Topic Segmentation using Large Language Models Master's Thesis Defense

Krishna Chaitanya Valaboju

Digital Engineering

Faculty of Media

11.06.2025

First Examiner : Prof. Dr. Benno Stein **Second Examiner** : Jun.-Prof. Dr. Jan Ehlers

Supervisor : Dr. Johannes Kiesel



Research Motivation

- Podcasts: 5 million+ episodes, informal & • conversational.
- No clear "chapters" = listeners waste time seeking key discussions.
- Unstructured + noisy transcripts make automatic \bullet segmentation hard.





YouTube Segmentation Parallel

- Video platforms (e.g., YouTube) already use "chapters" to break long videos into meaningful segments.
- Podcast transcripts: same idea, but only text. •
- Goal: Add chapter-like markers to audio via transcripts •



YouTube video segmentation with and without chapters, illustrating the difference in the progress bar.



Research Objectives



1. Build an automated pipeline to split podcast transcripts into topic segments. ĵ

3. Establish a robust evaluation: manual gold standard, F1, Pk, and WindowDiff.



2. Compare classical (TextTiling) vs. modern methods (LLM-similarity, Transformer BIO).



4. Identify the best approach for noisy, conversational data.



Topic Segmentation: Definition & Importance

What is Topic Segmentation?

- Process: Split a long transcript into coherent segments (distinct topics).
- Why? Transforms unstructured text into manageable,

semantically meaningful units.



- Longer documents such as magazine articles, news articles or books often talk about multiple topics
- Topic segmentation approaches attempt to break up these documents into subdocuments

Applications of Topic Segmentation

Information Retrieval

()

– Better search & retrieval

	P	
	-	٦
Ξ		I
-	_	
		ך ≡

Summarization

– Concise, segment-level

summaries.



Content Navigation

- Interactive navigation

(clickable chapters).

Recommendation

Systems

- Improved



recommendations (match

users to specific segments).

Challenges in Topic Segmentation

Lack of Clear Topic Boundaries Podcasts often exhibit informal and conversational speech, leading to 2 ambiguous or overlapping topic boundaries. \bigcirc E Long-form Content 8 Podcasts often span extensive durations, posing challenges for models with fixed context windows.

-

Data Noise

- Transcripts generated via ASR systems frequently contain errors due to factors like accents, background noise, or overlapping speech.

Conversational Nature

The dynamic and informal nature of -

podcast conversations adds

complexity to segmentation tasks.

Lack of Labeled Data

- The scarcity of annotated podcast
 - datasets hampers the development of
 - robust segmentation models.

Evolution of Topic Segmentation Techniques

Traditional Approaches

TextTiling (Hearst, 1997) and C99 (Choi, 2000) relied on lexical cohesion and fixed sliding windows to detect topic boundaries.

Neural Networks

RNNs and attention-based mechanisms improved capacity to capture long-range

dependencies in text.



Limitations of Traditional Rule-Based Methods



Semantic Relationships



Traditional rule-based methods typically depend on surface-level lexical overlap and frequency statistics. As a result, they overlook deeper semantic connections between sentences or paragraphs, such as paraphrasing, synonymy, and implicit references.

Context Awareness

Rule-based approaches generally operate on localized text windows using features like word distributions in fixed-length chunks or adjacent paragraphs and thus cannot model long-range dependencies.



Dynamic Vocabulary

podcast transcripts, speakers often employ colloquial language, domain-specific terminology, filler words, or abrupt changes in style. predetermined lexical cues cannot easily adapt to these

variations.



- In conversational domains like
- Rule-based methods relying on

Methodology Overview

Creating a Dataset for **Topic Segmentation**

Collection of podcast transcripts, preprocessing, and filtering to form a suitable dataset for analysis

Manual Annotation of **Topic Segments**

Development of annotation guidelines and execution of manual labeling to create a gold standard

Implementation of Automated Segmentation Methods **Development and** optimization of TextTiling, LLM-based, and Transformer-based

approaches

Evaluation and

performance using

Comparative Analysis

- Assessment of segmentation
- established metrics and
- detailed error analysis

Data Collection and Preprocessing



Process Details

- Collect 100 podcast audios \rightarrow transcribe via Whisper ulletASR.
- Filter: keep 30 episodes (10–30 min) ullet
- Minimal Cleaning: lowercase, normalize punctuation, and lemmatize. ullet



Podcast Duration Distribution



Duration Analysis

• (30 episodes, 10–30 min range)

Manual Annotation Workflow





Annotation Process

Overlap: transition sentences tagged in both segments

Annotation Example

Hierarchical Labeling Scheme (color-coded)

• <u>Main Topic</u>: primary segment

Eg: (Green segment in the figure)

