# A Transformer-based Approach for Arabic Offline Handwritten Text Recognition

*Saleh Momeni*[1] *Bagher BabaAli*[2]

[1,2] *School of Mathematics, Statistics, and Computer Science, University of Tehran, Tehran, Iran*
*\* E-mail: babaali@ut.ac.ir*

**Abstract:** Handwriting recognition is a challenging and critical problem in the fields of pattern recognition and machine learning, with applications spanning a wide range of domains. In this paper, we focus on the specific issue of recognizing offline Arabic handwritten text. Existing approaches typically utilize a combination of convolutional neural networks for image feature extraction and recurrent neural networks for temporal modeling, with connectionist temporal classification used for text generation. However, these methods suffer from a lack of parallelization due to the sequential nature of recurrent neural networks. Furthermore, these models cannot account for linguistic rules, necessitating the use of an external language model in the post-processing stage to boost accuracy. To overcome these issues, we introduce two alternative architectures, namely the Transformer Transducer and the standard sequence-to-sequence Transformer, and compare their performance in terms of accuracy and speed. Our approach can model language dependencies and relies only on the attention mechanism, thereby making it more parallelizable and less complex. We employ pre-trained Transformers for both image understanding and language modeling. Our evaluation on the Arabic KHATT dataset demonstrates that our proposed method outperforms the current state-of-the-art approaches for recognizing offline Arabic handwritten text.

## 1 Introduction

Optical Character Recognition (OCR) is a technology that enables automatic recognition of printed or handwritten text in scanned documents, converting it into machine-readable text. OCR is widely used in applications such as document digitization, search and translation of documents, and autonomous vehicles. OCR tasks typically involve two sub-problems: text detection and text recognition. The text detection algorithms aim to identify text samples in input images, while the text recognition module attempts to comprehend the patch contents and transcribe them into natural language tokens.

Despite significant advances in OCR technology for English text in recent years, other scripts such as Arabic remain relatively unexplored. The Arabic script is widely considered to be one of the most challenging scripts for OCR due to its cursive nature, which presents numerous obstacles, including difficulties with segmentation, character overlapping, and context dependency [27]. While previous research in Arabic handwriting recognition has largely focused on individual letters or isolated words [3, 16, 29], this approach may not be optimal for dense text, where it can be challenging to identify separate characters or words. To address this issue, this study proposes an end-to-end line-based OCR pipeline to extract Arabic handwritten text and determine the optimal Handwritten Text Recognition (HTR) model. The proposed pipeline aims to overcome the challenges associated with Arabic script by processing text at the line level. This approach considers the overall context of the handwritten text, which is particularly important in Arabic, where the shape of a character can vary depending on its position within a word. By processing text at the line level, the proposed pipeline can better handle the complexities of Arabic handwriting and improve recognition accuracy.

Current text recognition methods commonly use a Convolutional Neural Network (CNN) for feature extraction, followed by a Recurrent Neural Network (RNN) for capturing contextual information. However, RNNs are limited in terms of parallel processing during recognition due to their sequential nature. To overcome this limitation, we propose a Transformer-based pipeline that relies solely on the attention mechanism and does not use recurrence, allowing for more efficient and parallel image recognition.

Our model includes a separate text detection module, similar to ViT [14], which is responsible for identifying text instances in text line images. The input image is adjusted to a specific height while maintaining its aspect ratio, then broken into small patches of the same size and passed to the encoder as input. The encoder uses stacked self-attention layers to convert the input sequence into feature representations. The decoder then uses these encoded sequence features to attempt to decode the text content. We explore two decoding architectures: the Transducer decoder [18] and the cross-attention decoder [6].

Connectionist Temporal Classification (CTC) [19] is a widely used architecture in handwriting recognition systems, which assumes a monotonic alignment between input and output, simplifying the model. However, the output tokens are assumed to be independent in this architecture, which severely limits the ability of language modeling Consequently, CTC decoders are often paired with an external language model. The integration of an external language model has been observed to enhance the accuracy of the model considerably [13], which implies that decoding should incorporate language rules. The Transducer architecture elegantly addresses this limitation of CTC by incorporating a joiner network and considering both visual and language features when predicting the output. Initially, the Transducer architecture was used in conjunction with RNNs and was therefore referred to as a Recurrent Neural Network Transducer (RNN-T) [18]. However, this approach is not limited to RNNs. Recent research in text recognition has demonstrated promising results using transformer architecture [13, 22]. Motivated by these findings, we replaced the conventional RNN-T with the Transformer Transducer. This approach aligns with the current trend of substituting recurrent networks with attention modules. To the best of our knowledge, the use of Transformer Transducer architecture for text recognition has not been previously investigated. For our second architecture, we incorporated a standard Transformer decoder [36]. The Transformer decoder is the preferred decoder for many sequence-to-sequence (Seq2Seq) tasks, such as machine translation, and is suitable for sequentially processing images while also possessing the power of language modeling.

