# Topic Segmentation Using Generative Language Models

**Pierrre Mackenzie[1,2], Maya Shah[1,3], Patrick Frenett[1],**

[1]Adarga, [2] University of Edinburgh, [3] University of Exeter,
**Correspondence:** `lardet[dot]pierre[at]gmail.com`

## Abstract

Topic segmentation using generative Large Language Models (LLMs) remains relatively unexplored. Previous methods use lexical or semantic similarity between parts of a document to decide on boundaries but they lack the long range dependency and vast knowledge contained in LLMs. In this work we propose a new prompting strategy and compare to methods based on semantic similarity. We also support the adoption of a less commonly used evaluation metric: Boundary Similarity. Results show that LLMs can be more effective segmenters than existing methods, but issues remain to be solved before they can be relied upon for topic segmentation.

## 1 Introduction

### 1.1 Topic Segmentation Task Definition

Topic segmentation is the problem of dividing a string of text into constituent parts or 'segments'. Each segment should be semantically self-contained such that it is about one thing. The precise definition of a segment should be dependent upon the specific use cases. For this work, segment boundaries always lie on sentence boundaries.

We can interpret segmentation as a binary classification task. Given a list of input sentences $S$, of length $n$, the model must decide whether there exists a segment boundary between each pair of adjacent sentences. There are $n - 1$ possible boundaries, and therefore the solution space is $2^{n-1}$. Formally, the model must find a mapping $f$ from the list of sentences to a binary vector of length $n - 1$: $f(S) = \mathbf{y}$ where $\mathbf{y} = \{y_1, y_2 \dots, y_{n+1}\}$ and each of the $y_i \in \{0, 1\}$.

Relative to generative tasks such as summarisation, the space of possible solutions is much smaller, but the problem remains subjective as where a boundary should lie can be ambiguous. Frequently, humans cannot agree on a correct solution (Hearst, 1997).

### 1.2 Motivation

Topic segmentation can be important as a pre-processing step before some other NLP task, or can be important in its own right.

Tasks such as information retrieval, long-document summarisation and classification can all benefit from first being broken down by topic. For many open source or resource-constrained models, context windows limit the size of input for such tasks. This can be constraining especially for smaller models (Chung et al., 2024). Although context windows can be increased (Chen et al., 2023) and newer models are consistently increasing context lengths with ever-more-powerful GPUs, this alone does not fix all problems as LLMs do not fully utilise long context windows (Liu et al., 2024) (Petroni et al., 2020).

Furthermore, segmentation can be imporant for its own sake. One might use segments to generate a contents page for a long document, create individual summaries of segments within a document or provide titles for each segment. Alternatively, one might use segmentation to break down a long document into smaller, more manageable parts for a user to read or as a citation for RAG, in which the model must generate an answer to a question based on a segment of text.

Our application of interest is to use topic segmentation as both a pre-processing step for a document understanding pipeline and as an important step in its own right. This pipeline is used to extract structured information from unstructured text, and the segmentation step is used to break down the document into smaller, more manageable parts. When presented to a user, only semantically self-contained sections of a long document are presented to a user, which can be more easily understood and acted upon.

## 1.3 Related Work

(Xing, 2024) provides a recent, broad overview of topic segmentation. Their work discusses the history of segmentation, dialogue/hierarchical/multimodal segmentation, as well as exploration of large language models for topic segmentation.

Methods prior to the era of neural networks often relied on lexical similarity metrics. An influential framework introduced by (Hearst, 1997) in their paper *TextTiling* involves computing lexical similarity scores between adjacent sentences before boundaries are placed where similarity scores are lowest. A variety of such lexical similarity metrics were proposed (Galley et al., 2003; Eisenstein, 2009), but were proven to be too superficial compared to semantic similarity that came after it. Such a framework is still in use today but with semantic similarity metrics generated by embeddings from language models like SentenceBERT (Reimers and Gurevych, 2019).

