

FINAL CODE FOR SEGMENTATION:

```

1 pip install yt_dlp

Collecting yt_dlp
  Downloading yt_dlp-2025.2.19-py3-none-any.whl.metadata (171 kB)
    171.9/171.9 kB 11.6 MB/s eta 0:00:00
  Downloading yt_dlp-2025.2.19-py3-none-any.whl (3.2 MB)
    3.2/3.2 MB 56.9 MB/s eta 0:00:00
Installing collected packages: yt_dlp
Successfully installed yt_dlp-2025.2.19

1 import yt_dlp
2
3 def download_audio(video_url):
4     ydl_opts = {
5         'format': 'bestaudio/best',
6         'outtmpl': '/content/audio.mp3',
7         'cookiefile': '/content/cookies (3).txt', # Update this path
8         'postprocessors': [{
9             'key': 'FFmpegExtractAudio',
10            'preferredcodec': 'mp3',
11            'preferredquality': '192',
12        }],
13    }
14
15    with yt_dlp.YoutubeDL(ydl_opts) as ydl:
16        ydl.download([video_url])
17
18    return "/content/audio.mp3"
19
20 video_url = input("Enter YouTube video URL: ")
21 file_path = download_audio(video_url)
22 print(f"Downloaded audio: {file_path}")

Enter YouTube video URL: https://youtu.be/sK8SILOM37I
[youtube] Extracting URL: https://youtu.be/sK8SILOM37I
[youtube] sK8SILOM37I: Downloading webpage
[youtube] sK8SILOM37I: Downloading tv client config
[youtube] sK8SILOM37I: Downloading player 7d1d50a6
[youtube] sK8SILOM37I: Downloading tv player API JSON
[info] sK8SILOM37I: Downloading 1 format(s): 251
[download] Destination: /content/audio.mp3
[download] 100% of 39.37MiB in 00:00:03 at 13.06MiB/s
[ExtractAudio] Destination: /content/audio.mp3.mp3
Deleting original file /content/audio.mp3 (pass -k to keep)
Downloaded audio: /content/audio.mp3

1 import os
2
3 # Check the files in /content directory
4 print("Files in content:", os.listdir("/content"))
5
6 # Fix the file name if needed
7 old_path = "/content/audio.mp3.mp3"
8 new_path = "/content/audio.mp3"
9
10 if os.path.exists(old_path):
11     os.rename(old_path, new_path)
12     print(f"Renamed: {old_path} → {new_path}")
13 else:
14     print("Audio file not found!")

Files in content: ['.config', 'cookies (3).txt', 'audio.mp3.mp3', 'sample_data']
Renamed: /content/audio.mp3.mp3 → /content/audio.mp3

1 import IPython.display as ipd
2
3 # Set the correct file path
4 audio_path = "/content/audio.mp3"
5
6 # Play the audio with the correct sample rate (e.g., 44100 Hz)
7 ipd.Audio(audio_path, rate=44100)

```

```

1 import os
2 import whisper
3 from pydub import AudioSegment
4
5 def convert_audio_to_wav(audio_file):
6     """Converts an MP3 audio file to WAV format for better transcription."""
7     wav_file = "audio.wav"
8     if audio_file.endswith(".mp3"):
9         try:
10             AudioSegment.from_mp3(audio_file).export(wav_file, format="wav")
11             return wav_file
12         except Exception as e:
13             return f"Error converting to WAV: {str(e)}"
14     return audio_file # If already WAV, return as is
15
16 def transcribe_audio_with_whisper(audio_path):
17     """Transcribes audio using OpenAI's Whisper model."""
18     model = whisper.load_model("base") # Options: "tiny", "base", "small", "medium", "large"
19
20     print(" ♦ Transcribing with Whisper... (this may take some time)")
21
22     result = model.transcribe(audio_path)
23     return result["text"]
24
25 def save_transcript_to_file(transcript, filename="transcript.txt"):
26     """Saves the transcript to a text file."""
27     with open(filename, "w", encoding="utf-8") as file:
28         file.write(transcript)
29     print(f"✅ Transcript saved to {filename}")
30
31 # Example usage
32 if __name__ == "__main__":
33     audio_file = input("Enter the path to your audio file (MP3/WAV): ")
34
35     if not os.path.exists(audio_file):
36         print("❌ Error: File not found.")
37     else:
38         wav_file = convert_audio_to_wav(audio_file)
39
40         if "Error" in wav_file:
41             print(wav_file)
42         else:
43             transcript = transcribe_audio_with_whisper(wav_file)
44             print("\n ♦ Transcript:\n", transcript)
45
46             # Save transcript
47             save_transcript_to_file(transcript, "transcript.txt")
48
49             # Cleanup WAV file if it was converted
50             if wav_file != audio_file:
51                 os.remove(wav_file)

```

```

18         transcript = file.read()
19         return transcript
20     except Exception as e:
21         print(f"Error reading file: {e}")
22         return None
23
24 def split_into_sentences(transcript):
25     return sent_tokenize(transcript)
26
27 def calculate_cosine_similarity(sentences):
28     vectorizer = TfidfVectorizer()
29     tfidf_matrix = vectorizer.fit_transform(sentences)
30     cosine_sim = cosine_similarity(tfidf_matrix)
31     return cosine_sim
32
33 def segment_sentences(sentences, cosine_sim, threshold=0.5, min_sentences=5):
34     visited = [False] * len(sentences)
35     segments = []
36
37     for i in range(len(sentences)):
38         if not visited[i]:
39             segment = [sentences[i]]
40             visited[i] = True
41
42             for j in range(i + 1, len(sentences)):
43                 if not visited[j] and cosine_sim[i][j] >= threshold:
44                     segment.append(sentences[j])
45                     visited[j] = True
46
47             segments.append(segment)
48
49     # Merge smaller segments
50     merged_segments = []
51     temp_segment = []
52
53     for segment in segments:
54         temp_segment.extend(segment)
55         if len(temp_segment) >= min_sentences:
56             merged_segments.append(temp_segment)
57             temp_segment = []
58
59     if temp_segment:
60         if merged_segments:
61             merged_segments[-1].extend(temp_segment)
62         else:
63             merged_segments.append(temp_segment)
64
65     return merged_segments
66
67 def print_segments(segments):
68     for idx, segment in enumerate(segments, start=1):
69         print(f"\nSegment {idx}:")
70         for sentence in segment:
71             print(f" - {sentence}")
72
73 def save_segments_to_file(segments, output_file):
74     try:
75         with open(output_file, 'w', encoding='utf-8') as file:
76             for idx, segment in enumerate(segments, start=1):
77                 file.write(f"Segment {idx}:\n")
78                 for sentence in segment:
79                     file.write(f" - {sentence}\n")
80                 file.write("\n")
81             print(f"Segmented output saved to {output_file}")
82     except Exception as e:
83         print(f"Error saving file: {e}")
84
85 def process_transcript(file_path, threshold=0.5, min_sentences=5, output_file='segmented_output.txt'):
86     transcript = read_transcript(file_path)
87     if transcript is None:
88         return
89
90     sentences = split_into_sentences(transcript)
91     print("Sentences extracted:")
92     for i, sentence in enumerate(sentences):
93         print(f"S{i + 1}: {sentence}")
94
95     cosine_sim = calculate_cosine_similarity(sentences)
96
97     segments = segment_sentences(sentences, cosine_sim, threshold, min_sentences)
98
99     print("\nSegmented Sentences:")

```

```

100 print_segments(segments)
101
102 save_segments_to_file(segments, output_file)
103
104 file_path = '/content/transcript.txt' # Replace with your file path
105 output_file = '/content/segmented_output.txt'
106 process_transcript(file_path, threshold=0.15, min_sentences=5, output_file=output_file)

```

[nltk_data] Downloading package punkt to /root/nltk_data...

[nltk_data] Package punkt is already up-to-date!

[nltk_data] Downloading package punkt_tab to /root/nltk_data...

[nltk_data] Unzipping tokenizers/punkt_tab.zip.

[nltk_data] Error loading stopwords: Package 'stopword' not found in index

[nltk_data] Downloading package wordnet to /root/nltk_data...

[nltk_data] Downloading package omw-1.4 to /root/nltk_data...

[nltk_data] Downloading package averaged_perceptron_tagger to /root/nltk_data...

[nltk_data] Unzipping taggers/averaged_perceptron_tagger.zip.

Sentences extracted:

S1: So, sir, we know that India has seen a huge revolution with digital payments.

S2: We all thought that India is a place, at least the West thought that India is a place where many people do not get a square

S3: That was an narrative some 30 years ago.

S4: And not many are literate, people cannot read.

S5: But then we have now shown that digital payments, number one is India, while people thought that it wouldn't even come to top

S6: I think immediately after UPI the next big revolution, personally I think is in education.

S7: And the complete homework for this has happened in the form of NEP, the documentation of which many of us have read and reali

S8: So, followed by which we got NCRF framework done, which is a national creative framework, surrounding which we will be discus

S9: So, my question is, do you think NCRF plus NEP put together will be the next big revolution after UPI in India?

S10: Absolutely.

