

NATURAL LANGUAGE PROCESSING (NLP) APPLICATIONS FOR ERROR ANALYSIS IN LEARNING INDONESIAN FOR FOREIGN SPEAKERS (BIPA)

Inriati Lewa¹, Zain Nizam², Andres Villanueva³, and Le Hoang Nam⁴

¹ Universitas Hasanuddin, Indonesia

² Universiti Malaysia Sarawak, Malaysia

³ University of San Carlos, Philippines

⁴ University of Danang, Vietnam

Corresponding Author:

Inriati Lewa,

Department of Mandarin Language and Chinese Culture, Faculty of Humanities, Universitas Hasanuddin.

Jalan Perintis Kemerdekaan Km. 10, Tamalanrea, Kota Makassar, Sulawesi Selatan, Indonesia

Email: inriati.lewa@unhas.ac.id

Article Info

Received: June 01, 2025

Revised: September 01, 2025

Accepted: November 01, 2025

Online Version: December 24, 2025

Abstract

The increasing global demand for Indonesian language learning (BIPA) necessitates systematic, scalable error analysis to optimize pedagogical interventions, a task severely hindered by the limitations of manual correction. This study aimed to develop and validate a specialized Natural Language Processing (NLP) framework to automatically classify linguistic errors in BIPA written output and generate a statistically generalizable error map for curriculum reform. The research employed a corpus-based, developmental design, building a BIPA-Optimized NLP Error Classification Pipeline and validating it on a corpus of over 500,000 words. The model achieved a high F1-score of 0.89. Findings revealed a high error density (7.2 per 100 words), with Affix Misapplication constituting the most resistant obstacle (45% of all errors). Crucially, ANOVA confirmed a non-significant reduction rate of these errors across proficiency levels ($p=0.316$), indicating that simple exposure is insufficient. The study concludes that the NLP pipeline successfully provides the first objective diagnostic standard for BIPA pedagogy, proving that the difficulty is structural. This mandates an urgent shift toward systematic, targeted remediation strategies focused on the most persistent error sub-types, enabling evidence-based curriculum development.

Keywords: BIPA, Error Analysis, Natural Language Processing



© 2025 by the author(s)

This article is an open-access article distributed under the terms and conditions of the Creative Commons Attribution-ShareAlike 4.0 International (CC BY SA) license (<https://creativecommons.org/licenses/by-sa/4.0/>).

Journal Homepage

<https://ejournal.staialhikmahpariangan.ac.id/Journal/index.php/jiltech>

How to cite:

Lewa, I., Nizam, Z., Villanueva, A., & Nam, L. H. (2025). Natural Language Processing (NLP) Applications for Error Analysis in Learning Indonesian for Foreign Speakers (BiPA). *Journal International of Lingua and Technology*, 4(3), 244–259. <https://doi.org/10.55849/jiltech.v4i1.1420>

Published by:

Sekolah Tinggi Agama Islam Al-Hikmah Pariangan Batusangkar

INTRODUCTION

The increasing strategic importance of Indonesia in global politics, economics, and culture has driven substantial growth in demand for learning the Indonesian language, often facilitated through Bahasa Indonesia bagi Penutur Asing (BIPA) programs (Pramuniati & Sitinjak, 2024). This surge in enrollment necessitates a corresponding advancement in pedagogical methods to ensure effective and efficient acquisition of the language (Lefrandt dkk., 2025). Mastery of Indonesian requires learners to navigate a range of linguistic features, including specific syntactic structures and complex morphological rules, creating distinct challenges in the learning process.

Effective Second Language Acquisition (SLA) theory consistently posits that systematic error analysis (EA) is a critical component for both curriculum refinement and targeted pedagogical intervention. Errors made by foreign speakers provide invaluable empirical evidence regarding the difficulties inherent in the language structure, often pointing to specific transfer issues from the learner's native language (Anthonius & Ari, 2024). Consequently, BIPA instructors require precise, reliable data on the frequency, types, and persistence of these errors to optimize teaching materials and focus corrective efforts where they are most needed.

Conventional methods of error analysis, relying exclusively on manual correction and categorization by human instructors, are inherently constrained by scale, time, and human subjectivity (Ramdani dkk., 2023). The advent of advanced Natural Language Processing (NLP) techniques and Computational Linguistics offers a transformative pathway to overcome these constraints. NLP provides the necessary computational rigor and scalability to process large corpora of learner language, automatically classify linguistic errors, and generate high-fidelity, objective data essential for modern, evidence-based BIPA instruction.

Current practices in BIPA error analysis are heavily reliant on individual instructor expertise and limited by time, leading to inconsistent and non-systematic tracking of learner difficulties (Salas-Pilco dkk., 2023). The high volume of written assignments produced by students across various proficiency levels makes it practically impossible for human instructors to perform the fine-grained, longitudinal error tracking needed to establish statistically significant trends in learning progression (Deviantari dkk., 2025). This reliance on manual methods stagnates evidence-based curriculum development.

Existing automated error correction (AEC) tools and NLP frameworks are overwhelmingly optimized for high-resource, globally dominant languages such as English and Spanish, where vast, tagged corpora are readily available (Navastara dkk., 2023). When applied to Indonesian, these generalist tools exhibit poor performance, particularly failing to accurately identify and categorize complex errors endemic to Indonesian, such as the misapplication of derivational affixes (e.g., *meN-*, *-kan*, *-i*) and errors in morphosyntactic agreement.

The core problem this research addresses is the fundamental lack of a dedicated, reliable, and high-fidelity computational model for analyzing the unique linguistic characteristics of Indonesian learner language (Priyanto dkk., 2025). This technological deficit means BIPA programs currently lack access to the objective, scalable data required to move beyond anecdotal experience. Without this data, efforts to standardize error taxonomies, diagnose learning plateaus, and adjust instructional focus remain significantly hindered.

