
Digitizing Documents with a Neural Network Trained for Optical Character Recognition

Jason Albanus
ECE
jasona@vt.edu

Will Fete
ECE
willf00@vt.edu

Cameron Smith
CS
camerons03@vt.edu

1 Introduction

Handwriting recognition, the task of converting handwritten text images into machine-readable characters, has practical applications given the large number of scanned handwritten documents needing to be digitized. Although traditional optical character recognition (OCR) techniques have existed since the 1970s, their effectiveness is limited to printed text due to uniform spacing and font. On the other hand, handwritten text has a lot of challenges due to the considerable variety of handwriting, like stroke shape and size, spacing, and pressure of writing instruments [1]. The variability greatly complicates character segmentation and recognition. As a result, classic OCR is not reliable or effective for handwritten documents and requires a different approach.

Recent advancements in deep learning, specifically in convolutional neural networks (CNN), have shown better results in overcoming the limitations of OCR. No single universal CNN architecture exists as a perfect handwriting classifier. Moreover, no single CNN currently achieves high accuracy across different handwriting styles, which leaves a lot of room for further improvement.

This project aims to address the challenges of OCR by developing and comparing CNN models specifically designed for handwritten word recognition. To facilitate supervised learning, the IAM handwriting dataset from the University of Bern was chosen due to its structured pre-segmented format. The pre-segmentation of words allows for a decreased amount of pre-processing. However, there was still some pre-processing including removing corrupted or unclear images and punctuation marks to improve training efficiency and reliability.

Initially, we developed a baseline CNN model, which consisted of four convolution blocks designed to capture features from input images. However, early experiments revealed lots of limitations, including varying validation performance and poor accuracy. Motivated by these initial findings, we shifted focus and developed an improved model leveraging the ResNet18 architecture. Some of the main modifications, such as reducing convolutional steps and employing other strategies.

Comparative results clearly show the effectiveness of the second approach. The initial baseline CNN plateaued at a low accuracy of approximately 45%, indicating insufficient model complexity. Unlike the baseline model, the optimized CNN achieved training accuracy exceeding 95% and validation accuracy around 80%. Although this improved upon the first, the second model shows overfitting, showing that there needs to be further refinement.

Lastly, a simple graphical user interface (GUI) was created using Tkinter, allowing users to upload and classify handwritten images. Despite its simplicity, this front-end effectively demonstrates the potential practical application of the trained models.

33 2 Problem Description

34 The problem this project aims to solve is handwriting recognition. While millions of handwritten
35 documents have been scanned and stored as images, converting their contents into machine-readable
36 text offers significant advantages, such as easier search, storage, and analysis. However, manually
37 transcribing these documents is impractical at scale.

38 Optical character recognition (OCR), has existed since the 1970s; it is the process of turning text
39 from images to the respective characters. This conversion is fairly easy due to the predictable spacing
40 and similarity of digital fonts. The problem is that traditional OCR systems struggle with handwritten
41 text, instead handwriting recognition is required. Handwriting recognition is the process of turning
42 handwritten text from images to the correct characters. This is significantly more difficult than OCR
43 due to the high variance in human handwriting. The variability in shapes, styles, and writing pressure
44 complicates the task. Unlike text images, the unpredictable spacing and styles with handwritten
45 characters makes it more difficult to segment characters. While deep learning has produced promising
46 results, there is no universal model that performs reliably across all handwriting styles. That universal
47 solution is what this project aims to find, using convolutional neural networks designed and tuned to
48 best accomplish this task.

49 To train a model capable of recognizing handwriting, a high-quality dataset is needed. The IAM
50 dataset from the University of Bern was chosen for training and testing purposes as it was pre-
51 tokenized in terms of lines, sentences and words. To prevent biased results, punctuation was removed
52 as the images of them were blurred to the point that they were unrecognizable, even to humans.

