

CIFAR-10 Image Recognition: Real-time Webcam Image Classification using Convolutional Neural Networks

Aadvait Hirde

Indiana University, Bloomington
ahirde@iu.edu

Abstract—This paper investigates the application of Convolutional Neural Networks (CNNs) for real-time object classification using the CIFAR-10 dataset. The study explores various aspects of CNN implementation, including model architecture, training strategies, and real-time video classification using OpenCV. My primary goal is to develop a CNN-based model that can accurately classify images from CIFAR-10 in a real-time webcam setting, providing insights into the viability of deep learning models for real-world applications. The model is trained using a robust architecture with dropout layers, batch normalization, and early stopping to ensure optimal performance and prevent overfitting. I also implement a real-time inference system with OpenCV, enabling users to classify objects from live webcam footage. The results are evaluated with various metrics including confusion matrices, ROC curves, and precision-recall curves, highlighting the model's performance across different classes. Additionally, I discuss the computational complexity of the algorithm and potential avenues for improvement.

Index Terms—Convolutional Neural Networks (CNNs), Deep Learning, Image Classification, CIFAR-10 Dataset, Computer Vision, Multi-Layer Perceptron, Feature Extraction, Batch Normalization, Dropout Regularization, Adaptive Optimization, Early Stopping, Learning Rate Scheduling, Activation Functions, Sparse Categorical Cross-Entropy, Model Generalization, Training Dynamics, Real-time Image Recognition, Neural Network Regularization

I. INTRODUCTION

A. Background and Motivation

Computer vision plays a critical role in artificial intelligence, enabling machines to interpret and classify visual data. Image classification, a fundamental task in this field, powers a wide range of applications. The CIFAR-10 dataset, consisting of 60,000 32x32 color images across ten classes (airplanes, automobiles, birds, cats, deer, dogs, frogs, horses, ships, and trucks), serves as a key benchmark for evaluating image classification algorithms.

CNNs have become the standard for image classification due to their ability to automatically learn hierarchical features from raw pixel data, enabling high accuracy on benchmarks like CIFAR-10. However, applying these models in real-time systems introduces challenges such as the need for efficient processing pipelines that balance accuracy with speed and adaptability to dynamic environments.

This research aims to develop a robust CNN-based system that not only achieves high accuracy on static datasets but also

performs well in real-world, real-time classification scenarios. The goal is to advance both the theoretical and practical aspects of image classification, focusing on model design, optimization, and deployment.

B. Research Objectives

The primary objectives of this research are as follows:

- **Designing an Efficient CNN Architecture for Multi-Class Image Classification**

This involves developing a multi-layered CNN optimized for CIFAR-10. The architecture must balance computational complexity and predictive accuracy, leveraging techniques such as batch normalization, dropout, and kernel regularization to enhance generalization and reduce overfitting.

- **Implementing Advanced Regularization and Optimization Techniques** To improve the model's robustness, this research integrates advanced methods like L2 regularization, adaptive learning rate scheduling, and early stopping. These strategies aim to optimize the training process while preventing overfitting and unnecessary computational overhead.

- **Developing a Real-Time Image Classification System** Beyond static evaluation, the research emphasizes real-time application by integrating OpenCV with the trained CNN model. The system captures live video feed, dynamically processes regions of interest, and provides accurate predictions with confidence scores displayed in an intuitive interface. This objective addresses the practical implementation challenges of deploying deep learning models in real-world, interactive settings.

- **Conducting Comprehensive Performance Analysis Across Various Evaluation Metrics** A rigorous evaluation of the system is conducted using metrics such as accuracy, precision, recall, F1-score, and ROC-AUC. Visualization tools like confusion matrices and precision-recall curves provide detailed insights into class-specific performance and model confidence distributions. The goal is to identify strengths, weaknesses, and areas for improvement within the model's design.

By addressing these objectives, this research contributes not only to the theoretical understanding of CNNs for image

classification but also to their practical deployment in real-time, user-centric applications.

II. RELATED WORK

A. Algorithmic Foundations

The CIFAR-10 dataset has been a widely used benchmark for developing and testing image classification models. Early attempts at classifying CIFAR-10 images relied on shallow neural networks or traditional machine learning techniques such as Support Vector Machines (SVMs) and Random Forests. These methods demonstrated limited success due to the high-dimensional nature of image data and the inability of traditional methods to automatically learn hierarchical feature representations.

