

---

# MLP Coursework 2

---

s2617343

## Abstract

Deep neural networks have become the state-of-the-art in many standard computer vision problems thanks to their powerful representations and availability of large labeled datasets. While very deep networks allow for learning more levels of abstractions in their layers from the data, training these models successfully is a challenging task due to problematic gradient flow through the layers, known as vanishing/exploding gradient problem. In this report, we first analyze this problem in VGG models with 8 and 38 hidden layers on the CIFAR100 image dataset, by monitoring the gradient flow during training. We explore known solutions to this problem including batch normalization or residual connections, and explain their theory and implementation details. Our experiments show that batch normalization and residual connections effectively address the aforementioned problem and hence enable a deeper model to outperform shallower ones in the same experimental setup.

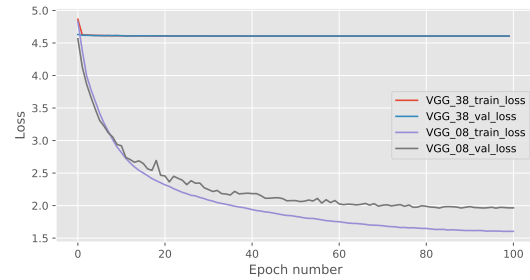
## 1. Introduction

Despite the remarkable progress of modern convolutional neural networks (CNNs) in image classification problems (Simonyan & Zisserman, 2014; He et al., 2016), training very deep networks is a challenging procedure. One of the major problems is the Vanishing Gradient Problem (VGP), a phenomenon where the gradients of the error function with respect to network weights shrink to zero, as they backpropagate to earlier layers, hence preventing effective weight updates. This phenomenon is prevalent and has been extensively studied in various deep neural networks including feedforward networks (Glorot & Bengio, 2010), RNNs (Bengio et al., 1993), and CNNs (He et al., 2016). Multiple solutions have been proposed to mitigate this problem by using weight initialization strategies (Glorot & Bengio, 2010), activation functions (Glorot & Bengio, 2010), input normalization (Bishop et al., 1995), batch normalization (Ioffe & Szegedy, 2015), and shortcut connections (He et al., 2016; Huang et al., 2017).

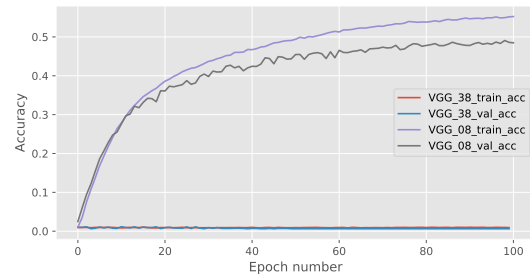
This report focuses on diagnosing the VGP occurring in the VGG38 model<sup>1</sup> and addressing it by implementing two standard solutions. In particular, we first study a “broken” network in terms of its gradient flow, L1 norm of gradients

---

<sup>1</sup>VGG stands for the Visual Geometry Group in the University of Oxford.



(a) Cross entropy error per epoch



(b) Classification accuracy per epoch

Figure 1. Training curves for VGG08 and VGG38 in terms of (a) cross-entropy error and (b) classification accuracy

with respect to its weights for each layer and contrast it to ones in the healthy and shallower VGG08 to pinpoint the problem. Next, we review two standard solutions for this problem, batch normalization (BN) (Ioffe & Szegedy, 2015) and residual connections (RC) (He et al., 2016) in detail and discuss how they can address the gradient problem. We first incorporate batch normalization (denoted as VGG38+BN), residual connections (denoted as VGG38+RC), and their combination (denoted as VGG38+BN+RC) to the given VGG38 architecture. We train the resulting three configurations, and VGG08 and VGG38 models on CIFAR100 (pronounced as ‘see far 100’) dataset and present the results. The results show that though separate use of BN and RC does mitigate the vanishing/exploding gradient problem, therefore enabling effective training of the VGG38 model, the best results are obtained by combining both BN and RC.