53 This problem has many challenges beyond what was mentioned above. To start off, we focused
54 on the recognition of the pre-segmented words from the dataset. This meant that the images used
55 needed to be normalized and padded to provide a uniform dimension for the model. Since the words
56 were pre-segmented, supervised learning can be utilized by providing the correct word transcriptions.
57 There are many limitations to this approach though, as sentences need to be converted word by word.
58 Additionally, limited training input meant that data augmentation techniques were needed to ensure
59 the resulting model would be applicable with unseen data. The model will struggle to recognize
60 words it has never seen in training due to it being trained on a limited vocabulary. Having described
61 the challenges of handwriting recognition, we now turn to the approach taken to solve them.

62 3 Methodology

63 3.1 Loading the Dataset

64 When preparing the data from the IAM dataset for model training and testing, care needed to be taken
65 with regards to the formatting decisions from the original provided file map. Some of the images in
66 the dataset had been corrupted, possibly from the original file transfer itself, thus each image had to
67 pass through a verification process. In addition, we found major sources of statistical noise and error
68 from the considerable amount of punctuation marks included in the dataset. To account for this, those
69 images were removed, and a final mapping was formed to pair each image's relative path with an
70 integer class label. Upon dataset initialization, no transformations were specified, as that was done
71 within the respective CNN model files as necessary. At load time, each image is converted to RGB,
72 and the transform is applied before pairing with the respective class label. Additionally, each image
73 is padded on the right with white pixels to ensure all items in the batch are of the same width. Once
74 initialized, the dataset is split into training and testing subsets with data loaders configured to provide
75 the necessary tensors needed by the networks.

76 3.2 Our First Model

77 The first convolutional neural network developed is a baseline CNN model for word classification. It
78 takes a variable-width 3-channel RGB word image and pushes it through four convolutional blocks.
79 The structure of each block consists of a convolutional layer, batch normalization, a rectified linear
80 unit (ReLU), and max-pooling. Specifically, the convolution layers utilize a 3 by 3 kernel with
81 padding set to maintain spatial dimensions. In addition, bias terms are disabled due to the batch
82 normalization after. The batch normalization standardizes the activations, stabilizes training, and
83 speeds up the convergence. After the ReLU enables the model to capture complex patterns, followed

84 by max-pooling to progressively down-sample feature maps in a 2 by 2 block. This decreases the
85 computational need and highlights dominant features.

86 After each block, channel depth increases, starting from 3 channels and going to 64 in the first block.
87 The network ends with a 512-channel feature map after the completion of the 4 blocks. The feature
88 depth allows the model to capture many important characteristics, including basic edges and texture.
89 The choice of four blocks was based on GPU hardware constraints. Further, after the completion of
90 the convolution stage, we convert the spatial feature maps into the format of a single column of fully
91 connected layers, regardless of image width and height. This is completed using adaptive average
92 pooling. It compresses the spatial dimensions to a single value per feature channel, which creates
93 a fixed-length representation of the image. This provides the flexibility needed for a wide variety
94 of image sizes. After the adaptive pooling is applied and the feature vector is a one-dimensional
95 tensor, a dropout is applied. The dropout is applied with a probability of 0.3 and serves as a way to
96 mitigate overfitting by randomly deactivating a subset of neuron connections during training. Finally,
97 the flattened vector is put through a single fully connected linear layer, which outputs the logits
98 corresponding to the word classes learned.

99 After the model was created, we created the training function. Training involved iterative updates to
100 the network parameters, where each epoch comprised a complete pass through the entire training
101 dataset. Images were grouped into small batches for computational efficiency and gradient estimation.
102 For the training, we chose the stochastic gradient descent (SGD) optimization with momentum set at
103 0.9. This enhanced the convergence stability and speed by smoothing the gradient directions. An
104 initial learning rate of 0.03 was chosen after some initial testing alongside a weight decay of 0.0001.
105 This further helps to prevent overfitting. To dynamically refine the learning rate, we implemented a
106 ReduceLROnPlateau, which allows for convergence on an optimal solution. During training, model
107 accuracy and loss were continuously monitored across both training and validation datasets. Each
108 epoch ended with evaluation across both training and validation datasets, which provided a direct
109 assessment of generalization performance.