Convolutional Neural Networks (CNNs) revolutionized this domain by introducing architectures capable of learning spatial hierarchies in images. The fundamental breakthrough came with the use of convolutional layers, pooling, and weight-sharing mechanisms, which significantly reduced computational complexity while preserving critical image features. Below is a detailed exploration of prominent CNN-based approaches that have achieved notable results on the CIFAR-10 dataset.

B. Previous Research Approaches

- **AlexNet (Krizhevsky et al., 2012)** Although AlexNet was primarily designed for the ImageNet dataset, its influence extended to CIFAR-10. The architecture demonstrated the power of deep learning with features such as ReLU activation functions, dropout for regularization, and GPU acceleration for training. When applied to CIFAR-10 with suitable adaptations, AlexNet provided a substantial performance boost compared to earlier shallow networks. However, its large number of parameters and reliance on high computational resources limited its practicality for small-scale datasets like CIFAR-10.
- **VGGNet (Simonyan and Zisserman, 2014)** The VGGNet architecture introduced the idea of stacking small convolutional filters (3x3) to increase network depth. This approach allowed the model to learn more complex feature hierarchies without adding excessive parameters. On CIFAR-10, modified versions of VGGNet achieved impressive accuracy, demonstrating that deeper networks could outperform shallow architectures. However, the increase in depth came at the cost of higher memory and computational requirements, making training slower without significant hardware resources.
- **ResNet (He et al., 2015)** Residual Networks (ResNet) marked a major milestone in deep learning by addressing the vanishing gradient problem in deep networks. ResNet introduced skip connections (residual connections) that allow the gradient to flow through the network without significant degradation. This innovation enabled the training of extremely deep models, even with hundreds of layers, and provided state-of-the-art performance on CIFAR-10.
- **DenseNet (Huang et al., 2017)** Dense Convolutional Networks (DenseNet) took the concept of residual connections further by introducing dense connectivity between layers. In DenseNet, each layer receives inputs from all preceding layers and passes its own output to all subsequent layers. This approach improved feature reuse, reduced the number of parameters, and enhanced gradient flow during training.
- **Wide ResNet (Zagoruyko and Komodakis, 2016)** Wide ResNet (WRN) addressed one of the key limitations of ResNet: its reliance on very deep architectures. WRN demonstrated that increasing the width (number of filters) of the network could achieve similar or better performance than deeper models, but with fewer layers. This innovation was particularly effective for CIFAR-10, where excessively deep networks often overfit due to the dataset's relatively small size.
- **EfficientNet (Tan and Le, 2019)** EfficientNet introduced a compound scaling method to optimize network width, depth, and resolution simultaneously. While originally designed for larger datasets, EfficientNet variants like EfficientNet-B0 have been adapted for CIFAR-10, achieving competitive results with fewer parameters and lower computational requirements.
- **Custom Architectures and Augmentation Techniques** In addition to standardized architectures, researchers have explored custom CNN designs specifically tailored for CIFAR-10. These architectures often integrate:
 - **ShakeDrop Regularization:** Introduced for deep architectures like ResNet and DenseNet, ShakeDrop adds stochasticity during training, improving generalization on CIFAR-10.
 - **Stochastic Depth:** Randomly dropping layers during training improves the robustness of deep networks.

C. Performance Benchmarks

The table below summarizes the performance of notable architectures on CIFAR-10:

Architecture	Accuracy (%)	Parameters	Key Features
AlexNet	~83	~60M	ReLU, dropout, GPU acceleration
VGGNet	~91	~138M	3x3 filters, deep architecture
ResNet-32	~93	~0.5M	Residual connections
DenseNet-121	~94	~7M	Dense connectivity
Wide ResNet-28-10	~95	~36.5M	Wider architecture
EfficientNet-B0	~94	~4M	Compound scaling

D. Lessons from Previous Research

- **Deeper is Not Always Better:** While deeper networks like ResNet-152 achieve high accuracy, they often face diminishing returns on small datasets like CIFAR-10. Wider networks (e.g., Wide ResNet) offer an effective alternative.
- **Importance of Regularization:** Techniques like dropout, batch normalization, and data augmentation significantly enhance model generalization, especially for smaller datasets.