Deep neural networks are known to require vast amounts of training data to attain a generalized model. In order to tackle this issue, we have developed a synthetic dataset comprising images

generated from Arabic texts along with their corresponding labels. Our models are pre-trained on this data prior to fine-tuning on the original dataset. Additionally, we initialize both the encoder and decoder using pre-trained Transformer models to leverage publicly available resources. Our experiments demonstrate that our approach surpasses the current state-of-the-art in the HTR task without the requirement for complex pre/post-processing steps. In summary, our work makes the following contributions:

1. We provide a comprehensive comparison of two distinct Transformer-based model design options to determine the best end-to-end architecture for the HTR task.

2. We investigate the performance of the Transformer Transducer in the HTR task. To the best of our knowledge, this architecture has not been previously employed for text recognition.

3. Our proposed model achieves a new state-of-the-art result on the benchmark dataset without the use of complicated pre/post-processing steps.

The remainder of this paper is organized as follows: we review related works in section 2 and introduce our proposed architectures in section 3. In section 4, we elaborate on synthetic data generation, implementation details, and the training scheme. Section 5 provides an analysis and comparison of the two proposed pipelines, as well as a comparative evaluation of our approach with other existing methods. Finally, we conclude our work in section 6.

## 2 Related Works

OCR is a challenging problem in which the input image is typically represented as a sequence that must be mapped to corresponding target labels. The challenge lies in the fact that the alignment between the input and output is often unknown. To address this challenge, CTC has emerged as a popular technique that enables training without the need for exact alignment information. Many text recognition frameworks rely on CTC, such as the end-to-end text recognition architecture proposed by Shi et al. [34], which utilizes a CNN for feature extraction and a RNN for learning spatial dependencies. Another approach, introduced by Bluche and Messina [10], employs gated convolution to enable context-sensitive feature extraction. Furthermore, previous studies have demonstrated the effectiveness of multidimensional RNNs in combination with a CTC decoder [9, 30, 37].

While the CTC architecture has been successful in OCR, it still has limitations in terms of language modeling, which can be improved by using an external language model as a post-processing step. To overcome this limitation, attention mechanisms have gained popularity in text recognition tasks [11, 17, 26, 32]. Chowdhury and Vig [11] introduced an RNN encoder-decoder model for HTR that utilized a CNN backbone. Michael et al. [26] further explored different attention mechanisms and positional encodings for improving text recognition accuracy.
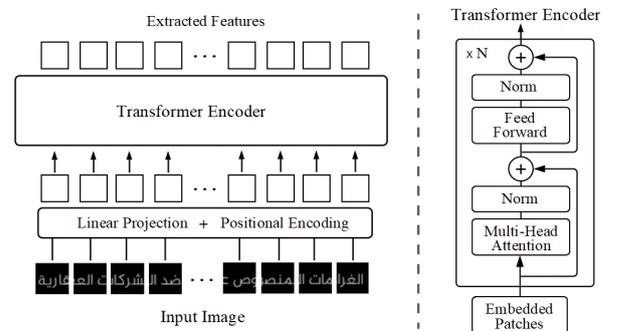
In recent years, the Transformer architecture [36] has gained significant traction in the domain of text recognition systems. Sheng et al. [33] were pioneers in utilizing the Transformer architecture for text recognition. They employed a modality-transform block to convert a 2D image to a 1D sequence, followed by a Transformer encoder-decoder for text recognition. The Transformer architecture's ability to model intricate language rules and facilitate greater parallelization has led to its widespread adoption, replacing traditional RNNs in many text recognition systems [8, 20, 21]. Baek et al. [5] proposed a four-stage framework to compare recognizer models and addressed inconsistencies between training and evaluation datasets. Diaz

et al. [13] conducted research on developing a universal architecture for handwritten and scene text recognition and compared various encoder/decoder combinations. The use of image Transformers [7, 23, 35] for feature extraction in text recognition tasks has also gained popularity, with Atienza [4] and Li et al. [22] among the early adopters of ViT-style Transformers, replacing the conventional CNN backbone. While Atienza [4] combined the vision Transformer with a CTC decoder, Li et al. [22] employed the Seq2Seq paradigm. The cross-attention architecture employed in our study is similar to the latter work, with the exception that we decompose each input image into a 1D sequence of patches, which allows for a more efficient handling of variable-length inputs.

