# Advancements In Deep Learning-Based Natural Language Processing For Indian Scripts: A Comprehensive Review

## Poonam Gaur

*Assistant Professor, Chandigarh Group of Colleges, Landran, Punjab- 140307*
*Email Id- bca.poonamgaur@gmail.com*

Recent advances in deep learning techniques have led to a considerable increase in the use of Natural Language Processing (NLP) for Indian languages. The diverse and complex nature of Indian scripts presents unique challenges, including data scarcity, code-mixing, morphological richness, and phonetic variations. This review explores the role of deep learning in overcoming these challenges, focusing on transformer-based models, transfer learning, and self-supervised learning. Recent innovations such as MuRIL, IndicBERT, and multilingual NLP frameworks have significantly improved machine translation, speech recognition, and sentiment analysis for Indian languages. However, several limitations remain, including computational costs, ethical concerns, and the lack of standardized evaluation benchmarks. This paper discusses these challenges, recent advancements, and future research directions to enhance the inclusivity and efficiency of NLP applications for Indian languages. To create scalable and reliable NLP solutions that are suited to India's linguistic diversity, cooperation between researchers, linguists, and legislators would be essential.

**Keywords:** Deep Learning, Transformer Models, Indian Languages, Natural Language Processing (NLP), and Multilingual NLP

## 1. Introduction

### 1.1 Background and Significance of NLP in Indian Scripts

The ability of natural language processing (NLP) to facilitate human-computer interaction by empowering machines to comprehend and interpret human languages has drawn a lot of attention recently. While global NLP research has predominantly focused on widely spoken languages such as English, there has been increasing interest in processing Indian languages due to their rich linguistic diversity and growing digital footprint. Given the complexities of Indian scripts—ranging from phonetic structures to syntactic variations—deep learning-based approaches have proven instrumental in enhancing NLP applications across multiple Indian languages.

A crucial aspect of Indian language NLP is the similarity among various scripts, which can be leveraged to improve performance in multiple NLP tasks. Aggarwal [1] explored this linguistic similarity and demonstrated its utility in various NLP applications, like named entity identification, sentiment analysis, and machine translation. The study highlights how transfer learning and cross-lingual models can enhance language processing capabilities for underrepresented Indian languages. Additionally, early work by Sowmya [2] addressed fundamental challenges in text input methods for Indian scripts, shedding light on the difficulties associated with keyboard layouts, transliteration, and character encoding. These challenges underscore the need for robust NLP models capable of handling complex script structures efficiently.

The development of NLP has also been characterized by improvements in computational methods.Shukla et al. [3] emphasized the growing role of deep learning in unlocking the potential of text and speech data, with an emphasis on low-resource and multilingual language processing. This aligns with Piotrowski [4], who provided insights into NLP techniques for historical texts, reinforcing the importance of context-aware models that can decipher script variations and linguistic evolution. Such historical perspectives offer valuable lessons in developing NLP solutions tailored to the ever-changing character of Indian languages.

New developments in NLP have also extended to applications in education and assessment. Abdullah et al. [5] proposed an automated deep learning model for evaluating handwritten answer scripts, integrating NLP techniques to assess linguistic structure and content quality. This work highlights the expanding role of NLP beyond traditional applications, demonstrating its potential in pedagogical and administrative domains.

Overall, NLP for Indian scripts presents both opportunities and challenges, requiring a blend of linguistic insights and technological innovations. The increasing adoption of deep learning and multilingual models promises to improve the accessibility and efficiency of language processing for India's vast linguistic landscape. Future research must continue addressing data scarcity, script complexity, and model interpretability to ensure robust and inclusive NLP solutions.

## 1.2 Challenges in Processing Indian Languages

Natural Language Processing (NLP) for Indian languages poses particular difficulties because of their great linguistic diversity, different scripts, and intricate structures. Unlike English and other widely studied languages, Indian languages belong to multiple linguistic families, such as Indo-Aryan and Dravidian, making standardization and resource development difficult. Furthermore, a lack of high-quality annotated datasets, morphological richness, and complex grammar rules further complicate NLP model development.

One of the major challenges in Indian language processing is machine translation. Singh et al. [6] provided an extensive review of machine translation systems for Indian languages, going at different modeling approaches and their drawbacks. The study identified major obstacles

such as syntactic divergence, scarcity of parallel corpora, and semantic ambiguity, which hinder the development of high-accuracy translation systems. Patel et al. [7] also highlighted similar issues, emphasizing the need for better linguistic resources and domain-specific models to improve translation quality.

Another significant issue is the availability of resources for low-resource Indian languages. Dongare [8] explored the process of corpus creation for such languages, identifying major bottlenecks, including the lack of digitized texts, orthographic variations, and dialectal differences. These challenges limit the training of robust NLP models, as deep learning techniques require large-scale datasets for effective learning. The scarcity of linguistic resources affects tasks like text categorization, named entity recognition, and speech recognition, which rely on well-structured language corpora.