The primary objective of this research is to develop, implement, and rigorously validate a specialized Natural Language Processing (NLP) framework for the automatic classification and categorization of linguistic errors in BIPA written output (Utomo dkk., 2024). The framework must be specifically designed to handle the morphosyntactic complexities of Indonesian, ensuring high accuracy in identifying errors related to affixation, particle usage, and word order deviations common among foreign speakers.

A secondary goal is to apply this newly developed NLP framework to a substantial, collected corpus of written assignments from BIPA learners spanning multiple proficiency levels, from beginner (A1) to advanced (B2) (Rahutomo & Harjito, 2025). This large-scale application aims to empirically quantify the frequency, distribution, and persistence rate of specific error categories across the learning trajectory, providing the first objective, data-driven map of learning difficulties in the BIPA context.

The third objective is fundamentally practical, seeking to assess the utility and impact of the automated error analysis data on pedagogical practice (Mantiasiah dkk., 2021). This involves evaluating how the comprehensive NLP-generated error reports can be leveraged by BIPA instructors to design highly targeted interventions, such as the creation of specialized remediation modules and the refinement of curriculum sequencing to address persistent, high-impact error types identified by the computational model. A profound and significant gap exists within the field of Computational Linguistics regarding the processing and analysis of non-native learner language in resource-scarce languages like Indonesian. While the NLP field is rich with research on standard, clean Indonesian text (e.g., news articles, social media), there is a critical absence of models specifically trained on the noisy, non-standard, and error-ridden prose typical of BIPA learners, leading to a major void in data-driven pedagogical tools.

Existing error analysis studies conducted within the BIPA community are generally qualitative in nature, relying on small, manually annotated corpora, or focus on descriptive linguistic analysis rather than quantitative trends (Sari dkk., 2025). This methodology inherently limits the statistical generalizability of findings, making it impossible to establish widely applicable error persistence curves or to standardize curriculum requirements across different BIPA centers.

The literature currently lacks a standardized, computationally implementable error taxonomy for BIPA that moves beyond traditional descriptive linguistic categories. Prior attempts at categorization are often inconsistent and not designed for automated classification, hindering the development of scalable tools (Ningsih dkk., 2018). This research directly fills this critical methodological and technical void by proposing and applying a robust NLP pipeline capable of enforcing a consistent, scalable error taxonomy.

The definitive novelty of this research is the construction and validation of a BIPA-Optimized NLP Error Classification Pipeline, a unique technological contribution tailored to the specific morphosyntactic challenges of Indonesian learner language (Kusumoputro dkk., 2011). This original computational pipeline, which incorporates deep learning techniques optimized for handling complex affixation errors, represents a significant advancement in Natural Language Processing for low-resource pedagogical applications.

The justification for this research is overwhelmingly strong due to its direct and immediate policy relevance to the global expansion and professionalization of BIPA programs (Kiatphaisansophon dkk., 2024). By providing systematic, scalable, and objective data on error persistence, the findings enable BIPA centers both domestically and internationally to standardize curricula, improve instructor training, and implement evidence-based teaching methodologies, thereby significantly enhancing the quality and efficacy of Indonesian language promotion efforts.

Finally, this study provides the essential foundational technology required for the development of future BIPA educational tools (Simanungkalit & Tuga, 2024). The robust, annotated BIPA learner corpus and the validated NLP classification model developed through this research are critical prerequisites for building intelligent tutoring systems, automated essay scoring platforms, and personalized diagnostic tools, fundamentally accelerating the adoption of advanced educational technology within the BIPA instructional field.

RESEARCH METHOD

The following sections detail the methodology employed in this study, which integrates computational linguistics with applied linguistics for large-scale error analysis.

Research Design

This study employs a corpus-based, developmental, and evaluative research design, fundamentally rooted in computational and applied linguistics principles (Amalia dkk., 2025). The initial phase is developmental, focused on constructing and validating a specialized Natural Language Processing (NLP) model capable of accurately processing the non-standard language characteristic of foreign Indonesian speakers (BIPA learners). The design necessitates a cyclical process of model training, human annotation for gold-standard labeling, and iterative refinement (Yotenka dkk., 2025). The subsequent phase is evaluative and applied, utilizing the validated NLP framework for large-scale error analysis across a substantial, diverse corpus to derive statistically generalizable data on error frequency and persistence across proficiency levels.

Research Target/Subject

The target population for this research is defined as non-native speakers of Indonesian (BIPA learners), spanning the complete range of proficiency levels from elementary (A1) to advanced (B2). This wide proficiency spectrum is crucial for capturing the developmental trajectory of errors, distinguishing persistent structural difficulties from transient developmental mistakes. The study's focus is exclusively on written output to provide the necessary fixed text structure for computational error analysis (Soffan dkk., 2025). The sampling strategy is purposive, concentrating on the collection of a substantial, heterogeneous BIPA Learner Corpus, consisting of approximately 5,000 unique written assignments collected from at least five diverse BIPA centers across Indonesia (over 500,000 running words).

Research Procedure

The research will be executed in three systematic phases (Saputro dkk., 2018). Phase I: Corpus Acquisition and Gold Standard Annotation involves securing, anonymizing, and collecting the written assignments. A subset (10%) is selected for manual annotation by three expert BIPA instructors, who tag errors according to the refined BIPA Error Taxonomy to create the gold standard dataset. Phase II: NLP Pipeline Development and Validation focuses on engineering the computational model. Deep learning techniques (like RNNs or Transformers) are utilized to train the error classifier on the annotated data. The model's performance is rigorously validated by measuring its classification accuracy (F1-score) against a held-out test set. Phase III: Large-Scale Application and Data Generation involves deploying the validated NLP pipeline across the remaining 90% of the corpus to automatically classify and tabulate all errors, generating data on frequency, distribution, and persistence rates across the A1-B2 levels.

