

TOPIC ANALYSIS OF TIMOR LESTE NEWS IN INDONESIA'S MASS MEDIA USING LATENT DIRICHLET ALLOCATION

**Expendito Pinto Ximenes, Luh Gede Surya Kartika, Eka Darmayanthi, Dora
Fanny Meidodga, Ni Kadek Ari Kesuma**

1(Dili Institute Of Technology, Timor Leste, expeditopinto27@gmail.com)

2(Informatics Progame, Hindu State University I Gusti Bagus Sugriwa Denpasar,
Indonesia, suryakartika@uhnsugriwa.ac.id)

3(Biro Academic, Hindu State University I Gusti Bagus Sugriwa Denpasar, Indonesia,
ekadarmayanthi32@gmail.com)

4(Informatics Progame, Hindu State University I Gusti Bagus Sugriwa Denpasar,
Indonesia, dorafanny06@gmail.com)

5(Informatics Progame, Hindu State University I Gusti Bagus Sugriwa Denpasar,
Indonesia, arikesumadevi@gmail.com)

ABSTRACT

This study analyzes news about Timor Leste in the Indonesian national mass media Detik.com using the Latent Dirichlet Allocation (LDA) method. This research is important to understand how the national mass media in Indonesia, especially Detik.com, report on Timor Leste. Timor-Leste is a neighboring country with a significant history and diplomatic relations with Indonesia. There is a gap between the existing literature, which often focuses on the political and historical links between the two countries, and the reality of everyday reporting covering various aspects of life in Timor-Leste. The aim of this study was to identify and analyse key topics in Timor-Leste reporting in Detik.com, in order to provide a more comprehensive picture of how Timor-Leste was presented to the Indonesian public. The methods used in this study include text mining, text preprocessing, and Latent Dirichlet Allocation (LDA). Using the LDA, we identified two main topics that emerged from articles discussing Timor Leste. The first topic (30.3% of tokens) focused on Timor-Leste's membership in ASEAN and relations with Indonesia, with keywords such as "asean", "leste", "timor", "year", "member", and "indonesia". The second topic (14.7% of tokens) revolved around local and sporting issues in Timor-Leste, with keywords such as "timor", "indonesia", "leste", "dili", "language", "minutes", "year", "goals", "east", and "second". The other eight topics deal with similar terms. The results showed that news about Timor-Leste in Detik.com covered aspects of politics and international relations as well as local developments and sports in the country. This research provides insight into how national media construct and convey information about Timor-Leste to the Indonesian public.

Keywords: Timor Leste, news, Latent Dirichlet Allocation, Detik.com

INTRODUCTION

The background of the research on the analysis of news topics about Timor Leste in the Indonesian National Mass Media is very relevant and important in the context of understanding the dynamics of bilateral relations between Indonesia and Timor Leste and their impact on public perception in Indonesia. Since Timor-Leste gained its independence in 2002, relations between the two countries have experienced significant developments in various fields, including political, economic, and social. Although the two countries have complicated histories, particularly with regard to East Timor's integration into Indonesia and the separation process that followed Timor-Leste's independence, their relations have improved in recent years. Indonesia has played an active role in supporting Timor Leste's development, especially in the fields of education, infrastructure, health, and agriculture. These cooperation programs

aim to help strengthen institutional capacity, improve social welfare, and accelerate development in Timor Leste (Ministry of Foreign Affairs of the Republic of Indonesia, 2024).

Although still relatively limited, economic cooperation between Indonesia and Timor-Leste has grown in recent years. These include Indonesian investment in Timor-Leste, bilateral trade, and efforts to enhance cooperation in sectors such as agriculture, energy, tourism, and trade (Taena et al., 2022; Taena & Afoan, 2020). Indonesia has provided support in strengthening Timor-Leste's security capacity, including through security personnel training, information exchange, and cooperation in law enforcement and border security (Vivi Pusvitasary, 2024). It aims to help East Timor strengthen stability and security in the country.

In the field of Education and culture, cultural and educational exchanges between Indonesia and Timor Leste have become an important part of bilateral relations (Ministry of Foreign Affairs of the Republic of Indonesia, 2023). Student, cultural, and academic exchange programs have helped strengthen relations between the two countries and enrich the experiences of people on both sides.