Additionally, the complexities of regional language processing have been extensively studied. Harish and Rangan [9] conducted a comprehensive survey on NLP techniques for Indian regional languages, pointing out issues such as code-mixing, lack of standardized spellings, and low adoption of Indian scripts in digital spaces. Their findings suggest that advancements in deep learning and linguistic adaptation techniques could bridge some of these gaps, but substantial efforts are still required in resource development and model optimization.

Overall, the challenges in processing Indian languages stem from both technical and linguistic factors. Addressing these issues requires a multi-faceted approach, including improved dataset creation, better machine translation models, and the integration of linguistic rules into deep learning architectures. Collaboration and ongoing research in this area will be essential for enhancing NLP applications and ensuring digital inclusivity for India's diverse linguistic landscape.

### 1.3 Role of Deep Learning in NLP
Natural Language Processing (NLP) has been transformed by deep learning, which allows models to comprehend and produce human language with astounding precision. Traditional rule-based and statistical NLP approaches struggled with linguistic complexities, but deep learning methods, particularly neural networks, have significantly improved tasks include speech recognition, sentiment analysis, and machine translation. Raaijmakers [10] highlighted how deep learning architectures, such as transformers, long short-term memory (LSTM) networks, and recurrent neural networks (RNNs), are advancing natural language processing. Through the acquisition of contextual representations from large volumes of textual data, these models have improved their language modeling capabilities.

Kamath et al. [11] discussed how deep learning has bridged the discrepancy between actual NLP applications and theoretical language knowledge. With the advent of transformer-based models like BERT and GPT, NLP systems can now achieve human-like language comprehension, making them highly effective in chatbots, automated text summarization, and question-answering systems. These developments keep expanding the realm of what is

conceivable in NLP, with deep learning driving innovations in multilingual processing and low-resource language modeling.

## 1.4 Objectives and Scope of the Review

**Objectives**

This review aims to critically analyze the Natural language processing (NLP) and deep learning for Indian scripts by addressing the following objectives:

1. To assess the present condition of NLP techniques applied to Indian languages, including their linguistic complexities and computational challenges.
2. To investigate the efficacy of deep learning architectures in Indian language processing, including Transformers, Long Short-Term Memory (LSTM), Recurrent Neural Networks (RNNs), and Convolutional Neural Networks (CNNs).
3. To assess the availability and limitations of datasets, linguistic corpora, and computational resources for NLP applications in Indian scripts.
4. To identify key challenges in machine translation, speech recognition, and text processing specific to Indian languages.
5. To highlight recent advancements, gaps, and future research directions for improving NLP performance in Indian linguistic contexts.

**Scope of the Review**

The scope of this review encompasses a broad spectrum of NLP applications for Indian languages, with a particular focus on deep learning-driven advancements. It includes:

- **Language Processing Techniques**: A comparative analysis of traditional rule-based methods, statistical approaches, and models based on deep learning for Indian scripts.
- Machine Learning and Deep Learning Models: Examination of state-of-the-art architectures like T5, GPT, and BERT specifically their adaptations for Indian languages.
- **Linguistic Challenges**: Discussion on code-mixing, morphological richness, syntactic complexity, and phonetic variations in Indian scripts.
- **Data Availability and Resource Constraints**: Evaluation of existing linguistic datasets, annotated corpora, and the role of **low**-resource language modeling in Indian NLP**.**
- Practical Applications**:** Exploration of machine translation systems**,** sentiment analysis, Speech-to-text and text-to-speech synthesis technologies in the Indian linguistic domain.
- **Research Gaps and Future Directions**: Identification of key limitations in current NLP research and potential avenues for improvement, particularly in deep learning approaches for underrepresented Indian languages.

## 2. Linguistic Characteristics of Indian Scripts

The diversity and complexity of Indian scripts pose significant computational difficulties for Natural Language Processing (NLP). Unlike Latin-based languages, Indian scripts exhibit unique characteristics such as ligatures, conjunct consonants, vowel diacritics, and script-specific variations, making text processing more intricate. Ghosh et al. [12] highlighted that script recognition in Indian languages requires sophisticated pattern analysis techniques due to the complex structure of characters and their contextual dependencies. The presence of multiple scripts each with distinct phonetic and morphological properties—further complicates the development of generalized NLP models.

Another major challenge is text rendering and type design for Indian scripts. Ghosh and Knuth [13] discussed early computational approaches for Indian script composition, emphasizing the difficulty in creating digital fonts and text formatting systems that accommodate intricate ligatures and orthographic rules. Additionally, Ishida [14] noted that the Unicode representation of Indian scripts involves multiple code points for a single character, making computational processing and string manipulation more demanding.

Pal and Chaudhuri [15] provided a comprehensive study of Indian scripts' character recognition, underscoring the difficulty of OCR (Optical Character Recognition) and handwriting recognition due to script variations and the absence of large annotated datasets. Furthermore, many Indian languages lack sufficient linguistic resources and annotated corpora, which hinders the training of deep learning models for NLP applications. Taking care of these challenges calls for the creation of script-specific pre-processing techniques, robust OCR systems, and models for machine learning tailored to the structural nuances of Indian languages.

