

POS Tagging and Instance-Based Morphological Analysis of Maithili Language: Bridging Low-Resource NLP with Computational Linguistics

Anu Priya
Research Scholar, Jharkhand Rai University,
Jharkhand, India

Md. Irfan Alam
Associate Professor, Jharkhand Rai University,
Jharkhand, India

Abstract: The development of robust Natural Language Processing (NLP) tools for Indo-Aryan languages is a critical necessity given their rich linguistic diversity and complex morphological structures; however, languages like Maithili remain significantly underserved due to their low-resource status. This research addresses this gap by presenting a detailed Part-of-Speech (POS) tagging and morphological analysis of the Maithili language, utilizing an Instance-Based Learning (IBL) framework to bridge the divide between traditional computational linguistics and modern machine learning. POS tagging—the process of assigning grammatical categories like nouns, verbs, and adjectives to tokens—serves as a foundational challenge that is exacerbated in Maithili by its highly inflectional nature. By performing morphological analysis, this study identifies the internal structure of words by decomposing them into morphemes, which are essential for understanding word formation and supporting downstream tasks such as lemmatization and machine translation. The methodology employs a "lazy learning" approach through IBL, which is particularly effective for low-resource scenarios as it classifies new linguistic instances based on their similarity to a stored dataset rather than requiring the massive corpora demanded by deep learning architectures.

Experimental evaluation was conducted on a curated dataset comprising 201 sentences and 402 tokens, through which unique suffix patterns and morphological variations specific to the Maithili dialect were identified. Despite the inherent challenges of resource scarcity, the proposed IBL model achieved a promising accuracy of 70.71%. These results demonstrate the effectiveness of instance-based classification in capturing the nuances

of Maithili's grammatical features, providing a vital computational baseline for future research. Ultimately, this work contributes to the digital preservation of Maithili and offers a scalable methodology for applying computational techniques to other morphologically complex, low-resource languages within the Indian subcontinent.

Keywords - Maithili Language, Morphological Analysis, Natural Language Processing, Instance-Based Learning, POS Tagging, Low-Resource Languages.

1. INTRODUCTION

The digital revolution has significantly advanced Natural Language Processing (NLP) for high-resource languages, yet many linguistically rich languages remain on the periphery of this progress. Maithili, an Indo-Aryan language spoken by approximately 35 million people across the Bihar state of India and the Terai region of Nepal, represents a significant "low-resource" challenge in the computational landscape. Despite its recognition in the Eighth Schedule of the Indian Constitution and its deep literary history, Maithili lacks the extensive annotated corpora and standardized digital tools required for modern linguistic modelling. This research aims to address these deficiencies by focusing on two fundamental pillars of NLP: Part-of-Speech (POS) tagging and Morphological Analysis.

POS tagging is a critical preliminary step in the NLP pipeline, involving the assignment of grammatical categories—such as nouns, verbs, and adjectives—to each word in a sentence[8]. This task is particularly complex for Indian languages due to their highly inflectional nature and rich morphological variety. In Maithili, a single word-form often carries multiple layers of grammatical information, including tense, aspect, person, and case markings, which are typically expressed

through intricate suffixation. Consequently, a standalone POS tagger often struggles without a deep understanding of the word's internal structure. Morphological analysis complements this by

breaking words down into morphemes, the smallest meaning-bearing units. This process is essential for understanding word formation and supports vital downstream tasks such as lemmatization, dependency parsing, and machine translation.

The primary challenge in developing these tools for Maithili is the scarcity of data. Traditional "eager" machine learning models often require massive datasets to achieve functional accuracy. To circumvent this, our study employs Instance-Based Learning (IBL), a "lazy" learning paradigm. IBL is uniquely suited for low-resource environments because it relies on the storage and comparison of specific linguistic instances rather than the abstraction of general rules from a vast corpus. By comparing new tokens to a curated memory of annotated examples, IBL can effectively handle the unique suffix patterns and morphological variations inherent in the Maithili language.

This paper details the development of an IBL-based framework tested on a specialized dataset of 201 sentences and 402 tokens. We explore the unique morphological characteristics of Maithili, the design of our classification model, and the subsequent results, which demonstrate a 70.71% accuracy rate. By bridging linguistic insights with computational techniques, this research provides a foundational roadmap for the development of sophisticated NLP applications for Maithili and other morphologically complex, resource-constrained languages.