Although there are still some border-related disputes between the two countries, efforts have been made to resolve the issue peacefully and based on international law. Timor-Leste and Indonesia have been engaged in dialogue and negotiations to reach a mutually beneficial agreement on borders and border security (Vivi Pusvitasary, 2024). Thus, although there are still some challenges and issues that need to be resolved, bilateral relations between Indonesia and Timor-Leste have experienced positive developments in recent years, focusing on development, economic cooperation, security, and cultural exchanges.

Mass media has an important role as a conveyor of information to the public. Reporting about Timor Leste in Indonesia's national mass media not only reflects the dynamics of bilateral relations between the two countries, but can also influence the perception and views of the Indonesian people towards Timor Leste and issues related to it. However, in reality, mass media reporting is often influenced by various factors, such as political, economic, and social interests (Pavelka, 2014). Therefore, to understand more deeply how Timor Leste is presented in Indonesia's national mass media coverage, a comprehensive analysis of the dominant topics in news coverage about Timor Leste is needed.

Through the analysis of news topics about Timor Leste in the Indonesian national mass media, various issues and themes that dominate the news can be identified, as well as narrative patterns that may affect public perception. This research can not only provide a deeper understanding of how Timor-Leste is presented in Indonesia's national mass media, but it can also be a basis for evaluating the extent to which mass media coverage reflects the reality of bilateral relations between the two countries and promotes better understanding between the two peoples.

Topic analysis research is a branch of text analysis that aims to identify and group the main themes or topics in a set of documents. One method that is often used for topic analysis is Latent Dirichlet Allocation (LDA). LDA is a generative model that assumes that a document is a mixture of several topics and that each topic is a distribution of words (Blei, 2011; Iparraguirre-Villanueva et al., 2023; Sutherland et al., 2020).

METHODS

In the context of this study, data will be taken from detik.com with the keyword "Timor Leste". Here are the steps in the implementation of this research:

The first step is news source determination. The initial stage is to select the news source to be taken for analysis. This research used Detik.com national mass media website as a source of news. The Indonesian national mass media website Detik.com chosen because it has HTML markers and relatively consistent elements on each news page. This makes it easier to do the scraping stages. In addition, according to (Annur, 2023; Cindy Mutia Annur, 2022), Detik.com

is one of the mass media that has the highest level of trust in Indonesia. So that data collection from Detik.com is considered feasible in this study.

The next step is the selection of scraping techniques. Scraping techniques are used to collect data from specified news sources (Krishna et al., 2022; Matta et al., 2020; Speckmann, 2021). It involves the use of a Python 3-language computer program to automatically extract text from the web pages of detik.com national mass media. The scraping technique used is to extract HTML elements (Hypertext Markup Language) that are relevant to the text taken. In this study, scraping of information: news title, news content, author, news date, and URL (Uniform Resource Locator) of the news.

Data Cleansing: Once the data is successfully retrieved, the next step is to clean the data. This involves removing irrelevant elements such as HTML tags, ads, or user comments. Data may also need to be normalized, for example, changing the date format or simplifying text formatting.

Tokenization and Text Processing. The cleaned text is then broken down into smaller units such as words or phrases. This process is called tokenization. Furthermore, the text can be further processed by performing steps such as the removal of stopwords (common words that do not give meaning), stemming (changing words into their basic form), or the removal of special characters. In addition to the use of the Stopwords Library in the Python Programming Language, some of the words added in Stopwords are:

```
['also', 'and', 'with', 'for', 'on', 'will', 'of', 'that', 'which',  
'content', 'continue', 'advertisementscroll to continue with  
content', 'say', 'do', 'become', 'word', 'by'].
```

The choice of words to omitted from the dataset is very important because it will affect the results. The more words added in Stopwords, the fewer final tokens earned.

