

Handwritten Text Recognition for Low Resource Languages

Sayantan Dey

Indian Institute of Technology Roorkee, Haridwar, Roorkee, 247667, Uttarakhand, India

Alireza Alaei

Faculty of Science and Engineering, Southern Cross University, Gold Coast, 4225, QLD, Australia

Partha Pratim Roy

Indian Institute of Technology Roorkee, Haridwar, Roorkee, 247667, Uttarakhand, India

Abstract

Despite considerable progress in handwritten text recognition, paragraph-level handwritten text recognition, especially in low-resource languages, such as Hindi, Urdu and similar scripts, remains a challenging problem. These languages, often lacking comprehensive linguistic resources, require special attention to develop robust systems for accurate optical character recognition (OCR). This paper introduces BharatOCR, a novel segmentation-free paragraph-level handwritten Hindi and Urdu text recognition. We propose a ViT-Transformer Decoder-LM architecture for handwritten text recognition, where a Vision Transformer (ViT) extracts visual features, a Transformer decoder generates text sequences, and a pre-trained language model (LM) refines the output to improve accuracy, fluency, and coherence. Our model utilizes a Data-efficient Image Transformer (DeiT) model proposed for masked image modeling in this research work. In addition, we adopt a RoBERTa architecture optimized for masked language modeling (MLM) to enhance the linguistic comprehension and generative capabilities of the proposed model. The transformer decoder generates text sequences from visual embeddings. This model is designed to iteratively process a paragraph image line by line, called implicit line segmentation. The proposed model was evaluated using our custom dataset ('Parimal Urdu') and ('Parimal Hindi'), introduced in this research work, as well as two public datasets. The

proposed model achieved benchmark results in the NUST-UHWR, PUCIT-OUHL, and Parimal-Urdu datasets, achieving character recognition rates of 96.24%, 92.05%, and 94.80%, respectively. The model also provided benchmark results using the Hindi dataset achieving a character recognition rate of 80.64%. The results obtained from our proposed model indicated that it outperformed several state-of-the-art Urdu text recognition methods.

Keywords: Segmentation-free Handwriting Recognition, Large Language Model, Transformer Network, Parimal Urdu and Hindi Datasets.

1. Introduction

Since the emergence of digital technology, optical character recognition (OCR) has been the subject of extensive research [30]. OCR aims to transform scanned document images into machine-readable text so that a vast amount of unstructured handwritten or printed text can be analyzed. Historically, handwritten text recognition, including Urdu and Hindi scripts, has been approached through segmentation at incremental levels of granularity—line, word, and character [11, 19, 25, 26]. Convolutional neural networks (CNNs) with long short-term memory (LSTM) networks were commonly used in the initial wave of deep learning-based OCR models [31, 32]. Using region proposals for better feature extraction, Region-Based CNNs (RCNNs) [31] enhanced this architecture. However, as these models were trained character-by-character rather than word-by-word, they had a significant drawback: they ultimately required a separate language model. As a result, Transformer-based OCR (TrOCR) and other transformer-based OCR models were proposed in the literature [17]. In this approach, a pre-trained language model was employed for text production and substitutes a pre-trained Vision Transformer in the CNN layers. While improving over its predecessor, it has faced the inherent challenge of accurately defining and isolating individual entities (lines, words, and characters) for effective segmentation at each level. In contrast, we propose a new segmentation-free text recognition approach that directly recognizes text from paragraph images without explicit segmentation during the training and decoding phases.

However, the leap to paragraph-level recognition introduces new complexities, such as varying text region/line numbers, diverse layout patterns and skews, and the imperative of defining a coherent reading order. To address these issues, we propose a model that adeptly navigates these challenges us-

ing a vision-language model for image-to-text generation, where a pre-trained image encoder extracts visual features using learned tokens. A transformer-based decoder then generates an initial text output that undergoes contextual refinement. Finally, a pre-trained language model is considered to enhance coherence and fluency, ensuring high-fidelity text generation for document understanding and OCR-based tasks. This approach circumvents the error compounding of traditional multi-stage segmentation-recognition approaches. Our approach distinguishes itself by processing entire paragraphs of Handwritten Urdu and Hindi text through a segmentation-free transformer, simplifying the recognition pipeline while enhancing the capability of the model to handle the linguistic and stylistic diversity of both scripts. The complex ligatures and cursive nature of Urdu, alongside the distinct conjunct formations and diacritic variations in Hindi, are effectively addressed within our unified framework. Furthermore, preliminary pre-training on paragraph images is incorporated to elevate recognition accuracy, paving the way for seamless paragraph-level recognition across both languages.

This paper presents a novel approach to handwritten Urdu and Hindi text recognition (HUHTR), moving beyond traditional segmentation-based methodologies in favor of an end-to-end segmentation-free text recognition framework. The contributions of this paper are three-fold:

1. **Introduction of a robust methodology:** This study introduces a novel segmentation-free framework for recognizing handwritten Urdu characters, words, and lines at the paragraph level.
2. **Inclusion of a Large Language Model:** We incorporate the RoBERTa base language model, which has been pre-trained using masked language modeling to improve the comprehension of Urdu and Hindi text, enabling more accurate interpretations of the subtle language features.
3. **New dataset created - ‘Parimal Urdu’ and ‘Parimal Hindi’:** Central to our research is the development of a new dataset composed of a diverse array of handwritten Urdu and Hindi pages. This dataset comprises 500 pages contributed by ten individuals from various age groups in each language, ensuring a wide variety of natural handwriting styles. Each contributor was carefully selected to represent different writing variations, making this dataset a valuable resource for handwriting recognition research.

