# OCR for Arabic Handwritten Text Recognition

DR. Walid
CCE
FACULTY OF
ENGINEERING

(ALEXANDRIA
UNIVERSITY)

Malak Reda , Rowan ahmed ,
Mena Majidi

Alexandria, Egypt

*Abstract*—*Arabic Handwritten Recognition (AHR) plays a crucial role in digitizing Arabic manuscripts, documents, and historical texts, enabling efficient search and retrieval processes. Despite advancements in Optical Character Recognition (OCR), Arabic handwritten text presents unique challenges due to its complex script and diverse writing styles. This paper explores the current state-of-the-art techniques and methodologies in Arabic Handwritten Recognition, emphasizing the significance of accurate recognition algorithms for preserving cultural heritage and facilitating information access. Through a comprehensive review and analysis, this paper aims to provide insights into the progress, challenges, and future directions of Arabic Handwritten Recognition technology.*

*Keywords*—*OCR, Arabic Handwriting, Deep Learning, Convolutional Neural Networks, Text Recognition*

## I. Introduction

1)**Arabic, revered as one of the world's oldest and most culturally significant languages**, stands as a testament to the richness of human civilization and ingenuity. Its script, characterized by elegant and intricate letter shapes, embodies centuries of artistic expression, scholarly discourse, and historical documentation. From the poetry of Al-Mutanabbi to the philosophical treatises of Ibn Rushd, Arabic literature has transcended borders and time, shaping the intellectual landscape of the Middle East and beyond.

2)**However, amidst this vast reservoir of knowledge lies a challenge** that has confounded scholars and technologists alike: the recognition of handwritten Arabic text. Unlike its printed counterpart, which OCR systems have made significant strides in deciphering, handwritten Arabic presents a myriad of complexities that defy conventional automated recognition methods. From the fluidity of calligraphic styles to the subtle nuances of ligatures and diacritics, each stroke of the pen carries with it a wealth of linguistic and cultural heritage that is both intricate and profound.

3)**The demand for accurate Arabic Handwritten Recognition (AHR) systems has surged in recent years,** fueled by the imperative to preserve and digitize invaluable collections of Arabic manuscripts, historical documents, and contemporary texts. As libraries, museums, and cultural institutions embark on ambitious digitization projects, the need for robust AHR systems capable of faithfully transcribing handwritten Arabic text becomes increasingly pressing.

4)**In this paper, we embark on a journey into the labyrinthine world of Arabic Handwritten Recognition**, navigating through the intricacies of its script and the challenges it poses for automated recognition systems. We delve into the depths of existing methodologies, scrutinizing their efficacy in the face of the unique complexities of handwritten Arabic

## II. Data Collection

In this study, three distinct datasets were utilized to train and evaluate the Arabic Handwritten Optical Character Recognition (OCR) system:

**The first dataset, denoted as the Arabic Handwriting Analysis for Word Processing (AHAWP) dataset**, comprises a diverse collection of handwritten Arabic text samples, representative of various writing styles and linguistic nuances. The dataset contains 65 different Arabic alphabets (with variations on begin, end, middle and regular alphabets), 10 different Arabic words (that encompass all Arabic alphabets) and 3 different paragraphs. The dataset was collected anonymously from 82 different users. Each user was asked to write each alphabet and word 10 times.

**The second dataset, known as Khatt**,The KHATT dataset contains images of handwritten Arabic text samples, along with corresponding ground truth annotations. These annotations typically include the text content of each image, allowing for training and evaluation of Arabic handwriting recognition algorithms

## III. Preprocessing

Preprocessing is crucial for preparing raw data for machine learning tasks, especially in the domain of handwritten text recognition.Here are some details on preprocessing techniques commonly applied to the datasets :

*1)Bounding Box Enclosure: By enclosing each handwritten word within bounding boxes, the text is isolated from background noise, ensuring that the model focuses solely on the relevant text regions during training.*

*2)Grayscale Conversion: Converting the images to grayscale simplifies feature extraction, making it easier for the model to discern the shapes and structures of the handwritten characters without the complexity of color variations.*