Neural networks have seen recent use for topic segmentation. The task can be conceived as a sequence labelling task. Naturally, BiLSTMs have been employed (Wang et al., 2016) as well as attention-based methods (Lukasik et al., 2020). Many methods leverage a hierarchical framework in which one model extracts sequential features from the text and another model constructs sentence boundaries based on these features. These models have included LSTMs (Koshorek et al., 2018; Badjatiya et al., 2018) or more recently transformers (Glava s and Somasundaran, 2020; Lo et al., 2021). The most recent work leverages pre-training for feature extraction followed by a dedicated transformers classifier. However, these methods may not take full advantage of the rich representations contained in the largest models as there is a bottleneck introduced by the hierarchical structure. We propose that an end-to-end method may be more effective.

LLMs have been the state of the art for a variety of NLP tasks for a number of years now such as as translation, summarisation and information extraction. However, there has been little research into the usage of LLMs for topic segmentation. A loss-based approach was proposed by (Feng et al., 2021) for annotation of dialogues and then extended by (Xing, 2024). This approach computes the mean NLL (Negative Log Likelihood) token-wise loss of a pre-trained model when predicting the tokens in a sentence. Boundaries are then placed where this loss is above some threshold, indicating that the sentence was hard to predict. While this approach is almost directly able to leverage the vast knowledge contained in LLMs and a tune-able threshold is appealing, it is also an arbitrary threshold, and the method relies on the assumption that all information about segment boundary location can be expressed by the next token prediction loss, even in nuanced cases.

The only work we found to explore the use of directly prompting an LLM for topic segmentation was (Xing, 2024). They find that prompting ChatGPT is the best dialogue segmentation model unless the input exceeds ChatGPT's limit. They propose two prompting methods: one which asks the LLM to return the original model with special characters delimiting the segment boundary, and a second in which the LLM is asked to return pairwise semantic coherence scores in the range of 0 to 1 for adjacent sentences. The first method falls short of our requirements as there is no guarantee that the model will return the original document unedited and is wasteful of tokens. The second method is methodologically flawed as the returned scores have no guarantee of directly corresponding to semantic coherence *for topic segmentation*, are very hard to interpret and there will still be a bottleneck in turning scores into segment boundaries. We propose a new prompting method which we believe may be more effective, less wasteful of tokens, and not limited by the maximum token limit of the model.

In this work, we compare our new prompting work to our own existing method which uses SentenceBERT embeddings and cosine similarity to decide on segment boundaries. We do not compare to other prompting methods nor methods based on hierarchites of transformers due to resource constraints. However, our focus is on whether our new prompting strategy is a feasible replacement for methods which used semantic similarity. Further work should more thoroughly compare our new prompting method to other prompting schema, to loss-based approaches and to a wider variety of hierarchical methods.

## 2 Method

### 2.1 Datasets

There are 4 datasets used in the experiments in this work: a small human-annotated dataset, a scrape of English wikipedia, a 'concatenated' wikipedia

scrape and a synthetic GPT3.5 generated dataset.

### 2.1.1 Human-Annotated Dataset

Lacking the resources to create a large dataset annotated by humans, this work uses a small dataset of 10 documents which was manually segmented. These documents are a mix of news articles, wikipedia articles, and miscellaneous documents such as podcast transcripts and scientific reports. This was intended to represent varying difficulties of segmentation and to provide examples which could be manually inspected to interpret segmenter behaviour. Only one annotator was used and the segments were created by hand. While these documents were long and the documents were segmented with care, further work should involve multiple annotators on a much larger dataset to ensure consistench and reliability.

### 2.1.2 Wikipedia Dataset

A plain text English wikipedia scrape[1] where headings are delimited by special characters was employed in our experiments. The articles were then automatically segmented based on these headings and filtered to remove articles with very few segments, too short segments, or too many artefacts such as tables and figures. After this process, there remained approximately 1000 segmented articles. We generate two versions of this dataset: one with headings removed, and one with headings included. We evaluate models on both versions of the dataset, but include only the version with headings removed in the final results.

### 2.1.3 Concatenated Wikipedia Dataset

We randomly sampled segments from the previous Wikipedia dataset and concatenated them in order to form new incoherent articles, with segments drawn from completely different domains. Intuitively, this dataset should be easier to segment as there are no semantic links between segments.