The remainder of the paper is outlined as follows: Section 2 provides an overview of the related work on Urdu text recognition. Section 3 details the

description of our proposed model. Section 4 discusses datasets, evaluation metrics, data augmentation, and experimental results and comparative analysis. Section 5 concludes the paper and presents future research directions.

2. Related Work

Due to various complexities and challenges, such as cursive text, complex ligatures, and varying writing styles, UHTR remains a difficult task [21]. Nevertheless, researchers have explored handwritten Urdu OCR in the literature [11, 18]. In the early stages, conventional machine learning models, such as SVM [22], were used for UHTR. However, a significant shift from traditional methods led to the introduction of holistic approaches utilizing hybrid CNN-RNN networks [15]. These modern methods employed CNNs to extract essential visual features from input images, which are then processed through Recurrent Neural Networks (RNNs) [13] to capture contextual information essential for transcription layers. Contemporary state-of-the-art Urdu OCR models incorporated various network architectures for feature extraction and sequential modeling. Bi-directional LSTM (BiLSTM) [12] networks and Multi-Dimensional LSTM (MDLSTM) [14] networks were particularly favored for sequential processing, with most models relying on a Connectionist Temporal Classification (CTC) layer for final output transcription.

In contrast, some models explored utilizing DenseNet [11] and GRU [16] networks, coupled with an attention-based decoding layer, marking a shift towards integrating attention mechanisms [11] to improve transcription accuracy. This trend of employing various deep learning architectures [10, 9, 8] mirrored the development trajectory observed in OCR technologies for languages with scripts similar to Urdu, such as Arabic. Moreover, to address the scarcity of extensive labeled datasets in this domain, enhanced prototypical network architectures were explored for few-shot recognition of handwritten Urdu characters, aiming to improve the classification performance with limited samples [34]. However, despite these advances, adapting Arabic OCR techniques to Urdu text often resulted in lower accuracy levels, emphasizing the unique challenges posed by Urdu script that were not fully addressed by existing methods. In addition, although each proposed model aspired to enhance Urdu OCR, a key limitation is the reliance on approaches designed for other languages without adequately tailoring them to accommodate the challenges inherent in Urdu script.

Research in handwritten Hindi character recognition has seen significant advancements through various approaches. Initial research focused on feature extraction techniques combined with conventional classifiers. For instance, studies employed k-means clustering and structural features to recognize handwritten Hindi characters [35]. Reddy and Babu [28] developed a handwritten Hindi character recognition system utilizing CNNs optimized with RMSprop and Adam, achieving high accuracy. In addition, Sharma and Ramakrishnan [25] introduced a handwritten Hindi character classification using a combination of global character and local sub-unit features, resulting in a recognition accuracy of 93.5%. Then, Chauhan et al. [26] proposed a script-independent deep learning network evaluated across multiple scripts, including Hindi, establishing new benchmarks with performance improvements up to 11%. Furthermore, a fusion-based hybrid-feature based on a combination of deep CNN features with handcrafted features has been explored, achieving a recognition accuracy of 98.7% [27]. The Integral Histogram of Oriented Displacement (IHOD) descriptor has also been used for handwritten Hindi script recognition, achieving high accuracy in character recognition and word spotting [36]. This method outperformed several CNN-based models and contributed a large dataset of handwritten Hindi word images. These studies collectively contributed to the advancement of Hindi handwritten character recognition, offering diverse methodologies to address the complexities inherent in processing handwritten Indic scripts. This oversight draws attention towards the missed opportunity to fully harness the potential inherent within this domain of research. Therefore, a refined approach that meticulously accounts for the unique complexities of low-resource script is necessary. Equally, the availability of more consistent and standardized datasets is indispensable to facilitate model comparisons between various OCR models.

3. Proposed Architecture

The block diagram of the novel architecture proposed for paragraph-level handwritten text recognition in this paper is shown in Figure 1. The proposed architecture begins by processing a handwritten paragraph image, which is divided into fixed-size patches and embedded into feature vectors. Positional encodings are also added to preserve spatial relationships, and a Vision Transformer (ViT) encoder extracts contextual visual features using multi-head self-attention. The extracted features are then flattened and pro-

jected into a sequence-compatible representation for the transformer decoder, which autoregressively generates text tokens using masked self-attention and cross-attention. The decoder aligns the generated tokens with ViT-extracted embeddings to ensure accurate transcription. A pre-trained language model (LM) further refines the generated text, enhancing fluency, coherence, and grammatical accuracy. The sequence is structured to maintain paragraph formatting and training is performed using cross-entropy loss combined with sequence-level optimization. Finally, the model provides a well-structured, refined transcription of the handwritten paragraph.

The proposed framework addresses the challenges of low-resource handwriting, such as varying cursive styles and inconsistent skews/slants, without relying on explicit segmentation. Moving away from traditional segmentation methods, the proposed model uses attention mechanisms to enhance text understanding and decoding, thereby improving efficiency and reducing recognition errors.

3.1. Pre-processing

In our proposed framework, we performed pre-processing across all datasets to maintain model consistency and optimal input quality. We implemented a down-sampling technique using bi-linear interpolation to halve the resolution of images to 150 dpi, consistently applying this specification across datasets. The down-sampled images were then resized to 224x224 for pre-training and 448x448 for fine-tuning. This uniform approach to pre-processing ensured that the model was trained and evaluated under consistent conditions, enhancing its ability to generalize effectively across diverse datasets.

3.2. Image Encoder

The proposed network uses DeiTForMaskedImageModeling, an adaptation of the Data-efficient Image Transformer (DeiT) constructed for complex advanced masked image modeling challenges [4]. This model is built on transformer architecture to extract intricate visual features from images.

