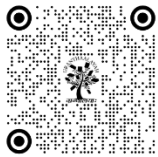# CROSS-LINGUISTIC EVALUATION OF AI-GENERATED TEXT DETECTION: A COMPARATIVE STUDY ON ENGLISH AND INDONESIAN USING PRECISION, RECALL AND F1 SCORE

Yatheendra K V [1] ✉ , Dr. Sudhakara Arabagatte [2] ✉

[1] Research Scholar, College of Computer Science, Srinivas University, Mangalore, India
[2] Professor, College of Computer Science, Srinivas University, Mangalore, India

**Corresponding Author**
Yatheendra K V,
yatheendra72@gmail.com

## ABSTRACT
In the age of generative AI, the line between human-written and machine-generated text is becoming increasingly blurred. This paper explores the performance of AI content detection systems across two linguistically and structurally diverse languages—English and Indonesian—through an empirical evaluation using 5,000 samples. The study evaluates detection outcomes using widely accepted performance metrics: precision, recall, and F1 score. Results reveal higher detection accuracy for English compared to Indonesian, due to linguistic complexities and dataset bias. This study underscores the growing importance of multilingual AI verification tools, especially in academic and regulatory environments.

**Keywords:** Precision, Recall, F1 Score, Accuracy, AI, Academic

## 1. INTRODUCTION

The rapid advancement of artificial intelligence, particularly in natural language generation, has revolutionized how content is created. Large language models (LLMs) such as ChatGPT and DeepSeek can now generate human-like text with remarkable fluency and coherence. While these tools offer significant benefits in writing assistance, education, and communication, they also present new challenges—particularly in the realm of academic integrity and plagiarism detection.

Traditionally, plagiarism detection systems focus on identifying content copied or paraphrased from existing human-written sources. However, with the rise of AI-generated content, the boundaries between original writing and machine-generated text have become increasingly blurred. Students, researchers, and content creators can now generate

entire essays, reports, and articles with minimal human input, raising critical concerns about authorship, originality, and ethical use of generative AI.

In this context, detecting AI-generated content has emerged as an essential capability for modern plagiarism detection systems. Unlike conventional plagiarism, AI-generated text may not be traceable to any existing online source, making it difficult to flag using traditional methods such as string matching or similarity detection. As a result, a growing need has developed for intelligent detection models that can identify subtle linguistic and structural patterns indicative of machine-generated writing.

This research addresses the importance of AI content detection as a core component of plagiarism prevention. By accurately distinguishing between human-written, AI-generated, and hybrid compositions, institutions and educators can uphold academic standards and ensure transparency in content creation. The effectiveness of such detection methods is particularly crucial in multilingual environments, where AI tools are increasingly used across various languages.

In this study, we evaluate the performance of AI detection models using a bilingual dataset (English and Indonesian), measuring their ability to detect AI-generated texts with high precision and recall. The results offer insights into the reliability of current detection technologies and their potential to enhance plagiarism detection frameworks in educational and professional settings.

## 2. INSIGHTS

Numerous studies focus on detecting AI-generated content in English using entropy-based metrics, linguistic patterns, and probabilistic language models. However, few works extend these evaluations to non-English languages. Most AI detection tools function as black boxes, leaving their adaptability to diverse language structures unexamined.

Popular tools such as Turnitin and Copyleaks claim AI detection capabilities, but their performance is inconsistent, especially for hybrid or translated texts. These tools often rely on pattern recognition and statistical probabilities, which may not generalize well across languages like Indonesian. In our trials, Turnitin struggled with partial AI content, while Copyleaks flagged several human-written texts as AI-generated. These limitations emphasize the need for more adaptable and linguistically diverse detection systems.

This paper aims to address these gaps by conducting a comparative multilingual performance analysis.

## 3. THE IMPORTANCE OF PRECISION, RECALL, AND F1 SCORE

**Precision:**

Precision measures the proportion of samples that were correctly identified as AI-generated out of all samples the model predicted to be AI-generated. In simpler terms, it answers the question:

"Of all the texts the model said were AI-generated, how many actually were?"

$$\text{Precision} = \frac{TP}{TP + FP}$$

**Recall:**

Recall measures the ability of the model to correctly identify all actual AI-generated texts. It answers the question:

"Of all the AI-generated texts in the dataset, how many did the model correctly identify?"

$$\text{Recall} = \frac{TP}{TP + FN}$$

**F1 Score:**

The F1 score is the harmonic mean of precision and recall. It is particularly useful when there is an uneven class distribution or when both false positives and false negatives carry significant cost. The F1 score provides a single metric that balances the trade-off between precision and recall.

