural Networks, Degradation Problem, Shortcut Connections

## 1. Introduction

### 1.1 Broad Background and Historical Context

In recent years, deep convolutional neural networks (CNNs) have become the cornerstone of modern computer vision, leading to unprecedented breakthroughs in tasks ranging from image classification to object detection and semantic segmentation. A central theme in the evolution of these architectures has been the pursuit of greater depth. The rationale is intuitive: deeper networks can learn a more complex hierarchy of features. The initial layers capture low-level features like edges and textures, subsequent layers compose these into mid-level features like object parts, and the final layers assemble them into high-level, abstract representations of entire objects and scenes [50]. This end-to-end, hierarchical feature learning is what gives deep models their power.

The history of deep learning is marked by a clear trend towards increasing network depth. Early models had only a handful of layers, but the landmark success of AlexNet [21] on the ImageNet challenge demonstrated the potential of deeper architectures. This was followed by a wave of even deeper models. The VGG networks [41] featured a simple, uniform architecture with up to 19 layers, showing that significant performance gains could be achieved simply by pushing depth. Concurrently, the GoogLeNet architecture [44], with its "Inception" modules, reached 22 layers, introducing structural innovations to manage computational cost while increasing depth and width. These models established that network depth is a critical dimension for improving performance on challenging visual recognition tasks. This has fueled a fundamental question within the research community: Is learning better networks as simple as stacking more layers?

1.2 Critical Literature Review

Answering this question has been far from straightforward. The path to deeper networks has historically been fraught with obstacles. The most widely recognized challenge is the problem of vanishing or exploding gradients [1, 9]. As gradients are backpropagated through many layers, they can either shrink exponentially towards zero or grow uncontrollably, making it impossible for the network to learn effectively. For a long time, this problem severely hampered the training of networks beyond a shallow depth. However, this issue has been largely mitigated by the development of more sophisticated initialization schemes [9, 37], which set the initial weights in a way that preserves signal variance, and the introduction of intermediate normalization layers. Batch Normalization (BN) [16], in particular, has become a transformative technique. By normalizing the inputs to each layer for each mini-batch, BN ensures that gradients remain in a healthy range, enabling networks with tens of layers to begin converging with standard stochastic gradient descent (SGD).

With the vanishing gradient problem largely addressed, a more subtle and perplexing issue was exposed: the degradation problem. As networks become deeper, their performance saturates and then begins to degrade rapidly. Counter-intuitively, this degradation is not a result of overfitting. In an overfitting scenario, a model performs well on the training data but poorly on unseen test data. With the degradation problem, the *training error itself* increases as more layers are added to a suitably deep model [11]. This implies that deeper models are fundamentally harder to optimize, even though they possess a larger parameter space.

This phenomenon is paradoxical. Consider a shallow, trained network and a deeper counterpart constructed by adding extra layers on top of it. A trivial solution exists for the deeper model that should allow it to perform at least as well as the shallow one: the added layers can simply learn an identity mapping, passing the output of the shallow network through unchanged, while the original layers retain their learned weights. The existence of this constructed solution means a deeper model's training error should be no higher than its shallower counterpart's. The empirical evidence of degradation suggests that standard optimization algorithms and network formulations struggle to find this trivial solution, or solutions of comparable quality, in a feasible amount of time.

The idea of using "shortcut connections" that skip one or more layers has a long history in the field of neural networks [2, 34, 49]. These connections provide an alternative path for gradient flow and information propagation. In some earlier works, intermediate layers were connected to auxiliary classifiers to help with gradient propagation in very deep networks [44, 24]. More recently, and developed concurrently with our work, "Highway Networks" [42, 43] introduced shortcut connections with gating functions. These gates, inspired by the Long Short-Term Memory (LSTM) cell [15], dynamically control the flow of information through the shortcut path and the transformation path. However, these gates are data-dependent and introduce additional parameters that must be learned.