S11: And why I think so is because in education the last policy came up many many years ago that was in 1986, which was slightly

S12: And thereafter so many changes have happened in the real world, so many changes have happened in the requirement of the indu

S13: However there were no corresponding changes which happened in education system.

S14: So, therefore I feel that this was the right time when we brought in the education policy 2020, honorable prime minister dec

S15: And we recently celebrated the fourth anniversary of NEP 2020.

S16: NCRF has been brought to implement the intent of 2020, 2020 is a policy, NEP 2020 is a policy.

S17: And for implementing a policy you need a framework.

S18: Now, why we call it a framework?

S19: We call it a framework because this is very flexible.

S20: This allows you all the innovation, the way you educate your students.

S21: Still it provides you the basic guidelines, the framework, the outer layer it provides.

S22: And that layer is mostly the enabling layer.

S23: That is such a enabling layer that it has broken the shackles which were there in the education sector.

S24: So, yes it is a big revolution and this is going to change the way we have been educating our kids and this will be game ch

S25: So, sir I think let us go with this example of let us say I started off living in a small 2BHK apartment and I slowly devel

S26: And one fine day you came and you changed my kitchen completely.

S27: I was using a bicycle, I moved to a scooter and a car and now we are asking me to fly and giving me an aircraft.

S28: NEP sounds more like that for me.

S29: How do I do the transition?

S30: I fear that I will crash if I use a aeroplane without training.

S31: I am talking about all the teachers in the country, all the schools in the country.

S32: We have been driving buses at max.

S33: Now we should fly how do we do this?

S34: Okay.

S35: Look at the requirement of the industry.

S36: The requirement of the industry has been moving very fast.

S37: The technology is emerging every day and the industry is moving with that speed.

S38: So, when a student is coming out of your institute and is going out in the market, he finds that whatever he has been taught

S39: When industry is moving that fast, when the requirement is moving that fast, don't you think it is important for us to chang

S40: Yes, how long can we wait?

S41: Yes, it will take a lot of effort for every one of us to adapt to this change.

S42: But this change is going to be not only beneficial but also very facilitated for all of us, very liberating for all of us.

S43: This is going to be highly liberative and choice based system.

S44: There are number of choices which are available to you, which are available to every student.

S45: Yes, when we introduce a new system, we have to really create new things, create new ways of doing things, learn something r

Keyword Generation

```
1 !pip install gensim spacy nltk
```

Collecting gensim

Downloading gensim-4.3.3-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (8.1 kB)

Requirement already satisfied: spacy in /usr/local/lib/python3.11/dist-packages (3.7.5)

Requirement already satisfied: nltk in /usr/local/lib/python3.11/dist-packages (3.9.1)

Requirement already satisfied: numpy<2.0,>=1.18.5 in /usr/local/lib/python3.11/dist-packages (from gensim) (1.26.4)

Collecting scipy<1.14.0,>=1.7.0 (from gensim)

Downloading scipy-1.13.1-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (60 kB)

60.6/60.6 kB 3.6 MB/s eta 0:00:00

Requirement already satisfied: smart-open>=1.8.1 in /usr/local/lib/python3.11/dist-packages (from gensim) (7.1.0)

Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.11 in /usr/local/lib/python3.11/dist-packages (from spacy) (3.0.12)

Requirement already satisfied: spacy-loggers<2.0.0,>=1.0.0 in /usr/local/lib/python3.11/dist-packages (from spacy) (1.0.5)

Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /usr/local/lib/python3.11/dist-packages (from spacy) (1.0.12)

Requirement already satisfied: cymem<2.1.0,>=2.0.2 in /usr/local/lib/python3.11/dist-packages (from spacy) (2.0.11)

Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /usr/local/lib/python3.11/dist-packages (from spacy) (3.0.9)

Requirement already satisfied: thinc<8.3.0,>=8.2.2 in /usr/local/lib/python3.11/dist-packages (from spacy) (8.2.5)

Requirement already satisfied: wasabi<1.2.0,>=0.9.1 in /usr/local/lib/python3.11/dist-packages (from spacy) (1.1.3)

Requirement already satisfied: srsly<3.0.0,>=2.4.3 in /usr/local/lib/python3.11/dist-packages (from spacy) (2.5.1)

```

Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in /usr/local/lib/python3.11/dist-packages (from spacy) (2.0.10)
Requirement already satisfied: weasel<0.5.0,>=0.1.0 in /usr/local/lib/python3.11/dist-packages (from spacy) (0.4.1)
Requirement already satisfied: typer<1.0.0,>=0.3.0 in /usr/local/lib/python3.11/dist-packages (from spacy) (0.15.2)
Requirement already satisfied: tqdm<5.0.0,>=4.38.0 in /usr/local/lib/python3.11/dist-packages (from spacy) (4.67.1)
Requirement already satisfied: requests<3.0.0,>=2.13.0 in /usr/local/lib/python3.11/dist-packages (from spacy) (2.32.3)
Requirement already satisfied: pydantic!=1.8,!1.8.1,<3.0.0,>=1.7.4 in /usr/local/lib/python3.11/dist-packages (from spacy) (2.11.6)
Requirement already satisfied: Jinja2 in /usr/local/lib/python3.11/dist-packages (from spacy) (3.1.6)
Requirement already satisfied: setuptools in /usr/local/lib/python3.11/dist-packages (from spacy) (75.1.0)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from spacy) (24.2)
Requirement already satisfied: langcodes<4.0.0,>=3.2.0 in /usr/local/lib/python3.11/dist-packages (from spacy) (3.5.0)
Requirement already satisfied: click in /usr/local/lib/python3.11/dist-packages (from nltk) (8.1.8)
Requirement already satisfied: joblib in /usr/local/lib/python3.11/dist-packages (from nltk) (1.4.2)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.11/dist-packages (from nltk) (2024.11.6)
Requirement already satisfied: language-data>=1.2 in /usr/local/lib/python3.11/dist-packages (from langcodes<4.0.0,>=3.2.0->spacy) (1.3.0)
Requirement already satisfied: annotated-types<=0.6.0 in /usr/local/lib/python3.11/dist-packages (from pydantic!=1.8,!1.8.1,<3.0.0->spacy) (0.7.0)
Requirement already satisfied: pydantic-core==2.27.2 in /usr/local/lib/python3.11/dist-packages (from pydantic!=1.8,!1.8.1,<3.0.0->spacy) (2.27.2)
Requirement already satisfied: typing-extensions>=4.12.2 in /usr/local/lib/python3.11/dist-packages (from pydantic!=1.8,!1.8.1,<3.0.0->spacy) (4.12.2)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy) (3.4.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy) (3.10.1)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy) (2.2.3)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy) (2025.1.31)
Requirement already satisfied: wrapt in /usr/local/lib/python3.11/dist-packages (from smart-open>=1.8.1->gensim) (1.17.2)
Requirement already satisfied: blis<0.8.0,>=0.7.8 in /usr/local/lib/python3.11/dist-packages (from thinc<8.3.0,>=8.2.2->spacy) (0.7.10)
Requirement already satisfied: confection<1.0.0,>=0.0.1 in /usr/local/lib/python3.11/dist-packages (from thinc<8.3.0,>=8.2.2->spacy) (0.0.4)
Requirement already satisfied: shellingham>=1.3.0 in /usr/local/lib/python3.11/dist-packages (from typer<1.0.0,>=0.3.0->spacy) (1.5.4)
Requirement already satisfied: rich>=10.11.0 in /usr/local/lib/python3.11/dist-packages (from typer<1.0.0,>=0.3.0->spacy) (13.9.4)
Requirement already satisfied: cloudpathlib<1.0.0,>=0.7.0 in /usr/local/lib/python3.11/dist-packages (from weasel<0.5.0,>=0.1.0->spacy) (0.19.0)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from Jinja2->spacy) (3.0.2)
Requirement already satisfied: marisa-trie>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from language-data>=1.2->langcodes<4.0.0,>=3.2.0->spacy) (1.3.0)
Requirement already satisfied: markdown-it-py>=2.2.0 in /usr/local/lib/python3.11/dist-packages (from rich>=10.11.0->typer<1.0.0,>=0.3.0->spacy) (3.0.0)
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in /usr/local/lib/python3.11/dist-packages (from rich>=10.11.0->typer<1.0.0,>=0.3.0->spacy) (2.19.1)
Requirement already satisfied: mdurl~0.1 in /usr/local/lib/python3.11/dist-packages (from markdown-it-py>=2.2.0->rich>=10.11.0->typer<1.0.0,>=0.3.0->spacy) (0.1.2)
Downloading gensim-4.3.3-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (26.7 MB)
26.7/26.7 MB 38.2 MB/s eta 0:00:00
Downloading scipy-1.13.1-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (38.6 MB)
38.6/38.6 MB 20.5 MB/s eta 0:00:00
Installing collected packages: scipy, gensim
Attempting uninstall: scipy
Found existing installation: scipy 1.14.1
Uninstalling scipy-1.14.1:

```

Using LDA

```

1 import re
2 import gensim
3 import spacy
4 import nltk
5 from gensim import corpora
6 from gensim.models import LdaModel
7 from nltk.corpus import stopwords
8
9 # Download stopwords if not available
10 nltk.download("stopwords")
11 stop_words = set(stopwords.words("english"))
12
13 # Load spaCy NLP model
14 nlp = spacy.load("en_core_web_sm")
15
16 # Function to preprocess text
17 def preprocess_text(text):
18     text = re.sub(r'\s+', ' ', text) # Remove extra spaces
19     doc = nlp(text.lower()) # Convert to lowercase and tokenize
20     tokens = [token.lemma_ for token in doc if token.is_alpha and token.lemma_ not in stop_words]
21     return tokens
22
23 # Read the transcript from file
24 input_file = "/content/segmented_output.txt" # Change this to your actual file name
25 with open(input_file, "r", encoding="utf-8") as file:
26     transcript = file.read()
27
28 # Preprocess each segment in the transcript
29 processed_texts = [preprocess_text(segment) for segment in transcript.split("\n") if segment.strip()]
30
31 # Create Dictionary and Corpus for LDA
32 dictionary = corpora.Dictionary(processed_texts)
33 corpus = [dictionary.doc2bow(text) for text in processed_texts]
34
35 # Train LDA Model
36 num_topics = 3 # Adjust based on requirement
37 lda_model = LdaModel(corpus, num_topics=num_topics, id2word=dictionary, passes=10)
38
39 # Display Extracted Topics
40 print("\nExtracted Topics:")

```

```
41 for idx, topic in lda_model.print_topics():
42     print(f"Topic {idx + 1}: {topic}")
```

```

➡ [nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!

```

Extracted Topics:

Topic 1: 0.066*"I" + 0.039*"segment" + 0.021*"say" + 0.017*"student" + 0.015*"think" + 0.012*"number" + 0.011*"see" + 0.011*"go" + 0.008*"study"

Topic 2: 0.030*"I" + 0.018*"ncrf" + 0.012*"framework" + 0.012*"learn" + 0.012*"absolutely" + 0.011*"way" + 0.009*"outcome" + 0.008*"work"

Topic 3: 0.021*"student" + 0.019*"skill" + 0.016*"go" + 0.014*"year" + 0.014*"time" + 0.012*"teach" + 0.011*"course" + 0.010*"teacher"

```

1 import re
2 import gensim
3 import spacy
4 import nltk
5 from gensim import corpora
6 from gensim.models import LdaModel
7 from nltk.corpus import stopwords
8
9 # Download stopwords if not available
10 nltk.download("stopwords")
11 stop_words = set(stopwords.words("english"))
12
13 # Load spaCy NLP model
14 nlp = spacy.load("en_core_web_sm")
15
16 # Function to preprocess text
17 def preprocess_text(text):
18     text = re.sub(r'\s+', ' ', text) # Remove extra spaces
19     doc = nlp(text.lower()) # Convert to lowercase and tokenize
20     tokens = [token.lemma_ for token in doc if token.is_alpha and token.lemma_ not in stop_words]
21     return tokens
22
23 # Read the transcript from file
24 input_file = "/content/segmented_output.txt" # Change this to your actual file name
25 with open(input_file, "r", encoding="utf-8") as file:
26     transcript = file.readlines()
27
28 # Preprocess each segment
29 processed_texts = [preprocess_text(segment) for segment in transcript if segment.strip()]
30
31 # Create Dictionary and Corpus for LDA
32 dictionary = corpora.Dictionary(processed_texts)
33 corpus = [dictionary.doc2bow(text) for text in processed_texts]
34
35 # Train LDA Model
36 num_topics = 5 # Adjust based on requirement
37 lda_model = LdaModel(corpus, num_topics=num_topics, id2word=dictionary, passes=10)
38
39 # Extract the dominant topic for each segment and save to file
40 output_file = "output_topics.txt"
41 with open(output_file, "w", encoding="utf-8") as out_file:
42     for i, bow in enumerate(corpus):
43         topic_distribution = lda_model.get_document_topics(bow)
44         dominant_topic = max(topic_distribution, key=lambda x: x[1])[0] # Get the most probable topic
45         topic_keywords = lda_model.print_topic(dominant_topic)
46
47         out_file.write(f"Segment {i + 1}:\n")
48         out_file.write(f"Topic: {topic_keywords}\n\n")
49
50 print(f"Extracted topics saved to {output_file}")
51

```


```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
Extracted topics saved to output/topics.txt
```

```
1 import re
2 import gensim
3 import spacy
4 import nltk
5 from gensim import corpora
6 from gensim.models import LdaModel
7 from nltk.corpus import stopwords
8
9 # Download stopwords if not available
10 nltk.download("stopwords")
11 stop_words = set(stopwords.words("english"))
12
13 # Load spaCy NLP model
14 nlp = spacy.load("en_core_web_sm")
```

```

15
16 # Function to preprocess text
17 def preprocess_text(text):
18     """Cleans and tokenizes text for topic modeling"""
19     text = re.sub(r'\s+', ' ', text) # Remove extra spaces
20     doc = nlp(text.lower()) # Convert to lowercase and tokenize
21     tokens = [token.lemma_ for token in doc if token.is_alpha and token.lemma_ not in stop_words]
22     return tokens
23
24 # Read the transcript file
25 input_file = "/content/segmented_output (2).txt" # Change this to your actual file name
26 with open(input_file, "r", encoding="utf-8") as file:
27     transcript = file.read()
28
29 # Split transcript into segments using a pattern like 'Segment X:' or other logical markers
30 segments = re.split(r'Segment \d+:', transcript)[1:] # Removes empty first split
31
32 # Preprocess each full segment
33 processed_texts = [preprocess_text(segment.strip()) for segment in segments if segment.strip()]
34
35 # Create Dictionary and Corpus for LDA
36 dictionary = corpora.Dictionary(processed_texts)
37 corpus = [dictionary.doc2bow(text) for text in processed_texts]
38
39 # Train LDA Model
40 num_topics = min(len(processed_texts), 5) # Set topic count based on available data
41 lda_model = LdaModel(corpus, num_topics=num_topics, id2word=dictionary, passes=10)
42
43 # Extract the top keyword for each full segment
44 output_file = "top_keywords.txt"
45 with open(output_file, "w", encoding="utf-8") as out_file:
46     for i, bow in enumerate(corpus):
47         topic_distribution = lda_model.get_document_topics(bow)
48         dominant_topic = max(topic_distribution, key=lambda x: x[1])[0] # Get the most probable topic
49         topic_keywords = lda_model.print_topic(dominant_topic)
50
51         # Extract the highest probability keyword
52         keywords = topic_keywords.split(" + ")
53         top_keyword = max(keywords, key=lambda x: float(x.split("*")[0])) # Get highest weight word
54         top_word = top_keyword.split("*")[1].replace("'", "").strip() # Extract the actual word
55
56         out_file.write(f"Segment {i + 1}: {top_word}\n")
57
58 print(f"Top keywords saved to {output_file}")

```

 [nltk_data] Downloading package stopwords to /root/nltk_data...
 [nltk_data] Package stopwords is already up-to-date!
 Top keywords saved to top_keywords.txt

```

1 import re
2 import gensim
3 import spacy
4 import nltk
5 from gensim import corpora
6 from gensim.models import LdaModel
7 from nltk.corpus import stopwords
8
9 # Download stopwords if not available
10 nltk.download("stopwords")
11 stop_words = set(stopwords.words("english"))
12
13 # Load spaCy NLP model
14 nlp = spacy.load("en_core_web_sm")
15
16 # Function to preprocess text
17 def preprocess_text(text):
18     """Cleans and tokenizes text for topic modeling"""
19     text = re.sub(r'\s+', ' ', text) # Remove extra spaces
20     doc = nlp(text.lower()) # Convert to lowercase and tokenize
21     tokens = [token.lemma_ for token in doc if token.is_alpha and token.lemma_ not in stop_words]
22     return tokens
23
24 # Read the transcript file
25 input_file = "/content/segmented_output (2).txt" # Change this to your actual file name
26 with open(input_file, "r", encoding="utf-8") as file:
27     transcript = file.read()
28
29 # Split transcript into segments using 'Segment X:' markers
30 segments = re.split(r'Segment \d+:', transcript)[1:] # Removes empty first split
31
32 # Preprocess each full segment
33 processed_texts = [preprocess_text(segment.strip()) for segment in segments if segment.strip()]

```



```

34
35 # Create Dictionary and Corpus for LDA
36 dictionary = corpora.Dictionary(processed_texts)
37 corpus = [dictionary.doc2bow(text) for text in processed_texts]
38
39 # Train LDA Model
40 num_topics = min(len(processed_texts), 5) # Set topic count based on available data
41 lda_model = LdaModel(corpus, num_topics=num_topics, id2word=dictionary, passes=10)
42
43 # Extract the top 3 keywords for each full segment
44 output_file = "top_keywords2.txt"
45 with open(output_file, "w", encoding="utf-8") as out_file:
46     for i, bow in enumerate(corpus):
47         topic_distribution = lda_model.get_document_topics(bow)
48         dominant_topic = max(topic_distribution, key=lambda x: x[1])[0] # Get the most probable topic
49         topic_keywords = lda_model.print_topic(dominant_topic)
50
51         # Extract top 3 highest probability keywords
52         keywords = topic_keywords.split(" + ")
53         top_keywords = sorted(keywords, key=lambda x: float(x.split("*")[0]), reverse=True)[:3]
54         top_words = [word.split("*")[1].replace("'", "").strip() for word in top_keywords]
55
56         out_file.write(f"Segment {i + 1}: {' '.join(top_words)}\n")
57
58 print(f"Top keywords saved to {output_file}")

