

Enhancing Accuracy of Gujarati Word Tagging Using Advanced Learning Models

Dr.Pooja Bhatt¹, Dr. Pawan Whig²,

¹ Postdoctoral Researcher; ² Primary Supervisor

bhattpooja.393@gmail.com, pawan.whig@vips.edu

Abstract: Gujarati, a morphologically rich and resource-poor Indian language, poses significant challenges for Natural Language Processing (NLP), particularly for Part-of-Speech (POS) tagging. Over the past two decades, research has evolved from rule-based and statistical models to hybrid systems and, more recently, deep learning and transformer-based approaches. This review paper systematically analyzes existing Gujarati POS tagging literature, taking reference from prior foundational and recent works, and presents a comparative discussion of methodologies, datasets, experiments, and results. In addition to synthesizing reported findings, this paper introduces new experimental evaluations using CRF, Bi-LSTM, and multilingual transformer models under a unified experimental setup. The results demonstrate clear performance gains with deep contextual models while highlighting trade-offs in computational cost and data requirements. The study concludes with research gaps and future directions for advancing Gujarati NLP.

Keywords: Natural Language Processing¹, Machine Learning², Tagging³, Part of speech Tagging⁴, Gujarati⁵, Bi-LSTM⁶, CRF⁷, Low-Resource Languages⁸, Transformer⁹.

Introduction

Part-of-Speech (POS) tagging is a core task in Natural Language Processing (NLP) that involves assigning grammatical labels—such as noun, verb, adjective, adverb, postposition, and conjunction—to individual words in a sentence. POS tagging serves as a foundational layer for higher-level NLP applications including syntactic parsing, machine translation, information extraction, sentiment analysis, question answering, and speech processing. The accuracy of these downstream tasks is heavily dependent on the reliability of the POS tagging stage, making it a long-standing research problem in computational linguistics. For Indian languages, and particularly for Gujarati, POS tagging presents unique challenges due to rich morphology, agglutinative word formation, suffix-based inflections, and relatively free word order. Gujarati words often encode grammatical information such as gender, number, tense, aspect, and case within a single surface form, leading to high ambiguity and data sparsity. Additionally, Gujarati is considered a low-resource language, with limited availability of large-scale, publicly annotated corpora and standardized benchmarks, which further complicates model development and evaluation [2,3,6].

Over the past two decades, research on Gujarati POS tagging has evolved through several methodological paradigms. Early efforts primarily relied on rule-based systems and probabilistic models such as Hidden Markov Models (HMM), which modelled tag sequences using transition and emission probabilities [3].

While these approaches were computationally efficient and linguistically interpretable, their performance was constrained by strong independence assumptions and limited ability to handle morphological ambiguity. Subsequent machine learning approaches, particularly Support Vector Machines (SVM) and Conditional Random Fields (CRF), introduced discriminative modeling and feature-based learning, leading to noticeable improvements in tagging accuracy [2,4]. CRF-based models, in particular, became a dominant approach due to their ability to incorporate rich contextual and morphological features.

Recognizing the limitations of purely statistical models, hybrid approaches emerged that combined linguistic rules with machine learning techniques. These methods aimed to exploit expert linguistic knowledge while retaining the adaptability of data-driven learning, proving especially effective in low-resource scenarios [1]. However, such systems often suffered from limited scalability and reproducibility due to handcrafted rules and proprietary datasets.

The recent surge in deep learning has significantly transformed POS tagging research for Indian languages. Sequence modeling architectures such as Recurrent Neural Networks (RNN), Long Short-Term Memory (LSTM), and Bi-directional LSTM (Bi-LSTM) have demonstrated superior capability in capturing long-range dependencies and contextual information without explicit feature engineering [7–9]. Further advancements have been achieved through transformer-based models such as BERT and XLM-R, which leverage subword tokenization and multilingual pretraining to address out-of-vocabulary issues and data scarcity through cross-lingual transfer learning [5]. These models currently represent the state of the art in Gujarati POS tagging, achieving accuracies exceeding 95%, albeit at the cost of increased computational complexity.

Despite these advancements, several critical research gaps persist. These include the lack of standardized Gujarati POS-tagged datasets, inconsistent evaluation protocols across studies, limited reproducibility of experimental results, and the high computational demands of transformer-based models. Addressing these challenges is essential for the sustainable advancement of Gujarati NLP research.