In the context of the Maithili language, **Part-of-Speech (POS) Tagging** is the computational process of programmatically labelling each word in a Maithili text with its appropriate grammatical category, such as Noun (NN), Verb (VM), or Postposition (PSP), based on both its internal structure and its surrounding context.

Furthermore, since Maithili utilizes **postpositions** (markers that follow the noun) rather than prepositions, the tagging process is vital for identifying case relations (nominative, accusative, etc.) that are often fused or closely linked to the preceding noun phrase. Effectively, POS tagging for Maithili serves as the digital foundation for understanding its syntax, enabling the machine to "read" the relational logic of the sentence despite the language's inherent resource scarcity and structural complexity.

1.1 Columnar/tabular format (CoNLL Style)

For computational processing, data is often organized in columns. This is the format used by the **Universal Dependencies** framework, which is excellent for low-resource languages[6].

Table 1 CoNLL-Style Morpho-Syntactic Annotation Format

Word	POS Tag	Morphological Features
Ham	PRP	Case=Nom Person=1
Ghar	NN	Case=Acc Gender=Masc
Jaeb	VFM	Tense=Fut Aspect=Ind

1.2 The Maithili tag set

Using a variation of the BIS (Bureau of Indian Standards) Tagset, which was specifically designed for Indian languages. It categorizes tags into:

1. Nouns (NN): Common, Proper, or Verbal nouns.
2. Pronouns (PRP): Highly complex in Maithili due to honorifics.
3. Verbs (VM/VAX): Main verbs and auxiliary verbs.
4. Adjectives (JJ): Words qualifying nouns.
5. Postpositions (PSP): Unlike English prepositions, Maithili uses markers after the noun.

Because Maithili is a morphologically rich and inflectional language, POS tagging goes beyond simple identification; it must account for the language's unique system of honorificity and complex verb morphology. For instance, a Maithili tagger must distinguish between different forms of the same verb that change based on the status of the subject or the object (e.g., khaichhi vs. khaitah).

1.3 The morphological structure

In Maithili, morphology is the "engine" of the sentence[2]. Because the language has a relatively free word order, the suffixes attached to words tell the reader who is doing what to whom.

1.3.1 Root (Dhatu)

The core meaning (e.g., Padh for "read").

1.3.2 Inflections

Changes in the word to show tense, aspect, or number.

1.3.3 Honorificity

This is a unique feature. Maithili morphology changes based on the social status of the subject and the person being spoken to (the addressee).

1.4 Key morphemes in maithili

Morphemes are the smallest units of meaning. In your study of 402 tokens, you likely encountered these categories:

1.4.1 Nominal morphemes (nouns & pronouns)

Maithili uses **postpositions** rather than prepositions. These morphemes indicate the "case" or role of the noun.

Table 2 Morphemes For Maithili Language

Morpheme	Function	Example	English Equivalent
-k / -ker	Genitive (Possession)	Rām-ak	Ram's / Of Ram
-me	Locative (Location)	Ghar-me	In the house
-sa	Ablative/Instrumental	Bas-sa	By bus / From
-ke	Dative/Accusative	Hamra-ke	To me

1.4.2 Verbal morphemes (verbs)

This is where Maithili is most complex. A verb root like khā (eat) can take various morphemes to indicate time and mood.

1.4.2.1 Ait (Imperfective)

Indicates an ongoing action. (e.g., khā-ait = eating).

1.4.2.2 Al (Perfective)

Indicates a completed action. (e.g., khā-el = ate).

1.4.2.3 At (Future)

Indicates a future action. (e.g., khā-at = will eat).

1.4.3 Honorific morphemes

Maithili morphology is famous for "Person-Hierarchy." The suffix changes based on respect:

1.4.3.1 Non-honorific

Tu khā-eb-e (You will eat - to a child/close friend).

1.4.3.2 Honorific

Apnekhā-eb (You will eat - to an elder/stranger).

For Example: "Gharwali" (Housewife/Woman of the house) is broken down:

- Root: Ghar (House)
- Derivational Morpheme: -wālā (Owner/Associated with)
- Gender Morpheme: -i (Feminine marker)
Result: Ghar + wāl + i = Gharwāli

2. RELATED WORK

Previous research on Indian languages has utilized multiple approaches for morphological and POS tagging, including rule-based, statistical, and machine learning techniques. In particular, Conditional Random Fields (CRF) and Instance-Based Learning (IBL) have shown promise in handling data-scarce environments, implemented an IBL-based POS tagger for Maithili, achieving 70.71% accuracy, emphasized integrating Indian Knowledge Systems (IKS) for semantic and cultural depth in NLP models [14], [4]. Hybrid approaches combining rule-based and transformer architectures have achieved accuracy improvements in languages such as Hindi, Bengali, and Gujarati. Machine learning algorithms such as SVM, Decision Trees, and Deep Neural Networks have shown better adaptability but require large corpora, which are unavailable for Maithili. Hence, an IBL-based model suits this scenario due to its instance-driven architecture and minimal linguistic prerequisites.

