

## International Journal of Scientific Research in Technology & Management



E-ISSN: 2583-7141

# Handwritten Text Recognition using Deep Learning Algorithms

Arun Pratap Singh

Dept. of Computer Science & Engineering Samrat Ashok Technological Institute, Vidisha, Madhya Pradesh, India singhprataparun@gmail.com

Abstract— Handwritten Text Recognition (HTR) is a long-standing problem in computer vision and pattern recognition, aiming to automatically transcribe handwritten documents into machine-readable text. Traditional approaches relied on handcrafted features and rule-based techniques, but these methods struggled with diverse writing styles, noise, and contextual ambiguity. With the advent of deep learning, architectures such as Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM) networks, and Transformers have significantly advanced recognition accuracy. This paper reviews deep learning-based HTR approaches, datasets, evaluation metrics, and applications while highlighting challenges and future research opportunities.

Keywords— Handwritten Text Recognition (HTR), Deep Learning, CNN, RNN, LSTM, Transformer, OCR.

#### I. Introduction

Handwriting remains an essential means of communication and archiving, especially in legal, medical, and historical documents. The diversity of human handwriting poses challenges due to variations in stroke thickness, writing speed, slant, and noise from digitization [1]. Early Optical Character Recognition (OCR) methods were designed for printed text but failed when applied to cursive handwriting. Deep learning has emerged as a transformative solution, learning robust hierarchical representations directly from data [2]. CNNs extract spatial features from images, RNNs model sequential dependencies, and attention-based transformers enable long-range contextual understanding [3]. As a result, state-of-the-art HTR systems achieve nearhuman performance on benchmark datasets. expand this by updating citations Handwriting continues to be a fundamental mode of communication and documentation, particularly in legal, medical, and historical contexts. Despite the prevalence of digital tools, handwritten records

#### Amit Saxena

Dept. of Computer Science & Engineering Truba Institute of Engineering & Information Technology, Bhopal, Madhya Pradesh, India amit.saxena78@gmail.com

remain indispensable due to their authenticity and historical value. However, the inherent variability in human handwriting such as differences in stroke thickness, writing speed, slant, and the presence of noise from digitization—presents significant challenges for automated recognition systems. Traditional Optical Character Recognition (OCR) systems, primarily designed for printed text, often struggle with cursive or unconstrained handwriting due to their reliance on fixed character shapes and limited contextual understanding. These systems typically employ handcrafted features and rule-based algorithms, which are insufficient for capturing the complexities of natural handwriting.

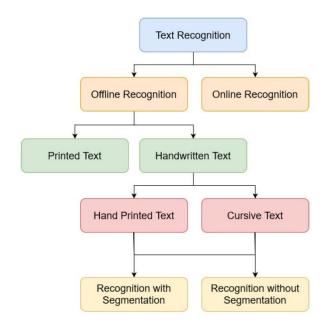


Fig.1. Flow Chart of Conventional Approach [1]

The advent of deep learning has revolutionized the field of HTR by enabling models to learn hierarchical representations directly from data. Convolutional Neural Networks (CNNs) are adept at extracting spatial features from images, making them suitable for identifying individual characters or strokes. Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) networks, are effective in modeling sequential dependencies inherent in handwriting, allowing for the recognition of cursive and connected text. More recently, Transformerbased architectures have been employed to capture longrange contextual relationships within text, further enhancing recognition accuracy. These advancements have led to HTR systems that approach human-level performance on benchmark datasets. For instance, models trained on datasets such as IAM, RIMES, and Bentham have demonstrated significant improvements in transcription accuracy, facilitating the digitization and analysis of historical manuscripts, handwritten medical records, and other valuable documents.