## 3. Literature Review

With research concentrating on topics like Named Entity identification (NER), machine translation, sentiment analysis, handwriting identification, and dependency parsing, natural language processing (NLP) for Indian languages has advanced dramatically over the last few decades. However, due to script diversity, phonetic variations, and resource limitations, Indian language NLP still faces many challenges compared to global languages like English.

This section presents a comprehensive literature review of major research contributions to NLP for Indian languages, highlighting key studies, methodologies, and findings.

**Table 1:** Key Literature in NLP for Indian Languages

| Author | Study Focus | Key Contributions | Ref |
|---|---|---|---|
| **Panchal & Shah (2024)** | Historical Evolution of NLP | Traces NLP advancements from rule-based methods to deep learning, focusing on Indic languages. | [16] |

| **Vinay (2024)** | NLP for Legal Documentation | Highlights challenges in processing legal texts in Indian languages, including syntactic ambiguity. | [17] |
|---|---|---|---|
| **Pandey & Nathani (2024)** | NER for Indian Languages | Examines ML, DL, and rule-based methods for NER; discusses data scarcity and linguistic diversity. | [18] |
| **Singh, Bhandari & Singh (2024)** | Punjabi Language Processing | Explores phonetic complexities and resource limitations in Punjabi NLP. | [19] |
| **Singh, Sharma & Chauhan (2023)** | Handwriting Recognition | Reviews CNN and RNN-based approaches for offline recognition of Indian scripts. | [20] |
| **Khurana et al. (2023)** | Trends in NLP Research | Identifies key challenges like low-resource settings and proposes transfer learning solutions. | [21] |
| **Guetari, Ayari & Sakly (2023)** | ML vs. DL in NLP | Compares ML and DL techniques, finding DL superior for semantic analysis and text generation. | [22] |
| **Joseph et al. (2023)** | NLP in Indian Languages | Surveys NLP resources, tools, and annotated corpora, stressing the need for language-specific embeddings. | [23] |
| **Kaladevi et al. (2022)** | Tamil NLP Evolution | Discusses hybrid rule-based and DL models for Tamil text-to-speech and script transformations. | [24] |
| **Sen et al. (2022)** | Bangla NLP Methods | Finds transformer-based models outperform ML in Bangla text classification and sentiment analysis. | [25] |
| **Marreddy et al. (2022)** | Telugu NLP Challenges | Highlights dataset creation, embeddings, and syntactic analysis issues; suggests transfer learning. | [26] |
| **Khanuja et al. (2022)** | Diversity & Inclusion in NLP | Identifies biases in Indian NLP models; suggests strategies for inclusivity. | [27] |
| **Rajendran et al. (2022)** | Tamil NLP Technologies | Reviews speech recognition, machine translation, and chatbot development; explores BERT applications. | [28] |
| **Sen et al. (2021)** | Bangla NLP Review | Compares ML vs. DL for Bangla NLP; highlights transformer and contextual embedding advancements. | [29] |
| **Awais et al. (2021)** | ML vs. DL in Activity Recognition | Finds DL superior for classifying activities in older adults using temporal prediction. | [30] |
| **Mukhamediev et al. (2021)** | ML to DL Evolution | Scientometric review of ML and DL across NLP, computer vision, and bioinformatics. | [31] |

| Harish & Rangan (2020) | Indian Language Processing | Surveys NER, sentiment analysis, and machine translation; emphasizes multilingual embeddings. | [32] |
|---|---|---|---|
| Kundu et al. (2020) | Indic Handwriting Recognition | Uses CNNs and attention-based models for word-level handwritten script recognition. | [33] |
| Rani & Kumar (2019) | Sentiment Analysis | Reviews lexicon-based, ML-based, and hybrid methods; discusses sarcasm detection challenges. | [34] |
| Guerdan et al. (2019) | DL vs. ML in Intelligence Prediction | Finds DL more effective than ML in cognitive intelligence prediction. | [35] |
| Gupta (2019) | Semantic Role Labeling | Shows DL outperforms rule-based methods for SRL in Indian languages. | [36] |
| Tandon (2018) | Dependency Parsing | Examines parsing techniques for Indian languages, covering transition- and graph-based approaches. | [37] |
| Menger et al. (2018) | ML vs. DL in Clinical Text | Finds DL significantly improves predictive accuracy for inpatient violence prediction. | [38] |