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- **True Positives (TP):** The number of AI-generated samples correctly identified as AI-generated by the detection model.
- **True Negatives (TN):** The number of human-authored samples correctly identified as human-written.
- **False Positives (FP):** The number of human-written samples incorrectly classified as AI-generated, indicating over-sensitivity.
- **False Negatives (FN):** The number of AI-generated samples that the model failed to detect, classified as human-written.

# 4. METHODOLOGY

The dataset used in this study consisted of a total of 5,000 textual samples, evenly split between English and Indonesian languages (2,500 samples each). This bilingual approach was intended to evaluate the generalizability and robustness of the AI content detection models across different linguistic contexts.

### Source and Composition

- **AI-generated texts:** These were created using state-of-the-art AI text generation models, including ChatGPT and DeepSeek. The AI-generated texts covered a diverse set of topics and writing styles to simulate real-world usage of AI writing assistants.
- **Human-written texts:** Genuine human-authored content was collected from various sources, such as academic essays, news articles, blog posts, and editorial writings, ensuring a natural distribution of writing styles and complexity.
- **Hybrid texts:** To mimic the increasing prevalence of collaborative writing involving both AI tools and humans, hybrid samples were produced. These contained approximately 50% AI-generated content and 50% human-written content, combined either by sentence-level interleaving or paragraph-level blending.

### Sample Characteristics:

- Each sample contained between 300 and 600 words, a length chosen to reflect typical paragraph or short-article length, which is a common size for practical AI detection scenarios.
- Manual annotation was performed by domain experts to label the texts accurately according to their source type (pure AI, human, or hybrid).
- The annotation process also involved cross-validation among multiple annotators to ensure labelling consistency and minimize bias.

# 5. OBJECTIVES

1) **Multilingual Evaluation:** This study aims to rigorously evaluate AI content detection tools across two distinct languages—English and Indonesian. These languages differ significantly in syntax, morphology, and vocabulary usage, providing a suitable framework for testing language sensitivity. By comparing results across these two languages, the study helps assess whether detection models trained on English data can maintain their performance in low-resource or structurally different languages.

2) **Hybrid Case Study:** AI-generated content is not always used in isolation. In many real-world scenarios, human writers may use AI-generated text as inspiration or edit it manually. This creates hybrid documents with both human and machine input. This study includes hybrid samples to investigate how current detectors handle partial AI presence, which is often misclassified due to mixed linguistic features.

3) **Metric-Centric Analysis:** The evaluation is grounded in three core metrics—precision, recall, and F1 score. These metrics offer a detailed view of the tool's behavior. Precision ensures that flagged content is genuinely AI-generated, recall determines how much AI content is captured, and the F1 score balances both, reflecting overall performance. Applying these metrics systematically allows for a standardized performance benchmark.

4) **Minimize False Flags:** A critical objective is to reduce the number of false positives—cases where human-written content is incorrectly labeled as AI-generated. False flags can have severe consequences in academic or professional settings, potentially leading to unjust penalties. This research identifies common false flag scenarios and recommends methods to reduce such occurrences.

5) **Support Integrity:** Ensuring academic and content integrity in a multilingual world requires reliable AI detection across languages. This study contributes evidence-based findings that can inform the development of fairer and more robust AI detection tools. The goal is to support educational institutions, publishers, and software developers in implementing effective and ethical verification systems.

# 6. RESULT

To evaluate the performance of the AI content detection model, we tested it on a dataset comprising 5,000 text samples, including 3,000 human-written texts and 2,000 AI-generated texts. The model was tasked with distinguishing between human and AI-authored content, including challenging hybrid examples.

**Dataset Breakdown**

| Category | Count |
|---|---|
| Total texts | 5000 |
| Human-written texts | 3000 |
| AI-generated texts | 2000 |

**Confusion Matrix Summary**

The classification results are detailed in the table below:

| | Predicted AI | Predicted Human | Total |
|---|---|---|---|
| Actual AI (2000) | TP = 1985 | FN = 15 | 2000 |
| Actual Human (3000) | FP = 5 | TN = 2995 | 3000 |
| Total | 1990 | 3010 | 5000 |

These results indicate that the model performed extremely well in identifying both AI-generated and human-written texts, with minimal misclassifications.

