

# MNIST Digit Recognition: Real-time Handwritten Digit Classification using Convolutional Neural Networks

Aadvait Hirde

*Indiana University, Bloomington*

ahirde@iu.edu

Darragh Buckley

*Indiana University, Bloomington*

buckled@iu.edu

Jared Gengnagel

*Indiana University, Bloomington*

jgengnag@iu.edu

Santiago Diaz

*Indiana University, Bloomington*

diazsd@iu.edu

**Abstract**—This research explores the development of a real-time handwritten digit recognition system using deep learning techniques, specifically focusing on the MNIST dataset. By implementing a neural network with TensorFlow and integrating OpenCV for live video processing, we demonstrate an end-to-end machine learning solution capable of instantaneously classifying handwritten digits with high accuracy. Our approach not only achieves over 99% classification accuracy but also provides a user-interactive interface for real-time digit recognition, bridging theoretical machine learning concepts with practical, interactive applications.

**Index Terms**—Handwritten Digit Recognition, Real-time Classification, Convolutional Neural Networks, MNIST Dataset, Deep Learning, TensorFlow, OpenCV, Machine Learning, Interactive User Interface, Computer Vision, Image Preprocessing, Neural Network Architecture, Pattern Recognition, Digit Classification, Live Video Processing, Model Performance, Empirical Analysis, Data Augmentation, Transfer Learning.

## I. INTRODUCTION

Handwritten digit recognition is a key challenge in both computer vision and machine learning, providing significant practical applications while also serving as a rigorous testing ground for algorithm development. The MNIST (Modified National Institute of Standards and Technology) dataset has long been a standard for benchmarking machine learning models, particularly neural networks, due to its balanced complexity and accessibility.

This research focuses on advancing both the theoretical and practical aspects of digit recognition. Specifically, our goals are as follows:

- **Develop a high-performance neural network for digit classification:** Our objective is to design a robust model that not only achieves high accuracy on the MNIST dataset but also generalizes well to real-world variations in handwriting and input quality.
- **Create an interactive, real-time digit recognition system:** In addition to static testing, we aim to develop a live system capable of interacting with users in real-time. This system could enable practical applications

such as automated form processing, accessibility tools, or educational platforms.

- **Provide a comprehensive empirical analysis of model performance:** A thorough evaluation will include performance metrics such as accuracy, precision-recall, and confusion matrices. We will also explore how the system performs under varying conditions, including changes in lighting, handwriting styles, and noise levels.
- **Explore the intersection of deep learning, computer vision, and interactive user interfaces:** By integrating deep learning and computer vision with user-friendly interfaces, we aim to enhance real-time classification and create interactive environments where users can adjust system parameters and view dynamic results.

Through these objectives, we aim to advance handwritten digit recognition systems by providing a robust, versatile, and user-oriented solution.

## II. RELATED WORK

### A. Algorithmic Foundations

The foundational work in handwritten digit recognition can be traced back to several key research efforts:

- **Lecun et al. (1998):** Pioneered the use of convolutional neural networks (CNNs) for digit recognition, introducing the original MNIST dataset. The authors demonstrated that CNNs, with their ability to learn local features and leverage spatial relationships, outperform traditional machine learning algorithms like fully connected neural networks and support vector machines on the MNIST task.
- **Simard et al. (2003):** Enhanced digit recognition techniques through advanced preprocessing and regularization methods. The researchers explored the impact of elastic distortions, a form of data augmentation, on improving the generalization capabilities of neural networks trained on MNIST. They showed that this simple yet effective technique could significantly boost model performance.

- **Wan et al. (2013):** Introduced dropout as a regularization technique to improve the generalization of neural networks. Dropout helps to prevent overfitting by randomly deactivating a subset of neurons during training, forcing the model to learn more robust features. The authors demonstrated the effectiveness of dropout on MNIST, achieving state-of-the-art results at the time.

### B. Previous Research Approaches

Over the years, researchers have employed various techniques for MNIST digit recognition, including:

- **Traditional machine learning algorithms:** Early approaches utilized classic algorithms such as support vector machines (SVMs) and random forests, which were able to achieve high accuracy on the MNIST dataset.
- **Shallow neural networks:** Before the advent of deep learning, researchers explored the use of simple neural network architectures with one or two hidden layers. While these models performed reasonably well, they were limited in their ability to capture complex patterns in the data.
- **Deep convolutional neural networks:** With the rise of deep learning, researchers have focused on designing more sophisticated CNN architectures that can effectively extract hierarchical features from the MNIST images. These models have consistently achieved state-of-the-art performance on the benchmark.
- **Ensemble methods:** Some researchers have investigated combining multiple classifiers, such as CNNs and SVMs, to create ensemble models that leverage the strengths of different approaches and further improve classification accuracy.

