

# Topic Modeling for Short Texts with Large Language Models

Tomoki Doi<sup>1</sup> Masaru Isonuma<sup>1,2</sup> Hitomi Yanaka<sup>1</sup>

<sup>1</sup> The University of Tokyo <sup>2</sup> The University of Edinburgh

{doi-tomoki701, hyanaka}@is.s.u-tokyo.ac.jp m.isonuma@ed.ac.uk

## Abstract

As conventional topic models rely on word co-occurrence to infer latent topics, topic modeling for short texts has been a long-standing challenge. Large Language Models (LLMs) can potentially overcome this challenge by contextually learning the meanings of words via pretraining. In this paper, we study two approaches to using LLMs for topic modeling: parallel prompting and sequential prompting. Input length limitations prevent LLMs from processing many texts at once. However, an arbitrary number of texts can be handled by LLMs by splitting the texts into smaller subsets and processing them in parallel or sequentially. Our experimental results demonstrate that our methods can identify more coherent topics than existing ones while maintaining the diversity of the induced topics. Furthermore, we found that the inferred topics cover the input texts to some extent, while hallucinated topics are hardly generated.

## 1 Introduction

Topic modeling is the classical task of discovering latent topics that best describe a set of documents (Blei et al., 2003; Churchill and Singh, 2022). Recently, while neural topic models have worked successfully on various kinds of long documents (Miao et al., 2017; Srivastava and Sutton, 2017; Deng et al., 2020), they have not been able to handle short texts, such as social media posts and news headlines (Li et al., 2016; Wu et al., 2022).

Large Language Models (LLMs), such as InstructGPT (Ouyang et al., 2022) and GPT-4 (OpenAI, 2023), have shown impressive results on various tasks by providing task instructions in a zero-shot manner (Wang et al., 2023; Kocóń et al., 2023). Since conventional topic models infer the topics of words by relying on word co-occurrence, they perform worse on short texts. In contrast, as LLMs contextually learn the meanings of words by pre-

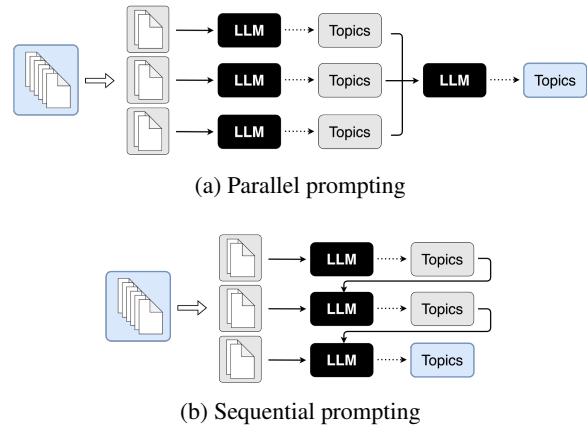


Figure 1: Topic modeling with LLMs by splitting a document set into subsets and prompting (a) in parallel or (b) sequentially.

training on massive text corpora, they could accurately infer the latent topics.

We propose two approaches to using LLMs for topic modeling: parallel prompting and sequential prompting (Figure 1). Due to the input length limitations of LLMs, an input document set must be split into smaller subsets, which are processed individually. Parallel prompting concurrently infers the topics of each subset and merges them to represent the topics of the whole document set. Sequential prompting processes each subset successively, updating the topics in every iteration. We assess our approaches across texts from various domains using multiple evaluation metrics.

The contributions of this study are as follows:

1. We propose parallel and sequential prompting methods for topic modeling using LLMs. Our methods can handle a large number of texts that cannot be processed in a single run due to the input length limitations of LLMs.
2. We validate the performance of our methods by comparing them with existing models and show that ours can identify more coherent topics than existing models while maintaining

the diversity of the induced topics.

3. We assess the document coverage and factuality of the inferred topics, due to concerns that LLMs may focus on only parts of documents or generate hallucinated topics. Evaluation results indicate that those concerns are negligible.

## 2 Background

Topic modeling is the task of identifying latent topics as a set of topic words representing each topic from a collection of documents (Blei et al., 2003). Topic modeling has conventionally been tackled with probabilistic models such as latent Dirichlet allocation (LDA, Blei et al., 2003). In recent years, however, neural models have come into widespread use due to their high performance (Srivastava and Sutton, 2017; Dieng et al., 2020; Grootendorst, 2022).

It is known that topic modeling for short texts is difficult for current topic models due to data sparsity (Li et al., 2016; Wu et al., 2022). TSCTM (Wu et al., 2022) is a current state-of-the-art neural topic model for short texts. This model addresses data sparsity by learning representations of documents using VQ-VAE (van den Oord et al., 2017), contrastive learning, and incorporation of data augmentation into the learning.

BERTopic (Grootendorst, 2022) uses a pre-trained encoder, Sentence-BERT (Reimers and Gurevych, 2019), to obtain clusters of documents and assigns topic words to each cluster by using a class-based TF-IDF procedure. Another related study is Stambach et al. (2023), in which LLMs are utilized to automatically evaluate topic quality. However, our study is the first to explore how well LLMs perform topic modeling.

## 3 Topic Modeling with LLMs

We introduce two approaches to performing topic modeling with LLMs: **parallel prompting** and **sequential prompting**. For these approaches, we apply common preprocessing, which involves randomly splitting a document set into subsets with the same size, smaller than the context length of the LLMs.