3.2.1 Our Second Model

For the improved CNN design, we started by importing a pretrained ResNet18 model to use as a feature extractor. It leverages the ImageNet weights, but includes several modifications from the default configuration to achieve better performance for the handwriting recognition task. Whereas the first convolution’s stride is typically 2 by 2 in the ResNet architecture, by reducing it to 1 by 1 we increase the model’s ability to preserve fine-grained details, which is valuable given the task at hand. In addition, the first block of the fourth, and final, residual stage has its main convolution and down-sample strides reduced from 2 by 2 to 1 by 1. Our aim was to prevent both premature loss of the fine-grained spatial detail and excessive down-sampling in the feature extraction. When character distinctions can be within a few pixels, these adjustments were necessary for the task to be performed well. Following the first stage’s convolution, a ReLU activation and 3 by 3 max pool are used before entering the four residual stages.

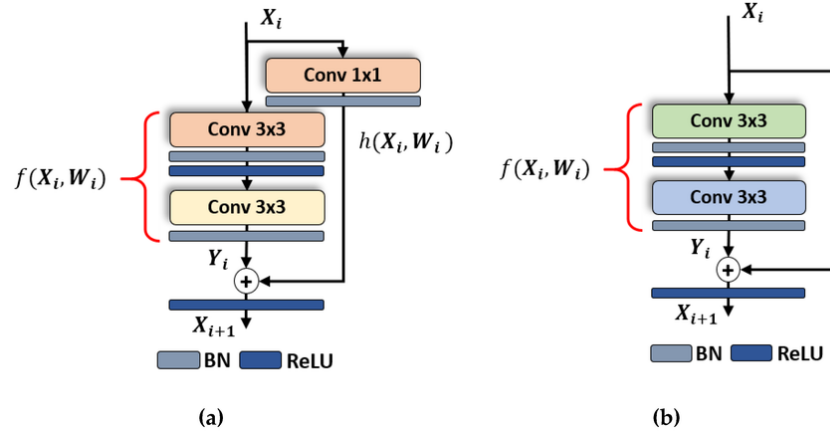


Figure 1: Residual Network Architecture (Obtained from ResearchGate [3]).

Within each residual stage, a 3 by 3 convolution, batch normalization, and ReLU activation sequence is applied twice. For the first residual stage, an identity skip is used (as seen in diagram (b) from Figure 1), however for the remaining three, to factor the changes for channel and spatial sizes, we used a 1 by 1 projection convolution (as seen in diagram (a) from Figure 1). The first stage maintains the input resolution and 64 input channels for its output. The second stage doubles the channel size to 128 with stride 2 down-sampling, and the third stage repeats that doubling to 256 with the same stride 2 down-sampling. Finally, the fourth stage outputs the 512 channels without further changes to the feature map dimensions. Every convolution is followed by a batch normalization and ReLU activation, so as to distribute the nonlinearities uniformly. ReLU Activation is used after the global pooling layer to ensure the classifier receives the raw, unbounded feature map. Using an average pool after the four residual layers, the spatial dimensions are reduced from 16 by 40 (output from the fourth layer) to 1 by 1.

Building upon this, we leverage an optimized DropBlock, as opposed to the conventional Dropout layer, for our regularization prior to classification. The DropBlock converts a scalar drop probability (0.1 for our model) into a spatial gamma, sampling a binary mask, and with max pooling, expands each dropped pixel into a 5 by 5 block. Using the masked feature map, we were able to force the network to distribute information across the map to mitigate over-reliance on any of the pen stroke fragments from the handwriting samples. For every call to the forward method, the DropBlock will flatten the 1 by 1 spatial dimensions to be fed into a single linear layer for classification. The linear layer will output one logit per class, each of which is passed to a Cross Entropy loss function (with applied softmax). As a result, our optimized model is well designed for the handwriting task, keeping the fine spatial details and regularizing against overfitting, with fine-tuned adjustments for memory efficiency.