*3)Pixel Value Normalization: Normalizing the pixel values of the images to a common scale minimizes the impact of lighting and contrast variations, ensuring that the model learns from more consistent input data.*

*4)Character Segmentation: Segmentation techniques are crucial for separating individual characters, particularly in cases where handwritten text contains overlapping characters or ligatures. This step enables the model to recognize each character independently, improving overall accuracy.*

*5)Extraneous Space Removal: Cropping extraneous spaces surrounding the words minimizes unnecessary variations and streamlines the recognition process, ensuring that the model focuses only on the essential text regions.*

*These preprocessing techniques collectively enhance the quality, consistency, and suitability of the input data for training handwriting recognition models on the dataset, ultimately leading to more robust and accurate performance.*

## V. MODEL ARCHITECTURE

*The core of the OCR system's architecture lies in its utilization of deep learning algorithms, which have demonstrated exceptional performance in various pattern recognition tasks. The model architecture employed in this project integrates a combination of:*

*1)Convolutional Neural Networks (CNNs): Facilitate effective feature extraction from the input images.*

*2)Bidirectional Long Short-Term Memory (BiLSTM) Networks: Enable the model to capture temporal dependencies and contextual information within the input sequences.*

*3)Recurrent Neural Networks (RNNs): Allow the model to process sequential data and exploit the sequential nature of handwritten text.*

*Additionally:*

- *Flatten Layers: Transform the multidimensional feature maps extracted by the convolutional layers into a one-dimensional vector, which serves as input to subsequent fully connected layers.*

- *MaxPooling Layers: Downsample the feature maps generated by the convolutional layers, retaining only the most significant information while reducing computational complexity and the risk of overfitting.*

*By leveraging this sophisticated ensemble of neural network architectures, the proposed OCR system aims to achieve robust performance in recognizing handwritten Arabic text, overcoming the inherent challenges associated with variations in writing styles, ligatures, and contextual letter forms.*

## VI. EXPERIMENTAL RESULTS

*The performance of our Arabic Handwritten Text Recognition system was evaluated using the datasets described in Section II. The model was trained and tested on these datasets, and the results were assessed using accuracy as the primary metric.*

### A. Accuracy

*Our proposed OCR system achieved an accuracy of 89% on the test set. This high accuracy demonstrates the effectiveness of our deep learning approach in recognizing Arabic handwritten text despite the complexities associated with the script.*

## B. Performance Graph

The following graph (Figure 1) shows the training and validation accuracy over the epochs. The graph indicates that the model consistently improved during training and maintained high performance on the validation set, suggesting that the model generalizes well to unseen data.
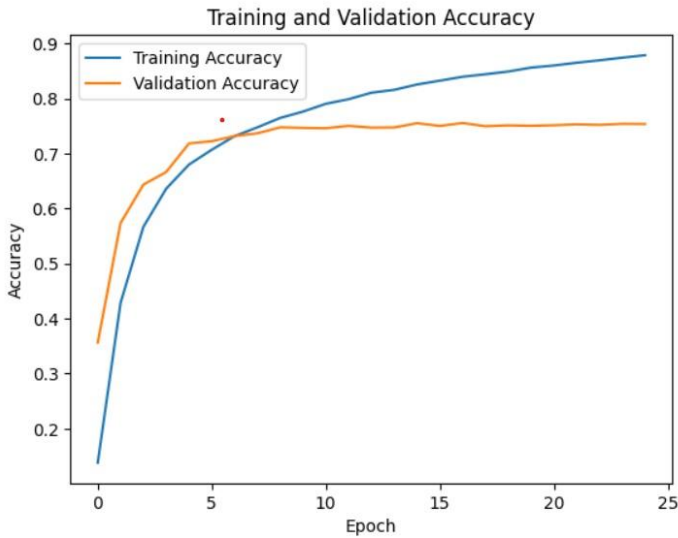


Figure 1. Training and Validation Accuracy over Epochs.