Text Analysis: Once the text data is processed, various text analysis techniques can be applied. This study used the Latent Dirichlet Allocation (LDA) model for text analysis. LDA makes it possible to identify topic patterns hidden within a corpus of text. LDA stands for Latent Dirichlet Allocation. This is a probabilistic model used to perform topic analysis on collections of text documents. LDA is one of the most commonly used natural language processing (NLP) methods for finding hidden or latent topics in a corpus of text (Blei, Ng, & Edu, 2003; Blei, Ng, & Jordan, 2003; Garg & Rangra, 2022).

The LDA model call in Python 3 is to use the following snippet of program code:

```
from pprint import pprint  
num_topics = 10  
lda_model=gensim.models.LdaMulticore(corpus=corpus, id2word=id2word,  
num_topics=num_topics)  
pprint(lda_model.print_topics())  
doc_lda = lda_model[corpus]
```

Using the above program code, the LDA model is asked to create as many as 10 topics from the existing corpus. In the LDA model, it is assumed that each document in the corpus is produced by a number of hidden topics, and each word in the document is associated with one or more of those topics. The purpose of LDA is to find the distribution of topics within the corpus and the distribution of words within each topic.

LDA works in a simple yet powerful way. It assumes that there are a number of topics in the corpus, and then assumes the probability distribution of the words in each topic as well as the probability distribution of the topics in each document. The model is then adjusted so as to produce the most likely distribution of topics within the corpus (Blei, Ng, & Edu, 2003). The

in news stories about Timor Leste. In LDA, "salient terms" refer to the most prominent or most representative words of each topic found by the model. In the process of LDA training, each topic is represented by a probability distribution of words, and words with a higher probability in that distribution are considered the most typical or most instrumental words in the topic. Therefore, "salient terms" are the words that appear most often or best represent the topic in topic analysis using LDA (Lotto et al., 2023). The salient term bar chart by default displays the 30 most prominent terms. Figure 2 shows that the words that appear look varied and relevant, it shows that the model has managed to find a meaningful topic in the corpus of text. By looking at the words that appear most often in each topic, we can find interesting patterns or insights in the corpus of text.