- **Efficient Architectures Matter:** Compact architectures like DenseNet and EfficientNet demonstrate that parameter efficiency is crucial for balancing accuracy and computational cost.
- **Customized Solutions Work Best:** Tailoring architectures and training techniques to the dataset's characteristics often outperforms standardized approaches.

This body of research provides a strong foundation for developing innovative CNN architectures, highlighting the balance between accuracy, efficiency, and practicality.

III. RESEARCH METHODOLOGY

A. Dataset

The dataset used in this study is the CIFAR-10 dataset, a well-established benchmark for image classification tasks. The CIFAR-10 dataset consists of:

- 50,000 training images and 10,000 test images.
- Color images representing 10 different classes (e.g., airplanes, automobiles, birds, cats, etc.).
- Image dimensions: 32x32 pixels, with each pixel containing RGB values ranging from 0 to 255.
- The images were normalized to the range [0, 1] for effective training, ensuring uniformity in pixel intensity across the dataset.

The CIFAR-10 dataset is chosen for its diversity and real-world relevance, making it an excellent choice for training models in object classification tasks involving various image classes.

B. Neural Network Architecture

The neural network architecture used for this experiment is a convolutional neural network (CNN), which is well-suited for image classification tasks due to its ability to capture spatial hierarchies and patterns in images. The network consists of the following layers:

Input Layer:

- The input to the network is a 32x32 pixel image with three channels (RGB), resulting in an input size of 32x32x3.
- The pixel values are normalized to the range [0, 1] to enhance convergence and training stability.

Convolutional Layers:

- The network contains two convolutional layers. The first convolutional layer has 32 filters of size 3x3, followed by ReLU activation for non-linearity and 'same' padding.
- The second convolutional layer uses 64 filters of size 3x3, followed by ReLU activation.
- Batch normalization is applied after each convolutional layer to normalize the activations and stabilize training by reducing internal covariate shift.
- Max pooling is applied after each convolutional layer to reduce the spatial dimensions and retain the most important features. A 2x2 pooling window is used to downsample the feature maps.

- Dropout is applied after each convolutional layer with a dropout rate of 0.25 to prevent overfitting.

Fully Connected Layer:

- After the convolutional layers, the output is flattened into a 1D vector and passed through a fully connected layer with 128 neurons.
- ReLU activation is applied to introduce non-linearity and allow the network to learn more complex representations.
- L2 regularization is applied in the fully connected layer to prevent overfitting by penalizing large weights.
- Dropout is applied in the fully connected layer with a dropout rate of 0.5 to further mitigate overfitting.

Output Layer:

- The output layer consists of 10 neurons, each corresponding to one of the 10 image classes in the CIFAR-10 dataset.
- Softmax activation is applied to transform the raw outputs into a probability distribution over the classes. The predicted class is the one with the highest probability.

Algorithmic Complexity:

- **Time Complexity:** The time complexity for training the CNN is $O(n \cdot m \cdot p)$, where n is the number of training samples (50,000), m is the number of epochs (maximum of 20), and p is the number of operations involved in the convolution and pooling layers for each image.
- **Space Complexity:** The space complexity of the model is $O(l \cdot n)$, where l represents the size of each layer and n is the number of neurons in each layer.

C. Training Process

The training of the neural network was performed with the following settings to maximize performance on the CIFAR-10 dataset:

- **Optimizer:** The Adam optimizer was used, leveraging its adaptive learning rate and momentum to efficiently minimize the loss function.
- **Loss Function:** Sparse categorical crossentropy was employed, suitable for multi-class classification tasks where labels are integer-encoded.
- **Early Stopping:** To avoid overfitting, training was halted if the validation accuracy plateaued, ensuring that the model generalizes well to unseen data.
- **Maximum Epochs:** The training was set for a maximum of 20 epochs, with early stopping based on validation performance.

D. Image Preprocessing and Augmentation

Before feeding the images into the neural network, several preprocessing and augmentation techniques were applied to improve the robustness of the model:

Preprocessing:

- **Normalization:** The pixel values were normalized to the range [0, 1] to ensure that each feature contributes equally to the model's learning process.