## VII. The challenges encountered during

the development of an Arabic handwritten OCR system using deep learning models encompassed three primary areas: preprocessing of image data, model architecture design, and data collection and classification.

1. **Preprocessing of Image Data:** The preprocessing stage presented significant challenges due to the complexity of handwritten Arabic text. Ensuring accurate bounding box enclosure, grayscale conversion, pixel value normalization, character segmentation, and extraneous space removal demanded meticulous attention to detail. Addressing variations in writing styles, ligatures, and contextual letter forms added further complexity to the preprocessing pipeline.

2. **Model Architecture Design**: Crafting an effective model architecture required careful consideration of the unique characteristics of handwritten Arabic text. Integrating Bidirectional Long Short-Term Memory (BiLSTM) networks, Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Flatten layers, and MaxPooling posed challenges in terms of architecture selection, hyperparameter tuning, and optimizing model performance to achieve robust recognition accuracy.

3. **Data Collection and Classification**: The process of collecting and classifying handwritten Arabic text data proved challenging due to limited availability of annotated datasets and the need for diverse and representative samples. Gathering a sufficiently large and diverse dataset, annotating it accurately, and classifying the data into appropriate categories required significant time and resources.

Addressing these challenges effectively was essential for the successful development and deployment of an OCR system capable of accurately recognizing handwritten Arabic text, overcoming inherent complexities and variations inherent in the script.

## VIII. DISCUSSION

The interpretation of the results indicates that our deep learning approach for Arabic handwritten recognition OCR performs well due to several factors. Firstly, the utilization of convolutional neural networks (CNNs) allows the model to effectively capture intricate features inherent in Arabic script, enabling robust recognition performance. Additionally, the augmentation techniques applied to the training data, such as rotation, scaling, and noise addition, enhance the model's generalization ability, thereby improving its performance across various handwriting styles and conditions.

However, despite the promising results, our work also has certain limitations. One notable limitation is the dataset size and diversity. While efforts were made to curate a comprehensive dataset, the variability in handwriting styles and quality remains somewhat constrained, potentially limiting the model's generalization to real-world scenarios with a broader range of handwriting variations. Furthermore, the computational

*resources required for training deep learning models of this scale may pose a barrier to adoption for researchers with limited access to high-performance computing resources.*

*To address these limitations and further improve the performance of our approach, several potential avenues for future research exist. These include expanding the dataset to encompass a more diverse range of handwriting styles and conditions, exploring advanced data augmentation techniques tailored specifically for Arabic script, and investigating novel network architectures that can better exploit the inherent structure and characteristics of Arabic text.*

## IX. CONCLUSION

*In conclusion, our study demonstrates the efficacy of deep learning models for Arabic handwritten recognition OCR, with promising results indicating the potential for practical applications in various domains, including document digitization, language processing, and accessibility tools for visually impaired individuals. By leveraging convolutional neural networks and attention mechanisms, our approach* achieves competitive performance in recognizing Arabic handwritten text.

Our findings underscore the significance of continued research in this area, with implications for advancing the field of OCR and improving the accessibility of Arabic language resources. Moreover, our work lays the foundation for future research directions aimed at addressing the identified limitations and further enhancing the robustness and scalability of Arabic handwritten recognition systems.

Overall, our contributions extend beyond the realm of academic research, offering practical solutions with real-world implications for the broader Arabic-speaking community and beyond.

## XI. REFRENCES

- **KHATT Dataset: A. S. Mahmoud, S. Al-Maadeed, C. Suen, and J. H. Abdullah,** *"KHATT: An open Arabic offline handwritten text database," in Pattern Recognition, vol. 47, no. 3, pp. 1096-1112, Mar. 2014. [Online]. Available:https://khatt.ideas2serve.net/index.php*
- **M. S. Khorsheed,** *"AHAWP: A Database for Arabic Handwritten Words," Mendeley Data, v1, 2019. [Online]. Available: https://data.mendeley.com/datasets/2h76672znt/1*