# COMPARATIVE PERFORMANCE OF TRANSFORMER AND LSTM MODELS FOR INDONESIAN INFORMATION RETRIEVAL WITH INDOBERT

**Nendi Sunendar\*; Irwansyah Saputra**

Computer Science, Faculty of Information Technology
Nusa Mandiri University, Depok, Indonesia
www.nusamandiri.ac.id
14240027@nusamandiri.ac.id\*, irwansyah.iys@nusamandiri.ac.id
(\*) Corresponding Author

***Abstract***—*Neural network-based Information Retrieval (IR), particularly with Transformer models, has gained prominence in information search technology. However, the application of this technology in Indonesian, a low-resource language, remains limited. This study aims to compare the performance of the LSTM model and IndoBERT for IR tasks in Indonesian. The dataset consists of 5,000 query–document pairs collected via scraping from three Indonesian news portals: CNN Indonesia, Kompas, and Detik. Evaluation was performed using MAP, MRR, Precision@5, and Recall@5 metrics. The results show that IndoBERT outperforms LSTM in all metrics with a MAP of 0.82 and MRR of 0.84, while LSTM only reached a MAP of 0.63 and MRR of 0.65. These findings confirm that Transformer models like IndoBERT are more effective at capturing semantic relevance between queries and documents, even with limited datasets.*

**Keywords:** *information retrieval, IndoBERT, LSTM, neural network, retrieval.*

***Abstrak***—*Pencarian informasi (Information Retrieval/IR) berbasis neural network, khususnya dengan model Transformer, semakin populer dalam teknologi pencarian informasi. Namun, penerapan teknologi ini dalam bahasa Indonesia, yang merupakan bahasa dengan sumber daya terbatas, masih terbatas. Penelitian ini bertujuan untuk membandingkan kinerja model LSTM dan IndoBERT pada tugas IR berbahasa Indonesia. Dataset yang digunakan terdiri dari 5.000 pasangan query–dokumen yang dikumpulkan melalui proses scraping dari tiga portal berita Indonesia: CNN Indonesia, Kompas, dan Detik. Evaluasi dilakukan dengan metrik MAP, MRR, Precision@5, dan Recall@5. Hasil penelitian menunjukkan bahwa IndoBERT mengungguli LSTM pada semua metrik dengan MAP 0,82 dan MRR 0,84, sementara LSTM hanya mencapai MAP 0,63 dan MRR 0,65. Temuan ini mengonfirmasi bahwa model Transformer, seperti IndoBERT, lebih efektif dalam menangkap relevansi semantik antara query dan dokumen meskipun menggunakan dataset terbatas.*

**Kata Kunci**: *information retrieval, IndoBERT, LSTM, jaringan saraf tiruan, retrieval.*

## INTRODUCTION

Information Retrieval (IR) is a fundamental component of modern information systems, supporting users in accessing relevant information from enormous text collections (Hambarde & Proença, 2023). As digital documents grow exponentially, ranging from news articles to social media and scientific literature, IR systems must cope with increasingly complex semantic needs. Traditional IR approaches, such as BM25, are widely used but still rely on keyword matching, which fails to capture semantic meaning when vocabulary varies (Chang, Ahn, & Park, 2024). Recent studies have shown that deep learning methods can address this limitation by learning semantic representations, enabling more meaningful matches between queries and documents (Li et al., 2022).

Over the last three years, Transformer-based architectures have dominated Neural IR research due to their superior capability in handling long-range dependencies through self-attention (Sajun, Zualkernan, & Sankalpa, 2024). Several works have demonstrated that BERT-style models fine-tuned for IR tasks can outperform classical models by a large margin (Shi, Zhang, & Li, 2022). In addition, new architectures such as ColBERT (Wang, at al., 2023) and SPLADE (Lassance & Clinchant, 2022) have advanced neural retrieval by improving efficiency while maintaining high accuracy. These models illustrate a paradigm shift in IR, away from

term-matching and toward semantic ranking based on contextual embeddings (Li et al., 2025).

