

<sup>1</sup>Ibadurrohman  
Irfan Fatani,

<sup>2</sup>Herry Irawan

· **Twitter, Instagram, Youtube Speak:  
Understanding Sentiments on LRT  
Jabodebek Services via Inset Lexicon,  
IndoBERT and BERTopic Approaches**



**Abstract:** - Rapid urbanization in the Jabodetabek region has led to an increased demand for public transportation. Responding to this need, the government has initiated the development of a new public transportation mode, namely the LRT Jabodebek. However, as a new public transportation mode, the LRT Jabodebek has both strengths and weaknesses in serving the community. Various public comments are expressed through social media platforms. To enhance service quality, it is crucial to pay attention to public comments. Therefore, a sentiment analysis is required to identify and delve into both positive and negative sentiments regarding the LRT Jabodebek service through comments on Twitter, Instagram, and Youtube. The methodology involves a combination of Lexicon-based, IndoBERT model, and BERTopic approaches to gain a deeper understanding of the Jabodebek LRT service trends. The study reveals that 55.9% of the 8,523 comments carry a negative sentiment, and the IndoBERT model achieves an accuracy of 85.97% in sentiment classification.

**Keywords:** IndoBERT, BERTopic, Lexicon-Based Approach, Light Rail Transit.

## I. INTRODUCTION

Indonesia is undergoing rapid urbanization, with its major cities experiencing significant growth [1], especially in metropolitan areas like Jabodetabek [2]. Therefore, the need for efficient public transportation in Jabodetabek is crucial in shaping the daily activities of the people [3]. Responding to this demand, the government introduced a new public transportation mode, namely Light Rail Transit (LRT) Jabodebek (Jakarta, Bogor, Depok, Bekasi), in August 2023 [4]. The aim is to fulfil the mobility needs of the population and contribute to the overall well-being of society [5].

The Role of LRT Jabodebek as a mode of transportation in service quality has positive implications. By providing comfort, timeliness, and safety, LRT Jabodebek can efficiently address the daily mobility needs of metropolitan residents [6]. The reduction in travel time, increased productivity, and decreased air pollution illustrates the tangible benefits of LRT Jabodebek to encourage people to switch to this public transportation option [7]. Despite the addition of a new public transportation mode that can improve mobilization with a quick travel time, the number of Jabodebek LRT users is not as substantial as it was during its initial operations [8]. This is allegedly because the community is perceived to be less satisfied when using this mode of transportation, which has various shortcomings in its service quality.

However, the long-term success of LRT Jabodebek relies on the quality of its services to provide a satisfying user experience. This enhancement is not merely an option but a necessity to ensure that LRT Jabodebek meets the expectations of the public, supporting more efficient daily lives and becoming the primary transportation choice. Therefore, promptly detecting and addressing issues related to services is a critical element [9]. Studies and analyses of the performance of LRT Jakarta and LRT in Klang Valley highlight efforts to enhance operational efficiency and services [6], [10]. This analysis underscores the urgency of service quality in public transportation and the continuous improvement needed to enhance the success of LRT Jabodebek [6].

The LRT Jabodebek, as a new mode of transportation, faces various challenges such as insufficient preparation for operational launch, recurring technical disruptions, and ensuring services alignment with user needs. These challenges have generated diverse public opinions regarding its operational services. The public's perception of LRT Jabodebek services, whether positive or negative, plays a crucial role in shaping improvements and

<sup>1</sup>School of Economic and Business, Telkom University,

<sup>2</sup>School of Economic and Business, Telkom University

<sup>1</sup>ibadirfan@student.telkomuniversity.ac.id, <sup>2</sup>herryir@telkomuniversity.ac.id

enhancements to this new mode of transportation [11]. Understanding diverse opinions in society is a complex task, involving exploration of various perspectives among users and non-users of LRT Jabodebek. Additionally, sentiments expressed by individuals on social media platforms serve as valuable data, reflecting diverse opinions, complaints, and aspirations related to LRT Jabodebek services [12].

