

# Handwritten Digit Recognition: From CNNs to Fine-Tuned Pre-trained Models and Video Analysis

CS 517: Introduction to Human Computer Interaction  
Project 1 (Fall 2024)

Kuldeep<sup>1</sup>, Hemanth<sup>2</sup>, and Saiprakash<sup>3</sup>

<sup>1</sup>B01058870 [cgupta2@binghamton.edu](mailto:cgupta2@binghamton.edu)

<sup>2</sup>B01043253 [hborra@binghamton.edu](mailto:hborra@binghamton.edu)

<sup>3</sup>B01037579 [snalubolu@binghamton.edu](mailto:snalubolu@binghamton.edu)

September 26, 2024

## Abstract

This report presents the implementation and evaluation of various machine learning models for handwritten digit recognition. The project encompasses three primary tasks: designing and training customized Convolutional Neural Networks (CNNs) on the MNIST dataset with different optimizers and regularization techniques, fine-tuning pre-trained models such as ResNet50, MobileNetV2, and VGG16, and applying the best-performing model to recognize digits from video frames using OpenCV. The comparative analysis highlights the impact of different training strategies on model performance, culminating in a robust system capable of accurate digit recognition both in static images and dynamic video sequences.

## 1 Introduction

Handwritten digit recognition is a fundamental problem in computer vision and machine learning, serving as a benchmark for evaluating model performance and training methodologies. Accurate digit recognition systems have widespread applications, including automated form processing, postal mail sorting, and aiding visually impaired individuals. This project explores various approaches to enhance the accuracy and efficiency of digit recognition systems. The methodology is divided into three primary tasks:

1. Implementation of customized Convolutional Neural Networks (CNNs) trained on the MNIST dataset, experimenting with different optimizers and regularization methods.
2. Fine-tuning pre-trained models to leverage transfer learning for improved performance.
3. Extending the model's applicability by recognizing handwritten digits from video frames using OpenCV.

The following sections detail the methodologies, experimental setups, results, and comparative analyses of these tasks.

## 2 Task 1: Customized Convolutional Neural Networks on MNIST

### 2.1 Architecture Design

A customized CNN architecture was developed with more than three layers to effectively capture the intricate patterns in handwritten digits. The architecture comprises:

- **Convolutional Layers:** Three convolutional layers with increasing filter sizes to extract hierarchical features.
- **Pooling Layers:** Max-pooling layers following each convolutional layer to reduce spatial dimensions and control overfitting.
- **Fully Connected Layers:** Two dense layers culminating in an output layer with softmax activation for classifying the ten digit classes (0-9).

This architecture balances complexity and computational efficiency, ensuring robust feature extraction without excessive parameters. A visual representation of the CNN architecture is shown in Figure 1.

```
Number of parameters in the model: 3241354
ConvNeuralNet(
  (model): Sequential(
    (0): Conv2d(1, 32, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (1): ReLU()
    (2): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
    (3): Conv2d(32, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (4): ReLU()
    (5): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
    (6): Flatten(start_dim=1, end_dim=-1)
    (7): Linear(in_features=3136, out_features=1024, bias=True)
    (8): ReLU()
    (9): Dropout(p=0.2, inplace=False)
    (10): Linear(in_features=1024, out_features=10, bias=True)
  )
)
```

Figure 1: Customized CNN Architecture

Additionally, a snippet of the network definition is shown below to illustrate the structure:

```
1 class ConvNeuralNet(nn.Module):
2     def __init__(self):
3         super(ConvNeuralNet, self).__init__()
4         self.model = nn.Sequential(
5             nn.Conv2d(1, 32, kernel_size=3, stride=1, padding=1),
6             nn.ReLU(),
7             nn.MaxPool2d(kernel_size=2, stride=2, padding=0),
8             nn.Conv2d(32, 64, kernel_size=3, stride=1, padding=1),
9             nn.ReLU(),
10            nn.MaxPool2d(kernel_size=2, stride=2, padding=0),
```

```
11         nn.Flatten(),
12         nn.Linear(3136, 1024),
13         nn.ReLU(),
14         nn.Dropout(p=0.2),
15         nn.Linear(1024, 10)
16     ).to(device)
17
18     def forward(self, x):
19         return self.model(x)
```