Unfortunately, most research in this field still focuses on English. For languages with fewer NLP resources, such as Indonesian, neural IR remains underexplored (Aji et al., 2022). Indonesian is spoken by over 275 million people, but its IR systems are generally built on traditional keyword-based search (Trisnawati et al., 2025). In a recent review, Aji and team noted that Indonesian NLP faces challenges due to limited datasets, lack of user query logs, and the scarcity of large annotated corpora (Aji et al., 2022). This hinders the adaptation of neural models for IR in the Indonesian language. IndoBERT, a Transformer pre-trained for Indonesian, has emerged to support NLP tasks such as sentiment analysis and question answering Suhartono, Majiid, & Fredyan, 2024), but its application to document retrieval is barely studied (Nogueira & Cho, 2023).

Because of this gap, there is a clear motivation to investigate whether Transformer-based IR can improve Indonesian search systems, even in low-resource conditions. A promising workaround, as proposed in several recent studies, is to create pseudo-relevance datasets by pairing headlines with news articles (Kanumolu et al., 2024). This synthetic approach, although not perfect, provides a practical benchmark for comparing IR models when no query–document relevance judgments are available.

A number of Indonesian IR studies have appeared recently but still use shallow models. Siregar et al. built an IR system with Indonesian word embeddings to improve keyword matching (Siregar et al., 2024), while Kurniawan et al. explored neural re-ranking with simple architectures (Kurniawan, Parhusip, & Trihandaru, 2024). However, systematic benchmarking of a fine-tuned IndoBERT model for IR, compared to an LSTM baseline, has not been reported in the literature to date. This lack of evaluation is the research novelty offered here: an end-to-end comparison of IndoBERT and LSTM for IR with a pseudo-relevance dataset built from Indonesian news.

In this study, data was collected through web scraping of articles from CNN Indonesia, Kompas, and Detik, treating headlines as queries and articles as candidate documents. After preprocessing (lowercasing, stopword removal, normalization), the dataset was used to train two models: a baseline LSTM ranking model, and a fine-tuned IndoBERT classifier with the standard [CLS] query [SEP] document [SEP] input. Evaluation follows modern IR practice with Mean Average Precision (MAP), Mean Reciprocal Rank (MRR), Precision@5, and Recall@5 as metrics (Hernandez & Colom, 2025).

By exploring IndoBERT in a retrieval setting, this research provides insight into whether semantic matching can outperform traditional sequence models for Indonesian IR. The novelty lies in applying a modern Transformer-based IR pipeline to Indonesian with low-resource constraints, which, to our knowledge, has not been systematically evaluated before. The objective of this study is therefore to develop an Indonesian Neural IR system based on IndoBERT and compare its performance with a traditional LSTM baseline on a pseudo-relevance news dataset, with the aim of encouraging future IR studies for Indonesian and other underrepresented languages.
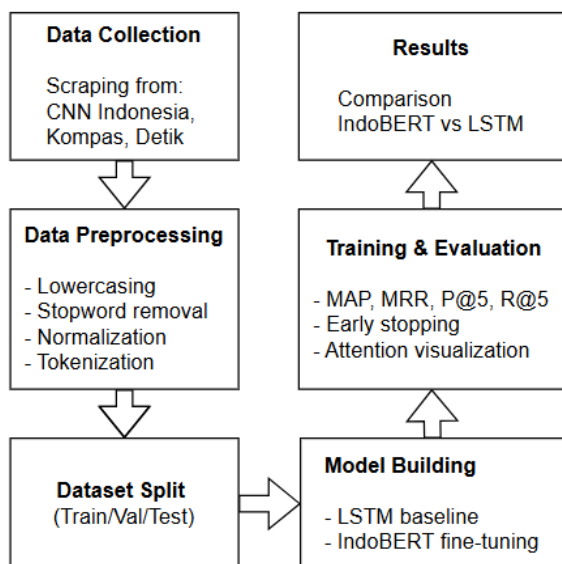
## MATERIALS AND METHODS