**Parallel Prompting** In the parallel prompting, LLMs identify topics for each subset in parallel by prompting the subset and the instruction of topic modeling. The topics of each subset are then

ID	Prompt
<b>Par<sub>TM</sub></b>	Write the results of simulating topic modeling for the following documents: [DOCS].
<b>Par<sub>Mrg</sub></b>	Write the results of merging the following topic modeling results: [TOPICS], [TOPICS], ...
<b>Seq<sub>TM</sub></b>	Write the results of simulating topic modeling for the following documents: [DOCS], Make the most use of the following topics: [TOPICS].

Table 1: Prompts for our methods. [DOCS] and [TOPICS] are replaced by a subset of documents and by previously identified topics, respectively.

Dataset	# of Documents	Text Length	Vocabulary Size
Tweet	2000	5.47	706
GoogleNewsT	11000	5.25	2376
StackOverFlow	19000	4.71	2544

Table 2: Dataset statistics. Each value is the average for five runs.

merged by LLMs. We use two kinds of prompts as shown in Table 1: (i) a **Par<sub>TM</sub>** prompt for parallel topic modeling for each subset, and (ii) a **Par<sub>Mrg</sub>** prompt for merging the topics from the results

**Sequential Prompting** In the sequential prompting, LLMs identify topics for each subset sequentially, considering the topics previously identified for the previous subset. We use the **Par<sub>TM</sub>** for the first subset, then use a **Seq<sub>TM</sub>** prompt in Table 1 for the other subsets. This prompt contains topics identified in the prior subset and instructions for referring to them.

## 4 Experiments

We investigate how well our methods perform topic modeling for short texts.

### 4.1 Dataset

We employ three tokenized datasets provided by Zhang et al. (2021): GoogleNewsT (Rakib et al., 2020), Tweet (Yin and Wang, 2016), and Stack-OverFlow.<sup>1</sup> Following Wu et al. (2022), the datasets are preprocessed as follows: (i) characters are converted to lower case; (ii) words with two or fewer letters are removed; (iii) words appearing fewer than five times are filtered out. We then split each preprocessed dataset into subsets for

<sup>1</sup><https://www.kaggle.com/competitions/predict-closed-questions-on-stack-overflow/data?select=train.zip>

Model	Tweet				GoogleNewsT				StackOverFlow			
	$K = 5$		$K = 15$		$K = 5$		$K = 15$		$K = 5$		$K = 15$	
	$C_v$	$TU$	$C_v$	$TU$	$C_v$	$TU$	$C_v$	$TU$	$C_v$	$TU$	$C_v$	$TU$
LDA	0.394	0.800	0.401	0.568	0.426	0.984	0.406	0.963	0.320	0.928	0.425	0.883
LDA <sub>Aug</sub>	0.445	0.968	0.436	0.856	0.411	0.984	0.381	0.981	0.360	0.920	0.508	0.952
TSCTM	0.393	<b>1.000</b>	0.467	0.997	0.333	<b>1.000</b>	0.374	<b>1.000</b>	0.244	<b>1.000</b>	0.313	<b>1.000</b>
TSCTM <sub>Aug</sub>	0.355	<b>1.000</b>	0.433	<b>1.000</b>	0.243	<b>1.000</b>	0.346	<b>1.000</b>	0.218	<b>1.000</b>	0.276	<b>1.000</b>
BERTopic	0.514	<b>1.000</b>	0.537	<b>1.000</b>	0.439	<b>1.000</b>	0.437	<b>1.000</b>	0.459	<b>1.000</b>	0.485	0.971
BERTopic <sub>Aug</sub>	0.535	<b>1.000</b>	0.526	<b>1.000</b>	0.412	<b>1.000</b>	0.417	<b>1.000</b>	0.460	<b>1.000</b>	0.489	0.955
<b>GPT-3.5</b> <sub>Par</sub>	0.476	0.992	0.532	0.900	0.571	0.960	0.535	0.913	0.312	0.864	0.496	0.913
<b>GPT-3.5</b> <sub>Seq</sub>	0.552	0.960	0.515	0.920	0.562	0.984	0.489	0.948	0.441	0.896	0.517	0.775
<b>GPT-4</b> <sub>Par</sub>	0.562	<b>1.000</b>	<b>0.576</b>	0.971	<b>0.618</b>	0.976	0.532	0.925	<b>0.466</b>	0.904	<b>0.571</b>	0.864
<b>GPT-4</b> <sub>Seq</sub>	<b>0.577</b>	0.992	0.551	0.976	0.556	0.944	<b>0.561</b>	0.963	0.318	0.744	0.532	0.853

Table 3: Topic coherence ( $C_v$ ) and diversity ( $TU$ ) results under 5 and 15 topics ( $K = 5$  and  $K = 15$ ). LLM<sub>Seq</sub> and LLM<sub>Par</sub> correspond to the parallel and sequential topic modeling methods with LLMs, respectively. MODEL<sub>Aug</sub> corresponds the performance of the model with data augmentation. The maximum  $TU$  is 1.000 when topic words are totally distinct from each other. The best scores are shown in **bold**.