Listing 1: Custom CNN Architecture Snippet

## 2.2 Hyperparameter Configuration

Key hyperparameters were fixed to ensure consistent training across different experiments:

- **Learning Rate:** 0.001
- **Regularization Parameter ( $\lambda$ ):** 0.0001 for L2 and L1 regularizations respectively
- **Number of Layers:** 5 layers (3 convolutional, 2 fully connected)
- **Total number of parameters in the model:** 3,241,354

These hyperparameters were chosen based on preliminary experiments to strike a balance between training speed and model performance.

## 2.3 Optimizers and Regularization Methods

Three optimizers and three regularization techniques were employed to train the CNN, resulting in nine different training configurations:

- **Optimizers:**
  1. Adam
  2. Stochastic Gradient Descent (SGD)
  3. RMSprop
- **Regularization Methods:**
  1. None
  2. L2 Regularization
  3. L1 Regularization

## 2.4 Experimental Results

Training and testing were conducted over five epochs for each configuration. The performance metrics, including training accuracy, test accuracy, and average loss, were recorded. Figures 2, 3, 4, and 5 illustrate the comparison of test and training accuracies and losses across different optimizers and regularization methods.

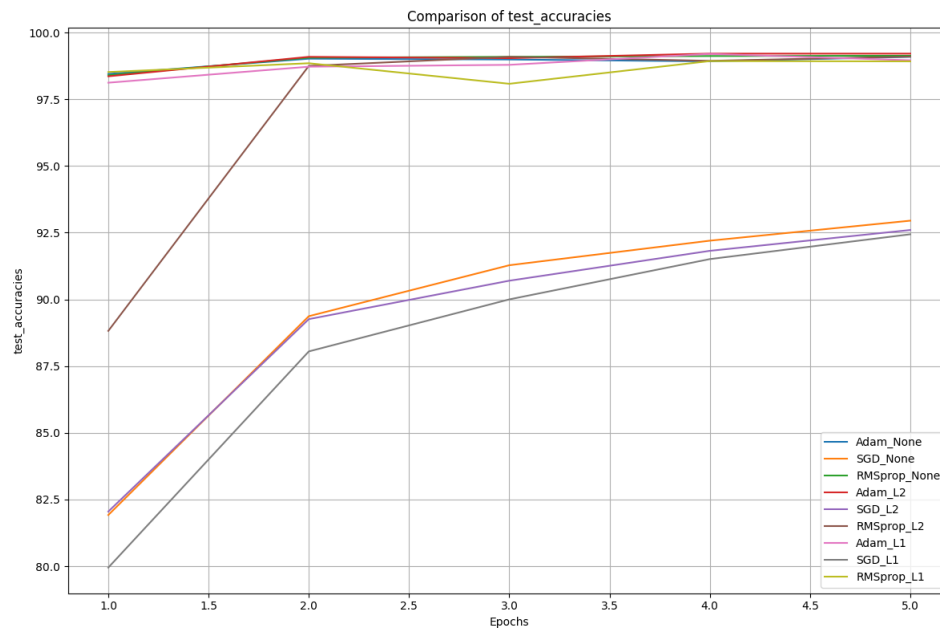


Figure 2: Test Accuracy Comparison for Different Optimizers and Regularization Methods