Our research builds upon these foundations while introducing a novel real-time, interactive implementation that integrates deep learning and computer vision techniques.

## III. RESEARCH METHODOLOGY

### A. Dataset

We utilized the MNIST dataset, a widely-used benchmark for handwritten digit recognition. The dataset comprises:

- 60,000 training images and 10,000 test images.
- Grayscale images of handwritten digits (0-9).
- Image dimensions: 28x28 pixels, with each pixel having a value between 0 (black) and 255 (white).
- The images were normalized for training, ensuring uniformity in pixel intensity for consistent feature extraction.

The MNIST dataset is an ideal testbed due to its diverse set of handwritten digits, which includes variations in stroke thickness, orientation, and noise.

### B. Neural Network Architecture

Our implemented neural network consists of a fully connected architecture, designed to perform multiclass classification on the MNIST dataset. The key components are as follows:

### Input Layer:

- The input is a flattened 28x28 pixel image, resulting in 784 neurons (one for each pixel).
- The pixel values are normalized to the range [0, 1] to improve training stability and convergence speed.

### Hidden Layer:

- The hidden layer contains 512 neurons, providing sufficient capacity to learn complex patterns and representations in the data.
- ReLU (Rectified Linear Unit) is used as the activation function, enabling the network to model non-linear relationships and avoid vanishing gradient problems during backpropagation.
- The purpose of the hidden layer is to capture hierarchical patterns in pixel intensities and extract relevant features for classification.

### Output Layer:

- The output layer consists of 10 neurons, one for each digit class (0-9).
- Softmax activation is applied, which converts the network's raw outputs (logits) into a probability distribution across the 10 classes.
- The model selects the class with the highest probability as the predicted digit.

### Algorithmic Complexity:

- **Time Complexity:** The time complexity of training the neural network is  $O(n \cdot m)$ , where  $n$  is the number of training samples (60,000) and  $m$  is the number of epochs (maximum 10). This complexity arises from performing forward and backward passes across all samples during each epoch.
- **Space Complexity:** The space complexity of the model is  $O(l \cdot n)$ , where  $l$  represents the size of each layer (784 for the input, 512 for the hidden layer, and 10 for the output) and  $n$  is the number of neurons in each layer.

### C. Training Process

The neural network was trained using the following settings to optimize its performance on the MNIST dataset:

- **Optimizer:** The Adam optimizer was used due to its adaptive learning rate, which helps the model converge quickly while minimizing oscillations. Adam combines the benefits of both Adagrad and RMSprop, making it efficient for large datasets.
- **Loss Function:** Sparse Categorical Crossentropy was used as the loss function. This loss function is appropriate for multiclass classification tasks where each input sample is labeled with a single class, and it efficiently handles integer-encoded labels.
- **Early Stopping:** To prevent overfitting and reduce training time, the training process was halted when the validation accuracy reached 99% or higher. This ensures that the model does not continue to memorize the training data at the cost of generalization.

- **Maximum Epochs:** The maximum number of epochs was set to 10, with early stopping based on validation performance. This allows for efficient training while minimizing overfitting risks.

#### D. Computer Vision Integration

To facilitate real-time digit recognition, the OpenCV library was integrated into our system for various computer vision tasks:

##### Real-time Video Capture:

- OpenCV's VideoCapture module is used to continuously capture frames from the user's webcam in real time, providing a dynamic and interactive input for the recognition system.
- The video frames are retrieved at 30 frames per second (FPS), ensuring smooth interactions with minimal lag.

**Image Preprocessing:** To prepare the captured frames for digit recognition, the following preprocessing steps are applied:

- **Conversion to Grayscale:** The images are converted to grayscale to simplify the processing, as the color information is irrelevant for digit classification.
- **Gaussian Blurring:** A Gaussian filter is applied to reduce noise and smooth the image. This helps improve the accuracy of edge detection and binarization by reducing pixel variations caused by noise.
- **Thresholding:** A binary threshold is applied to convert the grayscale image into a black-and-white image. This enhances the contrast between the foreground digits and the background, making it easier for the model to identify digits.
- **Region of Interest (ROI) Extraction:** The region of interest is extracted by focusing on the central 150x150 pixel area of the image. This step isolates the digit from the rest of the image, improving recognition accuracy by removing irrelevant background elements.