Model	Tweet				GoogleNewsT				StackOverFlow			
	$K = 5$		$K = 15$		$K = 5$		$K = 15$		$K = 5$		$K = 15$	
	$DC$	$Fa$	$DC$	$Fa$	$DC$	$Fa$	$DC$	$Fa$	$DC$	$Fa$	$DC$	$Fa$
LDA	<b>0.337</b>	<b>1.000</b>	0.561	<b>1.000</b>	0.488	<b>1.000</b>	0.664	<b>1.000</b>	<b>0.684</b>	<b>1.000</b>	<b>0.842</b>	<b>1.000</b>
LDA <sub>Aug</sub>	0.307	<b>1.000</b>	<b>0.579</b>	0.997	<b>0.531</b>	<b>1.000</b>	<b>0.763</b>	<b>1.000</b>	0.659	<b>1.000</b>	0.838	<b>1.000</b>
TSCTM	0.176	<b>1.000</b>	0.388	<b>1.000</b>	0.405	<b>1.000</b>	0.740	<b>1.000</b>	0.141	<b>1.000</b>	0.480	<b>1.000</b>
TSCTM <sub>Aug</sub>	0.187	<b>1.000</b>	0.331	0.987	0.309	<b>1.000</b>	0.608	0.979	0.419	0.888	0.441	0.888
BERTopic	0.293	<b>1.000</b>	0.471	<b>1.000</b>	0.433	<b>1.000</b>	0.748	<b>1.000</b>	0.656	<b>1.000</b>	0.796	<b>1.000</b>
BERTopic <sub>Aug</sub>	0.303	<b>1.000</b>	0.468	<b>1.000</b>	0.422	<b>1.000</b>	0.749	<b>1.000</b>	0.637	<b>1.000</b>	0.795	<b>1.000</b>
<b>GPT-3.5</b> <sub>Par</sub>	0.213	<b>1.000</b>	0.384	0.994	0.321	0.968	0.585	0.952	0.636	<b>1.000</b>	0.694	<b>1.000</b>
<b>GPT-3.5</b> <sub>Seq</sub>	0.197	0.984	0.335	0.967	0.334	0.975	0.583	0.954	0.479	<b>1.000</b>	0.689	0.994
<b>GPT-4</b> <sub>Par</sub>	0.241	<b>1.000</b>	0.402	<b>1.000</b>	0.392	<b>1.000</b>	0.661	0.995	0.578	<b>1.000</b>	0.754	<b>1.000</b>
<b>GPT-4</b> <sub>Seq</sub>	0.224	0.983	0.403	0.994	0.373	<b>1.000</b>	0.660	0.951	0.554	0.931	0.626	0.883

Table 4: Document coverage ( $DC$ ) and factuality ( $Fa$ ) results under 5 and 15 topics ( $K = 5$  and  $K = 15$ ). Since baseline models without data augmentation discover topics based only on documents, the factuality values are 1.000.

our methods, setting the size at  $1000^2$  and truncating the remaining example. Table 2 shows the final statistics of the datasets we use. Note that baseline models take the union of subsets as input, and each subset contains different examples for each run.

## 4.2 Model

We evaluate our approaches with GPT-3.5 (gpt-3.5-turbo-0125) and GPT-4 (gpt-4-0125-preview) provided by the OpenAI API.<sup>3</sup> For baseline models, we employ the three models mentioned in Section 2: LDA<sup>4</sup>, TSCTM<sup>4</sup>, and BERTopic.<sup>5</sup> Additionally, we report the results of each baseline model with data augmentation. Regarding data augmentation techniques and the hyperparameters of TSCTM, we follow the original settings that were used in

<sup>2</sup>In preliminary experiments, we checked the performance of our methods with subset sizes of 250, 500, and 1000. See Appendix A.3.

<sup>3</sup>In preliminary experiments, we also tried Llama 2 (Touvron et al., 2023), but we found that it was not sufficiently controllable for its output to be used in our approach. See Appendix A.2.

<sup>4</sup><https://github.com/BobXWu/TopMost>

<sup>5</sup><https://maartengr.github.io/BERTopic>

prior research (Wu et al., 2022).<sup>6</sup>

## 4.3 Evaluation

We evaluate the models under the condition that the number of topics is 5 or 15, and the number of topic words for each topic is 5. For evaluation metrics, we employ two widely used metrics for topic quality and two new metrics to assess possible issues of LLMs, i.e., the possibility of outputting topics reflecting only a very limited documents or hallucinated topics not included in documents. We run each model five times and report the average scores.

**Topic Coherence and Diversity** Following Wu et al. (2022), we calculate the coherence value<sup>7</sup> ( $C_v$ , Röder et al., 2015) with Wikipedia for topic coherence, and the topic uniqueness ( $TU$ , Nan et al., 2019) to assess the diversity in the inferred topics.

**Document Coverage** We are concerned that LLMs infer topics that reflect only a very limited

<sup>6</sup>The details can be found in Appendix B.1.