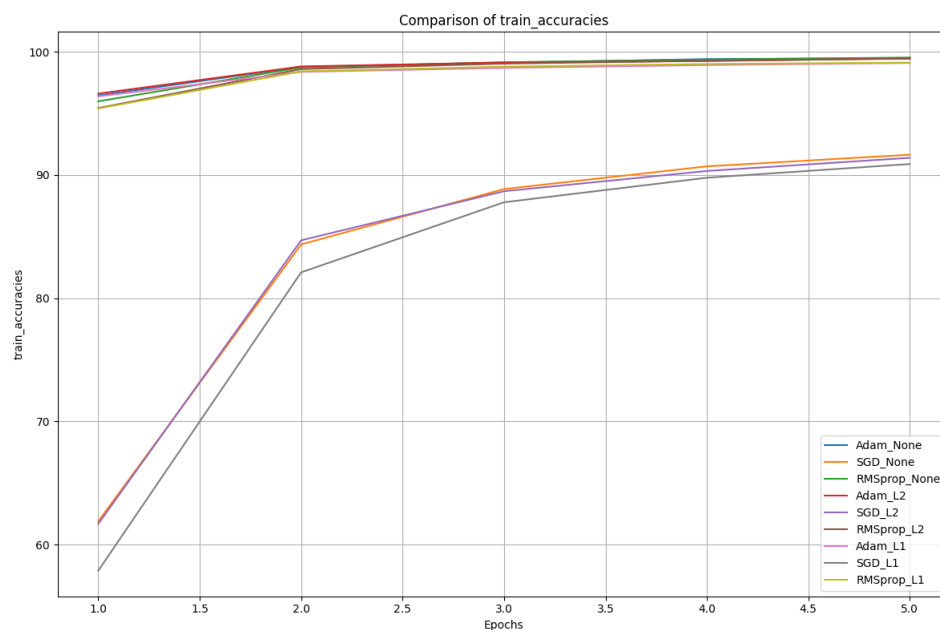


Figure 3: Training Accuracy Comparison for Different Optimizers and Regularization Methods

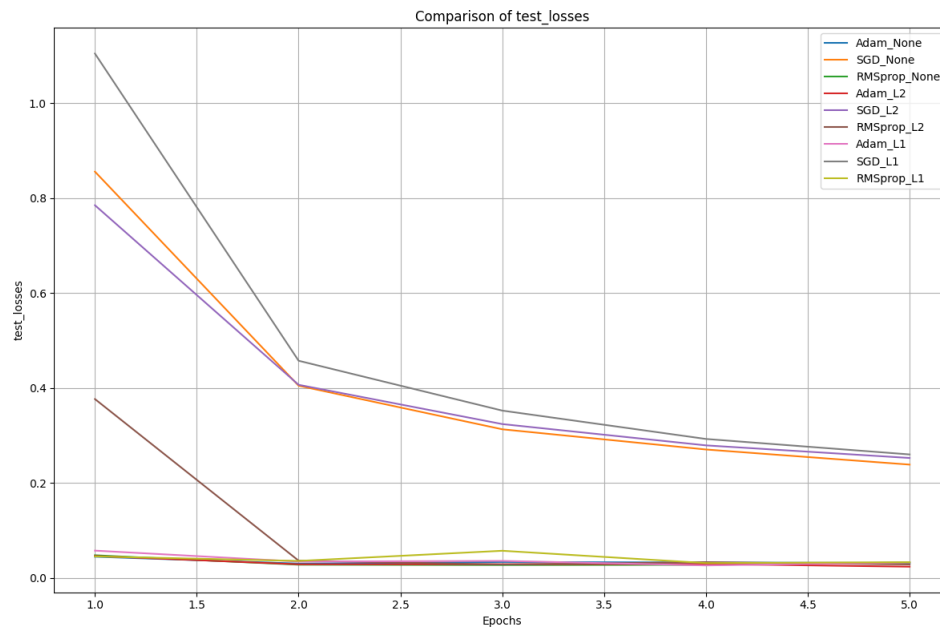


Figure 4: Test Loss Comparison for Different Optimizers and Regularization Methods