3.2.1. Knowledge Distillation

Our proposed model uses a data-efficient image transformer as the image encoder. A knowledge distillation technique [4], as shown in Figure 2, is used to pre-train this encoder. In knowledge distillation using a Vision Transformer [3], as the teacher model, the "distill token" emerges as a pivotal concept. This token is designed to capture and incorporate the distilled

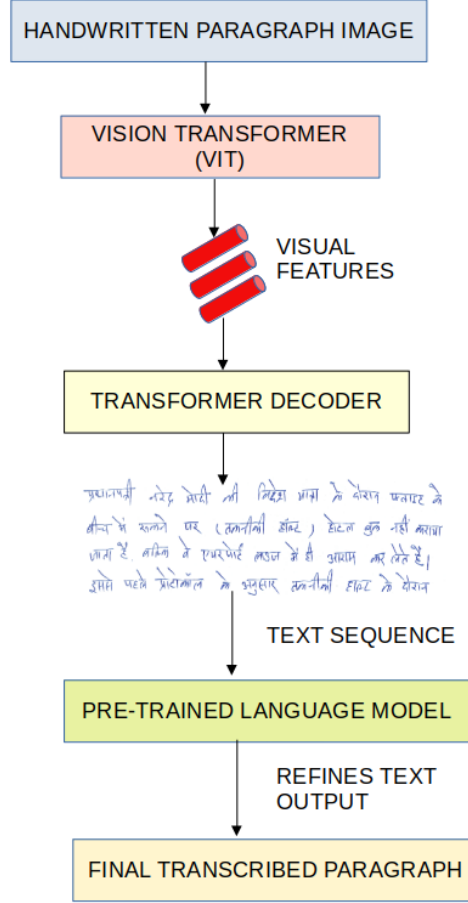


Figure 1: The architecture consists of an image encoder, a transformer decoder, and a pre-trained language model. The image encoder extracts visual features from paragraph images, the transformer decoder generates an initial text sequence, and the language model refines it for coherence and accuracy.

knowledge from the teacher model. During the training phase, it interacts with image patch embedding through the transformer’s self-attention mechanism to closely mirror the teacher’s output [1]. The distill token’s final state, refined through layers of attention, encapsulates the teacher’s learned representations, thus serving as a compact vessel for transferring knowledge to the student model. This method enhances the student model’s performance by imbuing it with the teacher’s insights without necessitating the computational heft of the teacher model. The distill token illustrates an

efficient strategy to use the complex patterns and relationships learned by larger models, facilitating the deployment of powerful yet lightweight models for various applications.

3.2.2. Image Encoder Architecture

The image encoder architecture processes resized images of 224x224 pixels to, a choice that balances the need for detail retention against computational efficiency. The core of our framework is built upon a deep transformer structure with 12 hidden layers and a hidden size of 768, designed to process complex visual data. It incorporates 12 attention heads to focus on multiple image features, enhancing feature extraction. The model processes the input images by breaking them into 16x16 pixel patches, treating them as sequences, thereby combining the transformer’s sequential processing capabilities with the spatial characteristics of images for more detailed analysis. We use the Gaussian Error Linear Unit (GELU) as the activation function and layer normalization to ensure stable learning dynamics. The model’s intermediate layer size of 3072 allows complex transformations in a higher-dimensional space. The encoder stride of 16 ensures efficient image processing. This DeiTForMaskedImageModeling model exemplifies a sophisticated application of transformers in computer vision, optimizing both performance and data efficiency.

3.3. Attention Mechanism

Attention mechanisms enable models to selectively focus on relevant parts of the input data, enhancing decision-making processes. Originally developed for natural language processing (NLP), these mechanisms help models dynamically prioritize different elements, such as words in a sentence. The concept has been expanded to applications in vision and audio processing. Cross-modal attention [5] is a further evolution of these mechanisms that allows models to integrate and focus on information from different data types, such as text and images. This invention encourages a more comprehensive understanding through the smooth integration of textual and visual features. It is especially valuable in tasks necessitating the synthesis of multi-modal inputs for predictions.

3.4. Language Model

In developing our proposed framework, we adopt a language transformer-based architecture [6] to improve the text sequence obtained from our proposed transformer decoder model. The language model was pre-trained for

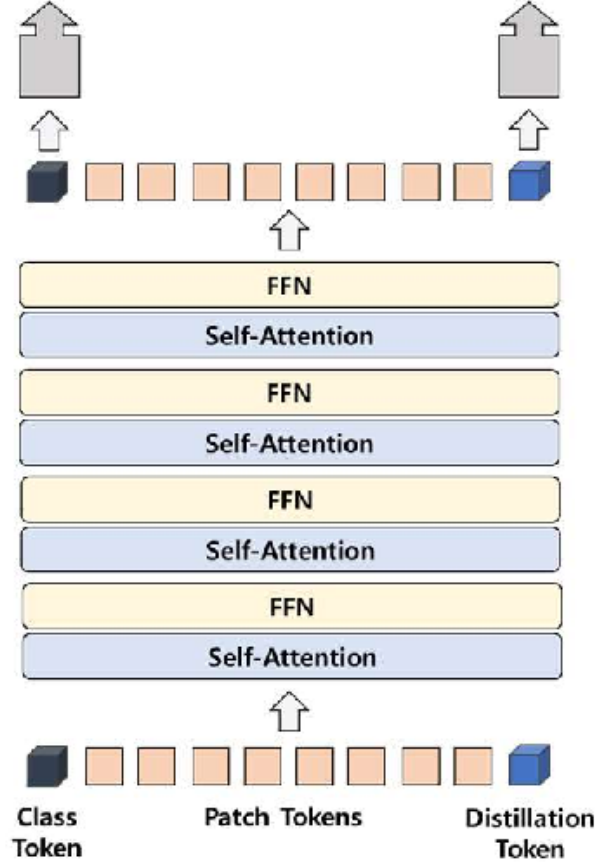


Figure 2: Our distillation procedure by including a new distillation token. It interacts with the class and patch tokens through the self-attention layers. This distillation token is employed similarly to the class token, except that on the network output, its objective is to reproduce the (hard) label predicted by the teacher instead of a true label. Both the class and distillation tokens input to the transformers are learned by back-propagation.