In this paper, we also explore the feasibility of replacing the conventional CTC decoder with a Transducer decoder in OCR models. The use of a Transducer decoder facilitates end-to-end training of the OCR model, thereby obviating the need for an external language model. The RNN-T loss, first proposed by Graves [18], is central to this end-to-end training. Although it has not been widely used in comparison, recent studies in the field of speech recognition have yielded promising results with the use of RNN-T loss, as demonstrated by Zhang et al. [39]. Inspired by this observation, we have employed the Transducer architecture in our HTR model.

## 3 Proposed Model Architecture

In the context of text recognition systems, two fundamental components are typically employed: a visual feature encoder that extracts salient features from textual images, and a decoder that transcribes the final output based on the extracted features. The present study is concerned with the problem of text transcription and, in particular, explores two distinct decoder choices while utilizing a fixed text detection module.



**Fig. 1**: The input image is divided into fixed-size patches. The patches are then embedded, provided with positional encodings, and fed to a standard Transformer encoder.

### 3.1 Encoder

As illustrated in Fig. 1, the encoder takes a grayscale image as input and resizes it to a fixed height while maintaining the aspect ratio. Since the transformer architecture cannot process images directly, the input image is partitioned into vertical patches. These patches are then flattened and processed through a linear layer, resulting in a sequence of vectors. To preserve positional information, absolute positional embeddings are added to the vectors, which are learned during the training process. The sequence of vectors is then passed through a sequence of identical transformer layers, each comprising a self-attention layer and a feed-forward layer. Following each of these layers, a normalizing layer is applied in conjunction with residual connections. The attention mechanism maps a set of queries, keys, and values to the output, which is computed as a weighted average

of the values. The weights are calculated based on the similarity between the key and its corresponding query:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V \qquad (1)$$

Where Q, K, and V are the matrices of the queries, keys, and values respectively, and d is the dimension of the input vector. The self-attention layer employs linear projection to derive the keys, queries, and values from the input sequence. In contrast to a single attention function, each attention module leverages multiple heads with distinct parameters, thereby enabling the model to gather information from distinct representation subspaces:
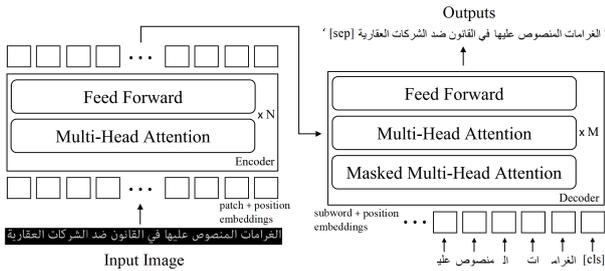
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$
$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \qquad (2)$$

The attention mechanism facilitates each encoder layer to encapsulate the information from the entire input sequence in a feature vector, bypassing the necessity for any recurrence.
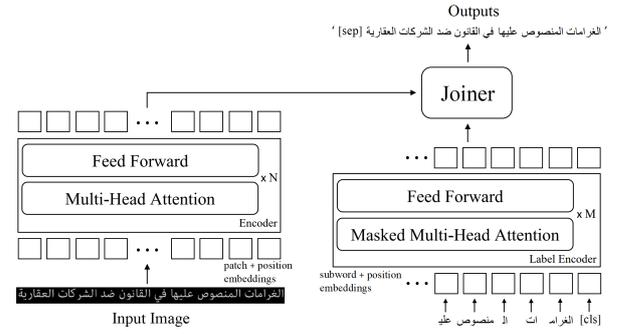
## 3.2 Decoder

The decoder component of the system takes in the encoded feature sequence and aims to convert the visual cues into corresponding natural language tokens. To achieve this goal, we investigate two distinct decoding methods. The proposed architectures for these decoding methods are illustrated in Fig. 2 and Fig. 3.



**Fig. 2**: Standard Transformer encoder-decoder architecture, where the decoder attends to the visual features using cross-attention.

*3.2.1 Transformer with Cross-Attention:* The Transformer decoder takes the encoded feature sequence and generates the output text sequence iteratively, utilizing the previously generated labels as additional input. Each input label is mapped to an embedding vector and positional embeddings are added to these vectors. The resulting embeddings are then passed through several identical transformer decoder layers, each comprising of self-attention and feed-forward layers, as well as a cross-attention layer with separate attention heads to distribute the weights on encoded features. In the cross-attention layer, the queries come from the decoder input, while the keys and values are from the encoder outputs. Attention masking is used during training to ensure that the model does not utilize the information of the tokens that are not yet predicted. To transform the decoder output into a probability distribution over the vocabulary, a linear layer with softmax activation is applied. Finally, the model is trained using cross-entropy loss.