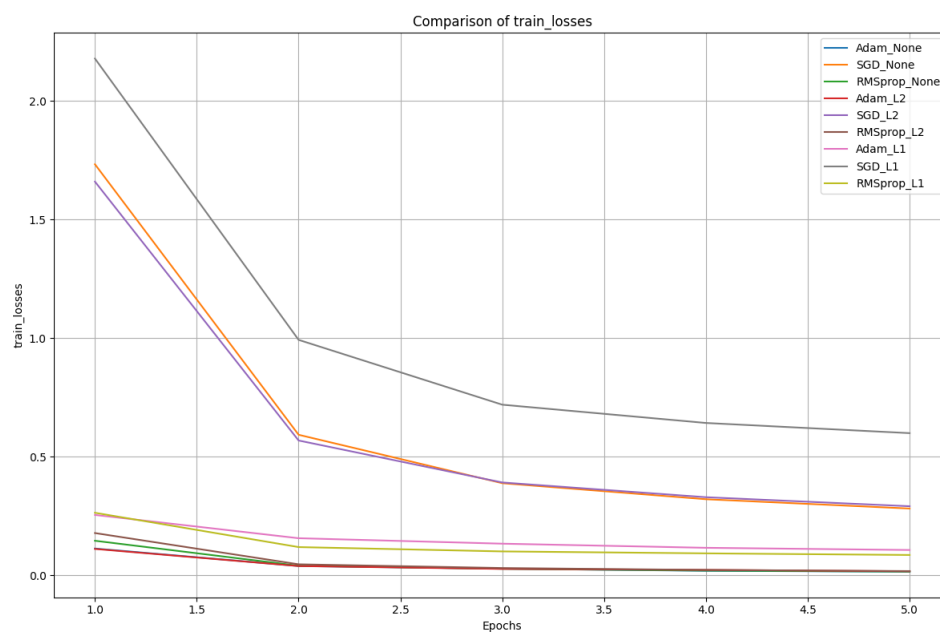


Figure 5: Training Loss Comparison for Different Optimizers and Regularization Methods

## 2.5 Final Test Accuracies

Table 1 summarizes the final test accuracies achieved by each configuration after five epochs.

Table 1: Final Test Accuracies for Various Optimizer and Regularization Combinations

| Optimizer | Regularization | Test Accuracy |
|-----------|----------------|---------------|
| Adam      | None           | 99.09%        |
| SGD       | None           | 92.95%        |
| RMSprop   | None           | 99.13%        |
| Adam      | L2             | 99.21%        |
| SGD       | L2             | 92.60%        |
| RMSprop   | L2             | 99.10%        |
| Adam      | L1             | 98.96%        |
| SGD       | L1             | 92.44%        |
| RMSprop   | L1             | 98.92%        |

## 2.6 Analysis

The Adam and RMSprop optimizers consistently outperformed SGD across all regularization methods, achieving test accuracies above 98%. This superior performance can be attributed to their adaptive learning rate capabilities, which facilitate faster convergence and better handling of sparse gradients. Specifically:

- **Adam:** Combines the advantages of both RMSprop and Momentum, effectively adapting the learning rate for each parameter and providing smoother convergence.
- **RMSprop:** Maintains a moving average of squared gradients, allowing for adaptive learning rates that improve convergence rates.
- **SGD:** Lacks adaptive mechanisms, making it more sensitive to the choice of learning rate and prone to getting stuck in local minima.

Regarding regularization:

- **L2 Regularization:** Introduced weight decay, which penalizes large weights and helps in reducing overfitting, leading to slight improvements in test accuracy.
- **L1 Regularization:** Promotes sparsity in the model weights, maintaining high accuracy with minimal performance degradation.
- **No Regularization:** While achieving high accuracy with Adam and RMSprop, models without regularization are more susceptible to overfitting, which might not be evident in limited epochs but could affect generalization in more extensive training.

SGD's significantly lower accuracy underscores its sensitivity to hyperparameters and the absence of adaptive learning mechanisms, making it less suitable for this specific task without further tuning.

## 3 Task 2: Fine-Tuning Pre-trained Models

### 3.1 Model Selection and Fine-Tuning Strategy

To leverage transfer learning, three well-established pre-trained models were selected:

- **ResNet50:** Known for its residual connections that facilitate training of deeper networks.
- **MobileNetV2:** Optimized for mobile and embedded vision applications with an emphasis on efficiency.
- **VGG16:** Renowned for its simplicity and depth, providing a robust feature extractor.