```

```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
Top keywords saved to top_keywords2.txt

```

Using KeyBERT

```
1 pip install keybert
```

```

Collecting keybert
  Downloading keybert-0.9.0-py3-none-any.whl.metadata (15 kB)
Requirement already satisfied: numpy>=1.18.5 in /usr/local/lib/python3.11/dist-packages (from keybert) (1.26.4)
Requirement already satisfied: rich>=10.4.0 in /usr/local/lib/python3.11/dist-packages (from keybert) (13.9.4)
Requirement already satisfied: scikit-learn>=0.22.2 in /usr/local/lib/python3.11/dist-packages (from keybert) (1.6.1)
Requirement already satisfied: sentence-transformers>=0.3.8 in /usr/local/lib/python3.11/dist-packages (from keybert) (3.4.1)
Requirement already satisfied: markdown-it-py>=2.2.0 in /usr/local/lib/python3.11/dist-packages (from rich>=10.4.0->keybert) (3.0.0)
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in /usr/local/lib/python3.11/dist-packages (from rich>=10.4.0->keybert) (2.11)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn>=0.22.2->keybert) (1.13.1)
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn>=0.22.2->keybert) (1.4.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn>=0.22.2->keybert) (3.5.0)
Requirement already satisfied: transformers<5.0.0,>=4.41.0 in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.3.8->keybert) (4.41.0)
Requirement already satisfied: tqdm in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.3.8->keybert) (4.67.1)
Requirement already satisfied: torch>=1.11.0 in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.3.8->keybert) (2.5.1)
Requirement already satisfied: huggingface-hub>=0.20.0 in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.3.8->keybert) (0.20.0)
Requirement already satisfied: Pillow in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.3.8->keybert) (11.1.0)
Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers) (3.16.1)
Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers) (2025.3.0)
Requirement already satisfied: packaging>=20.9 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers) (25.0)
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers) (6.0.2)
Requirement already satisfied: requests in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers) (2.32.3)
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers) (4.12.2)
Requirement already satisfied: mdurl>=0.1 in /usr/local/lib/python3.11/dist-packages (from markdown-it-py>=2.2.0->rich>=10.4.0->keybert) (0.1.2)
Requirement already satisfied: networkx in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (3.4.2)
Requirement already satisfied: Jinja2 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (3.1.4)
Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (12.4.127)
Requirement already satisfied: nvidia-cuda-runtime-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (12.4.127)
Requirement already satisfied: nvidia-cuda-cupti-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (12.4.127)
Requirement already satisfied: nvidia-cudnn-cu12==9.1.0.70 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (9.1.0.70)
Requirement already satisfied: nvidia-cublas-cu12==12.4.5.8 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (12.4.5.8)
Requirement already satisfied: nvidia-cufft-cu12==11.2.1.3 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (11.2.1.3)
Requirement already satisfied: nvidia-curand-cu12==10.3.5.147 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (10.3.5.147)
Requirement already satisfied: nvidia-cusolver-cu12==11.6.1.9 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (11.6.1.9)
Requirement already satisfied: nvidia-cusparse-cu12==12.3.1.170 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (12.3.1.170)
Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (2.21.5)
Requirement already satisfied: nvidia-nvtx-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (12.4.127)
Requirement already satisfied: nvidia-nvjitlink-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (12.4.127)
Requirement already satisfied: triton==3.1.0 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (3.1.0)
Requirement already satisfied: sympy==1.13.1 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (1.13.1)
Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from sympy==1.13.1->torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (1.3.0)
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.11/dist-packages (from transformers<5.0.0,>=4.41.0->sentence-transformers>=0.3.8->keybert) (2024.11.6)
Requirement already satisfied: tokenizers<0.22,>=0.21 in /usr/local/lib/python3.11/dist-packages (from transformers<5.0.0,>=4.41.0->sentence-transformers>=0.3.8->keybert) (0.21.0)
Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.11/dist-packages (from transformers<5.0.0,>=4.41.0->sentence-transformers>=0.3.8->keybert) (0.5.0)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from Jinja2->torch>=1.11.0->sentence-transformers>=0.3.8->keybert) (3.0.2)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests->huggingface-hub>=0.20.0->sentence-transformers>=0.3.8->keybert) (3.4.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests->huggingface-hub>=0.20.0->sentence-transformers>=0.3.8->keybert) (3.10.1)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests->huggingface-hub>=0.20.0->sentence-transformers>=0.3.8->keybert) (2.3.0)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests->huggingface-hub>=0.20.0->sentence-transformers>=0.3.8->keybert) (2025.11.12)
Downloading keybert-0.9.0-py3-none-any.whl (41 kB)
41.4/41.4 kB 1.9 MB/s eta 0:00:00

Installing collected packages: keybert
Successfully installed keybert-0.9.0


```



```

1 import re
2 import spacy
3 import nltk
4 from keybert import KeyBERT
5 from nltk.corpus import stopwords
6
7 # Download stopwords if not available
8 nltk.download("stopwords")
9 stop_words = set(stopwords.words("english"))
10
11 # Load spaCy NLP model
12 nlp = spacy.load("en_core_web_sm")
13
14 # Initialize KeyBERT model
15 kw_model = KeyBERT()
16
17 # Function to preprocess text
18 def preprocess_text(text):
19     """Cleans text and removes unnecessary characters."""
20     text = re.sub(r'\s+', ' ', text) # Remove extra spaces
21     doc = nlp(text.lower()) # Convert to lowercase and tokenize
22     tokens = [token.lemma_ for token in doc if token.is_alpha and token.lemma_ not in stop_words]
23     return " ".join(tokens) # Return as cleaned text
24
25 # Read the transcript file
26 input_file = "/content/segmented_output (2).txt" # Change this to your actual file name
27 with open(input_file, "r", encoding="utf-8") as file:
28     transcript = file.read()
29
30 # Split transcript into segments using 'Segment X:' markers
31 segments = re.split(r'Segment \d+:', transcript)[1:] # Removes empty first split
32
33 # Preprocess each segment
34 cleaned_segments = [preprocess_text(segment.strip()) for segment in segments if segment.strip()]
35
36 # Extract top 3 keywords using KeyBERT
37 output_file = "top_keywords.txt"
38 with open(output_file, "w", encoding="utf-8") as out_file:
39     for i, segment in enumerate(cleaned_segments):
40         keywords = kw_model.extract_keywords(segment, keyphrase_ngram_range=(1, 1), top_n=3)
41
42         # Get top 3 keywords
43         top_words = [word[0] for word in keywords]
44
45         out_file.write(f"Segment {i + 1}: {' '.join(top_words)}\n")
46
47 print(f"Top keywords saved to {output_file}")

```

 [nltk_data] Downloading package stopwords to /root/nltk_data...
 [nltk_data] Package stopwords is already up-to-date!
 /usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
 The secret `HF_TOKEN` does not exist in your Colab secrets.
 To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set it as :
 You will be able to reuse this secret in all of your notebooks.
 Please note that authentication is recommended but still optional to access public models or datasets.

warnings.warn(
modules.json: 100%	349/349 [00:00<00:00, 29.3kB/s]
config_sentence_transformers.json: 100%	116/116 [00:00<00:00, 7.47kB/s]
README.md: 100%	10.5k/10.5k [00:00<00:00, 726kB/s]
sentence_bert_config.json: 100%	53.0/53.0 [00:00<00:00, 3.12kB/s]
config.json: 100%	612/612 [00:00<00:00, 37.6kB/s]
model.safetensors: 100%	90.9M/90.9M [00:01<00:00, 111MB/s]
tokenizer_config.json: 100%	350/350 [00:00<00:00, 20.0kB/s]
vocab.txt: 100%	232k/232k [00:00<00:00, 4.20MB/s]
tokenizer.json: 100%	466k/466k [00:00<00:00, 19.9MB/s]
special_tokens_map.json: 100%	112/112 [00:00<00:00, 10.3kB/s]
config.json: 100%	190/190 [00:00<00:00, 8.24kB/s]
Top keywords saved to top_keywords.txt	

```

1 import yt_dlp
2
3 def download_audio(video_url):
4     ydl_opts = {
5         'format': 'bestaudio/best',
6         'outtmpl': 'linkedin_audio.mp3', # Save file with a meaningful name
7         'cookiefile': '/content/cookies (4).txt', # Ensure you have LinkedIn authentication
8         'postprocessors': [{
9             'key': 'FFmpegExtractAudio',
10            'preferredcodec': 'mp3',
11            'preferredquality': '192',
12        }],
13    }
14
15    with yt_dlp.YoutubeDL(ydl_opts) as ydl:
16        ydl.download([video_url])
17
18    return "linkedin_audio.mp3"
19
20 video_url = input("Enter LinkedIn video URL: ")
21 file_path = download_audio(video_url)
22 print(f"Downloaded audio: {file_path}")
23