*3.2.2 Transformer Transducer:* The Transducer decoder generates a probability distribution over the vocabulary at each time step by receiving the generated feature sequence as input and moving forward on it. This vocabulary incorporates a blank symbol that denotes the absence of output at a given time step. Unlike the CTC, this probability distribution is conditioned on



**Fig. 3**: Transformer Transducer architecture, incorporating a joiner network that combines information from the visual encoder and the label encoder to generate the final output.

the previously generated tokens. The Transducer decoder is composed of a label encoder and a joiner network. The label encoder resembles a conventional language model that leverages previous outputs to produce features to predict the subsequent label. The label encoder was initially implemented using RNNs in the model proposed by Graves [18]. However, in our study, we replaced the RNN with a Transformer architecture. Unlike the cross-attention architecture, the label encoder does not attend to the visual features. Instead, it encodes input labels into a sequence of features that contain semantic information. The last feature vector's information is then used to predict the next label, which is fed to the joiner network along with the visual features. The joiner network adds the predictor vector to the visual features and passes them through a linear layer with softmax activation, generating the probability distribution over the label space. The model can be trained end-to-end using the RNN-T loss. The Transducer model defines the probability of the target sequence $y$ given the input sequence $x$:

$$P(y \mid x) = \sum_{z \in \mathcal{Z}(x,y)} P(z \mid x) \qquad (3)$$

$$P(z \mid x) = \prod_i P(z_i \mid x, z_1, \dots, z_{i-1}) \qquad (4)$$

Where $z$ is an alignment between the input and the target sequence, and $\mathcal{Z}(x, y)$ is the set of all possible alignments between the two sequences. When utilizing the Transducer decoder, the alignment between the input and target sequences can vary depending on the generated tokens at each time step, leading to multiple possible alignments. As a result, the probability of the target sequence is computed by adding the probabilities of all possible alignments. The model's loss is then calculated as the sum of negative log probabilities over the training examples.

$$\text{loss} = -\sum_i \log P(y_i \mid x_i) \qquad (5)$$

However, directly computing the RNN-T loss by summing the probabilities of all possible alignments would be computationally inefficient. To address this issue, we use the forward-backward algorithm described in [18]. This algorithm takes advantage of the overlaps between alignments and stores intermediate computations that can be reused, resulting in improved computational efficiency.

## 3.3 Model Initialization

In order to enhance the accuracy of our model, we adopt an approach of initializing both the encoder and decoder with publicly available pre-trained models that have been trained on large-scale datasets. Specifically, we have employed the DeiT [35] model as the encoder and the asafaya-BERT [31] model as

the decoder. The DeiT model is a data-efficient image Transformer that is trained on the ImageNet dataset, and it introduces a distillation token to learn effectively from a teacher, resulting in competitive results on ImageNet without requiring external data. On the other hand, the asafaya-BERT model is the Arabic version of the BERT [12] model.
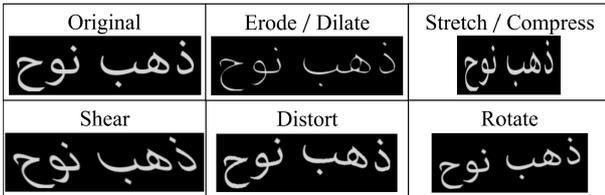
Despite the availability of pre-trained models, we have faced some challenges in initializing our model due to differences in the structure of our models and the pre-trained ones. For instance, the patch embedding layer in our encoder is different from the DeiT model, as the latter is designed to accept RGB images with a resolution of 384x384, while our encoder works with grayscale images of variable length. Moreover, the attention layers in our attention-based decoder that enable the encoder-decoder interactions are not present in the Arabic BERT. Furthermore, our Transducer decoder includes an additional joiner module. To overcome these differences, we have initialized all the parameters mentioned above randomly.

## 4    Experimental Setup

In this section, we begin by outlining the method employed for generating synthetic data to train the model. We then expound on the benchmark dataset, implementation details, and evaluation criteria used to assess the model's performance.

### 4.1    Synthetic Data For Pre-Training

Deep learning models require a substantial amount of training data to generate a model that can generalize well. To this end, we generated a synthetic dataset comprising 500,000 printed Arabic text-line images, along with their corresponding ground truth, using various open-source fonts. We used text from [15] to render the synthetic data. To improve the diversity and realism of the synthetic images, we applied various image processing techniques, such as shearing, rotation, distortion, erosion, and compression. Our models were trained on the synthetic dataset and fine-tuned on the original dataset to evaluate their performance.
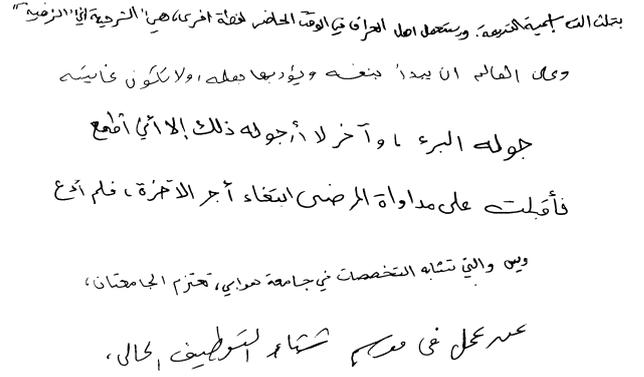