### 2.1.4 Synthetic Dataset

The final dataset was generated synthetically by GPT-3.5. The source data was a mix of CTC sentinel data[2], UN-Peacekeeping corpus[3] and an internal dataset at Adarga. These documents were segmented by querying OpenAI's API with the model `gpt-3.5-turbo-16k` using the overlapping

prompt schema defined in section 2.3. This dataset was used exclusively for fine-tuning the *FlanT5* (3.1.5) model. It was not used for evaluation as the results would be biased in favour of the generative models which generated or were trained on the data.

## 2.2 Evaluation

The methodology in evaluating the quality of a model's segmentation in this work follows Chris Fournier's *Evaluating Text Segmentation using Boundary Edit Distance* (Fournier, 2013). This work uses edited versions of the Boundary Similarity proposed along with associated information recall metrics: 'Boundary Precision' and 'Boundary Recall'.

The boundary edit distance algorithm pairs matches of segment boundaries between a reference segmentation and a hypothesised segmentation regardless of distance between matches. Exact matches score 1 and boundaries without matches score 0, whilst matches within some distance $n$ score partial points as a function of distance. For this work, the function always decreases score linearly as the boundary gets farther away from the true boundary. Boundary Similarity (B) is the mean score for all these matches, while Boundary Precision/Recall (BP/BR) measure the proportion of the hypothesis/reference boundaries which are matched, respectively.

A model is punished for both missing a boundary and for placing a boundary where there is none in a symmetric and intuitive manner. For further justification of the usage of the boundary edit distance algorithm as opposed to more traditional metrics such as WD and Pk (Pevzner and Hearst, 2002), see *Evaluating Text Segmentation using Boundary Edit Distance* (Fournier, 2013), or see our own investigations at segmenter evaluation metrics.

## 2.3 LLM-Based Text Segmentation

How can we get LLMs to output segment boundaries?

We might first consider passing in the input text as a prompt and asking the LLM to copy out the text, adding markers indicating where it has placed boundaries as is suggested by (Xing, 2024). However, not only is this wasteful of tokens, especially if the input text is long, but crucially, the LLM may fail to copy the input accurately, may change the formatting, and we found through qualitative experimentation that boundary placement

---

[1] https://www.kaggle.com/datasets/ltcmdrdata/plain-text-wikipedia-202011/data

[2] https://ctc.westpoint.edu/ctc-sentinel/

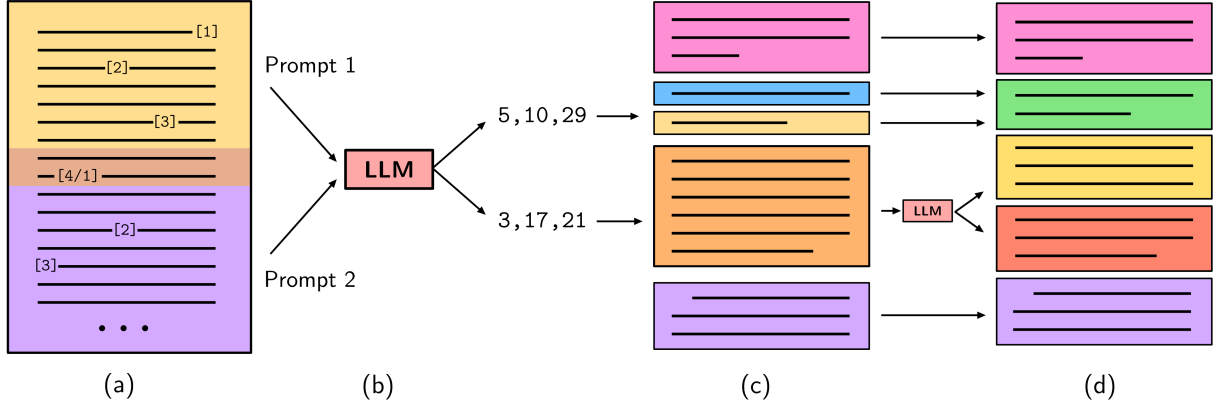[3] https://peacekeeping.un.org/en/reports

Figure 1: The overlapping and recursive prompting strategy for segmentation. In (a), a long document is split into overlapping sections with sentence boundaries enumerated. In (b), each section is segmented by the LLM. In (c), the segments are joined. Finally, in (d), each segment is validated to ensure it is not too long or too short. Segments that are too short are glued to neighbouring segments, and segments that are too long are recursively segmented.

was not any better than our final method. These problems are addressed by (Xing, 2024) through repeated prompting until the input and output sequence lengths match but this provides no guarantees. In our use case, guaranteeing that the input data would remain the same was essential so we opted for a different prompting strategy.

