HATFORMER: HISTORIC HANDWRITTEN ARABIC TEXT LINE RECOGNITION WITH TRANSFORMERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Arabic handwritten text recognition (HTR) is challenging, especially for historical texts, due to diverse writing styles and the intrinsic features of Arabic script. Additionally, Arabic handwriting datasets are smaller compared to English ones, making it difficult to train generalizable Arabic HTR models. To address these challenges, we propose HATFORMER, a transformer-based encoder-decoder architecture that builds on a state-of-the-art English HTR model. By leveraging the transformer's attention mechanism, HATFORMER captures spatial contextual information to address the intrinsic challenges of Arabic script through differentiating cursive characters, decomposing visual representations, and identifying diacritics. Our customization to historical handwritten Arabic includes an image processor for effective ViT information preprocessing, a text tokenizer for compact Arabic text representation, and a training pipeline that accounts for a limited amount of historic Arabic handwriting data. HATFORMER achieves a character error rate (CER) of 8.6% on the largest public historical handwritten Arabic dataset, with a 51% improvement over the best baseline in the literature. HATFORMER also attains a comparable CER of 4.2% on the largest private nonhistorical dataset. Our work demonstrates the feasibility of adapting an English HTR method to a low-resource language with complex, language-specific challenges, contributing to advancements in document digitization, information retrieval, and cultural preservation.

033 034

043

004

006

007 008 009

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

1 INTRODUCTION

Global archives contain hundreds of millions of manuscript pages written in the Arabic alphabet, primarily from the 19th and early 20th centuries, with around 25 million images from the Middle East and North Africa alone. For historians, the laborious process of sifting through these pages for relevant data is impractical due to time and resource constraints. Existing handwritten text recognition (HTR) systems for non-historical Arabic texts fail to effectively render these historical documents into a searchable format. Developing a dedicated HTR system for historical Arabic manuscripts would revolutionize digital humanities, enabling rapid data search and retrieval while facilitating the creation of advanced large language models for research, thus opening new avenues for historical and humanities scholarship.

This paper introduces HATFORMER, a transformer-based historical Arabic HTR system that lever-044 ages self-attention mechanisms to capture long-range dependencies, outperforming traditional HTR methods for complex scripts like Arabic. HTR systems such as Shi et al. (2016) have traditionally 046 relied on convolutional neural networks (CNNs) (LeCun et al., 1989) for feature extraction and re-047 current neural networks (RNNs) (Rumelhart et al., 1986) for text generation. However, RNN-based 048 methods often struggle to capture long-range dependencies, which are more crucial for handling 049 Arabic scripts than for English scripts. Recently, transformer (Vaswani et al., 2017) methods have shown to be promising for modern and historical English HTR tasks, with Li et al. (2023), Fujitake 051 (2024), and Parres & Paredes (2023) achieving state-of-the-art character error rates (CER) of 2.9%, 2.4%, and 2.7%, respectively. HATFORMER builds on the success of pretrained vision and text 052 transformers in HTR, introducing key adaptations to handle the intrinsic challenges of Arabic for more accurate recognition of historical text.

054 We will show through experimental verification that the inductive bias of the transformer's attention 055 mechanism effectively addresses the following three intrinsic challenges (Najam & Faizullah, 2023; 056 Faizullah et al., 2023) of Arabic script absent in English. First, Arabic is required to be written in cur-057 sive, making characters visually harder to distinguish. The attention mechanism allows the model to 058 better differentiate between connected characters. Second, Arabic characters are context-sensitive, meaning a character's shape can change depending on its position in a word and adjacent characters. Attention helps accurately decompose these visual representations by focusing on the relevant con-060 text within the sequence. Third, the Arabic language includes diacritics, which are markings above 061 or below characters that can completely alter the semantics of a word. Attention enables the model 062 to effectively identify diacritics by considering their contextual influence on surrounding characters. 063

064 In addition to the intrinsic challenges of Arabic scripts, Arabic handwritten datasets, especially historical ones, are significantly smaller than those available for languages like English. Many HTR 065 works (Li et al., 2023; Wigington et al., 2018; Zhang et al., 2019) focus on languages using the 066 modern Latin alphabet, such as English and French, where large amounts of training data are readily 067 available. Common datasets include IAM (Marti & Bunke, 2002), with over 1,500 handwritten 068 pages, and RIMES (Grosicki et al., 2024), which comprises a mix of handwritten and printed text 069 across approximately 12,500 pages. In contrast, the largest public dataset for handwritten Arabic (Saeed et al., 2024) consists of just over 1,600 pages, while another widely used dataset, KHATT 071 (Mahmoud et al., 2012), contains only 1,000 pages. A notable exception is the MADCAT (Lee et al., 072 2012; 2013a;b) dataset, which contains over 40,000 pages of handwritten Arabic. However, it is not 073 focused on historical writing, highlighting the limited availability of resources for historical texts.

074 We base our approach on TrOCR (Li et al., 2023) and leverage domain knowledge of the Arabic lan-075 guage to identify key factors in building an effective historical Arabic HTR system. We incorporate 076 a novel image preprocessor and synthetic dataset generator to enhance performance by minimizing 077 horizontal information loss and expanding the training dataset with realistic synthetic images. We perform extensive evaluation and cross-dataset experiments on HATFORMER. We will release the 079 image preprocessor, tokenizer, model weights, and source code for our HTR system, along with a detailed guide for researchers to interface our system with existing text detection packages for 081 page-level HTR evaluations and practical deployment. Additionally, we will release our dataset of realistic synthetic Arabic images and its generation source code, as well as provide an OCR Error Diagnostic App and its source code to benefit both machine learning and history studies researchers. 083 The contributions of our work are threefold. 084

- 1. Our proposed HATFORMER for historical Arabic HTR outperforms the state of the art across various Arabic handwritten datasets. It achieves a CER of 8.6% and 4.2% on the largest public and private handwritten Arabic datasets, respectively.
- 2. Our method has proven effective by leveraging the attention mechanism to address three intrinsic challenges of the Arabic language.
- 3. Our historical Arabic HTR system and OCR Error Diagnostic App will aid humanity researchers by automatically transcribing historical Arabic documents and debugging common recognition errors, thereby significantly enhancing the accessibility of these documents.

2 RELATED WORKS

085

087

089

090

091

092

094

095

096 Handwritten Text Recognition (HTR). Handcrafted features were historically used for optical character recognition (OCR) and HTR (Balm, 1970), but deep learning methods gradually took 098 over due to their improved performance. Common deep learning methods adopt the encoderdecoder paradigm where visual signals are encoded into a feature representation and the fea-100 ture is decoded for text generation. Graves & Schmidhuber (2008) proposed using a long short-101 term memory (LSTM) (Hochreiter & Schmidhuber, 1997) multidimensional recurrent neural net-102 work (MDRNN) (Graves et al., 2007) for feature extraction and a connectionist temporal classifi-103 cation (CTC) layer for decoding (Graves et al., 2006). Notably, Shi et al. (2016) introduced the 104 convolutional recurrent neural network (CRNN) architecture for OCR, where a CNN was used to 105 extract visual features from images, and a stacked bidirectional LSTM (BLSTM) (Graves & Schmidhuber, 2005; Graves et al., 2013) was used as the decoder. Puigcerver (2017); Wang & Hu (2017) 106 respectively adapted the encoder to use a CNN and modified recurrent convolutional neural net-107 work (RCNN) (Liang & Hu, 2015). Newer approaches (Michael et al., 2019; Wang et al., 2020)

attempt to incorporate the attention mechanism (Bahdanau et al., 2015) into the HTR pipeline. Co quenet et al. (2023) use the attention mechanism to perform full-page HTR, bypassing the need for
 line-level segmentation.