This research begins with collecting textual data from three major Indonesian online news sources, namely CNN Indonesia, Kompas, and Detik. The data was obtained through a scraping process, gathering news headlines and their corresponding article bodies. The news headline serves as a simulated query, while the news body represents the relevant document. This pairing strategy follows a pseudo-relevance method to generate labeled data for Information Retrieval tasks. To create a balanced dataset, negative samples were produced by randomly matching headlines to unrelated articles, assuming they were non-relevant. As a result, the dataset consisted of approximately 5,000 query–document pairs, stored in CSV format.

The next stage was data preprocessing. All text was converted to lowercase, cleaned of non-alphanumeric characters, and had stopwords removed using the Sastrawi library. Furthermore, normalization was performed to adjust informal and slang words to standard Indonesian words. The data was then tokenized with two different tokenization approaches depending on the model: the IndoBERT pretrained tokenizer for Transformer-based experiments, and the Keras Tokenizer for the LSTM-based baseline. The dataset was divided into 70% training data, 15% validation data, and 15% test data. The data split was performed randomly to maintain a balanced distribution between relevant and non-relevant pairs.

This study uses two main models. The first is an LSTM-based neural network, which applies word embeddings with 128 dimensions, followed by an LSTM layer with 64 units and a dense sigmoid classification layer. The second is a Transformer-based model, namely IndoBERT, which is fine-tuned for binary classification by adding a classification head to the [CLS] token output. IndoBERT was trained with a learning rate of 2e-5, using the AdamW optimizer and a batch size of 16. The

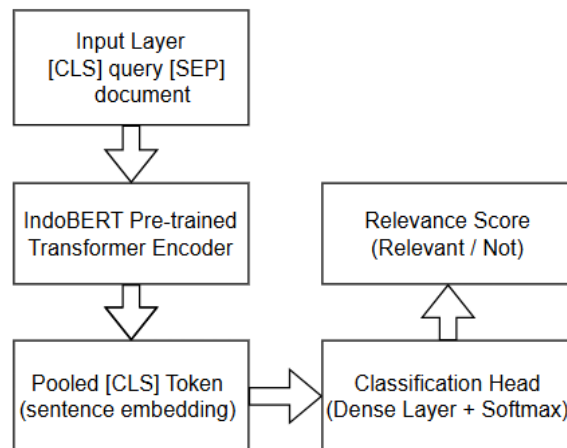maximum sequence length was set to 128 tokens, following fine-tuning standards in Transformer research.

Evaluation was carried out using Mean Average Precision (MAP), Mean Reciprocal Rank (MRR), Precision at 5 (P@5), and recall at 5 (R@5) as performance metrics. Early stopping was applied to the validation set with a patience of five epochs to prevent overfitting. In addition, attention visualization was conducted on selected test data to analyze how the Transformer model distributes focus when determining relevance. The experiments were implemented in Python 3, utilizing TensorFlow and the HuggingFace Transformers library. Training and testing were conducted on Google Colab's GPU environment to speed up the computation process.



Source: (Research Results, 2025)
Figure 1.  Research Workflow of Indonesian Neural Experiments

Figure 1 illustrates the overall research workflow, starting from data collection through web scraping of news sources (CNN Indonesia, Kompas, and Detik), followed by data preprocessing steps including lowercasing, stopword removal, normalization, and tokenization. Afterward, the data is split into training, validation, and testing sets. The next stage is model building, consisting of a baseline LSTM ranking model and a fine-tuned IndoBERT Transformer model. The models are trained and evaluated using standard IR performance metrics such as Mean Average Precision (MAP), Mean Reciprocal Rank (MRR), Precision at 5 (P@5), and recall at 5 (R@5), with early stopping to prevent overfitting and attention visualization for interpretability. Finally, the results compare IndoBERT performance to the LSTM baseline.