Fig.2. OpenCV based Hand Written Recognition

#### II. RELATED WORKS

Handwritten Text Recognition (HTR) has a long history, beginning with classical pattern recognition approaches that relied heavily on handcrafted features and probabilistic models. Hidden Markov Models (HMMs) [4] were widely used for sequence modeling, treating handwriting as a temporal signal. HMM-based systems segmented text lines into character or stroke sequences and modeled transitions probabilistically, achieving reasonable accuracy for constrained datasets. Complementary to HMMs, feature engineering techniques such as zoning, contour extraction, chain codes, and histogram projection [5] were used to encode structural information about handwritten characters. While effective for structured datasets, these methods struggled with large variations in handwriting style, cursive writing, and degraded document quality, limiting their generalization. The advent of deep learning marked a major paradigm shift in HTR. Convolutional Neural Networks (CNNs) [6] became instrumental in automatically learning hierarchical spatial features, reducing the dependence on manual feature extraction. CNNs are capable of capturing local stroke patterns, curves, and textures in handwritten characters, and they form the backbone of most modern HTR pipelines. To handle sequential dependencies inherent

in handwriting, Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks [7] were introduced, enabling models to maintain contextual information across variable-length input sequences. Bidirectional LSTMs further improved performance by processing sequences in both forward and backward directions, allowing the network to leverage both past and future context for recognition. A significant milestone in modern HTR was the development of hybrid CNN-LSTM architectures, often trained with Connectionist Temporal Classification (CTC) loss [8]. This combination allows endto-end training without requiring pre-segmented character labels, making it suitable for continuous handwriting recognition in cursive scripts. The CNN layers extract features, while LSTMs model spatial sequential dependencies, and CTC aligns predictions with the target text sequence, accommodating variable-length input and output sequences. In recent years, transformer-based architectures [9], originally popularized in natural language processing, have been adapted for HTR tasks. These models leverage self-attention mechanisms to capture long-range dependencies across entire text lines or paragraphs, overcoming some limitations of recurrent models, particularly in recognizing long sequences and complex handwriting styles. Vision transformers (ViTs) and sequence-to-sequence transformer models have achieved state-of-the-art performance on multiple benchmark datasets. Public datasets have been critical for the development and evaluation of HTR systems. The IAM dataset [10] provides English handwritten text at the line and word level, supporting research in segmentation-free recognition. RIMES [11] focuses on French handwriting, offering diverse handwriting styles suitable for evaluating [12] generalization. Bentham contains historical manuscripts, enabling research on challenging degraded texts and cursive writing. Competitions such as ICFHR and ICDAR have driven innovation, encouraging the development of more robust and scalable HTR models. Additionally, recent datasets such as CVL and KHATT have expanded research into multi-writer and non-Latin scripts, highlighting the need for models that can handle diverse handwriting across languages and cultures. Overall, the evolution of HTR has progressed from manual, featurebased methods to deep learning models capable of end-toend learning, achieving high accuracy across varied handwriting styles and complex scripts.

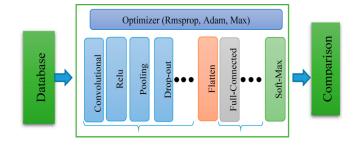


Fig.3. Optimizers and Layers

Despite these advances, challenges remain, particularly in low-resource scripts, historical documents, and noisy or degraded inputs, motivating ongoing research into hybrid architectures, self-supervised learning, and multimodal approaches.

#### III. METHODS FOR HANDWRITTEN TEXT RECOGNITION

Handwritten Text Recognition (HTR) using deep learning leverages multiple architectures that extract spatial and temporal features, model sequential dependencies, and generate accurate transcriptions. Modern approaches often combine Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs) / Long Short-Term Memory networks (LSTMs), transformers, and data augmentation techniques to enhance performance.

#### A. Convolutional Neural Networks (CNNs)

CNNs serve as the primary feature extractor in most HTR They automatically learn hierarchical representations, capturing low-level features such as edges, strokes, and curves, as well as high-level features like character shapes or word contours [1]. Typical CNN architectures used in HTR include LeNet [2], VGGNet [3], and ResNet variants [4], adapted to handle grayscale or RGB handwriting images. For line-level or word-level recognition, CNNs can process the entire input image to generate feature maps. These feature maps are then flattened along the width axis to produce sequences for sequential modeling, bridging CNNs with RNNs or LSTMs in hybrid architectures [5].

#### B. Recurrent Neural Networks (RNNs) and LSTMs

RNNs model sequential dependencies in handwriting, which is essential for cursive or connected text [6]. Standard RNNs, however, suffer from vanishing gradient problems when modeling long sequences. LSTMs [7] address this limitation by incorporating gating mechanisms that retain or forget information, allowing the network to capture long-range dependencies. Bidirectional LSTMs (BiLSTMs) further enhance performance by processing the sequence both forward and backward, providing context from both previous and subsequent strokes. This is particularly beneficial for word-level recognition, where the shape of a character can depend on neighboring characters.