```

Enter LinkedIn video URL: https://www.linkedin.com/posts/guillermoflor_elon-musk-wants-the-workplace-to-be-vibrant-activity-7208090955906707457
 [LinkedIn] Extracting URL: https://www.linkedin.com/posts/guillermoflor_elon-musk-wants-the-workplace-to-be-vibrant-activity-7208090955906707457: Downloading webpage
 ERROR: [LinkedIn] 7208090955906707457: Unable to extract video; please report this issue on <https://github.com/yt-dlp/yt-dlp/issues>

```

RegexNotFoundError                                Traceback (most recent call last)
/usr/local/lib/python3.11/dist-packages/yt_dlp/YoutubeDL.py in wrapper(self, *args, **kwargs)
    1636         try:
-> 1637             return func(self, *args, **kwargs)
    1638         except (CookieLoadError, DownloadCancelled, LazyList.IndexError, PagedList.IndexError):

```

11 frames

RegexNotFoundError: [LinkedIn] 7208090955906707457: Unable to extract video; please report this issue on <https://github.com/yt-dlp/yt-dlp/issues>, filling out the appropriate issue template. Confirm you are on the latest version using yt-dlp -U

During handling of the above exception, another exception occurred:

```

DownloadError                                Traceback (most recent call last)
/usr/local/lib/python3.11/dist-packages/yt_dlp/YoutubeDL.py in trouble(self, message, tb, is_error)
    1032         else:
    1033             exc_info = sys.exc_info()
-> 1034         raise DownloadError(message, exc_info)
    1035     self._download_retcode = 1
    1036

```

DownloadError: ERROR: [LinkedIn] 7208090955906707457: Unable to extract video; please report this issue on <https://github.com/yt-dlp/yt-dlp/issues>, filling out the appropriate issue template. Confirm you are on the latest version using yt-dlp -U

```

1 !pip install -U yt-dlp
2

```

Requirement already satisfied: yt-dlp in /usr/local/lib/python3.11/dist-packages (2025.2.19)

```

1 !yt-dlp --cookies-from-browser chrome "https://www.linkedin.com/posts/guillermoflor_elon-musk-wants-the-workplace-to-be-vibrant-acti

```

[LinkedIn] Extracting URL: https://www.linkedin.com/posts/guillermoflor_elon-musk-wants-the-workplace-to-be-vibrant-activity-7208090955906707457: Downloading webpage
 [LinkedIn] 7208090955906707457: Extracting cookies from chrome
 ERROR: could not find chrome cookies database in "/root/.config/google-chrome"

```

1 yt-dlp --cookies-from-browser chrome "https://www.linkedin.com/posts/guillermoflor_elon-musk-wants-the-workplace-to-be-vibrant-acti

```

File "<ipython-input-26-af5310eec6a0>", line 1
 yt-dlp --cookies-from-browser chrome "https://www.linkedin.com/posts/guillermoflor_elon-musk-wants-the-workplace-to-be-vibrant-activity-7208090955906707457-YQ__?utm_source=share&utm_medium=member_desktop&rcm=ACoAAC472oABr5i8229nuG8UuSMZV1GiCnsrmuE"
 ^
 SyntaxError: invalid syntax

TRANSCRIPT SEGMENTATION

```

1 import nltk
2 from nltk.corpus import stopwords
3 from nltk.tokenize import word_tokenize
4 from sklearn.feature_extraction.text import TfidfVectorizer
5 from sklearn.metrics.pairwise import cosine_similarity
6
7 # Ensure necessary resources are downloaded
8 nltk.download('punkt')
9 nltk.download('stopwords')
10
11 def preprocess_text(text):
12     stop_words = set(stopwords.words('english'))
13     words = word_tokenize(text.lower())
14     return [word for word in words if word.isalnum() and word not in stop_words]
15
16 def calculate_similarity(segment_words, new_sentence_words):
17     vectorizer = TfidfVectorizer()
18     combined_text = [" ".join(segment_words), " ".join(new_sentence_words)]
19     tfidf_matrix = vectorizer.fit_transform(combined_text)
20     return cosine_similarity(tfidf_matrix[0], tfidf_matrix[1])[0][0]
21
22 def segment_transcript(file_path, output_path, threshold=0.5):
23     with open(file_path, 'r', encoding='utf-8') as file:
24         lines = file.readlines()
25
26     lines = [line.strip() for line in lines if line.strip()]
27     segments = []
28     i = 0
29
30     while i < len(lines):
31         segment = lines[i:i+5]
32         segment_words = preprocess_text(" ".join(segment))
33
34         new_segment = segment[:]
35         i += 5
36
37         while i < len(lines):
38             new_sentence = lines[i]
39             new_sentence_words = preprocess_text(new_sentence)
40
41             similarity = calculate_similarity(segment_words, new_sentence_words)
42
43             if similarity >= threshold:
44                 new_segment.append(new_sentence)
45                 segment_words.extend(new_sentence_words)
46                 i += 1
47             else:
48                 break
49
50         segments.append(" ".join(new_segment))
51
52     # Start the next segment from the current sentence (not skipping lines)
53
54     with open(output_path, 'w', encoding='utf-8') as output_file:
55         for idx, segment in enumerate(segments):
56             output_file.write(f"Segment {idx + 1}:\n{segment}\n\n")
57
58     return segments
59
60 # Example usage
61 file_path = "/content/transcript (15).txt" # Path to the input .txt file
62 output_path = "segmented_transcript.txt" # Path to the output .txt file
63 segments = segment_transcript(file_path, output_path)
64 for idx, segment in enumerate(segments):
65     print(f"Segment {idx + 1} saved to file.")

```

→ Segment 1 saved to file.

```

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!

```

```

1 import nltk
2 from nltk.corpus import stopwords
3 from nltk.tokenize import word_tokenize
4 from sklearn.feature_extraction.text import TfidfVectorizer
5 from sklearn.metrics.pairwise import cosine_similarity
6
7 # Ensure necessary resources are downloaded
8 nltk.download('punkt')
9 nltk.download('stopwords')
10

```

```

11 def preprocess_text(text):
12     stop_words = set(stopwords.words('english'))
13     words = word_tokenize(text.lower())
14     return [word for word in words if word.isalnum() and word not in stop_words]
15
16 def calculate_similarity(reference_text, new_sentence):
17     vectorizer = TfidfVectorizer()
18     tfidf_matrix = vectorizer.fit_transform([" ".join(reference_text), new_sentence])
19     return cosine_similarity(tfidf_matrix[0], tfidf_matrix[1])[0][0]
20
21 def segment_transcript(file_path, output_path, threshold=1.0):
22     with open(file_path, 'r', encoding='utf-8') as file:
23         lines = [line.strip() for line in file.readlines() if line.strip()]
24
25     segments = []
26     i = 0
27
28     while i < len(lines):
29         segment = lines[i:i+5] # First 5 lines as an initial segment
30         segment_words = preprocess_text(" ".join(segment))
31
32         new_segment = segment[:]
33         i += 5 # Move forward
34
35         while i < len(lines):
36             new_sentence = lines[i]
37             new_sentence_words = preprocess_text(new_sentence)
38
39             similarity = calculate_similarity(segment_words, " ".join(new_sentence_words))
40
41             if similarity >= threshold:
42                 new_segment.append(new_sentence)
43                 i += 1
44             else:
45                 break # Start a new segment
46
47         segments.append(" ".join(new_segment))
48
49     with open(output_path, 'w', encoding='utf-8') as output_file:
50         for idx, segment in enumerate(segments):
51             output_file.write(f"Segment {idx + 1}:\n{segment}\n\n")
52
53     return segments
54
55 # Example usage
56 file_path = "/content/transcript (15).txt"
57 output_path = "segmented_transcript.txt"
58 segments = segment_transcript(file_path, output_path, threshold=0.6)
59
60 for idx, segment in enumerate(segments):
61     print(f"Segment {idx + 1} saved to file.")

```

```

→ Segment 1 saved to file.
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!