111 Transformers for HTR. Transformers (Vaswani et al., 2017) have recently been applied to 112 HTR with earlier works using architectures consisting of a CNN-feature-extractor encoder and 113 a transformer-encoder-decoder-hybrid decoder, which later was simplified to a transformer-only 114 encoder-decoder pair or a transformer-decoder-only architecture. Wick et al. (2021) proposed a 115 hybrid system that uses a CNN feature extractor and multiple encoder-decoder transformers for 116 bidirectional decoding. Li et al. (2023) proposed a transformer-only method utilizing pretrained 117 vision transformers (ViT) (Dosovitskiy et al., 2021), specifically BEiT (Bao et al., 2022), as its en-118 coder using raw pixels as input and text transformers, specifically RoBERTa (Liu et al., 2019), as the decoder. Fujitake (2024) proposed a transformer-decoder-only method, using GPT (Radford et al., 119 2018; 2019) in particular, with raw pixel inputs. However, the decoder-only method, in general, 120 requires more labeled data for training end-to-end, whereas a pretrained encoder could be used as 121 an initialization step for visual feature extraction. Our approach does not use a dedicated CNN fea-122 ture extractor and builds upon the transformer-only encoder-decoder architecture. ViTs have been 123 shown to outperform CNNs and can benefit from large-scale pretraining for downstream tasks with 124 low resources (Dosovitskiy et al., 2021), like Arabic HTR. 125

Arabic HTR. Arabic handwriting poses unique challenges to HTR systems, such as cursive writing, 126 connected letters, and context-dependent character shapes. One of the earliest approaches to Arabic 127 HTR is Graves & Schmidhuber (2008), which proposes using multidimensional LSTM (MDLSTM) 128 and CTC decoding. Shtaiwi et al. (2022); Lamtougui et al. (2023); Saeed et al. (2024) proposed using 129 a CNN and BLSTM architecture, with Shtaiwi et al. (2022); Saeed et al. (2024) based upon the Start, 130 Follow, Read network (Wigington et al., 2018). As with traditional English HTR, many Arabic HTR 131 systems are starting to use the transformer architecture. Mostafa et al. (2021) proposed a method that 132 combines a ResNet-101 (He et al., 2016) for feature extraction and an encoder-decoder transformer 133 for text prediction. Momeni & BabaAli (2024) proposed a system that solely uses transformers, 134 similar to Li et al. (2023), but also introduces transducers (Graves, 2012) for HTR, removing the 135 need for external postprocessing language models. We continue using transformers for HTR and leverage the most recent advancements to further improve recognition performance on Arabic texts. 136 137

- 138
- 130

3 BACKGROUND AND PRELIMINARIES

140 141 142

143

This section provides background information about the components that HATFORMER is built on.

Arabic-Character Encoding. Arabic characters can be efficiently represented in tokens for learn ing using byte-level byte pair encoding (BBPE) (Radford et al., 2019). It is a tokenization technique
 that compresses a string into a reversible compact representation by leveraging the UTF-8 encoding
 standard using a vocabulary dictionary. To train the vocabulary dictionary, it is initialized with all
 256 possible byte values as base tokens, allowing it to tokenize any Unicode character and eliminating the need for a task-specific vocabulary. They are then iteratively merged based on the most
 frequent token pairs in a corpus to form new tokens. This iterative process expands the vocabulary, allowing for more efficient encoding of frequent patterns.

151 **TrOCR.** The TrOCR framework (Li et al., 2023) for predicting text from images will be used as the 152 base architecture for this work. TrOCR employs a transformer-only encoder-decoder architecture, 153 specifically using a pretrained ViT as the encoder and a pretrained text transformer as the decoder. 154 The encoder takes an input image of shape $3 \times H_0 \times W_0$, which is resized to a fixed shape of 155 $3 \times H \times W$. The resized image is then decomposed into a sequence of $N = HW/P^2$ patches, 156 where each patch has a shape of $3 \times P \times P$. The encoder will use the patches with added positional 157 embeddings as input to generate encoder embeddings. The decoder employs masked attention on the 158 ground-truth text tokens to ensure it does not access more information during training than during 159 prediction. The ground-truth text tokens are then combined with the encoder embedding using crossattention. A linear layer projects the hidden states from the decoder to match the vocabulary size, 160 and the probabilities over the vocabulary are computed using the softmax function. Beam search is 161 used to generate the final output.



Figure 1: The architecture of HATFORMER. The input text-line image is processed by our BLOCK-PROCESSOR and the BEiT vision transformer. The ground-truth text string is tokenized using our Arabic BBPE tokenizer. The RoBERTa transformer is used for text prediction. HATFORMER addresses the three intrinsic challenges of Arabic scripts by leveraging attention and is able to work on smaller datasets with the help of our synthetic image training pipeline.

181 182 183

184 185

177

178

179

4 PROPOSED METHOD FOR HISTORICAL ARABIC HTR

In this section, we present HATFORMER, which tackles the unique challenges of Arabic handwriting
 recognition, particularly for historical documents. We describe the main components of our method,
 including an image processor for effective ViT information preprocessing, a text tokenizer for compact Arabic text representation, and a training pipeline that accounts for the limited availability of
 historic Arabic handwriting data.

Architecture and Unit of Analysis. Prediction for HTR involves recognizing and converting a text image into machine-readable characters. As illustrated in Figure 1, HATFORMER follows TrOCR's transformer encoder-decoder architecture for HTR text prediction. We focused on line-level images as in Li et al. (2023); Momeni & BabaAli (2024), which is more challenging than the word- and character-level predictions but less complex than the paragraph- and page-level predictions. This approach allows us to focus on Arabic HTR without the additional complexities of text document structure. HATFORMER can be easily integrated with existing layout detection methods, enabling full-page prediction capabilities.

198 199 200

4.1 BLOCKPROCESSOR FOR EFFECTIVE VIT INFORMATION PREPROCESSING

201 We introduce a BLOCKPROCESSOR to best prepare each text-line image for effective ViT compre-202 hension by applying image-processing insights and leveraging ViT's blocking and indexing behav-203 iors. The proposed BLOCKPROCESSOR works by first horizontally flipping a text-line image, then 204 standardizing its height to 64 pixels, and finally warping it to fill in the ViT's 384×384-pixel im-205 age container from left to right and top to bottom. The ViT's image container will allow up to six 206 nonoverlapping complete rows that are 384 pixels wide, accommodating line images with varying widths for up to 2,304 pixels. For shorter images, zeros will be padded. Figure 2(c) shows an out-207 put of the proposed BLOCKPROCESSOR respecting the input image's aspect ratio to allow potential 208 perfect reconstruction. In contrast, images in Figure 2(d), (b), and (f) show significant information 209 loss due to the direct use of ViT's image preprocessor. We provide analysis below and justify the 210 system design of BLOCKPROCESSOR. 211

Aspect Ratio. ViT resizes input images to 384×384 without respecting their original aspect ratios.
 This leads to an inefficient representation of text-line images from the Muharaf dataset, which has an average image width of 614 pixels after standardizing their heights to 64 pixels. A direct application of ViT image preprocessing will lead to horizontal compression of 1.6 times on average, losing the clarity of the strokes in the horizontal direction for confident recognition. Figure 2(d) shows a



Figure 2: Left: Our proposed BLOCKPROCESSOR respects the aspect ratio of (a) an original image and chunks it to fit within (c) a 384×384-pixel ViT image container. In contrast, the base ViT image processor naively resizes images to (d) fully occupy its square image container, resulting in (f) significant horizontal information loss of the vertical strokes when compared to (e) the raw version. **Right**: (g) Synthetic image generation pipeline. Realistic-looking text-line images are generated by randomly selecting words from a large Arabic corpus, rendering with a random font, paper background, and image augmentation.

