INTERNATIONAL JOURNAL OF RESEARCH IN EDUCATION AND

PSYCHOLOGY (IJREP) An International Peer Reviewed Journal http://ijrep.com/ SJIF Impact Factor 5.997 Vol.11, Issue 1 (January-March) 2025

RESEARCH ARTICLE





READING ARABIC TEXTS USING AI: CHALLENGES AND INNOVATIONS

Nashmiah Batel Hamed Alanazi

(Ministry of education, Al-Jouf education department, The fourth primary school in sakaka, Saudi Arabia.)

Email: ummfaisl5@gmail.com

https://doi.org/10.54513/IJREP.2025.11015

Abstract



Article Info: Article Received: 20-02-2025 Accepted on: 25-03-2025 Published online: 31-03-2025

The advancement of Artificial Intelligence (AI) in natural language processing (NLP) has significantly improved the ability to read and interpret Arabic texts. However, due to the complexity of Arabic script, including its rich morphology, diacritical system, and diverse dialects, AI models face unique challenges in accurately processing Arabic text. This paper explores the key difficulties in reading Arabic using AI, recent breakthroughs in Optical Character Recognition (OCR) and NLP, and future prospects for improving AI-driven Arabic text comprehension.

Keywords: Artificial Intelligence, Optical Character Recognition, Arabic texts, Reading, AI-driven Arabic text comprehension, difficulties, Improvements.

Author(s) retain the copyright of this article

Copyright © 2025 VEDA Publications

Author(s) agree that this article remains permanently open access under the terms of the Creative Commons Attribution License 4.0 International License (cc) EY

Nashmiah Batel Hamed Alanazi

1. Introduction

Reading Arabic text using AI is a growing area of interest in NLP and computer vision. From automated translation to voice assistants, AI has made substantial progress in understanding languages. However, Arabic presents distinct challenges due to its cursive script, varying orthographic styles, and contextual dependencies. This paper discusses the methods used in AI to process Arabic text, the hurdles that arise, and recent advancements in machine learning that help overcome these issues.

2. Challenges in Reading Arabic Text with AI

AI-based reading of Arabic texts encounters several technical and linguistic obstacles:

2.1. Complexity of Arabic Script

Arabic is written from right to left in a connected cursive script, meaning character shapes change depending on their position in a word. Unlike Latin-based scripts, there is no clear segmentation between letters, making it harder for AI to recognize individual characters.

2.2. Diacritics and Ambiguity

Arabic words often rely on diacritics (small marks above or below letters) to indicate vowel sounds, which affect meaning. However, diacritics are frequently omitted in modern texts, creating ambiguity. AI models must use contextual understanding to infer missing diacritics, a task requiring deep linguistic knowledge.

2.3. Dialectal Variation

Modern Standard Arabic (MSA) is used in formal writing, but spoken Arabic includes many regional dialects with significant differences in vocabulary, grammar, and pronunciation. AI models trained on MSA struggle with dialectal Arabic, requiring additional data and sophisticated training approaches.

2.4. Optical Character Recognition (OCR) Challenges

OCR is crucial for digitizing Arabic texts, yet Arabic OCR faces problems such as:

- Handwriting recognition difficulties
- Font variations and calligraphic styles
- Low-quality scanned texts with distortions

These factors complicate the accurate conversion of Arabic images into machine-readable text.

3. AI Approaches for Reading Arabic Texts

3.1. Machine Learning and Deep Learning Models

Deep learning techniques, particularly transformer-based architectures like BERT (Bidirectional Encoder Representations from Transformers) and its Arabic variants such as AraBERT and QARiB, have improved Arabic text processing (Alshammari,2023). These models use large-scale training data to enhance text recognition, translation, and summarization.

3.2. Optical Character Recognition (OCR) Solutions

Advancements in Arabic OCR involve deep neural networks, including convolutional neural networks (CNNs) and recurrent neural networks (RNNs), to improve character and word recognition. Google's Tesseract OCR and Arabic-specific OCR systems such as Sakhr and ABBYY FineReader have enhanced Arabic text digitization.(Hajj,2020)

3.3. Natural Language Processing (NLP) and Language Models

AI models trained on large Arabic datasets can handle tasks such as:

- Named Entity Recognition (NER) Identifying people, places, and organizations in Arabic text.
- Machine Translation Arabic-English translation with models like Google Translate and OpenAI's GPT.
- Speech-to-Text Processing AI-powered voice assistants (e.g., Siri and Google Assistant) recognize Arabic speech and convert it into written text.

4. Future Directions and Innovations

4.1. Improving Arabic NLP with Larger Datasets

Expanding high-quality Arabic datasets will enhance AI's ability to learn and adapt to different linguistic structures and dialects(Darwish.2020). Efforts such as the Arabic Gigaword corpus and OSCAR (Open Super-large Crawled Aggregated corpora) contribute to this goal.(Alshammari,2023)

4.2. Integration of Multimodal AI

Combining text, speech, and visual data will lead to more robust Arabic text comprehension. AI systems integrating OCR, speech recognition, and contextual understanding can improve performance in reading and translating Arabic texts.

4.3. Enhancing Zero-Shot and Few-Shot Learning

AI models trained in one dialect or language can generalize better to unseen Arabic dialects through few-shot or zero-shot learning techniques. This will reduce the need for massive labeled data.

5. Conclusion

This paper outlines the major aspects of AI-driven Arabic text reading and highlights potential future developments. Let me know if you need any modifications. AI has made remarkable progress in reading Arabic texts, but challenges remain due to the language's script complexity, diacritical ambiguity, and dialectal diversity. Advances in deep learning, NLP, and OCR continue to improve Arabic text recognition, making AI-powered Arabic reading applications more effective. Future research should focus on larger datasets, multimodal AI, and improved generalization techniques to further enhance AI's ability to read Arabic.

References

- Hajj, Hazem. Wissam Antoun, Fady Baly, Hazem
 •(2020). AraBERT: Transformer-based
 Model for Arabic Language Understanding. 28 Feb 2020 (v1), last revised 7 Mar 2021 (this version, v4)]. https://arxiv.org/abs/2003.00104
- Habash, N. Y. (2010). Introduction to Arabic natural language processing. (1 ed.) (Synthesis Lectures on Human Language Technologies).

https://doi.org/10.2200/S00277ED1V01Y201008HLT010

- Darwish, K., Habash, N., Abbas, M., Al-Khalifa, H.S., Al-Natsheh, H.T., El-Beltagy, S.R., Bouamor, H., Bouzoubaa, K., Cavalli-Sforza, V., El-Hajj, W., Jarrar, M., & Mubarak, H. (2020). A panoramic survey of natural language processing in the Arab world. Communications of the ACM, 64, 72 - 81.
- Alshammari, H., & EI-Sayed, A. (2023). AIRABIC: Arabic Dataset for Performance Evaluation of AI Detectors. 2023 International Conference on Machine Learning and Applications (ICMLA), 864-870.