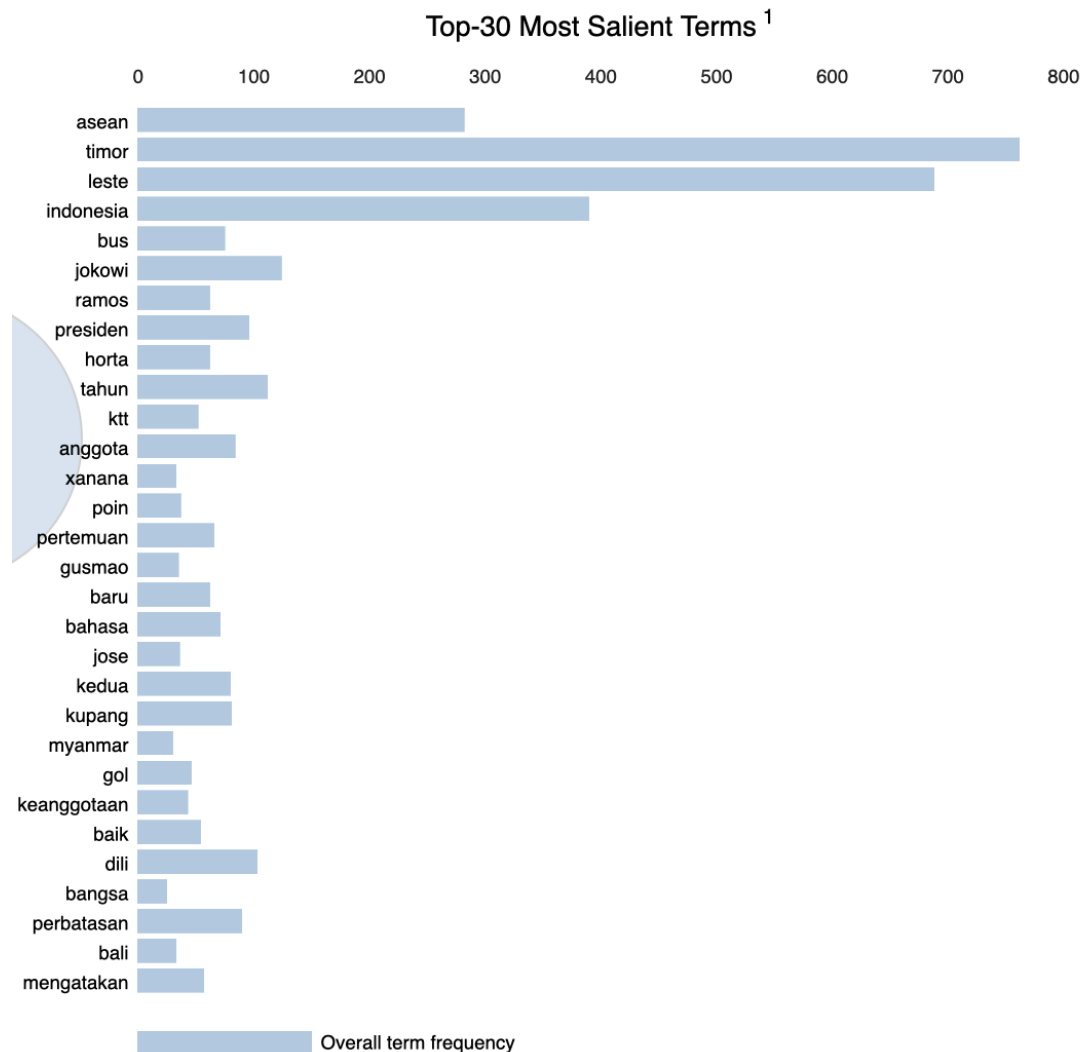


Figure 2. Top-30 Most Salient Terms

Figure 2 shows that the terms that appear the most are "Timor", "Leste", "Indonesia", "ASEAN", and "Jokowi". Of the 10 topics formed, the five most salient terms appear on most topics. This means that there are several reports that mention these words. Table 1 shows the topics and probability of each terms and the percentage of those topics against the overall token. Table 1 shows the most relevant or representative words of each topic in can be used to help understand what topics each group of documents or topics in the model represent. By looking at the most dominant words in each topic, we can identify the core or main theme of the topic.

Table 1. Probability of each term in each topic and percentage of token

Topics	Probability of each term	% of tokens
Topik 1	0.020*"asean" + 0.017*"leste" + 0.016*"timor" + 0.007*"tahun" + 0.006*"anggota" + 0.006*"indonesia" + 0.005*"baru" + 0.004*"dili" + 0.004*"ktt" + 0.004*"bus"	30,3% of Tokens
Topik 2	0.030*"timor" + 0.023*"indonesia" + 0.020*"leste" + 0.007*"dili" + 0.006*"bahasa" + 0.004*"menit" + 0.004*"tahun" + 0.004*"gol" + 0.004*"timur" + 0.004*"kedua"	14,7% of tokens
Topik 3	0.023*"timor" + 0.021*"leste" + 0.015*"asean" + 0.013*"indonesia" + 0.006*"jokowi" + 0.005*"xanana" + 0.005*"gusmao" + 0.005*"presiden" + 0.005*"ktt" + 0.004*"pertemuan"	10,7% of tokens
Topik 4	0.031*"timor" + 0.022*"leste" + 0.014*"indonesia" + 0.009*"asean" + 0.005*"anggota" + 0.004*"dua" + 0.004*"ramos" + 0.003*"horta" + 0.003*"gol" + 0.003*"laga"	9,8% of tokens
Topik 5	0.024*"leste" + 0.023*"timor" + 0.014*"asean" + 0.008*"ramos" + 0.008*"horta" + 0.008*"presiden" + 0.006*"indonesia" + 0.006*"jokowi" + 0.004*"baik" + 0.004*"kedua"	7,8% of tokens
Topik 6	0.037*"timor" + 0.034*"leste" + 0.018*"indonesia" + 0.008*"jokowi" + 0.007*"perbatasan" + 0.007*"asean" + 0.005*"bus" + 0.005*"dili" + 0.005*"kupang" + 0.004*"tahun"	6,4% of tokens
Topik 7	0.025*"leste" + 0.025*"timor" + 0.024*"asean" + 0.012*"indonesia" + 0.005*"anggota" + 0.005*"ktt" + 0.004*"poin" + 0.004*"presiden" + 0.003*"tenggara" + 0.003*"tahun"	5,8% of tokens
Topik 8	0.020*"leste" + 0.019*"timor" + 0.013*"asean" + 0.009*"tahun" + 0.007*"indonesia" + 0.006*"presiden" + 0.006*"jokowi" + 0.005*"ramos" + 0.005*"horta" + 0.005*"ktt"	5,4% of tokens
Topik 9	0.027*"timor" + 0.023*"leste" + 0.008*"asean" + 0.008*"bus" + 0.007*"indonesia" + 0.006*"kupang" + 0.004*"perbatasan" + 0.004*"bahasa" + 0.004*"kedua" + 0.003*"dili"	5,4% of tokens
Topik 10	0.030*"leste" + 0.029*"timor" + 0.014*"indonesia" + 0.012*"asean" + 0.006*"presiden" + 0.005*"jokowi" + 0.004*"baru" + 0.004*"tahun" + 0.004*"senjata" + 0.004*"poin"	3,7% of tokens