## 3.5 Research Gaps

1. **Limited Focus on Low-Resource Indian Languages:** Many studies focus on major Indian languages like Hindi, Tamil, and Bengali (e.g., [16], [19], [28]), but there is limited research on low-resource languages such as Konkani, Manipuri, or Tulu. Future work should develop NLP models and datasets for underrepresented Indian languages to improve inclusivity.
2. **Lack of Unified Frameworks for Multilingual NLP:** Studies have explored NLP for individual languages ([22], [26], [33]), but comparative and unified approaches remain underdeveloped. There is a need for multilingual frameworks that can handle multiple Indian scripts and languages simultaneously.
3. **Limited Use of Deep Learning for Script Recognition:** Several studies use machine learning for Indic script recognition ([12], [15], [20]), but deep learning-based approaches are not widely explored. More research is needed to integrate CNNs, RNNs, and transformer-based models for improving script recognition accuracy.
4. **Computational Challenges in Indian Language Processing:** Handling morphologically rich Indian languages ([17], [19], [29]) remains a major issue due to lack of robust NLP toolkits. Future research should work on efficient morphological analysers and retrained models specifically designed for Indian languages.
5. **Ethical Considerations and Bias in Indian NLP Models:** Studies like [27] highlight inclusion issues in NLP, but there is insufficient research on bias, fairness, and ethical AI

for Indian languages. Future work should investigate and mitigate biases in language models for more equitable NLP applications.

## 4. Deep Learning in NLP: An Overview

### 4.1 Evolution from Traditional NLP to Deep Learning
Deep learning-driven techniques have significantly replaced rule-based and statistical approaches in natural language processing (NLP). Early natural language processing (NLP) models were based on manually created rules and conventional machine learning methods like Conditional Random Fields (CRFs) and Hidden Markov Models (HMMs), which had trouble with long-range relationships and semantic understanding. Neural network-based models were introduced with the emergence of deep learning capable of capturing contextual meaning and hierarchical representations, significantly improving performance in language modeling, translation, and text generation ([39]). The development of architectures that successfully model sequential data, the availability of large-scale datasets, and improvements in processing power all contributed to this shift.

### 4.2 Key Deep Learning Architectures in NLP

### 4.2.1 Recurrent Neural Networks (RNNs)
RNNs were one of the earliest deep learning architectures used in NLP, created to manage sequential data by preserving historical data in a concealed state. They proved effective for tasks such as speech recognition and machine translation. However, the vanishing gradient issue plagues conventional RNNs, making it challenging to identify long-term dependencies in text sequences. [40]. More sophisticated recurrent architectures were created as a result of this restriction.

### 4.2.2 Long Short-Term Memory (LSTMs) and Gated Recurrent Units (GRUs)
LSTMs and GRUs were introduced to overcome the limitations of conventional RNNs. These models can better represent long-range interdependence by selectively remembering or forgetting information through the use of gating mechanisms. Because of their capacity to preserve context across lengthy sequences, LSTMs in particular have found extensive application in speech processing, language modeling, and text categorization. [41]. GRUs, being a more computationally efficient variant, have also been adopted in various NLP applications.

### 4.2.3 Transformer Models (BERT, GPT, etc.)
Transformers have revolutionized NLP by substituting recurrent connections includes mechanisms for self-attention, enabling parallel processing of entire sequences. Models like BERT (Bidirectional Encoder Representations from Transformers) and GPT (Generative Pre-trained Transformer) were created as a result of this breakthrough, and they have attained cutting-edge results in text production and comprehension.[42]. Unlike RNN-based models, Transformers do not suffer from vanishing gradients and can model long-range dependencies more effectively, making them the backbone of modern NLP applications.

### 4.2.4 Convolutional Neural Networks (CNNs) for NLP

CNNs have been effectively used for NLP applications including text classification and sentiment analysis, despite their primary use in image processing. By using convolutional filters, CNNs can capture local dependencies and hierarchical features in textual data. While not as dominant as Transformers, CNN-based models have demonstrated strong performance in areas requiring efficient feature extraction, particularly in low-resource settings[43].

### 4.3 Advantages of Deep Learning Over Classical NLP Techniques

Deep learning provides several advantages over traditional NLP approaches, primarily in its ability to learn complex representations directly from raw text. Unlike rule-based and statistical models that require extensive manual feature engineering, deep learning architectures automatically extract meaningful patterns from large datasets, significantly improving scalability and adaptability. Moreover, deep learning models, particularly Transformers, have demonstrated superior performance in capturing contextual dependencies, leading to more accurate and coherent language models [40]. Additionally, the ability to pre-train deep learning models on vast corpora has enabled transfer learning, which enables models to be optimized for certain tasks using a small amount of labeled data. Deep learning keeps pushing NLP innovation in spite of obstacles like high processing costs and the requirement for huge datasets, opening the door for more effective, interpretable, and flexible language models.

### 5. NLP Applications for Indian Scripts

Natural Language Processing (NLP) has seen notable progress in Indian language processing, driven by the need for improved accessibility, machine translation, sentiment analysis, and speech recognition. Given India's linguistic diversity, NLP models must address the unique challenges posed by morphologically rich languages, multiple scripts, and dialectal variations. While deep learning has enhanced the performance of various NLP tasks, there are still gaps in the data's availability, standardization, and model efficiency. This section explores key applications of NLP in Indian scripts, highlighting recent progress and challenges.

