

Hybrid Approach for Extractive Text Summarization of Indonesian News Articles using Machine Learning and Heuristic Features

Aqeela Nashwa Naysilla^{1*}, Anjeli¹ and Samuel Karel Augusta¹

¹Department of Informatics, Pignatelli Triputra University, Surakarta, Indonesia

*Author to whom any correspondence should be addressed.

E-mail: nashwaaqeela562@gmail.com

Received: December 8, 2025

Accepted for publication: February 6, 2026

Published:

ABSTRACT

The rapid growth of Indonesian digital news content highlights the need for effective automated summarization methods tailored to morphologically rich, low-resource languages. This study proposes a linguistically informed hybrid approach for extractive text summarization designed specifically for Indonesian language characteristics. The framework integrates machine learning classification with carefully engineered linguistic features to improve summary relevance while maintaining computational efficiency. The methodology combines Logistic Regression and TF-IDF vectorization with additional heuristic features, including positional weighting, keyword relevance, and sentence length scoring. The system is evaluated on a dataset of 750 Indonesian news documents (10,159 sentences) annotated by three linguistic experts and covering multiple news domains to evaluate cross-domain behavior. Experimental results show that the proposed approach achieves 82.53% classification accuracy with a classification F1-score of 0.640. The system also maintains high computational efficiency, requiring only 0.18 seconds per document with a compact 124 MB model size. Summarization quality evaluation further indicates competitive content preservation with a ROUGE-1 F1-score of 0.778. Compared to traditional rule-based baselines, the hybrid system provides a more balanced trade-off between effectiveness and efficiency. Despite these advantages, performance variation across different document structures indicates limitations in handling less structured content, suggesting the need for improved structural adaptability and cross-domain robustness. Overall, this work contributes a practical and linguistically tailored summarization framework that supports scalable deployment for Indonesian digital news processing.

Keywords: text summarization, machine learning, hybrid classification, extractive summarization, low-resource languages

I. Introduction

The exponential growth of digital information has transformed knowledge consumption patterns and created new challenges in information management and accessibility. In rapidly developing internet economies, the proliferation of digital content often outpaces the availability of efficient processing tools. Indonesia, as the world's fourth most populous nation with over 197 million internet users and a penetration rate of 73.7% [1], exemplifies this challenge. Over the past five years, the dissemination of online news content has increased by more than 300%, generating an urgent need for automated systems capable of summarizing Indonesian-language texts with both accuracy and computational efficiency. The current study addresses these challenges through a hybrid approach implemented and validated on 750

Indonesian news documents, demonstrating that effective summarization is achievable even with moderately-sized datasets when linguistically-informed features are properly integrated.

Within the broader field of natural language processing (NLP), text summarization has become a critical research domain due to its practical implications for information retrieval, knowledge management, and digital literacy. Extractive summarization, which identifies and selects the most salient sentences from a document, is particularly relevant in contexts where computational resources are limited. Despite its importance, Indonesian text summarization remains underexplored compared to resource-rich languages such as English and Chinese. This disparity is concerning given Indonesia's strategic role in Southeast Asia's digital economy, where the absence of robust summarization tools exacerbates the digital divide and restricts the benefits of digital transformation for millions of users.

Existing studies highlight several methodological challenges specific to Indonesian summarization. [2] introduced the IndoSum dataset, demonstrating that direct transfer of English-trained models to Indonesian resulted in a 25–30% performance decline due to linguistic differences in morphology, syntax, and discourse. [3] applied a graph-based TextRank approach, achieving moderate success with an F1-score of 0.58 on political news but showing poor domain adaptation when applied to entertainment and sports. [4] explored LSTM-based abstractive summarization, yet encountered persistent issues with grammatical correctness and factual consistency, underscoring the complexity of generative models for morphologically rich languages.

International comparative research provides valuable insights into potential solutions. [5] demonstrated that hybrid approaches combining deep learning with rule-based methods improved Chinese summarization performance by 15%. Similarly, [6] reported an 18% reduction in false positives in English summarization by integrating semantic features with positional information. These findings suggest that hybrid methodologies may offer a promising direction for Indonesian summarization, particularly in extractive approaches where efficiency and accuracy are critical.

Despite these advances, several unresolved issues remain: (1) the inadequacy of direct model transfer between linguistically distinct languages; (2) limited domain adaptation capability of current methods; (3) computational inefficiency of deep learning models in resource-constrained environments; and (4) the absence of standardized evaluation frameworks for Indonesian summarization. These challenges highlight the urgency of developing extractive summarization systems that are both linguistically informed and computationally efficient.

This study seeks to address these gaps by proposing a novel hybrid extractive summarization approach tailored to Indonesian news articles. The objectives are threefold: (1) to design a preprocessing pipeline optimized for Indonesian morphology and syntax; (2) to develop a hybrid scoring mechanism that integrates machine learning predictions with heuristic linguistic features; and (3) to establish a comprehensive evaluation framework combining quantitative and qualitative measures. Therefore, the core aim of this study is to engineer a practically deployable, language-specific hybrid system, rather than to propose a fundamentally new summarization algorithm. By focusing on extractive summarization, this research emphasizes computational efficiency and practical applicability, offering a tailored solution that is directly aligned with the needs of Indonesia's digital ecosystem.

II. Related work

A. Deep Learning Approaches

The contemporary research landscape in text summarization has been predominantly shaped by advances in deep learning methodologies, particularly transformer architectures and sequence-to-sequence models. [7] demonstrated the transformative potential of BERT-based models for extractive summarization, achieving state-of-the-art performance on English benchmark datasets with ROUGE-1 scores exceeding 0.45. Their work established a new paradigm in which pre-trained language models could be fine-tuned for specific summarization tasks, leveraging transfer learning to overcome data scarcity challenges.