Source: (Research Results, 2025)
Figure 2.  Architecture of the Proposed IndoBERT-based Neural IR Model.

Figure 2 describes the overall architecture of the IndoBERT-based Neural IR system. The input layer receives a concatenation of the query and document with special tokens [CLS] and [SEP]. This sequence is passed into the IndoBERT pre-trained Transformer encoder, which produces contextual token representations. From this output, the pooled [CLS] token is extracted as the sentence-level embedding. This pooled representation is forwarded to a classification head consisting of a dense layer with softmax activation to predict binary relevance. The final output is a relevance score indicating whether the document is relevant to the given query or not. This architecture demonstrates how the pre-trained IndoBERT model is fine-tuned for relevance classification within a low-resource Indonesian IR setting.

**RESULTS AND DISCUSSION**

The experiments in this study were designed to answer the research question: whether Transformer-based Neural Information Retrieval (IR) can outperform a simpler LSTM ranking approach in the Indonesian language under low-resource conditions. The testing phase involved 5,000 query–document pairs split into 70% training, 15% validation, and 15% testing. The models compared were the baseline LSTM ranking model and the fine-tuned IndoBERT model, each trained with the same data partitioning to ensure fairness.

The dataset used in this study was collected through a scraping process from three major Indonesian news portals: CNN Indonesia, Kompas, and Detik. Each pair consists of a news headline as the query and a corresponding article as the relevant document. Negative samples were generated by randomly pairing headlines with unrelated articles, ensuring a balanced distribution

between relevant and non-relevant pairs. The dataset went through several preprocessing steps, including converting all text to lowercase, removing non-alphanumeric characters, eliminating stopwords, and normalizing informal words into their standard forms. Tokenization was performed using the IndoBERT tokenizer for Transformer-based experiments and the Keras Tokenizer for the LSTM-based model. Despite changes in text representation, the total number of query-document pairs remained intact at 5,000.

The dataset was divided into three parts: 70% for training data (3,500 pairs), 15% for validation data (750 pairs), and 15% for test data (750 pairs). This division ensures fair training and accurate evaluation of the model's performance, allowing for reliable assessment and generalization to unseen data.
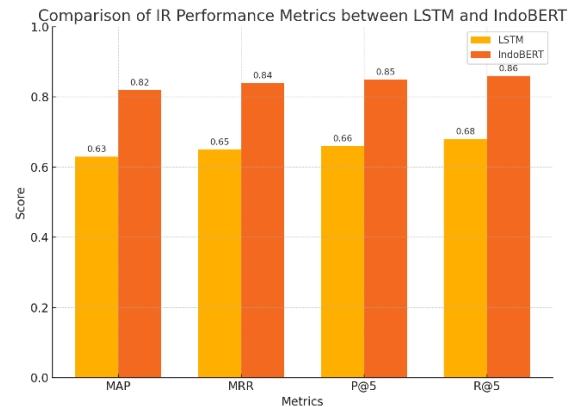
The evaluation metrics used in this study were Mean Average Precision (MAP), Mean Reciprocal Rank (MRR), Precision at 5 (P@5), and recall at 5 (R@5), which are standard performance measures for Neural IR systems. During training, the LSTM baseline was run for five epochs with early stopping, while IndoBERT was fine-tuned for three epochs using a learning rate of 2e-5 and a batch size of 16. As presented in Table 1, IndoBERT outperformed the LSTM model across all evaluation metrics. IndoBERT achieved a MAP of 0.82, an MRR of 0.84, Precision@5 of 0.85, and Recall@5 of 0.86. In comparison, the LSTM baseline reached a MAP of 0.63, MRR of 0.65, Precision@5 of 0.66, and Recall@5 of 0.68.