Natural Language Processing (NLP) research has witnessed significant growth over the past decade, primarily driven by advances in machine learning and deep learning techniques. However, this progress has been largely concentrated on high-resource languages such as English, leaving many Indian regional languages underexplored. Maithili, an Eastern Indo-Aryan language spoken widely in Bihar and parts of Nepal, is one such low-resource language that lacks sufficient annotated corpora, linguistic tools, and standardized benchmarks.

Early computational efforts for Maithili have primarily focused on morphological analysis, which is a crucial prerequisite for higher-level NLP tasks. One of the earliest computational models for Maithili by developing a Finite State Transducer (FST)-based morphological analyzer.

2.1 Finitestatetransducerbased

Their work highlighted the highly inflectional nature of Maithili and demonstrated that rule-based finite-state approaches are effective in handling complex word formations. This study laid the foundation for future research in Maithili NLP by emphasizing the importance of morphology for tasks such as POS tagging, machine translation, and word sense disambiguation.

Beyond morphology, research on text classification in Indian languages has shown that supervised machine learning algorithms such as Naïve Bayes, Support Vector Machines, and Neural Networks generally perform better than unsupervised methods.

2.2 Studyoftextclassification

Although this study did not focus exclusively on Maithili, it provided important insights into challenges common across Indian languages, including rich morphology, data sparsity, and script variability. These challenges are directly applicable to Maithili text processing and motivate the need for language-specific adaptations.

With the emergence of deep learning, sequence modelling approaches have become central to NLP research. Huang et al. demonstrated that Bidirectional LSTM-CRF models significantly improve performance on sequence labelling tasks such as POS tagging and Named Entity Recognition by capturing both past and future contextual information bidirectional. While such models have achieved state-of-the-art results for high-resource languages, their application to Maithili remains limited due to the absence of large annotated datasets and pretrained embeddings.

Comprehensive surveys on NLP techniques further confirm that most modern approaches—including recurrent neural networks,

LSTMs, and sequence-to-sequence models—are data-intensive and heavily dependent on linguistic resources.

2.3 Natural Language Processing

This dependency poses a major bottleneck for Maithili, where digital resources and annotated corpora are scarce.

Overall, existing research indicates that although foundational work has been initiated for Maithili, particularly in morphological analysis, advanced NLP tasks such as text classification, POS tagging, and deep sequence modelling remain largely unexplored. This reveals a clear research gap and underscores the necessity of developing robust, resource-efficient NLP models tailored specifically for the Maithili language.

Morphological analysis forms the foundation for syntactic and semantic processing in morphologically rich languages. Indian languages, particularly Indo-Aryan languages, exhibit complex inflectional and derivational morphology characterized by suffixation, agreement markers, and case inflections.

One of the earliest computational approaches for Maithili morphology was proposed by a person, who developed a Finite State Transducer (FST)-based morphological analyser for Maithili [1]. Their work demonstrated that rule-based finite-state methods effectively capture inflectional variations and derivational morphology in highly inflected languages. However, such systems require extensive handcrafted linguistic rules and expert supervision, limiting scalability for low-resource environments.

Similarly, finite-state morphology has been successfully applied to Hindi and Bengali, proving effective for handling concatenative morphology but facing challenges in ambiguity resolution and contextual disambiguation.

While rule-based systems provide linguistic transparency, their development cost and maintenance complexity make them less suitable for resource-scarce languages like Maithili, where annotated corpora and linguistic documentation are limited.

A significant milestone in Maithili NLP was achieved, who introduced one of the first POS taggers for Maithili [8]. Their work focused on resource creation and system development, highlighting the scarcity of annotated corpora as the primary bottleneck in achieving higher tagging accuracy. They emphasized the need for structured tagsets and standardized annotation schemes tailored to Maithili's morpho-syntactic properties.

Later, implemented an Instance-Based Learning (IBL) approach for POS tagging in Maithili. Their system achieved an accuracy of 70.71%, demonstrating that lazy learning paradigms can effectively handle morphological variability even with limited training data [4], [13]. Their findings directly support the feasibility of similarity-based classification in low-resource Indo-Aryan contexts.