The concept of solving problems by reformulating them in terms of residuals also has precedent in computer vision and other scientific domains. In image retrieval and classification, powerful shallow representations like VLAD [18] and Fisher Vectors [30] encode features as residual vectors relative to a dictionary. In numerical analysis, the widely used Multigrid method for solving partial differential equations reformulates the problem on multiple scales, where each scale is responsible for correcting the residual error from the coarser scale below it [3]. Similarly, hierarchical basis preconditioning methods rely on representing residual vectors between scales to accelerate convergence [45, 46]. These examples from different fields suggest a common principle: reformulating a problem to focus on learning a residual or a perturbation can often make optimization significantly more tractable.

1.3 Research Gap

Despite advances like Batch Normalization [16] that allow for the training of moderately deep networks, the fundamental degradation problem has remained a key barrier to exploring the full potential of network depth. Existing architectures and optimization methods appear ill-suited to learning identity mappings within a deep stack of non-linear layers, thereby preventing deeper models from achieving their theoretical performance potential. While Highway Networks [42, 43] presented a potential solution with gated shortcuts, they introduced extra parametric complexity and had not yet demonstrated consistent accuracy gains in extremely deep settings (e.g., over a hundred layers). There existed a clear need for a simple, general, and effective framework that could systematically address the degradation problem and enable the successful training of truly deep networks, unlocking the benefits of depth without encountering prohibitive optimization barriers.

1.4 Objectives and Hypotheses

Based on the identified research gap, this study was designed with the following objectives and hypotheses:

- **Objective 1:** To introduce and formally define a "deep

residual learning" framework aimed at easing the training of neural networks that are substantially deeper than those previously used.

- **Objective 2:** To empirically validate this framework by conducting extensive experiments on large-scale benchmark datasets, demonstrating its ability to overcome the degradation problem and achieve significant accuracy improvements from increased network depth.

- **Objective 3:** To leverage this framework to train extremely deep models for state-of-the-art performance on a suite of competitive image recognition tasks, including image classification, object detection, and localization.

Our investigation was guided by two primary hypotheses:

- **Hypothesis 1:** It is easier for a stack of non-linear layers to approximate a residual function (i.e., the difference between a target mapping and an identity mapping) than it is to approximate the original, unreferenced target mapping directly.

- **HypothesIS 2:** By explicitly reformulating network layers to learn residual functions using parameter-free identity shortcut connections, we can effectively address the degradation problem. This will allow for the successful training of networks that are hundreds or even thousands of layers deep, leading to lower training error and superior generalization performance compared to both shallower networks and conventional "plain" deep networks of the same depth.

## 2. Methods

2.1 Research Design

The core of our research methodology was a rigorous, comparative experimental design intended to isolate and evaluate the effects of our proposed residual learning framework. The design centered on a direct comparison between conventional "plain" deep neural networks and our proposed "residual" networks (ResNets). The primary independent variable was the network architecture—specifically, the presence or absence of shortcut connections that constitute the residual blocks.

To test our hypotheses, we systematically varied a second key variable: network depth. We designed and implemented pairs of plain and residual networks with matching depths and overall computational complexity. This included moderately deep networks (e.g., 18 and 34 layers) to directly observe the onset of the degradation

problem, as well as extremely deep networks (e.g., 50, 101, 152, and even over 1000 layers) to push the limits of the proposed framework.

The dependent variables were the model's performance metrics, including training error, validation/test error, and task-specific metrics like mean Average Precision (mAP). By analyzing the learning curves (training and validation error versus training iterations) for plain and residual models of increasing depth, we could directly assess whether our framework mitigates the degradation problem and enables accuracy gains from added depth.

Our evaluation was conducted on several large-scale, publicly available computer vision datasets. This ensures the reproducibility of our results and allows for a fair comparison against existing state-of-the-art methods. The design also included ablation studies to investigate the specific contributions of different components of our framework, such as the type of shortcut connection used.

2.2 Participants / Sample

As this research does not involve human subjects, this section details the datasets used for training and evaluation.

- **ImageNet 2012 Classification Dataset:** This is a large-scale dataset that has been the standard benchmark for image classification research [36]. It consists of 1,000 object classes, with approximately 1.28 million images for training, 50,000 for validation, and 100,000 for testing. The primary evaluation metrics are top-1 and top-5 error rates, which measure whether the correct class is within the model's top 1 or top 5 predictions, respectively. This dataset was our primary testbed for evaluating deep architectures and comparing them with prior work.