masked language modeling task, as shown in Figure 4, to enhance its linguistic capabilities. The model contains six hidden layers, and each layer is designed with a hidden size of 768 units and 12 attention heads. Each layer has multiple attention heads facilitating attention mechanisms across various input data segments. This configuration is tailored to process sequences with a maximum length of 512 tokens, supported by absolute position embedding to retain the syntactic and semantic order of language input. The architecture incorporates a 10 % dropout rate in the hidden layers and at-

tention probabilities to avoid overfitting and generalizing the model. The model’s activation function of choice, GELU, ensures smooth non-linearity and the ability to model complex relationships in the data. Our model benefits from a vocabulary size of 50,026 tokens to learn the intricacies of the Urdu language. This architectural design emphasizes our model’s capability to learn from masked linguistic contexts and predictively model language with remarkable accuracy and speed.

3.5. *Decoder*

The Transformer decoder, following an autoregressive structure [29], converts visual embeddings extracted by the Vision Transformer (ViT) into textual tokens for paragraph-level handwritten text recognition. The decoder comprises multi-head self-attention, cross-attention for alignment with ViT features, and feedforward layers to model long-range dependencies. During training, teacher forcing is employed, where the ground-truth token at each time step is provided as input for the next step instead of the previously generated token, ensuring stable convergence. The decoder is trained on paired image-text datasets using cross-entropy loss to minimize discrepancies between predicted and ground-truth sequences. To mitigate exposure bias, scheduled sampling is incorporated to gradually reduce reliance on teacher forcing. During inference, the model generates text autoregressively, predicting one token at a time while using its own outputs as subsequent inputs. Decoding strategies, such as beam search, nucleus sampling, and greedy decoding, enhance the quality of text generation. This approach ensures accurate, coherent, and contextually structured paragraph-level text recognition by leveraging self-attention, cross-attention, and language modeling techniques.

3.6. *Fine-tuning The Network*

The whole network, which consists of the image encoder, the transformer decoder, the language model, is fine-tuned using the Parimal Urdu and Hindi datasets, which consists of paragraph images. The whole architecture, when fine-tuning, is shown in Figure 1. During fine-tuning, the Image encoder processes a paragraph image and outputs visual features based on the visual tokens learned during pre-training. The transformer decoder outputs a sequence of text. The pre-trained language model then refines the output text contextually based on learned tokens and provides the final text output.

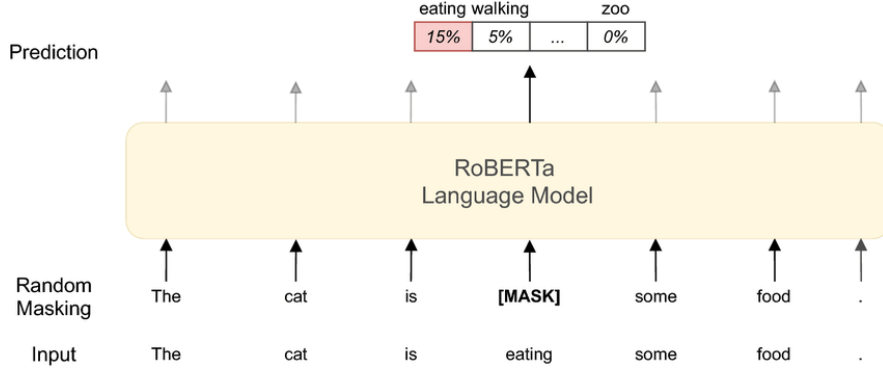


Figure 3: RoBERTa—masked language modeling with the input sentence: The cat is eating some food

4. Experimental Analysis and Discussion

4.1. Datasets

We used various datasets to pre-train our image encoder and language model, fine-tune them on handwritten documents and finally test the proposed framework.

Printed Datasets: To pre-train our DeiT model, we created a dataset of 21,000 images, as shown in Figure 5, from various online sources, focusing on the diversity of Urdu printed text. This dataset includes multiple font styles and layouts to comprehensively represent printed Urdu and Hindi text, which helps the model learn the intricacies of Urdu text recognition. The diversity of data in the training phase is crucial for helping the model generalize across different textual formats and provide a strong foundation for further fine-tuning for specialized tasks, ensuring robust initial training for enhanced performance in advanced applications.

Handwritten Paragraph Dataset (Parimal Urdu and Parimal Hindi): We also used the Parimal Urdu and Parimal Hindi dataset, which consists of 500 pages of handwritten Urdu and Hindi text as shown in Figure 6. It was contributed by 10 individuals from various age groups, each selected for their different handwriting styles shown in Figure 6. This variety introduces 10 unique styles into the dataset, providing a solid basis for evaluating the adaptability and effectiveness of our model with various natural writing variations. This dataset was crucial for the fine-tuning and evaluation phases.