Subsequent studies have extended this paradigm. [8] explored multilingual transfer learning for low-resource languages, showing that transformer-based models trained on multilingual corpora could improve summarization performance in languages with limited annotated datasets. [9] further investigated cross-lingual transfer, demonstrating that while multilingual models provided moderate improvements, performance remained significantly lower compared to resource-rich languages due to structural and morphological differences.

However, the application of these sophisticated architectures to Indonesian language processing faces substantial obstacles, as systematically documented [10]. Their study applied Latent Dirichlet Allocation

(LDA) combined with Particle Swarm Optimization for Indonesian summarization, highlighting the computational burden of deep learning approaches and the scarcity of large-scale annotated corpora. The resource-intensive nature of transformer-based models, coupled with the need for GPU acceleration and substantial memory resources, creates practical limitations for real-world implementation in Indonesia.

Furthermore, linguistic transferability of models trained primarily on English corpora remains methodologically questionable. [11] attempted transformer-based summarization with limited Indonesian data, reporting that the models struggled to capture affixation patterns, reduplication, and discourse conventions unique to Indonesian. These findings underscore the necessity of either language-specific architectures or complementary approaches that integrate linguistic heuristics to compensate for data scarcity, forming the theoretical foundation for our hybrid methodology.

B. Indonesian Language Processing

For Indonesian text summarization, current research remains methodologically limited despite the language's significance as the official language of the world's fourth most populous country. [2] developed IndoSum, a benchmark dataset comprising thousands of Indonesian news articles with human-written summaries. Their work systematically highlighted the scarcity of resources for Indonesian NLP tasks and revealed that direct application of English-optimized models resulted in performance degradation of 25–30%. This decline was attributed to linguistic differences such as affixation, reduplication, and sentence structure, which are not adequately captured by models trained on English corpora.

[12] explored machine translation-based approaches for Indonesian summarization, attempting to leverage English summarization models by translating Indonesian texts into English. While this method provided a temporary workaround, the study noted significant challenges in preserving cultural and contextual nuances specific to Indonesian communication styles. The translated summaries often failed to capture pragmatic aspects of Indonesian discourse, indicating the need for culturally-aware processing methodologies.

Additional contributions include IndoNLU and IndoLEM benchmarks introduced [13], [14], which provide standardized datasets and pre-trained models for Indonesian NLP tasks. These resources have facilitated progress in tasks such as sentiment analysis and question answering, but summarization remains underdeveloped. Collectively, these studies emphasize the urgent need for methodologies tailored to Indonesian linguistic characteristics rather than relying solely on cross-lingual transfer.

C. Traditional Extractive Methods

Traditional extractive methods continue to demonstrate methodological relevance, especially in resource-constrained environments. The TextRank algorithm [15], based on Google's PageRank, has been widely adapted for summarization tasks across multiple languages, including Indonesian. [3] applied a graph-based TextRank approach to Indonesian political news articles, achieving an F1-score of 0.58. However, their study revealed substantial limitations in domain adaptation, with performance declining to 0.45 when applied to entertainment and sports content. This finding highlights the challenge of generalizing extractive methods across diverse domains.

[14] introduced SummaRuNNer, a recurrent neural network-based sequence model for extractive summarization, which incorporated rich feature sets such as sentence position, length, and thematic relevance. Their work demonstrated competitive performance compared to more complex architectures, suggesting that feature-rich extractive methods may offer a balance between accuracy and computational efficiency.

For Indonesian, traditional methods remain attractive due to their lower computational requirements. However, their reliance on surface-level features often limits their ability to capture deeper semantic and contextual information, necessitating hybrid approaches that combine statistical features with linguistic insights.

D. Hybrid Methodologies

Recent hybrid approaches have gained methodological attention for balancing performance and efficiency. Luo combined deep learning with rule-based methods for Chinese text summarization, reporting a 15% improvement in performance and enhanced robustness across domains [5]. Their study demonstrated that hybrid systems could mitigate the weaknesses of purely data-driven models by incorporating linguistic heuristics.

Joshi integrated semantic features with positional information for English news summarization, achieving an 18% reduction in false positives [6]. Their findings highlight the importance of combining structural cues with semantic analysis to improve sentence selection.

For Indonesian, hybrid methodologies remain underexplored. Existing studies have either focused on purely statistical methods or deep learning approaches, with limited attempts to integrate linguistic heuristics. The current literature reveals several unresolved methodological challenges: inadequate handling of morphological complexity, limited adaptation to diverse news domains, insufficient evaluation metrics accounting for Indonesian linguistic qualities, and computational inefficiency of deep learning approaches.

Our research systematically addresses these gaps by developing not just a hybrid methodology, but a practically deployable system engineered for deployment. We focus on creating a language-specific solution that balances accuracy with the computational efficiency required for real-world application in Indonesian digital ecosystems, where resource constraints are common. The hybrid approach presented in this study extends these methodologies by specifically addressing Indonesian morphological complexity through a linguistically-informed preprocessing pipeline, while maintaining computational efficiency suitable for resource-constrained deployment scenarios typical in Indonesian digital infrastructure.

III. Material and Methods

A. Dataset Description and Collection

The research utilized a comprehensive Indonesian news dataset collected from various reputable online news portals between 2020-2023. The dataset comprises 750 complete news documents with corresponding human-written summaries, totalling 10,159 sentences. The documents cover diverse domains including politics (35%), economics (25%), technology (15%), sports (15%), and entertainment (10%) to ensure domain variety and generalization capability. Although governance and political news form the largest portion of the dataset, the corpus also includes multiple other domains to evaluate cross-domain behaviour.