- **CIFAR-10 Dataset:** This dataset consists of 60,000 32×32 color images in 10 classes, with 50,000 training images and 10,000 test images [20]. Due to its smaller size, it allows for faster experimentation, making it an ideal environment for controlled studies on the optimization behavior of extremely deep networks (e.g., >1000 layers) and for detailed analysis of layer responses.

- **PASCAL VOC 2007 & 2012 and MS COCO Datasets:** To assess the generalization capability of the features learned by our models, we used them as backbones for the downstream task of object detection. The PASCAL VOC datasets [5] are standard benchmarks for this task. The Microsoft COCO (Common Objects in Context) dataset [26] is a more recent and challenging benchmark with 80 object categories and more complex scenes. Performance on these datasets is

measured by mean Average Precision (mAP).

## 2.3 Materials and Apparatus

*Network Architectures*

Our experiments involved two main families of architectures: plain networks and residual networks.

- **Plain Networks:** Our plain network baselines were primarily inspired by the design philosophy of VGG nets [41]. These architectures are characterized by their simplicity and depth, consisting mainly of stacked 3×3 convolutional layers. We followed two main design rules: (i) for a given output feature map size, all layers have the same number of filters; and (ii) when the feature map size is halved (through a convolutional layer with a stride of 2), the number of filters is doubled to maintain roughly the same time complexity per layer. The networks end with a global average pooling layer followed by a 1000-way fully-connected layer with a softmax activation for classification.

- Residual Networks (ResNets): The residual networks were created by modifying the plain networks. The core innovation is the "residual block," which is formally defined as:
  $$y = F(x, \{W_i\}) + x \quad (1)$$

  Here, x and y are the input and output vectors of the block. The function $F(x, \{W_i\})$ represents the residual mapping to be learned by a few stacked layers (e.g., two or three convolutional layers). For a block with two layers, F can be expressed as $F = W_2 \sigma(W_1 x)$, where σ denotes the Rectified Linear Unit (ReLU) activation function [29] and biases are omitted for simplicity. The operation +x is performed by a shortcut connection and element-wise addition. A final ReLU activation is applied after the addition.

- **Shortcut Connections:** The implementation of the shortcut connection depends on the dimensions of x and F.

  - **Identity Shortcuts:** When the input and output dimensions are the same, the shortcut connection simply performs an identity mapping. This is the most common case and is parameter-free, adding no computational complexity to the model.

  - **Projection Shortcuts:** When the dimensions change (e.g., when the number of channels is increased), the identity shortcut is no longer sufficient. We explored two options: (A) The shortcut still performs an identity mapping, but extra zero entries are padded to the input to match the increased dimension. This option remains parameter-free. (B) A linear projection, implemented as a 1×1 convolution, is used in the shortcut connection to explicitly match the output dimension. This introduces additional parameters. For both options, when the shortcut spans feature maps of different sizes, a stride of 2 is used. Our experiments showed that option B yields slightly better results, but the economical, parameter-free option A is sufficient to address the degradation problem.

- **Bottleneck Architectures:** For training very deep networks (e.g., 50 layers or more) efficiently, we introduced a "bottleneck" design for the residual block. Instead of two stacked 3×3 layers, we use a stack of three layers: a 1×1, a 3×3, and another 1×1 convolution. The first 1×1 layer reduces the number of channels, creating a computational bottleneck for the expensive 3×3 layer. The final 1×1 layer then restores the number of channels. This design has a similar time complexity to the original two-layer block but allows for much greater depth with a manageable number of parameters.

*Software and Hardware*

All models were implemented and trained using the Caffe deep learning framework [19]. Training was accelerated using a distributed setup with multiple high-performance GPUs.

## 2.4 Data Collection Procedure

This section outlines the standardized training and fine-tuning procedures used across our experiments.

- **ImageNet Classification Training:**

  - **Preprocessing and Augmentation:** We followed standard practices [21, 41]. Images were resized so their shorter side was randomly sampled from the range [256, 480] for scale augmentation. A 224×224 crop was then randomly sampled from the image or its horizontal flip. The per-pixel mean, computed over the training set, was subtracted from each image. We also used the standard color augmentation scheme described in [21].