لوگوں کے لیے آپس میں ملتی ہیں اور پھر چلتی جاتی ہیں۔ اللہ نے جو چیزیں کو ان کی منزل کا پتہ دینے کا
 کیا سارا طریقہ عطا کیا ہے!
 جو چیزیں اپنے گروہ کے ساتھیوں کے ساتھ کسی کوئی جھگڑتی نہیں مگر دوسرے گروہ کے
 ساتھ جھگڑا بھی کرتی ہیں اور گروہ جماعت کے کھانے پینے کی چیزوں کو لوٹ لیتی ہیں حتیٰ کہ ان
 کے انڈوں کو بھی اٹھا لیتی ہیں لیکن ان سب کے باوجود اللہ نے ان کے دل میں بڑا ہی رحم عطا کیا
 ہے۔ لوٹے ہوئے انڈوں کو بریافتیں کرتی ہیں بلکہ انڈوں اور اس کے بچوں کو پال پوس کر اپنے
 گروہ میں شامل کر لیتی ہیں۔
 ایک حکماء ہے کہ جو کچھ کی آواز عرش تک جاتی ہے۔ یعنی خرب و عظام کی آواز خدا تک
 پہنچتی ہے۔ اسی لیے ہم اگھے بچوں کو بھی کچھ سے ضرور آواز کوئی تکلیف نہیں پہنچانی چاہیے۔
 کچھ تو یہ ہے کہ جو کچھ کی زندگی ہم سب کے لیے دوسرے جہت ہے۔ کیا ہم سب جو کچھ سے سبق
 حاصل کرنے کے لیے تیار ہیں!
 تمہارے اہل خانہ کے تمام افراد کو حسب مراتب سلام و دعا۔ اپنے اس خط کے جواب کا
 انتظار رہے جیسی سے کریں گے۔

اسی ہم نے دیکھا ہے کہ فرقہ کے فطری اور دائمی پہلو کیا ہیں۔ دوسرے معنی میں تمام
 نسل انسانی کو ایک فرقہ سے تعبیر کیا جاسکتا ہے۔ لیکن نسل انسانی کے جسے ایک
 دوسرے سے کئی لحاظ سے منقطع ہوئے ہیں جیسے: مقامی ملحدگی، زبانوں کا اختلاف
 مذہب، تعلیم، زندگی بسر کرنے کے طریقے اور بہت سے دوسرے حالات جو اس قسم
 کی ہم خدائی میں سامنے ہوتے ہیں جو حقیقی انسانی میل جول کے لیے ضروری ہے۔
 وہ لوگ بھی جو نسبتاً قریب اور ساتھ رہتے ہیں اکثر اوقات ایک دوسرے کے
 Prof-
 عدا مشابہت بہ اعتبار نوع کو میل جول کی بنیاد کی حیثیت سے پروردگار کو شکر
 Guiding- کا کتاب Principles of Sociology میں دیکھا جاسکتا ہے
 مگر ہمارا خیال ہے کہ اس نے نوع کی مشابہت کو غیر ضروری طور پر اہم قرار دیا ہے اور ہاں پر
 نوع کی مشابہت ہم خیالی سے ملحدہ ہے۔

Figure 4: Printed Urdu Text Paragraphs from different books of science and philosophy

1
 (1) جی دھلی، ۸۸، نو مہر: زبان: تاریخ میں ۹ نو مہر کی تاریخ: انراکھنڈ سے یوم
 تائیں سے طور پر درج ہے ایک علیحدہ انراکھنڈ کے مطالعہ کے کچھ سرائوں
 اضلاع سے جو انراکھنڈ کو رانا ۹ نو مہر کو ۲۰۰۰ کو سنائے ہوئے ہیں۔ راسن کے طور پر
 جو ۱۰۰۰ میں سنائے گئے ہیں ۲۰۰۰ تک اسے انراکھنڈ سے نام جانا تھا لیکن
 نام بدل کر انراکھنڈ کر دیا گیا۔ انراکھنڈ کے ۲۰۰۰ تک اسے انراکھنڈ سے نام جانا تھا لیکن
 اور متفرق میں نیپال سے متصل مغرب میں مہاجل پردیش اور جنوب میں
 انراکھنڈ میں اس کی سرحدی ریاستیں ہیں۔

جی دھلی، ۸۸، نو مہر: زبان: تاریخ میں ۹ نو مہر کی تاریخ: انراکھنڈ سے یوم
 تائیں سے طور پر درج ہے ایک علیحدہ انراکھنڈ کے مطالعہ کے کچھ سرائوں
 اضلاع سے جو انراکھنڈ کو رانا ۹ نو مہر کو ۲۰۰۰ کو سنائے ہوئے ہیں۔ راسن کے طور پر
 جو ۱۰۰۰ میں سنائے گئے ہیں ۲۰۰۰ تک اسے انراکھنڈ سے نام جانا تھا لیکن
 نام بدل کر انراکھنڈ کر دیا گیا۔ انراکھنڈ کے ۲۰۰۰ تک اسے انراکھنڈ سے نام جانا تھا لیکن
 اور متفرق میں نیپال سے متصل مغرب میں مہاجل پردیش اور جنوب میں
 انراکھنڈ میں اس کی سرحدی ریاستیں ہیں۔

Figure 5: Handwritten text paragraphs from two authors with quite different writing styles from the Parimal Urdu dataset.

Text Data:. In our research, the RoBERTa model was subjected to an extensive pre-training phase using a carefully curated corpus of Urdu text. This corpus, rich in diversity, consists of 20 million lines of text, amounting to approximately 200 million words, drawn from various sources, including news articles and magazines. This selection aimed to capture the linguistic richness and variability of the Urdu script. The sheer volume and variety of the dataset provided a solid foundation for the model, enabling it to acquire a deep understanding of Urdu syntax, semantics, and contextual nuances. Such comprehensive pre-training is crucial, as it significantly enhances the model’s performance across various NLP tasks.