In response to the challenges posed by the abundance of social media data, sentiment analysis becomes a relevant tool [13]. This research utilizes sentiment analysis to identify and delve into positive and negative sentiments regarding LRT Jabodebek services using comments on social media [14]. We combine Lexicon-based approaches [15], the application of IndoBERT specifically designed for Indonesian language analysis [16], and topic modelling techniques using BERTopic [17]. This research offers an innovative and holistic approach aiming to measure positive and negative sentiments of the community towards the Jabodebek LRT service using the InSet Lexicon and IndoBERT approaches. We analyzed frequently occurring topics for each polarity using BERTopic and assesses the dimensions of service quality improvement needed by the relevant stakeholders of the LRT Jabodebek. The results of this study can provide recommendations for LRT Jabodebek to obtain feedback for improving service quality.

## II. RELATED WORK

Several previous studies have explored sentiment analysis using the Lexicon-based approach for positive and negative labelling, the IndoBERT model approach, and topic modelling techniques using BERTopic. In [16], this study explores the growing significance of health service applications in societal life, they are employing the pre-training methodology of IndoBERT to evaluate user satisfaction with the Halodoc, Alodokter, and Klikdokter applications. The IndoBERT model delivers remarkable outcomes, boasting an accuracy score of 96%, recall of 96%, precision of 95%, and an F1 score of 95%. In [18], this research examines public sentiment towards PPKM policies in Indonesia during the pandemic using Twitter data. The study employs the BERT method with the IndoBERT model and compares it with SVM and Naive Bayes for sentiment analysis. They also utilize a lexicon approach to label data into three classes (positive, neutral, negative). The research findings indicate that IndoBERT achieves impressive evaluation metrics with an F-1 score of 84%, precision of 86%, and recall of 84%. In contrast, traditional methods such as SVM and Multinomial Naive Bayes show lower metrics. Zhang [19] conducted sentiment analysis research, amalgamating the advantages of lexicon-based and classification-based methods while addressing their shortcomings. Lexiconed BERT exhibited commendable performance in processing lengthy sentences, leading to the conclusion that the Lexiconed BERT approach effectively enhanced sentiment classification performance.

In [20], They compared the Latent Dirichlet Allocation (LDA) and BERTopic approaches to identify factors contributing to attitudes toward vaccines for three different vaccine brands. The research results indicate that BERTopic clustering outperforms LDA clustering. In [17], They also utilized the BERTopic model as the topic distribution for classifying Indonesian fake news. In [21], They also used BERTopic for topic modelling in the mHealth research domain.

## III. METHODOLOGY

This research involves several steps, starting from data collection, data preprocessing, and utilizing a lexicon-based approach along with the IndoBERT model. Ultimately, BERTopic technique is employed for topic modelling to identify positive and negative sentiments toward the Jabodebek LRT service.

### A. Data Collection

In the first phase is data collection process [22], we gathered opinion data using the effective web scraping tool APIFY to extract information from Twitter, Instagram, and Youtube [23]. The data collection occurred from September 1 to October 31, 2023, with the keyword "layanan" and "LRT Jabodebek" as well as posts about Jabodebek LRT services and we obtained a total of 8523 comments, including 3958 from Twitter, 1380 from Instagram, and 3185 from Youtube.

### B. Data Preprocessing

The second phase is data preprocessing, where the initial text data is prepared to facilitate easier processing by

algorithm [24]. First, we do case folding and convert slang terms to their normal forms [25]. The next step involves filtering by removing emojis, hashtags, punctuation, numbers, repeated words, and unwanted white spaces [11]. Subsequently, the process includes tokenization, where the text is broken down into small units [25]. Stopwords removal is then applied, and finally, stemming is implemented [26].

### C. *Lexicon-Based Approach*

In the third phase, data labelling employs a lexicon-based approach that utilizes lexical resources to determine the sentiment of words in the text [26]. In the application of this approach, we utilize the InSet Lexicon specifically designed for the Indonesian language to classify positive and negative sentiments [15]. This method aggregates sentiment scores for each word in the document using a predefined dictionary to match negative or positive lexicons. Subsequently, comments are labelled as positive (polarity > 0) or negative (polarity < 0) [12].