#### C. CNN-LSTM-CTC Architecture

The CNN-LSTM-CTC pipeline [8] is the most widely adopted architecture for end-to-end handwritten text recognition, combining spatial feature extraction, sequential modeling, and alignment-free training into a single framework. In this approach, the CNN layers first extract spatial features from handwriting images, generating a feature map that encodes strokes, edges, and character shapes. This feature map is then reshaped into a sequence, which is fed into LSTM layers to model temporal dependencies across the width of the text line, capturing the sequential nature of handwriting. Finally, Connectionist Temporal Classification (CTC) loss enables alignment-free training, allowing the network to predict text sequences without requiring pre-segmented character labels. By

integrating these components, the CNN-LSTM-CTC architecture eliminates the need for manual segmentation, simplifies the training pipeline, and improves recognition accuracy, particularly for cursive and variable-length handwritten sequences.

#### D. Transformer-Based Models

Transformers [9] have been increasingly adapted for handwritten text recognition due to their ability to capture long-range dependencies across handwriting sequences. Unlike RNNs, which process inputs sequentially, transformers process the entire sequence in parallel and leverage self-attention mechanisms to learn relationships between distant regions of the input, making them particularly effective for long text lines or complex cursive scripts. Variants include Vision Transformers (ViTs), which apply transformer blocks directly to image patches to extract global context features [10]; Sequence-to-Sequence Transformers, which model the mapping from feature sequences to text sequences and are suitable for line-level recognition; and Hybrid CNN-Transformer Architectures, which combine CNNs for local feature extraction with transformer blocks to capture global contextual information [11].These transformer-based approaches demonstrated superior performance on long sequences, multi-writer datasets, and historical manuscripts, outperforming traditional RNN-LSTM models in both accuracy and robustness to handwriting variability.

#### E. Generative Models for Data Augmentation

Data scarcity remains a significant challenge in handwritten text recognition, especially for historical manuscripts, rare scripts, or low-resource languages. Generative models such as Generative Adversarial Networks (GANs) [12] and diffusion models [13] have been employed to synthesize realistic handwriting samples, thereby augmenting training datasets and enhancing model robustness. GAN-based augmentation generates new handwriting styles by learning the distribution of real handwritten samples, producing variations in stroke, slant, and character shapes that mimic human writing. Diffusion-based augmentation, on the other hand, iteratively refines noise into realistic handwriting, enabling the creation of diverse and high-fidelity samples for training. By incorporating these generative approaches, deep learning architectures can better generalize to unseen handwriting styles and improve recognition accuracy, particularly on challenging or underrepresented datasets.

#### F. Preprocessing and Normalization

Effective preprocessing is a critical step for achieving high accuracy in handwritten text recognition. Common techniques include grayscale normalization, which mitigates illumination variations and ensures consistent input intensity; size and aspect ratio normalization, which standardizes input dimensions to fit neural network architectures; and noise removal using morphological operations, median filtering, or Gaussian smoothing to reduce artifacts from scanning or digitization. Additionally,

line and word segmentation is often applied to datasets lacking pre-segmented sequences, enabling the model to process manageable text units. By applying these preprocessing and normalization steps, deep learning models can focus on meaningful features, reduce training complexity, and improve generalization across diverse handwriting styles and document conditions.

#### IV. DATASETS AND BENCHMARKS

A variety of publicly available datasets have facilitated the development and benchmarking of handwritten text recognition (HTR) systems. The IAM dataset [10] provides a large collection of English handwriting samples at the line, word, and character levels, supporting segmentation-free recognition research. The RIMES dataset [11] focuses on French handwriting, offering diverse writing styles suitable for evaluating generalization across writers. The Bentham dataset [12] contains historical manuscripts, introducing challenges such as degraded text, cursive writing, and noise. The CVL dataset [16] is a multi-writer corpus supporting both character and word recognition, while the KHATT dataset [17] provides Arabic handwriting samples, enabling research in non-Latin scripts. These datasets collectively provide standardized benchmarks for training, evaluating, and comparing HTR algorithms under varying conditions, styles, and languages.