### 2.3.1 Prompting Method

We first annotate the text with indices between each sentence. Here is an example: 'Hello World. [1] The sky is blue. [2] The sun is is yellow. [3] The grass is green. [4] Machine learning is a rapidly evolving field'. We then ask the LLM to return a list of indices corresponding to boundaries. In the previous example, the ideal response might be '1,4'. In practice, the texts and list of segment boundaries are much longer.

We add a system prompt which describes the segmentation task, desired output format and primes the model to be a talented linguist. We also add a variety of short examples in line with the few-shot prompting technique (Brown et al., 2020) which improved performance. This did not necessarily reflect the fact that the LLM learned how to segment better. Instead, through manual testing, we suspect that it learned the ideal segment length and amount of information that should be contained within a segment, which was implicitly contained in the few-shot prompt (and also the testing datasets, therefore increasing performance). This suggests that different implicit definitions of a segment could be imposed by a few-shot prompt to a LLM, dependent upon use case.

The full method is shown in Figure 1 example prompt can be found in Appendix A.

### 2.3.2 Overlapping Prompts

A limitation also identified by (Xing, 2024) is that ChatGPT's segmentations were limited to the context window of the model. While our prompting strategy uses fewer tokens as the model need not copy out the input text, the same limitation applies. Therefore, we propose a simple overlapping prompt strategy to overcome this limitation.

Our prompting method works so long as the input text is within the context window of the LLM. At the time of experimentation, and with the use of gpt-3.5-turbo, we had a limit of 16k tokens. Many of our input texts exceeded this limit. Therefore, we need to split up these long documents into smaller chunks that can be processed by the LLM. However, this cannot be done by simply splitting the text at the nearest sentence before every 16k new tokens for two reasons. First, we do not know whether this sentence boundary should serve as a segment boundary. Second, the LLM loses valuable context which helps to choose where to place boundaries at the extremes of the 16k tokens.

Therefore, we instead send 16k prompts with some overlap. The overlap is calculated as twice the maximum segment length. In our experiments, we set a maximum segment length of 750 tokens, thus there is an overlap of 1500 tokens between prompts. We must then decide which boundaries to accept in this overlapping region. Given two generations which were prompted by 1500 overlapping tokens, we choose to accept the segment bound-

aries contained within the first 750 tokens of the overlapping section from the first generations and the boundaries in the final 750 tokens from the second prompt. We did not experiment with involving the responses from both outputs, but found no sign of worse performance towards segment boundaries.

While this method is wasteful of up to 1500 tokens per prompt, this is a small enough fraction of the 16k context that we were satisfied with the solution. If maximum segment lengths were much longer, the context may need to be limited more severely.

### 2.3.3 Segment Validation

We also performed some validation on the segments returned by the LLM. This primarily involved verifying that the returned segments are within a maximum and minimum segment length. Segments that are too short (for example, a model would sometimes return just a heading) were concatenated with another segment and segments that are too long were recursively segmented by the same model. This recursive segmentation was done with another prompt that asks the model to generate only one boundary. Again, we use a few-shot prompting strategy. This way, a segment that is too long will be split in two recursively until all segments are within the specified lengths.

## 3 Experiments

### 3.1 Models

#### 3.1.1 Baselines

We use two naive baselines as a point of reference. First, a segmenter which splits every $n$ sentences. Based on preliminary testing, we decided to split every 5 sentences which we call the *Split5* segmenter. We also define a *RandomF0.1* segmenter which splits at 10% of boundaries, placed randomly, with each potential boundary equally likely.

#### 3.1.2 BERT Segmenter

Our existing method generates a sequence of sentence similarities using sentence embeddings generated by (Reimers and Gurevych, 2019). Similarities are calculated as a weighted sum of the cosine similarities to the previous $n$ sentences. Ideal boundaries are then generated as troughs in the sequence of similarities, before post-processing of segments ensures that no segment is too long or too short, either in sentence length or in token length. We name this the *BERT* segmenter.