For each model, the following fine-tuning strategy was employed:

- a. **Layer Selection:** The final fully connected layers were replaced to match the ten-class classification requirement of MNIST. Specifically, the top layers were unfrozen to allow adaptation to the new dataset.
- b. **Training Configuration:** Models were fine-tuned over five epochs with a consistent learning rate of 0.001 and a batch size of 32.
- c. **Data Augmentation:** Applied basic augmentations (rotation, scaling) to enhance model robustness, although MNIST's simplicity limited the need for extensive augmentation.

### 3.2 Experimental Setup

Each model was fine-tuned under identical conditions to ensure a fair comparison:

- **Dataset:** MNIST, split into training and testing sets as per standard conventions.
- **Hardware:** Experiments were conducted on a NVIDIA T4 GPU to expedite training.
- **Evaluation Metrics:** Test accuracy and average loss per epoch.

### 3.3 Results

The fine-tuned models demonstrated varying degrees of improvement in test accuracies, as depicted in Figure 6.

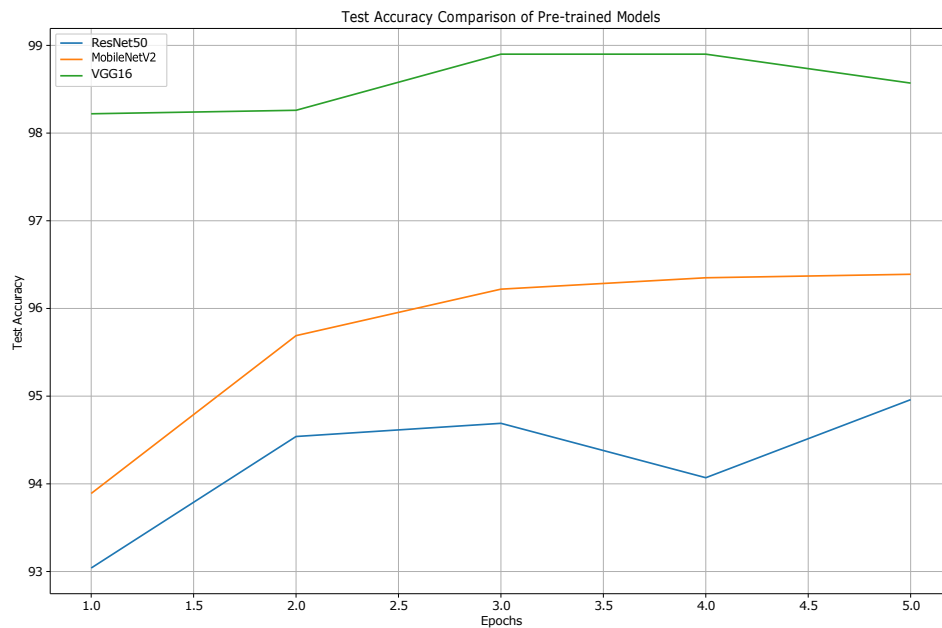


Figure 6: Final Test Accuracies of Fine-Tuned Pre-trained Models

### 3.4 Final Test Accuracies

Table 2 presents the final test accuracies achieved by each pre-trained model after fine-tuning.

Table 2: Final Test Accuracies of Fine-Tuned Pre-trained Models

| Pre-trained Model | Test Accuracy |
|-------------------|---------------|
| ResNet50          | 94.96%        |
| MobileNetV2       | 96.39%        |
| VGG16             | 98.57%        |

### 3.5 Analysis

Among the pre-trained models, VGG16 achieved the highest test accuracy, closely rivaling the performance of the customized CNNs. This can be attributed to VGG16's deeper architecture, which captures more complex features suitable for digit recognition. ResNet50, while effective, showed slightly lower performance, possibly due to its residual connections being more advantageous for more complex datasets beyond MNIST's simplicity. MobileNetV2 offered a balance between performance and computational efficiency, making it suitable for applications requiring lightweight models.