### D. *IndoBERT Model*

IndoBERT is a language model specifically trained and adapted for the Indonesian language. It is based on BERT (Bidirectional Encoder Representations from Transformers) and has undergone advanced training using Indonesian language datasets [18], this model is effective in sentiment analysis because this algorithm continues to be developed based on Natural Language Processing which is suitable for processing public opinion with big data. In this study, IndoBERT is employed for sentiment analysis on the LRT Jabodebek service. The labelled data is divided into three sets: training, validation, and testing, with an 80:10:10 distribution ratio [27]. After dividing the data proportions, we load the tokenizer from the IndoBERT model and conduct a vocabulary analysis of the pre-trained IndoBERT model [28]. Subsequently, we create functions to merge tokenization steps and add special tokens. Further, mapping functions are applied to produce a format compatible with the BERT model. The BERT model for classification is loaded from pre-trained models, compiled, and undergoes training for 10 epochs using training and validation data. After saving the trained model, evaluation on test data is conducted to obtain testing accuracy. Model predictions are then compared with actual labels using a confusion matrix [16]. This performance measurement helps evaluate how successful the model is in sentiment classification on review data using the IndoBERT model [18].

### E. *BERTopic*

This topic modelling employs BERT embedding vectors and c-TF-IDF to create compact topic clusters, enhancing the interpretability of topics [17]. By leveraging BERT embeddings, BERTopic generates richer and more contextual text representations, enhancing the accuracy of topic identification [20]. The outcomes of this process encompass two aspects: related topic sequences and associated probabilities. The topic sequences provide an understanding of the relationships between topics, while the probabilities indicate how likely a sentence is to belong to a particular topic [17].

## IV. RESULT

In this research, we obtained a total of 8523 comments about the LRT Jabodebek. This data constitutes raw information that underwent preprocessing to prepare it for analysis within the model algorithm. Before training the IndoBERT model, we performed the process of classifying the processed data into positive or negative classes using the InSet Lexicon method. Based on the process, 2891 positive sentiments and 3665 negative sentiments were obtained, comments with neutral sentiments were excluded as this research specifically focuses on positive and negative sentiments, indicating that positive reviews encompass 44.1% of the data, while negative reviews constitute 55.9%. This indicates that most sentiments expressed regarding the LRT Jabodebek service are negative. A higher percentage of negative sentiments reflects the public's perception of the LRT service, with most opinions being negative.

In the application of the IndoBERT model, the dataset was divided into three proportions to form the training data consisting of 5244 instances, validation data consisting of 656 instances, and test data consisting of 656 instances. The research utilized 10 epochs during the model training process. The result revealing that utilizing this approach led to favourable accuracy. This was substantiated by a noticeable upward trend in the curve.

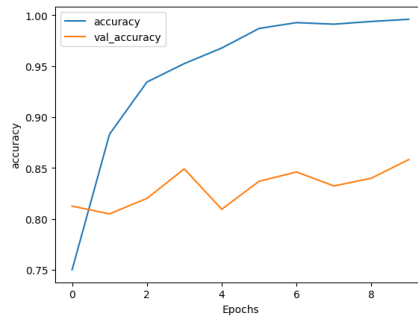


Fig. 1. Training accuracy history

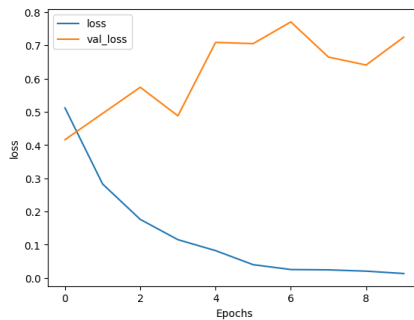


Fig. 2. Training loss history

According to Fig.1, as the number of iterations (epochs) increases in the model, the accuracy curve shows an upward trend to the right, and on the Fig.2 while the loss curve indicates a decreasing trend. This suggests that the model has been well-trained and is becoming more accurate in predictions. It is a positive indication that the model has undergone effective training and is ready to be tested on the test data. The research involves evaluating the performance of the trained model with test data using a confusion matrix, adding a level of detail in understanding the model's performance on the test data. The process of evaluating the model against the test data is explained in detail in the Fig. 3.