  - **Training Parameters:** All networks were trained from scratch using Stochastic Gradient Descent (SGD) with a mini-batch size of 256. The learning rate was initialized to 0.1 and was divided by 10

whenever the validation error plateaued. Models were trained for up to 60×104 iterations. We used a weight decay of 0.0001 and a momentum of 0.9. Following the practice in [16], we used Batch Normalization (BN) after each convolutional layer and before the activation function. We did not use dropout [14]. Weight initialization followed the method described in [13].

- **CIFAR-10 Classification Training:**

  - The training procedure was similar but adapted for the smaller dataset. Inputs were 32×32 images with the per-pixel mean subtracted. Data augmentation involved padding each side of an image with 4 pixels and then taking a random 32×32 crop from the padded image or its horizontal flip. Models were trained with a batch size of 128, starting with a learning rate of 0.1, which was divided by 10 at 32k and 48k iterations.

- **Object Detection Fine-tuning:**

  - For object detection tasks on PASCAL VOC and MS COCO, we used the Faster R-CNN framework [32]. Our ImageNet-trained ResNet models were used as the convolutional backbone, replacing the VGG-16 network used in the original Faster R-CNN paper. The models were then fine-tuned on the respective detection datasets. During fine-tuning, we fixed the parameters of the Batch Normalization layers, using the statistics computed on the ImageNet training set. This was done primarily to reduce memory consumption during the training of the detection models.

2.5 Data Analysis

- **Evaluation Metrics:**

  - **Classification:** Performance was measured using top-1 and top-5 error rates on the validation and test sets.

  - **Object Detection:** Performance was evaluated using mean Average Precision (mAP). For PASCAL VOC, we used mAP at an Intersection over Union (IoU) threshold of 0.5. For MS COCO, we report both mAP@0.5 and the standard COCO metric, which is the mAP averaged over IoU thresholds from 0.5 to 0.95 in steps of 0.05.

- **Testing Procedures:**

- For comparative classification studies, we used the standard 10-crop testing protocol, where predictions are averaged over 10 crops from an image (center, four corners, and their horizontal flips). For reporting our best results, we used a fully-convolutional testing approach, averaging classification scores over multiple scales of the input image.

- **Comparative and Ablation Analyses:**

  - The primary mode of analysis was the direct comparison of training and validation error curves for plain versus residual networks of varying depths. This allowed for a direct visualization and quantification of the degradation problem and the effectiveness of our proposed solution.

  - We conducted ablation studies on the ImageNet dataset to compare the performance of different shortcut connection strategies (zero-padding vs. projection) to understand their relative importance.

- **Layer Response Analysis:**

  - To gain insight into the internal dynamics of the networks, we analyzed the magnitude of layer responses. Specifically, we calculated the standard deviation of the outputs of each convolutional block (after BN but before the final ReLU/addition) across the entire test set. This allowed us to compare the response strength in plain vs. residual networks and across ResNets of different depths, providing evidence for our hypothesis that residual functions are generally closer to zero mappings.

## 3. Results

3.1 Preliminary Analyses: The Degradation Problem and its Resolution

Our initial experiments on the ImageNet dataset were designed to replicate the degradation problem with plain networks and to test whether our residual framework could resolve it. We trained an 18-layer and a 34-layer plain network. The results clearly demonstrated the degradation problem: the deeper 34-layer plain network exhibited a higher training error and validation error throughout the entire training process compared to its shallower 18-layer counterpart. This confirmed that despite having a solution space that is a superset of the shallower model, the 34-layer plain network was harder to optimize.

We then trained an 18-layer and a 34-layer residual network (ResNet) with architectures identical to the plain nets, except for the addition of identity shortcut connections. The results showed a complete reversal of the trend. The 34-layer ResNet consistently outperformed the 18-layer ResNet, achieving a significantly lower training error and a 2.8% lower top-1 error on the validation set.

A direct comparison between the 34-layer models was particularly revealing. The 34-layer ResNet achieved a top-1 validation error of 25.03%, a remarkable improvement over the 34-layer plain net's error of 28.54%. This 3.5% reduction in error stemmed directly from the ResNet's much lower training error, indicating that the residual learning formulation successfully addressed the optimization difficulties that plagued the plain network. This provided the first key piece of evidence that our framework was effective. Interestingly, the 18-layer plain net and 18-layer ResNet achieved comparable accuracy, though the ResNet converged faster in the early stages of training. This suggests that when a network is not overly deep, standard solvers are still capable of finding good solutions, but the residual formulation can still ease the optimization process.