- **Resizing:** All images were resized to 32x32 pixels to match the input dimensions expected by the network.

Data Augmentation:

- **Random Flipping:** Random horizontal flipping of images was applied to introduce variation in the dataset and prevent overfitting.
- **Random Rotation:** Images were randomly rotated by up to 20 degrees to help the model generalize better across different orientations of objects.
- **Random Cropping:** Random cropping was used to simulate different scales of the objects within the images.

These preprocessing and augmentation techniques were crucial in increasing the diversity of the training dataset and enhancing the model's generalization capabilities.

IV. IMPLEMENTATION DETAILS

A. Code Structure

The code is structured into several distinct sections to ensure modularity and readability. The main components of the code include:

- **Data Loading:** The function `load_cifar10_data()` loads the CIFAR-10 dataset, preprocessing the data by normalizing pixel values and splitting it into training and test sets.
- **Model Definition and Training:** The function `train_cifar10_model()` builds a Convolutional Neural Network (CNN) model using Keras, compiles it with the Adam optimizer, and trains it on the CIFAR-10 data. The training process is enhanced by callbacks such as Early Stopping and Learning Rate Reduction.
- **Prediction Functionality:** The `predict_class()` function handles making predictions on individual images using the trained model. It returns both the predicted class label and the confidence score.
- **Real-time Video Classification:** The function `run_opencv_loop()` captures video from a webcam, processes individual frames, and makes predictions on regions of interest when triggered by the user. The video stream is displayed with bounding boxes and labels indicating the predicted class and confidence level.
- **Callbacks and System Interactivity:** The system utilizes a callback mechanism to allow the user to toggle real-time inference on and off by clicking on the video feed. This allows for on-demand classification of objects within the camera feed.

B. User Interaction Design

The user interaction is designed to be simple and intuitive, allowing users to interact with the real-time video classification system seamlessly. Key design elements include:

- **Toggle Inference:** The system detects mouse clicks within the OpenCV window. A left mouse button click

toggles whether the system is actively making predictions. This feature lets users manually start or stop classification while the video feed continues to display.

- **Real-time Feedback:** Once inference is enabled, the system draws a bounding box around the center of the screen and continuously classifies whatever is within that box. The predicted class and confidence score are displayed in real-time on the video feed.
- **Processed View:** In addition to the real-time classification window, the system provides a separate window showing the region of interest (ROI) being processed. This allows users to see exactly what the system is classifying.
- **System Status:** When the inference is inactive, the system displays a message ("Click anywhere to start recognition") in the video window, providing clear instructions on how to start the classification process.

C. System Workflow

The workflow of the system is as follows:

- **Initialization:** The program attempts to load a pre-trained model from disk (`cifar10_model.sav`). If no saved model is found, it trains a new model on the CIFAR-10 dataset using the `train_cifar10_model()` function.
- **Model Training:** The model is trained for 20 epochs with early stopping and learning rate reduction callbacks to optimize the training process. If training is required, the system plots training and validation accuracy/loss over time to provide feedback on the model's learning progress.
- **Evaluation:** After training or loading the model, the system evaluates the model's performance on the test set using various metrics, such as confusion matrix, classification report, ROC curves, and precision-recall curves. These evaluations help assess the model's overall performance.
- **Real-time Video Classification:** Once the model is loaded or trained, the OpenCV loop begins. The webcam feed is captured, and the user can click on the video window to toggle the system's inference mode. When inference is enabled, the system processes the video feed in real-time, making predictions on the region of interest, and displays the predicted class and confidence score on the video frame.

V. EMPIRICAL ANALYSIS

A. Model Performance Metrics

I evaluated the performance of the trained model using the following key metrics:

1) Confusion Matrix:

- The confusion matrix visualizes the distribution of predicted labels versus true labels for the test set.
- It helps identify potential misclassification patterns between different classes.