<sup>7</sup><https://github.com/dice-group/Palmetto>

Model	Tweet		GoogleNewsT		StackOverFlow	
	$K = 15$		$K = 15$		$K = 15$	
	$Cv$	$DC$	$Cv$	$DC$	$Cv$	$DC$
GPT-3.5	0.532	0.366	0.517	0.569	0.464	0.634
GPT-4	0.580	0.395	0.523	0.665	0.519	0.747

Table 5: Average coherence ( $Cv$ ) and document coverage ( $DC$ ) of topics discovered by LLMs in parallel prompting without the merging process under 15 topics ( $K = 15$ ). For each subset, we take the average of the values in five runs.

documents. Thus, we propose the metric *document coverage*, which measures the extent to which discovered topics cover documents. Document coverage is defined as follows:

$$DC = \frac{\#(d_{ref} \text{ that contains at least one } w_{topic})}{\#(d_{ref})}$$

where  $d_{ref}$  is a document within the reference document collection, and  $w_{topic}$  is the topic word constituting the outputted topics. A higher  $DC$  means that discovered topics cover more reference documents. In this experiment, we use the preprocessed datasets without augmentation as references.

**Factuality** Another potential issue is hallucination, where topics discovered by LLMs may not be included in given documents. Therefore we introduce *factuality*, which measures the degree to which topic words are composed from the vocabulary in the reference documents. Factuality is defined as follows:

$$Fa = \frac{\#(w_{topic} \text{ present in at least one } d_{ref})}{\#(w_{topic})}$$

A higher  $Fa$  indicates that more topic words are composed from the vocabulary in the reference documents. Note that the factuality could be less than one in existing topic modeling with data augmentation due to word substitution using out-of-vocabulary words of the documents.

## 5 Results and Discussion

**Topic Quality** Table 3 shows that the topics discovered by our methods are relatively high-quality both in terms of coherence ( $Cv$ ) and diversity ( $TU$ ).<sup>8</sup> For coherence in particular, GPT-4 achieved the state-of-the-art performance in all settings, with up to 40 % improvement. For instance, the scores on GoogleNewsT have risen by 41% (from 0.439 to 0.618) and 28% (from 0.437 to 0.561), respectively, for each setting of the number of topics.

<sup>8</sup>Examples of topics are given in Appendix C.1.

**Document Coverage** Table 4 reports that LLMs showed relatively lower scores for document coverage ( $DC$ ) than the best baseline models. This means that the topics discovered by LLMs often cover fewer documents than those discovered by the baseline models. However, note that there is a trade-off between topic coherence ( $Cv$ ) and document coverage. For example, LDA<sub>Aug</sub> achieved the highest coverage on GoogleNewsT but showed the lowest coherence, with the exception of TSCTM and TSCTM<sub>Aug</sub>.

**Factuality** As shown in Table 4, LLMs showed lower scores for factuality ( $Fa$ ) than the baseline models, particularly those without augmentation. This indicates that some topic words output by LLMs are not included in the documents. However, their factuality loss was less than 5% in almost all settings. Furthermore, we analyzed these non-existent words and found that most were not problematic enough to mislead topic interpretation; these include synonyms, derivatives, and related words of the ones in the documents.<sup>9</sup> This suggests that LLMs do not generate hallucinated topics that would cause misinterpretation of the content.

**Parallel and Sequential Prompting** Table 3 and Table 4 show that the parallel prompting approach can identify topics with better coherence and document coverage than the sequential prompting one. To analyze the superior performance of the parallel approach, we calculated  $Cv$  and  $DC$  of topics before merging. Table 5 shows that  $Cv$  and  $DC$  scores before merging were worse than those of the parallel approach, demonstrating that the merging process can improve both their coherence and document coverage. On the other hand, we analyzed the transition of topics during the sequential approach and then observe that it tended to update the previously identified topic very little due to strict adher-

<sup>9</sup>Examples of non-existent words and analysis details are provided in Appendix C.2.



Model	#	Topics
BERTopic	#1	<b>kanye black thanksgiving west xbox</b>
	#2	china independence zone scotland air
	#3	hiv aarushi watkins ian woman
	#4	jellyfish robot seahorse flying methane
	#5	alzheimer brain infant risk gene
GPT-4 <sub>Par</sub>	#1	<b>kanye west kim kardashian parody</b>
	#2	<b>thanksgiving black friday shopping deal</b>
	#3	<b>xbox microsoft game console sale</b>
	#4	nokia lumia microsoft smartphone tablet
	#5	syria peace talk geneva conference
GPT-4 <sub>Seq</sub>	#1	<b>kanye west kim kardashian parody</b>
	#2	<b>black friday shopping thanksgiving deal</b>
	#3	<b>xbox game console playstation microsoft</b>
	#4	comet ison sun spectacular encounter
	#5	scottish independence salmond white paper

Table 6: Examples of topics discovered from GoogleNewsT when the number of topics and topic words is five, respectively. We have reordered the topics for illustrative purposes. **Bold topics** are mentioned in Section 5.

ence to our instructions, leading to lower document coverage compared with the parallel approach.<sup>10</sup>

**Qualitative Analysis** We conducted a qualitative analysis of the representative results that achieved the median topic coherence ( $C_v$ ) across five trials using the GoogleNewsT dataset under five topics and five topic words. Table 6 demonstrates that BERTopic, the best baseline model for  $C_v$ , has the potential to identify topics encompassing multiple themes, while our methods using LLMs discover highly consistent and distinct topics. For instance, topic #1 identified by BERTopic could be considered to contain three distinct themes (**Kanye West**, **Thanksgiving**, and **Xbox**), while GPT-4<sub>Par</sub> and GPT-4<sub>Seq</sub> effectively separated these into topics #1, #2, and #3, respectively.