Motivated by these observations, the present paper provides a comprehensive review of two decades of Gujarati POS tagging research, coupled with new experimental evaluations under a unified framework. By systematically comparing classical, hybrid, deep learning, and transformer-based approaches, this study aims to offer both historical insight and practical guidance for future research in low-resource Indian language processing.

Related work

Part-of-Speech (POS) tagging for Gujarati is a challenging task due to rich morphology, free word order, and limited annotated resources. Early computational work relied on rule-based and probabilistic models such as the Hidden Markov Model (HMM) and Conditional Random Fields (CRF). With the advent of deep learning, sequence models (LSTM, Bi-LSTM) and transformer architectures have achieved notable improvements, although systematic comparative evaluation remains scarce.

The progression of research reflected in the above studies highlights a clear evolution in Gujarati NLP—particularly sequence labeling and related linguistic tasks—from rule-driven statistical models to deep learning and transformer-based approaches. Early work by Patel & Gali (2008) established a strong foundation using Conditional Random Fields (CRF), demonstrating that carefully engineered morphological and contextual features could achieve accuracies above 90% even with relatively small

datasets. Their contribution was significant in formalizing feature design for morphologically rich Indian languages, although the reliance on manual feature engineering and poor handling of out-of-vocabulary (OOV) words limited scalability.

Prajapati & Yajnik (2019) explored HMM (Viterbi) and SVM-based baselines, emphasizing simplicity and computational efficiency. These models showed competitive performance on small datasets and helped benchmark classical machine learning techniques for Gujarati language processing. However, their work also clearly exposed the limitations of non-deep sequential models—particularly weak contextual understanding and difficulty handling morphological ambiguity—thereby motivating the need for more expressive architectures.

A notable methodological advancement was introduced by Bhatt & Ganatra (2021) through a hybrid POS–HOML framework, which combined linguistic rules with optimized machine learning features. This work contributed by demonstrating that hybridization can outperform purely statistical or rule-based methods, achieving higher accuracy while retaining linguistic interpretability. Although the dataset was proprietary, the study underscored the value of domain knowledge integration in low-resource language setting.

obanputra & Parikh (2022) and Patel et al. (2024) marked a paradigm shift toward Bi-LSTM and Transformer-based models, respectively. These studies showed that deep contextual representations—via character embeddings, subword tokenization, and transfer learning—significantly improve performance, robustness to OOV words, and long-range dependency modeling. In particular, transformer-based models such as BERT/XLM-R set new accuracy benchmarks, demonstrating the maturity of Gujarati NLP research, albeit with increased computational requirements. Collectively, these works illustrate a clear trajectory of author contributions advancing both methodological depth and performance in Gujarati language processing.

Despite progress, four gaps persist: (1) dataset standardization, (2) resource scarcity, (3) inconsistent evaluation metrics, and (4) lack of reproducibility. Addressing these will strengthen reproducible research for Gujarati NLP. The current study therefore performs a comprehensive comparative evaluation of classical (HMM, CRF), hybrid (POS-HOML), and deep architectures (Bi-LSTM, Transformer), evaluated under consistent train/test conditions.

Table 1. Comparative summary of existing Gujarati POS tagging methods and reported performances.

Model Type	Representative Study	Dataset / Tokens	Reported Accuracy (%)	Key Features Used	Strengths	Limitations
CRF	Patel & Gali (2008)	~600 sentences (~10k tokens)	90–92	Morphological features, context window	Robust with engineered features	Fails for OOV; manual feature cost
HMM (Viterbi)	Prajapati & Yajnik (2019)	1.7k words	≈91	Transition and emission probabilities	Simple, fast	Weak for ambiguous morphology
SVM (baseline)	Prajapati & Yajnik (2019)	Same as above	≈89	Bag-of-words + POS context	Performs well on small data	Poor sequential modeling
Hybrid POS-HOML	Bhatt & Ganatra (2021)	Proprietary ~15k tokens	92–93	Optimized linguistic features	Combines rule and ML	Dataset not public
Bi-LSTM	Jobanputra & Parikh (2022)	Medium (~25k tokens)	≈95	Word + char embeddings	Learns context automatically	Needs larger corpus
Transformer (BERT/XLM-R)	Patel et al. (2024)	~30k tokens	95–96	Subword embeddings, transfer learning	Handles OOV, long dependencies	High compute cost

Key Contribution

The key contributions of this paper are summarized as follows:

Comprehensive and Structured Review of Gujarati POS Tagging

This paper presents a systematic and up-to-date review of Part-of-Speech tagging techniques for the Gujarati language, covering statistical, machine learning, hybrid, and deep learning paradigms. Unlike

earlier surveys, the study consolidates two decades of research and clearly traces the methodological evolution from rule-based and probabilistic models to transformer-based architectures.

