

Deep Learning - Project 1

Team Deep Fried Learning

Rahil Singhi¹, Alper Mumcular², Divya Srinivasan³

New York University

¹rs9174@nyu.edu ²am14533@nyu.edu ³ds7852@nyu.edu

Abstract

In this project, we explore and introduce two ResNet architectures designed for efficient and accurate image classification on the CIFAR-10 dataset. Our first model, *LightResNet18_v2*, is a lightweight ResNet-18 variant optimized to stay under 5 million parameters, incorporating techniques like Squeeze-and-Excitation (SE) blocks and Stochastic Depth. Our second model, *RobustResNet64*, is a deeper architecture designed for enhanced generalization and performance, featuring pre-activation residual blocks and dropout regularization. Both models achieved over 96% test accuracy and secured a score of 0.87709 on the Kaggle leaderboard.

You can access the GitHub repository from this [link](#).

Overview

Residual Networks (ResNets) have become a powerful tool in deep learning for their ability to train extremely deep neural networks. The key strength of ResNets lies in their use of residual connections, also known as skip connections, which allow the network to bypass one or more layers. This technique effectively mitigates the vanishing gradient problem, enabling deeper models to converge efficiently. The residual block in a ResNet is characterized by the following operation:

$$\text{ReLU}(S(x) + F(x)) \quad (1)$$

where $S(x)$ is the identity (skip) connection and $F(x)$ is a series of convolutional, batch normalization, and activation layers. By combining these paths, ResNets excel in extracting hierarchical features, which makes them particularly effective in computer vision tasks.

1. CIFAR-10 Dataset

The CIFAR-10 dataset is a widely used benchmark in computer vision, consisting of 60,000 color images in 10 classes. Each image has a spatial resolution of 32×32 . There are 50,000 training images (split across five batches, each containing 10,000 images) and 10,000 test images:

- **Training set:** 50,000 images
- **Test set:** 10,000 images

Copyright © 2024, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

- **Classes:** airplane, automobile, bird, cat, deer, dog, frog, horse, ship, truck

However, for this project, we do *not* directly use the official CIFAR-10 test set as our final evaluation set. Instead, we use it as validation. We divide these 60,000 training images into a training subset and a validation subset in a way that best suits our experiments.

Each batch file is a Python “pickled” object containing:

- **data:** A $10,000 \times 3072$ numpy array of `uint8`, each row corresponding to one image in row-major format (the first 1024 entries are the red channel, then 1024 for green, and 1024 for blue).
- **labels:** A list of 10,000 integers in the range 0-9, indicating the class label for each image.

Additionally, the file `batches.meta` contains:

- **label_names:** A 10-element list giving names to the numerical labels (e.g., `label_names[0] == "airplane"`).

2. Custom Test Dataset

A custom test dataset is provided separately. It consists of images (also 32×32) with associated ID values. This dataset serves as the *final* evaluation set. Our goal is to apply our trained ResNet model to these images and output predictions in a CSV file containing two columns:

- **ID:** The ID of the test image.
- **Labels:** The predicted label corresponding to that image.

3. Findings

Through our experiments, we observed that carefully tuning the architecture’s parameters and hyperparameters was crucial in maximizing performance. By adjusting layer configurations, regularization techniques, and learning schedules, we achieved improved accuracy on CIFAR-10 while ensuring the model remained within the constraint of 5 million parameters.

Methodology

1. Modified ResNet-18

Our model, which we call *LightResNet18_v2*, follows the standard ResNet-18 macro-architecture but:

- Uses **Squeeze-and-Excitation** (SE) blocks to adaptively recalibrate channel-wise feature responses.
- Employs **Pre-activation** residual blocks for faster convergence.
- Uses **Stochastic Depth** to randomly drop residual connections (only during training), thereby acting as a strong form of regularization.