```

```

1 import nltk
2 from nltk.corpus import stopwords
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.metrics.pairwise import cosine_similarity
5
6 # Download stopwords if not already downloaded
7 nltk.download('stopwords')
8 nltk.download('punkt')
9
10 # Function to clean text by removing stopwords
11 def clean_text(text):
12     stop_words = set(stopwords.words('english'))
13     words = nltk.word_tokenize(text.lower()) # Tokenize and convert to lowercase
14     return ' '.join([word for word in words if word.isalnum() and word not in stop_words])
15
16 # Read transcript
17 def read_transcript(file_path):
18     with open(file_path, 'r', encoding='utf-8') as file:
19         return file.readlines()
20
21 # Function to group sentences into segments and save them
22 def segment_transcript(sentences, output_file):
23     segments = []
24     current_segment = sentences[:5] # First 5 lines as the initial segment

```

```

25     cleaned_segment_words = clean_text(' '.join(current_segment)) # Remove stopwords
26
27     for i in range(5, len(sentences)):
28         sentence = sentences[i]
29         cleaned_sentence = clean_text(sentence)
30
31         if not cleaned_sentence: # Skip empty sentences after stopwords removal
32             continue
33
34         # Calculate similarity
35         vectorizer = TfidfVectorizer().fit_transform([cleaned_segment_words, cleaned_sentence])
36         similarity = cosine_similarity(vectorizer)[0, 1]
37
38         if similarity > 0.5: # Threshold to decide if it belongs in the same segment
39             current_segment.append(sentence)
40         else:
41             segments.append(current_segment) # Save the current segment
42             current_segment = [sentence] + sentences[i+1:i+5] # New segment with next 4 lines
43             cleaned_segment_words = clean_text(' '.join(current_segment)) # Recalculate cleaned words
44
45     segments.append(current_segment) # Append the last segment
46
47     # Save segmented transcript to a new file
48     with open(output_file, 'w', encoding='utf-8') as file:
49         for idx, segment in enumerate(segments):
50             file.write(f"Segment {idx + 1}:\n")
51             file.write(" ".join(segment) + "\n")
52             file.write("-" * 50 + "\n")
53
54     return segments
55
56 # Example usage
57 if __name__ == "__main__":
58     transcript_lines = read_transcript("/content/transcript (15).txt")
59     segmented_text = segment_transcript(transcript_lines, "segmented_transcript.txt")
60
61     for idx, segment in enumerate(segmented_text):
62         print(f"Segment {idx + 1}:")
63         print(" ".join(segment))
64         print("-" * 50)

```

Segment 1:

So, sir, we know that India has seen a huge revolution with digital payments. We all thought that India is a place, at least the We

[nltk_data] Downloading package stopwords to /root/nltk_data...

[nltk_data] Package stopwords is already up-to-date!

[nltk_data] Downloading package punkt to /root/nltk_data...

[nltk_data] Package punkt is already up-to-date!

```

1 import nltk
2 from nltk.corpus import stopwords
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.metrics.pairwise import cosine_similarity
5
6 # Download stopwords if not already downloaded
7 nltk.download('stopwords')
8 nltk.download('punkt')
9
10 # Function to clean text by removing stopwords
11 def clean_text(text):
12     stop_words = set(stopwords.words('english'))
13     words = nltk.word_tokenize(text.lower()) # Tokenize and convert to lowercase
14     return ' '.join([word for word in words if word.isalnum() and word not in stop_words])
15
16 # Read transcript
17 def read_transcript(file_path):
18     with open(file_path, 'r', encoding='utf-8') as file:
19         return file.readlines()
20
21 # Function to group sentences into segments and save them
22 def segment_transcript(sentences, output_file):
23     segments = []
24     i = 0 # Pointer to track position in sentences
25
26     while i < len(sentences):
27         current_segment = sentences[i:i+5] # First 5 lines as the initial segment
28         cleaned_segment_words = clean_text(' '.join(current_segment)) # Remove stopwords
29         i += 5 # Move the pointer forward
30
31         while i < len(sentences):
32             cleaned_sentence = clean_text(sentences[i])

```

```

33
34     if not cleaned_sentence: # Skip empty sentences after stopwords removal
35         i += 1
36         continue
37
38     # Calculate similarity
39     vectorizer = TfidfVectorizer().fit_transform([cleaned_segment_words, cleaned_sentence])
40     similarity = cosine_similarity(vectorizer)[0, 1]
41
42     if similarity > 0.5: # Threshold to decide if it belongs in the same segment
43         current_segment.append(sentences[i])
44         i += 1 # Move to the next sentence
45     else:
46         break # Stop adding to this segment, start a new one
47
48     segments.append(current_segment)
49
50 # Save segmented transcript to a new file
51 with open(output_file, 'w', encoding='utf-8') as file:
52     for idx, segment in enumerate(segments):
53         file.write(f"Segment {idx + 1}:\n")
54         file.write("".join(segment) + "\n")
55         file.write("-" * 50 + "\n")
56
57 return segments
58
59 # Example usage
60 if __name__ == "__main__":
61     transcript_lines = read_transcript("/content/transcript (15).txt")
62     segmented_text = segment_transcript(transcript_lines, "segmented_transcript.txt")
63
64     for idx, segment in enumerate(segmented_text):
65         print(f"Segment {idx + 1}:")
66         print("".join(segment))
67         print("-" * 50)

```

↩ Segment 1:

So, sir, we know that India has seen a huge revolution with digital payments. We all thought that India is a place, at least the We

[nltk_data] Downloading package stopwords to /root/nltk_data...

[nltk_data] Package stopwords is already up-to-date!

[nltk_data] Downloading package punkt to /root/nltk_data...

[nltk_data] Package punkt is already up-to-date!

```

1 import nltk
2 from nltk.corpus import stopwords
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.metrics.pairwise import cosine_similarity
5
6 # Download stopwords if not already downloaded
7 nltk.download('stopwords')
8 nltk.download('punkt')
9
10 # Function to clean text by removing stopwords
11 def clean_text(text):
12     stop_words = set(stopwords.words('english'))
13     words = nltk.word_tokenize(text.lower()) # Tokenize and convert to lowercase
14     return ' '.join([word for word in words if word.isalnum() and word not in stop_words])
15
16 # Read transcript
17 def read_transcript(file_path):
18     with open(file_path, 'r', encoding='utf-8') as file:
19         return file.readlines()
20
21 # Function to group sentences into segments and save them
22 def segment_transcript(sentences, output_file):
23     segments = []
24     i = 0 # Pointer to track position in sentences
25
26     while i < len(sentences):
27         current_segment = sentences[i:i+5] # First 5 lines as the initial segment
28         cleaned_segment_words = clean_text(' '.join(current_segment)) # Remove stopwords
29         i += 5 # Move the pointer forward
30
31         while i < len(sentences):
32             cleaned_sentence = clean_text(sentences[i])
33
34             if not cleaned_sentence: # Skip empty sentences after stopwords removal
35                 i += 1
36                 continue
37

```

```

38         # Calculate similarity
39         vectorizer = TfidfVectorizer().fit_transform([cleaned_segment_words, cleaned_sentence])
40         similarity = cosine_similarity(vectorizer)[0, 1]
41
42         if similarity > 0.2: # Threshold to decide if it belongs in the same segment
43             current_segment.append(sentences[i])
44             i += 1 # Move to the next sentence
45         else:
46             break # Stop adding to this segment, start a new one
47
48         segments.append(current_segment)
49         i += 5 # Move the pointer forward to ensure the next segment starts correctly
50
51     # Save segmented transcript to a new file
52     with open(output_file, 'w', encoding='utf-8') as file:
53         for idx, segment in enumerate(segments):
54             file.write(f"Segment {idx + 1}:\n")
55             file.write("".join(segment) + "\n")
56             file.write("-" * 50 + "\n")
57
58     return segments
59
60 # Example usage
61 if __name__ == "__main__":
62     transcript_lines = read_transcript("/content/transcript (15).txt")
63     segmented_text = segment_transcript(transcript_lines, "segmented_transcript.txt")
64
65     for idx, segment in enumerate(segmented_text):
66         print(f"Segment {idx + 1}:")
67         print("".join(segment))
68         print("-" * 50)
69

```

➡ Segment 1:
 So, sir, we know that India has seen a huge revolution with digital payments. We all thought that India is a place, at least the We

[nltk_data] Downloading package stopwords to /root/nltk_data...
 [nltk_data] Package stopwords is already up-to-date!
 [nltk_data] Downloading package punkt to /root/nltk_data...
 [nltk_data] Package punkt is already up-to-date!

```

1 import nltk
2 from nltk.corpus import stopwords
3 from nltk.tokenize import sent_tokenize
4 from sklearn.feature_extraction.text import TfidfVectorizer
5 from sklearn.metrics.pairwise import cosine_similarity
6
7 # Download required NLTK data
8 nltk.download('stopwords')
9 nltk.download('punkt')
10
11 # Function to clean text by removing stopwords
12 def clean_text(text):
13     stop_words = set(stopwords.words('english'))
14     words = nltk.word_tokenize(text.lower()) # Tokenize and convert to lowercase
15     return ' '.join([word for word in words if word.isalnum() and word not in stop_words])
16
17 # Read transcript file
18 def read_transcript(file_path):
19     with open(file_path, 'r', encoding='utf-8') as file:
20         return sent_tokenize(file.read()) # Use sentence tokenizer instead of readlines
21
22 # Function to segment transcript
23 def segment_transcript(sentences, output_file):
24     segments = []
25     i = 0 # Pointer to track position
26
27     while i < len(sentences):
28         current_segment = [sentences[i]] # Start with first sentence
29         cleaned_segment = clean_text(sentences[i])
30         i += 1 # Move forward
31
32         while i < len(sentences):
33             cleaned_sentence = clean_text(sentences[i])
34
35             if not cleaned_sentence: # Skip empty cleaned sentences
36                 i += 1
37                 continue
38
39             # Vectorize segment and new sentence
40             vectorizer = TfidfVectorizer()