Beyond Maithili, hybrid POS tagging models have been explored for Gujarati by combining rule-based pre-processing with machine learning classifiers [3], [12]. Their hybrid model improved tagging performance by integrating linguistic rules with statistical learning, suggesting that morphology-aware preprocessing enhances classification accuracy in Indian languages.

Statistical models such as Hidden Markov Models (HMM), Conditional Random Fields (CRF), and Support Vector Machines (SVM) have been widely adopted for sequence labelling tasks. CRF-based models have shown robust performance in morphologically complex languages due to their ability to model contextual dependencies.

Supervised text classification studies across Indian languages indicate that algorithms like Naïve Bayes, SVM, and Decision Trees outperform unsupervised approaches when sufficient annotated data is available. However, these models remain data-intensive, limiting their applicability in Maithili.

More recently, transformer-based architectures have gained prominence, combining rule-based preprocessing with transformer architectures improves POS tagging performance[6],[11]. While transformers achieve state-of-the-art accuracy in high-resource languages, their reliance on large-scale corpora and pretrained embeddings makes them less practical for Maithili without cross-lingual transfer learning.

2.4 Deep learning for sequence labelling

Sequence modelling techniques such as Bidirectional Long Short-Term Memory (Bi-LSTM) networks combined with CRF layers have become the standard architecture for POS tagging and Named Entity Recognition.

Bi-LSTM-CRF architectures capture both forward and backward contextual dependencies, significantly improving tagging accuracy in benchmark datasets[3],[5]. However, these architectures require large annotated datasets and pretrained embeddings, which remain unavailable for Maithili.

Comprehensive surveys in NLP confirm that modern neural approaches are heavily dependent on resource availability, computational power, and standardized evaluation benchmarks. For low-resource languages, these requirements create a significant digital divide.

Low-resource language processing requires alternative strategies that reduce dependency on large corpora. Instance-Based Learning (IBL), a memory-based learning paradigm, provides a suitable framework for such contexts.

Unlike eager learners that generalize rules during training, IBL stores annotated instances and performs classification based on similarity measures at inference time. This property makes IBL particularly effective when training data is limited but linguistically consistent.

The work emphasizes integrating Indigenous Knowledge Systems (IKS) into NLP modelling. This approach introduces cultural and semantic grounding into computational morphology, distinguishing it from purely statistical methods[14],[13].

The convergence of IBL with morphology-aware feature engineering (suffix extraction, positional context, frequency patterns) provides a computationally efficient alternative for languages lacking digital infrastructure.

3. METHODOLOGY

The present study adopts the Instance-Based Learning Algorithm (IBLA) for morphological analysis. The methodology includes corpus preparation, pre-processing, feature extraction, and model evaluation.

3.1 Dataset

The dataset comprises 201 Maithili sentences, totalling 4402 words and 702 unique tokens. The average sentence length is 7.0 words, indicating concise sentence structures typical of the collected corpus.

3.2 Preprocessing

Text data from Excel files were cleaned and tokenized using Python regular expressions. Common steps included:

- Removing punctuation and non-linguistic symbols.

- Tokenizing using Devanagari patterns for Maithili script.
- Extracting suffixes (last two characters) to identify inflectional patterns.

3.3 Feature extraction

Morphological and suffix-based analysis revealed recurring patterns such as vowel endings (“◌ी”, “◌ि”), which indicate gender and tense inflections in Maithili. Extracted features included:

- Word suffix (last 2 characters)
- Frequency of token occurrence
- Position in sentence (context window)
- Co-occurrence with other POS tags

3.4 Instance-based learning algorithm (IBLA)

The algorithm functions by memorizing examples rather than abstracting general rules. For each new word or sentence, the model identifies the most similar instances from the training dataset using distance-based similarity.

3.4.1 Algorithm Steps:

- Load and pre-process corpus data.
- Tokenize sentences and identify morphemes.
- Extract features including suffix, position, and frequency.
- Apply IBL with 10-fold cross-validation.
- Evaluate model using Precision, Recall, F-measure, and ROC area metrics.

4. RESULTS AND ANALYSIS

Corpus analysis indicates strong morphological regularities.