ViT-resized text-line image and Figure 2(b) shows the image content if the resized image is rescaled back to its original shape. As the zoomed-in reconstruction block in Figure 2(f) reveals, the vertical strokes suffer the most severe blurring, making it difficult for any observer to confidently determine the exact thickness of a stroke at different vertical heights.

ViT Blocking & Indexing. In the BLOCKPROCESSOR, both horizontal flipping of text-line images and the standardization of their heights to 64 pixels are designed to better leverage ViT's blocking and indexing behaviors for more efficient transformer training. First, Arabic text-line images read from right to left, so flipping them horizontally can avoid representing the end of a sentence with beginning ViT image tokens. Even though positional embeddings will help with ordering, we opt not to add extra workload to the attention layers as it will potentially require more training data. Second, we standardize the line image to an integer multiple of the ViT's patch height of 16 pixels. This resizing choice ensures that when Arabic texts are written in the middle of a text line, the corresponding ViT image tokens with foreground text will always have similar indices. Without resizing the height to an integer multiple of 16 pixels, the boundary of a text line and the boundary of a ViT block will misalign at varying pixel counts for different rows. This will cause foreground text to appear in all ViT image tokens, increasing the learning complexity for the attention layers.

ARABIC BBPE TEXT TOKENIZER FOR COMPACT ARABIC TEXT REPRESENTATION 4.2

Text representation is integral for language modeling. Radford et al. (2019) showed the impact of using a byte-level representation for text with byte-level byte pair encoding (BBPE). This led to an optimal balance of token sequence length and vocabulary size. To efficiently represent Arabic text, we trained our own custom BBPE dictionary on a combined corpus from Abbas & Smaili (2005); Abbas et al. (2011); Saad & Alijla (2017). As the base BBPE dictionary from Radford et al. (2019) is skewed toward ASCII characters, our experiments show that Arabic text is represented with over 300% more tokens compared to our custom BBPE dictionary. The more compact representation from the custom BBPE dictionary results in a less complicated classification problem, resulting in higher accuracy along with our BLOCKPROCESSOR, as we will discuss in Section 5.5.

4.3 TRAINING ON REALISTIC SYNTHETIC AND REAL-WORLD TEXTLINE IMAGES

Our proposed method involves a two-stage training/fine-tuning process, i.e., training on a large synthetic dataset followed by fine-tuning on a real-world Arabic handwritten dataset.

274 Stage 1-Training on Large Synthetic Printed Dataset. To address the scarcity of historical 275 handwritten Arabic data and capture key intrinsic features of Arabic scripts, we first trained HAT-276 FORMER on a large dataset of one million synthetic text-line images. This approach mitigates the 277 impact posed by the limited availability of historical handwritten Arabic data. The synthetic images 278 contain all three inherent characteristics of Arabic, i.e., cursive writing, context-dependent character 279 shapes, and diacritics. This enables our system to learn these challenging characteristics of Arabic script before being trained on a downstream task. Synthetic training provides the necessary data 280 for the encoder to learn the visual features of Arabic, leading to more effective generalization. The 281 synthetic image generation pipeline will be described in Section 5.1 and shown in Figure 2. 282

283 Stage 2-Fine-Tuning on Real Handwritten Dataset. We fine-tune HATFORMER on real Arabic 284 handwritten datasets, primarily focusing on the Muharaf dataset containing 36,000 text-line images 285 due to its relevance to historical handwritten texts. To achieve strong performance on Arabic HTR, we leverage a technical insight for large-scale training from Hao et al. (2019); Mosbach et al. (2021). 286 Traditional machine learning theory suggests that when the validation loss flattens, the model has 287 converged, and no further learning occurs (Mohri et al., 2018; Jo, 2021). However, Mosbach et al. 288 (2021) demonstrated that transformers can continue to improve in task performance long after the 289 validation loss has plateaued. Mosbach et al. (2021) indicates that achieving a near-perfect training 290 loss can serve as a strong baseline for model performance. In Stage 2, we train past the plateau of 291 the validation loss and approach a near-perfect training loss while monitoring the validation CER as 292 the stopping criteria, which can take twice as long as the minimum validation loss. 293

2004

5 EXPERIMENTAL RESULTS

295 296 297

We present the experimental results for HATFORMER on three Arabic handwritten datasets and compare it with other Arabic HTR baselines. We also conduct ablation studies to assess the effectiveness of each component and analyze various parameters of HATFORMER.

299 300 301

298

5.1 SYNTHETIC & REAL-WORLD ARABIC DATASETS

Synthetic Stage 1 Training Dataset. For our Stage 1 training dataset, we generated 1,000,065
 synthetic images of Arabic text lines. We first randomly sampled between 1–20 words from an Arabic corpus containing 8.2 million words. The sampled words were then paired with one of 54
 Arabic text fonts on a background chosen from 130 paper background images and one of eight image augmentations to generate synthetic line images. Our ablation study in Section 5.5 will show that English OCR initialization is insufficient and synthetic Arabic training is required.

Arabic HTR Datasets. The Muharaf dataset (Saeed et al., 2024) is a collection of historical hand-309 written Arabic manuscripts that span from the early 19th century to the early 21st century. The 310 dataset contains over 36,000 text line images, which vary significantly in quality, from clear writing 311 on clean white backgrounds to illegible sentences on creased pages with ink bleeds. The KHATT 312 dataset (Mahmoud et al., 2012) is a collection of Arabic handwriting samples with over 6,600 seg-313 mented line images. All images have black text on a clean white background. The MADCAT 314 dataset (Lee et al., 2012; 2013a;b) is a collection of 740,000 handwritten Arabic line images created 315 under controlled writing conditions. All images have black text on a clean, white background. See 316 Appendix B for more detailed descriptions of each dataset.