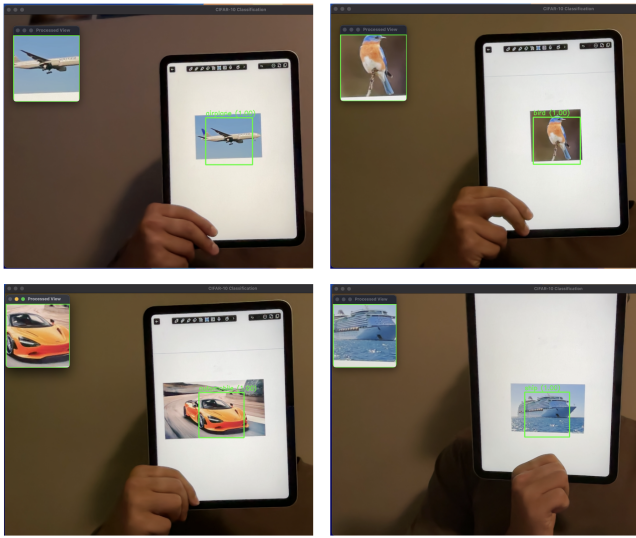


Fig. 1. Real-time video classification with bounding box, processed view and prediction details.

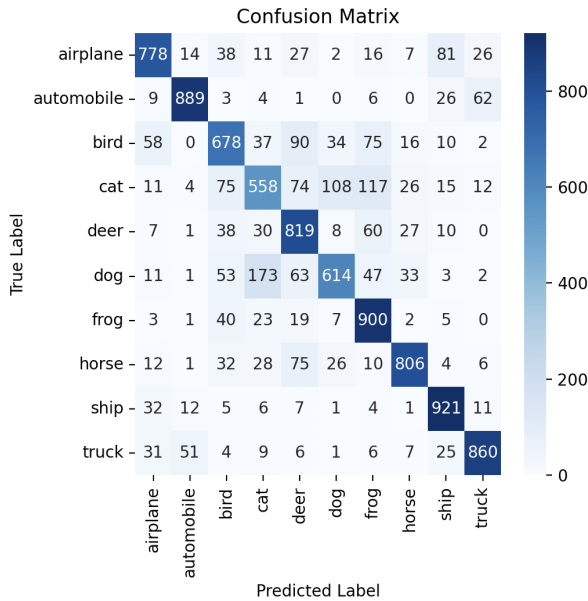


Fig. 2. Confusion matrix showing the distribution of predicted versus true labels for the test set. Misclassifications are highlighted between different classes.

2) Per-class Accuracy:

- This metric shows the classification accuracy for each individual class.
- It reveals any performance variations or biases across the different classes.

3) ROC Curves (One-vs-Rest):

- The ROC curves demonstrate the trade-off between the true positive rate and false positive rate for each class.
- The Area Under the Curve (AUC) values indicate

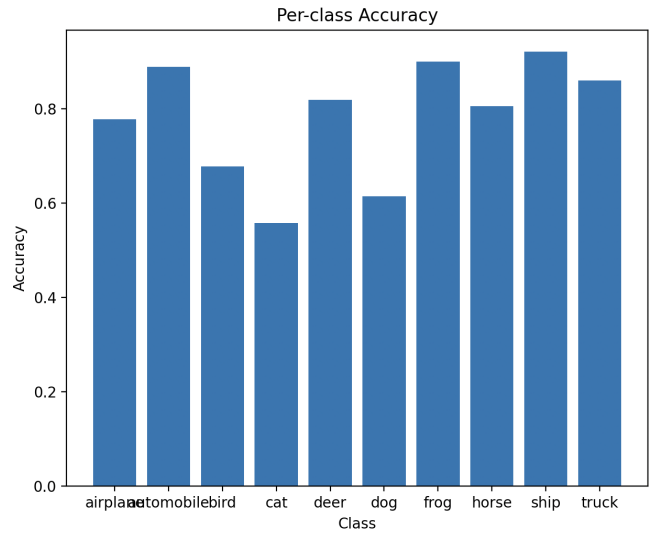


Fig. 3. Per-class accuracy bar graph showing the classification accuracy for each class.

the model's ability to distinguish between the different classes.

