

LDA Topic Modeling Analysis of Public Discourse on Indonesia's Free Nutritious Meals Program (MBG)

Cici Suhaeni*

Department of Statistics and Data Science
IPB University
Bogor, Indonesia
cici_suhaeni@apps.ipb.ac.id

Laily Nissa Atul Mualifah

Department of Statistics and Data Science
IPB University
Bogor, Indonesia
lailyatul@apps.ipb.ac.id

Hari Wijayanto

Department of Statistics and Data Science
IPB University
Bogor, Indonesia
hari@apps.ipb.ac.id

Article History

Received May 14th, 2025

Revised June 27th, 2025

Accepted June 29th, 2025

Published June, 2025

Abstract— This study investigates public discourse on Indonesia's Free Nutritious Meals (Makan Bergizi Gratis/MBG) program through Latent Dirichlet Allocation (LDA) topic modeling of YouTube comments. Filling a research gap on online public opinion regarding the MBG policy, this study identifies dominant themes and discursive patterns in public perception. A three-topic model, validated through coherence score evaluation and pyLDAvis visualization, reveals key topics: concerns over food prices and distribution, perceived benefits for children and society, and emotionally and politically driven reactions. The findings provide valuable insights into public opinion, while also highlighting challenges in processing Indonesian-language text, such as informal language and noisy data. This study contributes to understanding public perceptions of social policies in digital environments and recommends future research directions, including improved text preprocessing and alternative topic modeling approaches. By shedding light on online public discourse, this research informs policymakers and stakeholders about the effectiveness and potential areas for improvement in the MBG program.

Keywords— coherence score evaluation; discursive patterns; dominant themes; Makan Bergizi Gratis (MBG); text analytics

1 INTRODUCTION

In the era of digital transformation, Big Data is generated from diverse sources, and one of the most dominant forms is textual data, especially from social media platforms. Social media has transformed into a significant arena where public opinions are formed and expressed, with YouTube standing out as a rich platform for real-time public discourse. User-generated content on YouTube, especially through the comments section, offers an authentic lens into societal reactions, sentiments, and conversations surrounding pressing national issues [1].

The richness and volume of conversations on platforms like YouTube present both opportunities and challenges for researchers seeking to understand public sentiment and discourse. Unlike structured survey responses or official statements, user-generated content reflects spontaneous, diverse, and dynamic expressions of opinion. These characteristics make social media data a valuable yet complex source of information that requires advanced analytical methods to decode. To transform this unstructured textual data into meaningful insights, computational approaches such as Natural Language Processing (NLP) are essential.

Text data on such platforms can be analyzed using NLP techniques to uncover patterns and themes. Among the most effective and widely used techniques is topic modeling, which helps extract latent themes from a large corpus of text [2], [3]. Topic modeling is recognized as an effective analytical approach for uncovering latent patterns in textual data and is widely adopted across various disciplines because of its flexible and scalable nature [4], [5].

One foundational method in topic modeling is Latent Dirichlet Allocation (LDA) [6], which assumes that documents are a mixture of topics and each topic is a distribution over words [2], [7]. Due to its interpretability, flexibility, and simplicity, LDA continues to be a popular choice for analyzing text data. LDA and its variations have been the most frequently applied techniques across domains such as software engineering, communication, and digital economy research [8], [9]. It is also considered effective in extracting human-readable topics and generating semantically meaningful clusters in large datasets [10], [11].

With its ability to uncover latent themes and organize vast textual information into coherent topics, LDA offers a simple yet effective approach for analyzing large-scale social media data [6]. Its probabilistic framework enables the identification of semantically meaningful clusters, even from informal and unstructured user-generated content [7]. This makes LDA particularly suitable for extracting underlying themes from public comments on platforms such as YouTube. In the context of this study, LDA serves as an effective tool to explore the dominant narratives and public sentiments embedded in discussions surrounding the Free Nutritious Meals (MBG) program.

Applying topic modeling to public comments can be useful for evaluating public discourse, including reactions to public policies. As highlighted by Gandomi and Haider [12], big data analytics, including textual analysis, can inform decision-making by revealing the public's perspective. This

makes topic modeling an important tool in understanding public opinion toward newly launched government programs. Moreover, several studies emphasize the importance of statistical evaluation metrics, such as coherence score and perplexity, to assess the quality and validity of generated topics [9], [13], [14]. For better interpretability, visualization tools such as pyLDAvis are also widely used to present the topics and their distribution across the corpus [15].

Previous studies have applied LDA in various domains to extract public opinion and thematic patterns. For instance, LDA was used to analyze communication strategies in the 2024 Indonesian presidential election through scraped data from candidate-related sources [16]. Another study applied LDA to identify major topics in online news related to a local cosmetic brand, validating topic relevance using coherence scores and human judgment [17]. In the tourism sector, LDA was employed to explore discourse on Yogyakarta tourism based on Twitter data, supported by coherence and perplexity metrics for model evaluation [18]. On a broader scale, LDA was applied to predict research trends across 30 years of academic publications [19] and to analyze student feedback evaluations [20]. Additionally, studies on social media marketing and sustainability topics have used LDA to map long-term trends and thematic evolution [21],[22]. LDA has also been utilized to study flipped classroom adoption [23], management information systems [24], and aspect-based sentiment analysis using customized LDA models such as SS-LDA [25].

Despite the broad application of LDA across various fields, its utilization in analyzing public discourse around Indonesia's recent national policies remains limited. One such critical and timely issue is the Free Nutritious Meals (Makan Bergizi Gratis/MBG) program proposed by President Prabowo and Vice President Gibran, which aims to improve children's health and educational equity by providing free nutritious meals in schools. Intended to address nutritional disparities and reduce long-term educational inequalities, the MBG program has rapidly become a focal point of intense public debate. Discussions in the media and on social networks frequently highlight concerns over the feasibility of its nationwide implementation, budget allocation transparency, coordination across government institutions, and potential political motivations behind its inception. On social platforms such as YouTube, public reactions are notably polarized, reflecting a complex mix of support, skepticism, and criticism. Given this context, applying LDA topic modeling to analyze systematically user-generated comments can yield valuable insights into the dominant themes shaping public sentiment and discourse surrounding the MBG program.

Existing studies on MBG have predominantly utilized qualitative descriptive methods. One study explored how online media platforms such as Detik.com and VIVA.co.id frame the MBG policy, either in support or in opposition [26]. Another research examined the MBG program through the lens of Pancasila, particularly its alignment with the value of social justice [27]. Other works have discussed MBG from socio-political [28] and health-related perspectives [29], raised concerns about the program's sustainability and political motives [30], and explored its impact on student



motivation using a sociological lens [31]. Although a few studies have adopted NLP-based approaches, such as sentiment analysis using the Naive Bayes algorithm to classify YouTube comments on MBG, revealing a predominance of negative sentiments [32], no existing research has applied topic modeling techniques to explore public discourse on this issue, particularly through YouTube comment data. Given this gap, the application of LDA Topic Modeling as a simple and foundational effective topic modeling method to analyze YouTube comments related to MBG represents a meaningful starting point.

In response to this identified research gap, the purpose of this study is to initiate the use of topic modeling, specifically LDA, to analyze public discourse on the MBG program through YouTube comments, thereby providing an initial thematic map of public perceptions. As an exploratory study, this research serves as a preliminary step toward a broader and continuous research agenda on topic modeling applications for analyzing MBG-related discourse across various social media platforms.

2 METHOD

2.1 Data

The data used in this study consists of public comments retrieved from YouTube videos related to the Free Nutritious Meals (Makan Bergizi Gratis/MBG) program. Data collection was conducted using the YouTube Data API on March 25, 2025, with the keyword “makan bergizi gratis.” The scraping process targeted videos published between November 2024 and March 2025, a period during which the MBG program became a prominent topic in national discourse. A total of 17,396 comments were successfully collected from 50 videos across 12 different YouTube channels.