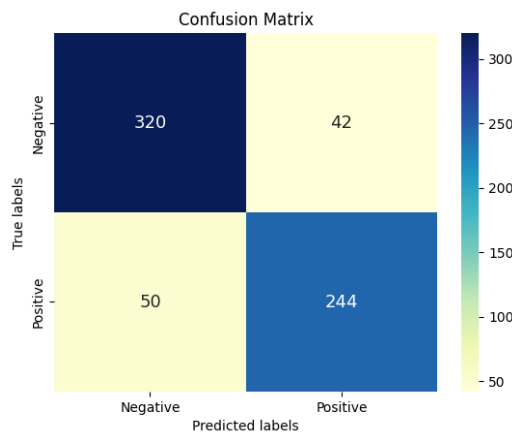


Fig. 3. Confusion matrix

According to Fig. 3, the IndoBERT model demonstrates excellent predictive capabilities for classifying sentiments related to the LRT Jabodebek service. Specifically, the values of true positive and true negative are significantly higher compared to false positive and false negative. Findings from this research reveal a commendable accuracy score of 85.97%, an F1-Score of 84.13%, a recall of 82.99%, and a Precision of 85.31%. Overall, the evaluation results indicate that the model exhibits a high level of accuracy, maintaining a balance between precision and recall. It also demonstrates proficiency in recognizing positive instances and providing accurate positive predictions for sentiment analysis of the LRT Jabodebek service.

Next, topic modelling analysis was conducted for each sentiment. The data used in the topic modelling analysis has undergone preprocessing and is labelled. We divided the topic modelling process for each positive and negative sentiment. This allows us to understand which topics are most frequently discussed in the LRT Jabodebek service for each sentiment by examining the frequency of recurring words.

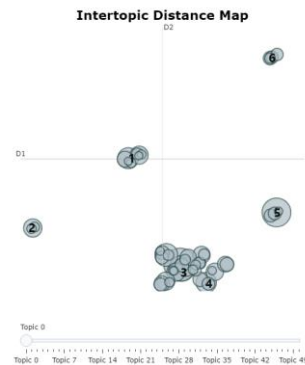


Fig. 4. Clustering positive

Based on the topic modelling analysis of positive sentiment towards LRT Jabodebek can be seen in Fig.4, there are several dimensions that reflect positive sentiments from the public:

Firstly, in cluster number 2 and 3, users appreciated the useful content and information presented by LRT Jabodebek's social media accounts on social media platforms. They highlighted the attention given to the relevance of public information needs related to LRT Jabodebek, showing a good response to the information delivery efforts made by the account, which reflects the service quality dimension of good responsiveness from LRT Jabodebek to user needs.

Secondly, in cluster number 1, 3, and 5, there is a high appreciation of the aspects of transport and modes that positively affect society, with a focus on the utilisation of high quality products produced locally by INKA companies. This reflects support for the development of transport infrastructure to improve people's quality of life, which demonstrates the service quality dimension of tangibility.

Third, in cluster number 3, there is public appreciation of LRT Jabodebek as an advancement of the Indonesian state, where LRT Jabodebek effectively contributes to serving daily activities that run smoothly in a crowded city, reflecting the service quality dimension of reliability of the services provided.

Fourthly, in cluster number 3 and 5, stakeholders in transport have strong beliefs and support study and development projects, showing commitment to maintain and improve service standards by playing an important role in ensuring optimal service quality, which is the assurance aspect of service quality.

Finally, in cluster number 3 and 6, Railmin, the LRT Jabodebek social media admin who interacts with the public through social media, provides an understanding response to users through their official accounts. They actively apologise and seek users' understanding when faced with undesirable situations, demonstrating their awareness of what users may experience, reflecting the empathy aspect of the service provided.

Thus, the analysis results show that responsiveness, tangibility, reliability, assurance, and empathy are important dimensions in creating positive perceptions of LRT Jabodebek services.

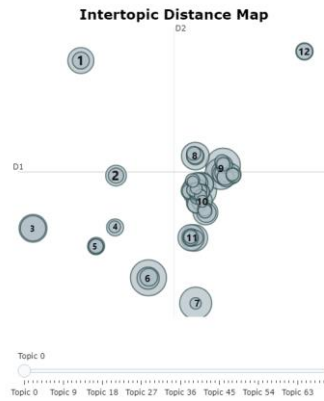


Fig. 5. Clustering negative

On the other hand, topic modeling for negative sentiment can be seen in Fig.5, there are several dimensions that reflect negative sentiments from the public:

Firstly, in cluster number 10 there is public dissatisfaction regarding train doors that are too short, which causes problems that cannot be expressed by the parties involved. This shows the need for better infrastructure improvements to improve user experience. This interpretation shows the dimension of service quality, namely tangibility.