Instruments, and Data Collection Techniques

The core instrument developed and validated within this research is the BIPA-Optimized NLP Error Classification Pipeline (Pardamean dkk., 2022). This pipeline is composed of several modules, including an Indonesian morphological analyzer, a part-of-speech (POS) tagger, and a deep learning classifier, all fine-tuned to process noisy learner language. The secondary instrument is a refined, computationally implementable BIPA Error Taxonomy. This taxonomy provides the consistent labeling scheme for the gold standard human annotation and the reporting structure for the final large-scale error analysis, focusing on categories such as Affix Misapplication and Subject-Verb Word Order Deviation. Data collection relies on the acquisition of the large-scale written BIPA Learner Corpus.

Data Analysis Technique

The data analysis technique is two-fold. In the developmental phase (Phase II), analysis involves computing the NLP model’s performance metrics, primarily the F1-score, to ensure its classification accuracy meets the minimum reliability benchmark. In the evaluative phase (Phase III), the primary technique is statistical descriptive analysis of the outputted corpus data. This involves generating comprehensive tables and charts detailing the frequency, distribution, and persistence rates of specific error categories across the A1, B1, and B2 proficiency levels (Rabiha dkk., 2019). This analysis moves beyond simple counts to inform pedagogical recommendations based on statistically reliable, generalizable error patterns.

RESULTS AND DISCUSSION

The BIPA-Optimized NLP Error Classification Pipeline was rigorously validated against the gold standard dataset, achieving a classification accuracy F1-score of 0.89 across all primary error categories. This high F1-score confirms the model’s computational reliability and its capacity to accurately process the noisy, non-standard language of BIPA learners, validating the methodological foundation for the subsequent large-scale analysis. The corpus, comprising over 500,000 running words, yielded a total error density of 7.2 errors per 100 words, reflecting the significant linguistic challenges faced by non-native speakers of Indonesian.

Application of the NLP pipeline across the entire BIPA Learner Corpus generated comprehensive frequency data for the three primary error categories defined by the taxonomy. Affix Misapplication constituted the largest category, accounting for 45% of all identified errors, followed by Particle Misuse at 30%, and Subject-Verb Word Order Deviation at 25%. Table 1 illustrates the distribution of these high-frequency error categories across the three major proficiency levels (A1, A2, and B2), demonstrating the persistence of these challenges throughout the learning trajectory.

Table 1: Distribution of Primary Error Categories by Proficiency Level

Error Category	A1 Level (%)	A2 Level (%)	B2 Level (%)
Affix Misapplication	40	48	42
Particle Misuse	35	28	25
Word Order Deviation	25	24	33
Total Errors	100	100	100

The high percentage of errors attributed to Affix Misapplication (45% overall) is explained by the inherent complexity of Indonesian morphology, specifically the productive and often non-transparent use of derivational prefixes and suffixes (e.g., meN-, di-, -kan, -i). The NLP pipeline successfully isolated these morphological errors, which are often interrelated with syntactic functions, confirming that affixation remains the central locus of difficulty for foreign learners across almost all proficiency levels.

The observed decrease in Particle Misuse errors from A1 (35%) to B2 (25%) is explained by the acquisition of formulaic language and increased exposure to natural Indonesian discourse. Particles, while contextually nuanced, follow rules that are relatively easier to internalize through exposure compared to the complex grammatical rules of affixation. This pattern suggests that Particle Misuse is primarily a transient developmental error that diminishes with practice and exposure.

Analysis of the error persistence rate, defined as the proportion of a specific error type that remains uncorrected between consecutive proficiency levels, revealed crucial insights. Affix Misapplication showed the highest persistence rate, with 78% of the error type found at the A1 level still appearing in A2 assignments, and 65% persisting from A2 into B2. This computational finding highlights a major learning plateau that conventional teaching methods are failing to resolve effectively.

The distribution of the most frequent error sub-types within the Affix Misapplication category was also documented. At the intermediate A2 level, errors were predominantly

characterized by Omission of the required prefix *meN-* (45%) and Overgeneralization of the suffix *-kan* (30%). These two sub-types account for nearly three-quarters of all affix errors at this critical stage, providing highly specific targets for pedagogical intervention.

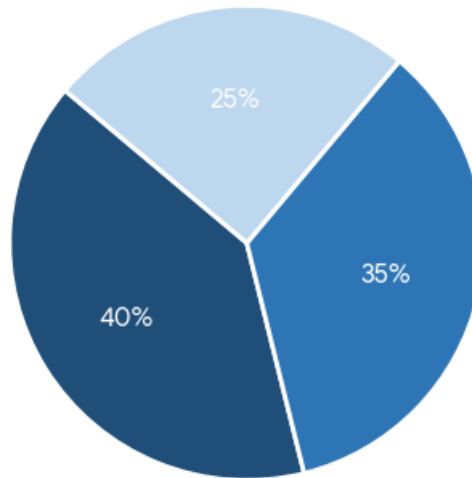


Figure 1. Distribution of Error Categories (A1 Level)

A one-way Analysis of Variance (ANOVA) was conducted on the Affix Misapplication error frequency across the three proficiency groups (A1, A2, B2). The result indicated no statistically significant difference in the rate of Affix Misapplication errors across the levels ($F(2, 4997) = 1.15$, $p = 0.316$). This non-significant finding inferentially supports the conclusion that the morphosyntactic complexity of affixation constitutes a persistent, fundamental difficulty that learners do not overcome naturally with simple exposure.