These YouTube comments serve as valuable data sources for capturing spontaneous public reactions and diverse opinions surrounding the MBG program. As such, they were used as the primary textual input for the topic modeling analysis conducted in this study.

2.2 Data Analysis

The data analysis in this study was carried out in three main stages: preprocessing, describing the preprocessed data, and performing LDA topic modeling (Fig. 1). The preprocessing stage involved preparing raw textual data into a clean and standardized format suitable for computational analysis. This was followed by descriptive analysis to explore patterns in the cleaned data, and finally, the Latent Dirichlet Allocation (LDA) algorithm was applied to identify the underlying topics discussed in YouTube comments related to the Free Nutritious Meals (MBG) program.

Stage 1: Preprocessing

To prepare the text data for topic modeling, a comprehensive preprocessing pipeline was implemented to clean and

normalize the comments collected from YouTube. The preprocessing steps included:

1. *Text Cleaning*: Regular expressions (re) were used to perform initial text cleaning. This involved converting all text to lowercase, removing URLs, numbers, punctuation, extra whitespaces, and non-ASCII characters. A custom regex pattern was also applied to remove non-alphabetical symbols, including emojis and currency signs, ensuring that only clean textual content was retained.
2. *Text Normalization*: The IndoNLP library was used for standardizing informal or slang words into their formal equivalents. This step also handled word elongation (e.g., “maaaaaakan” to “makan”). Additionally, a custom dictionary was integrated to improve normalization accuracy for domain-specific vocabulary.
3. *Tokenization*: Tokenization was performed using Python's built-in split() function, which separates each comment into individual words based on whitespace.
4. *Stopwords Removal*: Stopwords were removed using a combination of predefined stopwords from the IndoNLP library and a custom stopwords list tailored to the context of the MBG discussion.
5. *Additional Cleaning*: The pipeline removed words shorter than three characters, eliminated empty lines, and dropped duplicate entries to ensure the quality and uniqueness of the final dataset.

This preprocessing pipeline ensured that the data was clean, normalized, and ready for the topic modeling process using the LDA algorithm.

Stage 2: Describing the preprocessed data

Following the preprocessing stage, the cleaned dataset was explored to provide a descriptive overview of its structure and content before proceeding to topic modeling. This exploration aimed to understand patterns of public engagement with YouTube videos related to the Free Nutritious Meals (MBG) program. The analysis was conducted using Microsoft Excel for generating pivot tables and bar charts, while Python was utilized for word cloud visualization. The descriptive analysis consisted of the following components:

1. *Comment Distribution by Channel*
Using Excel's pivot table feature, the number of comments per YouTube channel was calculated and visualized as a bar chart. This helped identify which media sources attracted the most public attention on the MBG program.
2. *Video Distribution by Channel*
The number of videos per channel was also summarized through Excel, then visualized to show the extent of content production by each media outlet concerning the MBG topic.
3. *Top Five Videos and Channels by Comment Count*
The five most-commented videos were identified through simple sorting and tabulation. These were presented in a table, highlighting the content and channels that sparked the highest level of audience response.



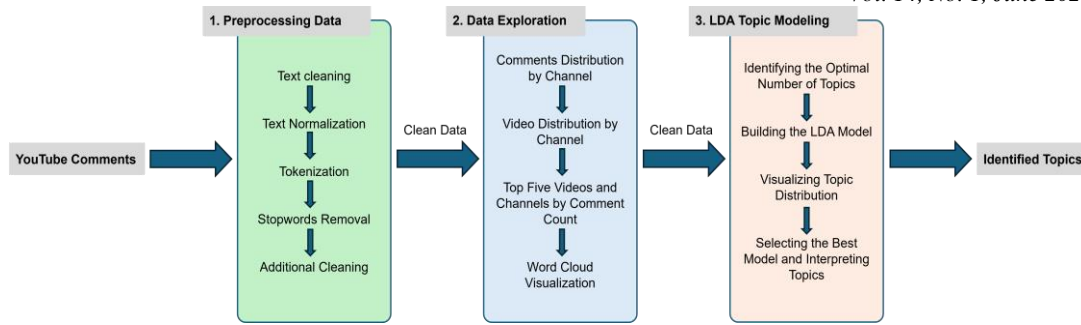


Figure 1. The research flow diagram

4. Word Cloud Visualization

A word cloud was generated using Python based on the preprocessed text data to provide an intuitive visual summary of the most frequently mentioned terms in the YouTube comments. This visualization served as a quick and accessible lens into the dominant vocabulary used by the public when discussing the MBG program.

Stage 3: LDA Topic Modeling

In this stage, Latent Dirichlet Allocation (LDA) was employed to extract latent thematic structures from the YouTube comment dataset related to the Free Nutritious Meals (MBG) program. LDA, first introduced by [7], is a generative probabilistic model designed to explain observations in large text corpora by inferring hidden topic distributions. The primary objective of LDA is to provide a compact representation of documents that captures the essential semantic relationships between words and documents, enabling tasks such as classification, summarization, novelty detection, and relevance assessment. LDA assumes that each document is a mixture of latent topics, and each topic is a distribution over words. The generative process for each document in a corpus involves the following steps [7]:

1. Choose the number of words $N \sim \text{Poisson}(\xi)$
2. Choose a topic distribution $\theta \sim \text{Dir}(\alpha)$
3. Choose a word $w_n \sim p(w_n | z_n, \beta)$, a multinomial distribution conditioned on topic z_n

Here, α is the parameter of the Dirichlet prior on the per-document topic distribution θ and β is the parameter of the topic-word distribution. The model infers the posterior distribution of the latent variables θ , z , and β given the observed data w , typically using approximate inference techniques such as Variational Bayes or Gibbs Sampling.

Input for the LDA model in this study consists of a preprocessed document-term matrix derived from the cleaned YouTube comments. Each document represents a single user comment, and each word corresponds to a term in the vocabulary after tokenization and stopword removal.

The output of the LDA model is a set of topics, each represented by a probability distribution over words, and a topic distribution for each document. These outputs were then

interpreted and labeled based on the top words associated with each topic. For evaluation, the use of coherence score and topic visualization will help ensure model quality and facilitate interpretation, following best practices from recent literature [33], [34].

The LDA topic modeling process in this study consisted of several key steps:

1. Identifying the Optimal Number of Topics

The first step involved determining the optimal number of topics by evaluating the coherence score across models with 2 to 10 topics. Topic coherence assesses the degree of semantic similarity among the top words in each topic, helping to differentiate interpretable topics from those that result from statistical noise [35]. A topic is considered coherent when its key terms support one another contextually [36]. In this study, we used the c_v coherence score proposed by [37], which combines a sliding window approach, NPMI, and cosine similarity to measure semantic consistency [36]. The topic number(s) with the highest coherence scores were selected as candidates for further modeling and interpretation.

2. Building the LDA Model

After identifying candidate topic numbers, LDA models were trained using the Gensim library. The Bag-of-Words (BoW) representation of the text was used, and model hyperparameters included passes=20 and alpha='auto' to ensure stability and convergence. Each model generated was examined for interpretability and topic separation.

3. Visualizing Topic Distribution

To understand the topic composition and inter-topic distance, the results of each LDA model were visualized using pyLDavis, an interactive tool that plots each topic as a circle based on its size and similarity to other topics. A good LDA model is indicated by three evaluation criteria of the topic modeling result:

- High coherence score,
- Clear separation between topics (no significant overlap),
- Semantically interpretable keywords within each topic.

This visual inspection was used alongside coherence scores to select the final model.



4. Selecting the Best Model and Interpreting Topics
From the candidate models, the best LDA model was selected based on the criteria mentioned in step 3.

Once the optimal model was selected, each topic was interpreted by examining its top representative words. This interpretation process involved human judgment and contextual understanding of the MBG-related discourse. The aim was to derive meaningful themes that reflect the diversity of public opinion found in the YouTube comment dataset.