## 6 Conclusion

In this study, we proposed two approaches to using LLMs for topic modeling: parallel prompting and sequential prompting. We implemented our methods on GPT-3.5 and GPT-4 and evaluated their performance on three datasets together with three existing topic models. In the evaluation, in addition to the well-known metrics for topic quality, we introduced two new metrics, document coverage and factuality, to assess the potential issues with LLMs reflecting only some documents or outputting hallucinated topics. The results showed that LLMs could find higher-quality topics than existing methods, and the impact of these issues was not

<sup>10</sup>Examples and further analysis are provided in Appendix C.3.

remarkable in practice. Future work will include improving our methods to enable topic assignment to each document.

## Acknowledgements

We thank the three anonymous reviewers for their helpful comments and suggestions, which improved this paper. This work was supported by JSPS KAKENHI Grant Number JP24H00809, Japan.

## References

- David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. **Latent dirichlet allocation**. *Journal of machine Learning research*, 3(Jan):993–1022.
- Rob Churchill and Lisa Singh. 2022. **The evolution of topic modeling**. *ACM Comput. Surv.*, 54(10s).
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. **BERT: Pre-training of deep bidirectional transformers for language understanding**. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Adji B. Dieng, Francisco J. R. Ruiz, and David M. Blei. 2020. **Topic Modeling in Embedding Spaces**. *Transactions of the Association for Computational Linguistics*, 8:439–453.
- Maarten Grootendorst. 2022. **Bertopic: Neural topic modeling with a class-based tf-idf procedure**. *Computing Research Repository*, arXiv:2203.05794. Version 1.

- Sosuke Kobayashi. 2018. [Contextual augmentation: Data augmentation by words with paradigmatic relations](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 452–457, New Orleans, Louisiana. Association for Computational Linguistics.
- Jan Kocoń, Igor Cichecki, Oliwier Kaszyca, Mateusz Kochanek, Dominika Szydło, Joanna Baran, Julita Bielaniec, Marcin Gruza, Arkadiusz Janz, Kamil Kanclerz, Anna Kocoń, Bartłomiej Koptyra, Wiktoria Mieszczenko-Kowszewicz, Piotr Miłkowski, Marcin Oleksy, Maciej Piasecki, Łukasz Radliński, Konrad Wojtasik, Stanisław Woźniak, and Przemysław Kazienko. 2023. [Chatgpt: Jack of all trades, master of none](#). *Information Fusion*, 99:101861.
- Chenliang Li, Haoran Wang, Zhiqian Zhang, Aixin Sun, and Zongyang Ma. 2016. [Topic modeling for short texts with auxiliary word embeddings](#). In *Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval*, page 165–174, New York, NY, USA. Association for Computing Machinery.
- Yishu Miao, Edward Grefenstette, and Phil Blunsom. 2017. [Discovering discrete latent topics with neural variational inference](#). In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, page 2410–2419.
- Feng Nan, Ran Ding, Ramesh Nallapati, and Bing Xiang. 2019. [Topic modeling with Wasserstein autoencoders](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6345–6381, Florence, Italy. Association for Computational Linguistics.
- OpenAI. 2023. [Gpt-4 technical report](#). *Computing Research Repository*, arXiv:2303.08774. Version 3.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744.
- Md Rashadul Hasan Rakib, Norbert Zeh, Magdalena Jankowska, and Evangelos Milios. 2020. [Enhancement of short text clustering by iterative classification](#). In *Natural Language Processing and Information Systems*, pages 105–117, Cham. Springer International Publishing.
- Nils Reimers and Iryna Gurevych. 2019. [SentenceBERT: Sentence embeddings using Siamese BERT-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Michael Röder, Andreas Both, and Alexander Hinneburg. 2015. [Exploring the space of topic coherence measures](#). In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, page 399–408, New York, NY, USA. Association for Computing Machinery.
- Akash Srivastava and Charles Sutton. 2017. [Autoencoding variational inference for topic models](#). In *International Conference on Learning Representations*.
- Dominik Stambach, Vilém Zouhar, Alexander Hoyle, Mrinmaya Sachan, and Elliott Ash. 2023. [Revisiting automated topic model evaluation with large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9348–9357, Singapore. Association for Computational Linguistics.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *Computing Research Repository*, arXiv:2307.09288. Version 2.
- Aaron van den Oord, Oriol Vinyals, and koray kavukcuoglu. 2017. [Neural discrete representation learning](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. [Document-level machine translation with large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16646–16661, Singapore. Association for Computational Linguistics.
- Xiaobao Wu, Anh Tuan Luu, and Xinshuai Dong. 2022. [Mitigating data sparsity for short text topic modeling by topic-semantic contrastive learning](#). In *Proceedings of the 2022 Conference on Empirical Methods*

*in Natural Language Processing*, pages 2748–2760, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Jianhua Yin and Jianyong Wang. 2016. [A model-based approach for text clustering with outlier detection](#). In *2016 IEEE 32nd International Conference on Data Engineering (ICDE)*, pages 625–636.

Dejiao Zhang, Feng Nan, Xiaokai Wei, Shang-Wen Li, Henghui Zhu, Kathleen McKeown, Ramesh Nallapati, Andrew O. Arnold, and Bing Xiang. 2021. [Supporting clustering with contrastive learning](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5419–5430, Online. Association for Computational Linguistics.