### 3.2 Main Findings: Scaling Depth and Generalizing Performance

*Ablation Study on Shortcut Connections*

We investigated the importance of the specific type of shortcut connection used when feature dimensions increase. We compared three options for the 34-layer ResNet: (A) parameter-free identity shortcuts with zero-padding for increasing dimensions; (B) projection shortcuts for increasing dimensions and identity shortcuts elsewhere; and (C) projection shortcuts for all connections. All three options significantly outperformed the 34-layer plain network. Option B (24.52% error) was slightly better than A (25.03% error), likely because the zero-padded dimensions in A do not facilitate residual learning. Option C (24.19% error) was marginally better than B, which we attribute to the extra parameters from the additional projection shortcuts. The small differences between these options indicated that projection shortcuts are not essential for solving the degradation problem. Consequently, we primarily used the more economical option B for deeper models to reduce complexity and model size, while noting the particular importance of parameter-free identity shortcuts for the efficiency of our bottleneck architecture.

*Performance of Very Deep Residual Networks*

To push the boundaries of depth, we developed deeper ResNets using our computationally efficient bottleneck architecture. We trained 50-layer, 101-layer, and 152-layer ResNets. The results demonstrated a clear and consistent benefit from increasing depth. The top-1 validation errors were 22.85% (ResNet-50), 21.75% (ResNet-101), and 21.43% (ResNet-152), respectively. We did not observe any degradation; instead, accuracy improved significantly with considerably increased depth.

Our single-model 152-layer ResNet achieved a top-5 validation error of 5.71%, a result that was already better than all previously reported ensemble results on the ImageNet benchmark. By forming an ensemble of six residual networks of varying depths, we achieved a final top-5 error of 3.57% on the ImageNet 2012 test set, securing a top position in the classification challenge.

*Generalization to Object Detection*

To evaluate the generalization power of the learned representations, we used our pre-trained ResNet-101 as a feature backbone within the Faster R-CNN object detection framework [32], replacing the standard VGG-16 backbone. The results were compelling. On the PASCAL VOC 2007 test set, using ResNet-101 improved the mAP from 73.2% (with VGG-16) to 76.4%. The improvement was even more dramatic on the more challenging MS COCO dataset. ResNet-101 boosted the mAP@[.5, .95] from 21.2% to 27.2%, a 6.0% absolute increase, which corresponds to a 28% relative improvement. This substantial gain, achieved without any changes to the detection methodology, is solely attributable to the superior feature representations learned by the much deeper residual network.

### 3.3 Exploratory Findings: Experiments on CIFAR-10

We conducted further experiments on the CIFAR-10 dataset to analyze the behavior of extremely deep networks in a more controlled setting.

*Training Extremely Deep Networks*

We trained a series of plain and residual networks with depths of 20, 32, 44, and 56 layers. As on ImageNet, the plain networks suffered from degradation, with training and test error increasing with depth. In contrast, the ResNets showed consistent improvement, with the 56-layer model outperforming the shallower ones.

Encouraged by this, we explored even more extreme depths. We successfully trained a 110-layer ResNet, which achieved a state-of-the-art error rate of 6.43% on the CIFAR-10 test set. Pushing the limits further, we constructed and trained a 1202-layer network. Remarkably, this model showed no optimization difficulty and was able to achieve a training error of less than 0.1%. Its test error was 7.93%. While this is a strong result, it is slightly worse than that of the 110-

layer network. We attribute this performance gap to overfitting. The 1202-layer model, with its 19.4 million parameters, was likely unnecessarily large for the small CIFAR-10 dataset, and we did not employ strong regularization techniques like dropout [14] or maxout [10] in this study. The key finding, however, is that our framework enabled the *optimization* of a network of this unprecedented depth.