145 3.3 Rationale for the Changes in Architecture

146 The first model we developed relied on the small four block convolutional network optimized
147 with SGD without transfer learning. Early experiments indicated to us that this model's validation
148 accuracy climbed slowly, reaching around 40-50% after many hours of training. To stay within the
149 computational limits imposed by the hardware at our disposal (NVIDIA 3070 TI desktop GPU),
150 development pivoted midway through the semester to start from scratch with the second architecture.
151 By using what we had learned from our original results, the revised model implemented several
152 optimization techniques with the goal of achieving an earlier plateau. This had the effect of reducing
153 the overall training time by several orders of magnitude on the same underlying hardware, allowing
154 a higher performing model to be created that would successfully achieve the results we had sought
155 after. While our original goal was to build a single model, after starting from scratch with the new
156 architecture, we did not want our previous work to go unrecognized. Thus, we made the decision to
157 re-frame our intentions towards a direct comparison between the residual neural network, and a more
158 traditional convolutional network architecture. Even with the under-trained model, it was clear the
159 plateau was approaching at a noticeably lower accuracy, hinting that our expectation was correct, and
160 the major improvements from the ResNet architecture massively contributed to its success.

161 3.4 User Interface

162 For our project's front-end, we used a simple Tkinter interface design which uses a file selection
163 dialog for the user to upload handwritten images. After an image is uploaded, we convert it to proper
164 grayscale and pass it through the same preprocessing pipeline used by the second model. In doing so,
165 the image is resized to the fixed height and width, with the same normalizations performed using
166 the ImageNet constants. The pretrained model file is loaded in with the respective label mapping,
167 specific to that model's training results. The inference is executed on said model to compute the
168 class probabilities from the image passed in by the user. The top ten probabilities are written to
169 console, with the most significant being displayed directly on the GUI. While this front-end was
170 slightly more rudimentary in design than what we had originally planned, it serves its purpose as a
171 way to adequately compare results.

172 4 Results

173 In this section the results of the models will be shown and discussed. There will be three main
174 implementations that will be evaluated in this section, the first model, the second model built upon
175 the ResNet architecture, and the simple front end. Starting off with the initial model, the training
176 function captured both the accuracy and loss for the validation and training datasets. This will be the
177 main results that will be analyzed for the first model. With the captured accuracies and losses, they
178 can be plotted against the epochs to show the model's progression over time. These results can be
179 seen below.

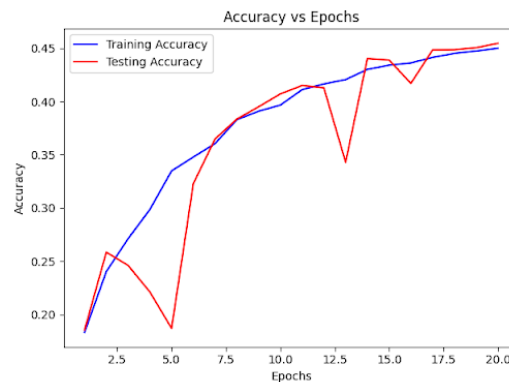


Figure 2: First Model - Plotted Accuracy.

180 As clearly seen in Figures 2 and 3, the performance of the first model was not great. The training
181 accuracy was initially very low, around 20%, before gradually increasing before plateauing at epoch

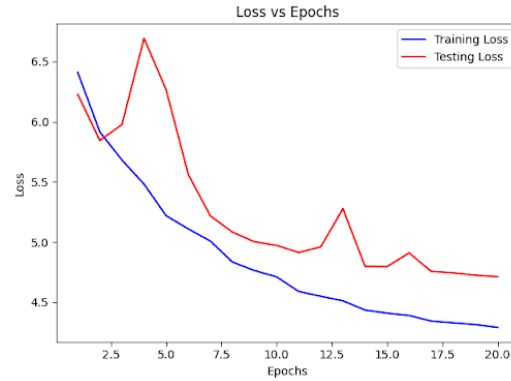


Figure 3: First Model - Plotted Loss.