3 RESULT AND DISCUSSION

This section presents the key findings of the study, organized into three main parts that reflect the sequential analytical process undertaken. First, it outlines the results of text preprocessing, which involved several steps such as case folding, stopword removal, normalization, and tokenization—essential steps to ensure the textual data was clean and suitable for computational analysis. Second, a descriptive analysis of YouTube comments is provided to highlight the volume, structure, and general characteristics of public engagement with the MBG program on the platform. This includes the frequency of comments, popular words, and emerging trends observed during preliminary exploration. Finally, the third part presents the results of topic modeling using the LDA method, offering insights into the underlying themes and discursive patterns present in over 16,000 comments. Together, these three components provide a comprehensive view of how the MBG program is being discussed, perceived, and contested within digital public discourse.

3.1 Preprocessing Results

The raw YouTube comments collected for this study contained various forms of noise, such as informal language, spelling variations, emojis, HTML artefacts, and non-standard expressions commonly found in social media discourse. To prepare the data for topic modeling, a multi-stage preprocessing pipeline was implemented to clean, normalize, and standardize the text.

Table 1 presents a comparison between the original and preprocessed versions of selected YouTube comments. As shown, the cleaning process significantly reduced textual noise by removing emojis, punctuation, HTML tags, and irrelevant formatting. These elements, while common in user-generated content, often carry minimal semantic value and can hinder accurate pattern detection during computational analysis. The resulting cleaned text is more focused, structurally uniform, and well-suited for downstream tasks like tokenization and topic modeling.

Table 1. Examples of Raw vs. Preprocessed Comments

No	Comment	Preprocessed Comment
----	---------	----------------------

1	Bagaimana dengan kajian awal tim pemerintah sendiri sebagai penyusun program? Apakah memiliki kapabilitas yang cukup untuk merancang kebijakan yang berbasis data, berdampak, dan dapat dipertanggungjawabkan? Mengapa mengandalkan hasil penelitian "independen" kemudian?	kajian tim pemerintah penyusun program memiliki kapabilitas merancang kebijakan berbasis data berdampak dipertanggungjawabkan mengandalkan hasil penelitian "independen"
2	🤔🤔🤔 para pembohong masih gk tau malu ya. Coba dipikirkan klu memang punya otak, apa hubungan keadaan skrg dgn dijalankannya program makan gratis?? Kenapa anda tidak memikirkan akibat dari mega korupsi yg ada di Pertamina, dan lain2 🤔🤔??	pembohong tau malu coba dipikirkan otak hubungan dijalankannya program makan gratis memikirkan akibat mega korupsi pertamina
3	Mbg belum di desa kami Bener kata ahok duit kasi orang tuanya kalau yang bikin yang lain udah bocor	MBG desa benar ahok uang kasih orang tuanya bikin bocor

The initial text cleaning was performed using regular expressions (regex), which efficiently handled basic normalization tasks such as converting all text to lowercase and removing URLs, numeric values, excessive punctuation, extra whitespaces, and non-ASCII characters. This step ensured the consistency of textual units across the dataset. Following this, tokenization was applied using Python's built-in `split()` function, which separated each comment into individual tokens (words) based on whitespace. Although straightforward, this method was sufficient for the study's goals, especially given the informal and concise nature of most YouTube comments. Future research could benefit from more advanced tokenization techniques, such as leveraging NLP libraries that consider context and compound words, to enhance text segmentation quality.

To handle the diverse range of informal expressions, a custom normalization dictionary was developed in addition to the IndoNLP library. This dictionary standardized common slang, misspellings, and elongated words. For example, terms like "gerati" and "geratis" were mapped to "gratis", and compound expressions such as "makan bergizi gratis" were standardized to the acronym "MBG." A sample of the custom normalization dictionary is shown in Table 2.

Table 2. Sample Entries from the Custom Text Normalization Dictionary

Informal Word	Normalized Word
gerati	gratis
geratis	gratis
makan bergizi gratis	MBG
makan bergizi	makanbergizi
mkan	makan
sd	sekolahdasar
duit	uang



orang tuanya	orangtua
dikasih	diberi
anaknya	anak

Stopwords were removed using a combination of the IndoNLP default stopword list and an extended custom list tailored to Indonesian social media language. This included colloquial terms such as *"tuh," "nih," "doang,"* and *"benarbrpak,"* which, while common in online conversation, contributed little to topic differentiation. A portion of this custom stopword list is presented in Table 3.

Final cleaning steps involved removing words shorter than three characters, eliminating empty rows, and dropping duplicate comments to ensure that unique, meaningful observations were retained. These preprocessing efforts were essential to improving the semantic quality of the corpus and minimizing the risk of noise-driven topic distortions. As a result, the dataset was transformed into a cleaner and more semantically coherent form, enabling a more accurate and interpretable topic modeling process in the next phase of the analysis.

3.2 Descriptive Analysis of YouTube Comments on MBG

Following the preprocessing phase, a descriptive analysis was conducted to better understand the structure and dynamics of the cleaned dataset. This stage aims to explore patterns of public engagement with YouTube content related to Indonesia's Free Nutritious Meals Program (MBG). Several exploratory visualizations were utilized to examine how discourse surrounding the MBG program unfolded across various media channels. These visualizations provide initial insights into the sources, intensity, and vocabulary of public responses before topic modeling is performed.

As shown in Fig. 2, the majority of comments were concentrated in a few high-traffic channels. The highest number of comments came from Liputan6 (6,226 comments), followed by Kompas.com (2,388) and KOMPASTV (2,182). In contrast, channels such as TVOneNews (37 comments) and CNBC Indonesia (36 comments) showed significantly lower engagement. This pattern indicates that public discourse on MBG was primarily driven by mainstream television-based news outlets with large subscriber bases and higher visibility on YouTube. The interaction volume suggests that audiences are more likely to comment on content from well-established and trusted media outlets. Consequently, the findings of this study largely reflect discussions taking place within the digital spaces of Indonesia's major news platforms.

Table 3. Sample Entries from the Custom Stopwords List

Custom Stopword		
nya	amp	bang
kau	pas	adek
bpk	banget	mah
gue	kalo	tuh



This article is distributed under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nc-nd/4.0/). See for details: <https://creativecommons.org/licenses/by-nc-nd/4.0/>

nak	nih	href
deh	biar	doang
dek	karna	namanya
gua	benarbrpak	sih
cuman	kayak	

Out of the initial 17,396 YouTube comments collected, a total of 16,228 comments were retained for analysis after the preprocessing stage. These comments were sourced from a diverse range of channels, including official media outlets and individual content creators who discussed the MBG program. The inclusion of both institutional and independent sources aimed to capture a wide spectrum of public discourse, from formal news narratives to grassroots reactions. By incorporating content from multiple channels, the dataset reflects not only the breadth of public engagement but also the diversity of opinion, tone, and framing related to the MBG initiative.

Fig. 3 shows the number of MBG-related videos uploaded by each channel, while Table 4 highlights the five most commented videos along with their respective channels. Although KOMPASTV and CNN Indonesia published multiple videos (14 and 10, respectively), the single video from Liputan6 received the highest number of comments, 6,226 entries, far surpassing others.

Interestingly, the total number of comments from just the top five videos amounted to 12,076, which represents approximately 74.4% of all comments in the dataset (16,228 comments). This reveals that the dataset used for topic modeling is highly dominated by reactions to a small number of viral videos, rather than a balanced distribution across various content.

In particular, the Liputan6 video, which features a student criticizing the taste of the MBG meal, appears to have sparked the strongest public reaction. This concentration of discourse suggests that the insights drawn from the topic modeling later in this study largely reflect netizens' responses to a handful of highly engaging videos, rather than the broader spectrum of coverage.