We inflate the base channel count from 32 to 40, with subsequent doubling ($40 \rightarrow 80 \rightarrow 160 \rightarrow 320$). This yields approximately 4.4 million trainable parameters. We used a lightweight ResNet-18 variant (`LightResNet18_v2`) designed to stay under 5 million parameters. The architecture follows the standard ResNet pattern of:

conv \rightarrow BN \rightarrow ReLU + Skip Connection

stacked across multiple layers, ending in a global average pool and a fully connected head.

Label Smoothing and Loss: We use a cross-entropy loss with a label smoothing factor of 0.1. When MixUp or CutMix is applied, we compute a weighted combination of losses for the interpolated labels.

Optimizer and Scheduler:

- **SGD** with Nesterov momentum (0.9) and weight decay (5×10^{-4}) is used for stable convergence.
- **Warmup + Cosine Decay:** We start with a *linear warmup* of the learning rate for the first few epochs, then switch to a *CosineAnnealingLR* schedule for the remainder of training (300 epochs).

2. RobustResNet64 Architecture - Best Model

Our model, termed `RobustResNet64`, is a ResNet-inspired architecture tailored for CIFAR-like image datasets. The model is constructed with the following design choices:

- Utilizes pre-activation residual blocks to improve convergence speed and gradient flow during backpropagation.
- Integrates optional dropout layers within residual blocks to enhance regularization and improve generalization performance.
- The base channel count is set to 32, progressively doubling across stages: $32 \rightarrow 64 \rightarrow 128$.
- Features a total of 62 layers, calculated as follows:
 - Initial convolution layer: 1 layer.
 - Three groups of residual blocks: Each group contains 10 blocks, with each block containing 2 convolution layers, resulting in $3 \times 10 \times 2 = 60$ layers.
 - Final fully connected (FC) layer: 1 layer.
- Employs **Batch Normalization** (BN) after each convolution layer to stabilize training and accelerate convergence.
- The model incorporates a **global average pooling** layer before the fully connected head to minimize parameter count and improve robustness.

Dropout and Regularization: To improve robustness against overfitting, we apply a dropout layer within each residual block, set at a probability of 0.1. This selectively drops activations during training, enhancing model generalization.

Optimizer and Learning Rate Scheduling: The model is trained using:

- **SGD** with Nesterov momentum (0.9) and a weight decay factor of 5×10^{-4} for stable convergence.
- A combination of Warmup followed by *CosineAnnealingLR* scheduling is applied to ensure a smooth transition into effective learning rates, stabilizing training across 300 epochs.

Total Trainable Parameters: `RobustResNet64` is designed to maintain an efficient parameter count, resulting in approximately 3.8 million trainable parameters. This ensures a balance between model complexity and effective learning capacity, suitable for CIFAR-like image datasets. It becomes faster to train compared to other model while providing higher accuracy.

3. Hyper-Parameter Tuning

Between the models and tests, we actively changed the values of 5 hyper-parameters, which are *depth*, *dropout*, *epoch*, *channel size*, and *filter size*.

- **Depth of the model:** The most fundamental difference between our 2 main architectures is the depth of the model. As we increased the number of layers, the parameter count came closer to 5 million.
- **Dropout:** Our best model consists of a basic pre-activation block with dropout as an option. This allows us to tune dropout as a hyper-parameter. During testing, this helped maintain generalization at acceptable levels.
- **Epoch:** During each attempt, we tested with different epoch numbers. The optimal epoch numbers prevented both underfitting and overfitting to the given data set.
- **Channel sizes:** The optimal size for the base channel increased between the models. This became a trade-off between computation and accuracy.
- **Channel Depth Scaling:** The `RobustResNet64` architecture increases channel depth across residual groups while maintaining a fixed 3×3 kernel size for all convolutional layers (except shortcut projections). This design allowed us to experiment with feature map resolution vs. channel capacity:
 - Initial layers (*base_channels* = 32) use fewer channels to capture low-level features.
 - Deeper layers progressively double channel capacity (*base_channels* \times 2, *base_channels* \times 4) to learn higher-level abstractions after downsampling.