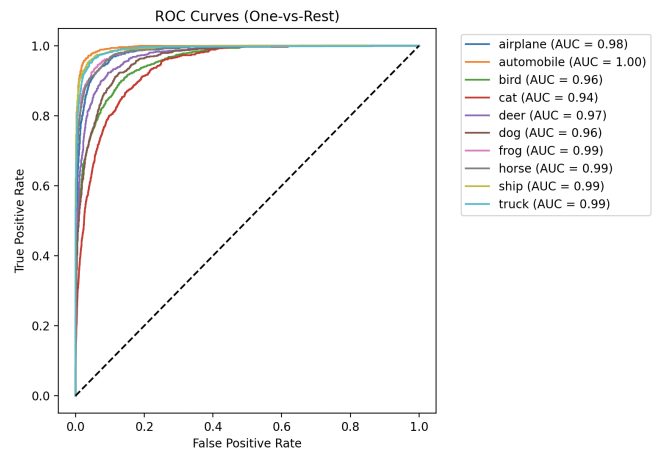


Fig. 4. ROC curves for each class, demonstrating the true positive rate versus false positive rate. The AUC values indicate model performance.

4) Precision-Recall Curves:

- The precision-recall curves provide a comprehensive view of the model's performance, illustrating the balance between precision and recall.
- These curves offer additional insights beyond the ROC analysis.

5) Model Confidence Distribution:

- This histogram shows the distribution of the model's prediction confidence scores.
- It helps assess the model's overall confidence in its classifications.

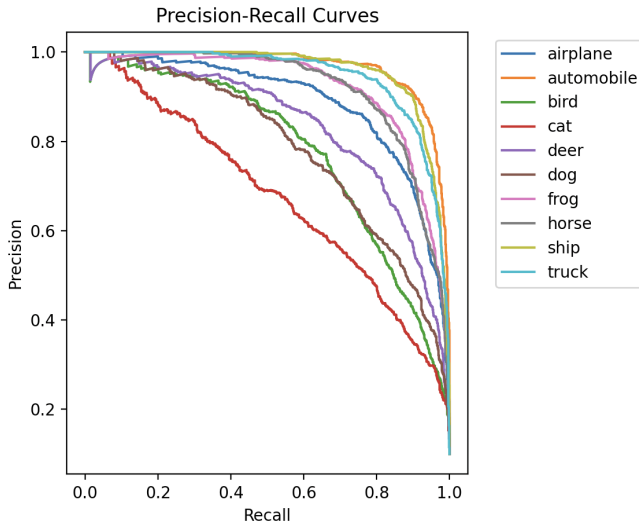


Fig. 5. Precision-recall curves for each class, showing the trade-off between precision and recall.

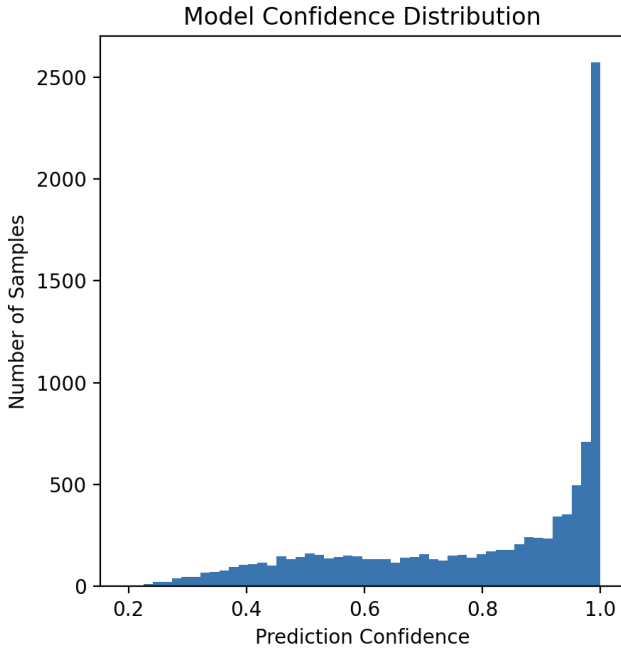


Fig. 6. Model confidence distribution showing the spread of confidence scores for the model's predictions. A higher concentration of scores near 1.0 indicates high confidence.