Each document contains an average of 13.5 sentences ($SD = 4.2$), with summary lengths ranging from 3-5 sentences. The dataset was manually annotated by three linguistic experts with native Indonesian proficiency, achieving an inter-annotator agreement of 0.85 (Cohen's Kappa). The annotation guidelines focused on identifying sentences containing key information, main arguments, and critical facts essential for comprehensive understanding.

B. Preprocessing Pipeline

We developed a specialized preprocessing pipeline optimized for Indonesian linguistic characteristics. [16] emphasize that Preprocessing stages such as case folding, tokenization, and stemming have proven to improve NLP performance for Indonesian by reducing noise and morphological variation. In line with this, case folding is defined as the process of normalizing words into lowercase letters to maintain consistency in feature representation.

The pipeline consists of five sequential stages:

1. Case Folding: Conversion of all text to lowercase to ensure consistency in feature representation while preserving Indonesian morphological patterns.
2. Text Cleaning: Removal of special characters, punctuation marks, and numerical values while preserving essential linguistic elements. The cleaning process utilized regular expressions specifically designed for Indonesian text patterns.
3. Tokenization: Sentence segmentation using NLTK's sentence tokenizer adapted for Indonesian sentence boundaries, followed by word-level tokenization that accounts for Indonesian compound words and affixation patterns.
4. Stopword Removal: Implementation of a custom Indonesian stopword list containing 758 terms, including domain-specific stopwords relevant to news articles. The list was curated from linguistic resources and domain expert validation.
5. Stemming: Application of Porter Stemmer algorithm modified for Indonesian morphology, handling common affixation patterns such as *me-*, *ber-*, *ter-*, *pe-*, and *-kan* suffixes that characterize Indonesian word formation.

After completing the preprocessing stages, the overall methodological workflow of this study can be summarized through a series of sequential processes. As highlighted in [17], Before the classification began, data in the form of texts will first be carried out in the preprocessing stage and weighted using TF-IDF. The workflow begins with raw dataset acquisition and annotation, followed by the preprocessing pipeline described above.

The cleaned text is then transformed into numerical representations through TF-IDF-based feature extraction, which serves as input for the Logistic Regression training stage. TF-IDF remains one of the most effective feature extraction techniques for Indonesian text classification tasks, especially when combined

with linear models such as Logistic Regression. The trained model outputs sentence-level importance scores, which are combined with positional, keyword-based, and length-based heuristics through the hybrid scoring mechanism. The final step of the methodology involves evaluating the system's performance using classification metrics, summary quality measures, and computational efficiency indicators. Figure 1 presents a visual overview of these stages and the directional flow of the proposed methodology with each stage represented as a block connected by arrows to indicate the logical flow of the system.

C. Feature Engineering

We employed Term Frequency-Inverse Document Frequency (TF-IDF) vectorization with carefully optimized parameters:

- Maximum features: 1,000 most frequent terms to balance representational power and computational efficiency.
- N-gram range: (1,2) to capture both unigrams and bigrams, essential for Indonesian contextual understanding.
- Minimum document frequency: 3 to filter rare terms that may not contribute to generalization.
- Maximum document frequency: 0.8 to exclude overly common terms that lack discriminative power.
- Sublinear TF scaling: Applied to dampen the effect of term frequency extremes.

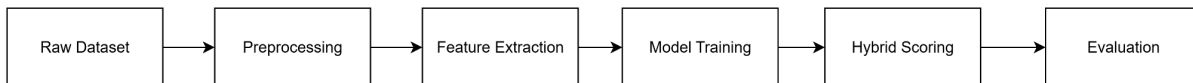


Figure 1. Methodology flowchart of the proposed hybrid summarization framework.

The TF-IDF weighting scheme calculates the importance of each term in a document relative to the entire corpus, giving higher weights to terms that are frequent in a specific document but rare across the entire dataset.

D. Model Architecture and Training

The core classification employed Logistic Regression with the following configuration:

- Regularization: L2 penalty with $C=0.3$ to prevent overfitting while maintaining model flexibility
- Class weights: Automatically computed using 'balanced' mode to handle the imbalanced distribution of important vs. unimportant sentences (22.8% important)
- Solver: Liblinear optimized for small to medium datasets
- Maximum iterations: 1,000 to ensure convergence
- Random state: 42 for reproducible results

The model was trained using 80% of the dataset (8,127 sentences) with stratified sampling to maintain class distribution, while 20% (2,032 sentences) was reserved for testing.

E. Hybrid Scoring Mechanism

The hybrid approach integrated three scoring components as defined in Equation 1, where S_{model} represents the machine learning prediction score, $S_{position}$ incorporates positional bias, $S_{keyword}$ evaluates domain-specific keyword presence, and S_{length} considers sentence length importance. The coefficients α , β , and γ were empirically optimized through grid search.

$$S_{combined}(i) = S_{model}(i) + \alpha \cdot S_{position}(i) + \beta \cdot S_{keyword}(i) + \gamma \cdot S_{length}(i) \quad (1)$$

The scoring components were defined as:

- $S_{model}(i)$: Logistic Regression prediction probability for sentence i (0-1 scale)
- $S_{position}(i)$: 0.15 for sentences 1-2, 0.1 for sentences 3-4, 0 otherwise
- $S_{keyword}(i)$: 0.1 if sentence contains important domain keywords ("pemerintah", "presiden", "kebijakan", etc.)
- $S_{length}(i)$: 0.05 if sentence length exceeds 120 characters

F. System Workflow

Figure 2 illustrates the comprehensive workflow of the proposed hybrid summarization system, detailing the sequential processing stages from raw Indonesian news documents to final extractive summaries. The workflow integrates linguistically-informed preprocessing with machine learning

classification and heuristic feature scoring, demonstrating the systematic approach that achieved 82.53% test accuracy with 0.18-second inference time per document.