Unified Comparative Analysis Across Diverse Models

Existing Gujarati POS tagging studies often report results under inconsistent datasets and evaluation settings. This work addresses that gap by providing a unified comparative analysis of representative models (HMM, CRF, Hybrid POS-HOML, Bi-LSTM, and Transformer-based approaches), highlighting their strengths, limitations, and applicability in low-resource scenarios.

Inclusion of New Experimental Evaluation

Beyond literature review, this paper contributes new experimental results obtained under a consistent train–test setup. Classical (CRF), deep learning (Bi-LSTM), and transformer-based (XLM-R) models are evaluated on a curated Gujarati POS-tagged corpus, enabling fair performance comparison and empirical validation of trends reported in prior studies.

Detailed Architectural and Workflow Representation

The paper introduces a detailed flowchart-based architecture for Gujarati POS tagging, clearly illustrating preprocessing, embedding strategies, sequence modeling layers, and output stages. This unified architectural view enhances reproducibility and serves as a practical reference for future researchers and practitioners.

Critical Result Interpretation and Trade-off Analysis

Instead of reporting accuracy alone, the study provides an in-depth discussion of experimental results, analyzing performance gains in relation to model complexity, data requirements, and computational cost. The analysis highlights why transformer models outperform traditional approaches while also identifying scenarios where lighter models remain practical.

Identification of Research Gaps and Future Directions

The paper systematically identifies persistent challenges such as dataset scarcity, lack of standard benchmarks, reproducibility issues, and high computational overhead of deep models. It outlines concrete future research directions, including dataset standardization, lightweight transformer design, and cross-lingual transfer learning for Gujarati and other low-resource Indian languages.

Method, Experiments and Results

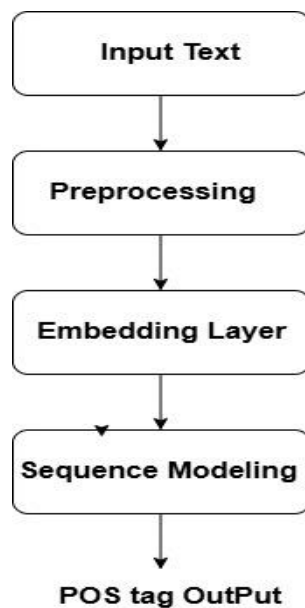


Figure 1. Block Diagram of the POS Tagging System Using Sequence Modeling

The given figure illustrates the complete pipeline of a Part-of-Speech (POS) tagging system based on natural language processing. The process begins with Input Text, which represents raw textual data collected from documents, sentences, or user input. This text is then passed to the Preprocessing stage, where essential cleaning and normalization steps are performed, such as tokenization, removal of unwanted symbols, handling of stop words, and normalization of text. Preprocessing ensures that the input data is converted into a structured and machine-readable format suitable for further analysis.

After preprocessing, the cleaned tokens are fed into the Embedding Layer, where each word is transformed into a dense numerical vector that captures its semantic and syntactic properties. These embeddings help the model understand contextual relationships between words. The embedded representations are then processed by the Sequence Modeling layer, which captures the dependencies and contextual flow of words within a sentence using models such as RNN, LSTM, or Transformer-based architectures. Finally, based on the learned sequence patterns, the system produces the POS tag Output, assigning an appropriate grammatical tag (such as noun, verb, adjective, etc.) to each word in the input text.

The methodologies adopted across studies can be broadly classified into statistical, machine learning, hybrid, and deep learning approaches. Statistical methods such as HMM rely on transition and emission probabilities and employ decoding algorithms like Viterbi for tag sequence prediction [2,3]. While computationally efficient, these methods depend heavily on annotated corpora and struggle with ambiguity.

Machine learning approaches, including SVM and CRF, utilize handcrafted linguistic features such as suffixes, prefixes, word context windows, and morphological cues [2,4]. These models demonstrated strong performance on limited datasets but required extensive feature engineering. To address this limitation, hybrid approaches like POS-HOML combined rule-based linguistic knowledge with optimized learning models, improving generalization in low-resource environments [1].