### 3.1.3 BERT-Graph Segmenter

(Costacurta, 2023) describes 'text tiling' (topic segmentation) also using BERT-generated similarity scores. This is instead followed by graph clustering to find the best segments and post-processing to ensure that the clusters return valid segments. Further details of this method can be found in the linked article. This model is called *BERTGraph* segmenter in our experiments. Code for this model was copied from the repository linked in the article.

### 3.1.4 GPT-3.5

OpenAI's `gpt-3.5-turbo-16k` model was queried using the prompting and segment validation strategy defined in Section 2 A. We call this model *GPT3.5*. We used the largest model available to us, which was the 16k token context window. Note that queries were made during the summer of 2023; the cheapest tier of model available from Open-AI has increased in capabilities since then. Given the nature of the task, we choose to use deterministic outputs from the model rather than sampling by setting topk $= 1$. We are looking for the best possible segmentations rather than a variety of possible segmentations. This also helps with reproducibility of segmentations.

### 3.1.5 Flan-T5-Finetuned

Our intended application of segmentation requires us to use our own model rather than an API. Therefore we fine-tune Google's Flan-T5 large (Chung et al., 2024) using LORA (Hu et al., 2022) on a combination of wiki, concatenated wiki and synthetic segmentations generated by GPT-3.5. We call this model *FlanT5*. We used the same deterministic output generation parameters as with *GPT3.5*. Because the model has been fine-tuned, we use a much shorter prompt with a brief instruction and no few-shot examples. Additionally, given the model's shorter context window, we have many more overlapping prompts per document.

### 3.2 Quantitative Results

Models were tested on the human-annotated dataset, the wikipedia dataset, the concatenated wikipedia dataset and a test-partition of the synthetic dataset. The fine-tuned model was not trained on this partition. We do not base our conclusions on results from the synthetic dataset as they would be biased in favor of the generative models which generated the segmentations. Indeed, while performance of the LLMs on the synthetic dataset

| | Boundary Similarity ($n = 2$) | | | | Boundary Precision and Recall | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Human | Wiki | Conc-Wiki | Synthetic | Human | | Wiki | | Conc-Wiki | |
| | | | | | BP | BR | BP | BR | BP | BR |
| *GPT3.5* | **0.38** | **0.25** | 0.29 | **0.35** | **0.51** | **0.60** | 0.36 | **0.55** | 0.42 | **0.63** |
| *FlanT5* | 0.25 | 0.24 | 0.41 | 0.33 | 0.38 | 0.46 | **0.43** | 0.37 | 0.65 | **0.63** |
| *BERTGraph* | 0.20 | 0.15 | 0.45 | 0.21 | 0.47 | 0.25 | 0.39 | 0.21 | 0.79 | 0.54 |
| *BERT* | 0.18 | 0.09 | **0.46** | 0.18 | 0.33 | 0.39 | 0.23 | 0.37 | **0.91** | 0.50 |
| *RandomF0.1* | 0.09 | 0.11 | 0.10 | 0.10 | 0.16 | 0.21 | 0.34 | 0.16 | 0.20 | 0.16 |
| *Split5* | 0.13 | 0.19 | 0.19 | 0.23 | 0.17 | 0.48 | 0.33 | 0.39 | 0.25 | 0.52 |

Table 1: Results for Boundary Similarity, Boundary Precision (BP) and Recall (BR) with $n = 2$ for each model and dataset. Best results for each metric are highlighted in **bold.**

is far from perfect (a boundary similarity of less than 0.5 with $n = 2$), it is disproportionately better than other models compared to other datasets.

We evaluate primarily using the previously discussed boundary similarity metric with $n = 2$. Results can be seen in Table 1. Evaluation was also run with boundary similarity $n = 5$ but the metric becomes noisier with higher values of $n$. Indeed, the naive baselines report better performance than some other models with high enough $n$.

Due to resource constraints, we could not test *GPT3.5* on the full wikipedia or full concatenated wikipedia datasets. Instead, we took the largest subset that fit within resource constraints. We evaluated all other models on the full datasets to verify that similar results are obtained. Full results on all evalutation metrics can be found online at the following links for both the smaller smaller and larger datasets.