## A Preliminary Experiments

In preliminary experiments, we tested different prompts and subset sizes to determine which maximize the performance of our methods.

### A.1 Prompts

We first considered the **Par<sub>TM</sub>** prompt and then proceeded to the **Par<sub>Mrg</sub>** and the **Seq<sub>TM</sub>** prompts.

**Par<sub>TM</sub>** We checked three kinds of prompts, which are shown in Table 7. Finally, we tentatively selected a **Direct** prompt as a **Par<sub>TM</sub>** prompt, which achieved the highest performance. We also considered the effects from inserting the following phrases, which were expected to improve scores for topic coherence, diversity, and document coverage, respectively.

**Cv** “NOTE: Make top words for each topic likely to occur together in the documents”

**TU** “NOTE: Make the top words unique across topics.”

**DC** “NOTE: Maximize the number of documents that contain at least one of the top words.”

However, we found that none of these can positively influence LLMs’ performance in our methods. Therefore, we selected a **Direct** prompt without phrase insertion as the **Par<sub>TM</sub>** prompt.

**Par<sub>Mrg</sub>** Regarding the **Par<sub>TM</sub>** prompt, we created a **Base Par<sub>Mrg</sub>** prompt, which has a similar structure to the **Par<sub>TM</sub>** (Table 8). We then considered the insertion of the following phrases:

**Goal** “We aim to identify topics for the entire document set by merging the topic modeling results for each subset.”

**Detail** “NOTE: Outputs should reflect the topics before merging as much as possible. Output should contain topics that often appear before merging and not have ones that don’t appear much before merging.”

Experimental results showed our methods performed the best when we inserted both the *Goal* phrase and the *Detail* phrase into the **Base Par<sub>TM</sub>**.

Consequently, we employed a **Base Par<sub>TM</sub>** prompt with both phrases as the **Par<sub>TM</sub>** prompt for the parallel approach.

**Seq<sub>TM</sub>** Similar to the prompt for parallel, we first created a simple **Base Seq<sub>TM</sub>** prompt for the sequential approach in Table 8, after which we validated the effect from inserting the following phrases.

**Goal** “We aim to identify topics for the entire document set by sequentially updating tentative topics identified from each subset, considering topics identified just before from another subset.”

**Detail** “NOTE: Outputs should be the same as the previous topics as much as possible. You can change them minimally only when the given documents don’t include them much, and a new topic needs to be added to describe the documents.”

We also found that the insertion of both of the above phrases was most effective at improving the performance of the sequential method. Thus, we utilized a **Base Seq<sub>TM</sub>** prompt that incorporates both phrases as the **Seq<sub>TM</sub>** prompt for the sequential approach.

### A.2 Llama 2

In preliminary experiments, we also tried using Llama-2-7b-chat<sup>11</sup> and Llama-2-13b-chat<sup>11</sup> as LLMs for our methods and found that it is difficult for Llama 2 (Touvron et al., 2023) to perform topic modeling regardless of the prompts and the subset size we use. Table 9 shows the outputs of Llama 2 when given the **Par<sub>TM</sub>** prompt with a subset size of 100 on GoogleNewsT. Llama 2 could not make adequate output for the number of topics and topic words in line with our instructions, while GPT-3.5 and GPT-4 could do so consistently under identical settings.

### A.3 Subset Size

We used 250, 500, and 1000 as options for the subset size. It would be difficult for the subset size to exceed 1000 due to the context length of GPT-3.5 (gpt-3.5-turbo-0125), which we planned to use for the main experiments.

We ran the parallel and the sequential methods with GPT-3.5 on GoogleNewsT for each subset size. Table 10 presents the average scores of each method for five runs. There was a tendency for

<sup>11</sup><https://huggingface.co/collections/meta-llama/llama-2-family-661da1f90a9d678b6f55773b>



topic coherence to improve as the subset size increased, but we could not discern any tendency for the other metrics. We ultimately selected 1000 as the subset size because the performance of each model was relatively high in all metrics under that setting.

Note, however, that using our proposed methods with the subset size of 250 or 500 could enable discovery of competitive or higher-quality topics compared with the existing models shown in Table 3 and Table 4. This suggests our methods could perform well regardless of the context length of LLMs applied them.

## B Experimental Details

### B.1 Implementation Details

We run TSCTM for 200 epochs. In the case without data augmentation, we run it with temperature as 0.5 and weight contrast as 1.0. In the case with data augmentation, we run it with temperature as 0.07, weight contrast as 3.0, and same quant as 0.001. For data augmentation, we apply WordNet<sup>12</sup> and Contextual Augmenter<sup>3</sup> (Kobayashi, 2018) with 30% word replacement, and filtered low-frequency words as in the preprocessing. Each Augmenter randomly replaces words in an input text with synonyms defined by WordNet and with words predicted by BERT (Devlin et al., 2019)<sup>13</sup>, respectively. We utilized the original configurations of gpt-3.5-turbo-0125, gpt-4-0125-preview, and BERTopic without modification.

### B.2 Examples of Prompts

Table 11 shows examples of prompts used in the experiment.

## C Result Details

### C.1 Examples of Topics

Following Wu et al. (2022), we randomly select some examples of topics identified by LDA, BERTopic, and our proposed methods with GPT-4.