317

318 5.2 EXPERIMENTAL CONDITIONS 319

We initialized our model from HuggingFace's trocr-base-stage1 334M parameter model. We use BEiT (Bao et al., 2022) and RoBERTa (Liu et al., 2019) as the encoder and decoder, respectively, since TrOCR (Li et al., 2023) empirically showed that they achieved the best CER performance. We used a batch size of 15 with a learning rate of 5×10^{-5} and linear warmup of 20,000 steps for synthetic Stage 1 training. For Stage 2 fine-tuning, we used a batch size of 30 with a learn-

Dataat	M	A		Training Data	Test Data	Model	CER (%)
Dataset	Model	Architecture	CER (%)↓	Muharaf	KHATT	Saeed et al. (2024)	38.5
Muharaf	Saeed et al. (2024)	CRNN	14.9	(Full)		Proposed Model	22.8
(Full)	Proposed Model	Transformer	11.7		MADCAT	Saeed et al. (2024)	30.5
Muharaf	Saeed et al. (2024)	CRNN	17.6			Proposed Model	21.6
Arabic Only)	Proposed Model	Transformer	86	Muharaf	KHATT	Saeed et al. (2024)	33.0
(Anable Only)	i toposed Model	mansionnei	0.0	(Arabic Only)		Proposed Model	27.5
KHATT	Saeed et al. (2024)	CRNN	14.1		MADCAT	Saeed et al. (2024)	28.9
	Lamtougui et al. (2023)	CRNN	19.9			Proposed Model	26.5
	Momeni & BabaAli (2024)	Transformer	18.5	KHATT	Muharaf	Saeed et al. (2024)	43.8
	Proposed Model	Transformer	15.4			Proposed Model	40.7
ΜΑΡΟΔΤ	Saeed et al. (2024)	CRNN	5.5		MADCAT	Saeed et al. (2024)	17.8
MINDON		CRIM	1.0			Proposed Model	18.1
	Shtaiwi et al. (2022)	CRNN	4.0	MADCAT	Muharaf	Saeed et al. (2024)	43.5
	Rawls et al. (2018)	CRNN	1.5 ¹			Proposed Model	41.4
	Proposed Model	Transformer	4.2		KHATT	Saeed et al. (2024)	17.8
Combined	Proposed Model	Transformer	15.3			Proposed Model	16.3

Table 1: Performance on Arabic Handwritten Datasets.



¹ Used the 2013 NIST OpenHART evaluation tools for computing

CER/WER, which involved normalizing certain diacritics.

342 343

ing rate of 10^{-4} and linear warmup of 2,000 steps. The warmup was followed by an inverse square root schedule for both Stage 1 training and Stage 2 fine-tuning. We trained on 2 to 4 A100 or H100 GPUs.

We did Stage 1 training on a train-validation-test dataset split of 90-9-1 for 1,000,065 synthetic line 346 images. We fine-tuned using a split of 85–15–5 for 25,767 line images from the Muharaf dataset; 347 the author recommended a 72-14-14 split for 6,687 line images from the KHATT dataset and a 72-348 18–10 split for 741,877 line images from the MADCAT dataset. We used traditional validation loss 349 as the early stopping criterion during Stage 1 of training, with a maximum of 5 epochs. However, 350 we used the overtraining technique during our Stage 2 fine-tuning and utilized the validation CER 351 for early stopping as explained in Section 4.3. 352

5.3 MAIN RESULTS

355 We compare the performance of HATFORMER against state-of-the-art baselines across the Muharaf, 356 KHATT, and MADCAT datasets. We evaluate the HTR performance using the character error 357 rate (CER) (Levenshtein, 1966), which is widely used for assessing the accuracy of OCR and HTR 358 systems (Neudecker et al., 2021). It is defined as CER = (S + D + I)/N, where S is the num-359 ber of substitutions, D is the number of deletions, I is the number of insertions, and N is the total 360 number of characters in the original text. The CER is based on the edit distance, which calculates 361 the number of aforementioned operations required to transform the predicted text into the original 362 text. We also performed cross-dataset comparisons to evaluate HATFORMER's ability to generalize across different datasets. 363

364 Table 1 reports the CER for HATFORMER and several existing baselines across the three datasets. An important note is that the only existing baseline for the Muharaf dataset is Saeed et al. (2024). 366 Since the source code for many existing Arabic HTR baseline models is not publicly available, 367 except Saeed et al. (2024), we compared our results to the reported numbers obtained from their 368 papers. For Saeed et al. (2024), we retrained their model on each dataset and with stage-1 synthetic 369 training for a fair comparison. It is important to note that the dataset splits used in these baselines may differ from those in our experiments, potentially affecting direct comparisons. Additionally, 370 we conducted experiments on two variants of Muharaf, the entire dataset and a subset containing 371 only Arabic characters. This allows us to investigate the impact of non-Arabic characters on HTR 372 performance. For clarity, our analysis will refer to the Arabic-only subset as Muharaf. 373

374 We first compared with CNN and RNN-based methods. HATFORMER achieves a CER of 8.6% 375 and 15.4% on the Muharaf and KHATT datasets, respectively, as compared to Saeed et al. (2024) who achieved a CER of 17.6% and 14.1%. Lamtougui et al. (2023) achieved a CER of 19.9% on the 376 KHATT dataset. These results indicate that the transformer architecture can significantly outperform 377 traditional HTR methods based on CRNNs with a 23-51% improvement in CER for handwritten

333 334

340 341

344

345

353

354

324



Figure 3: Self- and Cross-attention map visualizations. Yellow highlights areas of greater attention, 393 with attention maps overlaid onto the input image for easier comparison. Left: ViT encoder self-394 attention maps for selected patch tokens. The top of each column shows the relevant patch, followed 395 by attention maps showing what the transformer attends to as it progresses through its subsequent 396 layers. The leftmost column shows the attention for a diacritic patch. Red lines indicate the layer 397 cutoff where the attention association becomes too broad, as identified by our Arabic expert. **Right**: 398 RoBERTa decoder cross-attention maps for selected ground truth text tokens. Each row represents 399 consecutive text tokens, read from right to left, with the decoded token string above each map. 400 Tokens are annotated based on their type: red underlines indicate diacritic tokens, green underlines 401 denote subword tokens, and all other tokens correspond to full words, as identified by our Arabic language expert. The attention maps reveal the model's ability to attend to relevant image regions 402 for each token. It can handle a diverse range of text, from small diacritics to complex compounded 403 characters, demonstrating the model's ability to overcome the inherent challenges of Arabic script. 404

405

Arabic. This aligns with computer vision and natural language processing trends, where transformers are increasingly favored due to their superior ability to take care of contextual information.

While Saeed et al. (2024) slightly outperformed our method on the KHATT dataset, HATFORMER 409 remains comparable. We attribute this to differences in the dataset characteristics, which may favor 410 Saeed et al. (2024)'s hybrid architecture. Our results imply that CNNs and RNNs are no longer 411 required for HTR. We enable a fully transformer-based model that can surpass these hybrid archi-412 tectures by utilizing vision transformers as standalone feature extractors combined with a text trans-413 former decoder. We credit the effectiveness of our model to the attention mechanism, which allows 414 for learning contextual information critical for language modeling. Figure 3 illustrates how the at-415 tention mechanism captures character relationships. See Appendix A for a more detailed description 416 and analysis of the attention maps.

We also compared our approach with transformer-based methods. HATFORMER achieves a CER of 15.4%, compared with Momeni & BabaAli (2024), who achieved a CER of 18.5% on the KHATT dataset. Our 17% improvement in CER demonstrates the effectiveness of our preprocessing and overtraining methods. Our preprocessing pipeline mitigates information loss caused by horizontal image compression, resulting in a CER improvement discussed in Section 5.5, while our overtraining strategy establishes a strong baseline, ultimately leading to better performance.