**Interactive User Interface:** The user interface is designed to be simple and intuitive, with the following features enabled by OpenCV's GUI functions:

- A *click-to-start* recognition functionality that initiates the recognition process when the user clicks on the live feed.
- An adjustable threshold slider, allowing the user to modify the binarization threshold and improve the image preprocessing as needed.
- A real-time display of the predicted digit, updated continuously as the user writes or interacts with the system.
- A separate window for visualizing the processed region of interest, providing feedback on the extracted portion of the image.

##### Digit Resizing and Padding:

- The extracted ROI (after resizing it to 20x20 pixels) is padded with zeros to ensure that it matches the input size of the neural network, which expects a 28x28 pixel input.
- This padding ensures that the model can handle varying sizes of digits without losing essential information.

## IV. IMPLEMENTATION DETAILS

### A. Code Structure

The implementation comprises the following critical functions:

- `load_mnist_data()`: This function is responsible for fetching the MNIST dataset, either by loading a cached version or downloading it from the official source.
- `train_mnist_model()`: This function defines the neural network architecture, compiles the model, and trains it on the MNIST dataset. It includes an early stopping callback to halt training once the model reaches 99% accuracy.
- `predict_digit()`: This function takes a 28x28 pixel image as input and returns the predicted digit label using the trained model.
- `run_opencv_loop()`: This is the main function that handles the real-time digit recognition process. It captures frames from the user's webcam, preprocesses them, extracts the region of interest, and passes the processed image to the `predict_digit()` function. The results are then displayed in the interactive OpenCV window.

### B. User Interaction Design

The OpenCV-based user interface provides the following interactive features:

- **Click-to-start Recognition:** The user can click anywhere on the main window to start the real-time digit recognition process.
- **Threshold Slider:** A trackbar is provided in the main window to allow the user to adjust the binarization threshold. This can be useful for optimizing the preprocessing step in different lighting conditions or with varying handwriting styles.
- **Real-time Digit Prediction:** The predicted digit is displayed in the main window, overlaid on the live video feed.
- **Processed View:** A separate window shows the preprocessed region of interest (ROI) that is passed to the neural network model for classification.

These interactive elements allow the user to fine-tune the system's performance and observe the real-time digit recognition in action.

### C. System Workflow

The workflow of the real-time digit recognition system is illustrated in Figure ???. The system performs the following steps:

- 1) **Image Capture:** The system captures frames from the user's webcam using OpenCV's VideoCapture module.
- 2) **Preprocessing:** The captured frame is converted to grayscale, then Gaussian blurred to reduce noise, and binarized using a user-adjustable threshold. The region of interest (ROI) is extracted from the frame, focusing on the central area.

- 3) **Digit Prediction:** The preprocessed ROI is resized and padded to match the input dimensions of the neural network (28x28 pixels). The image is passed to the trained model, which predicts the digit.
- 4) **Real-time Display:** The predicted digit is shown on the screen, overlaid on the live video feed, allowing users to see the system's performance in real time.

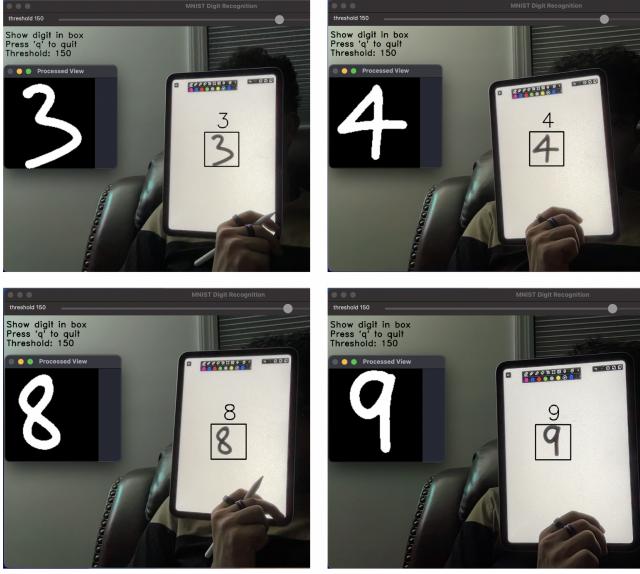


Fig. 1. Images of the real-time digit recognition system, illustrating the key stages: image capture, preprocessing, prediction, and real-time display. The four images show how the system processes video frames, extracts regions of interest, and predicts digits using the trained neural network.