**Fig. 4**: Illustration of image processing techniques used to produce synthetic data.

### 4.2    Dataset

The evaluation of our proposed method is conducted using the KHATT [25] dataset, which is publicly available and widely used in the field of Arabic handwriting recognition. The KHATT dataset comprises 1000 handwritten Arabic forms written by various writers from different backgrounds, including age, education levels, and countries. The dataset includes 2000 paragraph images with similar text and 2000 paragraph images with unique text, along with their corresponding ground truth. For the purposes of our experiments, we focus on the unique-text paragraphs, which consist of 6742 text lines after segmentation.

The KHATT dataset contains numerous text-line images that exhibit high skew and extraneous white regions. To prepare these images for height normalization, it is imperative to first remove any excess whitespace. In order to eliminate these additional spaces, we rotate the text-line images, using the same amount of rotation as their skew, but in the opposite direction. This allows



**Fig. 5**: Example text line images in the KHATT dataset which contains handwritten text with various styles.

us to search the image for black pixels from top to bottom, and discard any extra empty spaces before height normalization.

### 4.3    Implementation Details

In our approach, we utilize gray-scale images for training and evaluation. We normalize the input image to have a height of 64 pixels while preserving the aspect ratio. The encoder's patch size is set to 64x12, and we pad the inputs to a fixed length of 2400 to enable batch processing. We provide a summary of the network configuration in Table 1. We compare the accuracy by stacking different numbers of encoder and decoder layers. To account for differences in the hidden size of the encoder and decoder, we append a linear layer to the end of the encoder to project feature vectors to the decoder's dimension. We utilize word piece tokenization with a vocabulary size of 32000 to convert text lines to subword tokens.

**Table 1**  Transformer parameter setup.

| Parameter | Encoder | Decoder |
|---|---|---|
| Hidden Size | 768 | 512 |
| Intermediate Size | 3072 | 2048 |
| Number of Attention Heads | 12 | 8 |
| Dropout Ratio | 0.1 | 0.1 |

The models are implemented in PyTorch [28], and we initialize them using the $\text{DeiT}_{base}$ and asafaya-BERT$_{medium}$ models from the Hugging Face [38] repository. All the experiments in the paper are conducted on a Tesla T4 GPU with a 16 GB memory. The Adam optimizer is utilized with a batch size of 16 for training. During the first 5K steps, the learning rate linearly ramps up from 1e-7 to 1e-5, and then remains constant for 20K steps. Following that, it exponentially decays to 1e-6 over the next 5K steps. We fine-tune the models on the KHATT dataset for 10 epochs before evaluating them.

### 4.4    Evaluation Metric

The Character Error Rate (CER) is a widely used evaluation metric for OCR systems. It is defined as the Levenshtein distance between the predicted text and the ground truth, normalized by the number of characters in the ground truth. The Levenshtein distance between two strings is a measure of the difference between them, and it is calculated as the minimum number of character substitutions, deletions, and insertions required to transform one string into the other. More formally, the CER can be calculated as:

$$CER = \frac{s + i + d}{n} \tag{6}$$

4

where $s$ is the number of substitutions, $i$ is the number of insertions, $d$ is the number of deletions, and $n$ is the total number of characters in the ground truth.

# 5    Results

## 5.1    Ablation Study

In this section, we conduct a thorough ablation study of our models to investigate the impact of each architectural component on the final performance. We report on key metrics, including accuracy, latency, and the number of parameters to provide insight into the trade-offs associated with our approach.