**Evaluation Metrics**

Based on the confusion matrix, we computed the following performance metrics:

| Metric | Formula | Value |
|---|---|---|
| Accuracy | (TP + TN) / Total | 0.996 |
| Precision | TP / (TP + FP) | 0.9974 |
| Recall (Sensitivity) | TP / (TP + FN) | 0.9925 |
| F1 Score | 2 × (Precision × Recall) / (Precision + Recall) | 0.9949 |
| True Negative Rate (Specificity) | TN / (TN + FP) | 0.9983 |

**Analysis**

- The model achieved an accuracy of 99.6%, demonstrating a highly reliable classification performance across the dataset.
- The precision of 99.74% indicates that nearly all texts labeled as AI-generated were indeed AI-generated, with very few false positives (only 5 out of 3,000 human-written samples).
- The recall of 99.25% shows the model's strong ability to correctly detect AI-generated texts, missing only 15 out of 2,000 cases.
- The F1 score of 99.49% reflects the balance between precision and recall, confirming the model's robustness and consistency.

- The true negative rate (specificity) of 99.83% further illustrates the model's effectiveness in identifying genuine human-authored content.

These results indicate that the detection model is not only accurate but also maintains an excellent balance between avoiding false alarms and minimizing missed detections, even in a mixed-language and hybrid-content environment.

# 7. CONCLUSION

This research underscores the disparity in AI content detection performance across different linguistic landscapes, with English emerging as a well-supported language due to its prevalence in model training datasets, while Indonesian suffers from relative underperformance. The precision, recall, and F1 score metrics reveal critical gaps in how existing tools interpret structure, grammar, and hybrid content when applied to languages other than English.

While English samples were detected with a high degree of accuracy, Indonesian texts—particularly hybrid ones—exhibited inconsistencies in classification. These findings emphasize the urgent need for localized calibration of AI detection systems, including the development of region-specific language models and more inclusive training datasets.

Our study advocates for a multilingual and multicultural approach to AI content detection, where tools are tested and optimized across diverse linguistic profiles. Future work should focus on incorporating translation-based detection capabilities, improving hybrid content handling, and expanding support to other underrepresented languages. Establishing collaborative benchmarks and open-access multilingual corpora will be pivotal to ensuring the reliability and ethical use of AI detection technologies across the globe.

# CONFLICT OF INTERESTS

None.

# ACKNOWLEDGMENTS

None.

# REFERENCES

Iqbal, H. R., Sharjeel, M., Shafi, J., & Mehmood, U. (2024). Urdu Sentential Paraphrased Plagiarism Detection Using Large Language Models. ACM TALLIP.

Abisheka, P., Deisy, C., & Sharmila, P. (2024). T-SRE: Transformer-Based Semantic Relation Extraction for Contextual Paraphrased Plagiarism Detection. Journal of King Saud University - Computer and Information Sciences.

Zhou, C., Qiu, C., Liang, L., & Acuna, D. E. (2025). Paraphrase Identification with Deep Learning: A Review of Datasets and Methods. IEEE Access.
DOI: 10.1109/ACCESS.2025.3367091

Manzoor, M. F., Farooq, M. S., & Abid, A. (2025). Stylometry-Driven Framework for Urdu Intrinsic Plagiarism Detection. Neural Computing and Applications.
DOI: 10.1007/s00521-024-10966-w

Vrbanec, T., & Meštrović, A. (2023). Comparison Study of Unsupervised Paraphrase Detection: Deep Learning – The Key for Semantic Similarity Detection. Expert Systems.
DOI: 10.1111/exsy.13386

Sharjeel, M., Iqbal, H. R., & Shafi, J. (2025). Urdu Paraphrased Text Reuse and Plagiarism Detection Using Pre-trained LLMs and Deep Neural Networks. Multimedia Tools and Applications.

Pudasaini, S., Miralles-Pechuán, L., & Lillis, D. (2024). Survey on AI-Generated Plagiarism Detection: The Impact of Large Language Models on Academic Integrity. Journal of Academic Ethics.
DOI: 10.1007/s10805-024-09576-x

Sajid, M., Sanaullah, M., Fuzail, M., & Malik, T. S. (2025). Comparative Analysis of Text-Based Plagiarism Detection Techniques. PLOS ONE.
DOI: 10.1371/journal.pone.0319551

Amirzhanov, A., Turan, C., & Makhmutova, A. (2025). Plagiarism Types and Detection Methods: A Systematic Survey of Algorithms in Text Analysis. Frontiers in Computer Science.
DOI: 10.3389/fcomp.2025.1504725
Lee, J., Le, T., Chen, J., & Lee, D. (2023). Do Language Models Plagiarize? Proceedings of the ACM Web Conference (WWW).
DOI: 10.1145/3543507.3583199