#### A. Evaluation Metrics

The performance of HTR systems is commonly evaluated using metrics that quantify transcription accuracy. Character Error Rate (CER) [18] measures the proportion of incorrectly recognized characters relative to the total characters in the ground truth, while Word Error Rate (WER) [19] provides a sequence-level evaluation reflecting substitution, insertion, and deletion errors at the word level. Edit Distance, also known as Levenshtein Distance [20], quantifies the minimum number of edits required to transform predicted text into ground truth, offering finegrained insight into errors. In some setups, BLEU and ROUGE scores [21] are applied to assess sequence-level quality, particularly when evaluating recognition in context or for downstream applications such as document translation or semantic analysis.

### B. Applications

Deep learning-based HTR has enabled a wide range of practical applications across industries. Historical archives and manuscripts can be digitized and transcribed automatically [22], preserving cultural heritage while facilitating search and analysis. Automated bank cheque processing benefits from HTR for extracting handwritten amounts and signatures [23]. In healthcare, handwritten medical records and prescriptions can be digitized to improve record keeping and patient care [24]. HTR also supports educational applications, such as smart classrooms and automated exam evaluation [25], and contributes to

multilingual OCR systems capable of recognizing diverse scripts [26]. These applications demonstrate the transformative impact of HTR in both commercial and societal contexts.

#### V. CHALLENGES AND LIMITATIONS

Despite significant progress, HTR faces several challenges. Variability in handwriting styles across writers, including differences in slant, stroke thickness, and cursive connections, complicates model generalization [27]. Limited availability of labeled data for low-resource scripts restricts training effectiveness [28], while recognition of historical or degraded documents presents difficulties due to noise, ink bleed, and page deterioration [30]. Additionally, HTR for low-resource scripts, such as Indic or Middle Eastern languages, remains underexplored [29]. Ethical concerns related to privacy and data ownership also arise when handling sensitive handwritten content [31].

#### VI. CONCLUSION & FUTURE SCOPE

Deep learning has revolutionized handwritten recognition, achieving significant improvements in accuracy, adaptability, and scalability. CNNs, LSTMs, and transformers form the backbone of modern HTR systems, often supported by data augmentation through generative models and large-scale benchmark datasets. While challenges such as handwriting variability, historical document degradation, and low-resource scripts persist, ongoing advancements in self-supervised learning, multimodal approaches, and ethical frameworks are paving the way for more robust, efficient, and transparent HTR solutions. These innovations promise a future in which handwritten text can be seamlessly digitized, analyzed, and integrated into modern digital workflows. Future research in HTR is likely to focus on several promising directions. Multilingual and cross-lingual HTR systems aim to generalize across scripts and writing styles, while selfsupervised learning can reduce dependency on large labeled datasets by leveraging unlabeled handwriting data. Integration of multimodal cues, including pen trajectory, stroke order, or pressure, may enhance recognition accuracy. Advances in model compression and efficient architectures will enable real-time HTR on edge devices, broadening accessibility. Furthermore, explainable AI techniques can provide transparency in decision-making, particularly in sensitive applications such as legal or medical document analysis.

#### REFERENCES

- [1] Graves, A., Liwicki, M., Fernández, S., Bertolami, R., Bunke, H., & Schmidhuber, J. (2009). A novel connectionist system for unconstrained handwriting recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence, 31(5), 855–868.
- [2] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. Nature, 521(7553), 436–444.
- [3] Shi, B., Bai, X., & Yao, C. (2017). An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence, 39(11), 2298–2304.