### C.2 Examples of Topic Words Not Included in the Documents

Table 13 shows examples of words not included in the documents outputted in topic modeling on GoogleNewsT. The bold portion of the GPT-3.5 outputs are the names of entities (e.g., **brncos**,

**gree**, and **watson**) or words that do not exist in the real world (e.g., **dorffiefskee**). Such words are considered harmful because they may induce misinterpretation of topics. However, only a small number of such words were found, and most of them were synonyms, derivatives, or related words in the documents.

### C.3 Examples of the Processing

Table C.3 shows specific the concrete examples of topics identified for each subset and the final output to demonstrate the processing in our methods. In the parallel approach, we find that LLM reasonably merges topics from each subset. For instance, bold topics in each subset are merged into one topic in the final output, using words from both subsets. On the other hand, in the sequential approach the final output is the same as the topics for the first subset except for the one pair of bold words. This indicates that LLMs with the the sequential approach could too strictly retain topics from the previous subset, and thus they cannot output topics that sufficiently reflect the entire set.

## D Limitations

We do not thoroughly consider whether pre-training and instruction-tuning datasets of GPT-3.5 and GPT-4 might contain the datasets used in this study. Since topic modeling is an unsupervised task and we change the order of the samples randomly, we do not consider them able to utilize their knowledge about these datasets in our experiment.

<sup>12</sup><https://github.com/makcedward/nlpaug>

<sup>13</sup><https://huggingface.co/bert-base-uncased>

ID	Candidates for the Base Prompt Template
<b>Direct</b>	<p><b>Write the results of simulating topic modeling for the following documents</b>, each starting with "#."  Assume you will finally identify [NUM_TOPICS] topics and use 5 top words for each topic.  NOTE: Outputs must always be in the format "Topic k: word word word word word" and nothing else.</p> <p>""</p> <p>[DOCS]</p> <p>""</p>
<b>Indirect</b>	<p>Discover latent [NUM_TOPICS] topics in the following documents, each starting with "#."  For each topic, write 5 words extracted from input texts to show its meanings.  NOTE: Outputs must always be in the format "Topic k: word word word word word" and nothing else.</p> <p>""</p> <p>[DOCS]</p> <p>""</p>
<b>Direct<sub>reverse</sub></b>	<p>""</p> <p>[DOCS]</p> <p>""</p> <p><b>Write the results of simulating topic modeling for the above documents</b>, each starting with "#."  Assume you will finally identify [NUM_TOPICS] topics and use 5 top words for each topic.  NOTE: Outputs must always be in the format "Topic k: word word word word word" and nothing else.</p>

Table 7: Candidate prompts for Par<sub>TM</sub>. [DOCS] and [NUM\_TOPICS] are replaced by a subset of documents and by the number of topics.

ID	Base Prompt Template
<b>Base Par<sub>TM</sub></b>	<p>Write the results of merging the following topic modeling results for each subset of the document set.  Each result starts with "- n" and its topics start with "#"</p> <p>""</p> <p>- 1</p> <p>[TOPICS]</p> <p>- 2</p> <p>[TOPICS]</p> <p>- 3</p> <p>...</p> <p>""</p>
<b>Base Seq<sub>TM</sub></b>	<p>Write the results of simulating topic modeling for the following documents, each starting with "#."  Make the most use of the following topics previously identified from another set of documents, each starting with "Topic k:"</p> <p>""</p> <p>[TOPICS]</p> <p>""</p> <p>Assume you will finally identify [NUM_TOPICS] topics and use 5 top words for each topic.  NOTE: Outputs must always be in the format "Topic k: word word word word word" and nothing else.</p> <p>""</p> <p>[DOCS]</p> <p>""</p>

Table 8: Base prompts for the parallel and sequential methods. [DOCS], [TOPICS], and [NUM\_TOPICS] are replaced by a subset of documents, previously identified topics, and the number of topics, respectively.

Model size	Examples of Llama 2 Output
<b>7B</b>	Topic 1: Top words: relief, challenge, face Topic 2: Top words: welker, concussion, test Topic 3: Top words: live, stream, champion, league Topic 4: Top words: bargain, black, friday, shopping Topic 5: Top words: scotland, independence, white, paper Note: Each topic is represented by 5 top words, which are the most frequently occurring words in the given documents.
<b>13B</b>	Topic 1: Disasters and Relief Efforts Topic 2: Sports and Injuries Topic 3: Technology and Gadgets Topic 4: Politics and Leadership Topic 5: Entertainment and Celebrities

Table 9: Examples of Llama 2 outputs when we provide  $\text{Par}_{\text{TM}}$  on GoogleNewsT under the conditions that the number of topics and topic words is five and the subset size is 100.

Subset Size	$C_v$	$TU$	$DC$	$Fa$	Subset Size	$C_v$	$TU$	$DC$	$Fa$
250	0.531	0.936	<b>0.241</b>	<b>1.000</b>	250	0.524	0.976	<b>0.198</b>	<b>0.992</b>
500	<b>0.572</b>	0.896	<b>0.241</b>	<b>1.000</b>	500	0.529	<b>0.992</b>	0.193	0.976
1000	0.571	<b>0.960</b>	0.213	<b>1.000</b>	1000	<b>0.562</b>	0.984	0.197	0.984

(a) Parallel

(b) Sequential

Table 10: Results of the parallel and sequential methods under five topics on GoogleNewsT for subset sizes of 250, 500, and 1000. The best scores are shown in **bold**.