We see that *GPT3.5* outperforms all other models on the human-annotated dataset, the wikipedia dataset, and the synthetic dataset, with the other LLM *FlanT5* coming in second in the same datasets. The close performance between the models on the synthetic dataset implies that *FlanT5* is a good approximation of *GPT3.5* for this task, on boundaries generate by *GPT3.5*. However, the biggest performance gap is on the human-annotated dataset. Although this dataset is small, it represents a meaningful distribution shift and, unsurprisingly, the fine-tuned model is unable to generalize as well as the base model.

Interestingly, we see that both BERT-based models are superior in the purportedly easier task posed by the concatenated wikipedia dataset. We theorise that these models are better at finding clear bound-

aries between different domains, but struggle with more nuanced segment boundaries found in a news article, for example. Whereas, the LLMs are better at finding more nuanced segments that a human might have produced but are not significantly better at finding more clear-cut boundaries.

We also looked at precision and recall to better characterize the behavior of each model. Results can be found in the second half of Table 1. *GPT3.5* has the highest precision and recall in the human-annotated dataset and the highest recall on both the wiki and concatenated wiki datasets. This property of having high recall is true in general for the LLMs, even in the concatenated wiki dataset where the BERT models have higher overall boundary similarity. By contrast, precision scores are closer in general and much better for the BERT models on the concatenated wiki dataset. This suggests that in general the BERT models are more hesitant in placing boundaries but when they do they are more likely to be correct. It should be noted that this behaviour in both the LLMs and BERT models is subject to the prompt engineering and segment processing/validation procedures used.

### 3.3 Qualitative Evaluation

Manual testing and inspection of the segmentations produced by the models was also conducted, both during prompt engineering and after the models were trained.

Through manual inspection of segmentations with the human-annotated dataset, we found that *GPT3.5* generally found the boundaries which seemed reasonable from a human perspective, especially for simple documents like short news articles. The fine-tuned *FlanT5* model imitated this be-

havior, but was less consistent. BERT segmenters would find reasonable segments, but after the manual gluing and splitting procedure, would often lead to off-by-1 errors.

For documents with far more complex or nuanced text such as a podcast transcript, or with messy data like tables or artefacts from pdf to text conversions, *GPT3.5* would sometimes return indices with a regular pattern. For example, the LLM might return '$[1, 15, 22, \ldots, 76, 79, 82, 85, 88, \ldots, 184, 187, \ldots]$'. Often, the pattern would continue far beyond the number of sentences in the input indicating that the model became stuck in a regular pattern. Perhaps better prompt engineering, a more rigorous data-processing procedure, the use of newer models or the use of better generation parameters would help. However, our current approach was resource constrained but still required the ability to pass noisy documents to the model. A more thorough investigation of the logits computed by the model is required to understand how and when this occurs, and how to mitigate it.

## 4   Conclusion

Our work empirically compares generative LLMs with previous methods which use BERT embeddings and cosine similarity. In order to be more token-efficient and to provide guarantees that the original document will be unaltered, we propose a new overlapping prompt schema which centres on asking the LLM to return a list of indices corresponding to segment boundaries. We also support the use of boundary similarity and its associated information recall metrics as an evaluation for topic segmentation. Results indicate that LLMs can be more effective segmenters than existing methods where more nuanced segmentations are required. On the contrary, when the input is noisy or the segment boundaries are clear BERT-based methods are more reliable. Future work should focus on addressing highlighted issues with LLMs such as the regular patterns found in segmentations and conduct a more thorough comparison of different prompting methods or loss-based approaches. Lastly, larger human-annotated datasets of a range of media should be constructed to better assess generalisation capabilities.

## Ethics Statement

There were two sources of data in this work. Wikipedia data is available under under the Creative Commons Attribution-ShareAlike 4.0 International License (CC BY-SA). The other data sources are proprietary can cannot be made public.

All models used were access either by (paid) public API (*GPT3.5*) or are open source models (*FlanT5*, *BERT*, *BERTGraph*). The models were used in accordance with the terms of service of the respective providers.