Topic 1 has 30.3% of Tokens. A token is the smallest unit in natural language processing. In the context of text processing, tokens usually refer to individual words or punctuation marks that are considered separate units. The tokenization process is the process of dividing text into these tokens. For example, in the sentence "Citizens on the border of

Indonesia and Timor Leste, reveal..", the tokens are "citizens", "in", "borders", "Indonesia", "and", "Timor", "Leste", "disclose" and punctuation marks. Thus, tokens refer to separate units in text used for language analysis. So 30.3% of Tokens refers to the proportion of those tokens in the corpus of text attributed to Topic 1.

Looking at the terms in Topic 1, namely "asean", "leste", "timor", "year", "member", "Indonesia", "new", "dili", "summit", and "bus". It can be seen that topic 1 contains news about Timor Leste's membership in ASEAN (The Association of Southeast Asian Nations). If you look at the news in Detik.com, this topic appeared in more than 15 news stories starting from 2022. Topic 1 can also be related to reporting on the efforts of Indonesia and Timor Leste in encouraging the implementation of the five-point consensus at the ASEAN Laos Summit. News of the ASEAN Laos Summit consensus will begin in January 2024.

The next topic is Topic 2, which contains the terms "Timor", "Indonesia", "Leste", "Dili", "language", "minutes", "year", "goal", "east", "second". Judging from the composition of the term, topic 2 contains news about the Indonesian national team versus Timor Leste football match. A lot of news about this is contained in 2022-2023. Generally, the news on topic 2 is about the course of the 2023 U-20 match and the AFF U-23 Cup. The news conveyed was generally about the results obtained in the match, as well as the process of football matches that had occurred between Indonesia and Timor Leste.

Topic 3 contains the keywords "timor", "leste", "asean", "indonesia", "jokowi", "xanana", "gusmao", "president", "summit", and "meeting". Topic 3 is news related to the ASEAN Summit meeting. One of the news with this theme is the arrival of Xanana Gusmao to the ASEAN Summit in Jakarta in 2023, or about the actions, policies, and decisions of the Prime Minister of Timor Leste during the ASEAN Summit. The news that also contains the term in Topic 3 is about discussions that took place between the governments of Indonesia and Timor Leste regarding the right to attend all ASEAN meetings and summit sessions.

Topic 4 contains the words "Timor", "Leste", "Indonesia", "ASEAN", "Member", "Dua", "Ramos", "Horta", "Goal", "Match". The interesting thing here, is that manually there are no reports that contain the words "ramos" + "horta" and "goal" + "game" in the same report. This requires a more in-depth study of the evaluation of the Latent Dirichlet Allocation model. So that quality results can be obtained. Some possible model evaluations that can be considered for use are (Comotto, 2022):

Log Loss: can be used to measure the quality of probability predictions from models. It provides a measurement of how well the model estimates the correct probability for each data instance. The lower the log loss, the better the model, or;

Brier Score: can be used for general measurement to evaluate the quality of probability predictions from models. It measures the mean squared of the difference between the predicted probability and the actual probability, or;

Calibration Curve: This is a curve that plots the probability predicted by the model against the actual frequency of a positive class. A curve closer to the diagonal line indicates better calibration.