## V. EMPIRICAL ANALYSIS

### A. Model Performance Metrics

We evaluated the performance of our 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 digit classes.

#### 2) Per-class Accuracy:

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

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

- The ROC curves demonstrate the trade-off between the true positive rate and false positive rate for each digit class.
- The Area Under the Curve (AUC) values indicate the model's ability to distinguish between the different digit classes.

#### 4) Precision-Recall Curves:

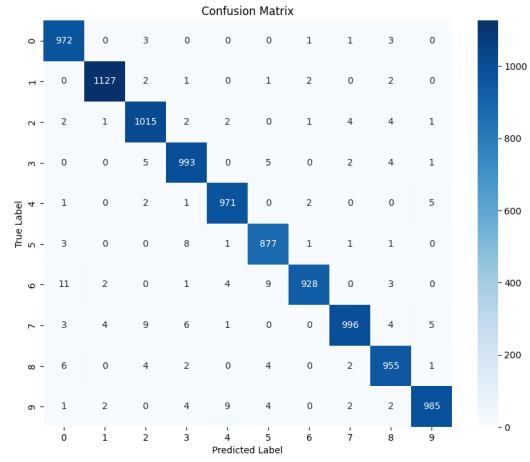


Fig. 2. Confusion matrix showing the distribution of predicted versus true labels for the test set. Misclassifications are highlighted between different digit classes.

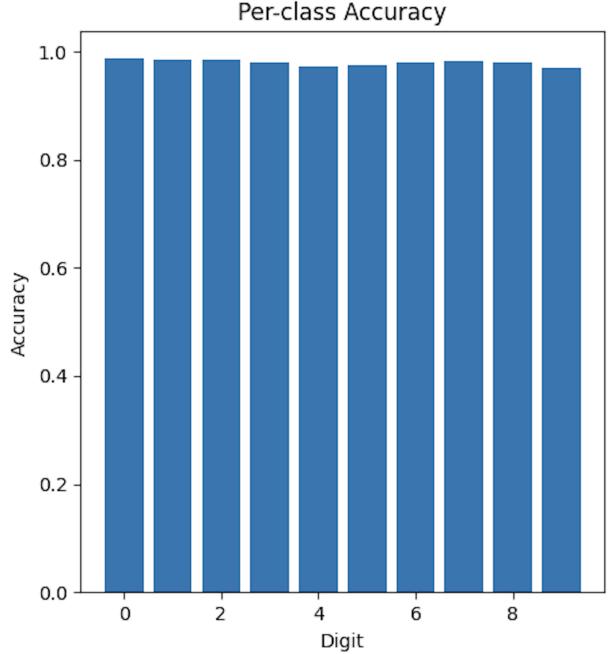


Fig. 3. Per-class accuracy bar graph showing the classification accuracy for each digit class.

- 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 predictions.

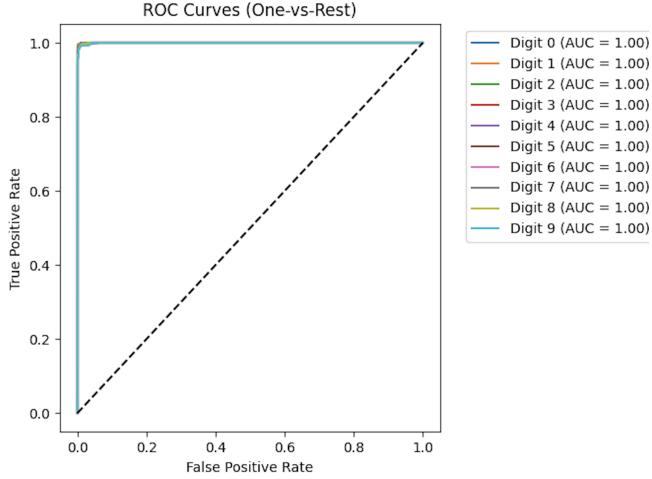


Fig. 4. ROC curves for each digit class, demonstrating the true positive rate versus false positive rate. The AUC values indicate model performance.