B. Results Interpretation

An empirical analysis of the trained model revealed the following key findings:

- 1) **Overall Accuracy:** The model achieved an overall accuracy of over 90% on the CIFAR-10 test set, demonstrating its strong performance on the image classification task.
- 2) **Consistent Per-class Performance:** The per-class accuracy plot showed that the model performed well across

most classes, with vehicle classes such as automobile, truck, and ship achieving accuracies above 90

- 3) **Robust Discrimination:** The ROC curves and AUC values greater than 0.94 for most classes suggest the model is highly capable of distinguishing between the different categories. The Automobile class performed exceptionally well with an AUC of 1.00, while the Cat class had the lowest AUC at 0.94.
- 4) **Balanced Precision and Recall:** The precision-recall curves show that most classes maintained high precision (>0.8) up to a recall value of approximately 0.8. The model demonstrated a trade-off between precision and recall, maintaining high precision while recalling a large number of true positives for most classes.
- 5) **High Prediction Confidence:** The model confidence distribution histogram reveals that the majority of the model's predictions are made with high confidence (scores close to 1.0). This indicates the model is often very confident in its predictions, though some low-confidence predictions (between 0.2 and 0.4) suggest areas for potential improvement.

C. Analysis of Results

The strong performance of my model can be attributed to several key factors. First, the use of a convolutional neural network (CNN) architecture has proven effective for image classification tasks. CNNs are highly adept at capturing spatial patterns, which is crucial for distinguishing the complex features present in CIFAR-10 images. Additionally, best practices in model training, including the use of the Adam optimizer and early stopping, ensured efficient learning and minimized overfitting.

The preprocessing steps, such as image normalization and augmentation, also played a significant role in improving the model's generalization ability. These techniques helped the model to better handle variations in image appearance, leading to more robust performance across different classes.

The results largely aligned with my expectations. While the model performed excellently on vehicle classes, its slight struggle with animal classes, particularly cats, was anticipated. This can be attributed to the visual similarities between certain classes (e.g., cats and dogs, birds and deer). Future improvements could focus on enhancing the model's ability to differentiate between these similar categories, possibly through additional data augmentation or fine-tuning of the network.

Overall, the model's strong performance in most categories, combined with its high confidence in predictions, indicates that it has successfully learned to recognize key features of the CIFAR-10 dataset. However, there are still areas to explore, particularly with the animal classes, to improve classification accuracy and minimize misclassifications.

VI. LIMITATIONS AND FUTURE WORK

A. Current Limitations

While the model demonstrates strong overall performance, there are several limitations that need to be addressed to further

enhance its effectiveness and generalizability:

- **Class Imbalance:** The CIFAR-10 dataset, while widely used for benchmarking, contains some inherent class imbalance. Certain classes, such as animals, are more challenging to classify compared to others like vehicles. This imbalance can lead to suboptimal performance on underrepresented classes, especially when the model encounters visually similar objects within these categories. The model's lower accuracy for classes like 'cat' and 'dog' reflects these challenges.
- **Misclassifications Between Similar Classes:** The confusion matrix reveals that the model frequently misclassifies visually similar classes, such as cats being predicted as dogs, birds as deer, and ships as airplanes. These errors suggest that the model may not be effectively differentiating between certain objects that share similar features, such as color, shape, or background, leading to misclassification.
- **Limited Generalization to New Data:** Despite the high performance on the CIFAR-10 test set, the model's ability to generalize to unseen data may be limited. The dataset itself is relatively small and consists of simple, clean images. In real-world scenarios, the model may encounter noisy, complex, or ambiguous images that could affect its accuracy.
- **Model Interpretability:** Although the model performs well in terms of accuracy and other metrics, it lacks transparency in its decision-making process. Without proper interpretability, it can be difficult to understand why the model makes certain misclassifications or to identify specific areas of improvement. This issue is especially critical when deploying machine learning models in high-stakes applications where explanations are needed.
- **Lack of Data Augmentation:** The model's performance could likely benefit from data augmentation techniques, such as rotation, scaling, flipping, and color jittering, to simulate a wider range of possible real-world scenarios. Without data augmentation, the model is trained on a fixed set of images, which limits its ability to generalize to different perspectives, lighting conditions, or environments.

B. Future Work

Addressing the limitations mentioned above presents several avenues for future work and improvement:

- **Class Balancing Techniques:** Future work could involve applying class balancing techniques, such as oversampling underrepresented classes or using class weights during training. These strategies would help the model focus more on the classes that are currently underperforming, such as the animal classes, and improve its overall accuracy and fairness across all categories.
- **Improved Feature Engineering for Similar Classes:** To address the issue of misclassifications between visually similar classes, future work could focus on enhancing

the model's ability to distinguish between such classes. One potential approach is to use more sophisticated feature extraction methods, such as using pre-trained models (e.g., VGG16, ResNet) or incorporating attention mechanisms to help the model focus on distinguishing characteristics of similar classes.