**Table 2**  CER on the KHATT dataset for models trained from scratch with different numbers of encoder and decoder layers.

| Model | Number of Hidden Layers | Parameter Size | CER(%) |
|---|---|---|---|
| Transformer-T | 8 encoder / 12 decoder | 129.0 M | 22.92 |
| | 10 encoder / 10 decoder | 136.9 M | 22.20 |
| | 12 encoder / 8 decoder | 144.7 M | 21.78 |
| Transformer with Cross-Attention | 8 encoder / 12 decoder | 141.6 M | 20.71 |
| | 10 encoder / 10 decoder | 147.4 M | 20.28 |
| | 12 encoder / 8 decoder | 153.1 M | 19.99 |

**Number of encoder and decoder layers analysis:** We begin by investigating the effect of the number of encoder and decoder layers on the model's accuracy. To do so, we compared three different configurations, each with the same total number of layers. All models were trained from scratch, and beam search was utilized to obtain the final predictions. The results are presented in Table 2. Our analysis indicated that the model with a higher number of encoder layers outperformed the others in terms of accuracy, suggesting that image understanding is more complex than language modeling in the HTR task. Based on these findings, we selected a configuration of 12 encoder layers and 8 decoder layers for our primary models.

**Table 3**  Comparison of Transformer Transducer and cross-attention based Transformer in terms of accuracy and speed on KHATT test set.

| Model | Parameter Size | Beam Width | Latency (ms) | CER(%) |
|---|---|---|---|---|
| Transformer-T | 144.7 M | K = 1 | 94.41 | 21.15 |
| | | K = 3 | 118.7 | 19.91 |
| | | K = 5 | 132.8 | 19.76 |
| Transformer with Cross-Attention | 153.1 M | K = 1 | 157.4 | 19.74 |
| | | K = 3 | 202.9 | 18.63 |
| | | K = 5 | 223.7 | 18.45 |

**Architecture comparison:** A comparison was conducted between the Transformer Transducer and the cross-attention Transformer architectures, where both models had identical layers, except for the cross-attention layers present in the Transformer decoder, resulting in a substantial increase of 8.4 million parameters in the cross-attention Transformer compared to the Transducer. It is noteworthy that the joiner network in the Transducer is equivalent to the linear layer at the end of the Transformer decoder, as the label encodings are added to the visual features instead of being concatenated. Evaluation was performed on the KHATT dataset, and the results are presented in Table 3. The Transformer Transducer achieved slightly lower accuracy than the cross-attention Transformer, but significantly outperformed the latter in terms of latency. The reported latencies represent the average time required by each model to recognize a single text line image in the KHATT test set. It was observed that the Transformer encoder took approximately 15.6 milliseconds to produce the visual features, and most of the latency was associated with the autoregressive nature of the decoders. In general, both architectures were competitive, but there were trade-offs. The cross-attention Transformer had a
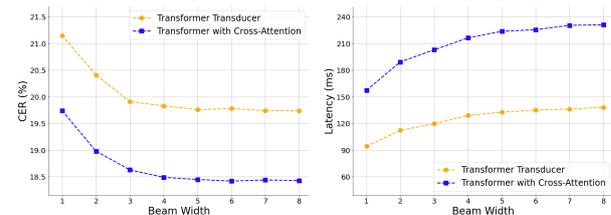
better CER, which was 1.31% lower than the Transformer Transducer. However, the Transformer Transducer was about 68% faster than its cross-attention counterpart.

**Table 4**  Impact of different model components on accuracy: in order to determine the significance of individual components, CER was assessed multiple times by toggling these components on and off.

| Model | Initialization | Beam Search | Augmentation | CER(%) |
|---|---|---|---|---|
| Transformer-T | ✗ | ✓ | ✓ | 21.78 |
| | ✓ | ✗ | ✓ | 21.15 |
| | ✓ | ✓ | ✗ | 20.69 |
| | ✓ | ✓ | ✓ | 19.76 |
| Transformer with Cross-Attention | ✗ | ✓ | ✓ | 19.99 |
| | ✓ | ✗ | ✓ | 19.74 |
| | ✓ | ✓ | ✗ | 19.30 |
| | ✓ | ✓ | ✓ | 18.45 |

**Importance of model initialization:** Pre-trained vision and language Transformers were employed as the initial point for our models. Pre-trained weight initialization allows the model to adapt quickly to new tasks with better performance. The significance of model initialization was assessed by training the model from scratch. When trained from scratch, a reduction of 1.54% in the cross-attention Transformer was observed, with a more notable decline of 2.02% in the Transformer Transducer. The results concluded that utilizing the pre-existing knowledge from pre-trained models significantly improved the models' accuracy.

**Importance of data augmentation:** Data augmentation is a widely used technique for reducing overfitting in machine learning models. In this work, various data augmentation techniques, such as shearing, rotation, blurring, and noise, were applied during the fine-tuning process. The primary objective of data augmentation was to generate slightly modified copies of the existing data, which could help to improve the generalization capability of the model. The experimental results illustrate that data augmentation plays a critical role in achieving good performance in our models. Specifically, the performance of the Transformer Transducer and cross-attention Transformer models decreased by 0.93% CER and 0.85% CER, respectively, when data augmentation was not applied. These findings emphasize the importance of data augmentation in our work and highlight the need for its inclusion in related applications.