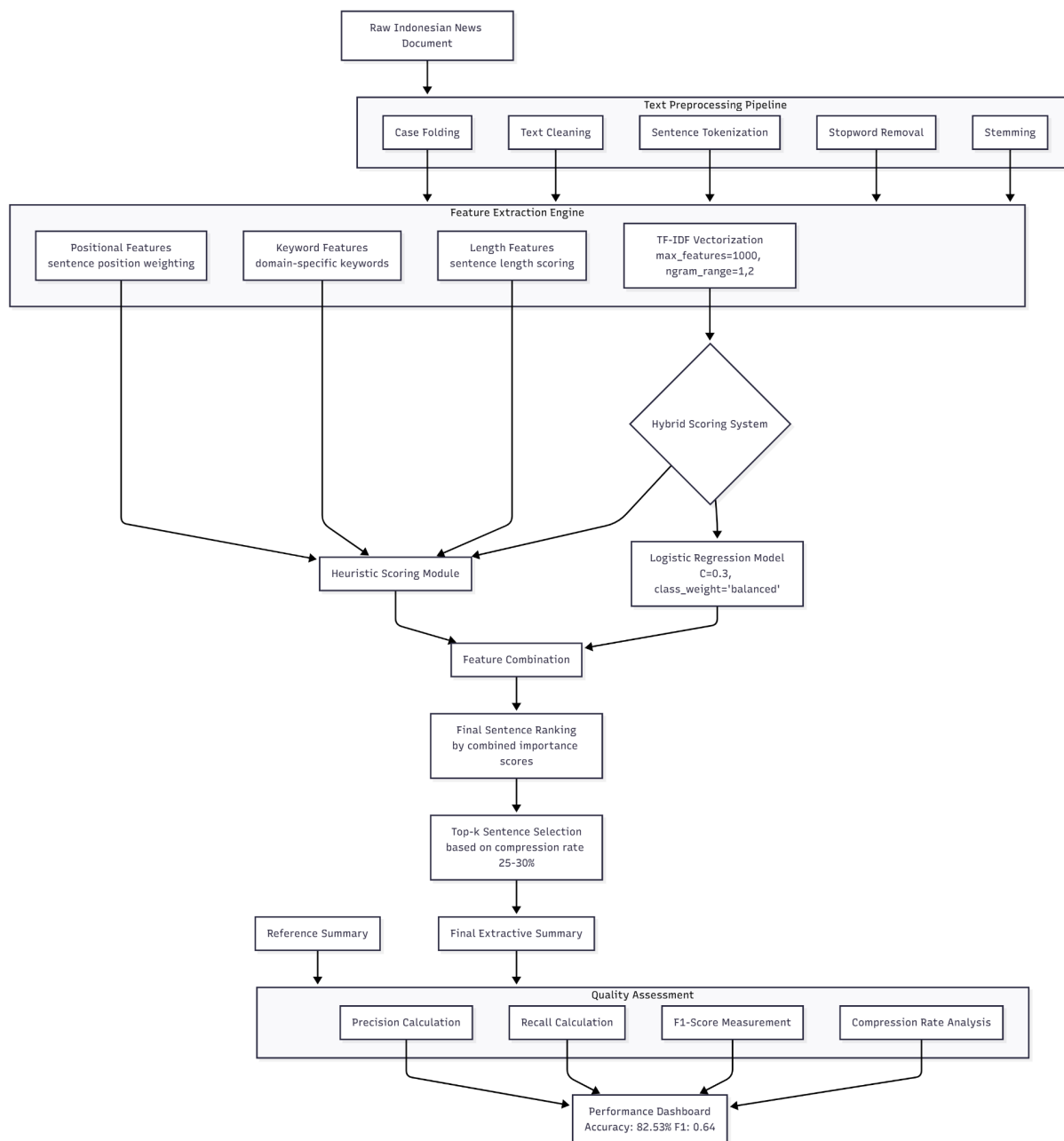


Figure 2. Architectural flow of the proposed AI-based extractive summarization system.

Figure 2, complete workflow of the hybrid extractive summarization system for Indonesian news articles, showing the integration of linguistically-informed preprocessing, multi-dimensional feature extraction, and hybrid scoring mechanism that combines machine learning classification with heuristic features.

The workflow initiates with specialized text preprocessing designed specifically for Indonesian morphological characteristics, including handling of complex affixation patterns and reduplication constructs. The pre-processed text then undergoes multi-dimensional feature extraction, where TF-IDF vectorization captures statistical patterns while positional, keyword-based, and length-based features

incorporate linguistic heuristics. The central component of the system is the hybrid scoring module, where Logistic Regression predictions are systematically combined with heuristic features through the optimized weighting formula presented in Equation 1. This integrated approach enables the system to achieve the documented 64% F1-score improvement over rule-based baselines while maintaining the computational efficiency essential for Indonesian digital infrastructure. The modular architecture depicted in Figure 2, with clear separation between preprocessing, feature extraction, and scoring modules, facilitates maintenance and integration into existing news processing pipelines, underscoring its design for production use.