Figure 3 shows the Intertopic Distance Map. This image is a visualization of the topic in two-dimensional space. The area of this topic circle is proportional to the number of words included in each topic in the dictionary. The circles are plotted using a multidimensional scaling algorithm based on the words contained in them, so that topics closer to each other have more words in common.



Figure 3 Intertopic Distance Map

Figure 3 shows that there are two intersecting topics, namely topic 5 and Topic 3. When two or more topics overlap on an intertopic distance map, this can indicate several things:

First, there are common themes between these topics. An overlap between Topic 3 and Topic 5 could indicate that there are several themes or concepts that are interrelated between the two topics. Topic 3 has the theme of the ASEAN Summit, while Topic 5 has the theme of Ramos Horta (President of Timor Leste). Both are possible to be united in the news about Ramos Horta and the ASEAN summit.

Secondly, the lack of a clear separation. Overlap in the intertopic distance map may indicate that the model has difficulty clearly separating different topics. This may happen if the corpus of text has a complex structure or if the topics are indeed very similar. It also highlights the importance of advanced analysis to understand the relationship between such topics. This could involve further exploration of key words representing overlapping topics and looking at larger contexts or patterns within the corpus of text. In addition, further research proposals using qualitative research models such as Framing Analysis. So it can be known whether the news tends to be monotonous or conveys similar information continuously.

Third, LDA models need further purification. Overlap between topics may indicate that the LDA model needs further refinement. This could involve adjusting model parameters, such as the number of topics or more sophisticated text preprocessing techniques, to improve separation between existing topics.

CONCLUSION

Research on topic analysis with the keyword "Timor Leste" in the Mass Media Detik.com using LDA has produced relevant topics. The main advantage of LDA is its ability to handle large text data sets without the need for manual annotations. The model is also unsupervised, so it can be used on unlabeled data. However, LDA also has challenges. One is to determine the optimal number of topics, as too many or too few topics can result in less accurate interpretations. In addition, the interpretation of the resulting topic can also be subjective and requires domain expertise to ensure the relevance and accuracy of the topic.

The most dominant topics appearing in news about Timor Leste in the national mass media Detik.com in Indonesia are in terms of the ASEAN Summit and the field of sports. Both of these appear on more than one topic with a term count of 30.5% of tokens in the data. The results showed that news about Timor-Leste in Detik.com covered aspects of politics and international relations as well as local developments and sports in the country. This research provides insight into how national media construct and convey information about Timor-Leste to the Indonesian public. For further research, the development of more sophisticated models and hybrid approaches that combine LDA with other techniques can provide more accurate and interpretable results, thus expanding the scope and application of topic analysis.