Public Handwritten Text Lines Dataset:. PUCIT-OHUL [23] dataset consists of multiple writing styles, different types of pen, ink types, text sizes

अक्सर अपने विवादित बयानों के कारण नर्चों में रहने वाली शास्त्री प्रताप सिंह ठाकुर ने एक बार फिर ऐसा बयान दिया है जो उन्हें विवादों के केंद्र में ले आया है। बुधवार को उन्होंने लोकसभा में महात्मा गांधी के हत्यारे नाथूराम गोडसे को दोबारा देशभक्त कह दिया। जिसपर विपक्ष काफी हंगामा कर रहा है। सरकार ने भी अंसद में अफाई देते हुए कहा कि यह बयान भिन्न है। वहीं अब आरक्षी प्रताप ने अपने बयान को लेकर द्विदर पर डीट करते हुए लिखा है कि मैंने कल ऊधम सिंह जी का अपमान नहीं अहा। उनका कहना है कि उनके बयान को गलत तरीके से लिया गया और वह स्वतंत्रता सेनानी ऊधम सिंह का जिक्र कर रही थी। उन्होंने लिखा, कभी-कभी झूठ का बवंडर इतना गहरा होता है कि दिन में भी रात लगाने लगती है। किंतु सूर्य अपना प्रकाश नहीं खोता। पलभर के बवंडर में लोग भ्रमित न हों। सूर्य का प्रकाश रुकाई है। अलग भरी है कि कल मैंने ऊधम सिंह जी का अपमान नहीं अहा। बस।

प्रधानमंत्री नरेद्र मोदी जी निदेश प्राप्त के दौरान प्रताप के बीच में खलने पर (तकनीकी हॉल) होटल बुक नहीं कराया जाया है। बकिंग के एयरपोर्ट लडज में ही आराम कर लेते हैं। इससे पहले प्रोटोकॉल के अनुसार तकनीकी हॉल के दौरान प्रधानमंत्री के सकेने के लिए होटल बुक होता था। पीएम मोदी अपनी लंबी निदेश प्राप्त के दौरान अलग गलमानों के लिए छोटे प्रानदंड के स्थापित कर रहे हैं। इसकी जानकारी खुद गृहमंत्री अमित शाह ने लोकसभा में एक नर्च के दौरान दी। बेपुखत प्रोटेक्शन ग्रुप (एसपीजी) संशोधन विधेयक पर नर्च के दौरान प्रभावित शाह ने कहा कि पीएम मोदी ने आज तक किसी भी देश में तकनीकी हॉल के दौरान उपपने लिए होटल बुक करने के निदेश नहीं दिए। एयरपोर्ट पर ही रुकते हैं, लडा नडाते हैं और विमान में स्थान भरने के बाद प्रिगल होते हैं।

Figure 6: Handwritten text paragraphs from two authors with quite different writing styles from the Parimal Hindi dataset.

and background types and colors. It was collected from 100 students between 20 and 24 years of age. A total of 479 pages of text were collected. Each page was scanned at 200 DPI and text lines were manually segmented. The dataset contains a total of 7401 text lines and 80,059 words. **NUST-UHWR** [18] dataset is obtained from various websites, including social networks and news websites, and contains 10,606 samples of handwritten Urdu text lines.

4.2. Data Augmentation

To improve accuracies obtained from our proposed model and prevent overfitting, we developed a comprehensive data augmentation strategy, which was implemented only during the training phase. This strategy includes a variety of techniques, such as resolution modification, perspective changes, elastic distortion, and adjustments in brightness and contrast, each with a 0.2 probability of use. This systematic and diversified augmentation approach significantly enhances the model's ability to generalize features and adapt to real-world data variations, which is crucial for robust pattern recognition performance.

4.3. Metrics

To evaluate the performance of our UHTR model, we considered three commonly used metrics, Character Error Rate (CER), Word Error Rate

(WER), and Line Error Rate (LER), which indicate the accuracy of the proposed model in recognizing individual characters, whole words, and Lines, respectively. These metrics are calculated using the Levenshtein distance (denoted by lev^d), which measures the difference between the ground truth (\mathbf{y}) and the recognized text ($\hat{\mathbf{y}}$), normalized by the total length of the ground truth ($\text{len}(\mathbf{y})$). This normalization process ensures that discrepancies in the shortest lines are proportionately weighted against those in longer segments, providing a fair metric across texts of varying lengths. The character, word and line recognition rates were calculated by subtracting their respective error rates from 100. The following equations were used to calculate CER, WER and LER.

$$\text{CER} = \frac{\sum_{i=1}^K \text{len}_i^y}{\sum_{i=1}^K \text{lev}_i^d(\hat{y}_i, y_i)} \quad (1)$$

$$\text{WER} = \frac{\sum_{i=1}^K \text{word count in } y_i}{\sum_{i=1}^K \text{lev}_i^d(\text{words in } \hat{y}_i, \text{words in } y_i)} \quad (2)$$

Here, K represents the total number of images in the dataset. For WER, punctuation characters are treated as independent words aligned with conventions.

$$\text{LER} = \frac{1}{N} \sum_{i=1}^N \text{lev}^d(\hat{y}_i, y_i) \quad (3)$$

Where N is the total number of sentences in the dataset and $\text{lev}(\hat{y}_i, y_i)$ is the Levenshtein distance between the predicted line \hat{y}_i and the ground truth line y_i for the i -th sentence. The denominator, $\sum_{i=1}^N 1$, counts the total number of sentences in the ground truth, providing a normalization factor for the error rate.

We also employed two metrics to assess our model’s text segmentation accuracy: Intersection over Union (IoU) and mean Average Precision (mAP). IoU measures the overlap of pixels classified as text versus the union of such pixels compared to the ground truth, normalized across images. mAP calculates average precision across IoU thresholds from 50% to 95% in 5% increments, weighted by pixel count for a dataset-wide score.