We found that scaling channel depth (rather than filter size) while retaining small 3×3 kernels improved feature extraction efficiency. Compared to architectures with larger filters (e.g., 5×5 or 7×7), this design reduced computational costs while achieving comparable accuracy, as hierarchical channel expansion better balances spatial and semantic learning.

4. Data Augmentation

To improve generalization, we tried an extensive data-augmentation pipeline:

- **RandomCrop and RandomHorizontalFlip:** Standard augmentations to introduce spatial variability.
- **ColorJitter and RandomErasing:** Adds further variance by randomly adjusting brightness, contrast, saturation, and hue; followed by occasional random patches in the image being erased.
- **RandAugment** (or *AutoAugment*): Introduces diverse transformations with controlled severity.
- **MixUp/CutMix Collation:** We randomly apply MixUp [1] or CutMix [2] per batch to blend images and labels. These methods encourage the model to be more robust to occlusions and interpolations of data.
- **RandomRotation:** Rotates the image by a randomly chosen degree from a specified interval.

Results

We trained our custom ResNet variant, **RobustResNet64**, on CIFAR-10 for a total of 300 epochs. In each epoch, we tracked two main metrics:

- **Training accuracy**, indicating how well the model fits the training set at that point.
- **Best validation accuracy so far**, a rolling metric that records the highest validation accuracy observed up to the current epoch (rather than just the accuracy of that epoch).

Together, these curves help illustrate the progression of generalization over time. The model leverages a range of regularization and augmentation methods (MixUp, CutMix, dropout, label smoothing, etc.), which is reflected in the gap between training and validation accuracy in the early epochs. Figure 1 plots both curves across all 300 epochs.

- **Rapid Early Improvement:** Both curves rise quickly within the first 10–15 epochs, suggesting fast initial learning of the basic features in CIFAR-10.
- **Steady Gains Thereafter:** After the initial surge, training accuracy continues climbing at a slower rate, while the best validation accuracy refines more gradually—eventually exceeding 96%.
- **Regularization Effects:** Label smoothing, MixUp, and CutMix appear to limit overfitting, as shown by the consistent gap between training and validation accuracy in the middle epochs.

Overall, this configuration—featuring pre-activation residual blocks, dropout, and advanced data augmentation—demonstrates strong generalization. The final model attains a best validation accuracy of about **96.93%**, Kaggle score of about **0.877** and contains only ≈ 3.80 million parameters, underscoring its efficiency on the CIFAR-10 benchmark.

The final configurations of our best model are:

- **Architecture:** RobustResNet64

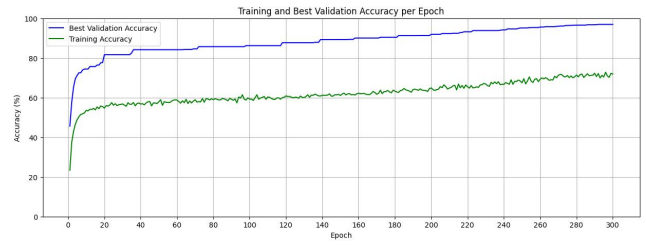


Figure 1: **Training vs. Best Validation Accuracy per Epoch.** The green curve shows per-epoch training accuracy, and the blue curve tracks the rolling best validation accuracy reached up to that epoch. Final convergence achieves roughly 96.93% validation accuracy.

Model	#Params (M)	N	B	C	Val Acc (%)
LightResNet18	2.81	4	2 (per group)	32, 64, 128, 256	93.12
LightResNet18.v2 (base=40)	4.39	4	2 (per group)	40, 80, 160, 320	94.35
LightResNet18.v2 (base=45)	5.00	4	2 (per group)	45, 90, 170, 340	94.88
RobustResNet64	3.80	3	10 (per group)	32, 64, 128	96.93

Table 1: **List of Model Configurations Experimented.** Comparison of four ResNet-based architectures in terms of parameter count (#Params), depth (N), block repetitions per group (B), and channel sizes (C). The final column reports the highest validation accuracy achieved by each model.