## 2. Identifying training problems of a deep CNN

### [ Question Figure 3 ]

Concretely, training deep neural networks typically involves



Figure 2. Gradient flow on VGG08

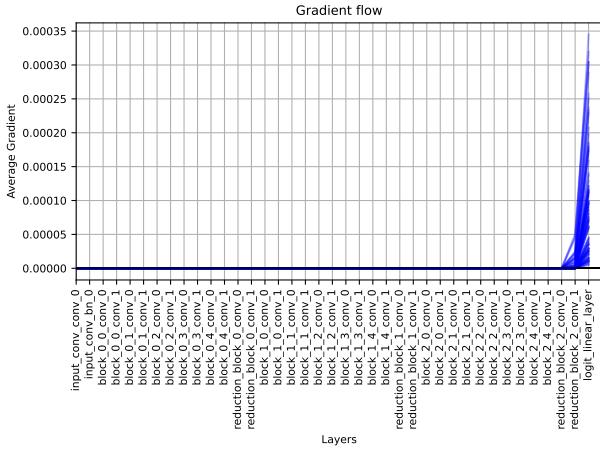


Figure 3. Gradient Flow on VGG38

three steps: forward pass, backward pass (or backpropagation algorithm (Rumelhart et al., 1986)) and weight update. The first step involves passing the input  $\mathbf{x}^{(0)}$  to the network and producing the network prediction and also the error value. In detail, each layer takes in the output of the previous layer and applies a non-linear transformation:

$$\mathbf{x}^{(l)} = f^{(l)}(\mathbf{x}^{(l-1)}; \mathbf{W}^{(l)}) \quad (1)$$

where  $(l)$  denotes the  $l$ -th layer in  $L$  layer deep network,  $f^{(l)}(\cdot, \mathbf{W}^{(l)})$  is a non-linear transformation for layer  $l$ , and  $\mathbf{W}^{(l)}$  are the weights of layer  $l$ . For instance,  $f^{(l)}$  is typically a convolution operation followed by an activation function in convolutional neural networks. The second step involves the backpropagation algorithm, where we calculate the gradient of an error function  $E$  (e.g. cross-entropy) for each layer’s weight as follows:

$$\frac{\partial E}{\partial \mathbf{W}^{(l)}} = \frac{\partial E}{\partial \mathbf{x}^{(L)}} \frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{x}^{(L-1)}} \cdots \frac{\partial \mathbf{x}^{(l+1)}}{\partial \mathbf{x}^{(l)}} \frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{W}^{(l)}}. \quad (2)$$

This step includes consecutive tensor multiplications between multiple partial derivative terms. The final step involves updating model weights by using the computed  $\frac{\partial E}{\partial \mathbf{W}^{(l)}}$

with an update rule. The exact update rule depends on the optimizer.

A notorious problem for training deep neural networks is the vanishing/exploding gradient problem (Bengio et al., 1993) that typically occurs in the backpropagation step when some of partial gradient terms in Eq. 2 includes values larger or smaller than 1. In this case, due to the multiple consecutive multiplications, the gradients *w.r.t.* weights can get exponentially very small (close to 0) or very large (close to infinity) and prevent effective learning of network weights.

Figures 2 and 3 depict the gradient flows through VGG architectures (Simonyan & Zisserman, 2014) with 8 and 38 layers respectively, trained and evaluated for a total of 100 epochs on the CIFAR100 dataset. **[By comparing Figure 2 and Figure 3, it is clear that VGG08 maintains significantly larger gradients across all layers (e.g., 0.025 at the bottom and 0.05-0.2 at the top). In contrast, VGG38 exhibits almost zero gradients beyond the bottom regression layer (0.00035) and dimensionality reduction block (0.000025), indicating a severe vanishing gradient problem.]**

**This occurs because gradients diminish exponentially during backpropagation due to the chain rule, leaving top-layer gradients effectively zero. As a result, parameter updates fail in most layers, halting learning. This is reflected in Figure 1, where VGG38’s accuracy remains at 0 throughout training, despite increasing epochs.]**

### 3. Background Literature

In this section we will highlight some of the most influential papers that have been central to overcoming the VGP in deep CNNs.

**Batch Normalization** (Ioffe & Szegedy, 2015) BN seeks to solve the problem of internal covariate shift (ICS), when distribution of each layer’s inputs changes during training, as the parameters of the previous layers change. The authors argue that without batch normalization, the distribution of each layer’s inputs can vary significantly due to the stochastic nature of randomly sampling mini-batches from your training set. Layers in the network hence must continuously adapt to these high variance distributions which hinders the rate of convergence gradient-based optimizers. This optimization problem is exacerbated further with network depth due to the updating of parameters at layer  $l$  being dependent on the previous  $l - 1$  layers.