423 HATFORMER achieves a CER of 4.2%, comparable to other baseline Arabic HTR models on the 424 MADCAT dataset. The 1.5% CER achieved by Rawls et al. (2018) may be due to several fac-425 tors, specifically text normalization during evaluation (Rawls et al., 2018), which can significantly 426 improve performance as shown in Section 5.5. MADCAT also presents many unique dataset-427 specific complexities and requires distinct preprocessing techniques, as highlighted by Abandah 428 & Al-Hourani (2018). Furthermore, as both KHATT and MADCAT are non-historical datasets, they pose a different set of challenges compared to the historical Arabic texts that are the main focus of 429 our work. While we included MADCAT and KHATT in our evaluation for a more complete com-430 parison with existing Arabic HTR systems, we did not specifically optimize for them, as our primary 431 goal is to enhance the performance of historical Arabic HTR.

Table 3: Ablation Study on Muharaf.

	Model	CER (%) \downarrow	Model	CER (%) \downarrow
(A)	Proposed Model	8.6		- (.) •
(B)	(A) - Overtraining	9.9	Best Base Model	8.6
(C1)	(B) - BLOCKPROCESSOR + TrOCR Processor	11.4	Domovo discritios	8.0
C2)	(B) - Modified text tokenizer + TrOCR Tokenizer	10.0	+ Remove unacritics	8.0
(D)	(B) - (C1) - (C2)	10.4	+ Remove without context	7.4
(E)	(D) - Synthetic Stage-1 fine-tuning	14.6	Demove with context	67
(F)	(E) - Pretrained weights	86.0	+ Remove with context	0.7

Table 4: Effects of Arabic normalization

postprocessing on Muharaf.

We combined the three handwritten datasets into a single large dataset to evaluate the model's performance across diverse handwriting styles. Using this combined dataset, HATFORMER achieved a CER of 15.3%. While this result is slightly worse than the individual dataset CERs, it reflects the challenge of adapting to significant image content and style variability across the Muharaf, KHATT, and MADCAT datasets, indicating that HATFORMER can still extract meaningful shared features even with the increased difficulty of combining datasets.

450 451 5.4 CROSS-DATASET EVALUATION

We conducted cross-dataset evaluations to explore the generalization ability of our model. Table 2 shows the results of cross-dataset evaluation. This table reveals the importance of using historical handwriting data for a strong general Arabic HTR model. While training on modern Arabic handwriting using either KHATT or MADCAT gives a high CER of 40% on Muharaf, training on historical Muharaf data gives a lower CER of 26% on modern Arabic handwriting. Hence, this shows that our model can perform well on the historic Muharaf handwriting and generalize to in-the-wild, unseen modern handwritten Arabic.

We also ran cross-dataset evaluations using Saeed et al. (2024)'s HTR system. As seen in Table 2, our proposed model outperforms their approach in every evaluation except one. Notably, when training on Muharaf (Arabic only) and testing on KHATT, our model outperforms (Saeed et al., 2024) by 16.7% at a CER of 27.5% compared to 33%, which further shows our model's capability to generalize better compared to other methods.

464 465

466

432

433

444

445

446

447

448

449

5.5 ABLATION STUDY

In our ablation study, we quantify the impact of each component of our model by starting with our
 best model and removing each component one at a time, as shown in Table 3 .

Baseline Model (A). Our baseline model achieved a CER of 8.6%. This served as the benchmark
 against which we compared the performance of the ablated models.

471
472
472
473
474 **Overtraining (B).** When we only trained to the minimum validation loss, we observed a slight increase in CER by 1.3%. This result is consistent with Hao et al. (2019); Mosbach et al. (2021), suggesting that our model was not fully trained.

BLOCKPROCESSOR and Modified Text Tokenizer. (C1) & (C2) & (D). When the BLOCKPRO-475 CESSOR and Arabic BBPE were added together, this led to a 0.5% CER improvement supporting 476 our ideas in Sections 4.1 and 4.2. Replacing TrOCR's image processor with our BLOCKPROCES-477 SOR led to a 0.4% CER improvement, whereas replacing the modified text tokenizer with TrOCR's 478 tokenizer led to a -1.0% CER performance change. This indicates that the BLOCKPROCESSOR 479 enhances image feature extraction. However, the modified text tokenizer struggles when paired with 480 TrOCR's processor due to the naive resizing, which discards essential features needed for predicting 481 compact Arabic token representations, as discussed in Section 4.1. The 1.1% CER improvement 482 observed due to the synergy when both components are combined highlights their complementary 483 roles: the BLOCKPROCESSOR enables richer feature extraction, while the Modified Text Tokenizer ensures compact and accurate Arabic text representation. This shows the importance of aligning 484 task-specific components to the target language, as their interaction can yield significant synergistic 485 effects beyond individual contributions.

486
487
488
488
489
489
480
480
480
480
480
481
482
483
484
484
484
485
486
486
486
487
488
488
488
489
488
489
489
489
480
480
480
480
480
481
481
482
482
483
484
484
484
485
486
486
487
488
488
488
489
488
489
488
489
489
488
489
489
489
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480
480

 Pretrained Weights (F). When the training of HATFORMER was started from randomly initialized
 weights, the model's performance plummeted to a CER of 86.0%. Despite the major differences between English and Arabic scripts, leveraging TrOCR's synthetic pretraining checkpoint for English
 OCR led to better results.

Arabic Specific Postprocessing Normalization. We leveraged Arabic domain knowledge to group 494 our model errors into normalization categories: replace without context, replace with context, 495 and remove diacritics. The replace without context category normalizes characters to a single 496 form that is phonetically similar and generally does not change the meaning of the word. The 497 replace with context category is where more aggressive normalization is applied. Characters that 498 are similar but can change the word's meaning are converted to a single form. Remove diacritics 499 is relevant to applications such as historical informational archival and search, where normalizing 500 certain characters into a single form is acceptable. Diacritics, in some cases, can be sparsely used and be removed in an Arabic OCR system. Table 4 shows that the model performance in terms of CER improves by 1.9% points or an additional $\sim 0.6\%$ per post-processing for each category. 502

503 504

505

5.6 FACTOR/SENSITIVITY STUDY

Table 5: Block Processor Comparison

We analyze the impact of various parameters on model
performance, with additional experiments in Appendix C.

Block Processor Methods. Several studies have explored dynamic aspect ratio image-processing approaches in vision-language models (Bavishi et al., 2023; Fadeeva et al., 2024; Dehghani et al., 2024). We compared our proposed BLOCKPROCESSOR with two notable methods:

Processor	CER (%) \downarrow		
Lee et al. (2023)	20.3		
Li et al. (2023)	11.4		
BLOCKPROCESSOR	8.6		

513 TrOCR (Li et al., 2023), which employs the standard ViT processing approach by resizing input 514 images to 384-by-384 pixels, and Pix2Struct (Lee et al., 2023), which scales input images while 515 preserving the aspect ratio to extract the maximum number of patches within a given sequence 516 length. Table 5 shows that our BLOCKPROCESSOR achieves the best CER of 8.6% on the Muharaf dataset. As discussed in Sections 4.1 and 5.5, TrOCR's processor suffers from information loss due 517 to its inefficient representation of resized images. While Pix2Struct addresses this by preserving the 518 aspect ratio, it introduces variability in the semantic meaning of patches, even when using absolute 519 2-dimensional positional embeddings. A patch may correspond to a character fragment in shorter 520 images, while a patch might represent an entire character in longer images. This inconsistency in 521 patch representation can negatively impact the model's ability to interpret and process the input. 522

522 523 524

525

6 CONCLUSION AND LIMITATIONS

In this paper, we have presented HATFORMER, a dedicated Arabic handwritten text recognition system harnessing the transformer's attention mechanism to address the unique challenges of the Arabic language. Our system integrates training methods with image and text processing techniques designed for Arabic HTR. Experiments show that HATFORMER outperforms baseline methods across multiple real-world datasets, highlighting the effectiveness of our approach.