Secondly, in cluster numbers 9 and 10, users felt frustrated with headway irregularities, especially at Dukuh station, which resulted in unpredictable waiting times. Incomplete schedules and uncertain travel times are also major concerns, highlighting the need for operational improvements and schedule regularity. This reflects the dimension of service quality, namely reliability.

Third, in cluster numbers 10 and 11 reflect the responsiveness dimension regarding criticism of officers to pay more attention to their duties and passenger safety, including priority seats and facilities on trains, showing the importance of improving service quality and security. In addition, in cluster number 10, passengers expect additional services, especially during night hours in Pancoran, indicating expectations for improved service quality at certain times.

Fourthly, cluster numbers 10 and 12 feel disappointed with tariffs which are considered expensive and the quality of service, which is considered poor, the public feels it is not appropriate regarding the payments that must be made. Apart from that, in cluster numbers 9 and 11, the public expressed concerns about infrastructure and felt a sense of fear or anxiety. These representations reflect an increase in service quality in the assurance dimension.

## V. CONCLUSIONS

In this research, data collection utilized web scraping via APIFY to extract public comments from social media platforms such as Twitter, Instagram, and Youtube, resulting in a dataset of 8,523 comments. Preprocessing was then applied to enhance the data's suitability for algorithmic application. Following preprocessing, the InSet Lexicon dictionary was employed to determine sentiment polarity related to the services of LRT Jabodebek. The labelled data was subsequently applied to IndoBERT, exhibiting robust classification performance with an accuracy score of 85.97%, F1 score of 84.13, Recall of 82.99%, and Precision of 85.31%. BERTopic topic modelling was applied to labelled data, revealing in the topic of modeling positive sentiment, interpreting users showing appreciation for informative content from LRT Jabodebek social media accounts, highlighting a good response to public information needs. Apart from that, there is appreciation for the contribution of the LRT Jabodebek in improving the quality of life of the community through high-quality transportation infrastructure, as well as strong support from stakeholders to maintain and improve service standards. In addition, Railmin, the LRT Jabodebek social media admin, showed an empathetic response to users through apologies and efforts to understand the user's situation. This shows the service quality dimensions of responsiveness, tangibility, reliability, assurance and empathy, parties related to the LRT Jabodebek can know the aspects that need to be maintained and improved regarding service quality. On the other hand, negative sentiments, shows several main aspects related to

the LRT Jabodebek service quality dimensions. First, there is a need to improve the infrastructure regarding train doors that are too short to improve user experience, highlighting the dimension of tangibility. Second, schedule uncertainty and irregular lead times indicate the need to improve operational and schedule reliability, reflecting the reliability dimension. Third, criticism of officers' duties and expectations for additional services indicate the importance of responsiveness and increased security, reflecting the responsiveness dimension. Finally, disappointment with high rates that do not correspond to the low quality of services provided indicates a need for improvement in the assurance dimension, while concerns about infrastructure reflect improvements in quality in the assurance dimension. By improving these aspects, the Jabodebek LRT service can improve the overall quality of service. This model holds the potential to be a valuable tool for various stakeholders, aiding LRT Jabodebek operators in enhancing their services based on public opinions on social media.

However, the research faces limitations, including potential inaccuracies in sentiment labelling using the InSet Lexicon method, especially when a word has multiple meanings. The next limitation is the data collection process, which was limited to a period of 2 months, thus insufficiently covering the variability of events and changes that occurred over time. Another limitation is related to the diverse nature of the dataset collected from various social media platforms, posing challenges for BERTopic to fully reflect the variations in opinions and perspectives. Future research could overcome these limitations by using more accurate lexicons or manual labelling with more accurate linguists and increasing training iterations to improve the accuracy of the classification and data scrapping process over longer time periods to cover a greater variety of data. Additionally, future research could explore developing the model with different case studies.