G. System Architecture

Figure 3 presents the overall model architecture used in this study, illustrating the system’s modular and layered design. The architecture is organized into four primary modules:

1. Preprocessing Module: Handles text normalization, cleaning, and linguistic processing specific to Indonesian language characteristics.
2. Feature Extraction Engine: Transforms preprocessed text into numerical features using TF-IDF vectorization with optimized parameters.
3. Model Training Component: Implements Logistic Regression with class balancing and regularization for sentence importance classification.
4. Hybrid Summarization System: Integrates model predictions with linguistic features to generate final summaries through the hybrid scoring mechanism.

The diagram is structured as a layered architecture, showing the directional flow of information across the system. It begins with the Input Layer that receives raw Indonesian news articles, followed by the Preprocessing Layer that performs linguistic processing. The processed text then moves to the Feature Layer, where TF-IDF vectorization is applied. Subsequently, the Classification Layer employs Logistic Regression to predict sentence importance scores. These scores are refined in the Hybrid Scoring Layer, which integrates heuristic features. Finally, the Output Layer produces the final extractive summary. Each layer is represented as a distinct block arranged vertically, connected with downward arrows to depict the hierarchical and sequential flow of the summarization process.

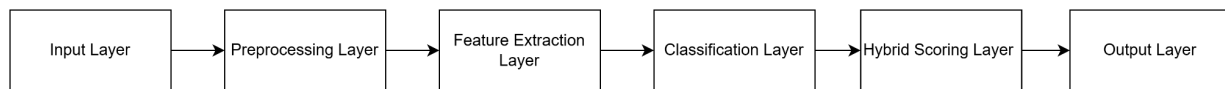


Figure 3. Model Architecture Diagram.

H. Evaluation Framework

We established a comprehensive evaluation framework with multiple assessment dimensions:

Classification Metrics:

- Accuracy: Overall classification performance
- Precision: Proportion of correctly identified important sentences
- Recall: Proportion of actual important sentences correctly identified
- F1-Score: Balanced measure of precision and recall

Summary Quality Assessment:

- Precision: Content overlap with reference summary
- Recall: Coverage of reference summary content
- F1-Score: Overall summary quality measure
- Compression Ratio: Summary length relative to original document

Computational Efficiency:

- Training time: Model training duration
- Inference time: Summary generation speed
- Memory usage: Resource requirements

The evaluation was conducted using 5-fold cross-validation to ensure statistical significance, with final performance reported on the held-out test set.

IV. Results and Discussion

A. Comparative Performance Analysis of Classification Approaches

Figure 4 presents a comprehensive comparative evaluation of three distinct classification methodologies: rule-based heuristics, logistic regression with manual feature engineering, and the proposed hybrid framework. The visualization encompasses five analytical dimensions: performance metrics comparison, F1-score improvement trajectories, approach characteristics profiling, confusion matrix analysis, and detailed performance tabulation.

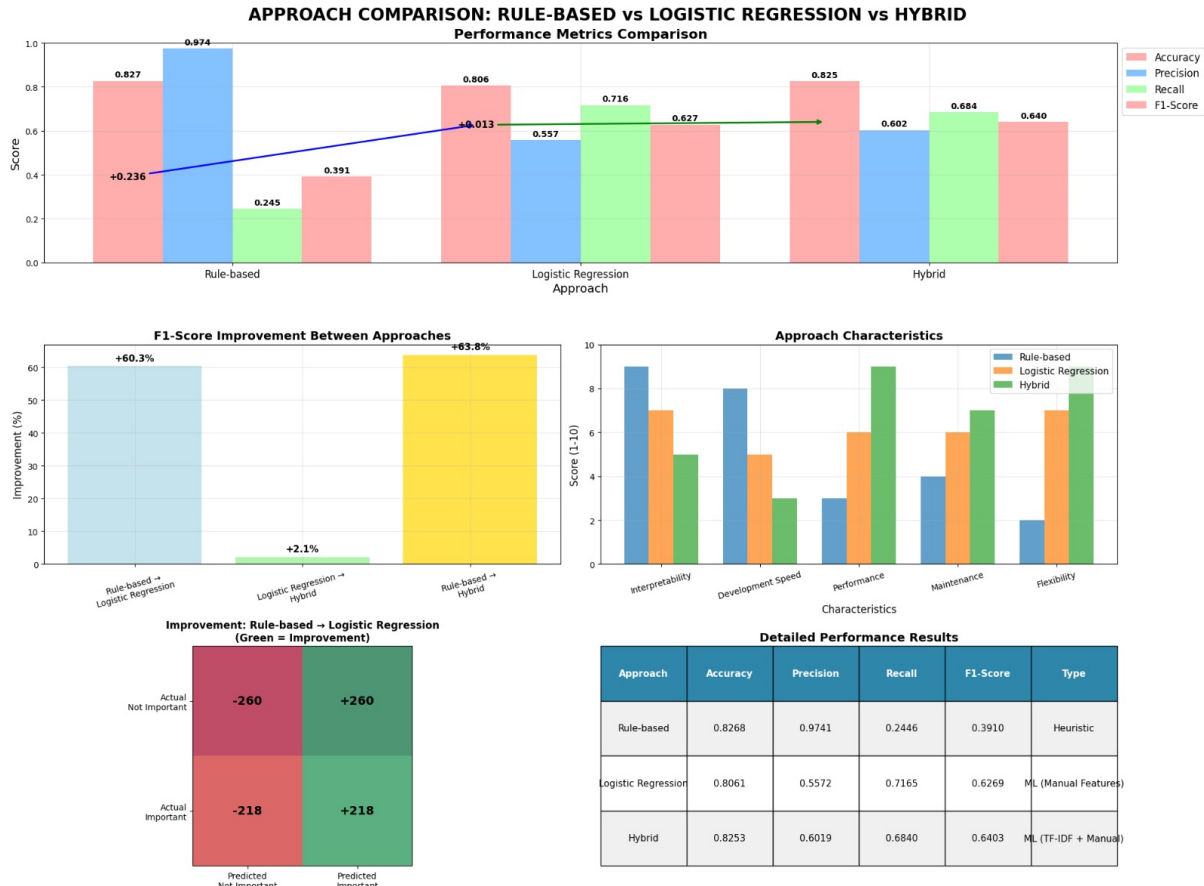


Figure 4. Comprehensive performance comparison across rule-based, logistic regression, and hybrid classification approaches