**Table 2:** Comparative Analysis of NLP Techniques

| Technique | Approach | Advantages | Challenges | Applications |
|-----------|----------|------------|------------|--------------|
| Rule-Based | Predefined linguistic rules | Explainability, good for small datasets | Limited scalability, needs manual updates | POS tagging, Named Entity Recognition |

| Machine Learning | Statistical models (SVM, CRF, HMM) | Automates learning, adaptable | Requires labeled data, feature engineering | Analysis of sentiment and voice recognition |
|---|---|---|---|---|
| Deep Learning | Neural Networks (RNN, CNN, Transformer) | Handles complex tasks, self-learning | Data-intensive, computationally expensive | Machine translation, chatbots |

## 5.1 Text Processing and Named Entity Recognition (NER)

Efficient text preprocessing is fundamental for Indian language NLP applications due to the complex morphology, free word order, and rich inflectional structures of these languages. Unlike English, many Indian languages lack standardized spelling conventions, making tokenization, stemming, and lemmatization challenging. Additionally, transliteration, where native scripts are written in Roman characters, further complicates text normalization. Named Entity Recognition (NER) is another critical component, used for tasks such as information retrieval, chatbots, and machine translation. However, existing NER models struggle with underrepresented languages, ambiguous entity boundaries, and limited annotated datasets. Future research must focus on improving preprocessing pipelines and developing high-quality, domain-specific annotated corpora to enhance model performance.

## 5.2 Machine Translation and Sentiment Analysis

The introduction of transformer-based neural machine translation (NMT) models has greatly enhanced machine translation between English and Indian languages. However, challenges persist due to structural differences between English and Indian languages, including variations in word order, morphological richness, and limited parallel corpora. The scarcity of high-quality bilingual datasets restricts the efficiency of translation systems for low-resource languages. Sentiment analysis, another crucial NLP application, faces issues such as code-switching, sarcasm detection, and dialectal variations. While deep learning-based sentiment analysis models perform well for high-resource languages, low-resource Indian languages require more annotated sentiment datasets and advanced context-aware models to improve accuracy.

## 5.3 Speech Processing and Optical Character Recognition (OCR)

Text-to-Speech (TTS) and Automatic Speech Recognition (ASR) technologies are essential for enabling voice-based digital interactions in Indian languages. Despite advances in ASR, challenges persist due to multiple dialects, tonal variations, and code-mixed speech. Many Indian languages lack large transcribed speech datasets, making it difficult to train accurate ASR models. Similarly, Optical Character Recognition (OCR) plays a crucial role in digitizing printed and handwritten Indian scripts, yet existing OCR models struggle with complex ligatures, conjunct characters, and non-standard fonts. The integration of deep learning, particularly convolutional and transformer-based architectures, has improved OCR accuracy,

but further enhancements are required to process handwritten text and degraded documents effectively.

## 6. Deep Learning Approaches for Indian Language NLP

Natural language processing (NLP) has greatly improved because to deep learning for Indian languages by improving text generation, translation, speech recognition, and sentiment analysisDespite these developments, managing code-mixed text, low-resource languages, and multilingual processing still present difficulties. This section discusses pretrained language models, transfer learning, data augmentation, and benchmark datasets for Indian language NLP.

### 6.1 Pretrained Language Models for Indian Scripts

Pretrained transformer models have transformed NLP for Indian languages by enabling better understanding of syntactic and semantic structures. Some key models include:

**Table 3:** Comparison of Pretrained Language Models for Indian NLP

| Model | Languages Supported | Architecture | Key Strengths | Limitations |
|---|---|---|---|---|
| MuRIL (Multilingual Representations for Indian Languages) | Hindi, Tamil, Telugu, Marathi, Bengali | BERT-based | Strong performance on multilingual NLP tasks | High computational cost |
| IndicBERT | 12 Indic languages | ALBERT-based | Efficient, lightweight | Moderate performance on complex tasks |
| Samanantar | 11 Indian languages | Transformer-based | High-quality translations | Requires large training data |
| AI4Bharat IndicTrans | Multiple Indian languages | Seq2Seq Transformer | State-of-the-art machine translation | Limited outside of translation |

These models significantly improve NLP activities such as speech recognition, sentiment analysis, and machine translation, but challenges remain in handling morphological complexity and dialectal variations.

### 6.2 Multilingual and Code-Mixed Language Processing

Indian languages frequently appear in code-mixed text (e.g., Hinglish, Tanglish), posing challenges for NLP models. Deep learning techniques help address language switching, transliteration, and spelling variations.

**Table 4:** Approaches for Handling Code-Mixed and Multilingual Text

| Approach | Methodology | Strengths | Challenges |
|---|---|---|---|
| Transliteration-Based Models | Convert words to a standard script before processing | Reduces ambiguity | Errors in script conversion |
| Multilingual Pretrained Models | Use models like XLM-R, MuRIL for mixed-language understanding | Strong generalization across languages | Requires large datasets |
| Hybrid Approaches | Combine rule-based and deep learning models | Improved performance on noisy data | Computationally expensive |

Many code-mixed language models rely on self-attention mechanisms and transfer learning to improve performance.