It is hence beneficial to embed the normalization of training data into the network architecture after work from LeCun *et al.* showed that training converges faster with this addition (LeCun et al., 2012). Through standardizing the inputs to each layer, we take a step towards achieving the fixed distributions of inputs that remove the ill effects of ICS. Ioffe and Szegedy demonstrate the effectiveness of their technique through training an ensemble of BN networks

which achieve an accuracy on the ImageNet classification task exceeding that of humans in 14 times fewer training steps than the state-of-the-art of the time. It should be noted, however, that the exact reason for BN’s effectiveness is still not completely understood and it is an open research question (Santurkar et al., 2018).

**Residual networks (ResNet)** (He et al., 2016) A well-known way of mitigating the VGP is proposed by He *et al.* in (He et al., 2016). In their paper, the authors depict the error curves of a 20 layer and a 56 layer network to motivate their method. Both training and testing error of the 56 layer network are significantly higher than of the shallower one.

[Generally, as network capacity increases, the model’s expressive power is expected to improve. However, larger-capacity models can suffer from overfitting, where the model learns data noise instead of general patterns. Additionally, differences in performance between networks of varying capacities may result from the vanishing gradient problem, as shown in Figure 1. Larger models also introduce more parameters, requiring higher computational resources and leading to slower convergence.

In Figure 1 from (He et al., 2016), the 56-layer network underperforms compared to the 20-layer network, highlighting the model degradation problem. This issue is not necessarily due to overfitting but could stem from optimization challenges, such as insufficient convergence speed or difficulty propagating gradients effectively in deeper networks.] .

Residual networks, colloquially known as ResNets, aim to alleviate VGP through the incorporation of skip connections that bypass the linear transformations into the network architecture. The authors argue that this new mapping is significantly easier to optimize since if an identity mapping were optimal, the network could comfortably learn to push the residual to zero rather than attempting to fit an identity mapping via a stack of nonlinear layers. They bolster their argument by successfully training ResNets with depths exceeding 1000 layers on the CIFAR10 dataset. Prior to their work, training even a 100-layer was accepted as a great challenge within the deep learning community. The addition of skip connections solves the VGP through enabling information to flow more freely throughout the network architecture without the addition of neither extra parameters, nor computational complexity.

## 4. Solution overview

### 4.1. Batch normalization

BN has been a standard component in the state-of-the-art convolutional neural networks (He et al., 2016; Huang et al., 2017). Concretely, BN is a layer transformation that is performed to whiten the activations originating from each layer. As computing full dataset statistics at each training iteration

would be computationally expensive, BN computes batch statistics to approximate them. Given a minibatch of  $B$  training samples and their feature maps  $X = (\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^B)$  at an arbitrary layer where  $X \in \mathbb{R}^{B \times H \times W \times C}$ ,  $H, W$  are the height, width of the feature map and  $C$  is the number of channels, the batch normalization first computes the following statistics:

$$\mu_c = \frac{1}{BWH} \sum_{n=1}^B \sum_{i,j=1}^{H,W} \mathbf{x}_{cij}^n \quad (3)$$

$$\sigma_c^2 = \frac{1}{BWH} \sum_{n=1}^B \sum_{i,j=1}^{H,W} (\mathbf{x}_{cij}^n - \mu_c)^2 \quad (4)$$

where  $c, i, j$  denote the index values for  $y, x$  and channel coordinates of feature maps, and  $\mu$  and  $\sigma^2$  are the mean and variance of the batch.

BN applies the following operation on each feature map in batch  $B$  for every  $c, i, j$ :

$$\text{BN}(\mathbf{x}_{cij}) = \frac{\mathbf{x}_{cij} - \mu_c}{\sqrt{\sigma_c^2 + \epsilon}} * \gamma_c + \beta_c \quad (5)$$

where  $\gamma \in \mathbb{R}^C$  and  $\beta \in \mathbb{R}^C$  are learnable parameters and  $\epsilon$  is a small constant introduced to ensure numerical stability.

At inference time, using batch statistics is a poor choice as it introduces noise in the evaluation and might not even be well defined. Therefore,  $\mu$  and  $\sigma$  are replaced by running averages of the mean and variance computed during training, which is a better approximation of the full dataset statistics.