4.4. Experimental results

We first used the Parimal Urdu dataset comprising 500 images of hand-written Urdu and Hindi paragraphs for experimentation. The dataset was methodically partitioned into training, testing, and validation subsets, with

Table 1: Recognition Rates (%) of the Proposed Model on Paragraph Urdu Data

Models	CRR		WRR		LRR	
	Val	Test	Val	Test	Val	Test
Proposed HWR	94.80	95.20	83.60	84.70	72.78	73.24

Table 2: Recognition Rates (%) of the Proposed Model on Paragraph Hindi Data

Models	CRR		WRR		LRR	
	Val	Test	Val	Test	Val	Test
Proposed HWR	80.64	78.20	70.60	67.65	54.78	57.24

respective proportions of 70%, 20%, and 10%. This split ensures a balanced approach, allowing for comprehensive training while retaining sufficient data to test and validate the models. The results were evaluated based on three key metrics, CRR, WRR and LRR. The results of these experiments are presented in Table 1 and Table 2.

4.5. Comparative Study

To compare the performance of our proposed model with state-of-the-art models we considered nine different methods from the literature. The experiments performed on two public datasets for Urdu text recognition, the NUST-UHWR and PUCIT-OHUL datasets, and the results were calculated based on CRR and WRR. The results obtained from the prior models and our proposed model On the NUST-UHWR dataset is shown in Table 3. From Table 3, it is clear that our proposed model performs better and sets a new benchmark. Table 4 shows a comparison study on the PUCIT-OHUL dataset with existing models. The results indicate that our proposed model outperforms existing models in both CRR and WRR. In the case of character recognition, the proposed model is superior to any existing models, showing that visual and textual tokens are mapped appropriately. The existing models based on CNNs do not perform well on Urdu text because they cannot identify the intricacies of Urdu handwritten text recognition. The models based on CNNs and transformers also suffer as they cannot process the contextual information in the text. It is important noting that as the proposed approach uses a language model to derive contextual information, and a multi-modal attention network to extract visual details, it improved

Table 3: Comparison of different models for offline Urdu handwritten text recognition on NUST-UHWR Dataset

Models	CRR%	WRR%
BLSTM [12]	72.60	61.37
Modified CRNN [13]	81.50	70.60
MDLSTM [14]	85.87	74.60
CNN-RNN [15]	86.75	76.05
BiGRU [16]	86.30	75.02
TrOCR [17]	79.88	68.53
Conv. Recursive [18]	92.75	81.88
Conv. Transformer [19]	94.03	83.13
Unified Arch. [20]	94.10	83.24
Proposed	96.24	85.3

Table 4: Comparison of different models for offline Urdu handwritten text recognition on PUCIT-OHUL dataset

Models	CRR%	WRR%
CNN-GRU	33.00	22.44
CNN-LSTM	32.04	21.42
CNN-BGRU	32.76	23.16
CNN-BLSTM [7]	32.02	22.23
SimpleHTR [9]	14.94	5.09
LineHTR [8]	30.05	18.19
CRNN [10]	32.86	24.18
CALText[11]	82.06	51.97
Proposed	92.05	80.54

handwritten Urdu text recognition. In addition, the proposed model suffers in word recognition compared to character recognition due to the use of inconsistent spacing, which is quite natural in Urdu scripts.

5. Conclusion and Future Scope

In segmentation-based approaches, the recognition of characters, words, and lines was severely affected due to wrong manual segmentation, as the Urdu script follows a cursive writing style. This study introduces a novel segmentation-free approach to recognizing handwritten Urdu datasets using multi-modal attention. This method is further refined by fine-tuning

with targeted datasets for handwritten Urdu recognition. In addition, a self-collected dataset, Parimal Urdu, is introduced in this research work. The proposed framework is evaluated on three different datasets, including Parimal, PUCIT-OHUL and NUST-UHWR. Notably, the new model achieved recognition accuracies that were higher than those of previous methods. In fact, it achieves state-of-the-art results on these datasets. Its implicit line segmentation process enables the recognition of inclined lines and paragraphs. For future work, this model can be extended to be applied to more diverse text recognition tasks, such as Patwari Urdu, Arabic, and Persian text recognition.

References

- [1] Vaswani A., Shazeer N., Parmar N., Uszkoreit J., Jones L., Gomez A. N., Polosukhin I. (2017). "Attention Is All You Need." In Advances in Neural Information Processing Systems (NeurIPS).
- [2] Devlin J., Chang M. W., Lee K., Toutanova K. (2019). "BERT: Pretraining of Deep Bidirectional Transformers for Language Understanding." In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL).
- [3] Dosovitskiy A., Beyer, L., Kolesnikov A., Weissenborn D., Zhai X., Unterthiner T., Houlsby N. (2021). "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale." Proceedings of the International Conference on Learning Representations (ICLR).
- [4] Hugo T. et al. (2021) "Training data-efficient image transformers and distillation through attention." International Conference on Machine Learning. PMLR.
- [5] Peng X., Zhu X. and Clifton D. A. (2023) "Multimodal learning with transformers: A survey." IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [6] Liu Y., et al. (2019) "RoBERTa: A robustly optimized BERT pretraining approach." arXiv preprint arXiv:1907.11692.