- **Blocks per group:** 10
- **Base channels:** 32
- **Dropout probability:** 0.1
- **Kernel Size:** 3
- **MixUp/CutMix probability:** 0.5
- **Batch Size:** 128
- **Optimizer:** SGD with Nesterov
- **Learning Rate:** 0.1
- **Momentum:** 0.9
- **Weight Decay:** 5e-4
- **Annealing Cosine:** Yes
- **Learning Rate Decay:** Cosine Annealing

By effectively combining ResNet’s strengths with CIFAR-10-specific optimizations, we achieved competitive accuracy while adhering to the project constraints.

Discussion

Our experiments confirm that a moderately deep ResNet-like architecture, when combined with strategic regularization and data augmentation, can achieve high accuracy on CIFAR-10 while staying under four million parameters. In particular, our RobustResNet64 design builds upon pre-activation residual blocks with an optional dropout layer. This adaptation proved beneficial for maintaining strong performance over longer training schedules (300 epochs), as it helps mitigate overfitting without severely hindering feature learning.

Another key factor in our training pipeline is the use of **MixUp** and **CutMix** via a custom collate function. By

blending or cutting-and-pasting images (and their corresponding labels) at each batch, we artificially enlarge the data distribution and encourage the model to become invariant to partial occlusions. In conjunction with **label smoothing**, these augmentation strategies reduce the tendency of the network to memorize training examples, as seen in the persistent gap between training accuracy and best validation accuracy. This gap indicates regularization is effectively preventing overfitting, while still allowing the model to push validation accuracy above 96%.

To further enhance generalization, we used **RandAugment**, which randomly selects transformation operations (e.g., rotations, translations, shearing). Although CIFAR-10 images are only 32×32 pixels, carefully tuned augmentations help broaden the effective training distribution. The resulting model is more robust and better able to handle variations in the test set.

During **inference**, a two-pass **test-time augmentation** (TTA) strategy is employed: the original and flipped versions of each image are passed through the network, then the two softmax outputs are averaged. Though relatively simple, TTA consistently yields small accuracy gains, pushing us toward the final result of around 96.93% on the validation set.

Finally, the model’s compact parameter footprint ($\approx 3.80\text{M}$) offers practical advantages for resource-limited scenarios, such as deployment on edge devices. Future work could explore additional strategies for compressing the network, such as pruning or quantization, to further reduce computational costs without sacrificing accuracy. Overall, these results underscore that with judicious selection of augmentation, dropout, and label smoothing, an appropriately sized residual network can deliver competitive performance on CIFAR-10 while remaining within a modest parameter budget.

Conclusion

In this project, we presented a `RobustResNet64` model trained on CIFAR-10, carefully balancing architectural depth and parameter constraints to remain under four million parameters while achieving a best validation accuracy of approximately 96.93%. Our approach combines standard ResNet concepts—pre-activation residual blocks and dropout—with advanced data augmentation methods such as MixUp, CutMix, and RandAugment, along with label smoothing. This synergy of architectural choices and regularization strategies proved crucial in preventing overfitting and steadily improving generalization performance over 300 epochs.

Moreover, inference-time horizontal flipping (test-time augmentation) added a slight but consistent accuracy gain. Our results show that a smaller, well-regularized convolutional network can compete with larger, more complex models on benchmarks like CIFAR-10. In future work, we plan to investigate whether these techniques generalize to higher-resolution tasks (e.g., ImageNet) and to examine how additional compression strategies (e.g., pruning, quantization) might further reduce the parameter and computational footprints without compromising accuracy.

Note: This report was generated with some assistance from GPT and DeepSeek.

References

- [1] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018.
- [2] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*, 2019.