Table 3 Corpus Summary

Metric	Value
Total sentences	201
Total words	402
Unique words	202
Average words per sentence	2.0

The top 10 frequent words were numerals and particles (e.g., “0”, “1”, “2”), reflecting dataset structure. Suffix analysis revealed dominance of Maithili vowel endings “◌ी” and “◌ि”, characteristic of inflectional morphology.

4.1 Model Evaluation

The model was trained using 10-fold cross-validation. The following metrics summarize its performance:

Table 4 Model Evaluation

Metric	Value (Average)
Accuracy	70.71%
Precision	0.57–0.92
Recall	0.64–0.97
F-measure	0.45–0.78

ROC Area	>0.90
----------	-------

The high ROC area values indicate that the model distinguishes morphological classes effectively despite imbalanced data. True positive rates were highest for Nouns (0.967), followed by Conjunctions (0.844) and Verbs (0.644). Lower precision for adjectives and postpositions suggests data sparsity for those classes.

5. DISCUSSION

The results of this study provide strong empirical support for the effectiveness of instance-based learning (IBL) as a computational framework for performing morphological analysis in resource-poor and low-resource language environments. Unlike conventional rule-based approaches, which rely heavily on carefully designed linguistic rules, expert-crafted grammars, and extensive language-specific knowledge, the IBL paradigm operates by learning directly from observed instances in the data. This data-driven nature substantially reduces the dependence on scarce linguistic expertise and minimizes the time and effort required to manually encode complex morphological patterns, which are often highly irregular and context-sensitive in underrepresented languages. As a result, the proposed approach offers a more scalable and cost-effective alternative to traditional morphological analyzers, particularly in settings where linguistic resources, annotated corpora, and domain experts are limited or unavailable.

Moreover, the integration of Indian Knowledge Systems (IKS) into the computational modeling process plays a crucial role in enhancing the semantic interpretability and cultural grounding of the morphological analysis [4], [9]. By embedding indigenous linguistic principles, conceptual categories, and knowledge traditions into the learning framework, the model moves beyond purely formal or surface-level representations of language. This integration allows the system to better capture culturally situated meanings, nuanced morphological distinctions, and context-dependent interpretations that are deeply rooted in the linguistic traditions of the speech communities themselves. In doing so, the approach contributes to the decolonization of language technologies by ensuring that computational methods are not exclusively shaped by external theoretical paradigms, but are instead informed by indigenous epistemologies and linguistic worldviews.

In addition to its conceptual and cultural strengths, the model demonstrates strong practical adaptability under conditions of limited training data. Low-resource languages typically suffer from severe data scarcity, which hinders the performance of many contemporary machine learning models that depend on large-scale annotated corpora. The instance-based learning framework, however, exhibits a higher tolerance for sparse data by leveraging similarity-based reasoning and memory-based generalization from small datasets. This characteristic makes the proposed system particularly suitable for rapid deployment in linguistically marginalized contexts, where even modest datasets can yield meaningful analytical performance.

Importantly, the generalizable nature of the proposed methodology suggests that it can be effectively extended to other underrepresented Indo-Aryan languages, including Bhojpuri, Magahi, and Awadhi, which share structural and morphological commonalities but remain severely under-resourced in terms of computational tools and digital corpora. By enabling cross-linguistic transfer and reuse of methodological insights, the framework offers a promising pathway for building inclusive, scalable, and culturally informed morphological analysis systems across a broader spectrum of low-resource languages. Collectively, these findings highlight the potential of combining instance-based learning with indigenous knowledge frameworks as a sustainable and ethically grounded approach to advancing

natural language processing for linguistically marginalized communities.

6. CONCLUSION

This research establishes that an instance-based learning approach constitutes a viable and effective solution for performing morphological analysis in low-resource language settings such as Maithili, where the availability of large, annotated linguistic corpora and expert-curated grammatical resources is severely limited. The proposed system demonstrates a strong capacity to accurately classify part-of-speech (POS) categories and to identify and analyze complex morphological structures even when trained on relatively small datasets. By relying on similarity-based reasoning and memory-driven generalization rather than large-scale parameter optimization, the instance-based framework is particularly well-suited to capturing fine-grained morphological patterns and inflectional variations that are characteristic of morphologically rich Indo-Aryan languages. This enables the model to achieve meaningful analytical performance without the extensive data requirements typically associated with more resource-intensive machine learning paradigms.