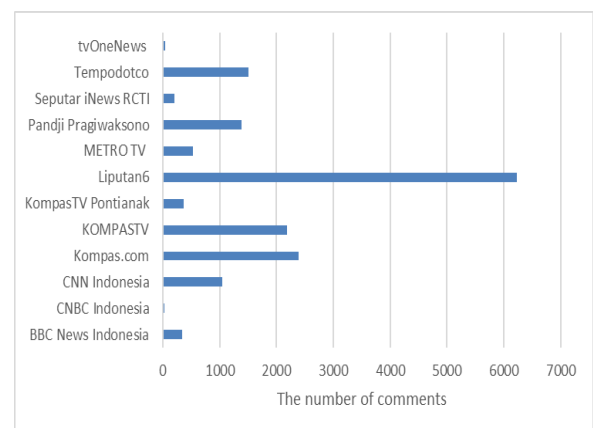
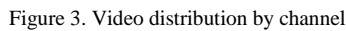


Figure 2. Comment distribution by channel



Channel	Video Title	Comments
Liputan6	Siswa SD di Palembang Kritik Menu Makan Bergizi Gratis, Rasanya Hambar Liputan 6	6226
Kompas.com	Nuel Nangis Dapat Makan Bergizi Gratis dari Tentara, Ternyata Ini Alasannya...	1848
Tempodotco	Omon-omon Makan Bergizi Gratis Explained	1382
Pandji Pragiwaksono	MEMBELA MAKAN BERGIZI GRATIS	1380
KOMPASTV	Ragam Impresi Siswa-Siswi SDN Larangan Sidoarjo Usai Disuguhkan Menu Makan Bergizi Gratis	1240
Total comments of the five top videos		12076

Key terms such as "*makan*" (eat), "*anak*" (child), "*gratis*" (free), "program", and "*makanan*" (food) appear most prominently, reflecting strong public attention toward the core concept of the Free Nutritious Meals initiative. These words suggest that the main elements of the policy have successfully captured public interest and become focal points of discussion. Their dominance in the word cloud implies a shared understanding or at least a repeated emphasis among users about the program's central objectives.



Figure 4. Word cloud of the most frequent terms

Table 5. Coherence scores for LDA models with various topic counts

The number of topics	Coherence Score
2	0.5721
3	0.5146
4	0.4845
5	0.4849
6	0.4322
7	0.4794
8	0.4880
9	0.5251
10	0.4783

In topic modeling, selecting the optimal number of topics involves balancing both statistical metrics and interpretability. Although the two-topic model demonstrated the highest coherence, it was considered insufficient to capture the thematic diversity of over 16,000 public comments. Conversely, the nine-topic model, despite its relatively high coherence score, exhibited overlapping topic areas based on the pyLDavis visualization, particularly among Topics 3 through 9, as shown in Fig. 5c. As noted in [37], coherence scores provide a valuable indicator of topic quality; however, semantic clarity and topic distinctiveness must also be examined visually. Supporting this view, visual tools like pyLDavis can be used to assess the interpretability and separation of topics [38]. Therefore, we argue that topic selection cannot rely solely on numerical values but must also consider how easily the topics can be meaningfully interpreted and differentiated by human judgment.

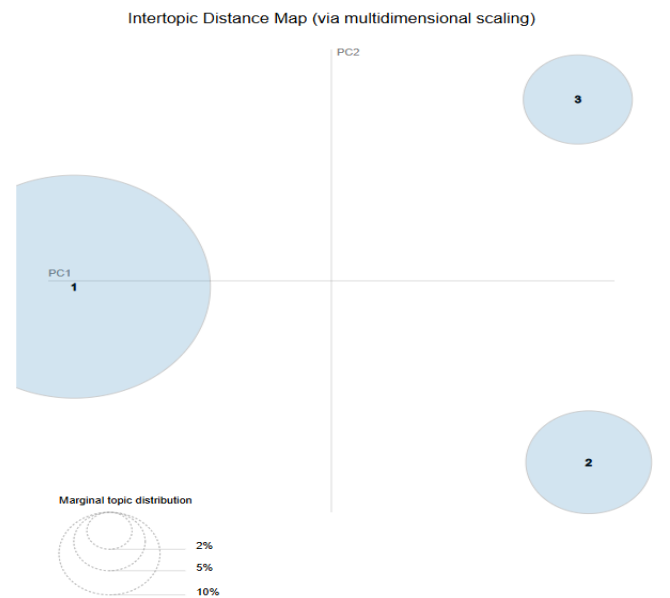
The three-topic model, presented in Fig. 5b, showed the best balance between coherence, topic separation, and interpretability. Compared to the overly generalized output of the two-topic model (Fig. 5a), the three-topic configuration provided more nuanced and distinct thematic groupings without the significant overlap found in the nine-topic model. The two-topic model tended to merge contrasting ideas under broad umbrellas, which made interpretation difficult, while the nine-topic model, although more granular, introduced redundancies and scattered subtopics. The three-topic structure, therefore, represented an optimal midpoint, offering sufficient detail while preserving thematic clarity. Based on this combination of quantitative coherence and qualitative visualization, the three-topic model was selected as the final model for further interpretation and analysis of public discourse on the MBG program. This decision is consistent with best practices in topic modeling research, which emphasize not only coherence scores but also human judgment in evaluating the clarity and usefulness of topic groupings.

To assess the quality and interpretability of the selected three-topic LDA model, the results were visualized using pyLDavis, an interactive tool designed to explore topic structure through both quantitative and qualitative dimensions. Fig. 6 displays the topic-specific visualizations for Topics 1, 2, and 3, respectively, with the lambda parameter (λ) set to 0.5. This setting balances word frequency and word uniqueness, aiding interpretation by highlighting

terms that are both common and exclusive to each topic. The visualization offers a spatial representation of topic distance and overlap, enabling intuitive assessment of the distinctness and relationships between topics. By examining the term relevance within each topic bubble, researchers can more accurately assign interpretive labels and understand how each topic reflects public sentiment toward different aspects of the MBG program. This step not only enhances transparency in topic interpretation but also supports reproducibility in computational discourse analysis.