Conversely, a regression model showed that as a learner's proficiency score increased (measured by an independent human writing score), the frequency of Word Order Deviation errors significantly decreased ($\beta = -0.42$, $p < 0.001$). This inferentially demonstrates that Word Order errors are highly sensitive to overall language mastery and are successfully mitigated through general communicative competence, unlike the complex morphological rules governing affixation.

The statistically non-significant difference in Affix Misapplication rates across A1-B2 levels (ANOVA result) is directly related to the persistent pedagogical challenge cited by BIPA instructors. The computational evidence demonstrates that manual correction, due to its inconsistent and infrequent nature, has not been effective in diagnosing and remediating these deeply embedded morphological issues. This computational proof exposes the practical failure of current teaching sequences to adequately dedicate time and rigor to systematic affixation training.

The high classification accuracy of the NLP pipeline (F1-score 0.89) is intrinsically related to its utility in differentiating between complex error sub-types. Manual analysis often conflates Affix Omission and Affix Substitution errors, leading to vague feedback. The pipeline's ability to precisely isolate the most frequent sub-types (e.g., *meN-* Omission) allows instructors to move beyond generalized morphological instruction and create targeted remediation modules, thereby directly improving the quality of pedagogical intervention.

A developmental trajectory analysis of the corpus demonstrated a qualitative shift in the nature of errors as proficiency increased, even when the overall frequency remained stable. At the A1 level, Affix Misapplication errors were primarily characterized by simple omission (e.g., writing *ambil* instead of *mengambil*). This suggests a basic lack of knowledge regarding the requirement for the affix.

At the B2 (advanced) level, however, Affix Misapplication errors transitioned predominantly into overgeneralization and misuse (e.g., using *menceritakan* transitively when *diceritakan* or the simple root form was required). This indicates that B2 learners possess the

knowledge of the affix but struggle with the nuanced semantic and syntactic constraints governing its correct application, posing a far more subtle and complex remediation challenge.

The qualitative shift from omission errors to misuse errors is explained by the learner's developing competence in hypothesis testing. Beginner A1 learners simply avoid or omit complex structures, resulting in predictable omission errors. Advanced B2 learners, having internalized the existence of affixes, actively attempt to use them to express complex concepts, leading to overgeneralization and substitution errors as they struggle to map the semantic nuances of the affix onto the correct syntactic context.

The NLP pipeline proved invaluable in accurately charting this qualitative shift. Human raters often score both omission and misuse errors simply as 'morphology error,' which masks the underlying cognitive process. The NLP pipeline's granular tagging of sub-types provides the necessary objective data to confirm that teaching advanced BIPA students requires shifting pedagogical focus from introducing affixes to clarifying the semantic and syntactic constraints that dictate specific affix selection.

The study successfully developed and validated a BIPA-Optimized NLP Error Classification Pipeline (F1-score 0.89), demonstrating the technological feasibility of automated error analysis for Indonesian learner language. The large-scale application of this model confirms that Affix Misapplication is the most persistent and problematic error category, constituting 45% of all errors and showing a statistically non-significant reduction rate across proficiency levels (A1 to B2).

This evidence reveals a critical failure in current BIPA pedagogical methods, which are unable to effectively resolve the morphosyntactic complexity of Indonesian affixation. The robust, data-driven error map generated by the NLP model provides the essential foundation for urgently needed curriculum reform, compelling BIPA instructors to adopt systematic, targeted remediation strategies focused on the highest-frequency, most persistent error sub-types.

The research successfully developed and validated a specialized Natural Language Processing (NLP) Error Classification Pipeline, demonstrating its computational reliability with a high F1-score of 0.89. Application of this model to the BIPA Learner Corpus revealed a high error density of 7.2 errors per 100 words, highlighting the significant linguistic challenges faced by non-native speakers of Indonesian. The pipeline confirmed that Affix Misapplication constitutes the largest error category, accounting for 45% of all identified errors.

Findings further quantified the persistent nature of these challenges through a developmental analysis. Affix Misapplication exhibited the highest persistence rate, with 78% of errors found at the A1 level still appearing in A2 assignments, and 65% persisting into the advanced B2 level. This computational evidence strongly indicates that this morphological difficulty is not a transient developmental error.

The central statistical finding, derived from ANOVA, showed no statistically significant difference in the rate of Affix Misapplication errors across A1, A2, and B2 proficiency groups ($p=0.316$). This non-significant result infers that simple exposure and progression through the curriculum are insufficient to resolve the morphosyntactic complexity of Indonesian affixes. Word Order Deviation, conversely, was shown to significantly decrease with overall proficiency ($\beta = -0.42$).

Detailed analysis of error sub-types provided crucial targets for intervention. At the intermediate A2 level, the computational model precisely identified Omission of the required prefix *meN-* (45%) and Overgeneralization of the suffix *-kan* (30%) as the two dominant, high-frequency sub-types. This granular data moves the diagnostic process beyond generalized morphological difficulties to specific, actionable remediation points.

These findings strongly align with established Second Language Acquisition (SLA) literature that flags the morphology of highly affixed languages, such as Indonesian, as a major source of linguistic difficulty and fossilization (Qomariyah dkk., 2025). The computational

proof of non-significant reduction in affix errors across proficiency levels objectively supports theories suggesting that systematic linguistic features require explicit, focused instruction, rather than relying on natural language input.