### 6.3 Low-Resource Language Processing and Transfer Learning

Many Indian languages suffer from limited annotated datasets, making NLP development challenging. Transfer learning addresses this issue by adapting pretrained models trained on high-resource languages. Fine-tuning models like MuRIL for specific languages enhances performance, while zero-shot learning allows models trained on related languages to generalize without direct training. Additionally, self-supervised learning utilizes large unlabeled corpora to develop language understanding, reducing dependence on manually labeled data. These techniques significantly improve NLP capabilities for low-resource Indian languages.

### 6.4 Data Augmentation and Annotation Strategies

Data scarcity remains a major challenge in Indian NLP, and data augmentation techniques help overcome this limitation by artificially increasing dataset diversity. Back translation generates variations for machine translation, while word replacement using resources like IndoWordNet improves NER and sentiment analysis. Text paraphrasing aids question answering and chatbot training, whereas text simplification enhances accessibility in education and assistive technologies. These augmentation methods help improve model generalization and robustness, making them essential for NLP applications in Indian languages.

### 6.5 Benchmark Datasets and Evaluation Metrics

To evaluate and compare NLP models, standardized datasets and assessment measures are necessary. The AI4Bharat IndicNLP Corpus supports NER, while Samanantar facilitates machine translation across multiple Indian languages. The ULCA Benchmark Dataset provides resources for speech recognition, translation, and NER. Evaluation metrics like BLEU score assess translation quality, F1-score measures NER and classification accuracy, and WER evaluates speech recognition efficiency. These resources enable consistent benchmarking, fostering advancements in Indian language NLP.

### 7. Challenges and Limitations

Despite significant advancements, NLP for Indian scripts faces several obstacles and constraints. The lack of extensive annotated corpora for many Indian languages is one of the main problems, which makes it challenging to train high-performance NLP models. This is especially problematic for low-resource languages such as Konkani, Manipuri, and Tulu, which remain underrepresented in existing datasets. Additionally, the phenomenon of code-mixing and dialectal variations—where speakers frequently switch between languages, such as in Hinglish (Hindi-English) or Tanglish (Tamil-English)—complicates text normalization, tokenization, and machine translation.

Another significant limitation is computational cost, as deep learning models, particularly transformer-based architectures, require extensive computational resources for training and inference. This makes their adoption challenging in resource-constrained environments, including academic research settings and small-scale industries. Furthermore, the script and phonetic complexity of Indian languages adds to the difficulty of tasks like Optical Character Recognition (OCR) and speech recognition. The presence of ligatures, conjunct consonants, and multiple character representations increases processing complexity compared to Latin-based scripts.

Additionally, the lack of standardized NLP frameworks for Indian languages results in fragmented research efforts, where models developed for one language may not generalize well to others. This lack of standardization is particularly problematic for machine translation and cross-lingual NLP applications. Another pressing concern is bias and ethical issues in NLP models, as these models often inherit biases from training data, leading to potential inaccuracies in sentiment analysis, automated decision-making, and digital assistants. Lastly, evaluation challenges persist due to the absence of standardized datasets and benchmarking metrics tailored for Indian languages, making it difficult to compare and assess model performance effectively.

## 8. Recent Advances and Future Directions
Recent advancements in deep learning have significantly improved NLP applications for Indian languages. Transformer-based models, including MuRIL (Multilingual Representations for Indian Languages) and IndicBERT have demonstrated remarkable success in multilingual text understanding, helping named entity identification, sentiment analysis, and machine translation. Another crucial development is the adoption of transfer learning and zero-shot learning, which allow pretrained models trained on the adaptation of high-resource languages for low-resource Indian languages, improving accuracy and efficiency.

Self-supervised learning techniques have emerged as a potential remedy for the issue of data scarcity by allowing models to pick up language patterns from large amounts of unannotated text. Moreover, data augmentation strategies, such as back translation, paraphrasing, and word replacement, have improved training data diversity, making NLP models more robust. Another key advancement is the rise of multilingual NLP frameworks, which integrate multiple Indian

languages within a single model, facilitating cross-lingual applications such as document classification and question-answering systems.

In addition to textual NLP applications, deep learning has led to significant improvements in speech processing and OCR for Indian scripts. Enhanced automatic speech recognition (ASR) systems now account for dialectal variations, while deep learning-powered OCR models have improved the accuracy of recognizing printed and handwritten Indian text. Furthermore, increasing attention is being given to bias mitigation strategies, as researchers work on ensuring fairness and reducing discrimination in NLP models. Another promising area of growth is conversational AI, where virtual assistants and chatbots are becoming more sophisticated in understanding and responding to queries in Indian languages, benefiting sectors such as customer service, education, and e-governance.