182 15 while at 45%. The training loss was very similar to the training accuracy except mirrored. The
 183 loss started pretty high before gradually decreasing to about 4.5 and flattening out around the 15th
 184 epoch. While the training loss and accuracy left a lot of room for improvement, the testing results
 185 were a main reason for the creation of the second model. As seen in the below figures, the testing loss
 186 and accuracy was very erratic and inconsistent. This pointed to the model being undertrained and a
 187 potential error in the implementation. Besides the dips and peaks, the model achieves similar results
 188 to the training sets in the validation phase. As discussed above, this combined with other factors such
 189 as computational time and power, led to the decision to move toward the second model.

190 Now the results of that second model can be shown, and compared with the first model to see whether
 191 the decision to create a new model was justified. The results will be shown the same as the first
 192 model, in terms of accuracy and loss versus epochs to show the second model's progression over
 193 time. These results can be seen below.

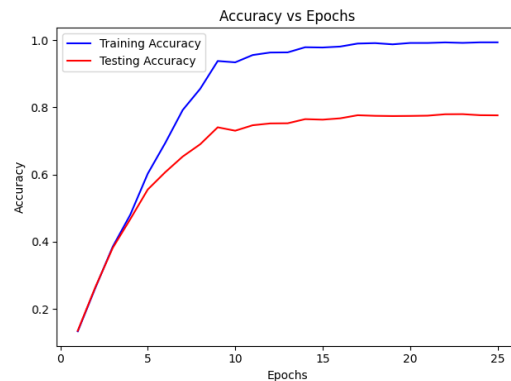


Figure 4: Second Model - Plotted Accuracy.

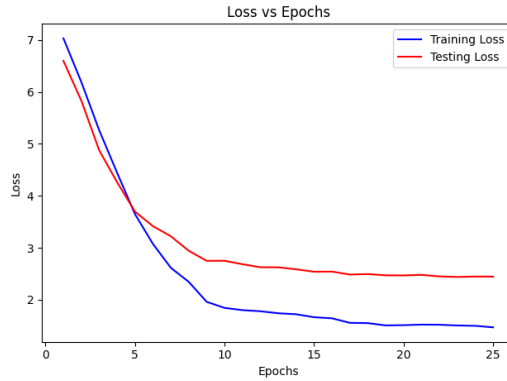


Figure 5: Second Model - Plotted Loss.

As seen in Figures 4 and 5, the results obtained are an improvement. The training accuracy starts low then rapidly increases to above 95% before flattening out. The loss is similar as it starts out around 7 before diving down to around 1. The model had the same behavior for the validation dataset, except the test accuracy was lower and the loss was higher at 80% and 3, respectively. Looking at the graphs, it is clear that this second model is overfitting, which can be due to high complexity, not enough data, and lack of enough regularization. While this second model led to better performance in terms of accuracy, loss, and consistency, the evidence of overfitting is important to recognize. With the results of both models, it is necessary to compare them to decide whether the decision to switch model architectures halfway through the project was justified.

To further test not only the second model, but the front end also, we can feed unseen data to the model through the front end. To do this, we created a small set of input data using our own handwriting of known words in the model's vocabulary. After selecting these images in the GUI, the results were promising. As seen in Figure 6 below, the model made several correct predictions, while the front end displayed the predicted word and its confidence. As you can see, it still is not perfect and limited by the training data size and architecture. The results of the front end testing were satisfactory and as expected. Overall, the results as the implementations discussed in the methodology were acceptable, with room for improvement in the future.