- [7] Hassan S., Irfan A., Mirza A., Siddiqi I. (2019) Cursive handwritten text recognition using bi-directional LSTMs: A case study on Urdu handwriting. In: 2019 International Conference on Deep Learning and Machine Learning in Emerging Applications (Deep-ML), pp. 67–72.
- [8] Lcm HT (2022) Line-level handwritten text recognition with TensorFlow. <https://github.com/lamhoangtung/LineHTR>.
- [9] Scheidl H (2022) Handwritten text recognition with TensorFlow. <https://github.com/githubharald/SimpleHTR>.
- [10] Shi B., Bai X., Yao C. (2016) An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE Trans Patt Anal Mach Intell* 39(11), pp. 2298–2304.
- [11] Tayaba A. and Khan N. (2023) "CALText: Contextual attention localization for offline handwritten text." *Neural Processing Letters*, 55(6), pp. 7227-7257.
- [12] Ul-Hasan A., et al. (2013) "Offline printed Urdu Nastaleeq script recognition with bidirectional LSTM networks." *Proceedings of the 12th International Conference on Document Analysis and Recognition*.
- [13] Baoguang S., Bai X. and Yao C. (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), pp. 2298-2304.
- [14] Graves, Alex, and Jürgen Schmidhuber. (2008) "Offline handwriting recognition with multi-dimensional recurrent neural networks." *Advances in neural information processing systems* 21.
- [15] Safarzadeh, Mohammadi V. and Jafarzadeh P. (2020) "Offline Persian handwriting recognition with CNN and RNN-CTC." *Proceedings of the 25th international computer conference, computer society of Iran (CS-ICC)*.
- [16] Chen L., Yan R., Peng L., Furuhashi A., Ding X. (2017) "Multi-layer Recurrent Neural Network based Offline Arabic Handwriting Recognition." *Proceedings of the 1st International Workshop on Arabic Script Analysis and Recognition (ASAR)*, pp. 6-10.

- [17] Li M. et al. (2023) "Trocrr: Transformer-based optical character recognition with pre-trained models." *Proceedings of the AAAI Conference on Artificial Intelligence*. 37(11).
- [18] Zia N., Naeem M. F., Raza S. M. K., Khan M. M., Ul-Hasan A., Shafait F. (2022) A Convolutional Recursive Deep Architecture for Unconstrained Urdu Handwriting Recognition. *Neural Computing and Applications*, 34(2), pp. 1635-1648.
- [19] Riaz N., Arbab H., Maqsood A., Nasir K. B., Ul-Hasan, A., Shafait, F.(2022) ConvTransformer Architecture for Unconstrained Off-Line Urdu Handwriting Recognition. *International Journal on Document Analysis and Recognition (IJDAR)*, volume 25, pp. 373–384.
- [20] Arooba M. et al. (2023) "A Unified Architecture for Urdu Printed and Handwritten Text Recognition." *International Conference on Document Analysis and Recognition*. Cham: Springer Nature Switzerland.
- [21] Satti DA, Saleem K. (2012) "Complexities and implementation challenges in offline urdu nastaliq OCR." *Proceedings of the Conference on Language and Technology*, pp 85–91
- [22] Sagheer M.W., He C.L., Nobile N., Suen C.Y. (2010) "Holistic Urdu handwritten word recognition using support vector machine." *Proceedings of the 20th International Conference on Pattern Recognition*, pp. 1900–1903.
- [23] Anjum T, Khan N (2020) An attention based method for offline handwritten Urdu text recognition. In: 2020 17th International Conference on Frontiers in handwriting recognition (ICFHR), pp 169–174.
- [24] D. Chaudhary and K. Sharma, "Hindi Handwritten Character Recognition using Deep Convolution Neural Network," 2019 6th International Conference on Computing for Sustainable Global Development (INDIA-Com), New Delhi, India, 2019, pp. 961-965.
- [25] Sharma, Anand, and A. G. Ramakrishnan. "Structural analysis of Hindi online handwritten characters for character recognition." *arXiv preprint arXiv:2310.08222* (2023).

- [26] Chauhan, V.K., Singh, S. & Sharma, A. HCR-Net: a deep learning based script independent handwritten character recognition network. *Multimed Tools Appl* 83, 78433–78467 (2024).
- [27] Rajpal, Danveer, et al. "A fusion-based hybrid-feature approach for recognition of unconstrained offline handwritten Hindi characters." *Future Internet* 13.9 (2021): 239.
- [28] Reddy, R. Vijaya Kumar, and U. Ravi Babu. "Handwritten Hindi character recognition using deep learning techniques." *International Journal of Computer Sciences and Engineering* 7.2 (2019): 1-7.
- [29] Radford, Alec, et al. "Improving language understanding by generative pre-training." (2018).
- [30] Rao, N. Venkata, et al. "OPTICAL CHARACTER RECOGNITION TECHNIQUE ALGORITHMS." *Journal of Theoretical & Applied Information Technology* 83.2 (2016).
- [31] Chauhan, Rahul, Kamal Kumar Ghanshala, and R. C. Joshi. "Convolutional neural network (CNN) for image detection and recognition." 2018 first international conference on secure cyber computing and communication (ICSCCC). IEEE, 2018.
- [32] Graves, Alex, and Alex Graves. "Long short-term memory." *Supervised sequence labelling with recurrent neural networks* (2012): 37-45.
- [33] Bao, Yu, et al. "Region-based CNN for logo detection." *Proceedings of the International Conference on Internet Multimedia Computing and Service*. 2016.
- [34] R. Sahay and M. Coustaty, "An Enhanced Prototypical Network Architecture for Few-Shot Handwritten Urdu Character Recognition," in *IEEE Access*, vol. 11, pp. 33682-33696, 2023.
- [35] A. Gaur and S. Yadav, "Handwritten Hindi character recognition using k-means clustering and SVM," 2015 4th International Symposium on Emerging Trends and Technologies in Libraries and Information Services, Noida, India, 2015, pp. 65-70.

- [36] Omayio, Enock Osoro, Sreedevi Indu, and Jeebananda Panda. "Word spotting and character recognition of handwritten Hindi scripts by Integral Histogram of Oriented Displacement (IHOD) descriptor." *Multimedia Tools and Applications* 83.1 (2024): 1379-1406.