(a) 2-topic model



(b) 3-topic model



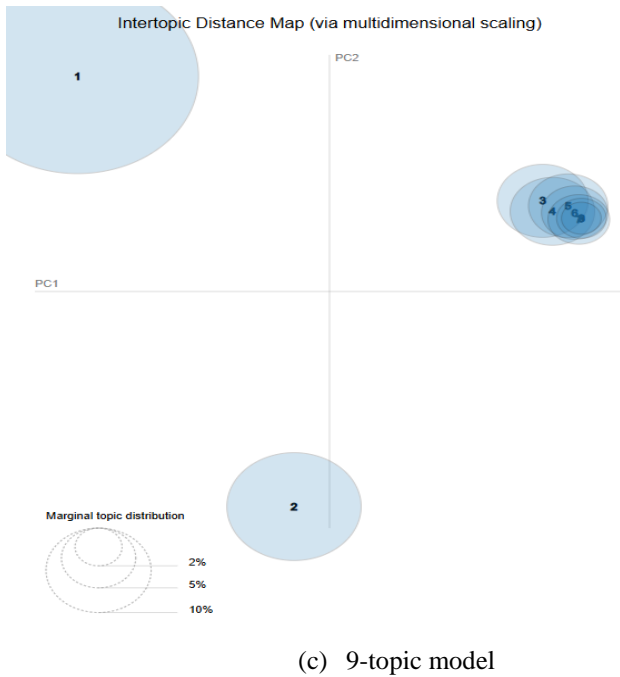
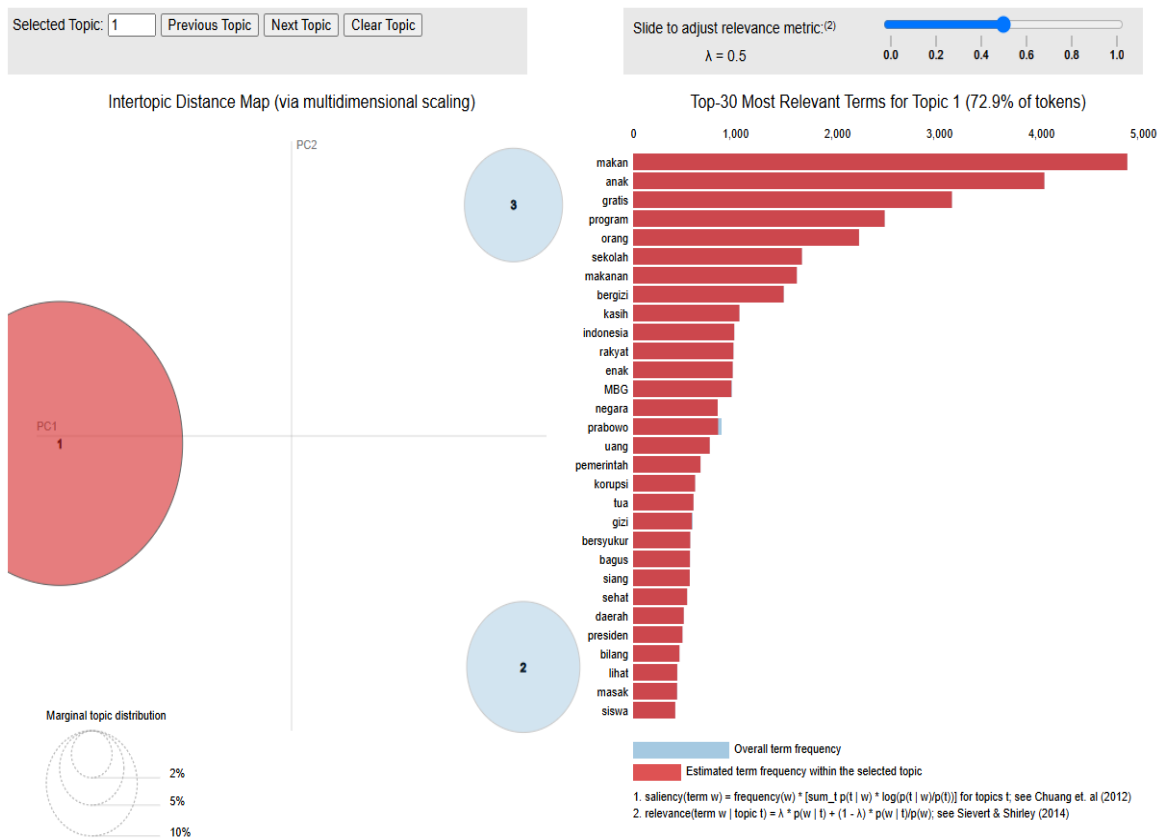


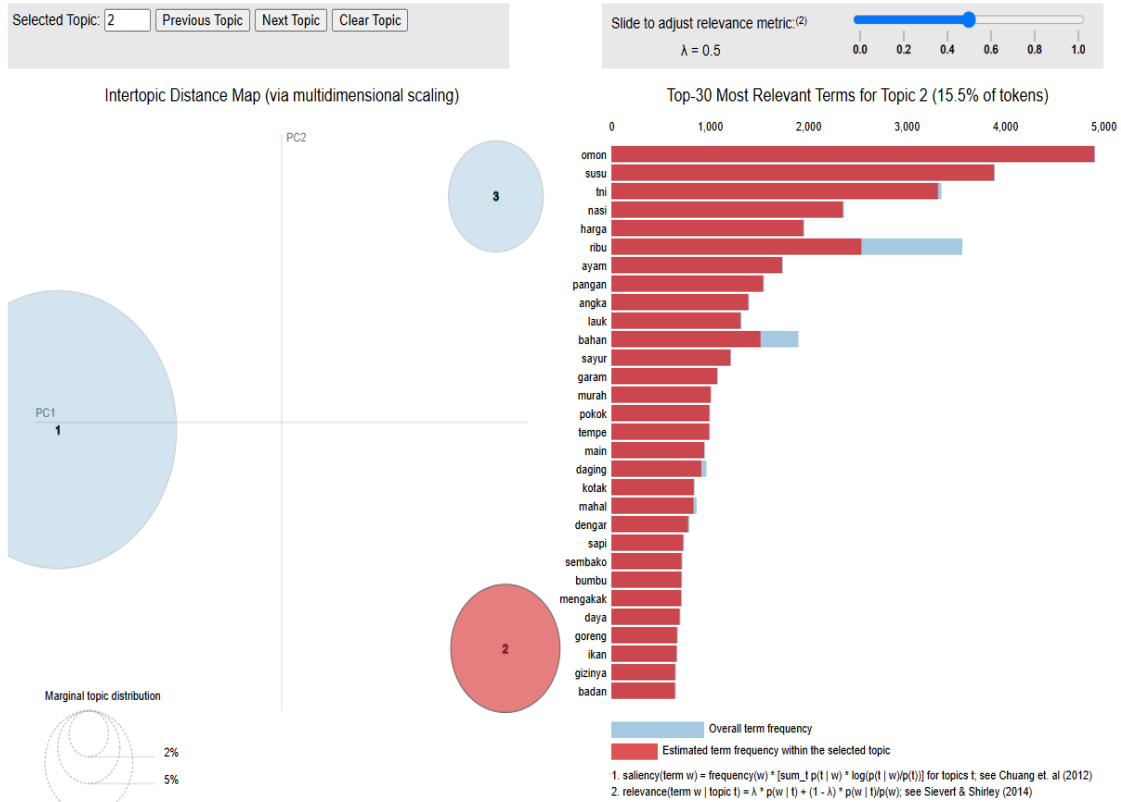
Figure 5. pyLDAvis visualization of the topic models

In pyLDAvis, λ determines how relevance scores are computed for terms within a topic [38]. This parameter directly affects the ordering of terms in the visualization by weighting two key factors: term frequency within the topic and term exclusivity across topics. When λ is set to 1, terms are ranked purely based on their probability within the selected topic. This means that common and frequently occurring words within that topic will dominate the list, even if they also appear in other topics. While this setting is useful for identifying the core vocabulary of a topic, it can sometimes obscure the terms that make the topic uniquely different from others. Conversely, when λ is set to 0, ranking prioritizes terms with the highest lift—those that are more unique to a particular topic compared to their frequency across all topics. Although this setting surfaces more distinctive terms, it may downplay high-frequency terms that are central to understanding the topic's overall content. Therefore, setting λ to 0.5 offers a balanced view by incorporating both term frequency and uniqueness, enabling a more nuanced interpretation. This compromise is particularly helpful for human analysts, as it supports the identification of both core concepts and distinguishing features, making the visualization more insightful and the topic labeling more accurate.