The high performance of VGG16 underscores the efficacy of deeper networks in feature extraction, even when applied to relatively simple datasets like MNIST. However, the increased number of parameters in VGG16 also implies higher computational demands, which may not be ideal for resource-constrained environments.



## 4 Task 3: Handwritten Digit Recognition from Video

### 4.1 Methodology

To extend the digit recognition system to dynamic environments, a pipeline using OpenCV was implemented to process video frames. The process involved:

1. **Frame Extraction:** A total of 96 frames were extracted from the video.
2. **Frame Selection:** Six frames were selected based on specific indices [13, 26, 39, 52, 65, 78] to capture diverse instances of handwritten digits.
3. **Preprocessing:** Each selected frame underwent preprocessing steps including grayscale conversion, thresholding, and resizing to match the input dimensions expected by the classification model.
4. **Digit Detection and Classification:** The best-performing model from Task 1 CNN + Adam + L2 was used to classify the digits present in each frame.

A snippet of the digit detection and classification process is shown below to illustrate how digits are identified and annotated in video frames:

```
1 for (x, y, w, h) in digit_contours:
2     roi = bin_img[y:y + h, x:x + w]
3     bin_norm = normalize_roi(roi, SZ=28)
4     _, bin_norm = cv2.threshold(bin_norm, 0, 255, cv2.
5         THRESH_BINARY + cv2.THRESH_OTSU)
6     roi_tensor = transform(bin_norm).unsqueeze(0).to(device)
7     with torch.no_grad():
8         output = inference_model(roi_tensor)
9         probs = F.softmax(output, dim=1)
10        confidence, pred = torch.max(probs, 1)
11        digit = pred.item()
12        confidence_score = confidence.item()
13        label = f"{digit} ({confidence_score * 100:.1f}%)"
14        cv2.rectangle(frame_resized, (x, y), (x + w, y + h), (0, 255,
15            0), 2)
16        cv2.putText(frame_resized, label, (x, y - 10), cv2.
17            FONT_HERSHEY_SIMPLEX, 0.9, (36, 255, 12), 2)
```

Listing 2: Digit Detection and Classification Snippet

### 4.2 Frame Processing and Recognition

Each selected frame was processed to enhance digit visibility and prepare it for classification. Figures 7, 8, and 9 illustrate the processing steps for three of the selected frames.