531 HATFORMER demonstrates significant progress in historical Arabic handwritten text recognition 532 but also has some limitations. As a text line recognition model, its performance relies on the quality 533 of line segmentations during real-world inference. Additionally, HATFORMER struggles with line 534 images exhibiting extreme slants without angle normalization, which can impact recognition accu-535 racy. The computational demands of the training process, particularly with the overtraining strategy, 536 pose challenges for institutions with extremely limited resources. Addressing this, future work could 537 explore parameter-efficient fine-tuning, such as low-rank adaptation (LoRA) (Hu et al., 2022) to enhance accessibility. These limitations point to key areas for improvement, including preprocessing 538 enhancements and optimization of training methods, to increase robustness and applicability across diverse contexts.

540 REFERENCES 541

548

553

554

555

565

566

567

569

- Gheith Abandah and Ahmad Al-Hourani. Challenges and preprocessing recommendations for 542 MADCAT dataset of handwritten Arabic documents. In International Congress on Image and 543 Signal Processing, BioMedical Engineering and Informatics, pp. 1–9. IEEE, 2018. 544
- Mourad Abbas and Kamel Smaili. Comparison of topic identification methods for Arabic language. 546 In International Conference on Recent Advances in Natural Language Processing, pp. 14–17, 547 2005.
- Mourad Abbas, Kamel Smaili, Daoud Berkani, et al. Evaluation of topic identification methods on 549 Arabic corpora. Journal of Digital Information Management, 9(5):185–192, 2011. 550
- 551 Dzmitry Bahdanau, KyungHyun Cho, and Yoshua Bengio. Neural machine translation by jointly 552 learning to align and translate. International Conference on Learning Representations, 2015.
 - GJ Balm. An introduction to optical character reader considerations. Pattern Recognition, 2(3): 151-166, 1970.
- 556 Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEiT: BERT pre-training of image transformers. International Conference on Learning Representations, 2022.
- 558 Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, 559 and Sağnak Taşırlar. Introducing our multimodal models. https://www.adept.ai/blog/ 560 fuyu-8b, 2023. 561
- 562 Denis Coquenet, Clément Chatelain, and Thierry Paquet. DAN: A segmentation-free document 563 attention network for handwritten document recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence, 45(7):8227-8243, 2023. 564
- Mostafa Dehghani, Basil Mustafa, Josip Djolonga, Jonathan Heek, Matthias Minderer, Mathilde Caron, Andreas Steiner, Joan Puigcerver, Robert Geirhos, Ibrahim M Alabdulmohsin, et al. Patch n'Pack: NaViT, a vision transformer for any aspect ratio and resolution. Advances in Neural 568 Information Processing Systems, 36, 2024.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 570 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-571 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at 572 scale. International Conference on Learning Representations, 2021. 573
- 574 Anastasiia Fadeeva, Philippe Schlattner, Andrii Maksai, Mark Collier, Efi Kokiopoulou, Jesse 575 Berent, and Claudiu Musat. Representing online handwriting for recognition in large visionlanguage models. arXiv preprint arXiv:2402.15307, 2024. 576
- 577 Safiullah Faizullah, Muhammad Sohaib Ayub, Sajid Hussain, and Muhammad Asad Khan. A survey 578 of OCR in Arabic language: Applications, techniques, and challenges. Applied Sciences, 13(7): 579 4584, 2023. 580
- Masato Fujitake. DTrOCR: Decoder-only transformer for optical character recognition. 581 In *IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 8025–8035, 2024. 582
- 583 Alex Graves. Sequence transduction with recurrent neural networks. International Conference of 584 Machine Learning Workshop, 2012. 585
- Alex Graves and Jürgen Schmidhuber. Framewise phoneme classification with bidirectional LSTM 586 networks. In IEEE International Joint Conference on Neural Networks, volume 4, pp. 2047–2052, 2005. 588
- 589 Alex Graves and Jürgen Schmidhuber. Offline handwriting recognition with multidimensional re-590 current neural networks. Advances in Neural Information Processing Systems, 21, 2008. 591
- Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist tem-592 poral classification: Labelling unsegmented sequence data with recurrent neural networks. International Conference on Machine Learning, pp. 369-376, 2006.

594 595	Alex Graves, Santiago Fernández, and Jürgen Schmidhuber. Multi-dimensional recurrent neural networks. In <i>International Conference on Artificial Neural Networks</i> , pp. 549–558, 2007.
590 597	Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep re- current neural networks. In <i>IEEE International Conference on Acoustics, Speech, and Signal</i>
599	<i>Processing</i> , pp. 6645–6649, 2013.
600	Emmanuèle Grosicki Matthieu Carré Edouard Geoffrois Emmanuel Augustin Françoise Preteux
601	and Ronaldo Messina. RIMES. https://doi.org/10.5281/zenodo.10812725.2024.
602	
603	Yaru Hao, Li Dong, Furu Wei, and Ke Xu. Visualizing and understanding the effectiveness of BERT.
604	Conference on Natural Language Processing pp 4143–4152 Association for Computational
605	Linguistics, 2019.
607	Keiming He Viengun Zhang, Shaeging Dan, and Jian Sun. Deen residual learning for image reason
608 600	nition. In <i>IEEE Conference on Computer Vision and Pattern Recognition</i> , pp. 770–778, 2016.
610 611	Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. <i>Neural Computation MIT-Press</i> , 1997.
612	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
613	and Weizhu Chen. LoRA: Low-rank adaptation of large language models. International Confer-
614	ence on Learning Representations, 2022.
615	Taeho Jo. Machine learning foundations. Springer, 2021.
617	Hicham Lamtougui, HE Mouhtahii, Hassan Fouadi, and Khalid Satori. An afficient hybrid model
618	for arabic text recognition. <i>Computers, Materials & Continua</i> , 74(2):2871–2888, 2023.
619	
620	Yann LeCun, Bernhard Boser, John Denker, Donnie Henderson, Richard Howard, Wayne Hubbard, and Lawrence Jackel, Handwritten digit recognition with a back propagation network. Advances
621	in Neural Information Processing Systems, 2, 1989.
622	
623	David Lee, Safa Ismael, Stephen Grimes, Dave Doermann, Stephanie Strassel, and Zhiyi Song.
624 625	LDC2012T15, 2012.
626	David Lee, Safa Ismael, Stephen Grimes, Dave Doermann, Stephanie Strassel, and Zhiyi Song.
627	MADCAT Phase 2 Training Set LDC2013T09. https://catalog.ldc.upenn.edu/
629	LDC2013T09, 2013a.
630	David Lee, Safa Ismael, Stephen Grimes, Dave Doermann, Stephanie Strassel, and Zhiyi Song.
631	MADCAT Phase 3 Training Set LDC2013T15. https://catalog.ldc.upenn.edu/
632	LDC2013T15, 2013b.
633	Kenton Lee, Mandar Joshi, Iulia Raluca Turc, Hexiang Hu, Fangyu Liu, Julian Martin Eisenschlos,
634	Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. Pix2Struct: Screen-
635	shot parsing as pretraining for visual language understanding. In International Conference on
636	<i>Machine Learning</i> , pp. 18893–18912, 2023.
639	Vladimir Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. Pro-
639	ceedings of the Soviet Physics Doklady, 1966.
640	Minghao Li, Tengchao Ly, Jingye Chen, Lei Cui, Yijuan Lu, Dinei Florencio, Cha Zhang, Zhou-
641	jun Li, and Furu Wei. TrOCR: Transformer-based optical character recognition with pre-trained
642	models. In AAAI Conference on Artificial Intelligence, volume 37, pp. 13094–13102, 2023.
643	Ming Liang and Xiaolin Hu. Recurrent convolutional neural network for object recognition. In
644	<i>IEEE Conference on Computer Vision and Pattern Recognition</i> , pp. 3367–3375, 2015.
645	Vinhan Liu Mula Ott Naman Gouel Lingfai Du Mandar Lash: Dansi Char, Orean L. Mil
645 647	Lewis, Luke Zettlemover, and Veselin Stovanov RoBERTa. A robustly optimized BERT pre-
047	training approach. arXiv preprint arXiv:1907.11692, 2019.