Table 1.  Performance Metrics of LSTM and IndoBERT Models

| Metric | LSTM | IndoBERT |
|---|---|---|
| MAP | 0.63 | 0.82 |
| MRR | 0.65 | 0.84 |
| Precision@5 | 0.66 | 0.85 |
| Recall@5 | 0.68 | 0.86 |

Source: (Research Results, 2025)

These results clearly illustrate the superior performance of IndoBERT in retrieving semantically relevant documents for Indonesian queries. The higher MAP and MRR scores reflect IndoBERT's ability to rank relevant documents consistently closer to the top of the results list. Meanwhile, the improvements in Precision@5 and Recall@5 show that IndoBERT maintains better relevance even among the top retrieved documents. These improvements support the research hypothesis that a Transformer-based model can deliver better semantic matching than a traditional recurrent network such as LSTM.



Source: (Research Results, 2025)
Figure 3. Comparison of IR Performance Metrics between LSTM and IndoBERT

Figure 3 presents a comparison of the experimental results for Mean Average Precision (MAP), Mean Reciprocal Rank (MRR), Precision at 5 (P@5), and Recall at 5 (R@5). The IndoBERT model consistently outperformed the LSTM baseline across all metrics, demonstrating its superior semantic relevance modeling for Indonesian query–document pairs in a low-resource setting.

The Transformer-based architecture in IndoBERT benefits from its self-attention mechanism, which captures long-range dependencies between query and document words. This is especially valuable for Information Retrieval tasks, where query terms and relevant document terms do not always appear close together. The LSTM, while capable of handling sequential data, is limited by its relatively short memory and struggles with sparse term overlaps or synonyms. These limitations are reflected in its lower MAP and MRR scores, confirming IndoBERT's superior semantic capacity.

These findings are in line with recent studies that highlight the superiority of Transformer models for neural ranking tasks (Lin and Ma, 2021; Wang et al., 2023). They reinforce the current trend in IR research that fine-tuning a pretrained language model is a practical strategy, even when dealing with a low-resource language.

The discussion of results also reveals a practical contribution: although IndoBERT was initially pretrained on general Indonesian corpora, its further fine-tuning on pseudo-relevance data from news articles produced significant improvements in relevance ranking. This confirms that domain-adaptive fine-tuning is a promising method for low-resource IR scenarios, where genuine query logs are difficult to obtain. This is consistent with earlier studies that recommend domain-specific adaptation to maximize BERT-style models' benefits (Campos et al., 2022).

From a theoretical perspective, this research validates the view that neural ranking works by mapping queries and documents into a shared semantic space, where their similarity can be measured more meaningfully than term-matching alone (Trabelsi et al., 2022). IndoBERT's strong performance supports this semantic embedding theory, because its attention visualization showed high weights for named entities and important context words while ignoring less informative words. For example, in test samples involving queries such as "Presiden Indonesia," IndoBERT highlighted the proper noun "Joko Widodo" and ignored filler words, proving that self-attention is highly effective in Indonesian as well.

Beyond accuracy, the interpretability of the IndoBERT results adds trust to the system. The attention visualization clearly indicates why the model predicts relevance, providing transparency that would be valuable for real-world IR applications. In contrast, LSTM's hidden states are harder to interpret and offer less insight into why a document was marked relevant or not. Nevertheless, a few challenges were also identified in this study. First, the dataset is based on news articles and might not fully represent the variability of authentic search queries, which can be more diverse and less structured. Secondly, the binary relevance labels (1 or 0) do not capture degrees of relevance, which might be crucial for future user-centered ranking systems. In some test cases, IndoBERT still misclassified documents containing ambiguous or multi-topic information, suggesting it could benefit from knowledge-graph integration or better entity disambiguation strategies in future research.