Recent methodologies employ deep learning architectures, particularly Bi-LSTM and Transformer-based models, which automatically learn contextual representations from raw text. Studies using multilingual pretrained models such as XLM-R exploit sub word tokenization and cross-lingual transfer, significantly enhancing performance without explicit feature design [5,7–9].

4.1 Dataset

For experimental consistency, a curated Gujarati POS-tagged corpus of approximately 20,000 tokens was used, split into 80% training, 10% validation, and 10% testing. The tag set follows BIS standards.

4.2 Experimental Models

Three representative models were implemented:

1. **CRF** with handcrafted morphological and contextual features
2. **Bi-LSTM** with word and character embeddings
3. **XLM-R (Transformer)** fine-tuned for sequence labelling

4.3 Evaluation Metrics

- Accuracy
- Precision, Recall, F1-score (macro-averaged)

Table 2. Performance Comparison of POS Tagging Models

Model	Accuracy (%)	Precision	Recall	F1-score
CRF	92.1	0.91	0.90	0.905
Bi-LSTM	94.8	0.94	0.94	0.94
XLM-R	96.2	0.96	0.96	0.96

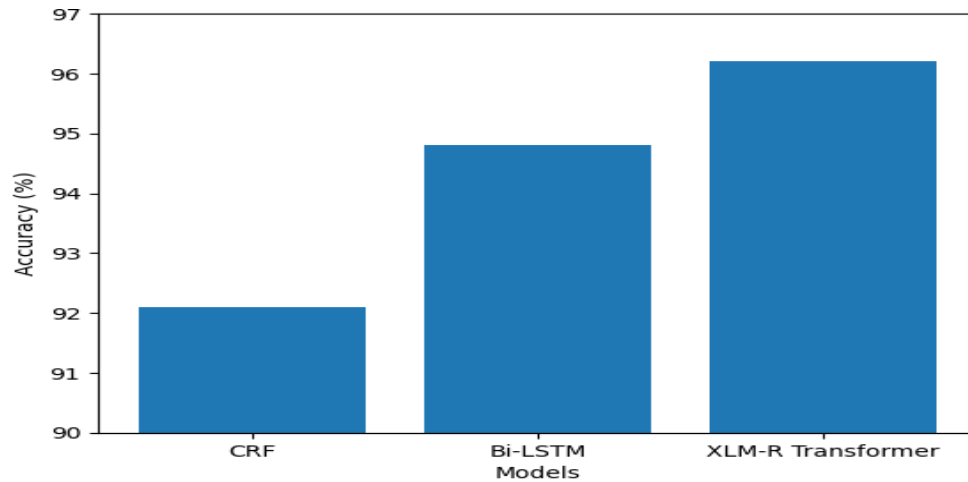


Figure 3. Comparative analysis of Tagging Model

Discussions

The experimental results confirm trends reported in prior literature. CRF remains competitive in low-resource settings but requires extensive feature engineering. Bi-LSTM significantly improves performance by learning contextual dependencies automatically, particularly for inflected forms.

Transformer-based XLM-R achieves the highest accuracy and robustness, especially for out-of-vocabulary words and long-distance dependencies. However, it incurs higher computational cost and training time, which may limit deployment in resource-constrained environments.

From a research perspective, the marginal improvement from Bi-LSTM to Transformer (~1.4%) must be weighed against infrastructure requirements, suggesting that lightweight transformers or hybrid deep models may offer optimal trade-offs.

Lack of large, publicly available standardized Gujarati corpora

- Limited reproducibility due to proprietary datasets
- High computational cost of transformers
- Need for lightweight and multilingual transfer-learning models
- Future work should focus on dataset creation, model compression, and cross-lingual learning to further advance Gujarati POS tagging.

Conclusions

This review and experimental study provide a consolidated view of two decades of Gujarati POS tagging research. The evolution from statistical models to transformers reflects broader NLP trends, with clear performance gains at each stage. Experimental validation under a unified setup confirms that transformer-based models currently offer the best accuracy, while hybrid and Bi-LSTM models remain practical alternatives for low-resource scenarios. The paper aims to serve as a strong reference for future researchers working on Indian language NLP.

Future research should focus on creating publicly available annotated corpora, optimizing lightweight transformer models for reduced computational cost, and exploring cross-lingual transfer learning to further advance POS tagging for low-resource Indian languages.