## REFERENCES

- [1] Monavia Ayu Rizaty, "Sebanyak 56,7% Penduduk Indonesia Tinggal di Perkotaan pada 2020," [databoks.katadata.co.id](http://databoks.katadata.co.id).
- [2] G. R. Pangaribuan and D. D. Purba, "The impact of LRT jabodebek in enforcing capability of the intercity transportation network in the greater Jakarta area," *Int J Adv Sci Eng Inf Technol*, vol. 10, pp. 828–836, 2020.
- [3] A. Madeppungeng, R. Said Rasul, and R. Irwanto, "RISK MANAGEMENT IN THE JABODEBEK LRT (LIGHT RAIL TRANSIT) DEVELOPMENT PROJECT, IN JAKARTA AND ITS NEIGHBOURING CITIES," 2022.
- [4] Railway Technology, "Indonesia's Jabodebek LRT to commence operations in mid-2023," [railway-technology.com](http://railway-technology.com).
- [5] F. Zhang, T. Song, X. Cheng, T. Li, and Z. Yang, "Transportation Infrastructure, Population Mobility, and Public Health," *Int J Environ Res Public Health*, vol. 20, no. 1, Jan. 2023, doi: 10.3390/ijerph20010751.
- [6] W. P. Pramudita and A. D. Nataadmadja, "Analysis of the performance of light rail transit (LRT) Jakarta as a transport demand management (TDM) strategy," in *IOP Conference Series: Earth and Environmental Science*, Institute of Physics, 2023. doi: 10.1088/1755-1315/1169/1/012021.
- [7] F. F. Rachman, R. Nooraeni, and L. Yuliana, "Public Opinion of Transportation integrated (Jak Lingko), in DKI Jakarta, Indonesia," in *Procedia Computer Science*, Elsevier B.V., 2021, pp. 696–703. doi: 10.1016/j.procs.2021.01.057.
- [8] Rio Adryawan, "Penyebab Utama LRT Jabodebek Sepi Penumpang," [economy.okezone.com](http://economy.okezone.com).
- [9] C. Sanchis-Pedregosa, J. A. D. Machuca, and M. D. M. González-Zamora, "Determinants of success in transport services outsourcing: Empirical study in Europe," in *International Journal of Logistics Management*, Emerald Group Publishing Ltd., 2018, pp. 261–283. doi: 10.1108/IJLM-09-2016-0207.
- [10] A. N. H. Ibrahim, M. N. Borhan, M. H. Osman, M. R. Mat Yazid, and M. Md. Rohani, "The Influence of Service Quality on User's Perceived Satisfaction with Light Rail Transit Service in Klang Valley, Malaysia," *Mathematics*, vol. 10, no. 13, Jul. 2022, doi: 10.3390/math10132213.
- [11] M. O. Pratama *et al.*, "The sentiment analysis of Indonesia commuter line using machine learning based on twitter data," in *Journal of Physics: Conference Series*, Institute of Physics Publishing, Apr. 2019. doi: 10.1088/1742-6596/1193/1/012029.
- [12] M. Bhardwaj, P. Mishra, S. Badhani, and S. K. Muttoo, "Sentiment analysis and topic modeling of COVID-19 tweets of India," *International Journal of System Assurance Engineering and Management*, 2023, doi: 10.1007/s13198-023-02082-0.