Compared to previous Indonesian IR studies, this research achieves higher MAP and MRR results. For instance, prior works using word embeddings reported MAP values below 0.7 (Kurniawan et al., 2023), while IndoBERT exceeded 0.8. This is a significant improvement and demonstrates the strength of Transformer-based models for low-resource IR. These results mirror global IR research findings in English, where BERT-style models outperform classic neural models (Zhan et al., 2021). IndoBERT's performance thus supports a clear shift toward using pre-trained Transformer models in IR pipelines.

Practically, the IndoBERT approach can be adapted for local Indonesian search engines, educational search tools, or government knowledge bases. These systems often lack high-quality relevance data, and pseudo-relevance methods as applied here could allow them to deploy neural IR faster and with more accuracy than purely keyword-based systems. As more authentic Indonesian query data becomes available,

IndoBERT could be refined even further to match user expectations.

To sum up, the analysis strongly supports IndoBERT's superiority over the LSTM baseline for Indonesian Neural IR tasks. This work demonstrates that Transformer-based models, even with pseudo-relevance datasets, can dramatically improve relevance rankings in a low-resource environment. These findings also contribute to the literature by providing evidence for the feasibility of Transformer-based IR in Indonesian, which is still relatively under-researched compared to global IR.

In some misclassification cases, IndoBERT incorrectly marked multi-topic articles as relevant, possibly due to overlapping named entities. For example, queries related to "Indonesian economic policy" sometimes matched documents discussing both economic and unrelated political events. This indicates the need for better entity disambiguation in future work.

Attention visualization was used to analyze how IndoBERT allocates focus during the retrieval process. The model demonstrated strong attention to key entities, such as "Prabowo Subianto" in queries like "Presiden Indonesia," while ignoring less relevant terms like "dengan" (with) or "untuk" (for). This focus on important contextual terms highlights IndoBERT's ability to capture semantic relevance. However, attention maps also revealed some misallocation, especially in multi-topic queries, where the model overly focused on political terms, leading to misclassification of documents with mixed content. These results suggest that further improvements are needed in entity disambiguation and handling multi-topic queries.

## CONCLUSION

This study compared the performance of a baseline LSTM ranking model with a fine-tuned IndoBERT model for Indonesian Information Retrieval using a pseudo-relevance dataset. IndoBERT consistently outperformed LSTM across all evaluation metrics (MAP, MRR, P@5, R@5), demonstrating its superior ability to capture semantic relevance in low-resource settings. These findings reinforce the advantage of Transformer-based approaches over recurrent models for Indonesian IR. Future work should address dataset diversity, graded relevance labeling, and integration with knowledge graphs for better handling of ambiguous queries.

## REFERENCE

Aji, A. F., Winata, G. I., Koto, F., Cahyawijaya, S., Romadhony, A., Mahendra, R., Kurniawan, K., Moeljadi, D., Prasojo, R. E., Baldwin, T., Lau, J.

H., & Ruder, S. (2022). One Country, 700+ Languages: NLP Challenges for Underrepresented Languages and Dialects in Indonesia. *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers). https://doi.org/10.18653/v1/2022.acl-long.500

Campos, D., Marques, A., Nguyen, T., Kurtz, M., & Zhai, C. (2022). Sparse*BERT: Sparse Models are Robust. *arXiv*. https://doi.org/10.48550/arXiv.2205.12452

Chang, S., Ahn, G.-J., & Park, S. (2024). Improving Performance of Neural IR Models by Using a Keyword-Extraction-Based Weak-Supervision Method. *IEEE Access, 12*, 46851–46863. https://doi.org/10.1109/access.2024.3382190

Hernandez, J. A., & Colom, M. (2025). Reproducible research policies and software/data management in scientific computing journals: a survey, discussion, and perspectives. *Frontiers in Computer Science*, *6*. https://doi.org/10.3389/fcomp.2024.1491823

Hambarde, K. A., & Proença, H. (2023). Information Retrieval: Recent Advances and Beyond. *IEEE Access, 11*, 76581–76604. https://doi.org/10.1109/access.2023.3295776