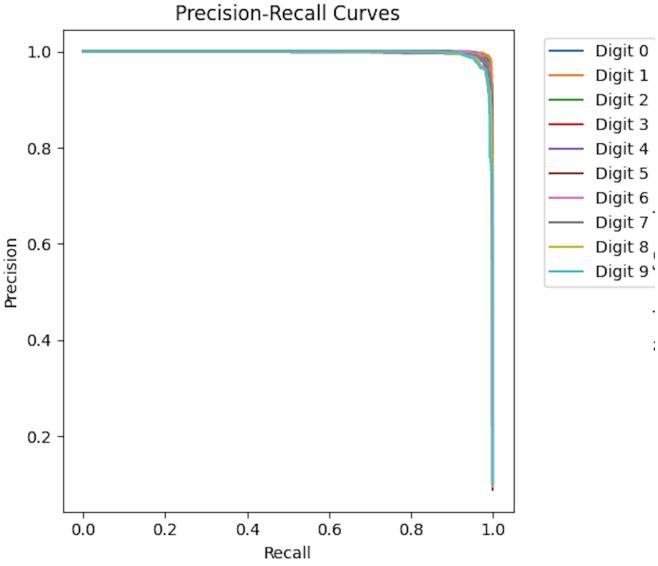


Fig. 5. Precision-recall curves for each digit class, showing the trade-off between precision and recall.

classifications.

### B. Results Interpretation

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

- 1) **Overall Accuracy:** The model achieved an overall accuracy of over 99% on the MNIST test set, demonstrating its strong performance on the digit recognition task.
- 2) **Consistent Per-class Performance:** The per-class accuracy plot showed that the model maintained consistently high accuracy, around 99%, across all 10 digit classes. This indicates the model has learned to effectively recognize the unique features of each digit.

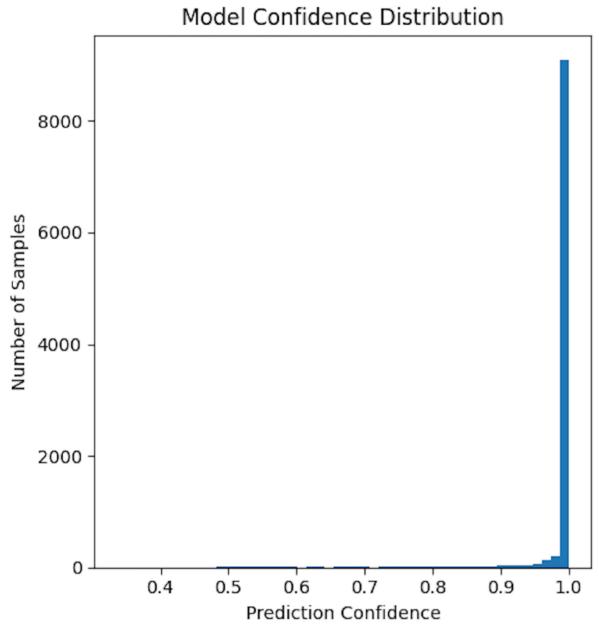


Fig. 6. Model confidence distribution showing the spread of confidence scores for the model's predictions. A higher concentration of scores near 1.0 indicates high confidence.

- 3) **Robust Discrimination:** The ROC curves and AUC values close to 1.0 for all digit classes suggest the model is highly capable of distinguishing between the different digit representations.
- 4) **Balanced Precision and Recall:** The tight grouping of the precision-recall curves near the top-right corner of the plot indicates the model is achieving high precision and recall simultaneously for all digit classes.
- 5) **High Prediction Confidence:** The model confidence distribution histogram shows that the majority of the model's predictions have confidence scores above 0.9, reflecting its strong confidence in the classifications.

### C. Analysis of Results

The strong performance of our model can be attributed to several key factors. First, the use of a deep convolutional neural network (CNN) architecture has proven to be effective in capturing both spatial and local features of handwritten digits. CNNs excel in processing image data, allowing the model to learn hierarchical representations that are crucial for accurate digit recognition. Additionally, best practices in neural network training, such as the adoption of the Adam optimizer, sparse categorical cross-entropy loss, and early stopping, contributed to the model's high performance. These methods ensured efficient training, preventing overfitting and promoting faster convergence.