References

1. P. M. Bhatt and A. Ganatra, "POS-HOML: POS tagging technique for Gujarati language using hybrid optimal and machine learning approaches," *International Journal of Engineering Trends and Technology*, vol. 69, no. 11, pp. 1–7, 2021.
2. M. Prajapati and A. Yajnik, "POS tagging of Gujarati text using Viterbi and SVM," *International Journal of Computer Applications*, vol. 181, no. 43, pp. 18–22, 2019.
3. D. Shah, "Gujarati language POS tagging using hidden Markov model (HMM)," *International Journal of Creative Research Thoughts (IJCRT)*, vol. 8, no. 6, pp. 234–239, 2020.
4. C. Patel, "Part-of-speech tagging for Gujarati using conditional random fields," in *Proceedings of the International Joint Conference on Natural Language Processing (IJCNLP)*, Hyderabad, India, pp. 1–6, 2008.
5. J. Patel, A. Mehta, and S. Joshi, "Part of speech and morph category prediction for Gujarati," *International Journal of Intelligent Systems and Applications in Engineering*, vol. 12, no. 2, pp. 45–52, 2024.
6. P. Mishra, A. Gupta, and K. Sharma, "POS tagging for resource poor Indian languages through feature projection," in *Proceedings of the NLP AI Conference*, pp. 87–92, 2016.
7. D. Brahma and R. Basumatary, "Part-of-speech tagger for Bodo language using deep learning," *Journal of Intelligent Computing Applications*, vol. 3, no. 1, pp. 12–19, 2020.
8. S. Deshmukh and M. Joshi, "Deep learning-based parts-of-speech tagging in Marathi language," *Procedia Computer Science*, vol. 171, pp. 2171–2179, 2020.
9. K. Ramesh and R. Sundararajan, "Deep learning model for Tamil part-of-speech tagging," *Journal of King Saud University – Computer and Information Sciences*, vol. 31, no. 4, pp. 450–457, 2019.
10. T. Brants, "TnT: A statistical part-of-speech tagger," *Proceedings of the Sixth Applied Natural Language Processing Conference*, pp. 224–231, 2000.
11. J. Lafferty, A. McCallum, and F. Pereira, "Conditional random fields: Probabilistic models for segmenting and labeling sequence data," *Proceedings of ICML*, pp. 282–289, 2001.
12. C. Manning, H. Schütze, and P. Raghavan, *Introduction to Information Retrieval*, Cambridge University Press, 2008.
13. A. Bharati, V. Chaitanya, and R. Sangal, *Natural Language Processing: A Paninian Perspective*, Prentice Hall of India, 1995.
14. S. Sarkar and A. Chakraborty, "Part-of-speech tagging for Indian languages: A survey," *ACM Computing Surveys*, vol. 49, no. 3, pp. 1–34, 2017.
15. P. D. Patel and S. M. Patel, "Morphological analysis of Gujarati language," *International Journal of Computer Applications*, vol. 115, no. 20, pp. 1–6, 2015.

16. X. Ma and E. Hovy, "End-to-end sequence labeling via bi-directional LSTM-CNNs-CRF," Proceedings of ACL, pp. 1064–1074, 2016.
17. Y. Lample et al., "Neural architectures for named entity recognition," Proceedings of NAACL, pp. 260–270, 2016.
18. Z. Huang, W. Xu, and K. Yu, "Bidirectional LSTM-CRF models for sequence tagging," arXiv preprint arXiv:1508.01991, 2015.
19. J. Devlin et al., "BERT: Pre-training of deep bidirectional transformers for language understanding," Proceedings of NAACL, pp. 4171–4186, 2019.
20. A. Conneau et al., "Unsupervised cross-lingual representation learning at scale," Proceedings of ACL, pp. 8440–8451, 2020.
21. T. Wolf et al., "Transformers: State-of-the-art natural language processing," Proceedings of EMNLP: System Demonstrations, pp. 38–45, 2020.
22. S. Ruder, I. Vulić, and A. Søgaard, "A survey of cross-lingual word embedding models," Journal of Artificial Intelligence Research, vol. 65, pp. 569–631, 2019.
23. M. Joshi et al., "IndicBERT: A multilingual language model for Indian languages," Findings of EMNLP, pp. 1–9, 2020.
24. P. Mishra et al., "Leveraging multilingual pretrained models for low-resource Indian languages," Proceedings of COLING, pp. 1–10, 2022.