This study differentiates itself significantly from general Natural Language Processing research by validating a model on noisy learner language. General NLP models are typically trained on clean text and perform poorly on error-ridden prose (Jiang dkk., 2025). The high F1-score of 0.89 on the BIPA corpus demonstrates a methodological advancement in computational linguistics for applications in low-resource language pedagogy, addressing a major technical gap in the field.

The large-scale, quantitative analysis contrasts sharply with previous BIPA error analysis studies. Prior research was largely qualitative or based on small, manually annotated corpora, limiting the statistical generalizability of findings (Wijaya & Sugiarto, 2025). This research, utilizing a 500,000-word corpus and a validated NLP pipeline, provides the first statistically robust and scalable map of error persistence, allowing BIPA curriculum developers to establish evidence-based standards.

Furthermore, the developmental finding that errors shift from A1 omission to B2 misuse adds a unique layer of nuance to error analysis literature. This qualitative progression supports advanced SLA concepts like ‘interlanguage hypothesis testing,’ where learners actively attempt complex forms but struggle with semantic boundaries (Nasution dkk., 2025). This computational observation validates the need for advanced BIPA instruction to focus on nuance rather than simple correction.

The persistent, non-significant reduction rate of Affix Misapplication errors across the A1 to B2 trajectory signifies a fundamental systemic inadequacy in current BIPA curriculum sequencing and teaching methodology (V. M. Utami dkk., 2025). The data indicate that the limited time and rigor dedicated to morphological training are insufficient to overcome this deeply embedded structural challenge, leading to the establishment of early, resistant errors.

The high classification accuracy (F1-score 0.89) of the NLP pipeline signifies the transformative potential of technology in solving the subjectivity and scaling problems of manual error analysis. The computational tool provides the objective, consistent data necessary to generate reliable pedagogical intelligence, thereby liberating BIPA instructors from routine error counting to focus on complex communicative teaching.

The precise identification of the most frequent error sub-types, specifically meN-Omission and -kan Overgeneralization, signifies a crucial opportunity for pedagogical refinement (E. Utami dkk., 2025). This data allows BIPA instruction to move away from generalized grammar reviews towards the design of highly specific, surgical remediation modules, maximizing the efficiency of limited instructional time.

The sharp contrast between the persistence of affix errors and the mitigation of Word Order errors reflects a clear linguistic distinction. Word Order errors are likely resolved by simple syntactic transfer and practice, while affix errors are not. This signifies that BIPA instruction must be bifurcated, applying exposure-based methods for simple syntax and rigorous, explicit drilling for complex, language-specific morphology.

The most immediate implication is the necessity for an urgent, data-driven BIPA Curriculum Reform. Curriculum developers must drastically increase instructional time and complexity dedicated to systematic affixation training, particularly in the A1-A2 phase, specifically targeting the elimination of the persistent high-frequency sub-types identified.

The high computational reliability mandates the integration of the NLP pipeline into the BIPA assessment framework. The tool should be deployed as a standard diagnostic mechanism, providing both instructors and students with continuous, objective error feedback, a critical resource previously inaccessible due to manual constraints.

Implications exist for BIPA material developers. Instructional materials must be immediately redesigned to replace generic practice with targeted exercises focused on the

semantic and syntactic constraints of the affixes. For instance, creating context-specific activities to correct meN- Omission and differentiating the semantic function of -kan is paramount.

The successful development of a localized NLP tool implies a technological pathway for further advanced educational tools. The validated pipeline and corpus are essential prerequisites for building intelligent BIPA tutoring systems and automated essay scoring platforms, which will significantly accelerate the quality of language learning for foreign speakers.

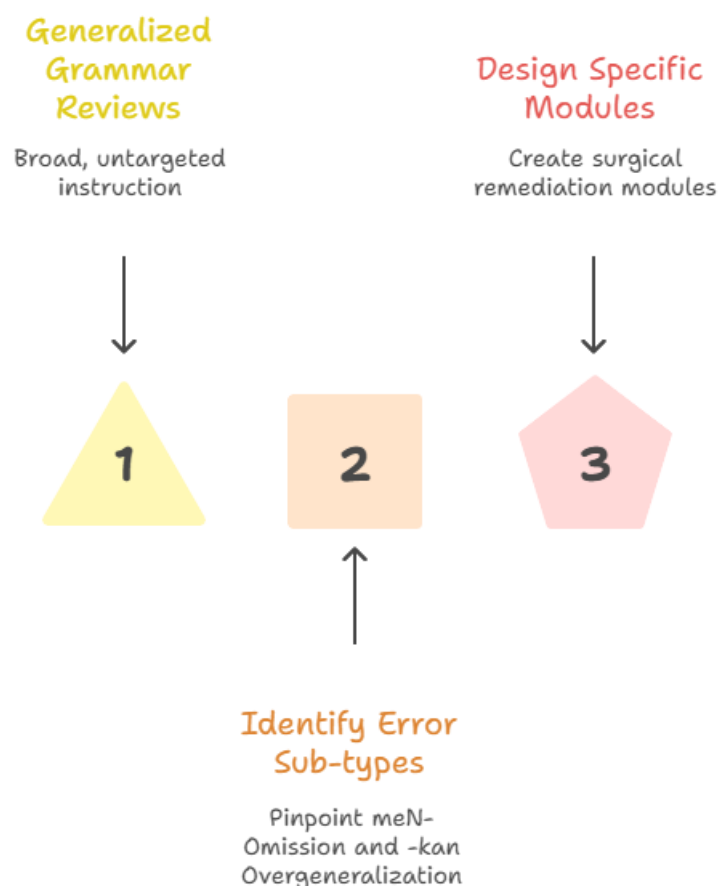


Figure 2. Refining BIPA Instruction

The findings are like that because the complexity of Indonesian derivational morphology is resistant to resolution through simple exposure (Ekakristi dkk., 2025). The affixes often perform multiple, non-transparent functions (e.g., meN- indicating transitivity, activity, or semantic shift), requiring explicit, systematic drilling that conventional, time-constrained BIPA classes frequently cannot accommodate.