*Analysis of Layer Responses*

To test our hypothesis that residual functions are easier to learn because they may be close to zero mappings, we analyzed the standard deviations of layer responses (outputs after BN, before nonlinearity). The results showed that for any given depth, ResNets consistently had smaller response magnitudes than their plain counterparts. Furthermore, we observed a clear trend within the ResNet family: as the networks got deeper (from ResNet-20 to 56 to 110), the average magnitude of layer responses decreased. This suggests that when more layers are available, individual residual blocks tend to learn functions that make smaller modifications to the signal being passed through the identity shortcut. This evidence supports our core motivation that the residual functions being learned are indeed small perturbations, making the optimization task simpler.

## 4. Discussion

### 4.1 Interpretation of Findings

Our results provide a clear and compelling answer to the question of whether learning better networks is as simple as stacking more layers. The answer is yes, provided that the network is constructed within a framework that allows it to be optimized effectively. The central finding of this study is that our deep residual learning framework provides such a mechanism. It consistently and effectively overcomes the degradation problem that has long been a barrier to training very deep neural networks.

The success of residual learning can be interpreted through the lens of optimization. The degradation problem suggests that it is difficult for multiple non-linear layers to learn an identity mapping. By reformulating the learning objective of a block of layers from fitting a desired mapping H(x) to fitting a residual mapping F(x)=H(x)−x, we provide a powerful form of preconditioning. The identity shortcut effectively provides a baseline mapping. If the optimal function for a given block is indeed close to an identity transformation, the optimizer can easily achieve this by driving the weights of the layers in the residual function F(x) toward zero. This is a much easier optimization target than fitting an identity function with a stack of non-linear layers. Our analysis of layer responses, which showed that

response magnitudes in ResNets are smaller than in plain nets and decrease with depth, lends strong support to this interpretation. The layers appear to be learning small corrective adjustments to the identity mapping, rather than entirely new representations.

The successful training of a 1202-layer network on CIFAR-10 is a landmark result. While the model overfitted the small dataset, the fact that it could be optimized at all is profound. It suggests that the limits of network depth are no longer primarily constrained by the optimization algorithm itself. Instead, the practical limits are now dictated by other factors, such as the size and complexity of the dataset, the availability of computational resources, and the need for effective regularization strategies to manage the vast capacity of such models.

### 4.2 Comparison with Literature

The performance of our ResNets represents a significant advance over previous state-of-the-art models. Compared to prominent architectures like VGG [41] and GoogLeNet [44], our models achieve superior accuracy simply by enabling greater depth in a computationally manageable way. Unlike the handcrafted Inception modules of GoogLeNet [44], which are designed to capture features at multiple scales, our approach is simpler and more general, relying on a single, repeatable architectural motif—the residual block. The performance gains we demonstrate are not the result of increased model width or complex connectivity patterns, but are almost entirely attributable to the depth we are able to achieve.

Our work is closely related to the concurrent development of Highway Networks [42, 43], which also utilize shortcut connections to ease information flow. However, there are crucial differences. Highway Networks employ gating mechanisms with learnable parameters to regulate the flow of information through the shortcut. In contrast, our identity shortcuts are parameter-free and never "close," ensuring that the signal always propagates. This design choice not only makes our models more economical but may also contribute to their robust performance. While Highway Networks provided an important step forward, our residual learning framework was the first to demonstrate consistent accuracy gains with extremely deep architectures of over 100 layers.

Our framework should be viewed as complementary to other techniques that facilitate deep network training. While methods like Batch Normalization [16] and improved weight initialization [13] are essential for addressing the vanishing gradient problem, they do not, on their own, solve the degradation problem. Our experiments consistently used BN, yet the plain networks still suffered from degradation. This indicates that residual learning addresses

a distinct and more fundamental optimization challenge related to the function classes that deep networks can easily learn. The combination of BN and residual learning proved to be a powerful pairing.

The dramatic improvements seen in object detection tasks when using ResNet backbones further underscore the value of our approach. The gains over VGG-16 within the same Faster R-CNN [32] system are not due to a new detection algorithm but are a direct consequence of the superior quality of the learned feature representations. This reinforces a core tenet of deep learning: better representations lead to better performance on a wide range of downstream tasks [8, 7].