648 Sabri A Mahmoud, Irfan Ahmad, Mohammad Alshayeb, Wasfi G Al-Khatib, Mohammad Tanvir 649 Parvez, Gernot A Fink, Volker Märgner, and Haikal El Abed. KHATT: Arabic offline handwritten 650 text database. In International Conference on Frontiers in Handwriting Recognition, pp. 449–454. 651 IEEE, 2012. 652 U-V Marti and Horst Bunke. The IAM-database: An english sentence database for offline handwrit-653 ing recognition. International Journal on Document Analysis and Recognition, 5:39–46, 2002. 654 655 Johannes Michael, Roger Labahn, Tobias Grüning, and Jochen Zöllner. Evaluating sequence-to-656 sequence models for handwritten text recognition. In IEEE International Conference on Docu-657 ment Analysis and Recognition, pp. 1286–1293, 2019. 658 Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. Foundations of machine learning. 659 MIT press, 2018. 660 661 Saleh Momeni and Bagher BabaAli. A transformer-based approach for Arabic offline handwritten 662 text recognition. Signal, Image and Video Processing, 18(4):3053-3062, 2024. 663 Marius Mosbach, Maksym Andriushchenko, and Dietrich Klakow. On the stability of fine-tuning 664 BERT: Misconceptions, explanations, and strong baselines. International Conference on Learn-665 ing Representations, 2021. 666 667 Aly Mostafa, Omar Mohamed, Ali Ashraf, Ahmed Elbehery, Salma Jamal, Ghada Khoriba, and 668 Amr S Ghoneim. OCFormer: A transformer-based model for Arabic handwritten text recognition. 669 In IEEE International Mobile, Intelligent, and Ubiquitous Computing Conference, pp. 182–186, 670 2021. 671 Rayyan Najam and Safiullah Faizullah. Analysis of recent deep learning techniques for Arabic 672 handwritten-text OCR and post-OCR correction. Applied Sciences, 13(13):7568, 2023. 673 674 Clemens Neudecker, Konstantin Baierer, Mike Gerber, Christian Clausner, Apostolos Antonacopou-675 los, and Stefan Pletschacher. A survey of ocr evaluation tools and metrics. In International 676 *Workshop on Historical Document Imaging and Processing*, pp. 13–18, 2021. 677 Daniel Parres and Roberto Paredes. Fine-tuning vision encoder-decoder transformers for handwrit-678 ing text recognition on historical documents. In International Conference on Document Analysis 679 and Recognition, pp. 253-268. Springer, 2023. 680 681 Joan Puigcerver. Are multidimensional recurrent layers really necessary for handwritten text recognition? In IEEE International Conference on Document Analysis and Recognition, volume 1, pp. 682 67-72, 2017. 683 684 Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language under-685 standing by generative pre-training. OpenAI blog, 2018. 686 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language 687 models are unsupervised multitask learners. OpenAI blog, 1(8):9, 2019. 688 689 Stephen Rawls, Huaigu Cao, Joe Mathai, and Prem Natarajan. How to efficiently increase resolution 690 in neural ocr models. In IEEE International Workshop on Arabic and Derived Script Analysis and 691 *Recognition*, pp. 140–144, 2018. 692 693 David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning internal representations by error propagation. In Parallel Distributed Processing: Explorations in the Microstructure of 694 Cognition, volume 1, pp. 318–362, 1986. 696 Motaz Saad and Basem Alijla. WikiDocsAligner: An off-the-shelf Wikipedia documents alignment 697 tool. In Palestinian International Conference on Information and Communication Technology, 2017. 699 Mehreen Saeed, Adrian Chan, Anupam Mijar, Joseph Moukarzel, Georges Habchi, Carlos Younes, 700 Amin Elias, Chau-Wai Wong, and Akram Khater. Muharaf: Manuscripts of handwritten arabic dataset for cursive text recognition. arXiv preprint arXiv:2406.09630, 2024.

702 Baoguang Shi, Xiang Bai, and Cong Yao. An end-to-end trainable neural network for image-based 703 sequence recognition and its application to scene text recognition. IEEE Transactions on Pattern 704 Analysis and Machine Intelligence, 39(11):2298–2304, 2016. 705 Reem E Shtaiwi, Gheith A Abandah, and Safaa A Sawalhah. End-to-end machine learning solution 706 for recognizing handwritten Arabic documents. In International Conference on Information and Communication Systems, pp. 180–185, 2022. 708 709 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, 710 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems, volume 30, pp. 5998-6008, 2017. 711 712 Jianfeng Wang and Xiaolin Hu. Gated recurrent convolution neural network for ocr. In Advances in 713 Neural Information Processing Systems, volume 30, 2017. 714 715 Tianwei Wang, Yuanzhi Zhu, Lianwen Jin, Canjie Luo, Xiaoxue Chen, Yaqiang Wu, Qianying Wang, and Mingxiang Cai. Decoupled attention network for text recognition. In AAAI Con-716 ference on Artificial Intelligence, volume 34, pp. 12216–12224, 2020. 717 718 Christoph Wick, Jochen Zöllner, and Tobias Grüning. Transformer for handwritten text recogni-719 tion using bidirectional post-decoding. In International Conference on Document Analysis and 720 Recognition, pp. 112-126. Springer, 2021. 721 Curtis Wigington, Chris Tensmeyer, Brian Davis, William Barrett, Brian Price, and Scott Cohen. 722 Start, follow, read: End-to-end full-page handwriting recognition. In European Conference on 723 Computer Vision, pp. 367–383, 2018. 724 725 Yaping Zhang, Shuai Nie, Wenju Liu, Xing Xu, Dongxiang Zhang, and Heng Tao Shen. Sequence-726 to-sequence domain adaptation network for robust text image recognition. In IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2740–2749, 2019. 727 728 729 А ATTENTION MAPS 730 731 We analyze the effectiveness of using self-attention for Arabic HTR, including visualizing the self-732 attention maps of our vision transformer encoder. We also visualize the cross-attention maps corre-

attention maps of our vision transformer encoder. we also visualize the cross-attention maps corresponding to a selected ground truth token.
Our visualization scheme for the vision transformer involves selecting a patch of interest in the image and then visualizing how it attends to other patches. In the heatmaps, a brighter color (yellow) indicates that the selected patch pays more attention to this patch. We accumulate the self-attention

indicates that the selected patch pays more attention to this patch. We accumulate the self-attention
heatmaps from the previous layers (taking into account residual connections) in the next layer to
get a more holistic view of the attention flow. In the following text, we discuss insights from the
attention maps and how our transformer model deals with the intricacies of Arabic handwriting and
addresses key Arabic script challenges that we outlined in Section 3.