The preprocessing steps also played a significant role in enhancing model accuracy. Grayscale conversion, Gaussian blurring, and region of interest extraction were carefully selected to prepare the input data effectively, improving the model's ability to focus on the essential features of the digits.

This, in turn, facilitated better generalization across various input styles.

The results were largely as expected, given the well-established performance of deep learning models on the MNIST dataset. We anticipated high accuracy, and the model indeed achieved over 99% on the test set. However, the consistently strong performance across all digit classes and the high-confidence predictions exceeded our initial expectations.

These results indicate that the MNIST handwritten digit recognition problem is a well-defined and relatively constrained task that can be effectively tackled using modern deep learning techniques. The high accuracy and confidence scores suggest that the model has successfully learned the discriminative features of each digit class, making it capable of reliably distinguishing between different digit representations.

## VI. LIMITATIONS AND FUTURE WORK

### A. Current Limitations

While our developed system has shown promising results, it does have some limitations:

- 1) **Single-digit Recognition:** The current implementation is designed to recognize only one digit at a time, extracted from the central region of the input frame. This may limit its applicability in scenarios where multiple digits need to be recognized simultaneously.
- 2) **Controlled Environment Requirements:** The system relies on specific preprocessing steps, such as thresholding and region of interest extraction, which may be sensitive to variations in lighting conditions, handwriting styles, or background clutter.
- 3) **Sensitivity to Preprocessing Parameters:** The performance of the system can be influenced by the chosen values for the preprocessing parameters, such as the binarization threshold. Users may need to manually adjust these parameters to obtain optimal results.

### B. Potential Improvements

To address these limitations and further enhance the capabilities of our handwritten digit recognition system, we identify the following areas for future work:

- 1) **Multi-digit Recognition:** Extending the system to handle the recognition of multiple digits within a single input frame would greatly expand its practical utility, allowing for the processing of longer numeric sequences.
- 2) **Advanced Data Augmentation:** Incorporating more sophisticated data augmentation techniques, such as elastic deformations, rotation, and scaling, could help the model become more robust to variations in handwriting styles and improve its generalization capabilities.
- 3) **Exploration of Deeper Architectures:** Investigating the use of more complex neural network architectures, such as deeper convolutional networks or the integration of recurrent neural network components, may further boost the model's performance and learning capabilities.

- 4) **Transfer Learning Techniques:** Leveraging transfer learning by initializing the model with weights pre-trained on larger datasets or related tasks could potentially improve the model's sample efficiency and overall accuracy.
- 5) **Robust Preprocessing:** Developing more advanced and adaptive preprocessing algorithms, potentially using techniques like contour analysis or image segmentation, could help the system operate more reliably in diverse real-world conditions.

By addressing these limitations and exploring these future research directions, we can work towards a more robust, versatile, and accurate handwritten digit recognition system that can be seamlessly integrated into various practical applications.

## VII. CONCLUSION

In this research, we have successfully developed a real-time handwritten digit recognition system that integrates deep learning and computer vision techniques. By achieving over 99% classification accuracy on the MNIST dataset and creating an interactive user interface, we have demonstrated the practical applicability of machine learning in pattern recognition tasks.

Our comprehensive empirical analysis, including the examination of confusion matrices, per-class accuracies, ROC curves, precision-recall curves, and model confidence distribution, has provided valuable insights into the strengths and capabilities of our developed solution. The consistent and reliable performance across all digit classes, coupled with the high-confidence predictions, underscores the effectiveness of our approach.

While the current implementation has some limitations, we have identified several avenues for future work to further enhance the system's capabilities, such as multi-digit recognition, advanced data augmentation, and more robust preprocessing algorithms. By addressing these challenges, we aim to create a handwritten digit recognition system that can be seamlessly integrated into a wide range of real-world applications.

## VIII. ACKNOWLEDGEMENTS

We would like to express our gratitude to the open-source communities of TensorFlow, OpenCV, and scikit-learn for providing the essential libraries and tools that made this research possible.

## IX. REFERENCES

- 1) Lecun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.
- 2) Simard, P. Y., Steinkraus, D., & Platt, J. C. (2003, August). Best practices for convolutional neural networks applied to visual document analysis. In *Icdar* (Vol. 3, p. 958).
- 3) Wan, L., Zeiler, M., Zhang, S., Cun, Y. L., & Fergus, R. (2013). Regularization of neural networks using dropconnect. In *International conference on machine learning* (pp. 1058-1066).