## 9. Conclusion
In conclusion, deep learning has significantly transformed Natural Language Processing for Indian languages, offering improved capabilities in machine translation, speech recognition, and sentiment analysis. However, several challenges remain, particularly in the areas of data availability, script complexity, and model efficiency. Addressing these challenges requires continued research and innovation in fields like data augmentation, self-supervised learning, and transfer learning, that can help bridge the resource gap for low-resource Indian languages.

The development of transformer-based models like MuRIL and IndicBERT has shown promising results, but further advancements are needed to enhance cross-lingual transferability and efficiency. Additionally, ensuring fairness and reducing biases in natural language processing models will be essential for the moral deployment of language technologies. As the field continues to evolve, collaboration between linguists, data scientists, policymakers, and technology developers will be essential to create more robust, inclusive, and effective NLP solutions for India's diverse linguistic landscape.

## References

1. Aggarwal, S. (2022). Exploiting Indian Languages' Similarity for Different NLP Applications (Doctoral dissertation, International Institute of Information Technology Hyderabad).
2. Sowmya, V. B. (2008). Text input methods for Indian languages (Doctoral dissertation, Master's Thesis, International Institute of Information Technology, Hyderabad).
3. Shukla, K., Vashishtha, E., Sandhu, M., & Choubey, P. R. (2023). Natural Language Processing: Unlocking the Power of Text and Speech Data. Xoffencerpublication.
4. Piotrowski, M. (2012). Natural language processing for historical texts. Morgan & Claypool Publishers.
5. Abdullah, A. S., Geetha, S., Aziz, A. A., & Mishra, U. (2024). Design of automated model for inspecting and evaluating handwritten answer scripts: A pedagogical approach with NLP and deep learning. Alexandria Engineering Journal, 108, 764-788.

6.  Singh, M., Kumar, R., & Chana, I. (2021). Machine translation systems for Indian languages: review of modelling techniques, challenges, open issues and future research directions. Archives of Computational Methods in Engineering, 28(4), 2165-2193.
7.  Dongare, P. (2024, May). Creating Corpus of Low Resource Indian Languages for Natural Language Processing: Challenges and Opportunities. In Proceedings of the 7th Workshop on Indian Language Data: Resources and Evaluation (pp. 54-58).
8.  Patel, R. N., Pimpale, P. B., & Sasikumar, M. (2019). Machine translation in Indian languages: challenges and resolution. Journal of Intelligent Systems, 28(3), 437-445.
9.  Harish, B. S., & Rangan, R. K. (2020). A comprehensive survey on Indian regional language processing. SN Applied Sciences, 2(7), 1204.
10. Raaijmakers, S. (2022). Deep learning for natural language processing. Simon and Schuster.
11. Kamath, U., Liu, J., & Whitaker, J. (2019). Deep learning for NLP and speech recognition (Vol. 84, p. 1). Cham, Switzerland: Springer.
12. Ghosh, D., Dube, T., & Shivaprasad, A. (2010). Script recognition—a review. IEEE Transactions on pattern analysis and machine intelligence, 32(12), 2142-2161.
13. Ghosh, P. K., & Knuth, P. K. (1983). An approach to type design and text composition in Indian scripts. Stanford: Department of Computer Science, Stanford University.
14. Ishida, R. (2002, September). An introduction to Indic scripts. In Proceedings of the 22nd Int. Unicode Conference (Vol. 8).
15. Pal, U., & Chaudhuri, B. B. (2004). Indian script character recognition: a survey. pattern Recognition, 37(9), 1887-1899.
16. Panchal, B. Y., & Shah, A. (2024, December). NLP Research: A Historical Survey and Current Trends in Global, Indic, and Gujarati Languages. In 2024 4th International Conference on Ubiquitous Computing and Intelligent Information Systems (ICUIS) (pp. 1263-1272). IEEE.
17. Vinay, S. B. (2024). Natural Language Processing for Legal Documentation in Indian Languages. International Journal of Natural Language Processing (IJNLP), 2(1), 1-11.
18. Pandey, P., & Nathani, B. (2024). State-of-art approach for Indian Language based on NER: Comprehensive Review.
19. Singh, G., Bhandari, R., & Singh, P. (2024, January). Advancing NLP for Punjabi Language: A Comprehensive Review of Language Processing Challenges and Opportunities. In 2024 2nd International Conference on Intelligent Data Communication Technologies and Internet of Things (IDCIoT) (pp. 1250-1257). IEEE.
20. Singh, S., Sharma, A., & Chauhan, V. K. (2023). Indic script family and its offline handwriting recognition for characters/digits and words: a comprehensive survey. Artificial Intelligence Review, 56(Suppl 3), 3003-3055.
21. Khurana, D., Koli, A., Khatter, K., & Singh, S. (2023). Natural language processing: state of the art, current trends and challenges. Multimedia tools and applications, 82(3), 3713-3744.
22. Guetari, R., Ayari, H., & Sakly, H. (2023). Computer-aided diagnosis systems: a comparative study of classical machine learning versus deep learning-based approaches. Knowledge and Information Systems, 65(10), 3881-3921.
23. Joseph, J., Lalithsriram, S. R., & Menon, N. (2023). Applications and Developments of NLP Resources for Text Processing in Indian Languages. Multilingual Digital Humanities.
24. Kaladevi, R., Revathi, A., & Manju, A. (2022). Analyzing the evolution of modern Tamil script for natural language processing. ECS Transactions, 107(1), 5219.
25. Sen, O., Fuad, M., Islam, M. N., Rabbi, J., Masud, M., Hasan, M. K., ... & Iftee, M. A. R. (2022). Bangla natural language processing: A comprehensive analysis of classical, machine learning, and deep learning-based methods. IEEE Access, 10, 38999-39044.