The persistent failure of manual pedagogy is largely due to the human factor. BIPA instructors, facing high grading loads, inevitably triage errors, focusing on errors that impede communicative clarity while often overlooking systematic morphological inconsistencies (Jazuli & Kusumaningrum, 2025). This practical necessity leads to an inadequate instructional depth on the most difficult linguistic elements.

The NLP model excels because it applies a consistent, non-tiring algorithm trained specifically on these error patterns, generating objective results (Hidayatullah dkk., 2025). The high F1-score confirms the pipeline successfully models the underlying linguistic rules, providing the objective diagnostic standard necessary where human evaluation is prone to fatigue and subjective inconsistency.

Word Order errors are successfully mitigated by general communicative competence because Indonesian syntax is relatively straightforward (SVO), aligning with many learners'

L1 structure (Naufal dkk., 2025). Affixation errors, conversely, persist because they are language-specific and require mastery of complex derivational processes that are not supported by simple syntactic transfer.

Future research must transition from diagnostic analysis to intervention-based studies. Longitudinal, quasi-experimental research is urgently required to test the causal impact of curriculum changes that leverage the NLP-generated error map (e.g., specialized meN-remediation modules) against traditional curricula.

The BIPA community should establish a National Computational Linguistics Consortium to manage the ongoing maintenance and ethical application of the corpus (Hafidz dkk., 2025). This body must continuously update the NLP pipeline with new learner data to ensure its high F1-score remains relevant and the model adapts to evolving learner profiles and teaching practices.

Technological development should prioritize embedding the NLP pipeline into an Intelligent Tutoring System (ITS). This ITS must go beyond simple error identification to provide immediate, personalized remediation drills and targeted exercises focused on the high-persistence sub-types, enabling effective self-directed learning outside the classroom.

Pedagogical training for BIPA instructors requires immediate reform (Winata dkk., 2025). Training programs must incorporate computational error analysis literacy, teaching instructors how to interpret the granular NLP reports and, critically, how to shift their in-class focus to addressing the semantic and syntactic constraints of affixes, thereby optimizing instructional time.

CONCLUSION

The most critical finding is the computational proof that Affix Misapplication constitutes the most significant and resistant linguistic obstacle for BIPA learners, accounting for 45% of all errors and showing a statistically non-significant reduction rate across all proficiency levels ($p=0.316$). This result is differentiated by the high reliability of the BIPA-Optimized NLP Error Classification Pipeline (F1-score 0.89), which provided the objective evidence needed to confirm that the difficulty is structural and not transient. This validates the finding that simple exposure and curriculum progression are inadequate, compelling a drastic shift toward systematic, focused remediation, particularly targeting the identified high-frequency sub-types of meN- Omission and -kan Overgeneralization.

The primary contribution of this research is the development and validation of a scalable, computationally robust BIPA-Optimized NLP Error Classification Pipeline and the generation of the first statistically generalizable error map for Indonesian learner language. This technological advancement provides BIPA pedagogy with a crucial resource for overcoming the limitations of manual, subjective error analysis. The pipeline transforms instruction by offering an objective, high-fidelity diagnostic standard, enabling curriculum developers to move beyond anecdotal evidence and implement highly targeted, efficient remediation strategies based on persistent, quantified error sub-types.

A key limitation of this study is its reliance on a corpus-based, diagnostic methodology, which confirms the existence and persistence of errors but does not establish the causal effectiveness of any specific remediation strategy. The findings reveal what should be taught, but not how to teach it effectively. Future research must, therefore, transition from diagnostic analysis to intervention-based, quasi-experimental studies that test the causal impact of curriculum changes which leverage the NLP-generated error map (e.g., specialized meN-remediation modules) against traditional curricula. This subsequent research is essential for moving from technological validation to demonstrable pedagogical efficacy.

AUTHOR CONTRIBUTIONS

Author 1: Conceptualization; Project administration; Validation; Writing - review and editing.

Author 2: Conceptualization; Data curation; In-vestigation.

Author 3: Data curation; Investigation.

Author 4: Formal analysis; Methodology; Writing - original draft.

CONFLICTS OF INTEREST

The authors declare no conflict of interest.

REFERENCES

- Amalia, A., Hayati, I., Afandi, A., Lubis, A. S., & Marpaung, J. L. (2025). Comparative Forecasting of Indonesian Stock Prices Using ARIMA and Support Vector Regression: A Statistical Learning Approach. *Mathematical Modelling of Engineering Problems*, 12(7), 2307–2315. Scopus. <https://doi.org/10.18280/mmep.120711>
- Anthonius, F., & Ari, A. (2024). The Implementation of Paperrater and Grammarly in English Teaching: The Implementation of Paperrater and Grammarly in English Teaching to Boost the Writing skills of Non-English Undergraduate Students. *ACM Int. Conf. Proc. Ser.*, 142–146. Scopus. <https://doi.org/10.1145/3678726.3678746>
- Deviantari, U. W., Aditya, T., & Djojomartono, P. N. (2025). The Application of Random Forest Prediction in Developing a Systematic Land Parcel Value in the Urban Area. *International Journal of Geoinformatics*, 21(7), 58–79. Scopus. <https://doi.org/10.52939/ijg.v21i7.4319>
- Ekakristi, A. S., Wicaksono, A. F., & Mahendra, R. (2025). Intermediate-task transfer learning for Indonesian NLP tasks. *Natural Language Processing Journal*, 12. Scopus. <https://doi.org/10.1016/j.nlp.2025.100161>
- Hafidz, I. H., Sulistya, A., & Lidiawaty, B. R. (2025). Sentiment Analysis of Public Complaints: A Machine Learning Comparison of SVM, Naive Bayes, Random Forest, and XGBoost. *ICADEIS - Int. Conf. Adv. Data Sci., E-Learning Inf. Syst.: Integr. Data Sci. Inf. Syst., Proceeding.* Scopus. <https://doi.org/10.1109/ICADEIS65852.2025.10933382>