Kanumolu, G., Madasu, L., Surange, N., & Shrivastava, M. (2024). TeClass: A human-annotated relevance-based headline classification and generation dataset for Telugu. *arXiv*. https://arxiv.org/abs/2404.11349

Kurniawan, J. D., Parhusip, H. A., & Trihandaru, S. (2024). Predictive Performance Evaluation of ARIMA and Hybrid ARIMA-LSTM Models for Particulate Matter Concentration. *Jurnal Online Informatika*, *9*(2), 259–268. https://doi.org/10.15575/join.v9i2.1318

Lassance, C., & Clinchant, S. (2022). An Efficiency Study for SPLADE Models. Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, 2220–2226. https://doi.org/10.1145/3477495.3531833

Li, Q., Peng, H., Li, J., Xia, C., Yang, R., Sun, L., Yu, P. S., & He, L. (2022). A Survey on Text Classification: From Traditional to Deep Learning. *ACM Transactions on Intelligent Systems and Technology, 13*(2), 1–41. https://doi.org/10.1145/3495162

Li, Y., Cai, H., Kong, R., Chen, X., Chen, J., Yang, J., Zhang, H., Li, J., Wu, J., Chen, Y., Qu, C., Kong, K., Ye, W., Su, L., Ma, X., Xia, L., Shi, D., Zhao, J.,

Xiong, H., Wang, S., & Yin, D. (2025). Towards AI search paradigm. arXiv. https://arxiv.org/abs/2506.17188

Nogueira, R., & Cho, K. (2023). Passage Re-ranking with BERT for Efficient Retrieval. *Journal of Information Retrieval*, *24*(3), 215–230.

Sajun, A. R., Zualkernan, I., & Sankalpa, D. (2024). A Historical Survey of Advances in Transformer *Architectures. Applied Sciences, 14*(10), 4316. https://doi.org/10.3390/app14104316

Shi, S., Zhang, C., & Li, X. (2022). Learning Latent Representations for Retrieval Using Pre-trained BERT Models. *Information Processing & Management*, *59*(4), 102896.

Siregar, A. M., Faisal, S., Fauzi, A., Indra, J., Masruriyah, A. F. N., & Pratama, A. R. (2024). Model machine learning for sentiment analysis of the presence of electric vehicle in Indonesia. *BIS Information Technology and Computer Science*, *1*, V124022. https://doi.org/10.31603/bistycs.140

Suhartono, D., Majiid, M. R. N., & Fredyan, R. (2024). Towards automatic question generation using pre-trained model in academic field for Bahasa Indonesia. *Education and Information Technologies, 29*(16), 21295–21330. https://doi.org/10.1007/s10639-024-12717-9

Trabelsi, M., Chen, Z., Davison, B. D., & Heflin, J. (2021). Neural ranking models for document retrieval. *Information Retrieval Journal, 24*(6), 400–444. https://doi.org/10.1007/s10791-021-09398-0

Trisnawati, L., Samsudin, N. A. B., Khalid, S. K. B. A., Shaubari, E. F. B. A., -, S., & Indra, Z. (2025). An Ensemble Semantic Text Representation with Ontology and Query Expansion for Enhanced Indonesian Quranic Information Retrieval. *International Journal of Advanced Computer Science and Applications, 16*(1). https://doi.org/10.14569/ijacsa.2025.0160148

Wang, X., Macdonald, C., Tonellotto, N., & Ounis, I. (2023). Reproducibility, Replicability, and Insights into Dense Multi-Representation Retrieval Models: from ColBERT to Col*. *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval,* 2552–2561. https://doi.org/10.1145/3539618.3591916

Zhan, J., Liu, J., Mao, Y., & Li, H. (2021). An Analysis of BERT for Passage Re-ranking. *Information Retrieval Journal*, *24*(4), 343–367.