• <u>Subtopics</u>: secondary segment

Eg: (purple and yellow segments in figure)

• <u>Ignore</u>: Refers to content that does not contribute to topic segmentation

Note: Sentences that serve as bridges between topics are included in both adjacent segments

Helio, and welcome to the IPP lotani-Palentine padcalit. (In Gravion, and today we'll be discussing the major coeffict best solution that establishes peace for both groups? So, how did this conflict arise? Well, the tarael-Palestine war started way back in 1948 when Jewish immigrants bombarded the state of stratil in an attempt to reciain the land they thought was dontfully theirs. Soon after the arrival of hundreds of lissall immigrants, they claimed largel as their own, provoking Palestinians. The Palestinians eventually started a war with the Israelis and lost and wouldn't stup. The Israelis were then known as the Sa-Day War. Ever s with barely any land, causing ion, and it is heir holy o live solution for both sides. meaning that th ion here to cite and model seacher ie closer together to orders be placed and only separates tible to being in danger from user. Another f Sub-top hanse may the will t tion? Another workdate that was thought tion for the israel and Palestin was the one-state solution. The two state solution is a solution that is important to all of us as Palestinians. The one-state solution is a resolution noticet, which ail rate state for both tonael and Palestine to share together. This means both process will be togethe equals, which also uestions. All right Kini, when the conflict? How do you think that has impacted the conflict tury The origins of th an conflict is an ongoing struggle be hanged some of their views on how they want to handle this conflict. I think the two state solution is the best solution s equally and lanore the best solution to create peace for both the israells and Palestinians. And thefa it for the israel Palestine podcast. See you next time. By

Segmentation Methods

TextTiling (Baseline)

장

 \bigcirc

 \bigcirc

- Lexical cohesion via fixed sliding windows.
- Detects sharp word-frequency dips as boundaries.

LLM-Based Topic Extraction with Similarity Thresholding method

- Use an LLM to propose topic labels & mpnet embeddings.
- Assign sentences when cosine similarity > θ .

Transformer-Based BIO Labeling method

- DistilBERT + CRF fine-tuned on BIO labels (B/I/O).
- Sentences tagged "Boundary" vs. "Inside" vs. "Ignore."

TextTiling (Baseline) Implementation

Implementation Workflow

- 1. <u>Pre-split</u>: Insert breaks every 5 sentences.
- 2. <u>Tokenize</u>: Split transcript into sentences.
- 3. <u>Slide windows of 5 sentences</u>: compute lexical cosine similarity.
- 4. Boundary = lowest similarity point.

House of Wax is a 1953 American warnercolor 3-D horror film about a disfigured sculptor who repopulates his destroyed wax museum by murdering people and using their wax-coated corpses as displays. Directed by Andre DeToth and starring Vincent Price, it is a remake of Warner Bros.' "Mystery of the Wax Museum" (1933), without the comic relief featured in the earlier film."House of Wax" was the first color 3-D feature from a major American studio and premiered just two days after the Columbia Pictures film "Man in the Dark", the first major-studio black-and-white 3-D feature.It was also the first 3-D film with stereophonic sound to be presented in a regular theater.

It premiered nationwide on April 10, 1953 and went out for a general release on April 25, 1953. In 1971, it was widely re-released to theaters in 3-D, with a full advertising campaign.Newly-struck prints of the film in Chris Condon's single-strip StereoVision 3-D format were used.

Another major re-release occurred during the 3-D boom of the early 1980s.In 2005, Warner Bros. distributed a new film also called "House of Wax", but its plot is very different from the one used in the two earlier films.The film starred Elisha Cuthbert, Chad Michael Murray, Paris Hilton and Jared Padalecki. This version received largely negative reviews from critics. In 2014, the film was deemed "culturally, historically, or aesthetically significant" by the Library of Congress and selected for preservation in the National Film Registry.



LLM-Based Topic Extraction with Similarity Thresholding

Implementation Workflow

- Ask LLM for ≤ 5 topic labels 1.
- Split the transcript into sentences. 2.
- Compute mpnet embeddings for each sentence & each topic. 3.
- Assign sentences to a topic when $\cos(\sin) > \theta$ ($\theta = 0.4$). 4.
- Iteratively reassign to refine boundaries. 5.



Bauhaus-Universität Weimar

Extract Topics

Generate Topic List

Segment Text

Compute Similarity

Transformer-Based Model with BIO Labeling

Implementation Workflow

- 1. <u>Manual BIO tags:</u> label each sentence B/I/O.
- 2. <u>Chunk transcripts</u>: 20 sent./chunk \rightarrow max 512 tokens.
- 3. Tokenize with DistilBERT + positional encodings.
- 4. Add CRF on top for sequence tag decoding.
- 5. Train via Leave-One-Out CV (30 folds).
- 6. Evaluate using F1, Pk, and WindowDiff.