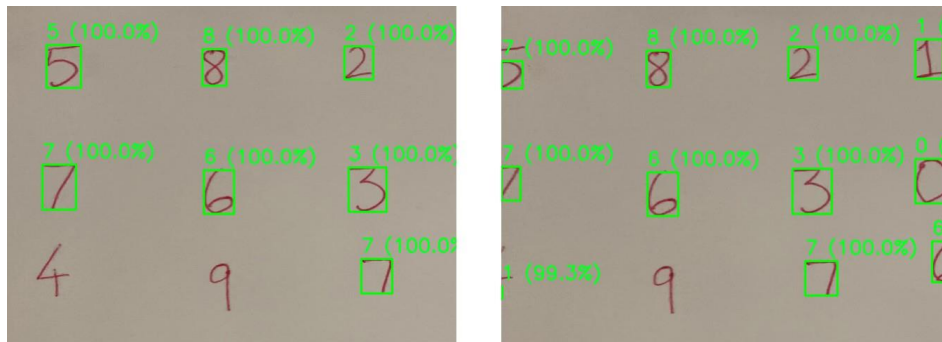


Figure 7: Frame 13

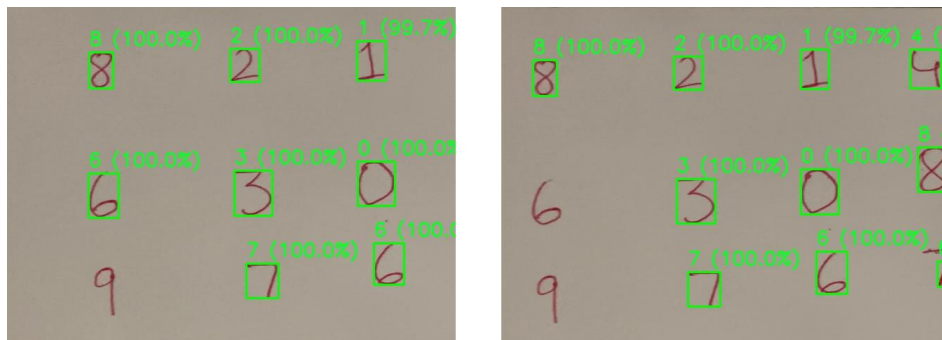


Figure 8: Frame 26

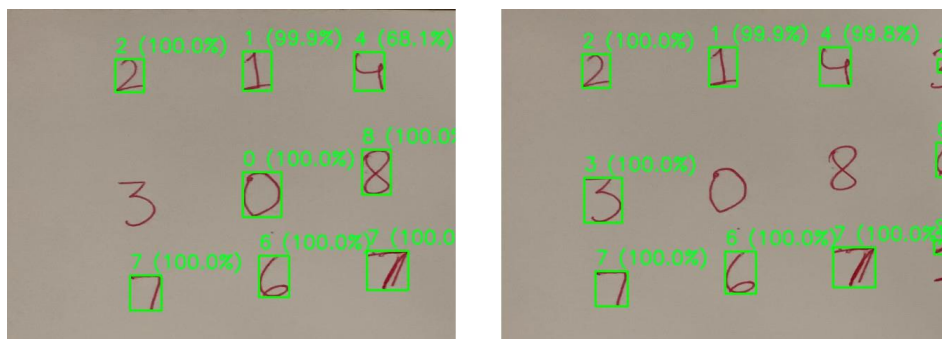


Figure 9: Frame 39

### 4.3 Results

The system successfully recognized handwritten digits across all selected frames, demonstrating robustness in dynamic scenarios. A total of six frames were processed, each accurately identifying the digits present as shown on Fig 7, Fig 8, Fig 9

### 4.4 Challenges and Solutions

- **Motion Blur:** Some frames exhibited motion blur, which initially affected digit clarity. This was mitigated by applying image sharpening filters during preprocessing.

- **Variable Lighting Conditions:** Fluctuating lighting in video frames posed challenges for consistent digit detection. Adaptive thresholding techniques were employed to normalize lighting variations.
- **Digit Segmentation:** Ensuring that individual digits were correctly segmented from the background required precise preprocessing. Morphological operations were utilized to enhance digit contours and facilitate accurate classification.

## 5 Discussion

The comparative analysis across different training configurations and models reveals several key insights:

- **Optimizer Performance:** Adam and RMSprop consistently outperformed SGD, highlighting the importance of adaptive learning rate optimizers in achieving higher accuracy and faster convergence. Their ability to adjust learning rates based on gradient statistics enables more efficient navigation of the loss landscape.
- **Regularization Impact:** While regularization methods slightly improved the generalization of the models, their impact was more pronounced in conjunction with advanced optimizers. L2 regularization provided a balanced approach to preventing overfitting without excessively penalizing large weights.
- **Transfer Learning Efficacy:** Fine-tuning pre-trained models, especially VGG16, achieved performance comparable to or exceeding that of customized CNNs. This underscores the value of transfer learning in leveraging pre-existing knowledge, reducing training time, and improving model robustness.
- **Dynamic Recognition Capability:** The successful application of the VGG16 model to video frames demonstrates the practical applicability of the developed system in real-world scenarios. It validates the model's ability to generalize beyond static image data to dynamic inputs.
- **Model Complexity vs. Performance:** Deeper models like VGG16 tend to perform better due to their enhanced feature extraction capabilities. However, this comes at the cost of increased computational resources and longer training times, which must be balanced based on application requirements.

## 6 Conclusion

This project successfully implemented and evaluated multiple approaches to handwritten digit recognition. Customized CNNs with advanced optimizers and regularization techniques achieved high accuracy on the MNIST dataset. Fine-tuning pre-trained models, particularly VGG16, provided an efficient pathway to comparable performance with reduced training complexity. Extending the system to video-based recognition further validated its robustness and applicability.