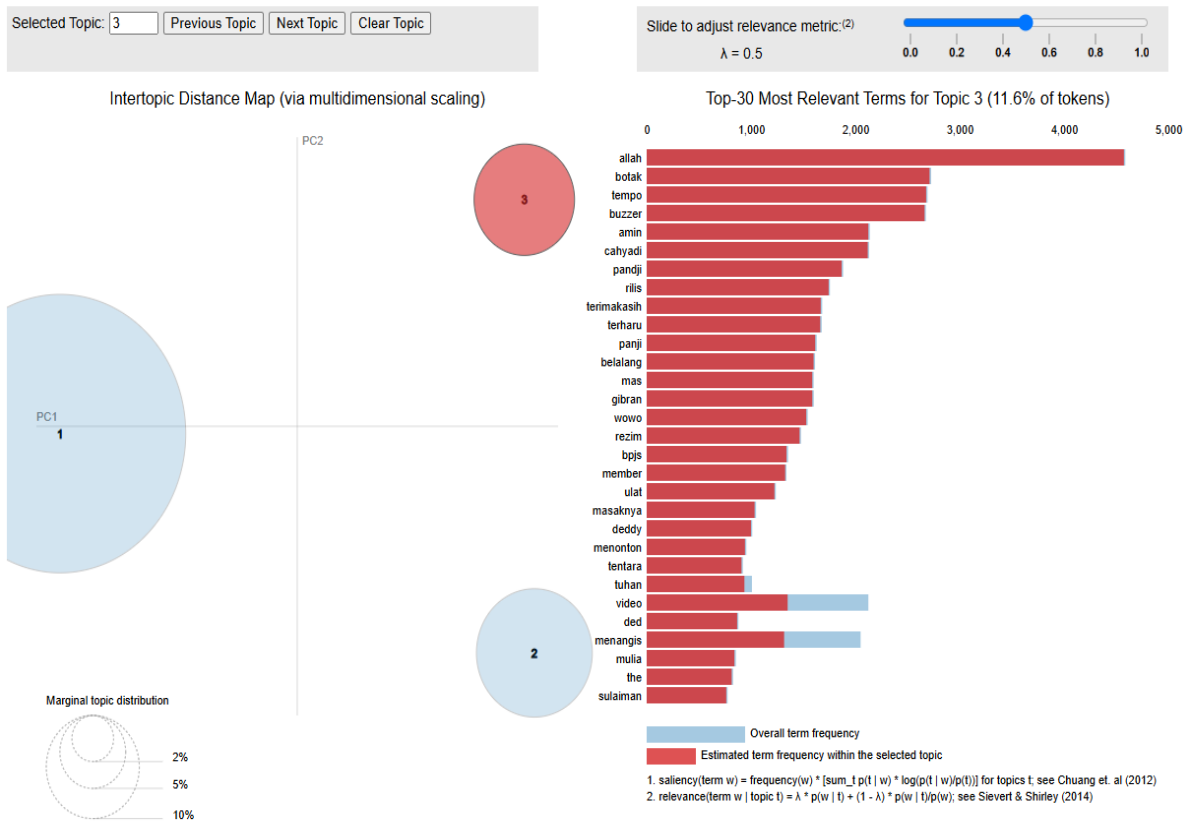


(a) Topic-1





(b) Topic-2



(c) Topic-3

Figure 6. PyLDAvis visualization



Each figure consists of two panels. The left panel shows the Intertopic Distance Map, where each circle represents a topic. The distance between circles reflects the semantic similarity between topics based on multidimensional scaling (MDS), and the size of the circle corresponds to the marginal topic distribution—the proportion of the entire corpus represented by each topic. Based on the bar charts, the topic distribution is as follows:

- Topic 1: 72.9% of tokens
- Topic 2: 15.5% of tokens
- Topic 3: 11.6% of tokens

This indicates that Topic 1 dominates the corpus, while Topic 3 represents the smallest share of discourse, though still thematically distinct.

The right panel of each figure presents the top 30 most relevant terms for the selected topic. The red bars indicate each word's frequency within the topic, while the blue bars show the same word's frequency in the overall corpus. This dual encoding helps reveal which words are semantically central and distinctive for each topic.

- Fig.6a (Topic 1) highlights words like *omon*, *susu*, *tni*, *harga*, and *pangan*, suggesting a focus on food distribution, logistics, and pricing concerns.
- Fig.6b (Topic 2) features terms such as *makan*, *anak*, *gratis*, *bergizi*, and *sekolah*, which reflect the core purpose of the MBG program.
- Fig.6c (Topic 3) includes emotionally and politically charged terms like *allah*, *buzzer*, *tempo*, and *terharu*, indicating a mix of public sentiment, media discourse, and political commentary.

These visualizations support the semantic clarity, topic separation, and interpretability of the model, justifying the selection of three topics as an optimal representation of public discourse surrounding the MBG program.

To deepen the understanding of each theme uncovered by the LDA model, the most representative keywords for each topic were compiled and interpreted semantically. Table 6 presents a summary of the three topics along with their general themes and concise interpretations, offering a clear view of how public discourse on the MBG program is distributed across different concerns and sentiments.

Table 6. Topic interpretation summary

Topic	General Theme	Interpretation
1	Food Prices and Distribution Concerns	Discussions revolve around food supply, pricing, and public skepticism toward MBG program execution.
2	Benefits of the Free Nutritious Meal Program	Public appreciation of the program's positive impact on schoolchildren and social welfare.
3	Social Reactions and Political Opinions	Emotional and political responses, including gratitude, criticism, and media/influencer involvement.



This article is distributed under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nc-nd/4.0/). See for details: <https://creativecommons.org/licenses/by-nc-nd/4.0/>

Topic 1 reflects public concerns over the logistics and affordability of food distribution under the MBG program. Keywords such as *susu* (milk), *nasi* (rice), *harga* (price), and *pangan* (food) point to a dominant theme around the basic components of the meals and their costs, indicating that many commenters were attentive to the practical aspects of what would be served. Moreover, terms like *ribu* (thousand), *angka* (numbers), and *TNI* (military) suggest discussions surrounding the program's budgeting scale and the involvement of institutional actors, such as the military, in operational roles. This reflects a broader concern about state capacity and transparency in executing large-scale public programs. Interestingly, the inclusion of the term *omon*, which appears to be a satirical or colloquial twist on *omong-omong* (idle talk), suggests a layer of public cynicism toward the promises made by political elites. Collectively, this topic captures critical discourse focused on the feasibility, cost, and trustworthiness of the program's implementation.

Topic 2 represents general public support and recognition of the MBG program's value, particularly for schoolchildren and their families. Keywords such as *makan* (eat), *anak* (child), *gratis* (free), *sekolah* (school), and *bergizi* (nutritious) strongly align with the program's stated objectives—providing accessible, healthy meals for students. The presence of words like *kasih* (care) and *program* reflects a positive framing of the policy as a form of social care and public service. Many comments associated with this topic likely come from individuals who view the initiative as a step forward in addressing child hunger and educational inequality. This topic illustrates a discourse grounded in gratitude, hope, and policy endorsement, showing that public sentiment is not solely critical, but also acknowledges the potential benefits of such welfare-driven initiatives.

Topic 3 is more complex and emotionally charged, capturing a blend of religious, emotional, and politicized reactions. Terms like *allah*, *amin*, *terimakasih*, and *terharu* suggest a tone of personal gratitude and spiritual reflection, possibly from users who emotionally resonated with the program's mission. At the same time, the presence of words such as *buzzer*, *tempo*, *botak*, and *pandji* indicate a layer of politicization, shaped by current political narratives, media framing, and influencer commentary. These terms hint at how online discussions about MBG are often intertwined with broader political debates, campaign messaging, or public skepticism toward political actors and media figures. This topic illustrates how public discourse can be shaped not only by policy content but also by the socio-political climate and the personalities driving public conversation.

Upon closer examination, some of these terms—particularly *pandji*, *panji*, or *daddy*—likely appear due to the inclusion of YouTube channel names or speaker identifiers within the comment texts. For instance, Pandji Pragiwaksono is a popular public figure whose video was part of the dataset. His name, repeated frequently in the comment section, may have skewed the topic content to reflect the identity of the uploader rather than the actual thematic focus of the discourse. This phenomenon is a common limitation in topic modeling based on user-generated content, where named entities or proper nouns—especially those tied to video titles, creators, or influencers—can disproportionately shape the

statistical distribution of words. Consequently, some topics may be more centered on individual personalities than on the underlying issues being discussed. This highlights the importance of incorporating manual filtering or named entity recognition (NER) during preprocessing, to differentiate between discussion-relevant keywords and metadata-like elements that may distort the semantic coherence of generated topics. Addressing this nuance ensures a more accurate interpretation of public discourse and prevents misattributing topic salience to factors external to the actual content of conversation.

This observation points to a limitation in the preprocessing stage, particularly the absence of filtering for channel names, speaker identifiers, and other non-content elements within the comment data. These elements, while structurally present in the dataset, do not represent the actual substance of the discourse. When such names are frequently mentioned, often because they are tied to the video title or uploader, they risk being interpreted by the model as significant terms, despite offering minimal thematic value.

Including these non-semantic terms introduces semantic noise, which can distort the accuracy of topic modeling. Instead of capturing meaningful patterns or genuine public sentiment, the model may inadvertently reflect patterns related to content structure or digital metadata. This can result in misleading topic interpretations, where emphasis appears to be placed on individuals or media channels rather than the issues discussed. Such distortions may reduce the validity of insights derived from the model, especially in studies aiming to map public opinion or discourse trends.