```



```

41     vectors = vectorizer.fit_transform([cleaned_segment, cleaned_sentence])
42     similarity = cosine_similarity(vectors)[0, 1]
43
44     if similarity > 0.9: # Threshold for grouping
45         current_segment.append(sentences[i])
46         cleaned_segment += " " + cleaned_sentence # Update segment content
47         i += 1 # Move to next sentence
48     else:
49         break # Stop adding, create a new segment
50
51     segments.append(current_segment)
52
53 # Save segmented transcript
54 with open(output_file, 'w', encoding='utf-8') as file:
55     for idx, segment in enumerate(segments):
56         file.write(f"Segment {idx + 1}:\n")
57         file.write(" ".join(segment) + "\n")
58         file.write("-" * 50 + "\n")
59
60     return segments
61
62 # Example usage
63 if __name__ == "__main__":
64     transcript_sentences = read_transcript("/content/transcript (15).txt")
65     segmented_text = segment_transcript(transcript_sentences, "segmented_transcript.txt")
66
67     for idx, segment in enumerate(segmented_text):
68         print(f"Segment {idx + 1}:")
69         print(" ".join(segment))
70         print("-" * 50)
71

```

```

[ntlk_data] Downloading package stopwords to /root/nltk_data...
[ntlk_data] Package stopwords is already up-to-date!
[ntlk_data] Downloading package punkt to /root/nltk_data...
[ntlk_data] Package punkt is already up-to-date!
Segment 1:
So, sir, we know that India has seen a huge revolution with digital payments.
-----
Segment 2:
We all thought that India is a place, at least the West thought that India is a place where many people do not get a square mill,
-----
Segment 3:
That was an narrative some 30 years ago.
-----
Segment 4:
And not many are literate, people cannot read.
-----
Segment 5:
But then we have now shown that digital payments, number one is India, while people thought that it wouldn't even come to top 50,
-----
Segment 6:
I think immediately after UPI the next big revolution, personally I think is in education.
-----
Segment 7:
And the complete homework for this has happened in the form of NEP, the documentation of which many of us have read and realized
-----
Segment 8:
So, followed by which we got NCRF framework done, which is a national creative framework, surrounding which we will be discussing
-----
Segment 9:
So, my question is, do you think NCRF plus NEP put together will be the next big revolution after UPI in India?
-----
Segment 10:
Absolutely.
-----
Segment 11:
And why I think so is because in education the last policy came up many many years ago that was in 1986, which was slightly tweek
-----
Segment 12:
And thereafter so many changes have happened in the real world, so many changes have happened in the requirement of the industry,
-----
Segment 13:
However there were no corresponding changes which happened in education system.
-----
Segment 14:
So, therefore I feel that this was the right time when we brought in the education policy 2020, honorable prime minister declared
-----
Segment 15:
And we recently celebrated the fourth anniversary of NEP 2020.
-----
Segment 16:
NCRF has been brought to implement the intent of 2020, 2020 is a policy, NEP 2020 is a policy.
-----
Segment 17:
And for implementing a policy you need a framework.
-----

```

```

Segment 18:
"""
1 import nltk
2 from nltk.tokenize import sent_tokenize
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.metrics.pairwise import cosine_similarity
5
6 # Download required data
7 nltk.download('punkt')
8
9 # Read transcript file
10 def read_transcript(file_path):
11     with open(file_path, 'r', encoding='utf-8') as file:
12         return sent_tokenize(file.read()) # Tokenize entire text into sentences
13
14 # Function to segment transcript using similarity
15 def segment_transcript(sentences, output_file, threshold=0.9):
16     segments = []
17     i = 0 # Pointer to track position
18
19     # Fit a single TF-IDF vectorizer on all sentences for consistency
20     vectorizer = TfidfVectorizer(stop_words="english")
21     sentence_vectors = vectorizer.fit_transform(sentences) # Vectorize entire text
22
23     while i < len(sentences):
24         current_segment = [sentences[i]] # Start with the first sentence
25         segment_indices = [i]
26         i += 1 # Move pointer forward
27
28         while i < len(sentences):
29             # Compute cosine similarity between last sentence of segment and next sentence
30             similarity = cosine_similarity(sentence_vectors[segment_indices[-1]], sentence_vectors[i])[0, 0]
31
32             if similarity >= threshold: # If similar, add to the segment
33                 current_segment.append(sentences[i])
34                 segment_indices.append(i) # Track index of included sentence
35                 i += 1
36             else:
37                 break # Stop and start a new segment
38
39         segments.append(current_segment) # Store the formed segment
40
41     # Save segmented transcript
42     with open(output_file, 'w', encoding='utf-8') as file:
43         for idx, segment in enumerate(segments):
44             file.write(f"Segment {idx + 1}:\n")
45             file.write(" ".join(segment) + "\n")
46             file.write("-" * 50 + "\n")
47
48     return segments
49
50 # Example usage
51 if __name__ == "__main__":
52     transcript_sentences = read_transcript("/content/transcript (15).txt")
53     segmented_text = segment_transcript(transcript_sentences, "segmented_transcript.txt")
54
55     for idx, segment in enumerate(segmented_text):
56         print(f"Segment {idx + 1}:")
57         print(" ".join(segment))
58         print("-" * 50)

```

```

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
Segment 1:
So, sir, we know that India has seen a huge revolution with digital payments.
-----
Segment 2:
We all thought that India is a place, at least the West thought that India is a place where many people do not get a square mill,
-----
Segment 3:
That was an narrative some 30 years ago.
-----
Segment 4:
And not many are literate, people cannot read.
-----
Segment 5:
But then we have now shown that digital payments, number one is India, while people thought that it wouldn't even come to top 50,
-----
Segment 6:
I think immediately after UPI the next big revolution, personally I think is in education.
-----
Segment 7:
And the complete homework for this has happened in the form of NEP, the documentation of which many of us have read and realized

```



```

-----
Segment 8:
So, followed by which we got NCRF framework done, which is a national creative framework, surrounding which we will be discussing
-----
Segment 9:
So, my question is, do you think NCRF plus NEP put together will be the next big revolution after UPI in India?
-----
Segment 10:
Absolutely.
-----
Segment 11:
And why I think so is because in education the last policy came up many many years ago that was in 1986, which was slightly tweaked
-----
Segment 12:
And thereafter so many changes have happened in the real world, so many changes have happened in the requirement of the industry,
-----
Segment 13:
However there were no corresponding changes which happened in education system.
-----
Segment 14:
So, therefore I feel that this was the right time when we brought in the education policy 2020, honorable prime minister declared
-----
Segment 15:
And we recently celebrated the fourth anniversary of NEP 2020.
-----
Segment 16:
NCRF has been brought to implement the intent of 2020, 2020 is a policy, NEP 2020 is a policy.
-----
Segment 17:
And for implementing a policy you need a framework.
-----
Segment 18:
Now, why we call it a framework?
-----
Segment 19:

```

```

1 # prompt: 1. Read the transcript
2 # 2. Separate each line
3 # 3. By default we have the first 5 lines in segment 1
4 # 4. Then we will remove the stop words from that segment and put the remaining words in a new list
5 # 5. Now we will start with the comparison with next lines one by one
6 # 6. Remove the stop words from the sentences below them and compare with those of the list from segment 1
7 # 7. If the stop words from the sentence 6 are similar to those of segment 1 (Finding the cosine similarity) Put them in segment 1
8 # 8. The next segment will consists of this sentence 6 and the next 4 chronological sentences should also be kept in segment 2.
9 # 9. Now we will repeat the steps 4-8 for the rest of the transcript.
10
11 import nltk
12 from nltk.tokenize import sent_tokenize
13 from sklearn.feature_extraction.text import TfidfVectorizer
14 from sklearn.metrics.pairwise import cosine_similarity
15 from nltk.corpus import stopwords
16
17 # Download required data
18 nltk.download('punkt')
19 nltk.download('stopwords')
20
21 # Function to clean text by removing stopwords
22 def clean_text(text):
23     stop_words = set(stopwords.words('english'))
24     words = nltk.word_tokenize(text.lower())
25     return ' '.join([word for word in words if word.isalnum() and word not in stop_words])
26
27 # Read transcript file
28 def read_transcript(file_path):
29     with open(file_path, 'r', encoding='utf-8') as file:
30         return sent_tokenize(file.read())
31
32 # Function to segment transcript using similarity
33 def segment_transcript(sentences, output_file, threshold=1):
34     segments = []
35     i = 0
36
37     # Fit a single TF-IDF vectorizer on all sentences for consistency
38     vectorizer = TfidfVectorizer(stop_words="english")
39     sentence_vectors = vectorizer.fit_transform(sentences)
40
41     while i < len(sentences):
42         current_segment = [sentences[i]]
43         segment_indices = [i]
44         i += 1
45
46         while i < len(sentences):
47             similarity = cosine_similarity(sentence_vectors[segment_indices[-1]], sentence_vectors[i])[0, 0]
48