- [13] A. Alamsyah, W. Rahmah, and H. Irawan, "SENTIMENT ANALYSIS BASED ON APPRAISAL THEORY FOR MARKETING INTELLIGENCE IN INDONESIA'S MOBILE PHONE MARKET," *J Theor Appl Inf Technol*, vol. 82, no. 2, 2015, [Online]. Available: [www.jatit.org](http://www.jatit.org)
- [14] G. K. Wadhvani, P. K. Varshney, A. Gupta, and S. Kumar, "Sentiment Analysis and Comprehensive Evaluation of Supervised Machine Learning Models Using Twitter Data on Russia–Ukraine War," *SN Comput Sci*, vol. 4, no. 4, Jul. 2023, doi: 10.1007/s42979-023-01790-5.
- [15] F. Koto and G. Y. Rahmaningtyas, "Inset lexicon: Evaluation of a word list for Indonesian sentiment analysis in microblogs," in *Proceedings of the 2017 International Conference on Asian Language Processing, IALP 2017*, Institute of Electrical and Electronics Engineers Inc., Jul. 2017, pp. 391–394. doi: 10.1109/IALP.2017.8300625.
- [16] H. Imaduddin, F. Yusfida A'la, and Y. S. Nugroho, "Sentiment Analysis in Indonesian Healthcare Applications using IndoBERT Approach," 2023. [Online]. Available: [www.ijacsa.thesai.org](http://www.ijacsa.thesai.org)
- [17] L. B. Hutama and D. Suhartono, "Indonesian Hoax News Classification with Multilingual Transformer Model and BERTopic," *Informatica (Slovenia)*, vol. 46, no. 8, pp. 81–90, 2022, doi: 10.31449/inf.v46i8.4336.
- [18] Fransiscus and A. S. Girsang, "Sentiment Analysis of COVID-19 Public Activity Restriction (PPKM) Impact using BERT Method," *International Journal of Engineering Trends and Technology*, vol. 70, no. 12, pp. 281–288, Dec. 2022, doi: 10.14445/22315381/IJETT-V70I12P226.
- [19] J. Zhang, "A Combination of Lexicon-Based and Classified-Based methods for Sentiment Classification based on Bert," in *IOP Conference Series: Earth and Environmental Science*, IOP Publishing Ltd, Mar. 2021. doi: 10.1088/1742-6596/1802/3/032113.
- [20] Y. Liu, J. Shi, C. Zhao, and C. Zhang, "Generalizing factors of COVID-19 vaccine attitudes in different regions: A summary generation and topic modeling approach," *Digit Health*, vol. 9, Jan. 2023, doi: 10.1177/20552076231188852.
- [21] M. Uncovska, B. Freitag, S. Meister, and L. Fehring, "Rating analysis and BERTopic modeling of consumer versus regulated mHealth app reviews in Germany," *NPJ Digit Med*, vol. 6, no. 1, Dec. 2023, doi: 10.1038/s41746-023-00862-3.
- [22] H. Irawan, G. Akmalia, and R. A. Masrury, "Mining tourist's perception toward Indonesia tourism destination using sentiment analysis and topic modelling," in *ACM International Conference Proceeding Series*, Association for Computing Machinery, Sep. 2019, pp. 7–12. doi: 10.1145/3361821.3361829.
- [23] A. W. K. Yeung *et al.*, "Are dental x-rays safe? Content analysis of English and Chinese YouTube videos," *Digit Health*, vol. 9, Jan. 2023, doi: 10.1177/20552076231179053.
- [24] V. W. Fitri, E. Cahyanti, and Q. G. Adiwijaya, "On The Feature Extraction For Sentiment Analysis of Movie Reviews Based on SVM," *2020 8th International Conference on Information and Communication Technology (ICoICT)*, 2020.
- [25] A. Muhariya, I. Riadi, Y. Prayudi, and I. A. Saputro, "Utilizing K-means Clustering for the Detection of Cyberbullying Within Instagram Comments," *Ingenierie des Systemes d'Information*, vol. 28, no. 4, pp. 939–949, Aug. 2023, doi: 10.18280/isi.280414.
- [26] C. Zong, R. Xia, and J. Zhang, *Text Data Mining*. Singapore: Springer Singapore, 2021. doi: 10.1007/978-981-16-0100-2.
- [27] A. Marpaung, R. Rismala, and H. Nurrahmi, "Hate Speech Detection in Indonesian Twitter Texts using Bidirectional Gated Recurrent Unit," in *KST 2021 - 2021 13th International Conference Knowledge and Smart Technology*, Institute of Electrical and Electronics Engineers Inc., Jan. 2021, pp. 186–190. doi: 10.1109/KST51265.2021.9415760.
- [28] B. V. Kartika, M. J. Alfredo, and G. P. Kusuma, "Fine-Tuned IndoBERT based model and data augmentation for Indonesian language paraphrase identification," *Revue d'Intelligence Artificielle*, vol. 37, no. 3, pp. 733–743, Jun. 2023, doi: 10.18280/ria.370322.