To address this issue, future studies should consider enhancing their preprocessing pipeline by implementing custom stopword lists or leveraging named entity recognition (NER) tools to identify and exclude non-content-specific terms. By systematically removing these elements, researchers can improve topic purity and ensure that the resulting themes are driven by authentic, content-based language rather than by repetitive mentions of media figures or channel labels. This refinement not only strengthens model interpretability but also enhances the methodological rigor of computational discourse analysis.

In summary, the topic modeling results reveal that public discourse surrounding the Free Nutritious Meals (MBG) program is characterized by three major themes: concerns about food pricing and distribution logistics, general appreciation of the program's benefits, and emotionally nuanced political and social reactions. These thematic groupings reflect the multidimensional nature of public response, where technical, emotional, and political dimensions coexist within online commentary. The first theme shows that many citizens are not only aware of the policy's intentions but also critically engaged with the practical challenges of implementation, such as cost, quality, and institutional capacity. The second theme highlights a layer of optimism and trust, especially among communities that may directly benefit from the policy, suggesting that welfare-oriented programs still resonate positively with segments of the public.

The third theme, rich with religious expression, gratitude, and political satire, reveals a more complex digital narrative, where public sentiment is shaped not only by the policy itself but also by broader socio-political undercurrents. This includes skepticism toward government agendas, influence from public figures, and ideological framing by online communities. Such emotionally driven and politically inflected discourse demonstrates how online platforms function as hybrid spaces: not just arenas of information exchange, but also emotional outlets and arenas of political engagement. The role of influencers, content creators, and even algorithms in amplifying certain voices or perspectives cannot be underestimated in this context.

By uncovering these latent themes, the analysis offers valuable insights into how large-scale social policies are perceived, negotiated, and contested in digital public spheres, particularly on platforms like YouTube. It highlights the utility of computational tools like LDA in capturing public sentiment on a scale while also emphasizing the need for contextual interpretation. These findings not only inform policymakers about the layers of public reaction but also contribute to the growing body of research at the intersection of text analytics and public policy discourse. The implications of these findings are further discussed in the concluding section.

4 CONCLUSION

This study aimed to identify dominant themes in public discourse surrounding Indonesia's Free Nutritious Meals (MBG) program by applying Latent Dirichlet Allocation (LDA) topic modeling to a dataset of YouTube comments. The analysis revealed three distinct topics: (1) food prices and distribution concerns, (2) public appreciation for the MBG program's benefits, and (3) emotionally and politically charged reactions involving social sentiment and media narratives. These findings underscore the multifaceted nature of public response to large-scale social policies, capturing both support and skepticism expressed through online platforms.

While the LDA model successfully extracted interpretable topics, the study also encountered limitations related to the text preprocessing stage, particularly due to the linguistic complexity and variability of informal Indonesian-language user comments. Challenges such as non-standard spelling, slang, repetition, and the presence of channel or speaker names (e.g., *Pandji*, *Deddy*) potentially affected topic purity and interpretability.

Building on the limitations identified in this study, future research should consider developing more robust and language-specific preprocessing techniques to better handle the informal and diverse linguistic characteristics of Indonesian social media data. Enhancements such as custom stopword lists, named entity recognition (NER), and semantic normalization tailored to Bahasa Indonesia are crucial for improving topic purity and interpretability. In addition, alternative topic modeling approaches—such as Non-negative Matrix Factorization (NMF), BERTopic, or GPT-based models—could be explored to provide comparative insights and address potential shortcomings of LDA. These



models offer various advantages, including improved topic coherence, contextual understanding, and zero-shot capabilities. Future studies may also expand the scope of analysis by incorporating data from other platforms such as Twitter, TikTok, or online news comment sections, or by integrating multimodal sources like video transcripts, hashtags, and visual content. Such advancements would contribute to a more comprehensive understanding of public discourse and the dynamic nature of opinion formation in the digital space.

AUTHOR'S CONTRIBUTION

In the research titled "*LDA Topic Modeling Analysis of Public Discourse on Indonesia's Free Nutritious Meals Program (MBG)*", Cici Suhaeni was responsible for data scraping, conducting the main analysis, and drafting the initial manuscript. Laily Nissa Atul Mualifah contributed by developing and refining the Python code and improving the overall structure and clarity of the writing. Hari Wijayanto guided the interpretation of results and supported the development of the discussion section. Through the collaborative efforts and contributions of all authors, this study was completed and is expected to contribute meaningfully to the understanding of public discourse on social policy in digital platforms.

COMPETING INTERESTS

Following the publication ethics of this journal, Cici Suhaeni, Laily Nissa Atul Mualifah, and Hari Wijayanto, as the authors of this journal article, declare that this journal article has no conflict of interest (COI) or competing interests (CI).