### 4.3 Strengths and Limitations

This study has several key strengths. First, the proposed residual learning framework is exceptionally **simple and general**. The core concept—the residual block with an identity shortcut—can be easily implemented in any standard deep learning library and incorporated into various network architectures. For identity shortcuts, it adds virtually no parameters or computational cost. Second, the framework is demonstrably **effective**. Our comprehensive experiments across multiple datasets and tasks show consistent, significant, and often state-of-the-art performance. Third, the framework successfully **enables extreme depth**, pushing the boundaries of trainable networks far beyond what was previously considered feasible and opening up new possibilities for architectural design.

Despite these strengths, our work has limitations. The **theoretical understanding** of why residual learning works so well is still incomplete. While our hypothesis about the ease of learning perturbations from an identity mapping is plausible and supported by empirical evidence, a more formal mathematical analysis of the optimization landscape of residual networks is an important area for future research [28].

Another limitation is the issue of **overfitting in extremely deep models**. Our experiment with the 1202-layer network highlighted that at such depths, model capacity can easily outstrip the information available in smaller datasets, leading to poor generalization. This study did not systematically investigate the interplay between residual learning and powerful regularization techniques like dropout [14] or maxout [10]. Combining these methods may be crucial for unlocking the full potential of extreme depth on a wider range of datasets.

Finally, while our bottleneck architecture improves **computational efficiency**, training very deep networks remains resource-intensive. Our research prioritized

demonstrating the limits of accuracy and depth, with less focus on developing models optimized for deployment in resource-constrained environments, a research direction pursued in works on model compression and efficiency [11].

### 4.4 Implications

The implications of this research are both practical and theoretical. On a practical level, deep residual networks have become a foundational architecture in computer vision. Pre-trained ResNet models are now standard backbones for a vast array of tasks, including object detection, instance segmentation [27], pose estimation, and facial recognition, often serving as the starting point for new research. Our work provided a clear and reliable recipe for building and training very deep, high-performing networks.

The theoretical implications are equally significant. This study helped shift the focus of deep learning optimization research. It demonstrated that beyond simply ensuring gradients can flow (addressing the vanishing gradient problem), the structure of the function space that a network must explore is critically important. Architectural choices, like the inclusion of residual connections, can act as a powerful form of preconditioning that dramatically simplifies the optimization task. This has inspired further research into how network architecture can be co-designed with optimization in mind.

### 4.5 Conclusion and Future Directions

In conclusion, this paper introduced a deep residual learning framework to address the degradation problem that hinders the training of very deep neural networks. By reformulating network layers to learn residual functions relative to their inputs via identity shortcut connections, we enabled the straightforward and effective training of networks of unprecedented depth. Our deep residual nets (ResNets) were shown to be easier to optimize than their "plain" counterparts, achieving consistent accuracy gains from greatly increased depth. This framework led to state-of-the-art results on the competitive ImageNet, PASCAL VOC, and MS COCO benchmarks for image classification and object detection.

Looking forward, several promising research avenues remain. A deeper theoretical analysis is needed to formally characterize the optimization landscape of residual networks and explain the mechanisms behind their ease of training. Another important direction is to explore the synergy between residual learning and other advanced architectural concepts, such as attention mechanisms, as well as its combination with stronger regularization methods to better harness the capacity of extremely deep models. Finally, the core principle of residual learning is generic. Investigating its application to other domains, such

as natural language processing with recurrent or transformer-based models [15], and to other data modalities like audio and time-series data, represents a fruitful area for future exploration.

References

[1] Y. Bengio, P. Simard, and P. 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] S. Gidaris and N. Komodakis. Object detection via a multi-region & semantic segmentation-aware cnn model. In ICCV, 2015.

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

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

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

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

[11] K. He and J. Sun. Convolutional neural networks at constrained time cost. In CVPR, 2015.

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

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

[14] 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.

[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 handwritten zip code recognition. Neural computation, 1989.

[23] Y. LeCun, L. Bottou, G. B. Orr, and K.-R. Muller. Efficient backprop. ¨ 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] S. Ren, K. He, R. Girshick, X. Zhang, and J. Sun. Object detection networks on convolutional feature maps. arXiv:1504.06066, 2015.

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

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

[36] 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.

[37] 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.

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

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

[40] 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.

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

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

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

[44] 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.

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

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

[47] 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.

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

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

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