## Limitations

This work is comparable to a part of the work contained in (Xing, 2024). However, this work is currently in pre-print as a PhD thesis which was made public in 2024. The research in this report is based on work from 2023, prior to the work from (Xing, 2024). This is why none of our experiments directly compare our method to their prompting methods or loss-based approaches. Our primary goal was to compare to the previously existing approach at Adarga, and to compare with other approaches available at the time. The code and datasets are no longer accessible to the authors to their proprietary nature and so we cannot rerun experiments on the same data with new prompting methods. We hope that future work can directly compare our prompting method with those proposed in (Xing, 2024) or loss-based approaches on larger datasets.

Results in this report are also subject to the subjective definition in a 'segment' as encoded by manual segmentation or by wikipedia heading placement. The conclusions in this paper may be subject to this implicit definition of a segment, but we are optimistic that methods and results presented here are valid with more flexible definitions of segments and across different genres of text.

This work is also limited by the size of the dataset used. The dataset used in this work is a relatively small subset of the full Wikipedia dataset. This is due to the computational resources available to us that were required to train and test the models. We hope that future work can be conducted on larger datasets, both for fine-tuning and testing.

## Acknowledgments

project to take place. The two interns (Pierre and Maya) were supervised by Patrick. Pierre and Maya would like thank Patrick and other members of the data science and engineering teams at Adarga for their guidance and support throughout the project.

## References

Pinkesh Badjatiya, Litton J. Kurisinkel, Manish Gupta, and Vasudeva Varma. 2018. Attention-based neural text segmentation. In *Advances in Information Retrieval*, pages 180–193, Cham. Springer International Publishing.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.

Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. 2023. Extending context window of large language models via positional interpolation. *Preprint*, arXiv:2306.15595.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2024. Scaling instruction-finetuned language models. *J. Mach. Learn. Res.*, 25:70:1–70:53.

Massimiliano Costacurta. 2023. Text tiling done right: Building solid foundations for your personal llm.

Jacob Eisenstein. 2009. Hierarchical text segmentation from multi-scale lexical cohesion. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 353–361, Boulder, Colorado. Association for Computational Linguistics.

Xiachong Feng, Xiaocheng Feng, Libo Qin, Bing Qin, and Ting Liu. 2021. Language model as an annotator: Exploring DialoGPT for dialogue summarization. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1479–1491, Online. Association for Computational Linguistics.

Chris Fournier. 2013. Evaluating text segmentation using boundary edit distance. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1702–1712, Sofia, Bulgaria. Association for Computational Linguistics.

Michel Galley, Kathleen R. McKeown, Eric Fosler-Lussier, and Hongyan Jing. 2003. Discourse segmentation of multi-party conversation. In *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics*, pages 562–569, Sapporo, Japan. Association for Computational Linguistics.

Goran Glava s and Swapna Somasundaran. 2020. Two-level transformer and auxiliary coherence modeling for improved text segmentation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34:7797–7804.

Marti A. Hearst. 1997. Texttiling: segmenting text into multi-paragraph subtopic passages. *Comput. Linguist.*, 23(1):33–64.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Omri Koshorek, Adir Cohen, Noam Mor, Michael Rotman, and Jonathan Berant. 2018. Text segmentation as a supervised learning task. 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 469–473, New Orleans, Louisiana. Association for Computational Linguistics.

Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173.

Kelvin Lo, Yuan Jin, Weicong Tan, Ming Liu, Lan Du, and Wray Buntine. 2021. Transformer over pre-trained transformer for neural text segmentation with enhanced topic coherence. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3334–3340, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Michal Lukasik, Boris Dadachev, Kishore Papineni, and Gonçalo Simões. 2020. Text segmentation by cross segment attention. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4707–4716, Online. Association for Computational Linguistics.

F Petroni, PSH Lewis, A Piktus, Tim Rocktäschel, Yuxiang Wu, AH Miller, and Sebastian Riedel. 2020. How context affects language models' factual predictions.

Lev Pevzner and Marti A. Hearst. 2002. A critique and improvement of an evaluation metric for text segmentation. *Computational Linguistics*, 28(1):19–36.

Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: 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.