REFERENCES

- [1] M. Thelwall, P. Sud, and F. Vis, "Commenting on YouTube videos: From guatemalan rock to El Big Bang," *J. Am. Soc. Inf. Sci.*, vol. 63, no. 3, pp. 616–629, Mar. 2012, doi: 10.1002/asi.21679.
- [2] D. M. Blei, "Probabilistic topic models," *Commun. ACM*, vol. 55, no. 4, pp. 77–84, Apr. 2012, doi: 10.1145/2133806.2133826.
- [3] P. DiMaggio, M. Nag, and D. Blei, "Exploiting affinities between topic modeling and the sociological perspective on culture: Application to newspaper coverage of U.S. government arts funding," *Poetics*, vol. 41, no. 6, pp. 570–606, Dec. 2013, doi: 10.1016/j.poetic.2013.08.004.
- [4] A. M. G. A., S. Robledo, and M. Zuluaga, "Topic Modeling: Perspectives From a Literature Review," *IEEE Access*, vol. 11, pp. 4066–4078, 2023, doi: 10.1109/ACCESS.2022.3232939.
- [5] I. Vayansky and S. A. P. Kumar, "A review of topic modeling methods," *Information Systems*, vol. 94, p. 101582, Dec. 2020, doi: 10.1016/j.is.2020.101582.
- [6] H. Jelodar et al., "Latent Dirichlet Allocation (LDA) and Topic modeling: models, applications, a survey," Dec. 06, 2018, arXiv: arXiv:1711.04305. doi: 10.48550/arXiv.1711.04305.
- [7] D. M. Blei, "Latent Dirichlet Allocation," *J. Mach. Learn. Res.*, vol. 3, pp. 993–1022, 2003.
- [8] C. C. Silva, M. Galster, and F. Gilson, "Topic modeling in software engineering research," *Empir Software Eng*, vol. 26, no. 6, p. 120, Nov. 2021, doi: 10.1007/s10664-021-10026-0.
- [9] U. Detthamrong et al., "Topic Modeling Analytics of Digital Economy Research: Trends and Insights," *J. Scientometric Res.*, vol. 13, no. 2, pp. 448–458, Aug. 2024, doi: 10.5530/jscires.13.2.35.
- [10] D. Yu and B. Xiang, "Discovering topics and trends in the field of Artificial Intelligence: Using LDA topic modeling," *Expert Systems with Applications*, vol. 225, p. 120114, Sep. 2023, doi: 10.1016/j.eswa.2023.120114.
- [11] R. Churchill and L. Singh, "The Evolution of Topic Modeling," *ACM Comput. Surv.*, vol. 54, no. 10s, pp. 1–35, Jan. 2022, doi: 10.1145/3507900.
- [12] A. Gandomi and M. Haider, "Beyond the hype: Big data concepts, methods, and analytics," *International Journal of Information Management*, vol. 35, no. 2, pp. 137–144, Apr. 2015, doi: 10.1016/j.ijinfomgt.2014.10.007.
- [13] A. Krishnan, "Exploring the Power of Topic Modeling Techniques in Analyzing Customer Reviews: A Comparative Analysis," Aug. 19, 2023, arXiv: arXiv:2308.11520. doi: 10.48550/arXiv.2308.11520.
- [14] A. Farea, S. Tripathi, G. Glazko, and F. Emmert-Streib, "Investigating the optimal number of topics by advanced text-mining techniques: Sustainable energy research," *Engineering Applications of Artificial Intelligence*, vol. 136, p. 108877, Oct. 2024, doi: 10.1016/j.engappai.2024.108877.
- [15] B. Gencoglu, M. Helms-Lorenz, R. Maulana, E. P. W. A. Jansen, and O. Gencoglu, "Machine and expert judgments of student perceptions of teaching behavior in secondary education: Added value of topic modeling with big data," *Computers & Education*, vol. 193, p. 104682, Feb. 2023, doi: 10.1016/j.compedu.2022.104682.
- [16] T. Irawan, L. Mutawalli, S. Fadli, and W. Bagye, "Topic Modelling Pola Komunikasi Pilpres 2024: Focus Web Scraping dan Latent Dirichlet Allocation," *J. Manaj. Inform. dan Sist. Inform.*, vol. 7, no. 2, pp. 186–194, 2024, https://doi.org/10.36595/misi.v7i2.1183.
- [17] E. Puspita, D. F. Shiddieq, and F. F. Roji, "Pemodelan Topik pada Media Berita Online Menggunakan Latent Dirichlet Allocation (Studi Kasus Merek Somethinc): Topic Modeling on Online News Media Using Latent Dirichlet Allocation (Case Study Somethinc Brand)," *MALCOM*, vol. 4, no. 2, pp. 481–489, Feb. 2024, doi: 10.57152/malcom.v4i2.1204.
- [18] Uray Nur Khadijah and Nuri Cahyono, "Analisis Topic Modelling Pariwisata Yogyakarta Menggunakan Latent Dirichlet Allocation (LDA)," *ijcs*, vol. 13, no. 4, Jul. 2024, doi: 10.33022/ijcs.v13i4.3816.
- [19] R. K. Gupta, R. Agarwalla, B. H. Naik, J. R. Evuri, A. Thapa, and T. D. Singh, "Prediction of research trends using LDA based topic modeling," *Global Transitions Proceedings*, vol. 3, no. 1, pp. 298–304, Jun. 2022, doi: 10.1016/j.gltp.2022.03.015.
- [20] O. Ishmael, E. Kiely, C. Quigley, and D. McGinty, "Topic Modelling using Latent Dirichlet Allocation (LDA) and Analysis of Students Sentiments," in 2023 20th International Joint Conference on Computer Science and Software Engineering (JCSSE), Phitsanulok, Thailand: IEEE, Jun. 2023, pp. 1–6. doi: 10.1109/JCSSE58229.2023.10201965.
- [21] J. H. Lee, J. Wood, and J. Kim, "Tracing the Trends in Sustainability and Social Media Research Using Topic Modeling," *Sustainability*, vol. 13, no. 3, p. 1269, Jan. 2021, doi: 10.3390/su13031269.
- [22] W.-R. Yang and H.-C. Yang, "Topic Modeling Analysis of Social Media Marketing using BERTopic and LDA," *The Journal of Industrial Distribution & Business*, vol. 13, no. 9, pp. 37–50, Sep. 2022, doi: 10.13106/JIDB.2022.VOL13.NO9.37.
- [23] O. Ozyurt, "Empirical research of emerging trends and patterns across the flipped classroom studies using topic modeling," *Educ Inf Technol*, vol. 28, no. 4, pp. 4335–4362, Apr. 2023, doi: 10.1007/s10639-022-11396-8.
- [24] H. Özköse, O. Ozyurt, and A. Ayaz, "Management Information Systems Research: A Topic Modeling Based Bibliometric Analysis," *Journal of Computer Information Systems*, vol. 63, no. 5, pp. 1166–1182, Sep. 2023, doi: 10.1080/08874417.2022.2132429.



- [25] B. Ozyurt and M. A. Akcayol, "A new topic modeling based approach for aspect extraction in aspect based sentiment analysis: SS-LDA," *Expert Systems with Applications*, vol. 168, p. 114231, Apr. 2021, doi: 10.1016/j.eswa.2020.114231.
- [26] M. E. D. Vanti, V. Octaviani, and M. Maryaningsih, "Analisis Framing Pemberitaan Program Makan Gratis Prabowo Subianto Di Media Online," *I. Kom. dan Adm. Pub.*, vol. 11, no. 1, Jun. 2024, doi: 10.37676/professional.v11i1.6396.
- [27] A. Andin et al., "Penerapan Nilai Pancasila Melalui Program Makan Bergizi Gratis," *Indonesian Journal of Education and Development Research*, vol. 3, no. 1, pp. 370–383, Dec. 2024, doi: 10.57235/ijedr.v3i1.4684.
- [28] A. Kiftiyah, F. A. Palestina, F. U. Abshar, and K. Rofiah, "Program Makan Bergizi Gratis (MBG) dalam Perspektif Keadilan Sosial dan Dinamika Sosial – Politik," *PJK*, vol. 5, no. 1, pp. 101–112, Apr. 2025, doi: 10.52738/pjk.v5i1.726.
- [29] R. Qomarrullah, S. Suratni, L. Wulandari S, and M. Sawir, "Dampak Jangka Panjang Program Makan Bergizi Gratis terhadap Kesehatan dan Keberlanjutan Pendidikan," *IJI_Publication*, vol. 5, no. 2, pp. 130–137, Mar. 2025, doi: 10.51577/ijipublication.v5i2.660.
- [30] P. A. Maharani, A. R. Namira, and T. V. Chairunnisa, "Peran Makan Siang Gratis Dalam Janji Kampanye Prabowo Gibran Dan Realisasinya," *JOLASOS*, vol. 1, no. 1, pp. 1–10, Jun. 2024, doi: 10.70656/jolasos.v1i1.79.
- [31] A. A. Merlinda and Yusmar Yusuf, "Analisis Program Makan Gratis Prabowo Subianto Terhadap Strategi Peningkatan Motivasi Belajar Siswa di Sekolah Tinjauan dari Perspektif Sosiologi Pendidikan," *RRJ*, vol. 7, no. 2, pp. 1364–1373, Jan. 2025, doi: 10.38035/trj.v7i2.1360.
- [32] B. Rahmatullah, S. A. Saputra, P. Budiono, and D. P. Wigandi, "Sentimen Analisis Makan Bergizi Gratis Menggunakan Algoritma Naive Bayes," *JifoTech*, vol. 5, no. 1, Mar. 2025, doi: 10.46229/jifotech.v5i1.978.
- [33] R. Egger and J. Yu, "A Topic Modeling Comparison Between LDA, NMF, Top2Vec, and BERTopic to Demystify Twitter Posts," *Front. Sociol.*, vol. 7, p. 886498, May 2022, doi: 10.3389/fsoc.2022.886498.
- [34] B. Yin and C.-H. Yuan, "Detecting latent topics and trends in blended learning using LDA topic modeling," *Educ Inf Technol*, vol. 27, no. 9, pp. 12689–12712, Nov. 2022, doi: 10.1007/s10639-022-11118-0.
- [35] K. Stevens, P. Kegelmeyer, D. Andrzejewski, and D. Buttler, "Exploring Topic Coherence over Many Models and Many Topics".
- [36] S. Mifrah, "Topic Modeling Coherence: A Comparative Study between LDA and NMF Models using COVID'19 Corpus," *IJATCSE*, vol. 9, no. 4, pp. 5756–5761, Aug. 2020, doi: 10.30534/ijatcse/2020/231942020.
- [37] M. Röder, A. Both, and A. Hinneburg, "Exploring the Space of Topic Coherence Measures," in *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, Shanghai China: ACM, Feb. 2015, pp. 399–408. doi: 10.1145/2684822.2685324.
- [38] C. Sievert and K. Shirley, "LDAvis: A method for visualizing and interpreting topics," in *Proceedings of the Workshop on Interactive Language Learning, Visualization, and Interfaces*, Baltimore, Maryland, USA: Association for Computational Linguistics, 2014, pp. 63–70. doi: 10.3115/v1/W14-3110.