Recent work has shown that BatchNorm has a more fundamental benefit of smoothing the optimization landscape during training (Santurkar et al., 2018) thus enhancing the predictive power of gradients as our guide to the global minimum. Furthermore, a smoother optimization landscape should additionally enable the use of a wider range of learning rates and initialization schemes which is congruent with the findings of Ioffe and Szegedy in the original BatchNorm paper (Ioffe & Szegedy, 2015).

### 4.2. Residual connections

Residual connections are another approach used in the state-of-the-art Residual Networks (He et al., 2016) to tackle the vanishing gradient problem. Introduced by He et. al. (He et al., 2016), a residual block consists of a convolution (or group of convolutions) layer, “short-circuited” with an identity mapping. More precisely, given a mapping  $F^{(b)}$  that denotes the transformation of the block  $b$  (multiple consecutive layers),  $F^{(b)}$  is applied to its input feature map  $\mathbf{x}^{(b-1)}$  as  $\mathbf{x}^{(b)} = \mathbf{x}^{(b-1)} + F(\mathbf{x}^{(b-1)})$ .

Intuitively, stacking residual blocks creates an architecture where inputs of each blocks are given two paths : passing through the convolution or skipping to the next layer. A residual network can therefore be seen as an ensemble

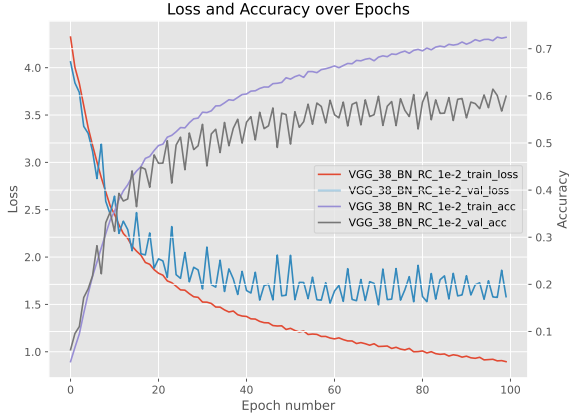


Figure 4. Training curves for VGG38 with Batch Normalisation and Residual Connection (Learning Rate 1e-2) in terms of cross-entropy error (left Y-axis) and classification accuracy (right Y-axis)

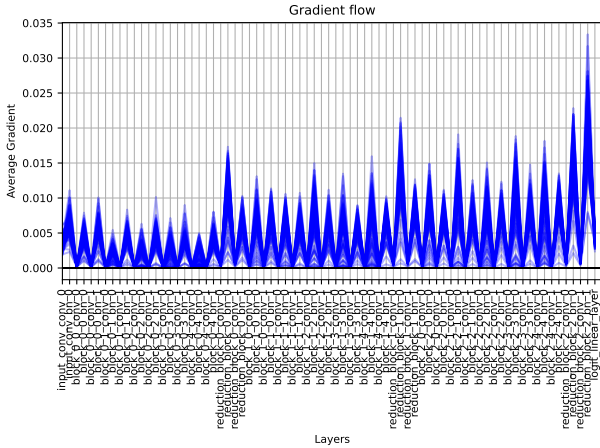


Figure 5. Gradient Flow on VGG38 with Batch Normalisation and Residual Connection (Learning Rate 1e-2)

model averaging every sub-network created by choosing one of the two paths. The skip connections allow gradients to flow easily into early layers, since

$$\frac{\partial \mathbf{x}^{(b)}}{\partial \mathbf{x}^{(b-1)}} = \mathbb{1} + \frac{\partial F(\mathbf{x}^{(b-1)})}{\partial \mathbf{x}^{(b-1)}} \quad (6)$$

where  $\mathbf{x}^{(b-1)} \in \mathbb{R}^{C \times H \times W}$  and  $\mathbb{1}$  is a  $\mathbb{R}^{C \times H \times W}$ -dimensional tensor with entries 1 where  $C$ ,  $H$  and  $W$  denote the number of feature maps, its height and width respectively. Importantly,  $\mathbb{1}$  prevents the zero gradient flow.