Load CSV Data
Raw data
Parse Annotations Convert to BIO tags
BIO-tagge
Data Chunking (20 sentences per chunk)
Chunks
PyTorch Dataset DataLoader creation
Batches
DistilBERT Tokenization & Embeddings
Embedding
Positional Embeddings
Input
$\fbox{\begin{tabular}{c} Model Training \\ (DistilBERT + CRF, LOOC \end{tabular}) \end{tabular}}$
Predictions
Evaluate Segmentation $(F1, P_k, WindowDiff)$
Results
Report Metrics (Avg. F1, Avg. P_k , Avg. W



Evaluation Metrics

F1 Score

F1 Score = $2 \times$ (Precision \times Recall)/(Precision + Recall) – Measures correct vs. missed boundaries.

Pk Metric

- Probability two sentences k apart are missegmented.
- –A lower Pk score indicates better performance, demonstrating fewer segmentation errors.

WindowDiff Metric

- Counts boundary count differences in a sliding window. – More sensitive to near-missed boundaries than Pk. - Similar to Pk, a lower WindowDiff score signifies better
- segmentation.



Quantitative Results

Method	F1 Score	Pk	WindowDi
TextTiling	0.53	0.44	0.45
LLM Similarity Thresholding	0.72	0.29	0.31
Transformer BIO Labeling	0.47	0.57	0.68

• **Best = LLM Similarity Thresholding method** (highest F1, lowest Pk & WD)





Comparative Performance Visualization



Performance Analysis

- The LLM-based Similarity Thresholding method substantially outperforms TextTiling on all reported metrics (F1 score of 0.72, Pk of 0.29, and WindowDiff of 0.31).
- Its reliance on contextual embeddings rather than raw lexical overlap allows it to better handle semantically subtle boundaries, such as transitions involving synonyms or related concepts.
- The Transformer-based BIO Labeling approach, which lacksquareachieved an F1 score of 0.47, a Pk of 0.57, and a WindowDiff of 0.68, performs variably in comparison to TextTiling.

Threshold Sensitivity in LLM-Based Method

Impact of Threshold Variation

- F1 peaks at $\theta \approx 0.4$ lacksquare
- Pk & WD minimized at $\theta \approx 0.4$ •





Threshold Sensitivity Analysis



Key Observations

- Peak in F1 Score Near 0.4 •
- Minimum in Pk and WindowDiff Around 0.4 •
- Sharp Changes After the Inflection Point ullet
- A threshold in the range of 0.3–0.4 appears to yield the most balanced performance.



Transformer BIO: Fold-Wise Performance (LOOCV)

Performance Trends Across Folds

- F1 varies from 0.30 to 0.85 across 30 folds •
- Pk errors peak in some folds (1.0) lacksquare
- WD errors also fluctuate widely. \bullet

Model highly sensitive to transcript structure & annotation.



Factors Influencing Segmentation Performance

$\overset{\circ}{\sim}$

Transcript Complexity

Transcripts with frequent digressions, informal discussions, or highly unstructured conversations led to lower scores. Such cases made it difficult for the model to learn stable topic boundaries.



Scripted vs. Spontaneous Speech

Folds with higher F1 scores corresponded to well-structured transcripts, such as scripted monologues or interview-style podcasts with clear topic transitions.



Annotation Inconsistencies

Some transcripts contained inconsistencies in manual annotations, where subtle transitions were labeled differently across transcripts. This impacted model training and resulted in lower F1 scores.



Topic Overlaps and Ambiguity

Transcripts in which multiple topics were discussed simultaneously, or where a single segment spanned multiple intertwined topics, tended to increase segmentation errors.



Key Findings

<u>[~]</u>	LLM-based The LLM-based flexibility	method superi	ority olding method der	nonstrated superior performance in balancing topic c
٢	C	Threshold sensitivity Optimal threshold selection (0.3-0.4) is critical for balancing over-segmentation and under-segmentation		
E Transcript s Well-structure results		structure impact red transcripts with clear transitions yielded better seg		
	Đ	17 ⊕		Contextual understanding Deep semantic representations outperformed surf features



coherence and

gmentation

rface-level lexical

Future Directions

£03

ե

R

 \bigoplus

Q

Automated Threshold Optimization

Develop adaptive thresholding mechanisms using reinforcement learning

Hybrid Models

Integrate semantic similarity and sequence labeling techniques

Weakly Supervised Learning

Explore semi-supervised methods to reduce reliance on labeled data

Domain Adaptation

test on various podcast genres

Real-World Deployment

Deploy in live podcast apps (chaptering & search)

Questions?



Thank you for your attention!