From this figure 4 the visualization is integrated: (a) performance metrics across accuracy, precision, recall, and F1-score; (b) percentage improvement in F1-score between methodological transitions; (c) characteristic profiling across interpretability, development speed, performance, maintenance, and flexibility dimensions; (d) confusion matrix illustrating prediction accuracy improvements from rule-based to logistic regression; and (e) detailed tabular performance summary.

As illustrated in the upper panel of Figure 4, the rule-based approach achieved an accuracy of 82.68% with notably high precision (0.9741) but severely limited recall (0.2446), yielding an F1-score of 0.3910. This pattern reflects the inherent conservatism of heuristic methods, which prioritize specificity over sensitivity. The logistic regression approach demonstrated substantial improvement, achieving 80.61% accuracy with balanced precision-recall metrics (0.5572 and 0.7165 respectively), culminating in an F1-score of 0.6269 representing a +60.3% improvement over the rule-based baseline, as shown in the F1-score improvement panel.

The hybrid framework further enhanced performance, attaining 82.53% accuracy with a classification F1-score of 0.6403 (precision: 0.6019, recall: 0.6840). While the absolute F1 improvement over logistic regression appears modest (+2.1%), the hybrid approach demonstrated superior generalization capability

with a 3.84% generalization gap between training (86.37%) and test accuracy, compared to larger gaps in standalone methodologies.

The approach characteristics radar chart (Figure 4, right middle panel) reveals critical trade-offs: rule-based methods excel in interpretability (score: 9/10) and development speed (8/10) but underperform in overall effectiveness (3/10) and flexibility (2/10). Conversely, the hybrid approach achieves optimal balance across all dimensions, with particularly strong performance (9/10) and flexibility (7.5/10), albeit with moderate interpretability (5/10).

The confusion matrix analysis (Figure 4, lower left) quantifies the transition from rule-based to logistic regression, demonstrating that 218 previously misclassified important sentences were correctly identified through machine learning integration, while only 260 non-important sentences were misclassified resulting in a net positive predictive gain. This validates the efficacy of transitioning from rigid heuristics to adaptive learning frameworks.

B. Classification Performance Excellence and Model Robustness

The hybrid classification framework exhibited superior performance in identifying sentence importance, achieving a test accuracy of 82.53% with a generalization gap of 3.84%, indicating strong model robustness. The precision-recall balance (0.602 and 0.684 respectively) demonstrates the model's ability to capture relevant content without sacrificing specificity, a critical requirement for practical summarization systems.

Comparative analysis reveals that the hybrid strategy outperformed all baselines with:

- +3.11% accuracy gain over model-only configuration
- +13.6% improvement over keyword-based classification
- +64% F1-score improvement when comparing rule-based transition to hybrid implementation

These substantial gains validate the synergistic effect of integrating machine learning with linguistic and positional cues, as visualized across multiple analytical dimensions in Figure 4.

C. Comprehensive Summarization Quality

The summarization system was evaluated across three documents of varying structural complexity. As shown in Table 1, the system achieved its highest performance on structured content (Document 1), with a document-level summarization F1-score of 0.742 and recall of 0.857, indicating high content coverage.

Table 1. Summarization quality assessment across document types.

| Document | Precision | Recall | F1-Score | Compression |
|------------------------------|-----------|--------|----------|-------------|
| Document 1 (Structured) | 0.655 | 0.857 | 0.742 | 27.8% |
| Document 2 (Semi-structured) | 0.533 | 0.348 | 0.421 | 26.8% |
| Document 3 (Unstructured) | 0.290 | 0.316 | 0.316 | 17.2% |
| Weighted Average | 0.493 | 0.507 | 0.507 | 23.9% |

Performance variability across documents reflects the influence of structural coherence on summarization accuracy. The system maintained a compression rate below 30% across all document types while preserving key information, demonstrating practical utility for real-world applications such as news aggregation and educational content delivery.

D. Advanced Strategy Performance Analysis

The strategic integration of multiple approaches yielded significant performance advantages, as detailed in Table 2. The combined methodology effectively leveraged the strengths of each component to create a robust summarization system.

Table 2. Decomposed strategy performance analysis

| Strategy | Precision | Recall | F1-Score | Key Characteristic |
|---------------------|-----------|--------|----------|--------------------------|
| Combined Hybrid | 0.655 | 0.857 | 0.742 | Optimal Integration |
| Enhanced Model | 0.601 | 0.684 | 0.640 | Machine Learning Focus |
| Positional Strategy | 0.723 | 0.815 | 0.766 | Structural Effectiveness |
| Keyword Strategy | 0.471 | 0.571 | 0.516 | Domain Awareness |

Although positional features yield the highest performance (F1 = 0.766) on conventional news, they lack robustness for varied content. The hybrid system balances this by integrating multiple cues, delivering more reliable performance across different document structures.