## 5. Experiment Setup

[ Question Figure 4 ]

[ Question Figure 5 ]

[ Question Table 1 ]

We conduct our experiment on the CIFAR100 dataset (Krizhevsky et al., 2009), which consists of 60,000 32x32

colour images from 100 different classes. The number of samples per class is balanced, and the samples are split into training, validation, and test set while maintaining balanced class proportions. In total, there are 47,500; 2,500; and 10,000 instances in the training, validation, and test set, respectively. Moreover, we apply data augmentation strategies (cropping, horizontal flipping) to improve the generalization of the model.

With the goal of understanding whether BN or skip connections help fighting vanishing gradients, we first test these methods independently, before combining them in an attempt to fully exploit the depth of the VGG38 model.

All experiments are conducted using the Adam optimizer with the default learning rate (1e-3) – unless otherwise specified, cosine annealing and a batch size of 100 for 100 epochs. Additionally, training images are augmented with random cropping and horizontal flipping. Note that we do not use data augmentation at test time. These hyperparameters along with the augmentation strategy are used to produce the results shown in Fig. 1.

When used, BN is applied after each convolutional layer, before the Leaky ReLU non-linearity. Similarly, the skip connections are applied from before the convolution layer to before the final activation function of the block as per Fig. 2 of (He et al., 2016). Note that adding residual connections between the feature maps before and after downsampling requires special treatment, as there is a dimension mismatch between them. Therefore in the coursework, we do not use residual connections in the down-sampling blocks. However, please note that batch normalization should still be implemented for these blocks.

### 5.1. Residual Connections to Downsampling Layers

[Adding residual connections to a network requires that the input and output of the block have exactly the same shape; otherwise, additive residual connections cannot be performed. However, downsampling layers typically alter the output shape, making it incompatible with direct addition of residual connections.]

We propose two solutions to address this issue. The first solution is to add a 1x1 convolution block to the residual connection to adjust the input shape, ensuring it matches the output shape. The advantage of this approach is its high flexibility, as the 1x1 convolution allows learning additional parameters to better retain input information. However, this method increases computational cost and model complexity.

The second solution is to add the same downsampling layer (e.g. average pooling) to the residual connection. This also aligns the input and output shapes, allowing residual addition. The advantage of this method is its simplicity and lack of additional computational overhead. However, the downside is that downsampling layers may lose important input information during the process.



| Model         | LR   | # Params | Train loss | Train acc | Val loss | Val acc |
|---------------|------|----------|------------|-----------|----------|---------|
| VGG08         | 1e-3 | 60 K     | 1.74       | 51.59     | 1.95     | 46.84   |
| VGG38         | 1e-3 | 336 K    | 4.61       | 00.01     | 4.61     | 00.01   |
| VGG38 BN      | 1e-3 | 339 K    | 1.70       | 52.16     | 1.84     | 48.08   |
| VGG38 RC      | 1e-3 | 336 K    | 1.33       | 61.52     | 1.84     | 52.32   |
| VGG38 BN + RC | 1e-3 | 339 K    | 1.26       | 62.99     | 1.73     | 53.76   |
| VGG38 BN      | 1e-2 | 339 K    | 1.70       | 52.28     | 1.99     | 46.72   |
| VGG38 BN + RC | 1e-2 | 339 K    | 0.90       | 72.43     | 1.58     | 59.92   |

Table 1. Experiment results (number of model parameters, Training and Validation loss and accuracy) for different combinations of VGG08, VGG38, Batch Normalisation (BN), and Residual Connections (RC), LR is learning rate.

Both methods enable residual connections in downsampling layers, but the choice depends on the specific requirements.] .

## 6. Results and Discussion

[From our experimental results, Batch Normalization (BN) and Residual Connections (RC) effectively alleviate the vanishing gradient problem and improve model performance. The specific effects can be observed in Table 1. First, by comparing VGG08 and VGG38, it is evident that deeper networks with larger capacity are more susceptible to the vanishing gradient problem. By adding BN and RC layers to VGG38, we observe a significant improvement in performance, as these methods successfully mitigate the gradient vanishing issue. Comparing the validation accuracy (val acc), we may infer that RC has a slightly greater impact than BN.