REFERENCES

- Annur, C. M. (2023). *Inilah Media yang Paling Dipercaya Warga Indonesia pada 2023, Ada Favoritmu?* <https://Databoks.Katadata.Co.Id/Datapublish/2023/06/15/Inilah-Media-Yang-Paling-Dipercaya-Warga-Indonesia-Pada-2023-Ada-Favoritmu>.
- Blei, D. M. (2011). *Introduction to Probabilistic Topic Models*. https://oar.princeton.edu/jspui/bitstream/88435/pr1bv3w/1/OA_IntroductionProbabilisticTopicModels.pdf
- Blei, D. M., Ng, A. Y., & Edu, J. B. (2003). Latent Dirichlet Allocation Michael I. Jordan. In *Journal of Machine Learning Research* (Vol. 3).
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3(4–5). <https://doi.org/10.7551/mitpress/1120.003.0082>
- Cindy Mutia Annur. (2022, June 6). *Tingkat Kepercayaan Responden terhadap Sejumlah Merek Media*. Databoks.
- Comotto, F. (2022). *Evaluation metrics: leave your comfort zone and try MCC and Brier Score*. <https://Towardsdatascience.Com/Evaluation-Metrics-Leave-Your-Comfort-Zone-and-Try-Mcc-and-Brier-Score-86307fb1236a#:~:Text=Brier%20Score%20is%20similar%20to,T%20have%20an%20upper%20bound>.
- Garg, M., & Rangra, P. (2022). Bibliometric Analysis of Latent Dirichlet Allocation. *DESIDOC Journal of Library and Information Technology*, 42(2). <https://doi.org/10.14429/djlit.42.2.17307>
- Iparraquirre-Villanueva, O., Sierra-Liñan, F., Salazar, J. L. H., Beltozar-Clemente, S., Pucuhuayla-Revatta, F., Zapata-Paulini, J., & Cabanillas-Carbonell, M. (2023). Search and classify topics in a corpus of text using the latent dirichlet allocation model. *Indonesian Journal of Electrical Engineering and Computer Science*, 30(1). <https://doi.org/10.11591/ijeecs.v30.i1.pp246-256>
- Krishna, N., Nayak, A., Badagan, S., & Jetty, C. (2022). A study on Web Scraping. *International Journal of Engineering Research in Computer Science and Engineering (IJERCSE)*, 9.

- Lotto, M., Hussain, I. Z., Kaur, J., Butt, Z. A., Cruvinel, T., & Morita, P. P. (2023). Analysis of Fluoride-Free Content on Twitter: Topic Modeling Study. *Journal of Medical Internet Research*, 25. <https://doi.org/10.2196/44586>
- Matta, P., Sharma, N., Sharma, D., Pant, B., & Sharma, S. (2020). Web Scraping: Applications and Scraping Tools. *International Journal of Advanced Trends in Computer Science and Engineering*, 9(5). <https://doi.org/10.30534/ijatcse/2020/185952020>
- Ministry of Foreign Affairs of the Republic of Indonesia. (2023). *Indonesia-Timor Leste Partnership: Scholarships for Outstanding Students*. <https://Kemlu.Go.Id/Portal/En/Read/5561/View/Indonesia-Timor-Leste-Partnership-Scholarships-for-Outstanding-Students>.
- Ministry of Foreign Affairs of the Republic of Indonesia. (2024). *President of Indonesia Receives Timor Leste PM Visit, Agree on Multiple Cooperation Agreements*. <https://Kemlu.Go.Id/Portal/En/Read/5708/Berita/President-of-Indonesia-Receives-Timor-Leste-Pm-Visit-Agree-on-Multiple-Cooperation-Agreements>.
- Pavelka, J. (2014). The Factors Affecting the Presentation of Events and the Media Coverage of Topics in the Mass Media. *Procedia - Social and Behavioral Sciences*, 140, 623–629. <https://doi.org/10.1016/j.sbspro.2014.04.482>
- Speckmann, F. (2021). Web Scraping. *Zeitschrift Für Psychologie*, 229(4). <https://doi.org/10.1027/2151-2604/a000470>
- Sutherland, I., Sim, Y., Lee, S. K., Byun, J., & Kiatkawsin, K. (2020). Topic modeling of online accommodation reviews via latent dirichlet allocation. *Sustainability (Switzerland)*, 12(5). <https://doi.org/10.3390/su12051821>
- Taena, W., & Afoan, F. (2020). Cross Border Tourism and Regional Development: Case Indonesia-Timor Leste Cross Border. *Ekulilibrium : Jurnal Ilmiah Bidang Ilmu Ekonomi*, 15(1). <https://doi.org/10.24269/ekulilibrium.v15i1.2330>
- Taena, W., Kase, M. S., & Afoan, F. (2022). The Externality and Sustainable Development Priority of Cross Border Tourism. *MIMBAR : Jurnal Sosial Dan Pembangunan*. <https://doi.org/10.29313/mimbar.v0i0.9422>
- Vivi Pusvitasary. (2024). Indonesian Diplomacy under the Leadership of Joko Widodo in Resolving Land Border Disputes with Timor-Leste. *Budi Luhur Journal of Strategic & Global Studies*, 2(1). <https://doi.org/10.36080/jsgs.v2i1.33>