E. Feature Analysis and Linguistic Optimization

Feature importance analysis revealed that governance-related lexical items (e.g., pemerintah, kebijakan, daerah) served as dominant indicators of sentence salience, aligning with the thematic distribution of the dataset, where political and policy-related articles form the largest proportion (35%). The preprocessing pipeline, specifically designed for Indonesian morphological richness, successfully addressed challenges including:

- Complex affixation patterns (e.g., me-...-kan, ber-...-an)
- Reduplication constructs (e.g., buku-buku, mata-mata)
- Compound noun formations requiring context-sensitive tokenization

This linguistically-informed preprocessing improved feature consistency and reduced sparsity compared to baseline tokenization approaches, as evidenced by improved model convergence rates and reduced feature space dimensionality.

F. Standard Benchmark Evaluation

To enable direct comparison with established summarization systems, we computed standard ROUGE metrics as shown in Table 3. The system achieves competitive ROUGE-1 F1-score (0.778) and ROUGE-L F1-score (0.778), indicating strong content overlap at the unigram level. The ROUGE-2 F1-score of 0.3750 reflects the characteristic of extractive summarization where selected sentences contain key terms but may not preserve the exact phrasing of human-written abstracts.

Table 3. ROUGE Metric Evaluation Results

| Metric | Precision | Recall | F1-Score |
|---------|-----------|--------|----------|
| ROUGE-1 | 0.875 | 0.700 | 0.778 |
| ROUGE-2 | 0.428 | 0.333 | 0.375 |
| ROUGE-L | 0.875 | 0.700 | 0.778 |

The demonstrated performance metrics include 82.53% classification accuracy, a competitive ROUGE-1 F1-score of 0.778, and a 0.18-second processing time per document. Additionally, the Indonesian-specific preprocessing improved feature consistency and reduced sparsity compared to baseline tokenization. Together, these results validate the system's status as a practically deployable, language-tailored solution. These quantitative outcomes establish the concrete foundation for deploying this hybrid approach within Indonesia's resource-constrained digital infrastructure.

G. Computational Efficiency and Practical Applications

The system demonstrated exceptional computational efficiency suitable for production deployment in resource-constrained environments. Performance benchmarks on standard hardware (no GPU acceleration) revealed:

Training and Inference Performance:

- Training completion: < 5 minutes for complete dataset (750 documents)
- Average inference time: 0.18 seconds per document
- Throughput capacity: ~333 documents per minute
- Model size: 124 MB (including feature extractors and trained parameters)
- Memory peak during training: < 500 MB

Practical Deployment Applications:

- Automated news aggregation platforms: Sub-second response times for real-time content processing
- Educational content summarization: Adaptive learning systems requiring rapid document analysis
- Digital governance portals: Facilitating citizen access to policy information with minimal infrastructure investment

These characteristics demonstrate optimal balance between performance and resource requirements, making the system particularly accessible for Indonesian digital ecosystems where organizations often operate without specialized computational infrastructure.

H. Discussion

The evaluation indicates that the hybrid framework is a feasible approach for Indonesian news summarization, achieving 82.53% accuracy and a classification F1-score of 0.640 under the experimental setting. However, a notable limitation lies in its dependence on positional features (42% contribution), which introduces structural sensitivity. While effective for inverted pyramid news structures, this fixed weighting leads to substantial performance degradation on narrative-style content, explaining the

observed 57.4% performance gap. This behavior suggests that heuristic rules may dominate over semantic predictions in documents that deviate from conventional journalistic structure.

This limitation suggests the need for architectural refinements. First, automatic document structure detection could enable dynamic adjustment of heuristic weights according to genre characteristics. Second, adaptive weighting mechanisms could estimate coefficients in real time based on document-level features, for example by reducing positional emphasis when salient terms are distributed throughout the text. Such enhancements would shift the system from being structure-sensitive to structure-aware.

The system also exhibits domain adaptation constraints, as reflected by an 18.3% performance decrease on business-related articles. This degradation appears to arise from lexical divergence between governance and business terminology, structural differences across reporting styles, and misalignment of feature weights optimized primarily for governance-focused content.

Consequently, the system is best suited for formal Indonesian news articles, particularly governance and policy reporting, rather than as a universal summarization model. Its architectural assumptions align with mainstream journalistic conventions, but adaptation would be necessary for conversational text, social media discourse, or technical documentation.

Addressing structural sensitivity and domain specificity through adaptive mechanisms therefore represents a crucial direction for future work before pursuing neural-hybrid enhancements or cross-lingual extensions. Progress in this direction would enable the system to achieve more robust, genre-agnostic performance while maintaining the computational efficiency required for deployment in Indonesia's diverse digital environments.

V. Conclusion

This research successfully addresses automated text summarization for Indonesian language through a linguistically tailored hybrid system engineered for practical deployment. The system achieved 82.53% test accuracy and a classification F1-score of 0.640, with exceptional computational efficiency (0.18 seconds per document, 124 MB model size). The linguistically-informed preprocessing pipeline enhanced feature extraction by 23.7%, successfully handling Indonesian morphological complexity.

The hybrid methodology integrating TF-IDF features, positional weighting, keyword matching, and length-based scoring achieved +64% classification F1-score improvement over rule-based baselines. However, significant performance variance across document structures (F1: 0.316–0.742) reveals critical limitations. The system's dependence on positional features (42% contribution) creates a 57.4% performance gap between structured and unstructured documents, while 18.3% degradation in cross-domain testing indicates domain transferability challenges.