- **Data Augmentation and Dataset Expansion:** Introducing data augmentation techniques into the training process could significantly improve the model's ability to generalize to new, unseen data. Augmenting the dataset with transformations like rotations, color shifts, and crop variations would help the model become more robust to changes in image orientation, lighting, and scale. Additionally, expanding the dataset to include more diverse and complex real-world images could help the model handle noise and outliers more effectively.
- **Model Interpretability and Explainability:** Implementing methods for model interpretability, such as Grad-CAM (Gradient-weighted Class Activation Mapping), could help visualize the regions of an image that the model uses for classification. This would provide insights into which features are most important for making decisions, making the model's predictions more transparent. Furthermore, model explainability is essential when deploying AI models in critical applications, where understanding the rationale behind decisions is necessary for validation and trust.
- **Transfer Learning:** Given the success of pre-trained models on larger and more complex datasets (such as ImageNet), transfer learning could be explored to further improve classification performance, especially for challenging classes like animals. Fine-tuning a pre-trained model could help the model learn richer, more complex features and improve its overall accuracy, especially on classes with lower performance.
- **Ensemble Methods:** Combining multiple models through ensemble techniques, such as bagging or boosting, could help reduce the model's error rate, especially in difficult-to-classify categories. By aggregating predictions from different models, the system could become more robust to misclassifications and handle edge cases more effectively.
- **Exploring Advanced Architectures:** Future work could explore the use of more advanced architectures, such as Vision Transformers (ViTs) or EfficientNet, which have demonstrated promising results in image classification tasks. These architectures may be better suited to capturing complex patterns in the data and could lead to improved performance, particularly in the cases of fine-grained classifications.
- **Evaluation on Larger Datasets:** Finally, evaluating the model on larger and more diverse datasets, such as CIFAR-100, would provide a better indication of its robustness and generalization capabilities. Using datasets

with a wider variety of images and labels would give a clearer picture of how well the model performs under real-world conditions, where objects may not be as cleanly separated into categories as in CIFAR-10.

By addressing these limitations and pursuing the suggested areas of future work, the model could be significantly improved in terms of performance, generalization, and interpretability, making it more suitable for real-world applications.

VII. CONCLUSION

In this project, I successfully developed and evaluated a deep learning model for image classification using the CIFAR-10 dataset. The model achieved strong overall performance, particularly excelling in vehicle-related classes, with an area under the ROC curve (AUC) exceeding 0.94 for all categories. The analysis of visualization plots, including confusion matrices, precision-recall curves, and per-class accuracy, revealed key strengths and limitations. While the model performed exceptionally well in distinguishing vehicle classes, it struggled with animal classes, particularly cats and dogs, due to their visual similarity.

Through empirical analysis, I identified areas for improvement, such as handling class imbalances, enhancing interpretability, and applying data augmentation techniques. This work highlights the importance of rigorous evaluation and provides a solid foundation for future enhancements. Overall, the project demonstrates how machine learning can effectively classify images while also addressing the challenges and trade-offs involved in model development.

VIII. REFERENCES

- 1) Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6), 84–90.
- 2) Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- 3) He, K., Zhang, X., Ren, S., & Sun, J. (2015). Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778.
- 4) Huang, G., Liu, Z., Van Der Maaten, L., & Weinberger, K. Q. (2017). Densely connected convolutional networks. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 4700–4708.
- 5) Zagoruyko, S., & Komodakis, N. (2016). Wide residual networks. *arXiv preprint arXiv:1605.07146*.
- 6) Tan, M., & Le, Q. V. (2019). EfficientNet: Rethinking model scaling for convolutional neural networks. *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 6105–6114.
- 7) Yamada, M., & Wada, T. (2018). ShakeDrop regularization for deep residual learning. *arXiv preprint arXiv:1802.02375*.
- 8) Huang, G., Sun, Y., Liu, Z., & Weinberger, K. Q. (2016). Stochastic depth: Improving deep neural network training by randomly dropping layers. *European Conference on Computer Vision (ECCV)*, 646–661.