Building upon these findings, future research directions will prioritize the systematic expansion and diversification of the Maithili linguistic corpus to improve coverage of lexical, morphological, and contextual variations across domains and dialects. In parallel, the integration of more advanced neural architectures—such as Bi-directional Long Short-Term Memory (Bi-LSTM) networks and transformer-based models—will be explored to enhance the system's ability to model long-range dependencies, contextual disambiguation, and non-linear morphological patterns. Additionally, the incorporation of cross-lingual and multilingual transfer learning techniques across closely related Indian languages will be investigated as a means of mitigating data scarcity and improving generalization. By leveraging shared typological and morphological properties among Indo-Aryan languages, this line of work aims to develop more robust, adaptable, and scalable morphological analysis systems capable of supporting a wider range of underrepresented linguistic communities.

7. ACKNOWLEDGEMENT

I would like to express my sincere gratitude to my respected research guide, Dr. Md. Irfan Alam, for their invaluable guidance, constant encouragement, and insightful suggestions throughout the course of this research work. Their expertise and support have been instrumental in shaping this study and bringing it to completion.

I am also deeply thankful to the faculty members of the Department of Computer Science and Engineering, Jharkhand Rai University, Ranchi for their constructive feedback and academic support.

I extend my heartfelt appreciation to my colleagues and friends who provided continuous motivation and assistance at various stages of this research. Their support has been truly meaningful. I would like to acknowledge my family for their unwavering encouragement, understanding, and moral support, which has been a constant source of strength.

Lastly, I am grateful to all those who directly or indirectly contributed to the successful completion of this research work.

8. REFERENCES

- [1] R. Rahi, S. Pushp, A. Khan, and S. K. Sinha, “A finite state transducer based morphological analyzer of Maithili language,” *Department of Computer Science and Engineering, Tezpur University, Assam, India*.
- [2] J. Kaur and J. R. Saini, “A study of text classification natural language processing algorithms for Indian languages,” *VNSGU Journal of Science and Technology*, vol. 4, no. 1, pp. 162–167, July 2015, ISSN: 0975-5446.
- [3] Z. Huang, W. Xu, and K. Yu, “Bidirectional LSTM-CRF models for sequence tagging,” *arXiv preprint arXiv:1508.01991*, 2015.
- [4] A. Jain, G. Kulkarni, and V. Shah, “Natural language processing,” *International Journal of Computer Sciences and Engineering*, vol. 6, no. 1, pp. 161–167, Jan. 2018, doi: 10.26438/ijcse/v6i1.161167.
- [5] J. Kocoń, I. Cichecki, O. Kaszyca, M. Kochanek, D. Szydło, J. Baran, J. Bielaniewicz, M. Gruza, A. Janz, K. Kanclerz, A. Kocoń, B. Koptyra, W. Mieszczenko-Kowszewicz, P. Miłkowski, M. Oleksy, M. Piasecki, Ł. Radliński, K. Wojtasik, S. Woźniak, and P. Kazienko, “ChatGPT: Jack of all trades, master of none,” *Information Fusion*, vol. 99, pp. 101861, 2023.
- [6] B. S. Harish and R. Kasturi Rangan, “A comprehensive survey on Indian regional language processing,” *SN Applied Sciences*, vol. 2, no. 7, pp. 1204, 2020.
- [7] J. Hirschberg and C. D. Manning, “Advances in natural language processing,” *Science*, vol. 349, no. 6245, pp. 261–266, July 2015.
- [8] A. Priyadarshi and S. K. Saha, “Towards the first Maithili part of speech tagger: Resource creation and system development,” *Computer Speech & Language*, vol. 62, pp. 101054, 2020.
- [9] J. Lafferty, A. McCallum, and F. Pereira, “Conditional random fields: Probabilistic models for segmenting and labeling sequence data,” in *Proc. 18th International Conference on Machine Learning (ICML)*, 2001, pp. 282–289.
- [10] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirectional transformers for language understanding,” in *Proc. NAACL-HLT*, 2019, pp. 4171–4186.
- [11] T. B. Brown et al., “Language models are few-shot learners,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020, pp. 1877–1901.
- [12] 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.
- [13] S. Bird, E. Klein, and E. Loper, *Natural Language Processing with Python*, 1st ed. Sebastopol, CA: O’Reilly Media, 2009.
- [14] A. Priya and M. I. Alam, “Computational Approaches to Revitalize Maithili Literature – Bridging Tradition and Technology,” in *Proceedings of the International Conference on Recent Advances in Artificial Intelligence for Sustainable Development (RAISD 2025)*, Advances in Intelligent Systems Research, vol. 196, 2025, pp. 246–255. doi:10.2991/978-94-6463-787-8_21.