Peilu Wang, Yao Qian, Frank K. Soong, Lei He, and Hai Zhao. 2016. Learning distributed word representations for bidirectional LSTM recurrent neural network. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 527–533, San Diego, California. Association for Computational Linguistics.

Linzi Xing. 2024. *Versatile neural approaches to more accurate and robust topic segmentation*. Ph.D. thesis, University of British Columbia.

## A  Prompting Strategy

The prompting strategy used in this work is a simple schema that is designed to be general and applicable to any LLM. The schema is as follows:

1. The LLM is prompted with the input text, with integers in square brackets delimiting the sentence boundaries, few-shot examples of the task, a short instruction and a system prompt.

2. Segments are validated. This means they must not be too long nor too short and that they do not contain too many punctuation marks as a proportion of the segment length.

3. Segments that are too long are recursively split into smaller segments through similar prompting strategy. This prompt asks the LLM to return a single segment boundary index.

4. This process is repeated until all segments are short enough.

5. Segments that are too short are merged with a neighbouring segment based on the semantic similarity to neighbouring sentences. This part could also be done via prompting but we found this unnecessary.

An example prompt is shown below. Note that this is not the exact prompt used in the experiments but a simplified version intended for illustrative purposes.

**System:**

You are an expert linguist and a master of nuance in the meaning of written text. You obey instructions. You do not hallucinate. You are not a chatbot. You are not a summariser.

**Prompt:**

You are given a document with sentence boundaries marked by square brackets. Your task is to segment the document into coherent parts. Return a list of indices corresponding to the segment boundaries of the document. This list should ONLY be a list of integers, for example '1, 3, 5'. Some examples are shown below.

Text:

It was a sunny day in the park. [1] The birds were singing. [2] The children were playing. [3] The adults were chatting. [4] The dogs were barking. [5] The sun was shining. [6] The day was perfect. [7] However, then the rain came. [8] The children ran for cover. [9] The adults laughed. [10] The dogs howled. [11] The sun disappeared. [12] The day was ruined. [13] Fortunately, the next day was sunny again. [14] But it was actually too hot! [15] The children were sweating. [16] The adults were fanning themselves.

Segments:

7, 13

...*[more examples]*...

Text:

The cat sat on the mat. [1] The dog sat on the floor. [2] The cat was black. [3] The dog was brown. [4] The cat was fluffy. [5] The dog was short-haired. [6] The cat was purring. [7] The dog was wagging its tail. [8] The cat was happy. [9] The dog was happy. [10] Then the cat went to London. [11] The dog went to Paris. [12] The cat saw the sights. [13] The dog saw the sights. [14] The cat ate fish and chips. [15] The dog ate croissants. [16] The cat drank tea. [17] The dog drank coffee. [18] The cat was happy. [19]

Segments:

**End Prompt**

We use a similar prompt for the recursive prompting mechanism with the same system prompt. For example:

**Recursive Prompt:**

You are given a document with sentence boundaries marked by square brackets. Your task is to

choose one segment boundary to split the document into two coherent parts. Return a single integer corresponding to the index of the segment boundary. This integer should be between 1 and the number of sentences in the document. Some examples are shown below.

Text:

The cat sat on the mat. [1] The cat was black. [2] The cat was fluffy. [3] The cat was purring. [4] The cat was happy. [5] On the other hand, the dog sat on the floor. [6] The dog was brown. [7] The dog was short-haired. [8] The dog was wagging its tail. [9] The dog was happy. [10]

Segment:

5

...[more examples]...

Text:

Jack and Jill went up the hill. [1] Jack fell down and broke his crown. [2] Jill came tumbling after. [3] This is a well known nursery rhyme that has been passed down through the generations. [4] It is a classic. [5] It is a favourite of many. [6] It is a favourite of mine. [7] It is a favourite of yours. [8] It is a favourite of everyone.

Segment:

3

**End Prompt**

These examples are not illustrative of the length or style of segmentations in our dataset, they merely serve to exemplify the prompting schema. The actual prompts used in the experiments were much longer and more complex, and included more examples which were more realistic. The system prompt was also more detailed and included more examples of what the model should not do, such as not repeating the same segment boundary multiple times, not exceeding the length of the input sentences and not getting stuck in a pattern of regular segment boundaries.