Furthermore, by varying the learning rate, we found that the BN layer is less sensitive to learning rate changes, while the RC layer is more sensitive. This implies that adjusting the learning rate can further optimize networks with RC layers. Figure 6 from (He et al., 2016) also demonstrates that BN effectively addresses gradient vanishing on its own, but the introduction of RC layers further enhances convergence on top of BN. Based on these experiments, we conclude that the best model configuration is VGG38 with BN and RC layers (learning rate of 1e-2). By analyzing its training curve (Figure 4) and gradient flow (Figure 5), we observe that this model mitigates the vanishing gradient problem well and achieves the highest validation accuracy of 59.92. However, there are signs of slight overfitting, and the validation curve exhibits greater volatility, possibly due to the model’s sensitivity to the learning rate.

To further improve model performance, I suggest the following experimental designs: 1. Add BN and RC layers to VGG08 to test whether their benefits extend to shallower networks. Given their success with VGG38, this could reveal whether these methods are universally applicable or more effective in deeper networks. 2. Incorporate appropriate regularization methods, such as Dropout or L1 norm penalties, to mitigate the overfitting observed in the current best-performing model.

To further investigate the effects of BN and RC layers, I propose the following experimental directions: 1. Evaluate whether the effects of BN and RC remain consistent across networks of varying depths, such as VGG08 and VGG38. 2. Explore the effect of learning rates on networks with only RC layers, given our observation that learning rate has minimal impact on networks with only BN layers. 3. Examine whether BN and RC independently enhance model performance or if there is an interaction between the two methods that amplifies their effectiveness.

These additional experiments will provide deeper insights into the behavior and benefits of BN and RC layers across different set up.] .

## 7. Conclusion

[In this paper, we primarily investigated the vanishing gradient problem in deep networks. We found that simply increasing the number of layers in a model can lead to the gradient failing to propagate effectively to the upper layers, thereby hindering parameter updates. To address this issue, we explored two solutions: Batch Normalization (BN) and Residual Connections (RC). Both methods were applied to the original VGG38 model, and extensive experimental results demonstrated that these methods effectively mitigate the vanishing gradient problem. Additionally, our results highlight the potential for further experiments to explore the specific effects of BN and RC layers on model performance.

Given that BN and RC layers have proven effective in addressing the vanishing gradient problem, we also observed that ResNet architectures fundamentally resolve the issue of network depth. By incorporating residual connections, ResNet enables deeper networks to maintain gradient flow. However, as noted in (He et al., 2016), the performance of very deep ResNets does not always surpass that of shallower ResNets. This suggests that network depth is no longer the primary limitation, and future research should instead focus on improving network optimization. Specifically, we recommend exploring methods to accelerate the convergence of large-capacity networks and address overfitting challenges.

In conclusion, while significant progress has been made

---

in mitigating gradient-related issues, optimizing deeper networks remains an open problem. Future work could include designing better regularization techniques, exploring adaptive learning rate strategies for large-scale models, and investigating alternative architectures that balance depth, capacity, and generalization ability.] .

## References

- Bengio, Yoshua, Frasconi, Paolo, and Simard, Patrice. The problem of learning long-term dependencies in recurrent networks. In *IEEE international conference on neural networks*, pp. 1183–1188. IEEE, 1993.
- Bishop, Christopher M et al. *Neural networks for pattern recognition*. Oxford university press, 1995.
- Glorot, Xavier and Bengio, Yoshua. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 249–256. JMLR Workshop and Conference Proceedings, 2010.
- He, Kaiming, Zhang, Xiangyu, Ren, Shaoqing, and Sun, Jian. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Huang, Gao, Liu, Zhuang, Van Der Maaten, Laurens, and Weinberger, Kilian Q. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- Ioffe, Sergey and Szegedy, Christian. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456. PMLR, 2015.
- Krizhevsky, Alex, Hinton, Geoffrey, et al. Learning multiple layers of features from tiny images. 2009.
- LeCun, Yann A, Bottou, Léon, Orr, Genevieve B, and Müller, Klaus-Robert. Efficient backprop. In *Neural networks: Tricks of the trade*, pp. 9–48. Springer, 2012.
- Rumelhart, David E, Hinton, Geoffrey E, and Williams, Ronald J. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- Santurkar, Shibani, Tsipras, Dimitris, Ilyas, Andrew, and Mądry, Aleksander. How does batch normalization help optimization? In *Proceedings of the 32nd international conference on neural information processing systems*, pp. 2488–2498, 2018.
- Simonyan, Karen and Zisserman, Andrew. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.