26. Marreddy, M., Oota, S. R., Vakada, L. S., Chinni, V. C., & Mamidi, R. (2022). Am I a resource-poor language? Data sets, embeddings, models and analysis for four different NLP tasks in telugu language. ACM Transactions on Asian and Low-Resource Language Information Processing, 22(1), 1-34.

27. Khanuja, S., Ruder, S., & Talukdar, P. (2022). Evaluating the diversity, equity and inclusion of nlp technology: A case study for indian languages. arXiv preprint arXiv:2205.12676.

28. Rajendran, S., Anand Kumar, M., Rajalakshmi, R., Dhanalakshmi, V., Balasubramanian, P., & Soman, K. P. (2022, November). Tamil NLP Technologies: Challenges, State of the Art, Trends and Future Scope. In International Conference on Speech and Language Technologies for Low-resource Languages (pp. 73-98). Cham: Springer International Publishing.

29. Sen, O., Fuad, M., Islam, M. N., Rabbi, J., Hasan, M. K., Fime, A. A., ... & Iftee, M. A. R. (2021). Bangla natural language processing: A comprehensive review of classical machine learning and deep learning based methods. CoRR.

30. Awais, M., Chiari, L., Ihlen, E. A., Helbostad, J. L., & Palmerini, L. (2021). Classical machine learning versus deep learning for the older adults free-living activity classification. Sensors, 21(14), 4669.

31. Mukhamediev, R. I., Symagulov, A., Kuchin, Y., Yakunin, K., & Yelis, M. (2021). From classical machine learning to deep neural networks: A simplified scientometric review. Applied Sciences, 11(12), 5541.

32. Harish, B. S., & Rangan, R. K. (2020). A comprehensive survey on Indian regional language processing. SN Applied Sciences, 2(7), 1204.

33. Kundu, S., Paul, S., Singh, P. K., Sarkar, R., & Nasipuri, M. (2020). Understanding NFC-Net: a deep learning approach to word-level handwritten Indic script recognition. Neural Computing and Applications, 32, 7879-7895.

34. Rani, S., & Kumar, P. (2019). A journey of Indian languages over sentiment analysis: a systematic review. Artificial Intelligence Review, 52, 1415-1462.

35. Guerdan, L., Sun, P., Rowland, C., Harrison, L., Tang, Z., Wergeles, N., & Shang, Y. (2019). Deep learning vs. classical machine learning: A comparison of methods for fluid intelligence prediction. In Adolescent Brain Cognitive Development Neurocognitive Prediction: First Challenge, ABCD-NP 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 13, 2019, Proceedings 1 (pp. 17-25). Springer International Publishing.

36. Gupta, A. (2019). Semantic Role Labeling for Indian languages (Doctoral dissertation, Ph. D. thesis, International Institute of Information Technology Hyderabad).

37. Tandon, J. (2018). Advancements in dependency parsing for indian languages. International Institute of Information Technology.

38. Menger, V., Scheepers, F., & Spruit, M. (2018). Comparing deep learning and classical machine learning approaches for predicting inpatient violence incidents from clinical text. Applied Sciences, 8(6), 981.

39. Khan, W., Daud, A., Khan, K., Muhammad, S., & Haq, R. (2023). Exploring the frontiers of deep learning and natural language processing: A comprehensive overview of key challenges and emerging trends. Natural Language Processing Journal, 4, 100026.

40. Torfi, A., Shirvani, R. A., Keneshloo, Y., Tavaf, N., & Fox, E. A. (2020). Natural language processing advancements by deep learning: A survey. arXiv preprint arXiv:2003.01200.

41. Kamath, U., Liu, J., & Whitaker, J. (2019). Deep learning for NLP and speech recognition (Vol. 84, p. 1). Cham, Switzerland: Springer.

42. Goyal, P., Pandey, S., & Jain, K. (2018). Deep learning for natural language processing. New York: Apress.

43.  Fahad, S. A., & Yahya, A. E. (2018, July). Inflectional review of deep learning on natural language processing. In 2018 international conference on smart computing and electronic enterprise (ICSCEE) (pp. 1-4). IEEE.