**Fig. 6**: The impact of beam search on accuracy and latency was assessed by gradually increasing the beam width while monitoring changes in CER.

**Influence of beam search:** we investigated the impact of beam search on the performance of our model. Beam search is an optimized version of the breadth-first search algorithm, which reduces the required memory by using a parameter called beam width. The beam width indicates the number of vertices stored in each step of the algorithm. The effect of beam width on CER and latency was evaluated, and the results were presented in Fig. 6. We observed a 1.39% CER improvement in the Transformer Transducer and a 1.29% improvement in the cross-attention Transformer by using beam search with a width of K=5, while the decoding time increased by 41%. Increasing the beam width further did not yield any significant improvement. Our findings revealed that beam search can significantly increase recognition accuracy, albeit at the cost of additional decoding time.

**Table 5** The results of the evaluations conducted on our models and existing approaches, utilizing the KHATT dataset.

| Model | Architecture | Recognition Level | External LM | CER(%) |
|---|---|---|---|---|
| Mahmoud et al. [24] | HMM | Character | No | 53.87 |
| Ahmad and Fink [1] | HMM | Character | No | 41.21 |
| Ahmad et al. [2] | MDLSTM + CTC | Character | No | 24.20 |
| Ours | Transformer-T | Subword | No | 19.76 |
|  | Transformer with Cross-Attention | Subword | No | 18.45 |

### 5.2 Comparative Evaluation

The comparison of our proposed models with existing approaches on the KHATT dataset is presented in Table 5. The results demonstrate that our model outperforms the previous methods and achieves a new state-of-the-art on the KHATT dataset without relying on any complex pre/post processing steps.

## 6 Conclusion

This study investigated two distinct end-to-end architectures for recognizing Arabic handwritten text lines: the Transformer Transducer and the standard Transformer architecture employing cross-attention. We utilized pre-trained text and image transformers to establish our models. Our proposed method is non-recurrent and open-vocabulary, and both architectures can model complex language dependencies without requiring an external language model. We found that both models are highly competitive, with the cross-attention Transformer achieving superior accuracy and the Transformer Transducer demonstrating faster processing speed. The experimental results indicate that our model achieves a new state-of-the-art performance on the KHATT benchmark.

## 7 References

1 Ahmad, I. and G. A. Fink (2019). Handwritten arabic text recognition using multi-stage sub-core-shape hmms. *International Journal on Document Analysis and Recognition (IJDAR) 22*(3), 329–349.

2 Ahmad, R., S. Naz, M. Z. Afzal, S. F. Rashid, M. Liwicki, and A. Dengel (2017). Khatt: A deep learning benchmark on arabic script. In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, Volume 7, pp. 10–14. IEEE.

3 Altwaijry, N. and I. Al-Turaiki (2021). Arabic handwriting recognition system using convolutional neural network. *Neural Computing and Applications 33*(7), 2249–2261.

4 Atienza, R. (2021). Vision transformer for fast and efficient scene text recognition. In *International Conference on Document Analysis and Recognition*, pp. 319–334. Springer.

5 Baek, J., G. Kim, J. Lee, S. Park, D. Han, S. Yun, S. J. Oh, and H. Lee (2019). What is wrong with scene text recognition model comparisons? dataset and model analysis. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4715–4723.

6 Bahdanau, D., K. Cho, and Y. Bengio (2014). Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

7 Bao, H., L. Dong, and F. Wei (2021). Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*.

8 Bleeker, M. and M. de Rijke (2019). Bidirectional scene text recognition with a single decoder. *arXiv preprint arXiv:1912.03656*.

9 Bluche, T. (2016). Joint line segmentation and transcription for end-to-end handwritten paragraph recognition. *Advances in neural information processing systems 29*.

10 Bluche, T. and R. Messina (2017). Gated convolutional recurrent neural networks for multilingual handwriting recognition. In *2017 14th IAPR international conference on document analysis and recognition (ICDAR)*, Volume 1, pp. 646–651. IEEE.

11 Chowdhury, A. and L. Vig (2018). An efficient end-to-end neural model for handwritten text recognition. *arXiv preprint arXiv:1807.07965*.

12 Devlin, J., M.-W. Chang, K. Lee, and K. Toutanova (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

13 Diaz, D. H., S. Qin, R. Ingle, Y. Fujii, and A. Bissacco (2021). Rethinking text line recognition models. *arXiv preprint arXiv:2104.07787*.