ID	Prompt Example
<b>Par<sub>TM</sub></b>	<p><b>Write the results of simulating topic modeling for the following documents</b>, each starting with "#."  Assume you will identify 5 topics and use 5 top words for each topic.  NOTE: Outputs must always be in the format "Topic k: word word word word word" and nothing else.  """"</p> <p># philippine typhoon relief effort face challenge  # wes welker concussion test bronco  # basel chelsea live stream champion league watch  ...  # discuss black friday shopping secret  """"</p>
<b>Par<sub>Mrg</sub></b>	<p>We aim to identify topics for the entire document set by merging the topic modeling results for each subset.  <b>Write the results of merging the following topic modeling results for each subset of the document set.</b>  Each result starts with "- n" and its topics start with "#"  """"</p> <p>- 1  # comet ison thanksgiving sun solar  # kanye west bound parody video  # nokia lumia release mobile device  # black friday shopping thanksgiving sale  # alec baldwin msnbc cancellation defends  ...  - 11  # nokia lumia sale december phone  # kanye west kim kardashian taylor  # black friday deal best sales  # irs rule political activity tax  # bronco patriot win game rivalry  """"</p> <p>Assume you will finally identify 5 topics and use 5 top words for each topic.  NOTE: Outputs should reflect the topics before merging as much as possible. Output should contain topics that often appear before merging and not have ones that don't appear much before merging.  NOTE: Outputs must always be in the format "Topic k: word word word word word" and nothing else.</p>
<b>Seq<sub>TM</sub></b>	<p>We aim to identify topics for the entire document set by sequentially updating tentative topics identified from each subset, considering topics identified just before from another subset.  Write the results of simulating topic modeling for the following documents, each starting with "#."  <b>Make the most use of the following topics previously identified from another set of documents, each starting with "Topic k:"</b>:  """"</p> <p>Topic 1: kanye west kim kardashian bound  Topic 2: xbox black friday cyber monday  Topic 3: hewlett packard nokia lumia company  Topic 4: dancing star finale winner season  Topic 5: syria peace talk china air  """"</p> <p>Assume you will finally identify 5 topics and use 5 top words for each topic.  NOTE: Outputs should be the same as the previous topics as much as possible. You can change them minimally only when the given documents don't include them much, and a new topic needs to be added to describe the documents.  NOTE: Outputs must always be in the format "Topic k: word word word word word" and nothing else.  """"</p> <p># spacex falcon launch attempt  # taylor swift princess gown winter white  # redbox instant window phone appears nokia exclusive  ...  # google backed company selling dna analysis kit ordered sale  """"</p>

Table 11: Examples of prompts used as Par<sub>TM</sub>, Par<sub>Mrg</sub>, and Seq<sub>TM</sub> for topic modeling on GoogleNewsT under five topics.

Model	Examples of Topics
LDA	xbox microsoft game patriot bronco nokia lumia oldboy launch google kobe bryant chelsea lakers basel
TSCTM	macy parade hanukkah thanksgiving travel china zone african japan johansson bronco patriot packer welker illinois
BERTopic	china zone air nsa porn methane ant emission fire burning thanksgiving friday black comet parade
GPT-4 <sub>Seq</sub>	wes welker nfl concussion game nokia lumia window phone december nfl season game player concussion
GPT-4 <sub>Par</sub>	san andreas mobile game release nokia lumia tablet smartphone launch thanksgivukkah hanukkah holiday feast rare

Table 12: Examples of topics discovered from GoogleNewsT under 15 topics.

Model	Examples of Topic Words Not Included in the Documents
TACTM <sub>Aug</sub>	twelvemonth sink railway blowout
<b>GPT-3.5<sub>Seq</sub></b>	<b>dorffiefskee broncos</b> patriots health advancement ocean guilty france legal attorney
<b>GPT-3.5<sub>Par</sub></b>	<b>gree watson</b> advertisement boat funding attorney declared refugees crash digital

Table 13: Examples of topic words not included in the documents when topic modeling on GoogleNewsT.

<p><b>Subset 1</b></p> <p>fishing fish bass fly report <b>superbowl commercial bowl super best</b> king speech oscar nomination award facebook privacy setting user change acai berry weight loss diet plan</p>	<p><b>Subset 1</b></p> <p>fishing commercial superbowl fly bass facebook privacy setting user <b>setting</b> king speech oscar nomination award berry acai weight diet loss christina aguilera national anthem super</p>
<p><b>Subset 2</b></p> <p>fishing fish fly book saltwater <b>superbowl commercial doritos pepsi volkswagen</b> king speech oscar nomination award best acai berry weight loss diet plan christina aguilera national anthem super bowl</p>	<p><b>Final Output</b></p> <p>fishing fly superbowl commercial bass facebook privacy setting user <b>security</b> king speech oscar nomination award acai berry weight diet loss christina aguilera national anthem super</p>
<p><b>Final Output</b></p> <p>fishing fish fly bass saltwater <b>superbowl commercial bowl pepsi doritos</b> king speech oscar nomination award acai berry weight loss diet health facebook privacy setting user change</p>	

(a) Parallel

(b) Sequential

Table 14: Topics identified for each subset and the final output by each method using GPT-4 on Tweet under five topics. **Bold words** are mentioned in Appendix C.3.