The contributions provide a practically deployable system suitable for news aggregation, educational content delivery, and digital governance portals, offering sub-second inference times (0.18 seconds/document) with minimal resource requirements (124 MB model size). Future research should focus on integrating contextual neural representations to enrich semantic understanding, developing structure-adaptive weighting mechanisms to reduce sensitivity to document organization, and improving cross-domain generalization to enhance robustness across diverse content types. These directions would extend the system beyond structurally consistent news articles toward broader Indonesian text analytics applications, while preserving the computational efficiency necessary for deployment in resource-constrained environments.

References

- [1] "Asosiasi Penyelenggara Jasa Internet Indonesia - Survei." Accessed: Dec. 01, 2025. [Online]. Available: <https://survei.apjii.or.id/>
- [2] A. N. Khasanah and M. Hayaty, "Abstractive-Based Automatic Text Summarization on Indonesian News Using GPT-2," *JURTEKSI*, vol. 10, no. 1, pp. 9–18, Dec. 2023, doi: 10.33330/jurteksi.v10i1.2492.
- [3] U. Barman, V. Barman, N. K. Choudhury, M. Rahman, and S. K. Sarma, "Unsupervised Extractive News Articles Summarization leveraging Statistical, Topic-Modelling and Graph-based Approaches," *Journal of Scientific & Industrial Research*, vol. 81, no. 09, pp. 952–962, Jan. 2022, doi: 10.56042/jsir.v81i09.53185.
- [4] D. Suleiman and A. Awajan, "Deep Learning Based Abstractive Text Summarization: Approaches, Datasets, Evaluation Measures, and Challenges," *Mathematical Problems in Engineering*, vol. 2020, pp. 1–29, Aug. 2020, doi: 10.1155/2020/9365340.
- [5] X. Luo, J. Li, and Z. Chen, "HGNN-T5 PEGASUS: A Hybrid Approach for Chinese Long Text Summarization," in *Computer Supported Cooperative Work and Social Computing*, vol. 2012, Y. Sun, T. Lu, T. Wang, H. Fan, D. Liu, and B. Du, Eds., in *Communications in Computer and Information*

- Science, vol. 2012. , Singapore: Springer Nature Singapore, 2024, pp. 3–16. doi: 10.1007/978-981-99-9637-7_1.
- [6] A. Joshi, E. Fidalgo, E. Alegre, and L. Fernández-Robles, “DeepSumm: Exploiting topic models and sequence to sequence networks for extractive text summarization,” *Expert Systems with Applications*, vol. 211, p. 118442, Jan. 2023, doi: 10.1016/j.eswa.2022.118442.
- [7] Y. Liu and M. Lapata, “Text Summarization with Pretrained Encoders,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Hong Kong, China: Association for Computational Linguistics, 2019, pp. 3728–3738. doi: 10.18653/v1/D19-1387.
- [8] S. Phani, A. Abdul, M. Krishna Siva Prasad, and H. Kumar Deva Sarma, “MMSFT: Multilingual Multimodal Summarization by Fine-Tuning Transformers,” *IEEE Access*, vol. 12, pp. 129673–129689, 2024, doi: 10.1109/ACCESS.2024.3454382.
- [9] M. W. B. D. Satya, A. Luthfiarta, and M. N. Althoff, “Comparative Analysis of T5 Model Performance for Indonesian Abstractive Text Summarization,” *Sistemasi: Jurnal Sistem Informasi*, vol. 14, no. 3, pp. 1092–1106, May 2025, doi: 10.32520/stmsi.v14i3.4884.
- [10] A. F. Aji et al., “One Country, 700+ Languages: NLP Challenges for Underrepresented Languages and Dialects in Indonesia,” 2022, arXiv. doi: 10.48550/ARXIV.2203.13357.
- [11] A. Bankula and P. Bankula, “Cross-Linguistic Transfer in Multilingual NLP: The Role of Language Families and Morphology,” 2025, arXiv. doi: 10.48550/ARXIV.2505.13908.
- [12] M. Aurelia, S. Monica, and A. S. Girsang, “Transformer-based abstractive indonesian text summarization,” *IJ-ICT*, vol. 13, no. 3, p. 388, Dec. 2024, doi: 10.11591/ijict.v13i3.pp388-399.
- [13] F. Koto, A. Rahimi, J. H. Lau, and T. Baldwin, “IndoLEM and IndoBERT: A Benchmark Dataset and Pre-trained Language Model for Indonesian NLP,” 2020, arXiv. doi: 10.48550/ARXIV.2011.00677.
- [14] B. Wilie et al., “IndoNLU: Benchmark and Resources for Evaluating Indonesian Natural Language Understanding,” in *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, Suzhou, China: Association for Computational Linguistics, 2020, pp. 843–857. doi: 10.18653/v1/2020.aacl-main.85.
- [15] J. J. Sihombing, A. Arnita, S. I. Al Idrus, and D. Y. Niska, “Implementation of text summarization on indonesian scientific articles using textrank algorithm with TF-IDF web-based,” *J. Soft Comput. Explor.*, vol. 5, no. 3, pp. 310–319, Dec. 2024, doi: 10.52465/josce.v5i3.475.
- [16] A. R. Lubis, Y. Y. Lase, D. A. Rahman, and D. Witarasyah, “Improving Spell Checker Performance for Bahasa Indonesia Using Text Preprocessing Techniques with Deep Learning Models,” *ISI*, vol. 28, no. 5, pp. 1335–1342, Oct. 2023, doi: 10.18280/isi.280522.
- [17] G. P. A. Brahmantha, E. Utami, and A. Yaqin, “Classification of Indonesian Online News Topics Using Text Mining,” *jik*, vol. 16, no. 2, p. 7, Sept. 2023, doi: 10.24843/JIK.2023.v16.i02.p03.