14 Dosovitskiy, A., L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.

15 El-Khair, I. A. (2016). 1.5 billion words arabic corpus. *arXiv preprint arXiv:1611.04033*.

16 El-Sawy, A., M. Loey, and H. El-Bakry (2017). Arabic handwritten characters recognition using convolutional neural network. *WSEAS Transactions on Computer Research 5*(1), 11–19.

17 Gao, Y., Y. Chen, J. Wang, and H. Lu (2017). Reading scene text with attention convolutional sequence modeling. *arXiv preprint arXiv:1709.04303*.

18 Graves, A. (2012). Sequence transduction with recurrent neural networks. *arXiv preprint arXiv:1211.3711*.

19 Graves, A., S. Fernández, F. Gomez, and J. Schmidhuber (2006). Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd international conference on Machine learning*, pp. 369–376.

20 Kang, L., P. Riba, M. Rusiñol, A. Fornés, and M. Villegas (2022). Pay attention to what you read: non-recurrent handwritten text-line recognition. *Pattern Recognition 129*, 108766.

21 Lee, J., S. Park, J. Baek, S. J. Oh, S. Kim, and H. Lee (2020). On recognizing texts of arbitrary shapes with 2d self-attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 546–547.

22 Li, M., T. Lv, L. Cui, Y. Lu, D. Florencio, C. Zhang, Z. Li, and F. Wei (2021). Trocr: Transformer-based optical character recognition with pre-trained models. *arXiv preprint arXiv:2109.10282*.

23 Liu, Z., Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo (2021). Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10012–10022.

24 Mahmoud, S. A., I. Ahmad, W. G. Al-Khatib, M. Alshayeb, M. T. Parvez, V. Märgner, and G. A. Fink (2014). Khatt: An open arabic offline handwritten text database. *Pattern Recognition 47*(3), 1096–1112.

25 Mahmoud, S. A., I. Ahmad, M. Alshayeb, W. G. Al-Khatib, M. T. Parvez, G. A. Fink, V. Märgner, and H. El Abed (2012). Khatt: Arabic offline handwritten text database. In *2012 International conference on frontiers in handwriting recognition*, pp. 449–454. IEEE.

26 Michael, J., R. Labahn, T. Grüning, and J. Zöllner (2019). Evaluating sequence-to-sequence models for handwritten text recognition. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pp. 1286–1293. IEEE.

27 Mostafa, A., O. Mohamed, A. Ashraf, A. Elbehery, S. Jamal, G. Khoriba, and A. S. Ghoneim (2021). Ocformer: A transformer-based model for arabic handwritten text recognition. In *2021 International Mobile, Intelligent, and Ubiquitous Computing Conference (MIUCC)*, pp. 182–186. IEEE.

28 Paszke, A., S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer (2017). Automatic differentiation in pytorch.

29 Pechwitz, M., S. S. Maddouri, V. Märgner, N. Ellouze, H. Amiri, et al. (2002). Ifn/enit-database of handwritten arabic words. In *Proc. of CIFED*, Volume 2, pp. 127–136. Citeseer.

30 Pham, V., T. Bluche, C. Kermorvant, and J. Louradour (2014). Dropout improves recurrent neural networks for handwriting recognition. In *2014 14th international conference on frontiers in handwriting recognition*, pp. 285–290. IEEE.

31 Safaya, A., M. Abdullatif, and D. Yuret (2020). Kuisail at semeval-2020 task 12: Bert-cnn for offensive speech identification in social media. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pp. 2054–2059.

32 Sajid, U., M. Chow, J. Zhang, T. Kim, and G. Wang (2021). Parallel scale-wise attention network for effective scene text recognition. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8. IEEE.

33 Sheng, F., Z. Chen, and B. Xu (2019). Nrtr: A no-recurrence sequence-to-sequence model for scene text recognition. In *2019 International conference on document analysis and recognition (ICDAR)*, pp. 781–786. IEEE.

34 Shi, B., X. Bai, and C. Yao (2016). An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE transactions on pattern analysis and machine intelligence 39*(11), 2298–2304.

35 Touvron, H., M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou (2021). Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pp. 10347–10357. PMLR.

36 Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin (2017). Attention is all you need. *Advances in neural information processing systems 30*.

37 Voigtlaender, P., P. Doetsch, and H. Ney (2016). Handwriting recognition with large multidimensional long short-term memory recurrent neural networks. In *2016 15th international conference on frontiers in handwriting recognition (ICFHR)*, pp. 228–233. IEEE.

38 Wolf, T., L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, et al. (2020). Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pp. 38–45.

39 Zhang, Q., H. Lu, H. Sak, A. Tripathi, E. McDermott, S. Koo, and S. Kumar (2020). Transformer transducer: A streamable speech recognition model with transformer encoders and rnn-t loss. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7829–7833. IEEE.