Cursive. Figure 3 demonstrates that when a patch containing a cursive line is selected, the attention
 map first highlights the relevant strokes of that character. This indicates that the network learns to
 distinguish individual characters and strokes before applying broader, global attention, effectively
 connecting relevant patches. The cross-attention map further shows that the model accurately iden tifies character boundaries. The model successfully segments the entire word in the image for tokens
 with complex cursive dependencies.

Context-Sensitive. A character in Arabic can take multiple forms depending on its position within
 a word and the surrounding characters. The cross-attention maps in Figure 3 demonstrate that even
 when words are split into multiple tokens, the model can accurately differentiate between word
 pieces and segment each part. These maps reveal that the model effectively learns the complex
 morphological rules of Arabic script and can distinguish between different positional forms of the
 same character.

Diacritics. Figure 3 demonstrates that the diacritic patch can attend to the character patches it is associated with. Importantly, both self-attention and cross-attention maps indicate that diacritic marks are not treated as noise but carry a strong signal. The self-attention maps reveal that the model can associate the relevant character corresponding with the diacritic. The cross-attention maps show
 the model correctly identifying the position of diacritics within the corresponding word. These maps
 highlight the network's ability to incorporate these small marks into the final token predictions rather
 than ignoring them.

Attention Maps Cutoff. To further evaluate the self-attention mechanism, our Arabic expert coauthor analyzed the progression of attention associations across layers. Specifically, our expert identified layers where the attention between a patch and other patches becomes excessively broad relative to the associated word, potentially diluting the model's focus on relevant features. These cutoff points are marked with red lines in the visualizations. This analysis provides valuable insights into how effectively the model maintains meaningful associations and highlights potential areas for improvement, particularly in leveraging Arabic-specific linguistic and structural knowledge.

B DATASETS

B.1 MUHARAF

The <u>Muharaf</u> dataset (Saeed et al., 2024) is a public collection of historical handwritten Arabic
manuscripts spanning from the early 19th century to the early 21st century. The dataset encompasses diverse document types, including personal letters, poems, dialogues, legal records, correspondences, and church documents. It consists of over 36,000 text line images, exhibiting significant variability in quality. These range from clear handwriting on clean white paper to highly degraded illegible text on creased pages with ink bleed-through. Fluent Arabic speakers scanned and transcribed the historical documents, ensuring high-quality annotations. This makes the Muharaf dataset a valuable resource for advancing research in historical handwriting recognition in Arabic.

B.2 KHATT

The <u>KHATT</u> dataset (Mahmoud et al., 2012) is a standard benchmark for Arabic HTR tasks. It is a public collection of modern Arabic handwriting samples comprising over 6,600 segmented line images. All images feature black text written on clean white backgrounds, ensuring consistent visual quality. The dataset was created under controlled conditions, where 1,000 participants transcribed 2,000 unique texts provided to them.

B.3 MADCAT

The <u>MADCAT</u> dataset (Lee et al., 2012; 2013a;b) is a proprietary dataset created by the Linguistic Data Consortium (LDC) to support the DARPA MADCAT Program. It comprises 740,000 modern handwritten Arabic line images created under controlled conditions with standardized writing speed, methodology, tool, and paper-type specifications. The text content was sourced from various digital mediums, including weblogs, newswires, and newsgroups. Each image features black text on a clean white background, ensuring high visual consistency. Due to its large size, the MADCAT dataset is a valuable resource for advancing Arabic HTR research.

B.4 SYNTHETIC

For our Stage 1 training dataset, we generated 1,000,065 synthetic images of Arabic text lines.
To create these, we randomly sampled between 1 and 20 words inclusive from an Arabic corpus comprising 8.2 million words, constructed by combining the datasets from Abbas & Smaili (2005);
Abbas et al. (2011); Saad & Alijla (2017). The selected words were rendered using one of 54 Arabic fonts and placed on a background randomly selected from a set of 130 paper background textures.
We source the Arabic fonts from freely available online websites. The 130 paper backgrounds are created from the Muharaf dataset by copying parts of the background image or created by using



Figure 4: (a) The impact of synthetic Stage-1 fine-tuning size on final HTR performance. A larger synthetic Stage-1 fine-tuning dataset allows for better generalization in terms of CER. (b) The CER and latency effect of inference beam size of our model on Muharaf. Using a larger beam size leads to a more accurate model but reduced speed. A beam width of three demonstrates a good trade-off between accuracy and computational speed. (c) The impact of inference length penalty of our model on Muharaf. A length penalty of 0.2 to 0.8 is preferred to achieve the best CER.

online paper texture images. Additionally, we applied one of eight image augmentation techniques: width distortion, height distortion, barrel distortion, left arc, right arc, left rotation, right rotation, or no distortion. We will release the realistic Arabic synthetic dataset and code to generate the images.

C ADDITIONAL FACTOR/SENSITIVITY STUDY

We analyze the impact of various parameters on model performance with further discussion and comparison of image processors in Section 5.6.¹

Stage-1 Synthetic Dataset Size. Figure 4(a) shows the impact of the synthetic Stage-1 training
 dataset size on the final performance of HATFORMER on the Muharaf dataset. As the size of the
 synthetic dataset increases, the CER decreases, demonstrating improved generalization. Specifically,
 datasets of 500k images and 1M images yield the best performance, with the CER dropping below
 10%. This trend suggests that a larger synthetic Stage-1 training dataset enhances the model's ability
 to effectively handle the inherent challenges of Arabic, ultimately leading to better CER performance
 in downstream HTR tasks.

Consecutive Whitespaces. The reported evaluation results throughout this paper were derived by
 removing consecutive whitespaces at test time. This is in line with the default CER score implementation in the HuggingFace evaluate library (v0.4.2). We observed that performing this normalization
 during the training stage instead of inference time leads to an additional 0.2% CER improvement.

Inference Beam Width. Figure 4(b) shows the effect of beam width on CER and generation speed.
The CER improves until the beam width is three and stabilizes beyond this point. Hence, we used a
beam width of three in our reported results. The inference speed was measured in tokens per second
over the Muharaf test set (total time / total tokens) on a single A10 (24GB) GPU with a batch size of
From the inference speeds we can see that our model can be used in a low-resource environment.

Inference Length Penalty. The length penalty parameter in beam search biases the generated output
 sequence length, where negative values encourage shorter sequences and positive values encourage
 longer ones. In Figure 4(c), we empirically show that on the Muharaf dataset, our model performs
 optimally using a length penalty between 0.2 and 0.8.

¹One notable parameter that was infeasible to study was decreasing the ViT patch size due to training computational complexity. Reducing patch size in ViT's results in a longer sequence and the attention mechanism requires quadratic cost $O(n^2)$ with respect to sequence length n.