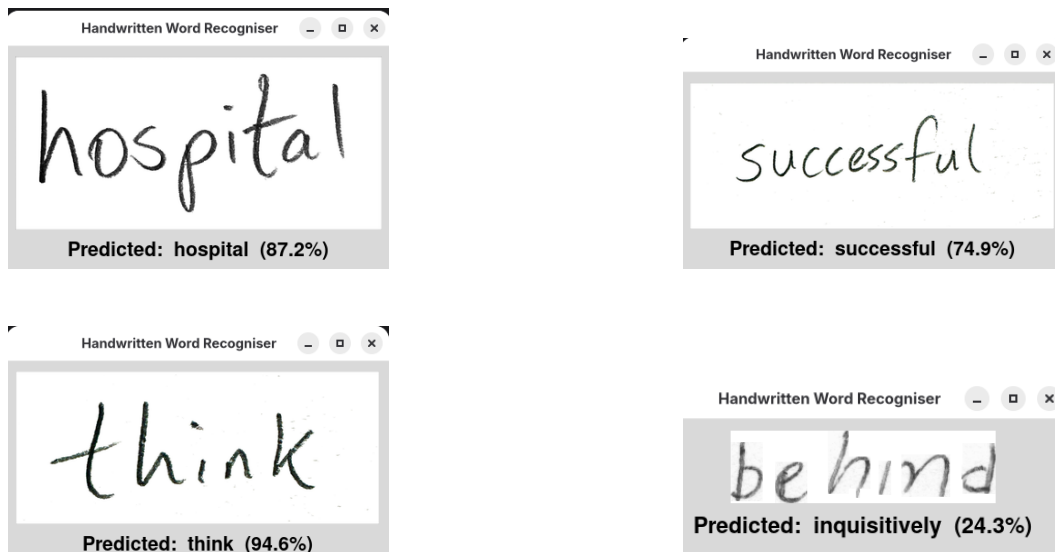


Figure 6: Four Examples from the GUI.

5 Conclusion

This work has shown us the capabilities that convolutional neural networks have at effectively addressing the many challenges imposed by handwriting recognition. There were many instances of success in addition to the setbacks, both of which provided practical exposure and insight. We were able to learn the importance of transfer learning in its ability to enhance accuracy and expedite convergence. In addition, understanding the delicate balancing act between model complexity against available computational resources, data filtering and augmentation to mitigate overfitting on a limited dataset, the effects that preserving fine-grained detail via early down-sampling had, and the improvements to generalization that were made with more advanced regularization techniques like DropBlock. Potential avenues to explore for future work include training the first model to completion on the ARC clusters provided by Virginia Tech, experimenting with attention-based architectures to handle full line and sentence recognition, expanding the dataset to broaden the model's applicability, and refining the front-end user interface for a more seamless interactive experience. Overall, this work helps to demonstrate the practicality of the methods of machine learning used for handwriting recognition tasks.

References

- [1] P. C. Vashist, A. Pandey and A. Tripathi, *A Comparative Study of Handwriting Recognition Techniques*, 2020 International Conference on Computation, Automation and Knowledge Management (ICCAKM), Dubai, United Arab Emirates, 2020, pp. 456-461, doi: 10.1109/ICCAKM46823.2020.9051464
- [2] Ahlawat, Savita et al. *Improved Handwritten Digit Recognition Using Convolutional Neural Networks (CNN)*. Sensors (Basel, Switzerland) vol. 20,12 3344. 12 Jun. 2020, doi:10.3390/s20123344
- [3] *Artificial Intelligence-Based Classification of Multiple Gastrointestinal Diseases Using Endoscopy Videos for Clinical Diagnosis* - Scientific Figure on ResearchGate. Available from: https://www.researchgate.net/figure/Residual-block-of-ResNet18-with-a-1-1-convolutional-mapping-based-residual-unit-and_fig3_334301817 [accessed 23 April 2025]
- [4] Retsinas, George, et al. *Best practices for a handwritten text recognition system*. International Workshop on Document Analysis Systems. Cham: Springer International Publishing, 2022.
- [5] Chowdhury, A., Vig, L.: "An efficient end-to-end neural model for handwritten text recognition" (2018)
- [6] He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.
- [7] Ogochukwu Patience Okechukwu, Godson Nnaeto Okechukwu, Okwuchukwu Ejike Chukwuogo, & Amanda Uloma Anyigor-Ogah. (2024). *Implementing handwritten text recognition using deep learning with TensorFlow: An MNIST dataset approach*. In World Journal of Advanced Engineering Technology and Sciences (Vol. 12, Issue 2, pp. 544–552). GSC Online Press. <https://doi.org/10.30574/wjaets.2024.12.2.0320>