```

```

49         if similarity >= threshold:
50             current_segment.append(sentences[i])
51             segment_indices.append(i)
52             i += 1
53         else:
54             break
55
56     segments.append(current_segment)
57
58     # Save segmented transcript
59     with open(output_file, 'w', encoding='utf-8') as file:
60         for idx, segment in enumerate(segments):
61             file.write(f"Segment {idx + 1}:\n")
62             file.write(" ".join(segment) + "\n")
63             file.write("-" * 50 + "\n")
64
65     return segments
66
67 # Example usage
68 if __name__ == "__main__":
69     transcript_sentences = read_transcript("/content/transcript (15).txt")
70     segmented_text = segment_transcript(transcript_sentences, "segmented_transcript.txt")
71
72     for idx, segment in enumerate(segmented_text):
73         print(f"Segment {idx + 1}:")
74         print(" ".join(segment))
75         print("-" * 50)
76

```

```

[!nltk_data] Downloading package punkt to /root/nltk_data...
[!nltk_data] Package punkt is already up-to-date!
[!nltk_data] Downloading package stopwords to /root/nltk_data...
[!nltk_data] Package stopwords is already up-to-date!
Segment 1:
So, sir, we know that India has seen a huge revolution with digital payments.
-----
Segment 2:
We all thought that India is a place, at least the West thought that India is a place where many people do not get a square mill,
-----
Segment 3:
That was an narrative some 30 years ago.
-----
Segment 4:
And not many are literate, people cannot read.
-----
Segment 5:
But then we have now shown that digital payments, number one is India, while people thought that it wouldn't even come to top 50,
-----
Segment 6:
I think immediately after UPI the next big revolution, personally I think is in education.
-----
Segment 7:
And the complete homework for this has happened in the form of NEP, the documentation of which many of us have read and realized
-----
Segment 8:
So, followed by which we got NCRF framework done, which is a national creative framework, surrounding which we will be discussing
-----
Segment 9:
So, my question is, do you think NCRF plus NEP put together will be the next big revolution after UPI in India?
-----
Segment 10:
Absolutely.
-----
Segment 11:
And why I think so is because in education the last policy came up many many years ago that was in 1986, which was slightly tweaked
-----
Segment 12:
And thereafter so many changes have happened in the real world, so many changes have happened in the requirement of the industry,
-----
Segment 13:
However there were no corresponding changes which happened in education system.
-----
Segment 14:
So, therefore I feel that this was the right time when we brought in the education policy 2020, honorable prime minister declared
-----
Segment 15:
And we recently celebrated the fourth anniversary of NEP 2020.
-----
Segment 16:
NCRF has been brought to implement the intent of 2020, 2020 is a policy, NEP 2020 is a policy.
-----
Segment 17:
And for implementing a policy you need a framework.
-----
Segment 18:
Now, why we call it a framework?

```

```
1 #Using LLM
```

```
1 pip uninstall bitsandbytes -y
```

⚠ WARNING: Skipping bitsandbytes as it is not installed.

```
1 pip install bitsandbytes --no-cache-dir
```

⚠ Requirement already satisfied: bitsandbytes in /usr/local/lib/python3.11/dist-packages (0.45.3)
 Requirement already satisfied: torch<3,>=2.0 in /usr/local/lib/python3.11/dist-packages (from bitsandbytes) (2.6.0+cu124)
 Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.11/dist-packages (from bitsandbytes) (2.0.2)
 Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (3.17.0)
 Requirement already satisfied: typing-extensions>=4.10.0 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (4.12.0)
 Requirement already satisfied: networkx in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (3.4.2)
 Requirement already satisfied: Jinja2 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (3.1.6)
 Requirement already satisfied: fsspec in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (2024.10.0)
 Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.127)
 Requirement already satisfied: nvidia-cuda-runtime-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.127)
 Requirement already satisfied: nvidia-cuda-cupti-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.127)
 Requirement already satisfied: nvidia-cudnn-cu12==9.1.0.70 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (9.1.0.70)
 Requirement already satisfied: nvidia-cublas-cu12==12.4.5.8 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.5.8)
 Requirement already satisfied: nvidia-cufft-cu12==11.2.1.3 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (11.2.1.3)
 Requirement already satisfied: nvidia-curand-cu12==10.3.5.147 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (10.3.5.147)
 Requirement already satisfied: nvidia-cusolver-cu12==11.6.1.9 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (11.6.1.9)
 Requirement already satisfied: nvidia-cusparse-cu12==12.3.1.170 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.3.1.170)
 Requirement already satisfied: nvidia-cusparselt-cu12==0.6.2 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (0.6.2)
 Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (2.21.5)
 Requirement already satisfied: nvidia-nvtx-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.127)
 Requirement already satisfied: nvidia-nvjitlink-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.127)
 Requirement already satisfied: triton==3.2.0 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (3.2.0)
 Requirement already satisfied: sympy==1.13.1 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (1.13.1)
 Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from sympy==1.13.1->torch<3,>=2.0->bitsandbytes) (1.3.0)
 Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from Jinja2->torch<3,>=2.0->bitsandbytes) (2.1.5)

```
1 from transformers import pipeline
```

```
1 from transformers import AutoModelForCausalLM, AutoTokenizer, pipeline
```

```
1 import torch
```

```
1 # Ensure GPU is available
```

```
2 assert torch.cuda.is_available(), "GPU is not available! Check your CUDA installation."
```

```
1 pipe = pipeline("text-generation", model="deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B")
```

⚠ config.json: 100% 679/679 [00:00<00:00, 16.8kB/s]
 model.safetensors: 100% 3.55G/3.55G [00:23<00:00, 163MB/s]
 generation_config.json: 100% 181/181 [00:00<00:00, 6.96kB/s]
 tokenizer_config.json: 100% 3.07k/3.07k [00:00<00:00, 193kB/s]
 tokenizer.json: 100% 7.03M/7.03M [00:00<00:00, 24.3MB/s]

```
1 with open("/content/transcript (15).txt", "r", encoding="utf-8") as file:
```

```
2     transcript_text = file.read()
```

```
1 inputs = [
2     {
3         "role": "system",
4         "content": "You are a helpful assistant tasked with segmenting a long video lecture transcript into clear, distinct topic sections."
5     },
6     {
7         "role": "user",
8         "content": f"You have been provided with a long video lecture transcript. Your task is to read the transcript, split it into 10 sections, and provide a summary of each section."
9     },
10    {
11        "role": "assistant",
12        "content": "Please segment the transcript provided below into clear sections without altering the original content. Ensure each section is clearly defined and easy to understand."
13    },
14    {
15        "role": "user",
16        "content": transcript_text
17    }
18 ]
```

```

1 a = pipe(inputs, max_length=500000)

1 Start coding or generate with AI.

1 #The Method we discussed

1 import nltk
2 import numpy as np
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.metrics.pairwise import cosine_similarity
5 from nltk.corpus import stopwords
6
7 nltk.download('stopwords')
8 nltk.download('punkt')
9 nltk.download('punkt_tab')
10 stop_words = set(stopwords.words('english'))
11
12 def preprocess_text(text):
13     words = nltk.word_tokenize(text.lower())
14     return ' '.join([word for word in words if word.isalnum() and word not in stop_words])
15
16 def segment_transcript(transcript):
17     lines = transcript.strip().split('\n')
18     segments = []
19     current_segment = lines[:5] # First 5 lines
20     processed_segment = [preprocess_text(line) for line in current_segment]
21     vectorizer = TfidfVectorizer()
22     segment_vectors = vectorizer.fit_transform(processed_segment)
23
24     i = 5
25     while i < len(lines):
26         sentence = lines[i]
27         processed_sentence = preprocess_text(sentence)
28         sentence_vector = vectorizer.transform([processed_sentence])
29
30         similarity_scores = cosine_similarity(sentence_vector, segment_vectors).flatten()
31         max_similarity = np.max(similarity_scores) if len(similarity_scores) > 0 else 0
32
33         if max_similarity > 0.5: # Threshold for similarity
34             current_segment.append(sentence)
35             processed_segment.append(processed_sentence)
36             segment_vectors = vectorizer.fit_transform(processed_segment)
37         else:
38             segments.append(current_segment) # Store completed segment
39             current_segment = [sentence] + lines[i+1:i+5] # Start new segment
40             processed_segment = [preprocess_text(line) for line in current_segment]
41             segment_vectors = vectorizer.fit_transform(processed_segment)
42             i += 4 # Skip next 4 lines as they are in the new segment
43
44         i += 1
45
46     if current_segment:
47         segments.append(current_segment) # Append last segment
48
49     return segments
50
51 # Example usage
52 transcript = """ So, sir, we know that India has seen a huge revolution with digital payments. We all thought that India is a place,
53 segmented_output = segment_transcript(transcript)
54
55 for idx, segment in enumerate(segmented_output, 1):
56     print(f"Segment {idx}:")
57     print("\n".join(segment))
58     print("-" * 50)
59
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data] Unzipping tokenizers/punkt_tab.zip.
Segment 1:
So, sir, we know that India has seen a huge revolution with digital payments. We all thought that India is a place, at least the Wes
-----

```

1 Start coding or [generate](#) with AI.