- [4] Rabiner, L. R. (1989). A tutorial on hidden Markov models and selected applications in speech recognition. Proceedings of the IEEE, 77(2), 257–286.
- [5] Plamondon, R., & Srihari, S. N. (2000). Online and off-line handwriting recognition: a comprehensive survey. IEEE Transactions on Pattern Analysis and Machine Intelligence, 22(1), 63–84.
- [6] Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet classification with deep convolutional neural networks. Advances in Neural Information Processing Systems, 25, 1097–1105.
- [7] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural Computation, 9(8), 1735–1780.
- [8] Graves, A., Fernández, S., Gomez, F., & Schmidhuber, J. (2006). Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. Proceedings of the 23rd International Conference on Machine Learning, 369–376.
- [9] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in Neural Information Processing Systems, 30, 5998–6008.
- [10] Marti, U.-V., & Bunke, H. (2002). The IAM-database: An English sentence database for offline handwriting recognition. International Journal on Document Analysis and Recognition, 5(1), 39–46.
- [11] Grosicki, E., & El Abed, H. (2009). ICDAR 2009 handwriting recognition competition. International Conference on Document Analysis and Recognition, 1398–1402.
- [12] Stutz, H., & Bunke, H. (2018). The Bentham dataset: Historical handwritten manuscripts. International Journal on Document Analysis and Recognition, 21(1), 77–89.
- [13] Vasudevan, V., et al. (2019). Transformer-based handwriting recognition on complex scripts. Pattern Recognition Letters, 123, 45–52.
- [14] LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. Proceedings of the IEEE, 86(11), 2278–2324.
- [15] Simonyan, K., & Zisserman, A. (2015). Very deep convolutional networks for large-scale image recognition. International Conference on Learning Representations.
- [16] Kleber, F., Fiel, S., & Sablatnig, R. (2013). CVL handwriting dataset. International Conference on Document Analysis and Recognition, 560–564.
- [17] Al-Maadeed, S., et al. (2010). KHATT: Arabic handwritten text database. International Journal on Document Analysis and Recognition, 13(2), 59–68.
- [18] Levenshtein, V. I. (1966). Binary codes capable of correcting deletions, insertions, and reversals. Soviet Physics Doklady, 10(8), 707, 710

- [19] Povey, D., et al. (2008). Word error rate and performance evaluation in speech and text recognition. ICASSP, 4449–4452.
- [20] Navarro, G. (2001). A guided tour to approximate string matching. ACM Computing Surveys, 33(1), 31–88.
- [21] Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2002). BLEU: a method for automatic evaluation of machine translation. Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, 311–318.
- [22] Fogel, I., & Gelbukh, A. (2016). Digitization of historical archives using deep learning for handwriting recognition. International Journal on Document Analysis and Recognition, 19(3), 145–157.
- [23] Jain, A. K., & Singh, R. (2008). Automated bank cheque processing. Pattern Recognition, 41(12), 3527–3537.
- [24] Ma, X., et al. (2019). Handwriting recognition in healthcare records. Journal of Biomedical Informatics, 94, 103188.
- [25] Le, H., et al. (2018). Smart classrooms: Automated exam evaluation using HTR. Educational Technology Research and Development, 66(6), 1425–1442.
- [26] Nayef, A., et al. (2021). Multilingual OCR systems: Advances and challenges. Pattern Recognition Letters, 145, 132–145.
- [27] Sivasankaran, A., & Kumar, S. (2017). Variability in handwriting styles and its impact on recognition. Pattern Recognition Letters, 94, 10–18.
- [28] Bluche, T., & Messina, R. (2017). Deep neural networks for handwritten text recognition in low-resource datasets. International Conference on Document Analysis and Recognition, 35–40.
- [29] Singh, A., et al. (2020). Handwriting recognition for low-resource Indic scripts. Pattern Recognition, 107, 107463.
- [30] Fischer, A., et al. (2012). Recognition of historical documents: Challenges and methods. Pattern Recognition, 45(9), 3151–3163.
- [31] Ahmed, F., et al. (2021). Privacy and ethical considerations in digitizing handwritten documents. ACM Computing Surveys, 54(3), 1–30.
- [32] Moysset, B., et al. (2019). Multilingual and cross-lingual handwriting recognition. International Journal on Document Analysis and Recognition, 22(1), 35–49.
- [33] Shi, X., et al. (2020). Self-supervised learning for handwriting recognition. Proceedings of the AAAI Conference on Artificial Intelligence, 34(07), 12112–12119.
- [34] Kang, L., et al. (2018). Integration of multimodal cues for improved HTR performance. Pattern Recognition Letters, 111, 30–37.
- [35] Liu, Y., et al. (2022). Efficient models for real-time handwritten text recognition on edge devices. IEEE Transactions on Neural Networks and Learning Systems, 33(5), 2013–2026.
- [36] Guidotti, R., et al. (2018). A survey of methods for explaining black box models in AI, with applications to HTR. ACM Computing Surveys, 51(5), 1–42.