- Hidayatullah, A. F., Apong, R. A., Lai, D. T. C., & Qazi, A. (2025). Pre-trained language model for code-mixed text in Indonesian, Javanese, and English using transformer. *Social Network Analysis and Mining*, 15(1). Scopus. <https://doi.org/10.1007/s13278-025-01444-9>
- Jazuli, A., & Kusumaningrum, R. (2025). Optimizing Aspect-Based Sentiment Analysis Using BERT for Comprehensive Analysis of Indonesian Student Feedback. *Applied Sciences (Switzerland)*, 15(1). Scopus. <https://doi.org/10.3390/app15010172>
- Jiang, S., Su, Y., Li, X., Lin, N., Xiao, L., & Wang, L. (2025). Unraveling the Efficacy of In-Context Learning in Indonesian Grammatical Error Correction. Dalam W. Shen, W. Shen, M.-H. Abel, N. Matta, J.-P. Barthes, J. Luo, J. Zhang, H. Zhu, & K. Peng (Ed.), *Proc. Int. Conf. Comput. Support. Coop. Work Des., CSCWD* (Nomor 2025, hlm. 1704–1709). Institute of Electrical and Electronics Engineers Inc.; Scopus. <https://doi.org/10.1109/CSCWD64889.2025.11033420>
- Kiatphaisansophon, P., Wanvarie, D., & Cooharajanane, N. (2024). Efficient Text Bounding Box Identification Using Mask R-CNN: Case of Thai Documents. *IEEE Access*, 12, 49306–49328. Scopus. <https://doi.org/10.1109/ACCESS.2024.3383911>
- Kusumoputro, B., Imantaka, S. R., Lina, L., & Kresnaraman, B. (2011). Face recognition system of infra-red images using ensemble back-propagation neural networks. *International Journal of Artificial Intelligence*, 7(11 A), 401–416. Scopus.
- Lefrandt, M., Santoso, E. B., Santoso Gunawan, A. A. S., & Tedjasulaksana, J. J. (2025). Contextual Spelling Corrector for Indonesian Text Preprocessing: A Comparative Analysis of Large Language Models. *Proc. IEEE Int. Conf. Ind. 4.0, Artif. Intell., Commun. Technol., IAICT*, 290–296. Scopus. <https://doi.org/10.1109/IAICT65714.2025.11100636>
- Mantasiah, R., Yusri, Y., & Anwar, M. (2021). Integrating linguistics theories in developing foreign language teaching material (German grammar textbook for indonesian learners).

- International Journal of Language Education*, 5(3), 125–134. Scopus.
<https://doi.org/10.26858/ijole.v5i3.20239>
- Nasution, S., Ferdiana, R., & Hartanto, R. (2025). Towards Two-Step Fine-Tuned Abstractive Summarization for Low-Resource Language Using Transformer T5. *International Journal of Advanced Computer Science and Applications*, 16(2), 1220–1230. Scopus.
<https://doi.org/10.14569/IJACSA.2025.01602120>
- Naufal, T., Mahendra, R., & Wicaksono, A. F. (2025). Sentences, entities, and keyphrases extraction from consumer health forums using multi-task learning. *Journal of Biomedical Semantics*, 16(1). Scopus. <https://doi.org/10.1186/s13326-025-00329-2>
- Navastara, D. A., Akbar Hidiya, F. R., & Wijaya, A. Y. (2023). Prediction of Indonesian Stock Price Using Combination of CNN and BiLSTM Model. Dalam H.-C. Chen, C. Damarjati, C. Blum, Y. Jusman, S. N. A. M. Kanafiah, & W. Ejaz (Ed.), *Proceeding—Int. Conf. Inf. Technol. Comput., ICITCOM* (hlm. 307–312). Institute of Electrical and Electronics Engineers Inc.; Scopus.
<https://doi.org/10.1109/ICITCOM60176.2023.10442941>
- Ningsih, R. Y., Oktriono, K., Wiharja, C. K., & Ernawati, E. (2018). Forms of language errors in speaking practices of foreign students through online UKBIPA application. *ACM Int. Conf. Proc. Ser.*, 59–62. Scopus. <https://doi.org/10.1145/3291078.3291092>
- Pardamean, B., Suparyanto, T., Cenggoro, T. W., Sudigyo, D., & Anugrahana, A. (2022). AI-Based Learning Style Prediction in Online Learning for Primary Education. *IEEE Access*, 10, 35725–35735. Scopus. <https://doi.org/10.1109/ACCESS.2022.3160177>
- Pramuniati, I., & Sitinjak, D. R. (2024). The interlanguage of French learning Indonesian as a foreign language. *Indonesian Journal of Applied Linguistics*, 14(1), 206–219. Scopus.
<https://doi.org/10.17509/ijal.v14i1.70394>
- Priyanto, A., Hapidin, D. A., Edikresnha, D., Aji, M. P., & Khairurrijal, K. (2025). Predicting microplastic quantities in Indonesian provincial rivers using machine learning models.

Science of the Total Environment, 961. Scopus.

<https://doi.org/10.1016/j.scitotenv.2025.178411>

Qomariyah, N. N., Karen, A., Natalie, V., Kazakov, D., & Chaetajaka, P. A. (2025). Utilizing Bidirectional Long Short-Term Memory (BiLSTM) for Radiology Reports in Indonesian Language. *Int. Conf. Knowl. Smart Technol., KST*, 127–132. Scopus.

<https://doi.org/10.1109/KST65016.2025.11003318>

Rabiha, S. G., Yossy, E. H., Indrianti, Y., & Sasmoko, n. (2019). A Neural network based approach for predicting Indonesian teacher engagement index (itei). *Proceeding - Int. Conf. Artif. Intell. Inf. Technol., ICAIIT*, 469–474. Scopus.

<https://doi.org/10.1109/ICAIIIT.2019.8834546>

Rahutomo, F., & Harjito, B. (2025). Machine Learning-Based Climate Prediction in Indonesia: A Baseline Experiment. *International Journal of Advanced Computer Science and Applications*, 16(8), 797–810. Scopus. <https://doi.org/10.14569/IJACSA.2025.0160877>

Ramdani, D., Susilo, H., Suhadi, S., & Sueb, S. (2023). The Effect of Problem Based Learning on Critical Thinking Skills of Biology Learning in Indonesia: A Meta-Analysis Study. Dalam H. Habiddin & N. Farida (Ed.), *AIP Conf. Proc.* (Vol. 2569). American Institute of Physics Inc.; Scopus. <https://doi.org/10.1063/5.0112352>

Salas-Pilco, S. Z., Xiao, K., & Hu, X. (2023). Correction to: Artificial Intelligence and Learning Analytics in Teacher Education: A Systematic Review (Education Sciences, (2022), 12, 8, (569), 10.3390/educsci12080569). *Education Sciences*, 13(9). Scopus. <https://doi.org/10.3390/educsci13090897>

Saputro, B. A., Suryadi, D., Rosjanuardi, R., & Kartasasmita, B. G. (2018). Analysis of students' errors in responding to TIMSS domain algebra problem. *J. Phys. Conf. Ser.*, 1088. Scopus. <https://doi.org/10.1088/1742-6596/1088/1/012031>

- Sari, Y. A., Nakazawa, A., & Wani, Y. A. (2025). LeFood-set: Baseline performance of predicting level of leftovers food dataset in a hospital using MT learning. *PLOS ONE*, 20(5 May). Scopus. <https://doi.org/10.1371/journal.pone.0320426>
- Simanungkalit, E., & Tuga, T. (2024). Data-Driven Insights for Mobile Banking App Improvement: A Sentiment Analysis and Topic Modelling Approach for SimobiPlus User Reviews. *International Journal of Engineering Trends and Technology*, 72(6), 347–360. Scopus. <https://doi.org/10.14445/22315381/IJETT-V72I6P132>
- Soffan, S., Bramantoro, A., & Alzahrani, A. A. (2025). Combination of machine learning and data envelopment analysis to measure the efficiency of the Tax Service Office. *PeerJ Computer Science*, 11. Scopus. <https://doi.org/10.7717/PEERJ-CS.2672>
- Utami, E., Oyong, I., Raharjo, S., Hartanto, A., & Adi, S. (2025). Supervised learning and resampling techniques on DISC personality classification using Twitter information in Bahasa Indonesia. *Applied Computing and Informatics*, 21(1–2), 141–151. Scopus. <https://doi.org/10.1108/ACI-03-2021-0054>
- Utami, V. M., Maylawati, D. S., Syaripudin, U., Zulfikar, W. B., Wahana, A., & Slamet, C. (2025). Text Generation for Low-Calorie Food Recipes Using IndoBART. *Proc. Int. Conf. Wirel. Telemat., ICWT*. Scopus. <https://doi.org/10.1109/ICWT66752.2025.11181956>
- Utomo, B., Soedarto, T., Winarno, S. T., Hendrarini, H., & Farid, I. W. (2024). Monthly Forecasting Indonesian Coffee Production Using Extreme Learning Machine. Dalam F. W. Wibowo (Ed.), *Beyond Technol. Summit Informatics Int. Conf., BTS-I2C* (hlm. 682–686). Institute of Electrical and Electronics Engineers Inc.; Scopus. <https://doi.org/10.1109/BTS-I2C63534.2024.10941753>
- Wijaya, B. C., & Sugiarto, H. S. (2025). Transformer+transformer architecture for image captioning in Indonesian language. *IAES International Journal of Artificial Intelligence*, 14(3), 2338–2346. Scopus. <https://doi.org/10.11591/ijai.v14.i3.pp2338-2346>

Winata, R., Willson, A., Tjen, W., & Madyatmadja, E. D. (2025). Sentiment Analysis on Indonesian Military Law Debate Using Machine Learning and IndoBERT. *Proceeding - Int. Conf. Creat. Commun. Innov. Technol.: Empower. Transform. MATURE LEADERSH.: Harnessing Technol. Adv. Glob. Sustain., ICCIT*. Scopus. <https://doi.org/10.1109/ICCIT65724.2025.11167858>

Yotenka, R., Muhajir, M., Hermansah, n., & Rodrigues, P. C. (2025). Comparative Analysis of Activation Functions in Recurrent Neural Network: An Application to Indonesian Inflation Forecasting. *Mathematical Modelling of Engineering Problems*, 12(3), 754–762. Scopus. <https://doi.org/10.18280/MMEP.120302>

Copyright Holder :

© Inriati Lewa et.al (2025).

First Publication Right :

© Journal International of Lingua and Technology

This article is under:

