

```

1 # The exclamation mark (!) is used to run shell commands directly from Jupyter Notebook or Google Colab.
2 # This command installs the SpeechRecognition library, which is used for converting speech to text in Python.
3 #!pip install speechrecognition

```

```

Collecting speechrecognition
  Downloading SpeechRecognition-3.14.1-py3-none-any.whl.metadata (31 kB)
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.11/dist-packages (from speechrecognition) (4.12.2)
Downloading SpeechRecognition-3.14.1-py3-none-any.whl (32.9 MB)
    32.9/32.9 MB 14.8 MB/s eta 0:00:00
Installing collected packages: speechrecognition
Successfully installed speechrecognition-3.14.1

```

```

1 # The exclamation mark (!) is used to run shell commands directly from Jupyter Notebook or Google Colab.
2 # This command installs the pydub library, which is used for audio processing tasks like converting, slicing, and merging audio files.
3 #!pip install pydub

```

```

Collecting pydub
  Downloading pydub-0.25.1-py2.py3-none-any.whl.metadata (1.4 kB)
  Downloading pydub-0.25.1-py2.py3-none-any.whl (32 kB)
Installing collected packages: pydub
Successfully installed pydub-0.25.1

```

```

1 # This command installs the openai-whisper library, which is a speech recognition model by OpenAI.
2 # Whisper is capable of transcribing audio into text and supports multiple languages.
3 # The exclamation mark (!) is used to run shell commands directly from Jupyter Notebook or Google Colab.
4 #!pip install openai-whisper

```

```

Collecting openai-whisper
  Downloading openai-whisper-20240930.tar.gz (800 kB)
    800.5/800.5 kB 12.4 MB/s eta 0:00:00
Installing build dependencies ... done
Getting requirements to build wheel ... done
Preparing metadata (pyproject.toml) ... done
Requirement already satisfied: numba in /usr/local/lib/python3.11/dist-packages (from openai-whisper) (0.61.0)
Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (from openai-whisper) (1.26.4)
Requirement already satisfied: torch in /usr/local/lib/python3.11/dist-packages (from openai-whisper) (2.5.1+cu124)
Requirement already satisfied: tqdm in /usr/local/lib/python3.11/dist-packages (from openai-whisper) (4.67.1)
Requirement already satisfied: more-itertools in /usr/local/lib/python3.11/dist-packages (from openai-whisper) (10.6.0)
Collecting tiktoken (from openai-whisper)
  Downloading tiktoken-0.9.0-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (6.7 kB)
Requirement already satisfied: triton>=2.0.0 in /usr/local/lib/python3.11/dist-packages (from openai-whisper) (3.1.0)
Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from triton>=2.0.0->openai-whisper) (3.17.0)
Requirement already satisfied: llvmlite<0.45,>=0.44.0dev0 in /usr/local/lib/python3.11/dist-packages (from numba->openai-whisper) (0.44.0)
Requirement already satisfied: regex>=2022.1.18 in /usr/local/lib/python3.11/dist-packages (from tiktoken->openai-whisper) (2024.11.1)
Requirement already satisfied: requests>=2.26.0 in /usr/local/lib/python3.11/dist-packages (from tiktoken->openai-whisper) (2.32.0)
Requirement already satisfied: typing-extensions>=4.8.0 in /usr/local/lib/python3.11/dist-packages (from torch->openai-whisper) (4.12.2)
Requirement already satisfied: networkx in /usr/local/lib/python3.11/dist-packages (from torch->openai-whisper) (3.4.2)
Requirement already satisfied: jinja2 in /usr/local/lib/python3.11/dist-packages (from torch->openai-whisper) (3.1.5)
Requirement already satisfied: fsspec in /usr/local/lib/python3.11/dist-packages (from torch->openai-whisper) (2024.10.0)
Collecting nvidia-cuda-nvrtc-cu12==12.4.127 (from torch->openai-whisper)
  Downloading nvidia_cuda_nvrtc_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cuda-runtime-cu12==12.4.127 (from torch->openai-whisper)
  Downloading nvidia_cuda_runtime_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cuda-cupti-cu12==12.4.127 (from torch->openai-whisper)
  Downloading nvidia_cuda_cupti_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cudnn-cu12==9.1.0.70 (from torch->openai-whisper)
  Downloading nvidia_cudnn_cu12-9.1.0.70-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cublas-cu12==12.4.5.8 (from torch->openai-whisper)
  Downloading nvidia_cublas_cu12-12.4.5.8-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cufft-cu12==11.2.1.3 (from torch->openai-whisper)
  Downloading nvidia_cufft_cu12-11.2.1.3-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-curand-cu12==10.3.5.147 (from torch->openai-whisper)
  Downloading nvidia_curand_cu12-10.3.5.147-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cusolver-cu12==11.6.1.9 (from torch->openai-whisper)
  Downloading nvidia_cusolver_cu12-11.6.1.9-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cuspars-cu12==12.3.1.170 (from torch->openai-whisper)
  Downloading nvidia_cuspars-cu12-12.3.1.170-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in /usr/local/lib/python3.11/dist-packages (from torch->openai-whisper) (2.21.5)
Requirement already satisfied: nvidia-nvtx-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch->openai-whisper) (12.4.127)
Collecting nvidia-nvjitlink-cu12==12.4.127 (from torch->openai-whisper)
  Downloading nvidia_nvjitlink_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Requirement already satisfied: sympy==1.13.1 in /usr/local/lib/python3.11/dist-packages (from torch->openai-whisper) (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->openai-whisper) (1.3.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests>=2.26.0->tiktoken->openai-whisper) (3.3.2)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests>=2.26.0->tiktoken->openai-whisper) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests>=2.26.0->tiktoken->openai-whisper) (2.31.2)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests>=2.26.0->tiktoken->openai-whisper) (2025.1.1)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from jinja2->torch->openai-whisper) (3.0.2)
Downloading tiktoken-0.9.0-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (1.2 MB)
    1.2/1.2 MB 47.7 MB/s eta 0:00:00
Downloading nvidia_cublas_cu12-12.4.5.8-py3-none-manylinux2014_x86_64.whl (363.4 MB)
    363.4/363.4 MB 1.5 MB/s eta 0:00:00
Downloading nvidia_cuda_cupti_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl (13.8 MB)
    13.8/13.8 MB 95.9 MB/s eta 0:00:00

```

```

1 '''# Import the AudioSegment class from the pydub library for audio processing tasks.
2 from pydub import AudioSegment
3
4 # Define a function to convert a WAV file to MP3 format.
5 # Parameters:
6 # - wav_file: The path to the input WAV file.
7 # - mp3_file: (Optional) The path to the output MP3 file. If not provided, it will be generated automatically.
8 def convert_wav_to_mp3(wav_file, mp3_file=None):
9     # Load the WAV file as an AudioSegment object.
10    audio = AudioSegment.from_wav(wav_file)
11
12    # If the output MP3 file name is not provided, generate it by replacing the .wav extension with .mp3.
13    if not mp3_file:
14        mp3_file = wav_file.rsplit('.', 1)[0] + ".mp3"
15
16    # Export the audio data as an MP3 file.
17    audio.export(mp3_file, format="mp3")
18
19    # Print a message indicating that the conversion is complete.
20    print(f"Conversion complete: {mp3_file}")
21
22 # Example usage of the function to convert a WAV file to MP3 format.
23 convert_wav_to_mp3("/content/converted_audio.wav")'''

1 '''# Import the whisper library for speech recognition and transcription tasks.
2 import whisper
3
4 # Define a function to transcribe audio files using the Whisper model.
5 # Parameters:
6 # - audio_file: The path to the input audio file (e.g., .mp3 or .wav).
7 # - model_size: The size of the Whisper model to use (default is "base").
8 # Other options include "tiny", "small", "medium", and "large" for different accuracy and speed trade-offs.
9 def transcribe_audio(audio_file, model_size="base"):
10    # Load the Whisper model of the specified size.
11    model = whisper.load_model(model_size)
12
13    # Transcribe the audio file and store the result.
14    result = model.transcribe(audio_file)
15
16    # Generate a file name for saving the transcription by replacing the original file extension with "_transcript.txt".
17    transcript_file = audio_file.rsplit('.', 1)[0] + "_transcript.txt"
18
19    # Save the transcription text to a file in UTF-8 encoding.
20    with open(transcript_file, "w", encoding="utf-8") as f:
21        f.write(result["text"])
22
23    # Print a message indicating where the transcription has been saved.
24    print(f"Transcription saved to {transcript_file}")
25
26    # Return the transcription text.
27    return result["text"]
28
29 # Example usage of the function:
30 # Specify the path to the audio file to be transcribed.
31 audio_path = "/content/converted_audio.mp3" # Change this to your actual file path
32
33 # Call the transcription function and print the transcript.
34 transcript = transcribe_audio(audio_path)
35 print("Transcript:", transcript)
36

1 '''# Install the pydub library, which is used for audio processing tasks like converting, slicing, and merging audio files.
2 # The exclamation mark (!) is used to run shell commands directly from Jupyter Notebook or Google Colab.
3 !pip install pydub
4
5 # Install the ffmpeg package, which is a powerful tool for handling audio and video files.
6 # ffmpeg is required by pydub to handle various audio formats (e.g., MP3 to WAV conversion).
7 # The 'apt install' command is used to install packages on Debian-based Linux systems, like Ubuntu and Google Colab.
8 !apt install ffmpeg
9

🔗 Requirement already satisfied: pydub in /usr/local/lib/python3.11/dist-packages (0.25.1)
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
ffmpeg is already the newest version (7:4.4.2-0ubuntu0.22.04.1).
0 upgraded, 0 newly installed, 0 to remove and 29 not upgraded.

1 '''# Install the moviepy library, which is used for video editing tasks like cutting, concatenating, and converting videos.
2 # It also supports extracting and manipulating audio from video files.

```

```

3 # The exclamation mark (!) is used to run shell commands directly from Jupyter Notebook or Google Colab.
4 !pip install moviepy
5

```

```

Requirement already satisfied: moviepy in /usr/local/lib/python3.11/dist-packages (1.0.3)
Requirement already satisfied: decorator<5.0,>=4.0.2 in /usr/local/lib/python3.11/dist-packages (from moviepy) (4.4.2)
Requirement already satisfied: imageio<3.0,>=2.5 in /usr/local/lib/python3.11/dist-packages (from moviepy) (2.37.0)
Requirement already satisfied: imageio_ffmpeg>=0.2.0 in /usr/local/lib/python3.11/dist-packages (from moviepy) (0.6.0)
Requirement already satisfied: tqdm<5.0,>=4.11.2 in /usr/local/lib/python3.11/dist-packages (from moviepy) (4.67.1)
Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.11/dist-packages (from moviepy) (1.26.4)
Requirement already satisfied: requests<3.0,>=2.8.1 in /usr/local/lib/python3.11/dist-packages (from moviepy) (2.32.3)
Requirement already satisfied: proglog<=1.0.0 in /usr/local/lib/python3.11/dist-packages (from moviepy) (0.1.10)
Requirement already satisfied: pillow>=8.3.2 in /usr/local/lib/python3.11/dist-packages (from imageio<3.0,>=2.5->moviepy) (11.1.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests<3.0,>=2.8.1->moviepy) (3.10)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests<3.0,>=2.8.1->moviepy) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests<3.0,>=2.8.1->moviepy) (2.2.3)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests<3.0,>=2.8.1->moviepy) (2025.1.1)

```

```

1 '''# Import the VideoFileClip class from the moviepy library for video processing tasks.
2 from moviepy.editor import VideoFileClip
3
4 # Define the path to the video file that you want to process.
5 video_path = "/content/videoplayback.mp4"
6
7 # Load the video file into a VideoFileClip object.
8 video = VideoFileClip(video_path)
9
10 # Extract the audio track from the video.
11 audio = video.audio
12
13 # Define the path and format for the extracted audio file.
14 # Here, we save the audio in MP3 format, but you can change it to WAV or other formats if needed.
15 audio_path = "extracted_audio.mp3"
16 audio.write_audiofile(audio_path)
17
18 # Print a message indicating that the audio extraction is complete and specify the file name.
19 print(f"Audio extracted and saved as {audio_path}")
20

```

```

WARNING:py.warnings:/usr/local/lib/python3.11/dist-packages/moviepy/video/io/sliders.py:61: SyntaxWarning: "is" with a literal. Did
you mean "is not"?
if event.key is 'enter':

```

```

MoviePy - Writing audio in extracted_audio.mp3
MoviePy - Done.
Audio extracted and saved as extracted_audio.mp3

```

```

1 import IPython.display as ipd
2
3 # Play the extracted audio
4 ipd.Audio("/content/extracted_audio.mp3")

```



```

1 from pydub import AudioSegment
2
3 # Load the audio file
4 audio = AudioSegment.from_mp3("/content/extracted_audio.mp3")
5
6 # Get the duration of the audio in milliseconds
7 duration = len(audio)
8
9 # Calculate the midpoint
10 midpoint = duration // 2
11
12 # Split the audio into two halves
13 first_half = audio[:midpoint]
14 second_half = audio[midpoint:]
15
16 # Export the two halves as separate files
17 first_half.export("/content/first_half.mp3", format="mp3")
18 second_half.export("/content/second_half.mp3", format="mp3")
19
20 print("Audio has been split into two halves successfully!")

```

```

Audio has been split into two halves successfully!

```



```

1 # prompt: Generate code to read the video url and extract the transcript of a video
2
3 # Install necessary libraries
4 !pip install youtube-transcript-api
5
6 from youtube_transcript_api import YouTubeTranscriptApi
7
8 def get_transcript(video_url):
9     """
10     Extracts the transcript from a YouTube video URL.
11
12     Args:
13         video_url: The URL of the YouTube video.
14
15     Returns:
16         A list of dictionaries, where each dictionary represents a segment of the transcript
17         with 'text' and 'start' keys, or None if an error occurs.
18     """
19
20     try:
21         # Extract video ID from URL
22         video_id = video_url.split("v=")[-1]
23         if "&" in video_id:
24             video_id = video_id[:video_id.index("&")]
25
26         # Get transcript
27         transcript = YouTubeTranscriptApi.get_transcript(video_id)
28         return transcript
29
30     except Exception as e:
31         print(f"Error: {e}")
32         return None
33
34 # Example usage
35 video_url = input("Enter the YouTube video URL: ")
36 transcript = get_transcript(video_url)
37
38 if transcript:
39     print("\nTranscript:")
40     for segment in transcript:
41         print(f"{segment['text']}")

```

1 #This is the final code

1 !pip install youtube_transcript_api

```

Collecting youtube_transcript_api
  Downloading youtube_transcript_api-0.6.3-py3-none-any.whl.metadata (17 kB)
Requirement already satisfied: defusedxml<0.8.0,>=0.7.1 in /usr/local/lib/python3.11/dist-packages (from youtube_transcript_api) (0.7.1)
Requirement already satisfied: requests in /usr/local/lib/python3.11/dist-packages (from youtube_transcript_api) (2.32.3)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests->youtube_transcript_api) (3.3.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests->youtube_transcript_api) (3.10.1)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests->youtube_transcript_api) (2.2.3)
Requirement already satisfied: certifi<2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests->youtube_transcript_api) (2024.7.4)
  Downloading youtube_transcript_api-0.6.3-py3-none-any.whl (622 kB)
    622.3/622.3 kB 12.8 MB/s eta 0:00:00
Installing collected packages: youtube_transcript_api
Successfully installed youtube_transcript_api-0.6.3

```

1 pip install pytube speechrecognition pydub

```

Collecting pytube
  Downloading pytube-15.0.0-py3-none-any.whl.metadata (5.0 kB)
Collecting speechrecognition
  Downloading SpeechRecognition-3.14.1-py3-none-any.whl.metadata (31 kB)
Collecting pydub
  Downloading pydub-0.25.1-py2.py3-none-any.whl.metadata (1.4 kB)
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.11/dist-packages (from speechrecognition) (4.12.2)
  Downloading pytube-15.0.0-py3-none-any.whl (57 kB)
    57.6/57.6 kB 3.2 MB/s eta 0:00:00
  Downloading SpeechRecognition-3.14.1-py3-none-any.whl (32.9 MB)
    32.9/32.9 MB 27.2 MB/s eta 0:00:00
  Downloading pydub-0.25.1-py2.py3-none-any.whl (32 kB)
Installing collected packages: pydub, speechrecognition, pytube
Successfully installed pydub-0.25.1 pytube-15.0.0 speechrecognition-3.14.1

```

```

1 import re
2 import urllib.parse

```

```

3 import requests
4 from youtube_transcript_api import YouTubeTranscriptApi
5 from pytube import YouTube
6 import speech_recognition as sr
7 from pydub import AudioSegment
8 import os
9
10 def extract_video_id(video_url):
11     """
12     Extracts the YouTube video ID from various URL formats.
13     """
14     parsed_url = urllib.parse.urlparse(video_url)
15     query_params = urllib.parse.parse_qs(parsed_url.query)
16
17     if "v" in query_params:
18         return query_params["v"][0]
19
20     match = re.search(r"(youtu\.be/|youtube\.com/embed/|youtube\.com/shorts/)([w-]+)", video_url)
21     if match:
22         return match.group(2)
23
24     return None
25
26 def download_audio(video_url):
27     """
28     Downloads the audio using yt-dlp with cookies and returns the file path.
29     """
30     try:
31         ydl_opts = {
32             'format': 'bestaudio/best',
33             'outtmpl': 'audio.%(ext)s',
34             'cookiefile': 'cookies (1).txt', # Use the exported cookies
35             'postprocessors': [{
36                 'key': 'FFmpegExtractAudio',
37                 'preferredcodec': 'mp3',
38                 'preferredquality': '192',
39             }],
40         }
41         with yt_dlp.YoutubeDL(ydl_opts) as ydl:
42             info = ydl.extract_info(video_url, download=True)
43             return "audio.mp3"
44     except Exception as e:
45         return f"Error downloading audio: {str(e)}"
46
47 def convert_audio_to_wav(audio_file):
48     """
49     Converts the downloaded MP3 audio to WAV format using pydub.
50     """
51     wav_file = "audio.wav"
52     try:
53         AudioSegment.from_mp3(audio_file).export(wav_file, format="wav")
54         return wav_file
55     except Exception as e:
56         return f"Error converting to WAV: {str(e)}"
57
58 def transcribe_audio(audio_path, chunk_length=30):
59     """
60     Splits audio into smaller chunks and transcribes each chunk separately.
61     Args:
62         audio_path (str): Path to the audio file.
63         chunk_length (int): Length of each chunk in seconds (default: 30).
64     Returns:
65         str: Transcribed text from the audio.
66     """
67     recognizer = sr.Recognizer()
68     audio = AudioSegment.from_wav(audio_path)
69     total_duration = len(audio) / 1000 # Convert to seconds
70     transcribed_text = []
71
72     print("Transcribing audio in chunks...")
73
74     # Split and transcribe audio in chunks
75     for start in range(0, int(total_duration), chunk_length):
76         end = min(start + chunk_length, int(total_duration))
77         chunk = audio[start * 1000:end * 1000] # Extract chunk in milliseconds
78         chunk.export("chunk.wav", format="wav") # Save chunk temporarily
79
80         with sr.AudioFile("chunk.wav") as source:
81             try:
82                 audio_data = recognizer.record(source)
83                 text = recognizer.recognize_google(audio_data)
84                 transcribed_text.append(text)

```

```

85         except sr.UnknownValueError:
86             transcribed_text.append("[Unintelligible]")
87         except sr.RequestError as e:
88             return f"Error with the speech recognition service: {str(e)}"
89
90     os.remove("chunk.wav") # Clean up temporary chunk file
91     return "\n".join(transcribed_text)
92
93 def get_transcript_unlisted(video_url):
94     """
95     Tries to fetch the transcript using youtube_transcript_api first,
96     then falls back to downloading and transcribing audio if necessary.
97     """
98     video_id = extract_video_id(video_url)
99     if not video_id:
100         return "Invalid YouTube URL."
101
102     # Try to fetch transcript using youtube_transcript_api
103     try:
104         transcript = YouTubeTranscriptApi.get_transcript(video_id)
105         return "\n".join([item['text'] for item in transcript])
106     except:
107         print("Transcript not available via API, attempting audio transcription...")
108
109     # Download and transcribe audio if no transcript is available
110     audio_file = download_audio(video_url)
111     if "Error" in audio_file:
112         return audio_file
113
114     wav_file = convert_audio_to_wav(audio_file)
115     if "Error" in wav_file:
116         return wav_file
117
118     transcription = transcribe_audio(wav_file)
119
120     # Cleanup temporary files
121     os.remove(audio_file)
122     os.remove(wav_file)
123
124     return transcription
125
126 # Example usage
127 if __name__ == "__main__":
128     video_url = input("Enter the YouTube video URL: ")
129     transcript = get_transcript_unlisted(video_url)
130     print("\nTranscript:\n", transcript)
131

```

Enter the YouTube video URL: <https://youtu.be/sK8SILOM37I>
 Transcript not available via API, attempting audio transcription...
 [youtube] Extracting URL: <https://youtu.be/sK8SILOM37I>
 [youtube] sK8SILOM37I: Downloading webpage
 [youtube] sK8SILOM37I: Downloading tv client config
 [youtube] sK8SILOM37I: Downloading player b191cf34
 [youtube] sK8SILOM37I: Downloading tv player API JSON
 [info] sK8SILOM37I: Downloading 1 format(s): 251
 [download] Destination: audio.webm
 [download] 100% of 39.37MiB in 00:00:05 at 6.69MiB/s
 [ExtractAudio] Destination: audio.mp3
 Deleting original file audio.webm (pass -k to keep)
 Transcribing audio in chunks...

Transcript:

so sorry we know that India has seen a huge Revolution with digital payments we all thought that India is a place at least the I think immediately after you play the next big revolution personally I think is an education and the complete homework for this do you think ncrf plus any people together will be the next big revolution after up in India absolutely and why I think so is because how many changes have happened in the real world so many changes have happened in the requirement of the industry requirement of July 2020 and we recently celebrated the 4th anniversary of NP 2020 ncrf has been brought to implement the intent of 2020 2020 this allows you all the Innovation the way you educate your kids you educate your students still it provides you the basic guide broken the shackles of which were there in the education sector right so yes it is a big Revolution and this is going to change I pull up my kitchen to add to my taste buds in a way that it's convenient for me and one fine day you came and you changed my kitchen we have been driving buses at Max now we should fly how do we do this look at the requirement of the industry requirement of the link that whatever he has been taught has no relevance to the real life world when industry is moving that fast when they require will you be not only beneficial but also very facilitated for all of us very liberating for all of us this is going to be highly new things create new ways of doing things learn something new but once we learn it there is no limit to Innovation and creativity creating our vision and Innovative Minds into the education sector where is we are applying it elsewhere everywhere else no we are to be fun for everyone and I can I can tell you that already a number of Institutions have adopted the any pain and CRF to vary the next step on how to implement an AP if I can request you to give me an elevator pitch for an EP and then an elevator pitch for next it allows for creditor of all learnings whether in academics or an experiential learning and all these three kinds and people who are already skilled or already in the professional area not there the experiential learning would play a big part increase of technology which has been created by single Department good question so therefore all of us we have to work in tandem all kinds of learning are being contractors including learning of soft skills employability skills life skills your hand skills go out it's all very flexible so that there's no Dropout there's no Dropout so these three things coupled with use of technology I want my son to be an engineer don't you think if we create a give me five approach to Credit Systems everybody will come to the turn off infrastructure India has today we created in last 1775 years we are going to double that infrastructure in next 10:11 years more number of other branches even liberal arts social sciences if I want to be an award Society I would need a proper mix of all

have you already seen the Fallout of this know you can see the photo you can see how many Engineers are there for Designing our l
[Unintelligible]
any of the new technology machines current any laser-based machines any automated operating machines robotic operations
robotic process and that is killing is equally important and this is important in multiple areas and therefore multidisciplinary
me and you want to know design a VTech in CSC syllabus or be taking AI syllabus that is Nip complaint under the ncrf framework h
are you teaching Teddy teaching Terry is not sufficient if you want the student to really understand and reply that concept pract
those horses are skill bass courses you divide every subject into 30 and its application how do you apply the theory and those ar
who is learning which is happening which which you are going through so then this looks like let's say my student stays for 1 ye
extension of BSC physics or BSC chemistry and you give him the actual knowledge of computer science right in the first year so th
that's unbelievable so that 50% of the time which means two full years and a btech program a person can stay outside the compass
are you learning outcomes and their alignment with the overall curricular structure and then once it comes back we have to test i
he gets the credits okay so here is where I have talked with some inhibitions about the entire setup where you are simply assumir
write although we can keep that a check it is not easy for us to keep some zones green some zones red operational cost for that v
what is the guarantee that the student is learning in the campus is there a is there an accident on some kind of a giant assessme
either online or with some time stamps it is it is being documented know what time is done and CVT has videographed every assessr
she claims that all right he is a good technician and he can repair any kind of car so she prepares you open the BMW engine and t
Julie appointed by the awarding body which Awards the certificate and then that video is kept forever so you imagine won't you fe
who is giving us I think that's a nail on the head where I think we all should pass for a moment and then think is our education
experimenting something like this I don't think we'll be damaging I think all the Institute should come out of their close-minde

```
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 4.3 MB/s eta 0:00:00
  Downloading yt_dlp-2025.2.19-py3-none-any.whl (3.2 MB)
    3.2/3.2 MB 32.6 MB/s eta 0:00:00

Installing collected packages: yt-dlp
Successfully installed yt-dlp-2025.2.19
```

```
1 import yt_dlp
```

```
1 from pydub import AudioSegment
2 import os
```


```
1 import re
2 from collections import Counter
3 import nltk
4 from nltk.tokenize import sent_tokenize
5 from sklearn.cluster import KMeans
6 from sklearn.feature_extraction.text import TfidfVectorizer
7
8 def get_transcript():
9     """Prompt user to input the transcript text."""
10    print("Enter the transcript (end input with 'END'):")
11    lines = []
12    while True:
13        line = input()
14        if line.strip().upper() == "END":
15            break
16        lines.append(line)
17    return "\n".join(lines)
18
19 def extract_topics(transcript, num_topics=5):
20     """Extract topics from the transcript using TF-IDF and clustering."""
21     sentences = sent_tokenize(transcript)
22     vectorizer = TfidfVectorizer(stop_words='english')
23     X = vectorizer.fit_transform(sentences)
24
25     kmeans = KMeans(n_clusters=num_topics, random_state=42, n_init=10)
26     kmeans.fit(X)
27
28     topic_sentences = {i: [] for i in range(num_topics)}
29     for i, label in enumerate(kmeans.labels_):
30         topic_sentences[label].append(sentences[i])
31
32     topics = {f"Topic {i+1}": " ".join(topic_sentences[i][:3]) for i in range(num_topics)}
33     return topics
34
35 def segment_transcript(transcript):
36     """Segment the transcript based on extracted topics."""
37     topics = extract_topics(transcript)
38
39     segmented_text = ""
40     for heading, example_text in topics.items():
41         segmented_text += f"\n### {heading}\n- {example_text}...\n"
42
43     return segmented_text
44
45 def main():
46     nltk.download('punkt') # Ensure required tokenizer is available
```



```

47 transcript = get_transcript()
48 segmented_text = segment_transcript(transcript)
49
50 output_file = input("Enter the desired output file name (or press Enter for default: segmented_transcript.txt): ")
51 output_file = output_file if output_file else "segmented_transcript.txt"
52
53 with open(output_file, "w", encoding="utf-8") as file:
54     file.write(segmented_text)
55
56 print(f"Segmentation complete! Output saved to {output_file}")
57
58 if __name__ == "__main__":
59     main()
60

```

 [nltk_data] Downloading package punkt to /root/nltk_data...
 [nltk_data] Unzipping tokenizers/punkt.zip.
 Enter the transcript (end input with 'END'):
 so sorry we know that India has seen a huge Revolution with digital payments we all thought that India is a place at least the best
 END

```

-----
LookupError                                Traceback (most recent call last)
<ipython-input-1-47d2f9668108> in <cell line: 0>()
    57
    58 if __name__ == "__main__":
--> 59     main()

```

```

----- 7 frames -----
/usr/local/lib/python3.11/dist-packages/nltk/data.py in find(resource_name, paths)
    577     sep = "*" * 70
    578     resource_not_found = f"\n{sep}\n{msg}\n{sep}\n"
--> 579     raise LookupError(resource_not_found)
    580
    581

```

```

LookupError:
*****
Resource punkt_tab not found.
Please use the NLTK Downloader to obtain the resource:

```

```

>>> import nltk
>>> nltk.download('punkt_tab')

```

For more information see: <https://www.nltk.org/data.html>

Attempted to load tokenizers/punkt_tab/english/

```

Searched in:
- '/root/nltk_data'
- '/usr/nltk_data'
- '/usr/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/share/nltk_data'
- '/usr/local/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/local/lib/nltk_data'

```

```

*****

```

```

1 import re
2 import nltk
3 from nltk.tokenize import sent_tokenize
4 from sklearn.cluster import KMeans
5 from sklearn.feature_extraction.text import TfidfVectorizer
6
7 def read_transcript(file_path):
8     """Read the transcript from a file."""
9     with open(file_path, "r", encoding="utf-8") as file:
10         return file.read()
11
12 def extract_topics(transcript, num_topics=5):
13     """Extract topics from the transcript using TF-IDF and clustering."""
14     sentences = sent_tokenize(transcript)
15     vectorizer = TfidfVectorizer(stop_words='english')
16     X = vectorizer.fit_transform(sentences)
17
18     kmeans = KMeans(n_clusters=num_topics, random_state=42, n_init=10)
19     kmeans.fit(X)
20
21     topic_sentences = {i: [] for i in range(num_topics)}
22     for i, label in enumerate(kmeans.labels_):
23         topic_sentences[label].append(sentences[i])
24
25     topics = {f"Topic {i+1}": " ".join(topic_sentences[i][:3]) for i in range(num_topics)}

```

```

26     return topics
27
28 def segment_transcript(transcript):
29     """Segment the transcript based on extracted topics."""
30     topics = extract_topics(transcript)
31
32     segmented_text = ""
33     for heading, example_text in topics.items():
34         segmented_text += f"\n### {heading}\n- {example_text}...\n"
35
36     return segmented_text
37
38 def main():
39     nltk.download('punkt') # Ensure required tokenizer is available
40     input_file = input("Enter the path of the transcript file: ")
41     transcript = read_transcript(input_file)
42     segmented_text = segment_transcript(transcript)
43
44     output_file = input("Enter the desired output file name (or press Enter for default: segmented_transcript.txt): ")
45     output_file = output_file if output_file else "segmented_transcript.txt"
46
47     with open(output_file, "w", encoding="utf-8") as file:
48         file.write(segmented_text)
49
50     print(f"Segmentation complete! Output saved to {output_file}")
51
52 if __name__ == "__main__":
53     main()
54

```

```

[ntlk_data] Downloading package punkt to /root/nltk_data...
[ntlk_data] Package punkt is already up-to-date!
Enter the path of the transcript file: /content/Kalsi sir transcript.txt
-----
LookupError                                Traceback (most recent call last)
<ipython-input-2-ec884cafd024> in <cell line: 0>()
    51
    52 if __name__ == "__main__":
--> 53     main()

```

7 frames

```

/usr/local/lib/python3.11/dist-packages/nltk/data.py in find(resource_name, paths)
    577     sep = "*" * 70
    578     resource_not_found = f"\n{sep}\n[msg]\n{sep}\n"
--> 579     raise LookupError(resource_not_found)
    580
    581

```

LookupError:

Resource `punkt_tab` not found.
Please use the NLTK Downloader to obtain the resource:

```
>>> import nltk
>>> nltk.download('punkt_tab')
```

For more information see: <https://www.nltk.org/data.html>

Attempted to load `tokenizers/punkt_tab/english/`

Searched in:

- '/root/nltk_data'
- '/usr/nltk_data'
- '/usr/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/share/nltk_data'
- '/usr/local/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/local/lib/nltk_data'

```

1 import re
2 import nltk
3 from nltk.tokenize import sent_tokenize
4 from sklearn.cluster import KMeans
5 from sklearn.feature_extraction.text import TfidfVectorizer
6
7 def read_transcript(file_path):
8     """Read the transcript from a file."""
9     with open(file_path, "r", encoding="utf-8") as file:
10         return file.read()
11
12 def extract_topics(transcript, num_topics=5):
13     """Extract topics from the transcript using TF-IDF and clustering."""
14     nltk.download('punkt') # Ensure required tokenizer is available
15     sentences = sent_tokenize(transcript)
16     vectorizer = TfidfVectorizer(stop_words='english')
17     X = vectorizer.fit_transform(sentences)
18
19     kmeans = KMeans(n_clusters=num_topics, random_state=42, n_init=10)
20     kmeans.fit(X)
21
22     topic_sentences = {i: [] for i in range(num_topics)}
23     for i, label in enumerate(kmeans.labels_):
24         topic_sentences[label].append(sentences[i])
25
26     topics = {f"Topic {i+1}": " ".join(topic_sentences[i][:3]) for i in range(num_topics)}
27     return topics
28
29 def segment_transcript(transcript):
30     """Segment the transcript based on extracted topics."""
31     topics = extract_topics(transcript)
32
33     segmented_text = ""
34     for heading, example_text in topics.items():
35         segmented_text += f"\n### {heading}\n- {example_text}...\n"
36
37     return segmented_text
38
39 def main():
40     nltk.download('punkt') # Ensure required tokenizer is available
41     input_file = input("Enter the path of the transcript file: ")
42     transcript = read_transcript(input_file)
43     segmented_text = segment_transcript(transcript)
44
45     output_file = input("Enter the desired output file name (or press Enter for default: segmented_transcript.txt): ")
46     output_file = output_file if output_file else "segmented_transcript.txt"
47
48     with open(output_file, "w", encoding="utf-8") as file:
49         file.write(segmented_text)
50
51     print(f"Segmentation complete! Output saved to {output_file}")
52
53 if __name__ == "__main__":
54     main()
55

```

```

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
Enter the path of the transcript file: /content/Kalsi sir transcript.txt
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
-----
LookupError                                Traceback (most recent call last)
<ipython-input-5-4884c9bae683> in <cell line: 0>()
     52
     53 if __name__ == "__main__":
--> 54     main()

-----
7 frames
/usr/local/lib/python3.11/dist-packages/nltk/data.py in find(resource_name, paths)
     577     sep = "*" * 70
     578     resource_not_found = f"\n{sep}\n{msg}\n{sep}\n"
--> 579     raise LookupError(resource_not_found)
     580
     581

LookupError:
*****
Resource punkt_tab not found.
Please use the NLTK Downloader to obtain the resource:

>>> import nltk
>>> nltk.download('punkt_tab')

For more information see: https://www.nltk.org/data.html

Attempted to load tokenizers/punkt_tab/english/

Searched in:
- '/root/nltk_data'
- '/usr/nltk_data'
- '/usr/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/share/nltk_data'
- '/usr/local/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/local/lib/nltk_data'
*****

```

```

1 import nltk
2 nltk.download('punkt')

```

```

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

```

```

1 import nltk
2 import os
3 import shutil
4
5 # Remove existing nltk_data directory
6 nltk_data_path = os.path.expanduser('~/.nltk_data')
7 if os.path.exists(nltk_data_path):
8     shutil.rmtree(nltk_data_path)
9
10 # Download necessary resources again
11 nltk.download('punkt')

```

```

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

```

```

1 import re
2 import nltk
3 from nltk.tokenize import sent_tokenize
4 from sklearn.cluster import KMeans
5 from sklearn.feature_extraction.text import TfidfVectorizer
6
7 def read_transcript(file_path):
8     """Read the transcript from a file."""
9     with open(file_path, "r", encoding="utf-8") as file:
10         return file.read()
11
12 def extract_topics(transcript, num_topics=5):
13     """Extract topics from the transcript using TF-IDF and clustering."""
14     nltk.download('punkt') # Ensure required tokenizer is available
15     sentences = sent_tokenize(transcript)
16     vectorizer = TfidfVectorizer(stop_words='english')
17     X = vectorizer.fit_transform(sentences)

```

```

18
19 kmeans = KMeans(n_clusters=num_topics, random_state=42, n_init=10)
20 kmeans.fit(X)
21
22 topic_sentences = {i: [] for i in range(num_topics)}
23 for i, label in enumerate(kmeans.labels_):
24     topic_sentences[label].append(sentences[i])
25
26 topics = {f"Topic {i+1}": " ".join(topic_sentences[i][:3]) for i in range(num_topics)}
27 return topics
28
29 def segment_transcript(transcript):
30     """Segment the transcript based on extracted topics."""
31     topics = extract_topics(transcript)
32
33     segmented_text = ""
34     for heading, example_text in topics.items():
35         segmented_text += f"\n### {heading}\n- {example_text}...\n"
36
37     return segmented_text
38
39 def main():
40     nltk.download('punkt') # Ensure required tokenizer is available
41     input_file = input("Enter the path of the transcript file: ")
42     transcript = read_transcript(input_file)
43     segmented_text = segment_transcript(transcript)
44
45     output_file = input("Enter the desired output file name (or press Enter for default: segmented_transcript.txt): ")
46     output_file = output_file if output_file else "segmented_transcript.txt"
47
48     with open(output_file, "w", encoding="utf-8") as file:
49         file.write(segmented_text)
50
51     print(f"Segmentation complete! Output saved to {output_file}")
52
53 if __name__ == "__main__":
54     main()

```

```

[ nltk_data ] Downloading package punkt to /root/nltk_data...
[ nltk_data ] Package punkt is already up-to-date!
Enter the path of the transcript file: /content/Kalsi sir transcript.txt
[ nltk_data ] Downloading package punkt to /root/nltk_data...
[ nltk_data ] Package punkt is already up-to-date!

```

```

-----
LookupError                                Traceback (most recent call last)
<ipython-input-7-4884c9bae683> in <cell line: 0>()
    52
    53 if __name__ == "__main__":
--> 54     main()

```

```

----- 7 frames -----
/usr/local/lib/python3.11/dist-packages/nltk/data.py in find(resource_name, paths)
    577     sep = "*" * 70
    578     resource_not_found = f"\n{sep}\n{msg}\n{sep}\n"
--> 579     raise LookupError(resource_not_found)
    580
    581

```

```

LookupError:
*****

```

Resource `punkt_tab` not found.
Please use the NLTK Downloader to obtain the resource:

```

>>> import nltk
>>> nltk.download('punkt_tab')

```

For more information see: <https://www.nltk.org/data.html>

Attempted to load `tokenizers/punkt_tab/english/`

```

Searched in:
- '/root/nltk_data'
- '/usr/nltk_data'
- '/usr/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/share/nltk_data'
- '/usr/local/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/local/lib/nltk_data'

```

```

*****

```

```

1 import nltk
2 nltk.download('punkt')

```

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

```
1 import nltk
2 print(nltk.data.path)
```

```
[ '/root/nltk_data', '/usr/nltk_data', '/usr/share/nltk_data', '/usr/lib/nltk_data', '/usr/share/nltk_data', '/usr/local/share/nltk_data', '/usr/local/lib/nltk_data' ]
```

```
1 nltk.data.path.append('/usr/local/share/nltk_data')
```

```
1 import shutil
2 import os
3 import nltk
4
5 nltk_data_path = os.path.expanduser('~/.nltk_data')
6 if os.path.exists(nltk_data_path):
7     shutil.rmtree(nltk_data_path)
8
9 nltk.download('punkt')
```

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

```
1 import os
2 print(os.listdir(nltk.data.find("tokenizers")))
```

```
-----
IndexError                                Traceback (most recent call last)
<ipython-input-12-82277fac31aa> in <cell line: 0>()
      1 import os
----> 2 print(os.listdir(nltk.data.find("tokenizers")))

/usr/local/lib/python3.11/dist-packages/nltk/data.py in find(resource_name, paths)
    554
    555     # Identify the package (i.e. the .zip file) to download.
--> 556     resource_zipname = resource_name.split("/")[1]
    557     if resource_zipname.endswith(".zip"):
    558         resource_zipname = resource_zipname.rpartition(".")[0]

IndexError: list index out of range
```

```
1 import nltk
2 print(nltk.data.path)
```

```
[ '/root/nltk_data', '/usr/nltk_data', '/usr/share/nltk_data', '/usr/lib/nltk_data', '/usr/share/nltk_data', '/usr/local/share/nltk_data', '/usr/local/lib/nltk_data' ]
```

```
1 import os
2 nltk_data_path = os.path.expanduser("~/.nltk_data/tokenizers")
3 if os.path.exists(nltk_data_path):
4     print(os.listdir(nltk_data_path))
5 else:
6     print("NLTK tokenizer directory not found!")
```

```
NLTK tokenizer directory not found!
```

```
1 import os
2
3 nltk_data_path = os.path.expanduser("~/.nltk_data/tokenizers")
4 if not os.path.exists(nltk_data_path):
5     os.makedirs(nltk_data_path)
6
7 print(f"Created directory: {nltk_data_path}")
```

```
Created directory: /root/nltk_data/tokenizers
```

```
1 import nltk
2 nltk.download('punkt')
```

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

```

1 import os
2
3 nltk_data_path = os.path.expanduser("~/nltk_data/tokenizers")
4 if os.path.exists(nltk_data_path):
5     print(os.listdir(nltk_data_path))
6 else:
7     print("NLTK tokenizer directory still missing!")

```

↗ []

```

1 import nltk
2 import os
3
4 # Set a specific NLTK data directory
5 nltk_data_dir = os.path.expanduser("~/nltk_data")
6 if not os.path.exists(nltk_data_dir):
7     os.makedirs(nltk_data_dir)
8
9 # Append the path manually in case it's not detected
10 nltk.data.path.append(nltk_data_dir)
11
12 # Download the 'punkt' tokenizer
13 nltk.download('punkt', download_dir=nltk_data_dir)
14
15 # Verify that the tokenizer exists now
16 tokenizer_path = os.path.join(nltk_data_dir, "tokenizers")
17 print(f"Checking for tokenizers in: {tokenizer_path}")
18 print(os.listdir(tokenizer_path) if os.path.exists(tokenizer_path) else "Tokenizer directory not found!")

```

↗ Checking for tokenizers in: /root/nltk_data/tokenizers
 []
 [nltk_data] Downloading package punkt to /root/nltk_data...
 [nltk_data] Package punkt is already up-to-date!

```

1 import nltk
2 from nltk.data import find
3
4 try:
5     punkt_path = find('tokenizers/punkt')
6     print(f"'punkt' found at: {punkt_path}")
7 except LookupError:
8     print("Could not find 'punkt' in nltk data.")

```

↗ Could not find 'punkt' in nltk data.

```

1 import shutil
2 import os
3
4 nltk_data_path = os.path.expanduser("~/nltk_data")
5
6 if os.path.exists(nltk_data_path):
7     shutil.rmtree(nltk_data_path)
8
9 print("Deleted existing NLTK data. Ready for fresh installation.")

```

↗ Deleted existing NLTK data. Ready for fresh installation.

```

1 import nltk
2 import os
3
4 nltk_data_path = os.path.expanduser("~/nltk_data")
5 os.makedirs(nltk_data_path, exist_ok=True)
6
7 nltk.download('punkt', download_dir=nltk_data_path)
8
9 print("Downloaded 'punkt' to:", nltk_data_path)

```

↗ Downloaded 'punkt' to: /root/nltk_data
 [nltk_data] Downloading package punkt to /root/nltk_data...
 [nltk_data] Package punkt is already up-to-date!

```

1 import os
2
3 tokenizer_path = os.path.expanduser("~/nltk_data/tokenizers")
4 print("Checking for tokenizers in:", tokenizer_path)
5 print(os.listdir(tokenizer_path) if os.path.exists(tokenizer_path) else "Tokenizer directory not found!")

```

↗ Checking for tokenizers in: /root/nltk_data/tokenizers
 Tokenizer directory not found!


```

1 # prompt: segment the transcript from the path with relevant topics
2
3 import nltk
4 from nltk.tokenize import sent_tokenize
5 from sklearn.feature_extraction.text import TfidfVectorizer
6 from sklearn.cluster import KMeans
7 import os
8
9 def read_transcript(file_path):
10     """Read the transcript from a file."""
11     with open(file_path, "r", encoding="utf-8") as file:
12         return file.read()
13
14 def extract_topics(transcript, num_topics=5):
15     """Extract topics from the transcript using TF-IDF and clustering."""
16     try:
17         sentences = sent_tokenize(transcript)
18         vectorizer = TfidfVectorizer(stop_words='english')
19         X = vectorizer.fit_transform(sentences)
20
21         kmeans = KMeans(n_clusters=num_topics, random_state=42, n_init=10)
22         kmeans.fit(X)
23
24         topic_sentences = {i: [] for i in range(num_topics)}
25         for i, label in enumerate(kmeans.labels_):
26             topic_sentences[label].append(sentences[i])
27
28         topics = {f"Topic {i+1}": " ".join(topic_sentences[i][:3]) for i in range(num_topics)}
29         return topics
30     except Exception as e:
31         print(f"An error occurred during topic extraction: {e}")
32         return None
33
34 def segment_transcript(transcript):
35     """Segment the transcript based on extracted topics."""
36     topics = extract_topics(transcript)
37     if topics is None:
38         return "Topic extraction failed."
39
40     segmented_text = ""
41     for heading, example_text in topics.items():
42         segmented_text += f"\n### {heading}\n- {example_text}...\n"
43
44     return segmented_text
45
46 def main():
47     nltk.download('punkt', quiet=True) # Download punkt if not present, suppress output
48     input_file = input("Enter the path of the transcript file: ")
49
50     if not os.path.exists(input_file):
51         print(f"Error: File '{input_file}' not found.")
52         return
53
54     transcript = read_transcript(input_file)
55     segmented_text = segment_transcript(transcript)
56
57     output_file = input("Enter the desired output file name (or press Enter for default: segmented_transcript.txt): ")
58     output_file = output_file if output_file else "segmented_transcript.txt"
59
60     with open(output_file, "w", encoding="utf-8") as file:
61         file.write(segmented_text)
62
63     print(f"Segmentation complete! Output saved to {output_file}")
64
65 if __name__ == "__main__":
66     main()
67

```

Enter the path of the transcript file: /content/Kalsi sir transcript.txt

An error occurred during topic extraction:

Resource `punkt_tab` not found.

Please use the NLTK Downloader to obtain the resource:

```
>>> import nltk
>>> nltk.download('punkt_tab')
```

For more information see: <https://www.nltk.org/data.html>

Attempted to load `tokenizers/punkt_tab/english/`

Searched in:

- '/root/nltk_data'
- '/usr/nltk_data'

```

- '/usr/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/share/nltk_data'
- '/usr/local/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/local/lib/nltk_data'
- '/usr/local/share/nltk_data'
- '/root/nltk_data'

```

Enter the desired output file name (or press Enter for default: segmented_transcript.txt):
Segmentation complete! Output saved to segmented_transcript.txt

1 #Another approach

1 pip install transformers nltk

```

Requirement already satisfied: transformers in /usr/local/lib/python3.11/dist-packages (4.48.3)
Requirement already satisfied: nltk in /usr/local/lib/python3.11/dist-packages (3.9.1)
Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from transformers) (3.17.0)
Requirement already satisfied: huggingface-hub<1.0,>=0.24.0 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.28.1)
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.11/dist-packages (from transformers) (1.26.4)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from transformers) (24.2)
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.11/dist-packages (from transformers) (6.0.2)
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.11/dist-packages (from transformers) (2024.11.6)
Requirement already satisfied: requests in /usr/local/lib/python3.11/dist-packages (from transformers) (2.32.3)
Requirement already satisfied: tokenizers<0.22,>=0.21 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.21.0)
Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.5.3)
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.11/dist-packages (from transformers) (4.67.1)
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: fsspec>=2023.5.0 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.24.0->transf
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.24.0->transf
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (3.10)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (2.3.0)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (2025.1.1)

```

```


1 import os
2 import nltk
3 from transformers import pipeline
4 from nltk.tokenize import sent_tokenize
5 import torch
6
7 # Download the Punkt tokenizer for sentence splitting if not already downloaded
8 nltk.download('punkt')
9
10 # Define file paths
11 input_file = '/content/Kalsi sir transcript.txt' # Replace with your transcript file name
12 output_file = 'segmented_transcript.txt'
13
14 # Load the summarization model for heading generation
15 summarizer = pipeline("summarization", model="facebook/bart-large-cnn", device=-1 if torch.cuda.is_available() else -1)
16
17 # Read the transcript
18 with open(input_file, 'r', encoding='utf-8') as file:
19     transcript = file.read()
20
21 # Split transcript into sentences
22 sentences = sent_tokenize(transcript)
23
24 # Segment transcript into chunks of approximately 500 words for manageable heading generation
25 chunk_size = 300
26 chunks = []
27 current_chunk = []
28
29 for sentence in sentences:
30     current_chunk.append(sentence)
31     if sum(len(s.split()) for s in current_chunk) >= chunk_size:
32         chunks.append(' '.join(current_chunk))
33         current_chunk = []
34
35 if current_chunk:
36     chunks.append(' '.join(current_chunk))
37
38 # Generate headings and organize content
39 segmented_content = ""
40 for idx, chunk in enumerate(chunks, 1):
41     if not chunk.strip():
42         continue # Skip empty chunks
43
44     # Limit chunk to 1024 tokens if it's too large

```

```

45     tokens = tokenizer(chunk)['input_ids']
46     if len(tokens) > 1024:
47         chunk = tokenizer.decode(tokens[:1024])
48
49     # Generate a heading using summarization
50     heading = summarizer(chunk, max_length=10, min_length=3, do_sample=False)[0]['summary_text']
51     segmented_content += f"### {heading}\n\n{chunk}\n\n"
52
53 # Save segmented content to a new file
54 with open(output_file, 'w', encoding='utf-8') as file:
55     file.write(segmented_content)
56
57 print(f"Segmented transcript saved to {output_file}")
58

```

 [nltk_data] Downloading package punkt to /root/nltk_data...
 [nltk_data] Package punkt is already up-to-date!
 Device set to use cpu
 Token indices sequence length is longer than the specified maximum sequence length for this model (3126 > 512). Running this sequen
 Segmented transcript saved to segmented_transcript.txt

```

1 import nltk
2 nltk.download('all')

```

 [nltk_data] Downloading collection 'all'
 [nltk_data] |
 [nltk_data] | Downloading package abc to /root/nltk_data...
 [nltk_data] | Unzipping corpora/abc.zip.
 [nltk_data] | Downloading package alpino to /root/nltk_data...
 [nltk_data] | Unzipping corpora/alpino.zip.
 [nltk_data] | Downloading package averaged_perceptron_tagger to
 [nltk_data] | /root/nltk_data...
 [nltk_data] | Unzipping taggers/averaged_perceptron_tagger.zip.
 [nltk_data] | Downloading package averaged_perceptron_tagger_eng to
 [nltk_data] | /root/nltk_data...
 [nltk_data] | Unzipping
 [nltk_data] | taggers/averaged_perceptron_tagger_eng.zip.
 [nltk_data] | Downloading package averaged_perceptron_tagger_ru to
 [nltk_data] | /root/nltk_data...
 [nltk_data] | Unzipping
 [nltk_data] | taggers/averaged_perceptron_tagger_ru.zip.
 [nltk_data] | Downloading package averaged_perceptron_tagger_rus to
 [nltk_data] | /root/nltk_data...
 [nltk_data] | Unzipping
 [nltk_data] | taggers/averaged_perceptron_tagger_rus.zip.
 [nltk_data] | Downloading package basque_grammars to
 [nltk_data] | /root/nltk_data...
 [nltk_data] | Unzipping grammars/basque_grammars.zip.
 [nltk_data] | Downloading package bcp47 to /root/nltk_data...
 [nltk_data] | Downloading package biocreative_ppi to
 [nltk_data] | /root/nltk_data...
 [nltk_data] | Unzipping corpora/biocreative_ppi.zip.
 [nltk_data] | Downloading package bllip_wsj_no_aux to
 [nltk_data] | /root/nltk_data...
 [nltk_data] | Unzipping models/bllip_wsj_no_aux.zip.
 [nltk_data] | Downloading package book_grammars to
 [nltk_data] | /root/nltk_data...
 [nltk_data] | Unzipping grammars/book_grammars.zip.
 [nltk_data] | Downloading package brown to /root/nltk_data...
 [nltk_data] | Unzipping corpora/brown.zip.
 [nltk_data] | Downloading package brown_tei to /root/nltk_data...
 [nltk_data] | Unzipping corpora/brown_tei.zip.
 [nltk_data] | Downloading package cess_cat to /root/nltk_data...
 [nltk_data] | Unzipping corpora/cess_cat.zip.
 [nltk_data] | Downloading package cess_esp to /root/nltk_data...
 [nltk_data] | Unzipping corpora/cess_esp.zip.
 [nltk_data] | Downloading package chat80 to /root/nltk_data...
 [nltk_data] | Unzipping corpora/chat80.zip.
 [nltk_data] | Downloading package city_database to
 [nltk_data] | /root/nltk_data...
 [nltk_data] | Unzipping corpora/city_database.zip.
 [nltk_data] | Downloading package cmudict to /root/nltk_data...
 [nltk_data] | Unzipping corpora/cmudict.zip.
 [nltk_data] | Downloading package comparative_sentences to
 [nltk_data] | /root/nltk_data...
 [nltk_data] | Unzipping corpora/comparative_sentences.zip.
 [nltk_data] | Downloading package comtrans to /root/nltk_data...
 [nltk_data] | Downloading package conll2000 to /root/nltk_data...
 [nltk_data] | Unzipping corpora/conll2000.zip.
 [nltk_data] | Downloading package conll2002 to /root/nltk_data...
 [nltk_data] | Unzipping corpora/conll2002.zip.
 [nltk_data] | Downloading package conll2007 to /root/nltk_data...

```
1 from transformers import AutoTokenizer
2
3 tokenizer = AutoTokenizer.from_pretrained("t5-small")
```



tokenizer_config.json: 100%

2.32k/2.32k [00:00<00:00, 191kB/s]

spiece.model: 100%

792k/792k [00:00<00:00, 15.5MB/s]

tokenizer.json: 100%

1.39M/1.39M [00:00<00:00, 16.3MB/s]

```

1 import os
2 import re
3 import nltk
4 from transformers import pipeline, AutoTokenizer
5 from nltk.tokenize import sent_tokenize
6 import torch
7
8 # Download the Punkt tokenizer for sentence splitting if not already downloaded
9 nltk.download('punkt')
10
11 # Define file paths
12 input_file = '/content/Kalsi sir transcript.txt' # Replace with your transcript file name
13 output_file = 'segmented_transcript.txt'
14
15 # Load the summarization model and tokenizer for heading generation
16 summarizer = pipeline("summarization", model="facebook/bart-large-cnn", device=0 if torch.cuda.is_available() else -1)
17 tokenizer = AutoTokenizer.from_pretrained("facebook/bart-large-cnn", use_fast=True)
18
19 # Read and clean the transcript
20 with open(input_file, 'r', encoding='utf-8') as file:
21     transcript = file.read()
22
23 # Clean special characters and excessive spaces
24 transcript = re.sub(r'^\x00-\x7F|+', ' ', transcript) # Remove non-ASCII characters
25 transcript = re.sub(r'\s+', ' ', transcript).strip() # Replace multiple spaces with a single space
26
27 # Split transcript into paragraphs based on double newlines or fallback to sentence tokenization
28 paragraphs = transcript.split('\n\n') if '\n\n' in transcript else [transcript]
29
30 # Create clearer and more coherent chunks
31 chunk_size = 300 # Limit chunk size to ~400 tokens for clarity and coherence
32 chunks = []
33
34 for paragraph in paragraphs:
35     sentences = sent_tokenize(paragraph)
36     current_chunk = []
37
38     for sentence in sentences:
39         current_chunk.append(sentence)
40         # Check token length instead of word count for more reliable chunking
41         token_count = len(tokenizer(' '.join(current_chunk))['input_ids'])
42         if token_count >= chunk_size:
43             chunks.append(' '.join(current_chunk))
44             current_chunk = []
45
46     # Append remaining sentences in the current chunk
47     if current_chunk:
48         chunks.append(' '.join(current_chunk))
49
50 # Generate headings and organize content
51 segmented_content = ""
52 for idx, chunk in enumerate(chunks, 1):
53     if not chunk.strip() or len(chunk.split()) < 5:
54         continue # Skip empty or very short chunks
55
56     # Limit chunk to 1024 tokens if it's too large
57     tokens = tokenizer(chunk)['input_ids']
58     if len(tokens) > 1024:
59         chunk = tokenizer.decode(tokens[:1024])
60
61     # Debugging print to track token info
62     print(f"Processing chunk {idx}/{len(chunks)} with {len(tokens)} tokens.")
63
64     # Generate a heading using summarization
65     heading = summarizer(chunk, max_length=10, min_length=3, do_sample=False)[0]['summary_text']
66     segmented_content += f"### {heading.strip()}\n\n{chunk.strip()}\n\n"
67
68 # Save segmented content to a new file
69 with open(output_file, 'w', encoding='utf-8') as file:
70     file.write(segmented_content)
71
72 print(f"Segmented transcript saved to {output_file}")

```

```

[ntlk_data] Downloading package punkt to /root/ntlk_data...
[ntlk_data] Package punkt is already up-to-date!
Device set to use cpu
Processing chunk 1/2 with 2950 tokens.
-----
IndexError                                Traceback (most recent call last)
<ipython-input-37-f582a8017948> in <cell line: 0>()
     63
     64     # Generate a heading using summarization
--> 65     heading = summarizer(chunk, max_length=10, min_length=3, do_sample=False)[0]['summary_text']
     66     segmented_content += f"### {heading.strip()}\n\n{chunk.strip()}\n\n"
     67

-----
16 frames
/usr/local/lib/python3.11/dist-packages/torch/nn/functional.py in embedding(input, weight, padding_idx, max_norm, norm_type,
scale_grad_by_freq, sparse)
    2549     # remove once script supports set_grad_enabled
    2550     _no_grad_embedding_renorm_(weight, input, max_norm, norm_type)
-> 2551     return torch.embedding(weight, input, padding_idx, scale_grad_by_freq, sparse)
    2552
    2553

IndexError: index out of range in self

1 import os
2 import re
3 import nltk
4 from transformers import pipeline, AutoTokenizer
5 from nltk.tokenize import sent_tokenize
6 import torch
7
8 # Download the Punkt tokenizer for sentence splitting if not already downloaded
9 nltk.download('punkt')
10
11 # Define file paths
12 input_file = '/content/Kalsi sir transcript.txt' # Replace with your transcript file name
13 output_file = 'segmented_transcript.txt'
14
15 # Load the summarization model and tokenizer for heading generation
16 summarizer = pipeline("summarization", model="facebook/bart-large-cnn", device=0 if torch.cuda.is_available() else -1)
17 tokenizer = AutoTokenizer.from_pretrained("facebook/bart-large-cnn", use_fast=True)
18
19 # Re-download the tokenizer files to fix potential corruption
20 tokenizer.save_pretrained("./tokenizer")
21 tokenizer = AutoTokenizer.from_pretrained("./tokenizer")
22
23 # Read and clean the transcript
24 with open(input_file, 'r', encoding='utf-8') as file:
25     transcript = file.read()
26
27 # Clean special characters and excessive spaces
28 transcript = re.sub(r'[\x00-\x7F]+', ' ', transcript) # Remove non-ASCII characters
29 transcript = re.sub(r'\s+', ' ', transcript).strip() # Replace multiple spaces with a single space
30
31 # Split transcript into paragraphs based on double newlines or fallback to sentence tokenization
32 paragraphs = transcript.split('\n\n') if '\n\n' in transcript else [transcript]
33
34 # Create clearer and more coherent chunks
35 chunk_size = 400 # Limit chunk size to ~400 tokens for clarity and coherence
36 chunks = []
37
38 for paragraph in paragraphs:
39     sentences = sent_tokenize(paragraph)
40     current_chunk = []
41
42     for sentence in sentences:
43         current_chunk.append(sentence)
44         # Check token length instead of word count for more reliable chunking
45         token_count = len(tokenizer(' '.join(current_chunk))['input_ids'])
46         if token_count >= chunk_size:
47             chunks.append(' '.join(current_chunk))
48             current_chunk = []
49
50 # Append remaining sentences in the current chunk
51 if current_chunk:
52     chunks.append(' '.join(current_chunk))
53
54 # Generate headings and organize content
55 segmented_content = ""
56 for idx, chunk in enumerate(chunks, 1):
57     if not chunk.strip() or len(chunk.split()) < 5:
58         continue # Skip empty or very short chunks

```

```

59
60 # Tokenize the chunk
61 tokens = tokenizer(chunk)['input_ids']
62
63 # Split chunk into smaller parts if it exceeds 1024 tokens
64 if len(tokens) > 1024:
65     sub_chunks = [tokens[i:i + 1024] for i in range(0, len(tokens), 1024)]
66 else:
67     sub_chunks = [tokens]
68
69 # Process each sub-chunk separately
70 for sub_idx, sub_chunk in enumerate(sub_chunks, 1):
71     # Decode sub-chunk back to text
72     sub_chunk_text = tokenizer.decode(sub_chunk, skip_special_tokens=True)
73
74     # Debugging print to track token info
75     print(f"Processing chunk {idx}. {sub_idx}/{len(chunks)} with {len(sub_chunk)} tokens.")
76
77     # Generate a heading using summarization
78     try:
79         heading = summarizer(sub_chunk_text, max_length=10, min_length=3, do_sample=False)[0]['summary_text']
80         segmented_content += f"### {heading.strip()}\n\n{sub_chunk_text.strip()}\n\n"
81     except IndexError as e:
82         print(f"Skipping a chunk due to error: {e}")
83         continue
84
85 # Save segmented content to a new file
86 with open(output_file, 'w', encoding='utf-8') as file:
87     file.write(segmented_content)
88
89 print(f"Segmented transcript saved to {output_file}")

```

```

[ntlk_data] Downloading package punkt to /root/nltk_data...
[ntlk_data] Package punkt is already up-to-date!
Device set to use cpu
Processing chunk 1.1/2 with 1024 tokens.
Skipping a chunk due to error: index out of range in self
Processing chunk 1.2/2 with 1024 tokens.
Skipping a chunk due to error: index out of range in self
Processing chunk 1.3/2 with 902 tokens.
Processing chunk 2.1/2 with 1024 tokens.
Skipping a chunk due to error: index out of range in self
Processing chunk 2.2/2 with 1024 tokens.
Skipping a chunk due to error: index out of range in self
Processing chunk 2.3/2 with 1024 tokens.
Skipping a chunk due to error: index out of range in self
Processing chunk 2.4/2 with 482 tokens.
Segmented transcript saved to segmented_transcript.txt

```

```

1 import os
2 import re
3 import nltk
4 from transformers import pipeline, AutoTokenizer, AutoModel
5 from sentence_transformers import SentenceTransformer, util
6 from nltk.tokenize import sent_tokenize
7 import torch
8
9 # Download the Punkt tokenizer for sentence splitting if not already downloaded
10 nltk.download('punkt')
11
12 # Define file paths
13 input_file = '/content/Kalsi sir transcript.txt' # Replace with your transcript file name
14 output_file = 'segmented_transcript.txt'
15
16 # Load summarization and zero-shot classification models
17 summarizer = pipeline("summarization", model="facebook/bart-large-cnn", device=0 if torch.cuda.is_available() else -1)
18 classifier = pipeline("zero-shot-classification", model="facebook/bart-large-mnli")
19 tokenizer = AutoTokenizer.from_pretrained("facebook/bart-large-cnn", use_fast=True)
20
21 # Load sentence-transformers model for semantic similarity
22 embedder = SentenceTransformer('all-MiniLM-L6-v2')
23
24 # Read and clean the transcript
25 with open(input_file, 'r', encoding='utf-8') as file:
26     transcript = file.read()
27
28 # Clean special characters and excessive spaces
29 transcript = re.sub(r'[^\x00-\x7F]+', ' ', transcript) # Remove non-ASCII characters
30 transcript = re.sub(r'\s+', ' ', transcript).strip() # Replace multiple spaces with a single space
31
32 # Split transcript into sentences
33 sentences = sent_tokenize(transcript)
34

```



```

35 # Create chunks based on semantic similarity and concept shifts
36 chunk_size = 300
37 chunks = []
38 current_chunk = []
39 current_embedding = None
40
41 for sentence in sentences:
42     # Embed the current sentence
43     sentence_embedding = embedder.encode(sentence, convert_to_tensor=True)
44
45     # Check if this is the first sentence
46     if current_embedding is None:
47         current_embedding = sentence_embedding
48         current_chunk.append(sentence)
49         continue
50
51     # Calculate similarity with the last chunk's embedding
52     similarity = util.pytorch_cos_sim(current_embedding, sentence_embedding).item()
53
54     # If similarity is low or chunk size exceeds limit, start a new chunk
55     if similarity < 0.6 or len(tokenizer(' '.join(current_chunk))['input_ids']) >= chunk_size:
56         chunks.append(' '.join(current_chunk))
57         current_chunk = [sentence]
58         current_embedding = sentence_embedding
59     else:
60         current_chunk.append(sentence)
61         current_embedding = (current_embedding + sentence_embedding) / 2 # Update embedding with average
62
63 # Append remaining chunk if exists
64 if current_chunk:
65     chunks.append(' '.join(current_chunk))
66
67 # Generate headings and organize content
68 segmented_content = ""
69 for idx, chunk in enumerate(chunks, 1):
70     if not chunk.strip() or len(chunk.split()) < 5:
71         continue # Skip empty or very short chunks
72
73     # Generate a heading using zero-shot classification for topic detection
74     topics = ["Technology", "Education", "AI and Machine Learning", "Research", "Career Guidance", "Software Development", "Leadership"]
75     classification = classifier(chunk, candidate_labels=topics, multi_label=False)
76     main_topic = classification['labels'][0]
77
78     # Generate a concise heading based on main topic
79     heading = summarizer(chunk, max_length=15, min_length=5, do_sample=False)[0]['summary_text']
80     segmented_content += f"### {main_topic}: {heading.strip()}\n\n{chunk.strip()}\n\n"
81
82 # Save segmented content to a new file
83 with open(output_file, 'w', encoding='utf-8') as file:
84     file.write(segmented_content)
85
86 print(f"Segmented transcript saved to {output_file}")
87

```

```

[!nltk_data] Downloading package punkt to /root/nltk_data...
[!nltk_data] Unzipping tokenizers/punkt.zip.
/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 .
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(
config.json: 100%                               1.58k/1.58k [00:00<00:00, 30.1kB/s]
model.safetensors: 100%                         1.63G/1.63G [00:21<00:00, 227MB/s]
generation_config.json: 100%                    363/363 [00:00<00:00, 19.6kB/s]
vocab.json: 100%                                899k/899k [00:00<00:00, 1.98MB/s]
merges.txt: 100%                                456k/456k [00:00<00:00, 11.3MB/s]
tokenizer.json: 100%                             1.36M/1.36M [00:00<00:00, 18.9MB/s]
Device set to use cpu
config.json: 100%                               1.15k/1.15k [00:00<00:00, 72.2kB/s]
model.safetensors: 100%                         1.63G/1.63G [00:09<00:00, 180MB/s]
tokenizer_config.json: 100%                     26.0/26.0 [00:00<00:00, 1.69kB/s]
vocab.json: 100%                                899k/899k [00:00<00:00, 17.6MB/s]
merges.txt: 100%                                456k/456k [00:00<00:00, 30.5MB/s]
tokenizer.json: 100%                             1.36M/1.36M [00:00<00:00, 38.5MB/s]
Device set to use cpu
modules.json: 100%                              349/349 [00:00<00:00, 33.9kB/s]
config_sentence_transformers.json: 100%          116/116 [00:00<00:00, 10.1kB/s]
README.md: 100%                                10.7k/10.7k [00:00<00:00, 799kB/s]
sentence_bert_config.json: 100%                 53.0/53.0 [00:00<00:00, 3.47kB/s]
config.json: 100%                               612/612 [00:00<00:00, 42.5kB/s]
model.safetensors: 100%                         90.9M/90.9M [00:00<00:00, 180MB/s]
tokenizer_config.json: 100%                     350/350 [00:00<00:00, 26.8kB/s]
vocab.txt: 100%                                 232k/232k [00:00<00:00, 17.4MB/s]
tokenizer.json: 100%                             466k/466k [00:00<00:00, 25.0MB/s]
special_tokens_map.json: 100%                   112/112 [00:00<00:00, 6.13kB/s]
config.json: 100%                               190/190 [00:00<00:00, 17.9kB/s]

-----
LookupError                                Traceback (most recent call last)
<ipython-input-1-8e707ad1f898> in <cell line: 0>()
    31
    32 # Split transcript into sentences
1 #Transcript With timestamp
    35 # Create chunks based on semantic similaritv and concept shifts
1 import re
2 import urllib.parse
3 import requests
4 from youtube_transcript_api import YouTubeTranscriptApi
5 from pytube import YouTube
6 import speech_recognition as sr
7 from pydub import AudioSegment
8 import os
9
10 def extract_video_id(video_url):
11     """
12     Extracts the YouTube video ID from various URL formats.
13     """
14     parsed_url = urllib.parse.urlparse(video_url)
15     query_params = urllib.parse.parse_qs(parsed_url.query)
16
17     if "v" in query_params:
18         return query_params["v"][0]
19
20     match = re.search(r"(youtu\.be/|youtube\.com/embed/|youtube\.com/shorts/)([w-]+)", video_url)
21     if match:
22         return match.group(2)
23
24     return None
25
26 def download_audio(video_url):

```

```

27 """
28 Downloads the audio using yt-dlp with cookies and returns the file path.
29 """
30 try:
31     ydl_opts = {
32         'format': 'bestaudio/best',
33         'outtmpl': 'audio.%(ext)s',
34         'cookiefile': 'cookies (1).txt', # Use the exported cookies
35         'postprocessors': [{
36             'key': 'FFmpegExtractAudio',
37             'preferredcodec': 'mp3',
38             'preferredquality': '192',
39         }],
40     }
41     with yt_dlp.YoutubeDL(ydl_opts) as ydl:
42         info = ydl.extract_info(video_url, download=True)
43         return "audio.mp3"
44 except Exception as e:
45     return f"Error downloading audio: {str(e)}"
46
47 def convert_audio_to_wav(audio_file):
48     """
49     Converts the downloaded MP3 audio to WAV format using pydub.
50     """
51     wav_file = "audio.wav"
52     try:
53         AudioSegment.from_mp3(audio_file).export(wav_file, format="wav")
54         return wav_file
55     except Exception as e:
56         return f"Error converting to WAV: {str(e)}"
57
58 def transcribe_audio(audio_path, chunk_length=30):
59     """
60     Splits audio into smaller chunks and transcribes each chunk separately.
61     Args:
62         audio_path (str): Path to the audio file.
63         chunk_length (int): Length of each chunk in seconds (default: 30).
64     Returns:
65         list: List of dictionaries containing transcribed text and timestamps.
66     """
67     recognizer = sr.Recognizer()
68     audio = AudioSegment.from_wav(audio_path)
69     total_duration = len(audio) / 1000 # Convert to seconds
70     transcribed_segments = []
71
72     print("Transcribing audio in chunks...")
73
74     # Split and transcribe audio in chunks
75     for start in range(0, int(total_duration), chunk_length):
76         end = min(start + chunk_length, int(total_duration))
77         chunk = audio[start * 1000:end * 1000] # Extract chunk in milliseconds
78         chunk.export("chunk.wav", format="wav") # Save chunk temporarily
79
80         with sr.AudioFile("chunk.wav") as source:
81             try:
82                 audio_data = recognizer.record(source)
83                 text = recognizer.recognize_google(audio_data)
84                 transcribed_segments.append({
85                     "start": start,
86                     "end": end,
87                     "text": text
88                 })
89             except sr.UnknownValueError:
90                 transcribed_segments.append({
91                     "start": start,
92                     "end": end,
93                     "text": "[Unintelligible]"
94                 })
95             except sr.RequestError as e:
96                 return f"Error with the speech recognition service: {str(e)}"
97
98     os.remove("chunk.wav") # Clean up temporary chunk file
99     return transcribed_segments
100
101 def get_transcript_unlisted(video_url):
102     """
103     Tries to fetch the transcript using youtube_transcript_api first,
104     then falls back to downloading and transcribing audio if necessary.
105     """
106     video_id = extract_video_id(video_url)
107     if not video_id:
108         return "Invalid YouTube URL."

```

```

109
110 # Try to fetch transcript using youtube_transcript_api
111 try:
112     transcript = YouTubeTranscriptApi.get_transcript(video_id)
113     # Add 'end' time to each segment
114     for segment in transcript:
115         segment["end"] = segment["start"] + segment["duration"]
116     return transcript # Return transcript with timestamps
117 except:
118     print("Transcript not available via API, attempting audio transcription...")
119
120 # Download and transcribe audio if no transcript is available
121 audio_file = download_audio(video_url)
122 if "Error" in audio_file:
123     return audio_file
124
125 wav_file = convert_audio_to_wav(audio_file)
126 if "Error" in wav_file:
127     return wav_file
128
129 transcription = transcribe_audio(wav_file)
130
131 # Cleanup temporary files
132 os.remove(audio_file)
133 os.remove(wav_file)
134
135 return transcription
136
137 # Example usage
138 if __name__ == "__main__":
139     video_url = input("Enter the YouTube video URL: ")
140     transcript = get_transcript_unlisted(video_url)
141
142     if isinstance(transcript, list):
143         print("\nTranscript with Timestamps:")
144         for segment in transcript:
145             print(f"{segment['start']} - {segment['end']}: {segment['text']}")
146     else:
147         print("\nTranscript:\n", transcript)

```

Enter the YouTube video URL: <https://youtu.be/sK8SILOM37I>

Transcript with Timestamps:

0.08 - 5.92: so sir we know that India has seen a
3.04 - 8.280000000000001: huge Revolution with digital payments we
5.92 - 9.679: all thought that India is a place at
8.28 - 12.559: least the West thought that India is a
9.679 - 14.799: place where uh many people do not get a
12.559 - 17.039: square meal right that was the narrative
14.799 - 21.358999999999998: some 30 years ago and not many are
17.039 - 24.840000000000003: literate people cannot read um but then
21.359 - 26.679000000000002: we have now shown that digital payments
24.84 - 29.32: number one is India while people thought
26.679 - 30.759: that it wouldn't even come to top 50 let
29.32 - 32.88: alone number one
30.759 - 34.16: I think immediately after UPA the next
32.88 - 37.84: big
34.16 - 40.599999999999994: revolution personally I think is in
37.84 - 43.28: education and the complete homework for
40.6 - 45.68: this has happened in the form of NEP the
43.28 - 48.079: documentation of which many of us have
45.68 - 52.16: read and realize that it is too good to
48.079 - 53.64: be true followed by which we got ncrf
52.16 - 55.68: framework done which is the national
53.64 - 57.44: credit framework surrounding which we'll
55.68 - 60.48: be discussing
57.44 - 63.76: today sir my question is do you think
60.48 - 67.158999999999999: ncrf plus NEP put together will be the
63.76 - 69.2: next big revolution after UPA in
67.159 - 72.0: India
69.2 - 77.36: absolutely and why I think
72.0 - 80.119: so is because in education the last
77.36 - 82.56: policy came up many many years ago that
80.119 - 87.079: was in
82.56 - 87.079000000000001: 1986 which was slightly tweaked in
87.72 - 93.2: 1982 and thereafter so so many changes
90.68 - 96.399: have happened in the real
93.2 - 99.119: world so many changes have happened in
96.399 - 102.320000000000001: the requirement of the
99.119 - 104.2: industry requirement of Manpower for the
102.32 - 106.199999999999999: industry however there were no
104.2 - 107.479: corresponding changes which happened
106.2 - 111.0: education
107.479 - 113.32: system so therefore I feel that this was
111.0 - 116.52: the right time when we brought
113.32 - 121.798999999999999: in the education policy 2020 honorable

```

116.52 - 125.19999999999999: prime minister declared it open in July
121.799 - 129.399: 2020 and we recently celebrated the
125.2 - 133.4: fourth anniversary of NEP 2020 ncrf has
129.399 - 136.8: been brought to implement the intent of
133.4 - 139.36: 2020 2020 is a policy the NP 2020 is a
136.8 - 141.44: policy and for implementing a policy you
139.36 - 144.48000000000002: need a
141.44 - 149.04: framework now why we call it a framework
144.48 - 151.16: we call it a framework because this is
149.04 - 154.35999999999999: very flexible
151.16 - 159.4: this allows you all the
154.36 - 161.84: Transcription the way you educate your kids

```

```
1 pip install youtube_transcript_api
```

```

Collecting youtube_transcript_api
  Downloading youtube_transcript_api-0.6.3-py3-none-any.whl.metadata (17 kB)
Requirement already satisfied: defusedxml<0.8.0,>=0.7.1 in /usr/local/lib/python3.11/dist-packages (from youtube_transcript_api) (0.7.1)
Requirement already satisfied: requests in /usr/local/lib/python3.11/dist-packages (from youtube_transcript_api) (2.32.3)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests->youtube_transcript_api) (3.4.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests->youtube_transcript_api) (3.10.1)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests->youtube_transcript_api) (2.2.3)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests->youtube_transcript_api) (2025.1.1)
  Downloading youtube_transcript_api-0.6.3-py3-none-any.whl (622 kB)
    622.3/622.3 kB 25.5 MB/s eta 0:00:00
Installing collected packages: youtube_transcript_api
Successfully installed youtube_transcript_api-0.6.3

```

```
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 9.7 MB/s eta 0:00:00
  Downloading yt_dlp-2025.2.19-py3-none-any.whl (3.2 MB)
    3.2/3.2 MB 47.6 MB/s eta 0:00:00
Installing collected packages: yt-dlp
Successfully installed yt-dlp-2025.2.19

```

```
1 pip install pytube
```

```

Collecting pytube
  Downloading pytube-15.0.0-py3-none-any.whl.metadata (5.0 kB)
  Downloading pytube-15.0.0-py3-none-any.whl (57 kB)
    57.6/57.6 kB 4.1 MB/s eta 0:00:00
Installing collected packages: pytube
Successfully installed pytube-15.0.0

```

```
1 pip install SpeechRecognition
```

```

Collecting SpeechRecognition
  Downloading SpeechRecognition-3.14.1-py3-none-any.whl.metadata (31 kB)
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.11/dist-packages (from SpeechRecognition) (4.12.2)
  Downloading SpeechRecognition-3.14.1-py3-none-any.whl (32.9 MB)
    32.9/32.9 MB 17.4 MB/s eta 0:00:00
Installing collected packages: SpeechRecognition
Successfully installed SpeechRecognition-3.14.1

```

```
1 pip install pydub
```

```

Collecting pydub
  Downloading pydub-0.25.1-py2.py3-none-any.whl.metadata (1.4 kB)
  Downloading pydub-0.25.1-py2.py3-none-any.whl (32 kB)
Installing collected packages: pydub
Successfully installed pydub-0.25.1

```

```

1 import re
2 import urllib.parse
3 import requests
4 from youtube_transcript_api import YouTubeTranscriptApi
5 from pytube import YouTube
6 import speech_recognition as sr
7 from pydub import AudioSegment
8 import os
9
10 def extract_video_id(video_url):
11     """
12     Extracts the YouTube video ID from various URL formats.
13     """
14     parsed_url = urllib.parse.urlparse(video_url)
15     query_params = urllib.parse.parse_qs(parsed_url.query)
16
17     if "v" in query_params:

```

```

18     return query_params["v"][0]
19
20     match = re.search(r"(youtu\.be/|youtube\.com/embed/|youtube\.com/shorts/)([^\w-]+)", video_url)
21     if match:
22         return match.group(2)
23
24     return None
25
26 def download_audio(video_url):
27     """
28     Downloads the audio using yt-dlp with cookies and returns the file path.
29     """
30     try:
31         ydl_opts = {
32             'format': 'bestaudio/best',
33             'outtmpl': 'audio.%(ext)s',
34             'cookiefile': 'cookies (1).txt', # Use the exported cookies
35             'postprocessors': [{
36                 'key': 'FFmpegExtractAudio',
37                 'preferredcodec': 'mp3',
38                 'preferredquality': '192',
39             }],
40         }
41         with yt_dlp.YoutubeDL(ydl_opts) as ydl:
42             info = ydl.extract_info(video_url, download=True)
43             return "audio.mp3"
44     except Exception as e:
45         return f"Error downloading audio: {str(e)}"
46
47 def convert_audio_to_wav(audio_file):
48     """
49     Converts the downloaded MP3 audio to WAV format using pydub.
50     """
51     wav_file = "audio.wav"
52     try:
53         AudioSegment.from_mp3(audio_file).export(wav_file, format="wav")
54         return wav_file
55     except Exception as e:
56         return f"Error converting to WAV: {str(e)}"
57
58 def transcribe_audio(audio_path, chunk_length=30):
59     """
60     Splits audio into smaller chunks and transcribes each chunk separately.
61     Args:
62         audio_path (str): Path to the audio file.
63         chunk_length (int): Length of each chunk in seconds (default: 30).
64     Returns:
65         list: List of dictionaries containing transcribed text and timestamps.
66     """
67     recognizer = sr.Recognizer()
68     audio = AudioSegment.from_wav(audio_path)
69     total_duration = len(audio) / 1000 # Convert to seconds
70     transcribed_segments = []
71
72     print("Transcribing audio in chunks...")
73
74     # Split and transcribe audio in chunks
75     for start in range(0, int(total_duration), chunk_length):
76         end = min(start + chunk_length, int(total_duration))
77         chunk = audio[start * 1000:end * 1000] # Extract chunk in milliseconds
78         chunk.export("chunk.wav", format="wav") # Save chunk temporarily
79
80         with sr.AudioFile("chunk.wav") as source:
81             try:
82                 audio_data = recognizer.record(source)
83                 text = recognizer.recognize_google(audio_data)
84                 transcribed_segments.append({
85                     "start": start,
86                     "end": end,
87                     "text": text
88                 })
89             except sr.UnknownValueError:
90                 transcribed_segments.append({
91                     "start": start,
92                     "end": end,
93                     "text": "[Unintelligible]"
94                 })
95             except sr.RequestError as e:
96                 return f"Error with the speech recognition service: {str(e)}"
97
98     os.remove("chunk.wav") # Clean up temporary chunk file
99     return transcribed_segments

```

```

100
101 def get_transcript_unlisted(video_url):
102     """
103     Tries to fetch the transcript using youtube_transcript_api first,
104     then falls back to downloading and transcribing audio if necessary.
105     """
106     video_id = extract_video_id(video_url)
107     if not video_id:
108         return "Invalid YouTube URL."
109
110     # Try to fetch transcript using youtube_transcript_api
111     try:
112         transcript = YouTubeTranscriptApi.get_transcript(video_id)
113         # Add 'end' time to each segment
114         for segment in transcript:
115             segment["end"] = segment["start"] + segment["duration"]
116         return transcript # Return transcript with timestamps
117     except:
118         print("Transcript not available via API, attempting audio transcription...")
119
120     # Download and transcribe audio if no transcript is available
121     audio_file = download_audio(video_url)
122     if "Error" in audio_file:
123         return audio_file
124
125     wav_file = convert_audio_to_wav(audio_file)
126     if "Error" in wav_file:
127         return wav_file
128
129     transcription = transcribe_audio(wav_file)
130
131     # Cleanup temporary files
132     os.remove(audio_file)
133     os.remove(wav_file)
134
135     return transcription
136
137 def save_transcript_to_file(transcript, filename="transcript.txt"):
138     """
139     Saves the transcript to a text file.
140     Args:
141         transcript (list or str): The transcript to save.
142         filename (str): The name of the output file.
143     """
144     with open(filename, "w", encoding="utf-8") as file:
145         if isinstance(transcript, list):
146             for segment in transcript:
147                 file.write(f"{segment['start']} - {segment['end']}: {segment['text']}\n")
148         else:
149             file.write(transcript)
150     print(f"Transcript saved to {filename}")
151
152 # Example usage
153 if __name__ == "__main__":
154     video_url = input("Enter the YouTube video URL: ")
155     transcript = get_transcript_unlisted(video_url)
156
157     if isinstance(transcript, list):
158         print("\nTranscript with Timestamps:")
159         for segment in transcript:
160             print(f"{segment['start']} - {segment['end']}: {segment['text']}")
161     else:
162         print("\nTranscript:\n", transcript)
163
164     # Save transcript to a text file
165     save_transcript_to_file(transcript, "transcript.txt")

```

Enter the YouTube video URL: <https://youtu.be/sK8SILOM37I>

Transcript with Timestamps:
0.08 - 5.92: so sir we know that India has seen a
3.04 - 8.280000000000001: huge Revolution with digital payments we
5.92 - 9.679: all thought that India is a place at
8.28 - 12.559: least the West thought that India is a
9.679 - 14.799: place where uh many people do not get a
12.559 - 17.039: square meal right that was the narrative
14.799 - 21.358999999999998: some 30 years ago and not many are
17.039 - 24.840000000000003: literate people cannot read um but then
21.359 - 26.679000000000002: we have now shown that digital payments
24.84 - 29.32: number one is India while people thought
26.679 - 30.759: that it wouldn't even come to top 50 let
29.32 - 32.88: alone number one
30.759 - 34.16: I think immediately after UPA the next
32.88 - 37.84: big



34.16 - 40.599999999999994: revolution personally I think is in
 37.84 - 43.28: education and the complete homework for
 40.6 - 45.68: this has happened in the form of NEP the
 43.28 - 48.079: documentation of which many of us have
 45.68 - 52.16: read and realize that it is too good to
 48.079 - 53.64: be true followed by which we got ncrf
 52.16 - 55.68: framework done which is the national
 53.64 - 57.44: credit framework surrounding which we'll
 55.68 - 60.48: be discussing
 57.44 - 63.76: today sir my question is do you think
 60.48 - 67.158999999999999: ncrf plus NEP put together will be the
 63.76 - 69.2: next big revolution after UPA in
 67.159 - 72.0: India
 69.2 - 77.36: absolutely and why I think
 72.0 - 80.119: so is because in education the last
 77.36 - 82.56: policy came up many many years ago that
 80.119 - 87.079: was in
 82.56 - 87.079000000000001: 1986 which was slightly tweaked in
 87.72 - 93.2: 1982 and thereafter so so many changes
 90.68 - 96.399: have happened in the real
 93.2 - 99.119: world so many changes have happened in
 96.399 - 102.320000000000001: the requirement of the
 99.119 - 104.2: industry requirement of Manpower for the
 102.32 - 106.199999999999999: industry however there were no
 104.2 - 107.479: corresponding changes which happened
 106.2 - 111.0: education
 107.479 - 113.32: system so therefore I feel that this was
 111.0 - 116.52: the right time when we brought
 113.32 - 121.798999999999999: in the education policy 2020 honorable
 116.52 - 125.199999999999999: prime minister declared it open in July
 121.799 - 129.399: 2020 and we recently celebrated the
 125.2 - 133.4: fourth anniversary of NEP 2020 ncrf has
 129.399 - 136.8: been brought to implement the intent of
 133.4 - 139.36: 2020 2020 is a policy the NP 2020 is a
 136.8 - 141.44: policy and for implementing a policy you
 139.36 - 144.480000000000002: need a
 141.44 - 149.04: framework now why we call it a framework
 144.48 - 151.16: we call it a framework because this is
 149.04 - 154.359999999999999: very flexible
 151.16 - 159.4: this allows you all the
 154.36 - 161.84: Innovation the way you educate your kids

```
1 import json
2 from sklearn.cluster import KMeans
3 from sentence_transformers import SentenceTransformer
4 import numpy as np
5
6 def load_transcript_from_file(file_path):
7     """
8     Loads the transcript from a file.
9     Args:
10         file_path (str): Path to the transcript file.
11     Returns:
12         list: List of dictionaries with 'start', 'end', and 'text' keys.
13     """
14     try:
15         with open(file_path, "r", encoding="utf-8") as file:
16             if file_path.endswith(".json"):
17                 transcript = json.load(file)
18             else:
19                 # For plain text files, assume each line is in the format: [start] - [end]: [text]
20                 transcript = []
21                 for line in file:
22                     match = re.match(r"(\d+\.\d+) - (\d+\.\d+): (.+)", line.strip())
23                     if match:
24                         start, end, text = match.groups()
25                         transcript.append({
26                             "start": float(start),
27                             "end": float(end),
28                             "text": text
29                         })
30                 return transcript
31     except Exception as e:
32         print(f"Error loading transcript: {e}")
33         return None
34
35 def segment_transcript_semantically(transcript, num_clusters=5):
36     """
37     Segments the transcript into semantically meaningful clusters.
38     Args:
39         transcript (list): List of dictionaries with 'start', 'end', and 'text' keys.
40         num_clusters (int): Number of clusters to create (default: 5).
41     Returns:
42         list: List of dictionaries with 'start', 'end', and 'text' for each segment.
43     """
44     # Extract sentences and their timestamps
```

```

45 sentences = [segment["text"] for segment in transcript]
46 timestamps = [(segment["start"], segment["end"]) for segment in transcript]
47
48 # Load a pre-trained sentence embedding model
49 model = SentenceTransformer('all-MiniLM-L6-v2')
50 sentence_embeddings = model.encode(sentences)
51
52 # Perform K-Means clustering
53 kmeans = KMeans(n_clusters=num_clusters, random_state=42)
54 clusters = kmeans.fit_predict(sentence_embeddings)
55
56 # Group sentences by cluster
57 clustered_segments = {}
58 for i, cluster in enumerate(clusters):
59     if cluster not in clustered_segments:
60         clustered_segments[cluster] = []
61     clustered_segments[cluster].append({
62         "start": timestamps[i][0],
63         "end": timestamps[i][1],
64         "text": sentences[i]
65     })
66
67 # Merge segments in each cluster
68 segmented_transcript = []
69 for cluster, segments in clustered_segments.items():
70     start_time = segments[0]["start"]
71     end_time = segments[-1]["end"]
72     combined_text = " ".join([segment["text"] for segment in segments])
73     segmented_transcript.append({
74         "start": start_time,
75         "end": end_time,
76         "text": combined_text
77     })
78
79 # Sort segments by start time
80 segmented_transcript.sort(key=lambda x: x["start"])
81
82 return segmented_transcript
83
84 def save_segmented_transcript(segmented_transcript, output_path="segmented_transcript.txt"):
85     """
86     Saves the segmented transcript to a text file.
87     Args:
88         segmented_transcript (list): List of dictionaries with 'start', 'end', and 'text'.
89         output_path (str): The path to save the output file.
90     """
91     with open(output_path, "w", encoding="utf-8") as file:
92         for segment in segmented_transcript:
93             file.write(f"{segment['start']} - {segment['end']}: {segment['text']}\n")
94     print(f"Segmented transcript saved to {output_path}")
95
96 # Example usage
97 if __name__ == "__main__":
98     # Ask the user for the transcript file path
99     transcript_path = input("Enter the path to the transcript file: ")
100     output_path = input("Enter the path to save the segmented transcript (default: segmented_transcript.txt): ") or "segmented_tra
101
102     # Load the transcript from the file
103     transcript = load_transcript_from_file(transcript_path)
104     if not transcript:
105         print("Failed to load transcript. Exiting.")
106         exit()
107
108     # Segment the transcript semantically
109     segmented_transcript = segment_transcript_semantically(transcript, num_clusters=3)
110
111     # Print the segmented transcript
112     print("\nSegmented Transcript:")
113     for segment in segmented_transcript:
114         print(f"{segment['start']} - {segment['end']}: {segment['text']}")
115
116     # Save the segmented transcript to a file
117     save_segmented_transcript(segmented_transcript, output_path)

```



Enter the path to the transcript file: /content/transcript.txt
Enter the path to save the segmented transcript (default: segmented_transcript.txt):

Segmented Transcript:

0.08 - 3444.2: so sir we know that India has seen a all thought that India is a place at least the West thought that India is a place
3.04 - 3436.7599999999998: huge Revolution with digital payments we documentation of which many of us have framework done which is 1
37.84 - 3444.2000000000003: education and the complete homework for so is because in education the last education in the education
Segmented transcript saved to segmented_transcript.txt

```

1 import json
2 import re
3 from bertopic import BERTopic
4 from sentence_transformers import SentenceTransformer
5 from sklearn.feature_extraction.text import CountVectorizer
6
7 def load_transcript_from_file(file_path):
8     """
9     Loads the transcript from a file.
10    Args:
11        file_path (str): Path to the transcript file.
12    Returns:
13        list: List of dictionaries with 'start', 'end', and 'text' keys.
14    """
15    try:
16        with open(file_path, "r", encoding="utf-8") as file:
17            if file_path.endswith(".json"):
18                transcript = json.load(file)
19            else:
20                # For plain text files, assume each line is in the format: [start] - [end]: [text]
21                transcript = []
22                for line in file:
23                    match = re.match(r"(\d+\.\d+) - (\d+\.\d+): (.+)", line.strip())
24                    if match:
25                        start, end, text = match.groups()
26                        transcript.append({
27                            "start": float(start),
28                            "end": float(end),
29                            "text": text
30                        })
31                return transcript
32    except Exception as e:
33        print(f"Error loading transcript: {e}")
34        return None
35
36 def segment_transcript_with_topics(transcript, num_topics=5):
37     """
38     Segments the transcript into topics using BERTopic.
39    Args:
40        transcript (list): List of dictionaries with 'start', 'end', and 'text' keys.
41        num_topics (int): Number of topics to create (default: 5).
42    Returns:
43        list: List of dictionaries with 'start', 'end', 'text', and 'topic' for each segment.
44    """
45    # Extract sentences and their timestamps
46    sentences = [segment["text"] for segment in transcript]
47    timestamps = [(segment["start"], segment["end"]) for segment in transcript]
48
49    # Load a pre-trained sentence embedding model
50    embedding_model = SentenceTransformer('all-MiniLM-L6-v2')
51
52    # Initialize BERTopic
53    vectorizer_model = CountVectorizer(stop_words="english") # Remove stopwords for better topic modeling
54    topic_model = BERTopic(embedding_model=embedding_model, vectorizer_model=vectorizer_model, nr_topics=num_topics)
55
56    # Fit the model and transform the sentences
57    topics, _ = topic_model.fit_transform(sentences)
58
59    # Assign topics to each sentence
60    segmented_transcript = []
61    for i, (start, end) in enumerate(timestamps):
62        segmented_transcript.append({
63            "start": start,
64            "end": end,
65            "text": sentences[i],
66            "topic": int(topics[i])
67        })
68
69    return segmented_transcript, topic_model
70
71 def save_segmented_transcript(segmented_transcript, topic_model, output_path="segmented_transcript.txt"):
72     """
73     Saves the segmented transcript with topics to a text file.
74    Args:
75        segmented_transcript (list): List of dictionaries with 'start', 'end', 'text', and 'topic'.
76        topic_model (BERTopic): The trained BERTopic model.
77        output_path (str): The path to save the output file.
78    """
79    # Get topic labels

```

```

80 topic_info = topic_model.get_topic_info()
81 topic_labels = {row["Topic"]: row["Name"] for _, row in topic_info.iterrows()}
82
83 with open(output_path, "w", encoding="utf-8") as file:
84     for segment in segmented_transcript:
85         topic_label = topic_labels.get(segment["topic"], "Unknown Topic")
86         file.write(f"{segment['start']} - {segment['end']} [{topic_label}]: {segment['text']}\n")
87     print(f"Segmented transcript saved to {output_path}")
88
89 # Example usage
90 if __name__ == "__main__":
91     # Ask the user for the transcript file path
92     transcript_path = input("Enter the path to the transcript file: ")
93     output_path = input("Enter the path to save the segmented transcript (default: segmented_transcript.txt): ") or "segmented_tra
94
95     # Load the transcript from the file
96     transcript = load_transcript_from_file(transcript_path)
97     if not transcript:
98         print("Failed to load transcript. Exiting.")
99         exit()
100
101     # Segment the transcript into topics
102     segmented_transcript, topic_model = segment_transcript_with_topics(transcript, num_topics=5)
103
104     # Print the segmented transcript
105     print("\nSegmented Transcript with Topics:")
106     for segment in segmented_transcript:
107         print(f"{segment['start']} - {segment['end']} [Topic {segment['topic']}]: {segment['text']}")
108
109     # Save the segmented transcript to a file
110     save_segmented_transcript(segmented_transcript, topic_model, output_path)

```

Enter the path to the transcript file: /content/transcript.txt
Enter the path to save the segmented transcript (default: segmented_transcript.txt):

Segmented Transcript with Topics:

```

0.08 - 5.92 [Topic 1]: so sir we know that India has seen a
3.04 - 8.280000000000001 [Topic -1]: huge Revolution with digital payments we
5.92 - 9.679 [Topic 1]: all thought that India is a place at
8.28 - 12.559 [Topic 1]: least the West thought that India is a
9.679 - 14.799 [Topic -1]: place where uh many people do not get a
12.559 - 17.039 [Topic -1]: square meal right that was the narrative
14.799 - 21.358999999999998 [Topic -1]: some 30 years ago and not many are
17.039 - 24.840000000000003 [Topic -1]: literate people cannot read um but then
21.359 - 26.679000000000002 [Topic -1]: we have now shown that digital payments
24.84 - 29.32 [Topic 1]: number one is India while people thought
26.679 - 30.759 [Topic 1]: that it wouldn't even come to top 50 let
29.32 - 32.88 [Topic 1]: alone number one
30.759 - 34.16 [Topic 1]: I think immediately after UPA the next
32.88 - 37.84 [Topic 1]: big
34.16 - 40.599999999999994 [Topic 1]: revolution personally I think is in
37.84 - 43.28 [Topic 0]: education and the complete homework for
40.6 - 45.68 [Topic 2]: this has happened in the form of NEP the
43.28 - 48.079 [Topic -1]: documentation of which many of us have
45.68 - 52.16 [Topic -1]: read and realize that it is too good to
48.079 - 53.64 [Topic 2]: be true followed by which we got ncrf
52.16 - 55.68 [Topic 2]: framework done which is the national
53.64 - 57.44 [Topic 0]: credit framework surrounding which we'll
55.68 - 60.48 [Topic -1]: be discussing
57.44 - 63.76 [Topic 1]: today sir my question is do you think
60.48 - 67.158999999999999 [Topic 2]: ncrf plus NEP put together will be the
63.76 - 69.2 [Topic 1]: next big revolution after UPA in
67.159 - 72.0 [Topic 1]: India
69.2 - 77.36 [Topic 1]: absolutely and why I think
72.0 - 80.119 [Topic 0]: so is because in education the last
77.36 - 82.56 [Topic -1]: policy came up many many years ago that
80.119 - 87.079 [Topic -1]: was in
82.56 - 87.079000000000001 [Topic 1]: 1986 which was slightly tweaked in
87.72 - 93.2 [Topic 1]: 1982 and thereafter so so many changes
90.68 - 96.399 [Topic 1]: have happened in the real
93.2 - 99.119 [Topic 1]: world so many changes have happened in
96.399 - 102.320000000000001 [Topic -1]: the requirement of the
99.119 - 104.2 [Topic 0]: industry requirement of Manpower for the
102.32 - 106.199999999999999 [Topic 0]: industry however there were no
104.2 - 107.479 [Topic 1]: corresponding changes which happened
106.2 - 111.0 [Topic 0]: education
107.479 - 113.32 [Topic -1]: system so therefore I feel that this was
111.0 - 116.52 [Topic 1]: the right time when we brought
113.32 - 121.798999999999999 [Topic 0]: in the education policy 2020 honorable
116.52 - 125.199999999999999 [Topic -1]: prime minister declared it open in July
121.799 - 129.399 [Topic -1]: 2020 and we recently celebrated the
125.2 - 133.4 [Topic 2]: fourth anniversary of NEP 2020 ncrf has
129.399 - 136.8 [Topic 0]: been brought to implement the intent of
133.4 - 139.36 [Topic -1]: 2020 2020 is a policy the NP 2020 is a
136.8 - 141.44 [Topic 0]: policy and for implementing a policy you
139.36 - 144.480000000000002 [Topic 1]: need a
141.44 - 149.04 [Topic 2]: framework now why we call it a framework

```

```
144.48 - 151.16 [Topic 2]: we call it a framework because this is
149.04 - 154.35999999999999 [Topic -1]: very flexible
151.16 - 159.4 [Topic 1]: this allows you all the
```

```
1 pip install bertopic
```

```
Collecting bertopic
  Downloading bertopic-0.16.4-py3-none-any.whl.metadata (23 kB)
Requirement already satisfied: hdbscan>=0.8.29 in /usr/local/lib/python3.11/dist-packages (from bertopic) (0.8.40)
Requirement already satisfied: numpy>=1.20.0 in /usr/local/lib/python3.11/dist-packages (from bertopic) (1.26.4)
Requirement already satisfied: pandas>=1.1.5 in /usr/local/lib/python3.11/dist-packages (from bertopic) (2.2.2)
Requirement already satisfied: plotly>=4.7.0 in /usr/local/lib/python3.11/dist-packages (from bertopic) (5.24.1)
Requirement already satisfied: scikit-learn>=0.22.2.post1 in /usr/local/lib/python3.11/dist-packages (from bertopic) (1.6.1)
Requirement already satisfied: sentence-transformers>=0.4.1 in /usr/local/lib/python3.11/dist-packages (from bertopic) (3.4.1)
Requirement already satisfied: tqdm>=4.41.1 in /usr/local/lib/python3.11/dist-packages (from bertopic) (4.67.1)
Requirement already satisfied: umap-learn>=0.5.0 in /usr/local/lib/python3.11/dist-packages (from bertopic) (0.5.7)
Requirement already satisfied: scipy>=1.0 in /usr/local/lib/python3.11/dist-packages (from hdbscan>=0.8.29->bertopic) (1.13.1)
Requirement already satisfied: joblib>=1.0 in /usr/local/lib/python3.11/dist-packages (from hdbscan>=0.8.29->bertopic) (1.4.2)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.1.5->bertopic) (2.9.0)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.1.5->bertopic) (2025.1)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.1.5->bertopic) (2025.1)
Requirement already satisfied: packaging>=6.2.0 in /usr/local/lib/python3.11/dist-packages (from plotly>=4.7.0->bertopic) (24.2)
Requirement already satisfied: packaging in /usr/local/lib/python3.11/dist-packages (from plotly>=4.7.0->bertopic) (24.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn>=0.22.2.post1->bertopic) (3.5.0)
Requirement already satisfied: transformers<5.0.0,>=4.41.0 in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.4.1->bertopic) (4.41.0)
Requirement already satisfied: torch>=1.11.0 in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.4.1->bertopic) (2.5.1)
Requirement already satisfied: huggingface-hub>=0.20.0 in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.4.1->bertopic) (0.20.0)
Requirement already satisfied: Pillow in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.4.1->bertopic) (10.4.0)
Requirement already satisfied: numba>=0.51.2 in /usr/local/lib/python3.11/dist-packages (from umap-learn>=0.5.0->bertopic) (0.61.0)
Requirement already satisfied: pynndescent>=0.5 in /usr/local/lib/python3.11/dist-packages (from umap-learn>=0.5.0->bertopic) (0.5.10)
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: 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) (2.32.0)
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) (4.12.0)
Requirement already satisfied: llvmlite<0.45,>=0.44.0dev0 in /usr/local/lib/python3.11/dist-packages (from numba>=0.51.2->umap-learn) (0.44.0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas>=1.1.5->bertopic) (1.17.0)
Requirement already satisfied: networkx in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic) (3.4.2)
Requirement already satisfied: Jinja2 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic) (3.1.4)
Collecting nvidia-cuda-nvrtc-cu12==12.4.127 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cuda_nvrtc_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cuda-runtime-cu12==12.4.127 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cuda_runtime_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cuda-cupti-cu12==12.4.127 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cuda_cupti_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cudnn-cu12==9.1.0.70 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cudnn_cu12-9.1.0.70-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cublas-cu12==12.4.5.8 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cublas_cu12-12.4.5.8-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cufft-cu12==11.2.1.3 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cufft_cu12-11.2.1.3-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-curand-cu12==10.3.5.147 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_curand_cu12-10.3.5.147-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cusolver-cu12==11.6.1.9 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cusolver_cu12-11.6.1.9-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cuspars-cu12==12.3.1.170 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cuspars-cu12-12.3.1.170-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic) (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>=0.4.1->bertopic) (1.11.0)
Collecting nvidia-nvjitlink-cu12==12.4.127 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_nvjitlink_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Requirement already satisfied: triton==3.1.0 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic) (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>=0.4.1->bertopic) (1.13.1)
```

```
1 from sentence_transformers import SentenceTransformer
2 from sklearn.cluster import KMeans
3 import re
4
5 # Load pre-trained model for sentence embeddings
6 model = SentenceTransformer('all-MiniLM-L6-v2')
7
8 # Function to extract sentences and starting timestamps
9 def extract_sentences_and_timestamps(file_path):
10     sentences = []
11     timestamps = []
12     # Enhanced regex pattern for range-based timestamps
13     pattern = re.compile(r"(\d{1,2}(\.?|\d+)?)\s*-\s*(\d{1,2}(\.?|\d+)?):\s*(.*)")
14
15     try:
16         with open(file_path, 'r') as file:
17             for line in file:
18                 line = line.strip()
19                 if line:
20                     match = pattern.match(line)
21                     if match:
22                         start_time = float(match.group(1))
```

```

23         # Convert start time to hh:mm:ss format
24         hours = int(start_time // 3600)
25         minutes = int((start_time % 3600) // 60)
26         seconds = start_time % 60
27         timestamp = f"{hours:02}:{minutes:02}:{seconds:05.2f}"
28         timestamps.append(timestamp)
29         sentences.append(match.group(3).strip())
30     else:
31         print(f"⚠️ Unmatched line format: {line}") # Debug print for unmatched lines
32 except FileNotFoundError:
33     print(f"❌ Error: The file '{file_path}' was not found.")
34     return [], []
35
36 if not sentences:
37     print("⚠️ Warning: No sentences found. Check if the file is empty or if the format is incorrect.")
38 return sentences, timestamps
39
40 # Extract sentences and timestamps
41 input_file = 'transcript.txt'
42 sentences, timestamps = extract_sentences_and_timestamps(input_file)
43
44 # Debug: Print extracted sentences and timestamps
45 print("\nExtracted Sentences:", sentences)
46 print("Extracted Timestamps:", timestamps)
47
48 # Check if sentences are empty
49 if not sentences:
50     print("❌ Error: No valid sentences were extracted. Exiting the script.")
51 else:
52     # Generate sentence embeddings
53     embeddings = model.encode(sentences)
54
55     # Debug: Check if embeddings are generated correctly
56     if embeddings.size == 0:
57         print("❌ Error: No embeddings were generated. Check if the input text is valid.")
58     else:
59         # Determine number of clusters (segments)
60         num_clusters = 3 # You can adjust this based on your transcript length
61         kmeans = KMeans(n_clusters=num_clusters, random_state=0)
62         labels = kmeans.fit_predict(embeddings)
63
64         # Create segmented output with timestamps
65         segments = [[] for _ in range(num_clusters)]
66         for i, label in enumerate(labels):
67             segments[label].append(f"[{timestamps[i]}] {sentences[i]}")
68
69         # Save segmented output to a new file
70         output_file = 'segmented_transcript.txt'
71         with open(output_file, 'w') as file:
72             for idx, segment in enumerate(segments):
73                 file.write(f"--- Segment {idx + 1} ---\n")
74                 file.write("\n".join(segment))
75                 file.write("\n\n")
76
77         print(f"✅ Segmented transcript saved to {output_file}")

```



⚠️ Unmatched line format: 96.399 - 102.32000000000001: the requirement of the
 ⚠️ Unmatched line format: 99.119 - 104.2: industry requirement of Manpower for the
 ⚠️ Unmatched line format: 102.32 - 106.19999999999999: industry however there were no
 ⚠️ Unmatched line format: 104.2 - 107.479: corresponding changes which happened
 ⚠️ Unmatched line format: 106.2 - 111.0: education
 ⚠️ Unmatched line format: 107.479 - 113.32: system so therefore I feel that this was
 ⚠️ Unmatched line format: 111.0 - 116.52: the right time when we brought
 ⚠️ Unmatched line format: 113.32 - 121.79899999999999: in the education policy 2020 honorable
 ⚠️ Unmatched line format: 116.52 - 125.19999999999999: prime minister declared it open in July
 ⚠️ Unmatched line format: 121.799 - 129.399: 2020 and we recently celebrated the
 ⚠️ Unmatched line format: 125.2 - 133.4: fourth anniversary of NEP 2020 ncrf has
 ⚠️ Unmatched line format: 129.399 - 136.8: been brought to implement the intent of
 ⚠️ Unmatched line format: 133.4 - 139.36: 2020 2020 is a policy the NP 2020 is a
 ⚠️ Unmatched line format: 136.8 - 141.44: policy and for implementing a policy you
 ⚠️ Unmatched line format: 139.36 - 144.48000000000002: need a
 ⚠️ Unmatched line format: 141.44 - 149.04: framework now why we call it a framework
 ⚠️ Unmatched line format: 144.48 - 151.16: we call it a framework because this is
 ⚠️ Unmatched line format: 149.04 - 154.35999999999999: very flexible
 ⚠️ Unmatched line format: 151.16 - 159.4: this allows you all the
 ⚠️ Unmatched line format: 154.36 - 161.84: Innovation the way you educate your kids
 ⚠️ Unmatched line format: 159.4 - 165.44: you educate your
 ⚠️ Unmatched line format: 161.84 - 168.48: students still it provides you the basic
 ⚠️ Unmatched line format: 165.44 - 171.959: guidelines the the framework the outer
 ⚠️ Unmatched line format: 168.48 - 174.84: layer it provides and that layer is
 ⚠️ Unmatched line format: 171.959 - 177.84: mostly the enabling
 ⚠️ Unmatched line format: 174.84 - 180.76: layer that is such an enabling layer
 ⚠️ Unmatched line format: 177.84 - 183.48: that it has
 ⚠️ Unmatched line format: 180.76 - 185.879: broken the shackles which were there in
 ⚠️ Unmatched line format: 183.48 - 189.48: the education



```

⚠ Unmatched line format: 185.879 - 191.959: sector right so yes it is a big
⚠ Unmatched line format: 189.48 - 195.23899999999998: Revolution and this is going to change
⚠ Unmatched line format: 191.959 - 197.68: the way we have been educating our kids
⚠ Unmatched line format: 195.239 - 201.08: and this will be Game Changer uh for
⚠ Unmatched line format: 197.68 - 204.56: rest of the world as well so sir I think
⚠ Unmatched line format: 201.08 - 206.36: um uh let's go with this example of
⚠ Unmatched line format: 204.56 - 209.799: let's say I started off living in a
⚠ Unmatched line format: 206.36 - 212.04000000000002: small 2bhk apartment M and I slowly
⚠ Unmatched line format: 209.799 - 214.15900000000002: developed my kitchen to add to my taste
⚠ Unmatched line format: 212.04 - 216.76: but in a way that it's convenient for
⚠ Unmatched line format: 214.159 - 219.39999999999998: me right and one fine day you came and
⚠ Unmatched line format: 216.76 - 221.879: you changed my kitchen completely I was
⚠ Unmatched line format: 219.4 - 226.0: using a bicycle I moved to a scooter and
⚠ Unmatched line format: 221.879 - 228.2: a car and now you're asking me to fly
⚠ Unmatched line format: 226.0 - 230.879: and you're giving me an aircraft NEP
⚠ Unmatched line format: 228.2 - 232.83999999999997: sounds more like that for me how do I do
⚠ Unmatched line format: 230.879 - 234.92: the transition I fear that I will I will
⚠ Unmatched line format: 232.84 - 236.239: crash if I use a aeroplane without
⚠ Unmatched line format: 234.92 - 238.35999999999999: training I'm talking about all the
⚠ Unmatched line format: 236.239 - 241.12: teachers in the in the country all the
⚠ Unmatched line format: 238.36 - 244.0: schools in the country we have been
⚠ Unmatched line format: 241.12 - 245.59900000000002: driving buses at Max now we should fly
⚠ Unmatched line format: 244.0 - 250.0: how do we do this
⚠ Unmatched line format: 245.599 - 252.23899999999998: okay look at the requirement of the
⚠ Unmatched line format: 250.0 - 253.72: industry requirement of the industry has
⚠ Unmatched line format: 252.239 - 258.4: been moving very
⚠ Unmatched line format: 253.72 - 262.32: fast the technology is emerging every
⚠ Unmatched line format: 258.4 - 265.03999999999996: day and the industry is moving with that

1 from sentence_transformers import SentenceTransformer
2 from sklearn.cluster import KMeans
3 import re
4
5 # Load pre-trained model for sentence embeddings
6 model = SentenceTransformer('all-MiniLM-L6-v2')
7
8 # Function to extract sentences and starting timestamps
9 def extract_sentences_and_timestamps(file_path):
10     sentences = []
11     timestamps = []
12     # Enhanced regex pattern for range-based timestamps
13     pattern = re.compile(r"(\d{1,2}(\?:\.\d+)?)\s*-\s*(\d{1,2}(\?:\.\d+)?):\s*(.*)")
14
15     try:
16         with open(file_path, 'r') as file:
17             for line in file:
18                 line = line.strip()
19                 if line:
20                     match = pattern.match(line)
21                     if match:
22                         start_time = float(match.group(1))
23                         # Convert start time to hh:mm:ss format
24                         hours = int(start_time // 3600)
25                         minutes = int((start_time % 3600) // 60)
26                         seconds = start_time % 60
27                         timestamp = f"{hours:02}:{minutes:02}:{seconds:05.2f}"
28                         timestamps.append(timestamp)
29                         sentences.append(match.group(3).strip())
30                     else:
31                         print(f"⚠ Unmatched line format: {line}") # Debug print for unmatched lines
32     except FileNotFoundError:
33         print(f"❌ Error: The file '{file_path}' was not found.")
34         return [], []
35
36 if not sentences:
37     print("⚠ Warning: No sentences found. Check if the file is empty or if the format is incorrect.")
38 return sentences, timestamps
39
40 # Extract sentences and timestamps
41 input_file = 'transcript.txt'
42 sentences, timestamps = extract_sentences_and_timestamps(input_file)
43
44 # Debug: Print extracted sentences and timestamps
45 print("\nExtracted Sentences:", sentences)
46 print("Extracted Timestamps:", timestamps)
47
48 # Check if sentences are empty
49 if not sentences:
50     print("❌ Error: No valid sentences were extracted. Exiting the script.")
51 else:
52     # Generate sentence embeddings
53     embeddings = model.encode(sentences)
54
55     # Debug: Check if embeddings are generated correctly

```



```

56 if embeddings.size == 0:
57     print("❌ Error: No embeddings were generated. Check if the input text is valid.")
58 else:
59     # Determine number of clusters based on sentence count
60     # Here, we aim for topic-based segmentation, so we adjust clusters based on transcript length
61     num_clusters = max(1, len(sentences) // 5) # Example: 1 cluster per 5 sentences
62     kmeans = KMeans(n_clusters=num_clusters, random_state=0)
63     labels = kmeans.fit_predict(embeddings)
64
65     # Create segmented output with timestamps
66     segments = [[] for _ in range(num_clusters)]
67     for i, label in enumerate(labels):
68         segments[label].append(f"[{timestamps[i]}] {sentences[i]}")
69
70     # Save segmented output to a new file
71     output_file = 'segmented_transcript.txt'
72     with open(output_file, 'w') as file:
73         for idx, segment in enumerate(segments):
74             file.write(f"--- Topic Segment {idx + 1} ---\n")
75             file.write("\n".join(segment))
76             file.write("\n\n")
77
78     print(f"✅ Topic-segmented transcript saved to {output_file}")

```



⚠️ Unmatched line format: 96.399 - 102.32000000000001: the requirement of the
 ⚠️ Unmatched line format: 99.119 - 104.2: industry requirement of Manpower for the
 ⚠️ Unmatched line format: 102.32 - 106.19999999999999: industry however there were no
 ⚠️ Unmatched line format: 104.2 - 107.479: corresponding changes which happened
 ⚠️ Unmatched line format: 106.2 - 111.0: education
 ⚠️ Unmatched line format: 107.479 - 113.32: system so therefore I feel that this was
 ⚠️ Unmatched line format: 111.0 - 116.52: the right time when we brought
 ⚠️ Unmatched line format: 113.32 - 121.79899999999999: in the education policy 2020 honorable
 ⚠️ Unmatched line format: 116.52 - 125.19999999999999: prime minister declared it open in July
 ⚠️ Unmatched line format: 121.799 - 129.399: 2020 and we recently celebrated the
 ⚠️ Unmatched line format: 125.2 - 133.4: fourth anniversary of NEP 2020 ncrf has
 ⚠️ Unmatched line format: 129.399 - 136.8: been brought to implement the intent of
 ⚠️ Unmatched line format: 133.4 - 139.36: 2020 2020 is a policy the NP 2020 is a
 ⚠️ Unmatched line format: 136.8 - 141.44: policy and for implementing a policy you
 ⚠️ Unmatched line format: 139.36 - 144.48000000000002: need a
 ⚠️ Unmatched line format: 141.44 - 149.04: framework now why we call it a framework
 ⚠️ Unmatched line format: 144.48 - 151.16: we call it a framework because this is
 ⚠️ Unmatched line format: 149.04 - 154.35999999999999: very flexible
 ⚠️ Unmatched line format: 151.16 - 159.4: this allows you all the
 ⚠️ Unmatched line format: 154.36 - 161.84: Innovation the way you educate your kids
 ⚠️ Unmatched line format: 159.4 - 165.44: you educate your
 ⚠️ Unmatched line format: 161.84 - 168.48: students still it provides you the basic
 ⚠️ Unmatched line format: 165.44 - 171.959: guidelines the the framework the outer
 ⚠️ Unmatched line format: 168.48 - 174.84: layer it provides and that layer is
 ⚠️ Unmatched line format: 171.959 - 177.84: mostly the enabling
 ⚠️ Unmatched line format: 174.84 - 180.76: layer that is such an enabling layer
 ⚠️ Unmatched line format: 177.84 - 183.48: that it has
 ⚠️ Unmatched line format: 180.76 - 185.879: broken the shackles which were there in
 ⚠️ Unmatched line format: 183.48 - 189.48: the education
 ⚠️ Unmatched line format: 185.879 - 191.959: sector right so yes it is a big
 ⚠️ Unmatched line format: 189.48 - 195.23899999999998: Revolution and this is going to change
 ⚠️ Unmatched line format: 191.959 - 197.68: the way we have been educating our kids
 ⚠️ Unmatched line format: 195.239 - 201.08: and this will be Game Changer uh for
 ⚠️ Unmatched line format: 197.68 - 204.56: rest of the world as well so sir I think
 ⚠️ Unmatched line format: 201.08 - 206.36: um uh let's go with this example of
 ⚠️ Unmatched line format: 204.56 - 209.799: let's say I started off living in a
 ⚠️ Unmatched line format: 206.36 - 212.04000000000002: small 2bhk apartment M and I slowly
 ⚠️ Unmatched line format: 209.799 - 214.15900000000002: developed my kitchen to add to my taste
 ⚠️ Unmatched line format: 212.04 - 216.76: but in a way that it's convenient for
 ⚠️ Unmatched line format: 214.159 - 219.39999999999998: me right and one fine day you came and
 ⚠️ Unmatched line format: 216.76 - 221.879: you changed my kitchen completely I was
 ⚠️ Unmatched line format: 219.4 - 226.0: using a bicycle I moved to a scooter and
 ⚠️ Unmatched line format: 221.879 - 228.2: a car and now you're asking me to fly
 ⚠️ Unmatched line format: 226.0 - 230.879: and you're giving me an aircraft NEP
 ⚠️ Unmatched line format: 228.2 - 232.83999999999997: sounds more like that for me how do I do
 ⚠️ Unmatched line format: 230.879 - 234.92: the transition I fear that I will I will
 ⚠️ Unmatched line format: 232.84 - 236.239: crash if I use a aeroplane without
 ⚠️ Unmatched line format: 234.92 - 238.35999999999999: training I'm talking about all the
 ⚠️ Unmatched line format: 236.239 - 241.12: teachers in the in the country all the
 ⚠️ Unmatched line format: 238.36 - 244.0: schools in the country we have been
 ⚠️ Unmatched line format: 241.12 - 245.59900000000002: driving buses at Max now we should fly
 ⚠️ Unmatched line format: 244.0 - 250.0: how do we do this
 ⚠️ Unmatched line format: 245.599 - 252.23899999999998: okay look at the requirement of the
 ⚠️ Unmatched line format: 250.0 - 253.72: industry requirement of the industry has
 ⚠️ Unmatched line format: 252.239 - 258.4: been moving very
 ⚠️ Unmatched line format: 253.72 - 262.32: fast the technology is emerging every
 ⚠️ Unmatched line format: 258.4 - 265.03999999999996: day and the industry is moving with that

```

1 from sentence_transformers import SentenceTransformer
2 from sklearn.cluster import KMeans
3 import re
4

```

```

5 # Load pre-trained model for sentence embeddings
6 model = SentenceTransformer('all-MiniLM-L6-v2')
7
8 # Function to extract sentences and starting timestamps
9 def extract_sentences_and_timestamps(file_path):
10     sentences = []
11     timestamps = []
12     # Enhanced regex pattern for range-based timestamps
13     pattern = re.compile(r"(\d{1,2}(\?:\.\d+)?)\s*-\s*(\d{1,2}(\?:\.\d+)?):\s*(.*)")
14
15     try:
16         with open(file_path, 'r') as file:
17             for line in file:
18                 line = line.strip()
19                 if line:
20                     match = pattern.match(line)
21                     if match:
22                         start_time = float(match.group(1))
23                         # Convert start time to hh:mm:ss format
24                         hours = int(start_time // 3600)
25                         minutes = int((start_time % 3600) // 60)
26                         seconds = start_time % 60
27                         timestamp = f"{hours:02}:{minutes:02}:{seconds:05.2f}"
28                         timestamps.append(timestamp)
29                         sentences.append(match.group(3).strip())
30                     else:
31                         print(f"⚠ Unmatched line format: {line}") # Debug print for unmatched lines
32     except FileNotFoundError:
33         print(f"❌ Error: The file '{file_path}' was not found.")
34         return [], []
35
36     if not sentences:
37         print(f"⚠ Warning: No sentences found. Check if the file is empty or if the format is incorrect.")
38     return sentences, timestamps
39
40 # Extract sentences and timestamps
41 input_file = 'transcript.txt'
42 sentences, timestamps = extract_sentences_and_timestamps(input_file)
43
44 # Debug: Print extracted sentences and timestamps
45 print("\nExtracted Sentences:", sentences)
46 print("Extracted Timestamps:", timestamps)
47
48 # Check if sentences are empty
49 if not sentences:
50     print(f"❌ Error: No valid sentences were extracted. Exiting the script.")
51 else:
52     # Generate sentence embeddings
53     embeddings = model.encode(sentences)
54
55     # Debug: Check if embeddings are generated correctly
56     if embeddings.size == 0:
57         print(f"❌ Error: No embeddings were generated. Check if the input text is valid.")
58     else:
59         # Determine number of clusters dynamically based on conceptual shifts
60         num_clusters = max(1, len(sentences) // 5) # Example: 1 cluster per 5 sentences
61         kmeans = KMeans(n_clusters=num_clusters, random_state=0)
62         labels = kmeans.fit_predict(embeddings)
63
64         # Create segmented output with timestamps
65         segments = [[] for _ in range(num_clusters)]
66         for i, label in enumerate(labels):
67             segments[label].append(f"[{timestamps[i]}] {sentences[i]}")
68
69         # Save segmented output to a new file
70         output_file = 'segmented_transcript.txt'
71         with open(output_file, 'w') as file:
72             for idx, segment in enumerate(segments):
73                 file.write(f"--- Conceptually Relevant Segment {idx + 1} ---\n")
74                 file.write("\n".join(segment))
75                 file.write("\n\n")
76
77         print(f"✅ Conceptually relevant segmented transcript saved to {output_file}")

```



- ⚠ Unmatched line format: 96.399 - 102.32000000000001: the requirement of the
- ⚠ Unmatched line format: 99.119 - 104.2: industry requirement of Manpower for the
- ⚠ Unmatched line format: 102.32 - 106.19999999999999: industry however there were no
- ⚠ Unmatched line format: 104.2 - 107.479: corresponding changes which happened
- ⚠ Unmatched line format: 106.2 - 111.0: education
- ⚠ Unmatched line format: 107.479 - 113.32: system so therefore I feel that this was
- ⚠ Unmatched line format: 111.0 - 116.52: the right time when we brought
- ⚠ Unmatched line format: 113.32 - 121.79899999999999: in the education policy 2020 honorable
- ⚠ Unmatched line format: 116.52 - 125.19999999999999: prime minister declared it open in July

⚠ Unmatched line format: 121.799 - 129.399: 2020 and we recently celebrated the
 ⚠ Unmatched line format: 125.2 - 133.4: fourth anniversary of NEP 2020 ncrf has
 ⚠ Unmatched line format: 129.399 - 136.8: been brought to implement the intent of
 ⚠ Unmatched line format: 133.4 - 139.36: 2020 2020 is a policy the NP 2020 is a
 ⚠ Unmatched line format: 136.8 - 141.44: policy and for implementing a policy you
 ⚠ Unmatched line format: 139.36 - 144.48000000000002: need a
 ⚠ Unmatched line format: 141.44 - 149.04: framework now why we call it a framework
 ⚠ Unmatched line format: 144.48 - 151.16: we call it a framework because this is
 ⚠ Unmatched line format: 149.04 - 154.35999999999999: very flexible
 ⚠ Unmatched line format: 151.16 - 159.4: this allows you all the
 ⚠ Unmatched line format: 154.36 - 161.84: Innovation the way you educate your kids
 ⚠ Unmatched line format: 159.4 - 165.44: you educate your
 ⚠ Unmatched line format: 161.84 - 168.48: students still it provides you the basic
 ⚠ Unmatched line format: 165.44 - 171.959: guidelines the the framework the outer
 ⚠ Unmatched line format: 168.48 - 174.84: layer it provides and that layer is
 ⚠ Unmatched line format: 171.959 - 177.84: mostly the enabling
 ⚠ Unmatched line format: 174.84 - 180.76: layer that is such an enabling layer
 ⚠ Unmatched line format: 177.84 - 183.48: that it has
 ⚠ Unmatched line format: 180.76 - 185.879: broken the shackles which were there in
 ⚠ Unmatched line format: 183.48 - 189.48: the education
 ⚠ Unmatched line format: 185.879 - 191.959: sector right so yes it is a big
 ⚠ Unmatched line format: 189.48 - 195.23899999999998: Revolution and this is going to change
 ⚠ Unmatched line format: 191.959 - 197.68: the way we have been educating our kids
 ⚠ Unmatched line format: 195.239 - 201.08: and this will be Game Changer uh for
 ⚠ Unmatched line format: 197.68 - 204.56: rest of the world as well so sir I think
 ⚠ Unmatched line format: 201.08 - 206.36: um uh let's go with this example of
 ⚠ Unmatched line format: 204.56 - 209.799: let's say I started off living in a
 ⚠ Unmatched line format: 206.36 - 212.04000000000002: small 2bhk apartment M and I slowly
 ⚠ Unmatched line format: 209.799 - 214.15900000000002: developed my kitchen to add to my taste
 ⚠ Unmatched line format: 212.04 - 216.76: but in a way that it's convenient for
 ⚠ Unmatched line format: 214.159 - 219.39999999999998: me right and one fine day you came and
 ⚠ Unmatched line format: 216.76 - 221.879: you changed my kitchen completely I was
 ⚠ Unmatched line format: 219.4 - 226.0: using a bicycle I moved to a scooter and
 ⚠ Unmatched line format: 221.879 - 228.2: a car and now you're asking me to fly
 ⚠ Unmatched line format: 226.0 - 230.879: and you're giving me an aircraft NEP
 ⚠ Unmatched line format: 228.2 - 232.83999999999997: sounds more like that for me how do I do
 ⚠ Unmatched line format: 230.879 - 234.92: the transition I fear that I will I will
 ⚠ Unmatched line format: 232.84 - 236.239: crash if I use a aeroplane without
 ⚠ Unmatched line format: 234.92 - 238.35999999999999: training I'm talking about all the
 ⚠ Unmatched line format: 236.239 - 241.12: teachers in the in the country all the
 ⚠ Unmatched line format: 238.36 - 244.0: schools in the country we have been
 ⚠ Unmatched line format: 241.12 - 245.59900000000002: driving buses at Max now we should fly
 ⚠ Unmatched line format: 244.0 - 250.0: how do we do this
 ⚠ Unmatched line format: 245.599 - 252.23899999999998: okay look at the requirement of the
 ⚠ Unmatched line format: 250.0 - 253.72: industry requirement of the industry has
 ⚠ Unmatched line format: 252.239 - 258.4: been moving very
 ⚠ Unmatched line format: 253.72 - 262.32: fast the technology is emerging every

```
1 pip install nltk
```

```

➡ Requirement already satisfied: nltk in /usr/local/lib/python3.11/dist-packages (3.9.1)
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: tqdm in /usr/local/lib/python3.11/dist-packages (from nltk) (4.67.1)

```

```
1 pip install texttiling_segmenter.py
```

```

➡ ERROR: Could not find a version that satisfies the requirement texttiling_segmenter.py (from versions: none)
ERROR: No matching distribution found for texttiling_segmenter.py

```

```


1 import nltk
2 from nltk.corpus import stopwords
3 from nltk.tokenize import sent_tokenize, word_tokenize
4 from nltk.tokenize.texttiling import TextTilingTokenizer
5 from google.colab import files # For file upload
6
7 # Ensure you have the necessary NLTK data
8 nltk.download('punkt')
9 nltk.download('stopwords')
10
11 def preprocess_text(text):
12     """
13     Preprocess the text by tokenizing, lowercasing, and removing stopwords.
14     """
15     stop_words = set(stopwords.words('english'))
16     sentences = sent_tokenize(text)
17     words = [word_tokenize(sentence.lower()) for sentence in sentences]
18     words = [[word for word in sentence if word.isalnum() and word not in stop_words] for sentence in words]
19     return sentences, words
20
21 def segment_transcript(file_path):
22     """
23     Read a transcript from a file, preprocess it, and apply TextTiling to segment it.

```

```

24 """
25 # Read the file
26 with open(file_path, 'r', encoding='utf-8') as file:
27     transcript = file.read()
28
29 # Preprocess the text
30 sentences, words = preprocess_text(transcript)
31
32 # Apply TextTiling
33 tt = TextTilingTokenizer()
34 segments = tt.tokenize(transcript)
35
36 # Print the segments
37 for i, segment in enumerate(segments):
38     print(f"Segment {i+1}: \n{segment}\n")
39
40 # Upload the file to Colab
41 uploaded = files.upload()
42
43 # Get the file name
44 file_name = list(uploaded.keys())[0]
45
46 # Segment the transcript
47 segment_transcript(file_name)

```

 [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!
 No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving Kalsi sir transcript.txt to Kalsi sir transcript (1).txt

```

-----
StopIteration                                Traceback (most recent call last)
/usr/local/lib/python3.11/dist-packages/nltk/tokenize/texttiling.py in _create_token_table(self, token_sequences, par_breaks)
    242         try:
--> 243             current_par_break = next(pb_iter) # skip break at 0
    244         except StopIteration as e:

```

StopIteration:

The above exception was the direct cause of the following exception:

```

-----
ValueError                                Traceback (most recent call last)
-----
3 frames
/usr/local/lib/python3.11/dist-packages/nltk/tokenize/texttiling.py in _create_token_table(self, token_sequences, par_breaks)
    243             current_par_break = next(pb_iter) # skip break at 0
    244         except StopIteration as e:
--> 245             raise ValueError(
    246                 "No paragraph breaks were found(text too short perhaps?)"
    247             ) from e


```

ValueError: No paragraph breaks were found(text too short perhaps?)

```

1 import nltk
2 nltk.download('punkt')
3 nltk.download('punkt_tab')

```

 [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.
 True

```


1 import nltk
2 from nltk.corpus import stopwords
3 from nltk.tokenize import sent_tokenize, word_tokenize
4 from nltk.tokenize.texttiling import TextTilingTokenizer
5 from google.colab import files
6
7 # Download required NLTK resources
8 nltk.download('punkt')
9 nltk.download('punkt_tab')
10 nltk.download('stopwords')
11
12 def preprocess_text(text):
13     """
14     Preprocess the text by tokenizing, lowercasing, and removing stopwords.
15     """
16     stop_words = set(stopwords.words('english'))

```

```

17 sentences = sent_tokenize(text)
18 words = [word_tokenize(sentence.lower()) for sentence in sentences]
19 words = [[word for word in sentence if word.isalnum() and word not in stop_words] for sentence in words]
20 return sentences, words
21
22 def insert_paragraph_breaks(text, sentences_per_paragraph=5):
23     """
24     Insert paragraph breaks after every `sentences_per_paragraph` sentences.
25     """
26     sentences = sent_tokenize(text)
27     paragraphs = []
28     for i in range(0, len(sentences), sentences_per_paragraph):
29         paragraph = " ".join(sentences[i:i + sentences_per_paragraph])
30         paragraphs.append(paragraph)
31     return "\n\n".join(paragraphs)
32
33 def segment_transcript(file_path):
34     """
35     Read a transcript from a file, preprocess it, and apply TextTiling to segment it.
36     """
37     # Read the file
38     with open(file_path, 'r', encoding='utf-8') as file:
39         transcript = file.read()
40
41     # Insert paragraph breaks
42     transcript_with_breaks = insert_paragraph_breaks(transcript)
43
44     # Preprocess the text
45     sentences, words = preprocess_text(transcript_with_breaks)
46
47     # Apply TextTiling
48     tt = TextTilingTokenizer()
49     segments = tt.tokenize(transcript_with_breaks)
50
51     # Print the segments
52     for i, segment in enumerate(segments):
53         print(f"Segment {i+1}:\n{segment}\n")
54
55 # Upload the file to Colab
56 uploaded = files.upload()
57
58 # Get the file name
59 file_name = list(uploaded.keys())[0]
60
61 # Segment the transcript
62 segment_transcript(file_name)

```

 [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] Package punkt_tab is already up-to-date!
 [nltk_data] Downloading package stopwords to /root/nltk_data...
 [nltk_data] Package stopwords is already up-to-date!
 No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving Kalsi sir transcript.txt to Kalsi sir transcript (2).txt

```

-----
StopIteration                                Traceback (most recent call last)
/usr/local/lib/python3.11/dist-packages/nltk/tokenize/texttiling.py in _create_token_table(self, token_sequences, par_breaks)
    242         try:
--> 243             current_par_break = next(pb_iter) # skip break at 0
    244         except StopIteration as e:

```

StopIteration:

The above exception was the direct cause of the following exception:

```

ValueError                                Traceback (most recent call last)
-----
3 frames
/usr/local/lib/python3.11/dist-packages/nltk/tokenize/texttiling.py in _create_token_table(self, token_sequences, par_breaks)
    243         current_par_break = next(pb_iter) # skip break at 0
    244     except StopIteration as e:
--> 245         raise ValueError(
    246             "No paragraph breaks were found(text too short perhaps?)"
    247         ) from e

```

ValueError: No paragraph breaks were found(text too short perhaps?)

```

1 import re
2 from sklearn.feature_extraction.text import CountVectorizer
3 from sklearn.decomposition import LatentDirichletAllocation

```

```

4 import numpy as np
5
6 def preprocess_text(text):
7     """
8     Preprocess the text by splitting it into sentences and cleaning.
9     """
10    # Split into sentences
11    sentences = re.split(r'(?<=[.!?])\s+', text)
12    # Remove short sentences (optional)
13    sentences = [s.strip() for s in sentences if len(s.split()) > 5]
14    return sentences
15
16 def identify_topics(sentences, num_topics=5):
17     """
18     Identify topics in the text using Latent Dirichlet Allocation (LDA).
19     """
20    # Create a document-term matrix
21    vectorizer = CountVectorizer(stop_words='english', max_df=0.95, min_df=2)
22    doc_term_matrix = vectorizer.fit_transform(sentences)
23
24    # Fit LDA model
25    lda = LatentDirichletAllocation(n_components=num_topics, random_state=42)
26    lda.fit(doc_term_matrix)
27
28    # Get topic distribution for each sentence
29    topic_distribution = lda.transform(doc_term_matrix)
30    topic_labels = np.argmax(topic_distribution, axis=1)
31
32    # Get top words for each topic
33    feature_names = vectorizer.get_feature_names_out()
34    topics = []
35    for topic_idx, topic in enumerate(lda.components_):
36        top_words = [feature_names[i] for i in topic.argsort()[-5:]]
37        topics.append(f"Topic {topic_idx + 1}: {' '.join(top_words)}")
38
39    return topic_labels, topics
40
41 def segment_transcript(input_file, output_file, num_topics=5):
42     """
43     Read a transcript from a text file, analyze it, and generate segmented output.
44     """
45    # Read the input file
46    with open(input_file, 'r', encoding='utf-8') as file:
47        transcript = file.read()
48
49    # Preprocess the text
50    sentences = preprocess_text(transcript)
51
52    # Identify topics
53    topic_labels, topics = identify_topics(sentences, num_topics)
54
55    # Group sentences by topic
56    segments = {}
57    for i, label in enumerate(topic_labels):
58        if label not in segments:
59            segments[label] = []
60        segments[label].append(sentences[i])
61
62    # Write the segmented output to a file
63    with open(output_file, 'w', encoding='utf-8') as file:
64        for label, segment_sentences in segments.items():
65            file.write(f"Segment {label + 1}: {topics[label]}\n")
66            file.write("\n".join(segment_sentences) + "\n\n")
67
68    # Input and output file paths
69    input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
70    output_file = "segmented_transcript.txt" # Output file path
71
72    # Segment the transcript
73    segment_transcript(input_file, output_file, num_topics=5)
74
75    print(f"Segmented transcript saved to {output_file}")

```

```

-----
ValueError                                Traceback (most recent call last)
<ipython-input-13-2d212f11d399> in <cell line: 0>()
    71
    72 # Segment the transcript
--> 73 segment_transcript(input_file, output_file, num_topics=5)
    74
    75 print(f"Segmented transcript saved to {output_file}")

-----
3 frames -----
/usr/local/lib/python3.11/dist-packages/sklearn/feature_extraction/text.py in fit_transform(self, raw_documents, y)
    1384         min_doc_count = min_df if isinstance(min_df, Integral) else min_df * n_doc
    1385         if max_doc_count < min_doc_count:
-> 1386             raise ValueError("max_df corresponds to < documents than min_df")
    1387         if max_features is not None:
    1388             X = self._sort_features(X, vocabulary)

ValueError: max_df corresponds to < documents than min_df

1 import re
2 from sklearn.feature_extraction.text import CountVectorizer
3 from sklearn.decomposition import LatentDirichletAllocation
4 import numpy as np
5
6 def preprocess_text(text):
7     """
8     Preprocess the text by splitting it into sentences and cleaning.
9     """
10    # Split into sentences
11    sentences = re.split(r'(?<=[.!?])\s+', text)
12    # Remove short sentences (optional)
13    sentences = [s.strip() for s in sentences if len(s.split()) > 5]
14    return sentences
15
16 def identify_topics(sentences, num_topics=5):
17     """
18     Identify topics in the text using Latent Dirichlet Allocation (LDA).
19     """
20    # Create a document-term matrix
21    vectorizer = CountVectorizer(stop_words='english', max_df=1.0, min_df=1)
22    doc_term_matrix = vectorizer.fit_transform(sentences)
23
24    # Fit LDA model
25    lda = LatentDirichletAllocation(n_components=num_topics, random_state=42)
26    lda.fit(doc_term_matrix)
27
28    # Get topic distribution for each sentence
29    topic_distribution = lda.transform(doc_term_matrix)
30    topic_labels = np.argmax(topic_distribution, axis=1)
31
32    # Get top words for each topic
33    feature_names = vectorizer.get_feature_names_out()
34    topics = []
35    for topic_idx, topic in enumerate(lda.components_):
36        top_words = [feature_names[i] for i in topic.argsort()[-5:]]
37        topics.append(f"Topic {topic_idx + 1}: {' '.join(top_words)}")
38
39    return topic_labels, topics
40
41 def segment_transcript(input_file, output_file, num_topics=5):
42     """
43     Read a transcript from a text file, analyze it, and generate segmented output.
44     """
45    # Read the input file
46    with open(input_file, 'r', encoding='utf-8') as file:
47        transcript = file.read()
48
49    # Preprocess the text
50    sentences = preprocess_text(transcript)
51
52    # Check if there are enough sentences for topic modeling
53    if len(sentences) < num_topics:
54        raise ValueError(f"Not enough sentences ({len(sentences)}) for {num_topics} topics.")
55
56    # Identify topics
57    topic_labels, topics = identify_topics(sentences, num_topics)
58
59    # Group sentences by topic
60    segments = {}
61    for i, label in enumerate(topic_labels):

```

```

62     if label not in segments:
63         segments[label] = []
64     segments[label].append(sentences[i])
65
66     # Write the segmented output to a file
67     with open(output_file, 'w', encoding='utf-8') as file:
68         for label, segment_sentences in segments.items():
69             file.write(f"Segment {label + 1}: {topics[label]}\n")
70             file.write("\n".join(segment_sentences) + "\n\n")
71
72 # Input and output file paths
73 input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
74 output_file = "segmented_transcript.txt" # Output file path
75
76 # Segment the transcript
77 segment_transcript(input_file, output_file, num_topics=5)
78
79 print(f"Segmented transcript saved to {output_file}")

```



```

-----
ValueError                                Traceback (most recent call last)
<ipython-input-14-39c7be34e700> in <cell line: 0>()
    75
    76 # Segment the transcript
--> 77 segment_transcript(input_file, output_file, num_topics=5)
    78
    79 print(f"Segmented transcript saved to {output_file}")

<ipython-input-14-39c7be34e700> in segment_transcript(input_file, output_file, num_topics)
    52     # Check if there are enough sentences for topic modeling
    53     if len(sentences) < num_topics:
--> 54         raise ValueError(f"Not enough sentences ({len(sentences)}) for {num_topics} topics.")
    55
    56     # Identify topics

ValueError: Not enough sentences (2) for 5 topics.

```

```

1 import re
2 from sklearn.feature_extraction.text import CountVectorizer
3 from sklearn.decomposition import LatentDirichletAllocation
4 import numpy as np
5
6 def preprocess_text(text):
7     """
8     Preprocess the text by splitting it into sentences and cleaning.
9     """
10    # Split into sentences
11    sentences = re.split(r'(?<=[.!?])\s+', text)
12    # Remove short sentences (optional)
13    sentences = [s.strip() for s in sentences if len(s.split()) > 5]
14    return sentences
15
16 def identify_topics(sentences, num_topics=2):
17     """
18     Identify topics in the text using Latent Dirichlet Allocation (LDA).
19     """
20    # Create a document-term matrix
21    vectorizer = CountVectorizer(stop_words='english', max_df=1.0, min_df=1)
22    doc_term_matrix = vectorizer.fit_transform(sentences)
23
24    # Fit LDA model
25    lda = LatentDirichletAllocation(n_components=num_topics, random_state=42)
26    lda.fit(doc_term_matrix)
27
28    # Get topic distribution for each sentence
29    topic_distribution = lda.transform(doc_term_matrix)
30    topic_labels = np.argmax(topic_distribution, axis=1)
31
32    # Get top words for each topic
33    feature_names = vectorizer.get_feature_names_out()
34    topics = []
35    for topic_idx, topic in enumerate(lda.components_):
36        top_words = [feature_names[i] for i in topic.argsort()[-5:]]
37        topics.append(f"Topic {topic_idx + 1}: {'', ' '.join(top_words)}")
38
39    return topic_labels, topics
40
41 def segment_transcript(input_file, output_file):
42     """
43     Read a transcript from a text file, analyze it, and generate segmented output.
44     """
45    # Read the input file

```



```

46 with open(input_file, 'r', encoding='utf-8') as file:
47     transcript = file.read()
48
49 # Preprocess the text
50 sentences = preprocess_text(transcript)
51
52 # Dynamically adjust the number of topics
53 num_topics = min(len(sentences), 5) # Use up to 5 topics or the number of sentences, whichever is smaller
54 if num_topics < 1:
55     raise ValueError("Not enough sentences for topic modeling.")
56
57 print(f"Using {num_topics} topics for segmentation.")
58
59 # Identify topics
60 topic_labels, topics = identify_topics(sentences, num_topics)
61
62 # Group sentences by topic
63 segments = {}
64 for i, label in enumerate(topic_labels):
65     if label not in segments:
66         segments[label] = []
67     segments[label].append(sentences[i])
68
69 # Write the segmented output to a file
70 with open(output_file, 'w', encoding='utf-8') as file:
71     for label, segment_sentences in segments.items():
72         file.write(f"Segment {label + 1}: {topics[label]}\n")
73         file.write("\n".join(segment_sentences) + "\n\n")
74
75 # Input and output file paths
76 input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
77 output_file = "segmented_transcript.txt" # Output file path
78
79 # Segment the transcript
80 segment_transcript(input_file, output_file)
81
82 print(f"Segmented transcript saved to {output_file}")

```

→ Using 2 topics for segmentation.
Segmented transcript saved to segmented_transcript.txt

```

1 import re
2 from nltk.tokenize.texttiling import TextTilingTokenizer
3 from nltk.corpus import stopwords
4 from nltk.tokenize import sent_tokenize
5 from nltk import download
6
7 # Download necessary NLTK data
8 download('punkt')
9 download('stopwords')
10
11 def preprocess_text(text):
12     """
13     Preprocess the text by removing stopwords and tokenizing into sentences.
14     """
15     stop_words = set(stopwords.words('english'))
16     sentences = sent_tokenize(text)
17     cleaned_sentences = []
18     for sentence in sentences:
19         words = [word.lower() for word in re.findall(r'\b\w+\b', sentence) if word.lower() not in stop_words]
20         cleaned_sentences.append(" ".join(words))
21     return cleaned_sentences
22
23 def segment_transcript(input_file, output_file):
24     """
25     Segment the transcript using the TextTiling algorithm.
26     """
27     # Read the input file
28     with open(input_file, 'r', encoding='utf-8') as file:
29         transcript = file.read()
30
31     # Preprocess the text
32     sentences = preprocess_text(transcript)
33
34     # Combine sentences into a single string for TextTiling
35     text = "\n\n".join(sentences)
36
37     # Initialize TextTilingTokenizer
38     tt = TextTilingTokenizer()
39
40     # Segment the text
41     try:

```

```

42     segments = tt.tokenize(text)
43 except ValueError as e:
44     print(f"TextTiling failed: {e}")
45     print("Falling back to sentence-based segmentation.")
46     segments = ["\n".join(sentences)]
47
48 # Write the segmented output to a file
49 with open(output_file, 'w', encoding='utf-8') as file:
50     for i, segment in enumerate(segments):
51         file.write(f"Segment {i + 1}:\n{segment}\n\n")
52
53 # Input and output file paths
54 input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
55 output_file = "segmented_transcript.txt" # Output file path
56
57 # Segment the transcript
58 segment_transcript(input_file, output_file)
59
60 print(f"Segmented transcript saved to {output_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!
Segmented transcript saved to segmented_transcript.txt

```

```

1 import re
2 from nltk.tokenize.texttiling import TextTilingTokenizer
3 from nltk.corpus import stopwords
4 from nltk.tokenize import sent_tokenize, word_tokenize
5 from sklearn.feature_extraction.text import TfidfVectorizer
6 from nltk import download
7
8 # Download necessary NLTK data
9 download('punkt')
10 download('stopwords')
11
12 def preprocess_text(text):
13     """
14     Preprocess the text by removing stopwords and tokenizing into sentences.
15     """
16     stop_words = set(stopwords.words('english'))
17     sentences = sent_tokenize(text)
18     cleaned_sentences = []
19     for sentence in sentences:
20         words = [word.lower() for word in word_tokenize(sentence) if word.isalnum() and word.lower() not in stop_words]
21         cleaned_sentences.append(" ".join(words))
22     return cleaned_sentences
23
24 def extract_topic(segment, top_n=3):
25     """
26     Extract the top keywords from a segment to use as the topic label.
27     """
28     vectorizer = TfidfVectorizer(max_features=10)
29     tfidf_matrix = vectorizer.fit_transform([segment])
30     feature_names = vectorizer.get_feature_names_out()
31     tfidf_scores = tfidf_matrix.toarray()[0]
32     top_keywords = [feature_names[i] for i in tfidf_scores.argsort()[-top_n:][::-1]]
33     return " ".join(top_keywords)
34
35 def segment_transcript(input_file, output_file):
36     """
37     Segment the transcript using the TextTiling algorithm and assign topic labels.
38     """
39     # Read the input file
40     with open(input_file, 'r', encoding='utf-8') as file:
41         transcript = file.read()
42
43     # Preprocess the text
44     sentences = preprocess_text(transcript)
45
46     # Combine sentences into a single string for TextTiling
47     text = "\n\n".join(sentences)
48
49     # Initialize TextTilingTokenizer
50     tt = TextTilingTokenizer()
51
52     # Segment the text
53     try:
54         segments = tt.tokenize(text)
55     except ValueError as e:
56         print(f"TextTiling failed: {e}")

```

```

57     print("Falling back to sentence-based segmentation.")
58     segments = ["\n".join(sentences)]
59
60     # Write the segmented output to a file with topic labels
61     with open(output_file, 'w', encoding='utf-8') as file:
62         for i, segment in enumerate(segments):
63             # Extract the topic for the segment
64             topic = extract_topic(segment)
65             file.write(f"Segment {i + 1} - Topic: {topic}\n")
66             file.write(f"{segment}\n\n")
67
68 # Input and output file paths
69 input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
70 output_file = "segmented_transcript.txt" # Output file path
71
72 # Segment the transcript
73 segment_transcript(input_file, output_file)
74
75 print(f"Segmented transcript saved to {output_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!
Segmented transcript saved to segmented_transcript.txt

```

```

1 import nltk
2 from nltk.corpus import stopwords
3 from nltk.tokenize import word_tokenize
4 from nltk.stem import WordNetLemmatizer
5 import string
6
7 nltk.download('punkt')
8 nltk.download('punkt_tab')
9 nltk.download('stopwords')
10 nltk.download('wordnet')
11 nltk.download('omw-1.4')
12
13 def preprocess(text):
14     # Tokenize
15     tokens = word_tokenize(text.lower())
16
17     # Remove punctuation and stopwords
18     tokens = [word for word in tokens if word.isalnum() and word not in stopwords.words('english')]
19
20     # Lemmatize
21     lemmatizer = WordNetLemmatizer()
22     tokens = [lemmatizer.lemmatize(word) for word in tokens]
23
24     return tokens
25
26 # Example transcript
27 transcript = """
28 Your transcript text goes here. It can be multiple paragraphs long.
29 This is just an example to show how the preprocessing works.
30 """
31
32 # Preprocess the transcript
33 processed_transcript = preprocess(transcript)

```

```

[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] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data] Package omw-1.4 is already up-to-date!

```

```

1 # Create a dictionary and corpus
2 dictionary = Dictionary([processed_transcript])
3 corpus = [dictionary.doc2bow(processed_transcript)]
4
5 # Train an LDA model
6 lda_model = LdaModel(corpus, num_topics=5, id2word=dictionary, passes=15)
7
8 # Tokenize the transcript into sentences
9 sentences = sent_tokenize(transcript)
10
11 # Segment the transcript based on topics

```

```

12 segments = []
13 for sentence in sentences:
14     bow = dictionary.doc2bow(preprocess(sentence))
15     topic_distribution = lda_model.get_document_topics(bow)
16     dominant_topic = max(topic_distribution, key=lambda x: x[1])[0]
17     segments.append((sentence, dominant_topic))
18
19 # Write the segments to a text file
20 with open("segmented_transcript_lda.txt", "w") as output_file:
21     for sentence, topic in segments:
22         output_file.write(f"Topic {topic}: {sentence}\n\n")
23
24 print("Segmentation complete. Output saved to 'segmented_transcript_lda.txt'.")

1 from gensim.corpora import Dictionary # Import Dictionary from gensim.corpora

1 from gensim.models import LdaModel

1 from nltk.tokenize import word_tokenize, sent_tokenize

1 # Preprocessing function
2 def preprocess(text):
3     tokens = word_tokenize(text.lower())
4     tokens = [word for word in tokens if word.isalnum() and word not in stopwords.words('english')]
5     lemmatizer = WordNetLemmatizer()
6     tokens = [lemmatizer.lemmatize(word) for word in tokens]
7     return tokens
8
9 # Load the transcript from a file
10 with open("Kalsi sir transcript.txt", "r") as file:
11     transcript = file.read()
12
13 # Preprocess the transcript
14 processed_transcript = preprocess(transcript)
15
16 # Create a dictionary and corpus
17 dictionary = Dictionary([processed_transcript])
18 corpus = [dictionary.doc2bow(processed_transcript)]
19
20 # Train an LDA model
21 lda_model = LdaModel(corpus, num_topics=5, id2word=dictionary, passes=15)
22
23 # Tokenize the transcript into sentences
24 sentences = sent_tokenize(transcript)
25
26 # Segment the transcript based on topics
27 segments = []
28 for sentence in sentences:
29     bow = dictionary.doc2bow(preprocess(sentence))
30     topic_distribution = lda_model.get_document_topics(bow)
31     dominant_topic = max(topic_distribution, key=lambda x: x[1])[0]
32     segments.append((sentence, dominant_topic))
33
34 # Write the segments to a text file
35 with open("segmented_transcript_lda.txt", "w") as output_file:
36     for sentence, topic in segments:
37         output_file.write(f"Topic {topic}: {sentence}\n\n")
38
39 print("Segmentation complete. Output saved to 'segmented_transcript_lda.txt'.")

```

➞ Segmentation complete. Output saved to 'segmented_transcript_lda.txt'.

1 pip install bertopic

➞ Collecting bertopic

Downloading bertopic-0.16.4-py3-none-any.whl.metadata (23 kB)

Requirement already satisfied: hdbscan>=0.8.29 in /usr/local/lib/python3.11/dist-packages (from bertopic) (0.8.40)

Requirement already satisfied: numpy>=1.20.0 in /usr/local/lib/python3.11/dist-packages (from bertopic) (1.26.4)

Requirement already satisfied: pandas>=1.1.5 in /usr/local/lib/python3.11/dist-packages (from bertopic) (2.2.2)

Requirement already satisfied: plotly>=4.7.0 in /usr/local/lib/python3.11/dist-packages (from bertopic) (5.24.1)

Requirement already satisfied: scikit-learn>=0.22.2.post1 in /usr/local/lib/python3.11/dist-packages (from bertopic) (1.6.1)

Requirement already satisfied: sentence-transformers>=0.4.1 in /usr/local/lib/python3.11/dist-packages (from bertopic) (3.4.1)

Requirement already satisfied: tqdm>=4.41.1 in /usr/local/lib/python3.11/dist-packages (from bertopic) (4.67.1)

Requirement already satisfied: umap-learn>=0.5.0 in /usr/local/lib/python3.11/dist-packages (from bertopic) (0.5.7)

Requirement already satisfied: scipy>=1.0 in /usr/local/lib/python3.11/dist-packages (from hdbscan>=0.8.29->bertopic) (1.13.1)

Requirement already satisfied: joblib>=1.0 in /usr/local/lib/python3.11/dist-packages (from hdbscan>=0.8.29->bertopic) (1.4.2)

Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.1.5->bertopic) (2025.1)

Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.1.5->bertopic) (2025.1)

Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.1.5->bertopic) (2025.1)

Requirement already satisfied: tenacity>=6.2.0 in /usr/local/lib/python3.11/dist-packages (from plotly>=4.7.0->bertopic) (9.0.0)

```

Requirement already satisfied: packaging in /usr/local/lib/python3.11/dist-packages (from plotly>=4.7.0->bertopic) (24.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn>=0.22.2.post1->bertopic) (3.3.0)
Requirement already satisfied: transformers<5.0.0,>=4.41.0 in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.4.1->bertopic) (4.41.0)
Requirement already satisfied: torch>=1.11.0 in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.4.1->bertopic) (2.4.0)
Requirement already satisfied: huggingface-hub>=0.20.0 in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.4.1->bertopic) (0.20.0)
Requirement already satisfied: Pillow in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=0.4.1->bertopic) (10.4.0)
Requirement already satisfied: numba>=0.51.2 in /usr/local/lib/python3.11/dist-packages (from umap-learn>=0.5.0->bertopic) (0.60.0)
Requirement already satisfied: pynndescent>=0.5 in /usr/local/lib/python3.11/dist-packages (from umap-learn>=0.5.0->bertopic) (0.5.10)
Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers) (3.13.1)
Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers) (2024.10.0)
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) (2.32.0)
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) (4.12.0)
Requirement already satisfied: llvmlite<0.44,>=0.43.0dev0 in /usr/local/lib/python3.11/dist-packages (from numba>=0.51.2->umap-learn) (0.43.0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas>=1.1.5->bertopic) (1.16.0)
Requirement already satisfied: networkx in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic) (3.3)
Requirement already satisfied: Jinja2 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic) (3.1.4)
Collecting nvidia-cuda-nvrtc-cu12==12.4.127 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cuda_nvrtc_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cuda-runtime-cu12==12.4.127 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cuda_runtime_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cuda-cupti-cu12==12.4.127 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cuda_cupti_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cudnn-cu12==9.1.0.70 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cudnn_cu12-9.1.0.70-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cublas-cu12==12.4.5.8 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cublas_cu12-12.4.5.8-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cufft-cu12==11.2.1.3 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cufft_cu12-11.2.1.3-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-curand-cu12==10.3.5.147 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_curand_cu12-10.3.5.147-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-curand-cu12==10.3.5.147 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_curand_cu12-10.3.5.147-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cusolver-cu12==11.6.1.9 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cusolver_cu12-11.6.1.9-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cusparse-cu12==12.3.1.170 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_cusparse_cu12-12.3.1.170-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers) (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) (12.4.127)
Collecting nvidia-nvjitlink-cu12==12.4.127 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_nvjitlink_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Requirement already satisfied: triton==3.1.0 in /usr/local/lib/python3.11/dist-packages (from torch>=1.11.0->sentence-transformers) (3.1.0)

```

```

1 from bertopic import BertTopic
2
3 # Load the transcript from a text file
4 input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
5 with open(input_file, "r") as file:
6     transcript = file.read()
7
8 # Split the transcript into sentences
9 sentences = transcript.split('.') # Simple split for demonstration
10 sentences = [s.strip() for s in sentences if s.strip()]
11
12 # Train BertTopic model
13 topic_model = BertTopic()
14 topics, _ = topic_model.fit_transform(sentences)
15
16 # Get topic information
17 topic_info = topic_model.get_topic_info()
18
19 # Save the topics to a file
20 output_file = "bertopic_clustered_transcript.txt" # Replace with your desired output file path
21 with open(output_file, "w") as f:
22     for topic_id in topic_info['Topic']:
23         if topic_id != -1: # Ignore outliers
24             f.write(f"Topic {topic_id}:\n")
25             for sentence, assigned_topic in zip(sentences, topics):
26                 if assigned_topic == topic_id:
27                     f.write(f" - {sentence}\n")
28             f.write("\n")
29
30 print(f"Clustering complete. Output saved to '{output_file}'.")

```



modules.json: 100%	349/349 [00:00<00:00, 22.5kB/s]
config_sentence_transformers.json: 100%	116/116 [00:00<00:00, 9.46kB/s]
README.md: 100%	10.5k/10.5k [00:00<00:00, 864kB/s]
sentence_bert_config.json: 100%	53.0/53.0 [00:00<00:00, 4.68kB/s]
config.json: 100%	612/612 [00:00<00:00, 36.2kB/s]
model.safetensors: 100%	90.9M/90.9M [00:01<00:00, 105MB/s]
tokenizer_config.json: 100%	350/350 [00:00<00:00, 22.0kB/s]
vocab.txt: 100%	232k/232k [00:00<00:00, 4.22MB/s]
tokenizer.json: 100%	466k/466k [00:00<00:00, 17.6MB/s]
special_tokens_map.json: 100%	112/112 [00:00<00:00, 6.22kB/s]
config.json: 100%	190/190 [00:00<00:00, 8.95kB/s]

ValueError Traceback (most recent call last)

```
<ipython-input-15-12c776d8ef88> in <cell line: 0>()
    12 # Train BERTopic model
    13 topic_model = BERTopic()
--> 14 topics, _ = topic_model.fit_transform(sentences)
    15
    16 # Get topic information
```

⬆ 5 frames

```
/usr/local/lib/python3.11/dist-packages/numpy/core/_methods.py in _amax(a, axis, out, keepdims, initial, where)
    39 def _amax(a, axis=None, out=None, keepdims=False,
    40           initial=NoValue, where=True):
--> 41     return umr_maximum(a, axis, None, out, keepdims, initial, where)
    42
    43 def _amin(a, axis=None, out=None, keepdims=False,
```

ValueError: zero-size array to reduction operation maximum which has no identity

```
1 from bertopic import BERTopic
2 import nltk
3
4 # Download the Punkt tokenizer for sentence splitting
5 nltk.download('punkt')
6
7 # Load the transcript from a text file
8 input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
9 with open(input_file, "r") as file:
10     transcript = file.read()
11
12 # Split the transcript into sentences using NLTK's sent_tokenize
13 sentences = nltk.sent_tokenize(transcript) # More robust sentence splitting
14 sentences = [s.strip() for s in sentences if s.strip()] # Remove empty strings
15
16 # Debugging: Print the sentences to verify
17 print("Sentences extracted from the transcript:")
18 for i, sentence in enumerate(sentences):
19     print(f"{i + 1}: {sentence}")
20
21 # Check if sentences are empty
22 if not sentences:
23     raise ValueError("No sentences found in the transcript. Please check the input file.")
24
25 # Train BERTopic model
26 topic_model = BERTopic()
27 topics, _ = topic_model.fit_transform(sentences)
28
29 # Get topic information
30 topic_info = topic_model.get_topic_info()
31
32 # Save the topics to a file
33 output_file = "bertopic_clustered_transcript.txt" # Replace with your desired output file path
34 with open(output_file, "w") as f:
35     for topic_id in topic_info['Topic']:
36         if topic_id != -1: # Ignore outliers
37             f.write(f"Topic {topic_id}:\n")
38             for sentence, assigned_topic in zip(sentences, topics):
39                 if assigned_topic == topic_id:
40                     f.write(f" - {sentence}\n")
41             f.write("\n")
42
43 print(f"Clustering complete. Output saved to '{output_file}'.")
```

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

Sentences extracted from the transcript:

1: so sorry we know that India has seen a huge Revolution with digital payments we all thought that India is a place at least the I think immediately after you play the next big revolution personally I think is an education and the complete homework for this do you think ncrf plus any people together will be the next big revolution after up in India absolutely and why I think so is because how many changes have happened in the real world so many changes have happened in the requirement of the industry requirement of July 2020 and we recently celebrated the 4th anniversary of NP 2020 ncrf has been brought to implement the intent of 2020 2020 this allows you all the Innovation the way you educate your kids you educate your students still it provides you the basic guide: broken the shackles of which were there in the education sector right so yes it is a big Revolution and this is going to change I pull up my kitchen to add to my taste buds in a way that it's convenient for me and one fine day you came and you changed my kitchen we have been driving buses at Max now we should fly how do we do this look at the requirement of the industry requirement of the link that whatever he has been taught has no relevance to the real life world when industry is moving that fast when they require will you be not only beneficial but also very facilitated for all of us very liberating for all of us this is going to be highly new things create new ways of doing things learn something new but once we learn it there is no limit to Innovation and creativity creating our vision and Innovative Minds into the education sector where is we are applying it elsewhere everywhere else no we are to be fun for everyone and I can I can tell you that already a number of Institutions have adopted the any pain and CRF to vary the next step on how to implement an AP if I can request you to give me an elevator pitch for an EP and then an elevator pitch for not it allows for creditor of all learnings whether in academics or killing or an experiential learning and all these three kinds and people who are already skilled or already in the professional area not there the experiential learning would play a big part increase of technology which has been created by single Department good question so therefore all of us we have to work in tandem all kinds of learning are being contractors including learning of soft skills employability skills life skills your hand skills you go out it's all very flexible so that there's no Dropout there's no Dropout so these three things coupled with use of technology I want my son to be an engineer don't you think if we create a give me five approach to Credit Systems everybody will come to the turn off infrastructure India has today we created in last 1775 years we are going to double that infrastructure in next 10:11 years more number of other branches even liberal arts social sciences if I want to be an award Society I would need a proper mix of all have you already seen the Fallout of this know you can see the photo you can see how many Engineers are there for Designing our [Unintelligible]

any of the new technology machines current any laser-based machines any automated operating machines robotic operations robotic process and that is killing is equally important and this is important in multiple areas and therefore multidisciplinary me and you want to know design a VTech in CSC syllabus or be taking AI syllabus that is Nip complaint under the ncrf framework how are you teaching Teddy teaching Terry is not sufficient if you want the student to really understand and reply that concept practical those horses are skill based courses you divide every subject into 30 and its application how do you apply the theory and those are who is learning which is happening which which you are going through so then this looks like let's say my student stays for 1 year extension of BSC physics or BSC chemistry and you give him the actual knowledge of computer science right in the first year so that that's unbelievable so that 50% of the time which means two full years and a btech program a person can stay outside the compass are you learning outcomes and their alignment with the overall curricular structure and then once it comes back we have to test if he gets the credits okay so here is where I have talked with some inhibitions about the entire setup where you are simply assuming write although we can keep that a check it is not easy for us to keep some zones green some zones red operational cost for that what is the guarantee that the student is learning in the campus is there a is there an accident on some kind of a giant assessment either online or with some time stamps it is it is being documented know what time is done and CVT has videographed every assessment she claims that all right he is a good technician and he can repair any kind of car so she prepares you open the BMW engine and if Julie appointed by the awarding body which Awards the certificate and then that video is kept forever so you imagine won't you feel who is giving us I think that's a nail on the head where I think we all should pass for a moment and then think is our education experimenting something like this I don't think we'll be damaging I think all the Institute should come out of their close-minded [Unintelligible]

play nursery and she says it lightly and the answer lies in this new framework and know how many times as you rightly say that so it can be adopted but colleges that are approved by a city let's say a state Technical University may not know these rules of course of course a price them go and talk to them and then educate them but over and Beyond this there are some subtle problems that I play by 50% or I'll engage my faculty for the betterment of my students at a level which need not necessarily be teaching in the picture Affiliated to these 1200 University leaving all the institutions that is Institute of national importance but they have about 4 2: 3 4 4

do you think this ratio is sufficient to teach a technical subject know once we are sending our students out again these teachers it's not easy to create a project which is outcome based creating that project itself is going to take a lot of time similarly the [Unintelligible]

challenge dead and then they're going to be no nutrition I suppose it is really against the students should not disincentivize what clothes should not be in the classroom but outside the classroom evaluation and teacher would be much happier teacher would

```

1 from bertopic import BERTopic
2 import nltk
3
4 # Download the Punkt tokenizer for sentence splitting
5 nltk.download('punkt')
6
7 # Load the transcript from a text file
8 input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
9 with open(input_file, "r") as file:
10     transcript = file.read()
11
12 # Check if the transcript is empty
13 if not transcript.strip():
14     raise ValueError("The input file is empty. Please provide a valid transcript.")
15
16 # Split the transcript into sentences using NLTK's sent_tokenize
17 sentences = nltk.sent_tokenize(transcript) # More robust sentence splitting
18 sentences = [s.strip() for s in sentences if s.strip()] # Remove empty strings
19
20 # Debugging: Print the sentences to verify
21 print("Sentences extracted from the transcript:")
22 for i, sentence in enumerate(sentences):
23     print(f"{i + 1}: {sentence}")
24
25 # Check if sentences are empty
26 if not sentences:
27     raise ValueError("No sentences found in the transcript. Please check the input file.")
28
29 # Train BERTopic model

```


[illegible]

It is where we find you can get your accounting education at not paying on Google on the course completely at fees like a paper on

```
1 from bertopic import BERTopic
2 import nltk
3
4 # Download the Punkt tokenizer for sentence splitting
5 nltk.download('punkt')
6
7 # Load the transcript from a text file
8 input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
```



```

9 with open(input_file, "r") as file:
10     transcript = file.read()
11
12 # Check if the transcript is empty
13 if not transcript.strip():
14     raise ValueError("The input file is empty. Please provide a valid transcript.")
15
16 # Split the transcript into sentences using NLTK's sent_tokenize
17 sentences = nltk.sent_tokenize(transcript) # More robust sentence splitting
18 sentences = [s.strip() for s in sentences if s.strip()] # Remove empty strings
19
20 # Debugging: Print the sentences to verify
21 print("Sentences extracted from the transcript:")
22 for i, sentence in enumerate(sentences):
23     print(f"{i + 1}: {sentence}")
24
25 # Check if sentences are empty
26 if not sentences:
27     raise ValueError("No sentences found in the transcript. Please check the input file.")
28
29 # Train BERTopic model
30 topic_model = BERTopic()
31 topics, _ = topic_model.fit_transform(sentences)
32
33 # Get topic information
34 topic_info = topic_model.get_topic_info()
35
36 # Save the topics to a file
37 output_file = "bertopic_clustered_transcript.txt" # Replace with your desired output file path
38 with open(output_file, "w") as f:
39     for topic_id in topic_info['Topic']:
40         if topic_id != -1: # Ignore outliers
41             f.write(f"Topic {topic_id}:\n")
42             for sentence, assigned_topic in zip(sentences, topics):
43                 if assigned_topic == topic_id:
44                     f.write(f" - {sentence}\n")
45             f.write("\n")
46
47 print(f"Clustering complete. Output saved to '{output_file}'.")

```

so we picked up the critical learning outcomes the generic learning outcomes which every student must gain as a graduate so The (play level descriptors based on a particular soft skills you will get a particular level in the NCR that course will be related there are six kinds of thinking itself there is adaptive thinking computational thinking critical thinking then you have creative and we have to integrate those skills in the courses these skills cannot be just caught on a standalone basis yes you can explain to talk and interact with students why are we really bothered about that cognitive time the teacher is one who is dedicated his it'll be a very thoughtful process cognitive process through which he will also be enjoying along with the students he will feel

ValueError Traceback (most recent call last)

<ipython-input-18-ed92283479ab> in <cell line: 0>()

```

29 # Train BERTopic model
30 topic_model = BERTopic()
--> 31 topics, _ = topic_model.fit_transform(sentences)
32
33 # Get topic information

```

⤴ 5 frames

/usr/local/lib/python3.11/dist-packages/numpy/core/_methods.py in _amax(a, axis, out, keepdims, initial, where)

```

39 def _amax(a, axis=None, out=None, keepdims=False,
40           initial=_NoValue, where=True):
--> 41     return umr_maximum(a, axis, None, out, keepdims, initial, where)
42
43 def _amin(a, axis=None, out=None, keepdims=False,

```

ValueError: zero-size array to reduction operation maximum which has no identity

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

Sentences extracted from the transcript:

1: so sorry we know that India has seen a huge Revolution with digital payments we all thought that India is a place at least the I think immediately after you play the next big revolution personally I think is an education and the complete homework for this do you think ncrf plus any people together will be the next big revolution after up in India absolutely and why I think so is because how many changes have happened in the real world so many changes have happened in the requirement of the industry requirement of July 2020 and we recently celebrated the 4th anniversary of NP 2020 ncrf has been brought to implement the intent of 2020 2020 this allows you all the Innovation the way you educate your kids you educate your students still it provides you the basic guide: broken the shackles of which were there in the education sector right so yes it is a big Revolution and this is going to change I pull up my kitchen to add to my taste buds in a way that it's convenient for me and one fine day you came and you changed my kitchen we have been driving buses at Max now we should fly how do we do this look at the requirement of the industry requirement of the link that whatever he has been taught has no relevance to the real life world when industry is moving that fast when they require will you be not only beneficial but also very facilitated for all of us very liberating for all of us this is going to be highly new things create new ways of doing things learn something new but once we learn it there is no limit to Innovation and creativity creating our vision and Innovative Minds into the education sector where is we are applying it elsewhere everywhere else no we are to be fun for everyone and I can I can tell you that already a number of Institutions have adopted the any pain and CRF to vary in next step on how to implement an AP if I can request you to give me an elevator pitch for an EP and then an elevator pitch for not it allows for creditor of all learnings whether in academics or killing or an experiential learning and all these three kinds and people who are already skilled or already in the professional area not there the experiential learning would play a big part increase of technology which has been created by single Department good question so therefore all of us we have to work in tandem all kinds of learning are being contractors including learning of soft skills employability skills life skills your hand skills you go out it's all very flexible so that there's no Dropout there's no Dropout so these three things coupled with use of technology I want my son to be an engineer don't you think if we create a give me five approach to Credit Systems everybody will come to the turn off infrastructure India has today we created in last 1775 years we are going to double that infrastructure in next 10:11 years more number of other branches even liberal arts social sciences if I want to be an award Society I would need a proper mix of all have you already seen the Fallout of this know you can see the photo you can see how many Engineers are there for Designing our [Unintelligible]

any of the new technology machines current any laser-based machines any automated operating machines robotic operations robotic process and that is killing is equally important and this is important in multiple areas and therefore multidisciplinary me and you want to know design a VTech in CSC syllabus or be taking AI syllabus that is Nip complaint under the ncrf framework how are you teaching Teddy teaching Terry is not sufficient if you want the student to really understand and reply that concept practical those horses are skill based courses you divide every subject into 30 and its application how do you apply the theory and those are who is learning which is happening which which you are going through so then this looks like let's say my student stays for 1 year extension of BSC physics or BSC chemistry and you give him the actual knowledge of computer science right in the first year so that that's unbelievable so that 50% of the time which means two full years and a btech program a person can stay outside the compass are you learning outcomes and their alignment with the overall curricular structure and then once it comes back we have to test if he gets the credits okay so here is where I have talked with some inhibitions about the entire setup where you are simply assuming write although we can keep that a check it is not easy for us to keep some zones green some zones red operational cost for that what is the guarantee that the student is learning in the campus is there a is there an accident on some kind of a giant assessment either online or with some time stamps it is it is being documented know what time is done and CVT has videographed every assessment she claims that all right he is a good technician and he can repair any kind of car so she prepares you open the BMW engine and if Julie appointed by the awarding body which Awards the certificate and then that video is kept forever so you imagine won't you feel who is giving us I think that's a nail on the head where I think we all should pass for a moment and then think is our education experimenting something like this I don't think we'll be damaging I think all the Institute should come out of their close-minded [Unintelligible]

play nursery and she says it lightly and the answer lies in this new framework and know how many times as you rightly say that so it can be adopted but colleges that are approved by a city let's say a state Technical University may not know these rules of course of course a price them go and talk to them and then educate them but over and Beyond this there are some subtle problems that I play by 50% or I'll engage my faculty for the betterment of my students at a level which need not necessarily be teaching in the picture Affiliated to these 1200 University leaving all the institutions that is Institute of national importance but they have about 4 2: 3 4 4

do you think this ratio is sufficient to teach a technical subject know once we are sending our students out again these teachers it's not easy to create a project which is outcome based creating that project itself is going to take a lot of time similarly the [Unintelligible]

challenge dead and then they're going to be no nutrition I suppose it is really against the students should not disincentivize what clothes should not be in the classroom but outside the classroom evaluation and teacher would be much happier teacher would

```
1 from bertopic import BERTopic
2 import nltk
3
4 # Download the Punkt tokenizer for sentence splitting
5 nltk.download('punkt')
6
7 # Load the transcript from a text file
8 input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
9 with open(input_file, "r") as file:
10     transcript = file.read()
11
12 # Check if the transcript is empty
13 if not transcript.strip():
14     raise ValueError("The input file is empty. Please provide a valid transcript.")
15
16 # Split the transcript into sentences using NLTK's sent_tokenize
17 sentences = nltk.sent_tokenize(transcript) # More robust sentence splitting
18 sentences = [s.strip() for s in sentences if s.strip()] # Remove empty strings
19
20 # Debugging: Print the sentences to verify
21 print("Sentences extracted from the transcript:")
22 for i, sentence in enumerate(sentences):
23     print(f"{i + 1}: {sentence}")
24
25 # Check if sentences are empty
26 if not sentences:
27     raise ValueError("No sentences found in the transcript. Please check the input file.")
28
29 # Train BERTopic model
```

[illegible]

1. **Registration:** The first step in the process is to register the business with the appropriate state and federal agencies. This involves obtaining a tax ID number, registering for sales tax, and obtaining any necessary licenses or permits.

https://colab.research.google.com/drive/1bS3tCUKU1o3xt_AWZA_oZYU_aaDKLEr#scrollTo=8T19YzuSFP5T&printMode=true 56/65

```

[!nltk_data] Downloading package punkt to /root/nltk_data...
[!nltk_data] Package punkt is already up-to-date!
Sentences extracted from the transcript:
1: so sorry we know that India has seen a huge Revolution with digital payments we all thought that India is a place at least th
I think immediately after you play the next big revolution personally I think is an education and the complete homework for this
do you think ncrf plus any people together will be the next big revolution after up in India absolutely and why I think so is bec
how many changes have happened in the real world so many changes have happened in the requirement of the industry requirement of
July 2020 and we recently celebrated the 4th anniversary of NP 2020 ncrf has been brought to implement the intent of 2020 2020 tl
this allows you all the Innovation the way you educate your kids you educate your students still it provides you the basic guide:
broken the shackles of which were there in the education sector right so yes it is a big Revolution and this is going to change i
pull up my kitchen to add to my taste buds in a way that it's convenient for me and one fine day you came and you changed my kitc
we have been driving buses at Max now we should fly how do we do this look at the requirement of the industry requirement of the
link that whatever he has been taught has no relevance to the real life world when industry is moving that fast when they require
will you be not only beneficial but also very facilitated for all of us very liberating for all of us this is going to be highly
new things create new ways of doing things learn something new but once we learn it there is no limit to Innovation and creativi
creating our vision and Innovative Minds into the education sector where is we are applying it elsewhere everywhere else no we ar
to be fun for everyone and I can I can tell you that already a number of Institutions have adopted the any pain and CRF to varyir
next step on how to implement an AP if I can request you to give me an elevator pitch for an EP and then an elevator pitch for ne
it allows for creditor of all learnings weather in academics orange killing or an experiential learning and all these three kind:
and people who are already skilled or already in the professional area not there the experiential learning would play a big part
increase of technology which has been created by single Department good question so therefore all of us we have to work in tander
all kinds of learning are being contractors including learning of soft skills employability skills life skills your hand skills y
go out it's all very flexible so that there's no Dropout there's no Dropout so these three things coupled with use of technology
I want my son to be an engineer don't you think if we create a give me five approach to Credit Systems everybody will come to the
turn off infrastructure India has today we created in last 1775 years we are going to double that infrastructure in next 10:11 ye
more number of other branches even liberal arts social sciences if I want to be an award Society I would need a proper mix of al
have you already seen the Fallout of this know you can see the photo you can see how many Engineers are there for Designing our l
[Unintelligible]
any of the new technology machines current any laser-based machines any automated operating machines robotic operations
robotic process and that is killing is equally important and this is important in multiple areas and therefore multidisciplinary
me and you want to know design a VTech in CSC syllabus or be taking AI syllabus that is Nip complaint under the ncrf framework he
are you teaching Teddy teaching Terry is not sufficient if you want the student to really understand and reply that concept pract
those horses are skill bass courses you divide every subject into 30 and its application how do you apply the theory and those ar
who is learning which is happening which which you are going through so then this looks like let's say my student stays for 1 ye
extension of BSC physics or BSC chemistry and you give him the actual knowledge of computer science right in the first year so th
that's unbelievable so that 50% of the time which means two full years and a btech program a person can stay outside the compass
are you learning outcomes and their alignment with the overall curricular structure and then once it comes back we have to test i
he gets the credits okay so here is where I have talked with some inhibitions about the entire setup where you are simply assumir
write although we can keep that a check it is not easy for us to keep some zones green some zones red operational cost for that i
what is the guarantee that the student is learning in the campus is there a is there an accident on some kind of a giant assessme
either online or with some time stamps it is it is being documented know what time is done and CVT has videographed every assessr
she claims that all right he is a good technician and he can repair any kind of car so she prepares you open the BMW engine and i
Julie appointed by the awarding body which Awards the certificate and then that video is kept forever so you imagine won't you fe
who is giving us I think that's a nail on the head where I think we all should pass for a moment and then think is our education
experimenting something like this I don't think we'll be damaging I think all the Institute should come out of their close-minder
[Unintelligible]
play nursery and she says it lightly and the answer lies in this new framework and know how many times as you rightly say that si
it can be adopted but colleges that are approved by a city let's say a state Technical University may not know these rules of cou
of course a price them go and talk to them and then educate them but over and Beyond this there are some subtle problems that I i
play by 50% or I'll engage my faculty for the betterment of my students at a level which need not necessarily be teaching in the
picture Affiliated to these 1200 University leaving all the ionis that is Institute of national importance but they have about 4
2: 3 4 4
do you think this ratio is sufficient to teach a technical subject know once we are sending our students out again these teacher:
it's not easy to create a project which is outcome based creating that project itself is going to take a lot of time similarly tl
[Unintelligible]
challenge dead and then they're going to be no nutrition I suppose it is really against the students should not disincentivize re
what clothes should not be in the classroom but outside the classroom evaluation and teacher would be much happier teacher would

```

```

1 from openai import OpenAI
2 from sklearn.cluster import KMeans
3 import numpy as np
4 import nltk
5
6 # Initialize the OpenAI client
7 client = OpenAI(api_key="sk-proj-hTAcJedPA-L6e5dEsETts2Smkea8aobk4gRC5FmPUY82gFFquqZItpbxjGBIVPata9A2upge_rT3B1bkFJfAYUkHc61RQEERO_
8
9 # Load the transcript from a text file
10 input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
11 with open(input_file, "r") as file:
12     transcript = file.read()
13
14 # Check if the transcript is empty
15 if not transcript.strip():
16     raise ValueError("The input file is empty. Please provide a valid transcript.")
17
18 # Split the transcript into sentences using NLTK's sent_tokenize
19 nltk.download('punkt')
20 sentences = nltk.sent_tokenize(transcript) # More robust sentence splitting
21 sentences = [s.strip() for s in sentences if s.strip()] # Remove empty strings
22
23 # Debugging: Print the sentences to verify
24 print("Sentences extracted from the transcript:")
25 for i, sentence in enumerate(sentences):
26     print(f"{i + 1}: {sentence}")
27
28 # Check if sentences are empty
29 if not sentences:

```

```

30     raise ValueError("No sentences found in the transcript. Please check the input file.")
31
32 # Generate embeddings using OpenAI's text-embedding-ada-002 model
33 def get_embeddings(texts):
34     response = client.embeddings.create(
35         input=texts,
36         model="text-embedding-ada-002"
37     )
38     return [item.embedding for item in response.data]
39
40 # Get embeddings for all sentences
41 sentence_embeddings = get_embeddings(sentences)
42
43 # Apply K-Means clustering
44 num_clusters = 5 # Number of topics
45 kmeans = KMeans(n_clusters=num_clusters)
46 kmeans.fit(sentence_embeddings)
47 cluster_labels = kmeans.labels_
48
49 # Group sentences by cluster
50 clustered_sentences = {i: [] for i in range(num_clusters)}
51 for sentence, label in zip(sentences, cluster_labels):
52     clustered_sentences[label].append(sentence)
53
54 # Save the clusters to a file
55 output_file = "openai_clustered_transcript.txt" # Replace with your desired output file path
56 with open(output_file, "w") as f:
57     for cluster_id, sentences_in_cluster in clustered_sentences.items():
58         f.write(f"Cluster {cluster_id}:\n")
59         for sentence in sentences_in_cluster:
60             f.write(f" - {sentence}\n")
61         f.write("\n")
62
63 print(f"Clustering complete. Output saved to '{output_file}'.")

```

↕ 1 frames

```

/usr/local/lib/python3.11/dist-packages/openai/lib/_old_api.py in __call__(self, *_args, **_kwargs)
37
38     def __call__(self, *_args: Any, **_kwargs: Any) -> Any:
---> 39         raise APIRemovedInV1(symbol=self._symbol)
40
41

```

APIRemovedInV1:

You tried to access `openai.Embedding`, but this is no longer supported in `openai>=1.0.0` - see the README at <https://github.com/openai/openai-python> for the API.

You can run ``openai migrate`` to automatically upgrade your codebase to use the 1.0.0 interface.

Alternatively, you can pin your installation to the old version, e.g. ``pip install openai==0.28``

A detailed migration guide is available here: <https://github.com/openai/openai-python/discussions/742>

https://colab.research.google.com/drive/1bS3tCUKU1o3xt_AWZA_oZYU_ggDKLEr#scrollTo=8T19YzuSFP5T&printMode=true 59/65

```

1 from sentence_transformers import SentenceTransformer
2 from sklearn.cluster import KMeans
3 import numpy as np
4 import nltk
5
6 # Load a pre-trained sentence embedding model
7 model = SentenceTransformer('all-MiniLM-L6-v2')
8
9 # Load the transcript from a text file
10 input_file = "/content/Kalsi sir transcript.txt" # Replace with your input file path
11 with open(input_file, "r") as file:
12     transcript = file.read()
13
14 # Check if the transcript is empty
15 if not transcript.strip():
16     raise ValueError("The input file is empty. Please provide a valid transcript.")
17
18 # Split the transcript into sentences using NLTK's sent_tokenize
19 nltk.download('punkt')
20 sentences = nltk.sent_tokenize(transcript) # More robust sentence splitting
21 sentences = [s.strip() for s in sentences if s.strip()] # Remove empty strings
22
23 # Debugging: Print the sentences to verify
24 print("Sentences extracted from the transcript:")
25 for i, sentence in enumerate(sentences):
26     print(f"{i + 1}: {sentence}")
27
28 # Check if sentences are empty
29 if not sentences:
30     raise ValueError("No sentences found in the transcript. Please check the input file.")
31
32 # Generate embeddings using Sentence Transformers
33 sentence_embeddings = model.encode(sentences)
34
35 # Apply K-Means clustering
36 num_clusters = 5 # Number of topics
37 kmeans = KMeans(n_clusters=num_clusters)
38 kmeans.fit(sentence_embeddings)
39 cluster_labels = kmeans.labels_
40
41 # Group sentences by cluster
42 clustered_sentences = {i: [] for i in range(num_clusters)}
43 for sentence, label in zip(sentences, cluster_labels):
44     clustered_sentences[label].append(sentence)
45
46 # Save the clusters to a file
47 output_file = "clustered_transcript.txt" # Replace with your desired output file path
48 with open(output_file, "w") as f:
49     for cluster_id, sentences_in_cluster in clustered_sentences.items():
50         f.write(f"Cluster {cluster_id}:\n")
51         for sentence in sentences_in_cluster:
52             f.write(f" - {sentence}\n")
53         f.write("\n")
54
55 print(f"Clustering complete. Output saved to '{output_file}'.")

```



```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
Sentences extracted from the transcript:
1: so sorry we know that India has seen a huge Revolution with digital payments we all thought that India is a place at least the
I think immediately after you play the next big revolution personally I think is an education and the complete homework for this
do you think ncrf plus any people together will be the next big revolution after up in India absolutely and why I think so is bec
how many changes have happened in the real world so many changes have happened in the requirement of the industry requirement of
July 2020 and we recently celebrated the 4th anniversary of NP 2020 ncrf has been brought to implement the intent of 2020 2020 t
this allows you all the Innovation the way you educate your kids you educate your students still it provides you the basic guide:
broken the shackles of which were there in the education sector right so yes it is a big Revolution and this is going to change i
pull up my kitchen to add to my taste buds in a way that it's convenient for me and one fine day you came and you changed my kitc
we have been driving buses at Max now we should fly how do we do this look at the requirement of the industry requirement of the
link that whatever he has been taught has no relevance to the real life world when industry is moving that fast when they require
will you be not only beneficial but also very facilitated for all of us very liberating for all of us this is going to be highly
new things create new ways of doing things learn something new but once we learn it there is no limit to Innovation and creativi
creating our vision and Innovative Minds into the education sector where is we are applying it elsewhere everywhere else no we ar
to be fun for everyone and I can I can tell you that already a number of Institutions have adopted the any pain and CRF to varyir
next step on how to implement an AP if I can request you to give me an elevator pitch for an EP and then an elevator pitch for ne
it allows for creditor of all learnings weather in academics orange killing or an experiential learning and all these three kind:
and people who are already skilled or already in the professional area not there the experiential learning would play a big part
increase of technology which has been created by single Department good question so therefore all of us we have to work in tander
all kinds of learning are being contractors including learning of soft skills employability skills life skills your hand skills y
go out it's all very flexible so that there's no Dropout there's no Dropout so these three things coupled with use of technology
I want my son to be an engineer don't you think if we create a give me five approach to Credit Systems everybody will come to the
turn off infrastructure India has today we created in last 1775 years we are going to double that infrastructure in next 10:11 ye
more number of other branches even liberal arts social sciences if I want to be an award Society I would need a proper mix of al
have you already seen the Fallout of this know you can see the photo you can see how many Engineers are there for Designing our l
[Unintelligible]
any of the new technology machines current any laser-based machines any automated operating machines robotic operations
robotic process and that is killing is equally important and this is important in multiple areas and therefore multidisciplinary
me and you want to know design a VTech in CSC syllabus or be taking AI syllabus that is Nip complaint under the ncrf framework h
are you teaching Teddy teaching Terry is not sufficient if you want the student to really understand and reply that concept pract
those horses are skill bass courses you divide every subject into 30 and its application how do you apply the theory and those ar
who is learning which is happening which which you are going through so then this looks like let's say my student stays for 1 ye
extension of BSC physics or BSC chemistry and you give him the actual knowledge of computer science right in the first year so th
that's unbelievable so that 50% of the time which means two full years and a btech program a person can stay outside the compass
are you learning outcomes and their alignment with the overall curricular structure and then once it comes back we have to test i
he gets the credits okay so here is where I have talked with some inhibitions about the entire setup where you are simply assumir
write although we can keep that a check it is not easy for us to keep some zones green some zones red operational cost for that v
what is the guarantee that the student is learning in the campus is there a is there an accident on some kind of a giant assessme
either online or with some time stamps it is it is being documented know what time is done and CVT has videographed every assessr
she claims that all right he is a good technician and he can repair any kind of car so she prepares you open the BMW engine and i
Julie appointed by the awarding body which Awards the certificate and then that video is kept forever so you imagine won't you fe
who is giving us I think that's a nail on the head where I think we all should pass for a moment and then think is our education
experimenting something like this I don't think we'll be damaging I think all the Institute should come out of their close-minder
[Unintelligible]
play nursery and she says it lightly and the answer lies in this new framework and know how many times as you rightly say that si
it can be adopted but colleges that are approved by a city let's say a state Technical University may not know these rules of cou
of course a price them go and talk to them and then educate them but over and Beyond this there are some subtle problems that I i
play by 50% or I'll engage my faculty for the betterment of my students at a level which need not necessarily be teaching in the
picture Affiliated to these 1200 University leaving all the ionis that is Institute of national importance but they have about 4
2: 3 4 4
do you think this ratio is sufficient to teach a technical subject know once we are sending our students out again these teacher:
it's not easy to create a project which is outcome based creating that project itself is going to take a lot of time similarly tl
[Unintelligible]
challenge dead and then they're going to be no nutrition I suppose it is really against the students should not disincentivize r
what clothes should not be in the classroom but outside the classroom evaluation and teacher would be much happier teacher would
```

```
1 from transformers import AutoTokenizer, AutoModelForSequenceClassification
2
3 # Load the tokenizer and model
4 tokenizer = AutoTokenizer.from_pretrained("WebOrganizer/TopicClassifier")
5 model = AutoModelForSequenceClassification.from_pretrained(
6     "WebOrganizer/TopicClassifier",
7     trust_remote_code=True,
8     use_memory_efficient_attention=False
9 )
10
11 # Read the content from the uploaded text file
12 file_path = '/content/Kalsi sir transcript.txt' # Replace with your actual file name
13 with open(file_path, 'r') as file:
14     web_page = file.read()
15
16 # Tokenize the input text
17 inputs = tokenizer([web_page], return_tensors="pt")
18
19 # Get the model's predictions
20 outputs = model(**inputs)
21
22 # Compute the probabilities and get the predicted topic
23 probs = outputs.logits.softmax(dim=-1)
24 predicted_topic = probs.argmax(dim=-1).item()
25
26 # Print the predicted topic
27 print(f"Predicted Topic: {predicted_topic}")
28 # The predicted topic will correspond to the label (e.g., 5 for "Hardware")
```

have you been engaged with the best of the engineering colleges that has already been done and now we are also making a part of 1

destroy world we can pick quickly they can pick quickly if they want to implement every solution is there in Seattle and if some
 tokenizers from 100%
 can you sing happy, I'm going to the class teaching and coming out except that or
 special tokens map,son: 100%
 anything over on top of this is difficult for us because it's going to be time-consuming for us and we also have a research mand
 so far we can't
 however there would be certain people who would like to work with them for them either here or in the industry anyway so we have
 on the days from 100%
 good morning assume I'm walking into a bachelor's program knowing good amount of physics it is BSN physics bsnbcs I don't want to
 morning outcomes beforehand why should you be forced to go to the same course and waste your time why can't the time be used for
 configuration.py
 integrity is very important that's an important parameter right that's the most important parameter here for the first time ugc :
 and this provides that the first semester course student cannot have that basic advantage of RPL hence he's not he cannot say that you c
 day of difficulty 100%
 course up to 50% of the time or up to 50% of the credits can be earned either through online courses or through RPM apprenticesh
 in favor of the operation but as of today this is what I have been provided it is significant number of me things that I can affi
 modeling.py
 institution who are not working around that they will teach the entire engineering to the student in the first one and a half ye
 do students have to be placed in highly repeated companies like Google Microsoft article some of the top companies in the world :
 model references: 100%
 a half years very nice model anyone can afford model weather so a c
 this is never integrated into the curriculum properly which I think is very important but currently model Dimension being in pl

so we picked up the critical learning outcomes the generic learning outcomes which every student must gain as a graduate so The (n
 nlay level descriptors based on a particular soft skills you will get a particular level in the MCR that course will be related ,

```
1 from transformers import AutoTokenizer, AutoModelForSequenceClassification
2
3 # Load the tokenizer and model
4 tokenizer = AutoTokenizer.from_pretrained("WebOrganizer/TopicClassifier")
5 model = AutoModelForSequenceClassification.from_pretrained(
6     "WebOrganizer/TopicClassifier",
7     trust_remote_code=True,
8     use_memory_efficient_attention=False
9 )
10
11 # Read the content from the uploaded text file
12 input_file_path = '/content/Kalsi sir transcript.txt' # Replace with your actual file name
13 output_file_path = 'predicted_segments.txt' # Output file name
14
15 # Open the input file and read its content
16 with open(input_file_path, 'r') as file:
17     content = file.read()
18
19 # Split the content into segments (e.g., by paragraphs)
20 segments = content.split('\n\n') # Adjust the delimiter as needed
21
22 # Process each segment and write the results to the output file
23 with open(output_file_path, 'w') as output_file:
24     for segment in segments:
25         if segment.strip(): # Skip empty segments
26             # Tokenize the segment
27             inputs = tokenizer([segment], return_tensors="pt")
28
29             # Get the model's predictions
30             outputs = model(**inputs)
31
32             # Compute the probabilities and get the predicted topic
33             probs = outputs.logits.softmax(dim=-1)
34             predicted_topic = probs.argmax(dim=-1).item()
35
36             # Write the segment and its predicted topic to the output file
37             output_file.write(f"Segment:\n{segment}\n")
38             output_file.write(f"Predicted Topic: {predicted_topic}\n")
39             output_file.write("-" * 50 + "\n") # Separator for readability
40
41 print(f"Predicted segments have been saved to {output_file_path}")
```

/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(
 Predicted segments have been saved to predicted_segments.txt

```
1 import nltk
2 from nltk.tokenize import sent_tokenize
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.cluster import AgglomerativeClustering
5 import re
6 import numpy as np
7
8 # Download necessary NLTK data (if you haven't already)
```

```

9 try:
10     nltk.data.find("tokenizers/punkt")
11 except LookupError:
12     nltk.download("punkt")
13
14
15 def segment_text(file_path, num_clusters=5, clustering_method='agglomerative'):
16     """
17     Segments a text file into sections based on topic clustering,
18     using either K-Means or Agglomerative Clustering.
19
20     Args:
21         file_path (str): The path to the text file.
22         num_clusters (int): The number of topic clusters to create.
23         clustering_method (str): 'kmeans' or 'agglomerative' (default: 'agglomerative').
24
25     Returns:
26         dict: A dictionary where keys are topic labels (cluster numbers) and
27               values are lists of sentences belonging to that topic.
28     """
29
30     try:
31         with open(file_path, 'r', encoding='utf-8') as file:
32             text = file.read()
33     except FileNotFoundError:
34         return "Error: File not found."
35     except Exception as e:
36         return f"Error reading file: {e}"
37
38     # 1. Sentence Tokenization
39     sentences = sent_tokenize(text)
40
41     # 2. Text Vectorization (TF-IDF)
42     vectorizer = TfidfVectorizer(stop_words='english', max_df=0.7) # Remove common English stop words
43     tfidf_matrix = vectorizer.fit_transform(sentences)
44
45     # 3. Clustering
46     if clustering_method == 'kmeans':
47         from sklearn.cluster import KMeans # Import here to avoid unnecessary dependency if not used
48
49         kmeans = KMeans(n_clusters=num_clusters, random_state=42, n_init=10) # Set random_state for reproducibility
50         kmeans.fit(tfidf_matrix)
51         clusters = kmeans.labels_
52
53     elif clustering_method == 'agglomerative':
54         # Convert sparse matrix to dense matrix
55         tfidf_matrix_dense = tfidf_matrix.toarray()
56
57         agglomerative = AgglomerativeClustering(n_clusters=num_clusters, linkage='ward') # 'ward' linkage minimizes variance
58         clusters = agglomerative.fit_predict(tfidf_matrix_dense)
59
60     else:
61         return "Error: Invalid clustering method. Choose 'kmeans' or 'agglomerative'."
62
63     # 4. Segment Creation
64     segmented_text = {}
65     for i, cluster in enumerate(clusters):
66         if cluster not in segmented_text:
67             segmented_text[cluster] = []
68         segmented_text[cluster].append(sentences[i])
69
70     return segmented_text
71
72
73 def write_segmented_text_to_file(segmented_text, output_file="segmented_text.txt"):
74     """
75     Writes the segmented text to a file, with each topic as a section.
76
77     Args:
78         segmented_text (dict): The dictionary containing segmented text.
79         output_file (str): The name of the output file (default: "segmented_text.txt").
80     """
81     try:
82         with open(output_file, 'w', encoding='utf-8') as outfile:
83             for cluster, sentences in segmented_text.items():
84                 outfile.write(f"Topic {cluster + 1}:\n") # Topics are numbered starting from 1
85                 for sentence in sentences:
86                     outfile.write(sentence + "\n")
87                 outfile.write("\n") # Add a blank line between topics
88             print(f"Segmented text written to {output_file}")
89     except Exception as e:
90         print(f"Error writing to file: {e}")

```

```

91
92
93 def clean_filename(filename):
94     """
95     Sanitizes a filename by removing or replacing invalid characters.
96
97     Args:
98         filename (str): The original filename.
99
100    Returns:
101        str: A cleaned filename.
102    """
103    # Replace spaces with underscores
104    filename = filename.replace(" ", "_")
105    # Remove any characters that are not alphanumeric, underscores, or periods
106    filename = re.sub(r"[^a-zA-Z0-9_.]", "", filename)
107    return filename
108
109
110 def main():
111     file_path = input("Enter the path to the text file: ")
112     num_topics = int(input("Enter the desired number of topics (clusters): ")) # Get the desired number of topics from the user
113
114     clustering_method = input("Enter the clustering method ('kmeans' or 'agglomerative'): ").lower() # Get the desired clustering
115
116     segmented_data = segment_text(file_path, num_topics, clustering_method)
117
118     if isinstance(segmented_data, str): # Handle error messages from segment_text
119         print(segmented_data)
120         return
121
122     # Extract filename from path for a more descriptive output filename
123     file_name_base = file_path.split('/')[-1].split('.')[0] # Get filename without extension
124     cleaned_filename = clean_filename(file_name_base)
125     output_file_name = f"{cleaned_filename}_segmented.txt"
126
127     write_segmented_text_to_file(segmented_data, output_file_name)
128
129
130 if __name__ == "__main__":
131     main()
132

```

Enter the path to the text file: /content/Kalsi sir transcript.txt
Enter the desired number of topics (clusters): 2
Enter the clustering method ('kmeans' or 'agglomerative'): agglomerative
Segmented text written to Kalsi_sir_transcript_segmented.txt

```

1 import numpy as np
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 from sklearn.metrics.pairwise import cosine_similarity
4 import networkx as nx
5
6 # Load and preprocess the text file
7 file_path = '/content/Kalsi sir transcript.txt' # Replace with your text file pat
8 with open(file_path, 'r', encoding='utf-8') as file:
9     paragraphs = file.read().split('\n\n')
10
11 # Convert paragraphs to TF-IDF vectors
12 vectorizer = TfidfVectorizer()
13 tfidf_matrix = vectorizer.fit_transform(paragraphs)
14
15 # Compute cosine similarity matrix
16 cosine_sim = cosine_similarity(tfidf_matrix)
17
18 # Create a graph based on similarity thresholds
19 csim_threshold = 0.2 # Compare similarity value
20 graph = nx.Graph()
21 for i in range(len(paragraphs)):
22     for j in range(i + 1, len(paragraphs)):
23         if cosine_sim[i, j] > csim_threshold:
24             graph.add_edge(i, j, weight=cosine_sim[i, j])
25
26 # Detect themes using connected components
27 themes = [list(component) for component in nx.connected_components(graph)]
28
29 # Generate a summary based on themes
30 summary = []
31 for theme in themes:
32     theme_text = " ".join([paragraphs[i] for i in theme])
33     summary.append(theme_text)
34

```

```

35 # Display themes and summary
36 for idx, theme in enumerate(summary):
37     print(f"Theme {idx + 1}:\n{theme[:500]}...\n{'-'*50}")
38
39
40
41 import numpy as np
42 from sklearn.feature_extraction.text import TfidfVectorizer
43 from sklearn.metrics.pairwise import cosine_similarity
44 import networkx as nx
45
46 # Load and preprocess the text file
47 file_path = '/content/Kalsi sir transcript.txt' # Replace with your text file path
48 with open(file_path, 'r', encoding='utf-8') as file:
49     paragraphs = [p.strip() for p in file.read().split('\n\n') if p.strip()]
50
51 # Convert paragraphs to TF-IDF vectors
52 vectorizer = TfidfVectorizer()
53 tfidf_matrix = vectorizer.fit_transform(paragraphs)
54
55 # Compute cosine similarity matrix
56 cosine_sim = cosine_similarity(tfidf_matrix)
57
58 # Create a graph based on similarity thresholds
59 csim_threshold = 0.1 # Compare similarity value
60 graph = nx.Graph()
61 for i in range(len(paragraphs)):
62     for j in range(i + 1, len(paragraphs)):
63         if cosine_sim[i, j] > csim_threshold:
64             graph.add_edge(i, j, weight=cosine_sim[i, j])
65
66 # Detect themes using connected components
67 themes = [list(component) for component in nx.connected_components(graph)]
68
69 # Generate a summary based on themes
70 summary = []
71 for theme in themes:
72     # Extract a representative paragraph for each theme (first paragraph in the theme)
73     theme_text = " ".join([paragraphs[i] for i in theme])
74     summary.append(theme_text)
75
76 # Save themes to a text file
77 with open('themes_output.txt', 'w', encoding='utf-8') as out_file:
78     for idx, theme in enumerate(summary):
79         out_file.write(f"Theme {idx + 1}:\n{theme}\n{'-'*50}\n")
80
81 # Print a summary of themes
82 print(f"Detected {len(themes)} themes. Check 'themes_output.txt' for details.")
83 for idx, theme in enumerate(summary):
84     print(f"Theme {idx + 1} (Preview):\n{theme[:300]}...\n{'-'*50}")
85

```

➦ Detected 0 themes. Check 'themes_output.txt' for details.

```
1 print("Cosine Similarity Matrix:\n", cosine_sim)
```

➦ Cosine Similarity Matrix:
[[1.]]

```
1 print(f"Graph Info: Nodes = {graph.number_of_nodes()}, Edges = {graph.number_of_edges()}")
```

➦ Graph Info: Nodes = 0, Edges = 0

```
1 pip install networkx
```

➦ Requirement already satisfied: networkx in /usr/local/lib/python3.11/dist-packages (3.4.2)

```

1 from transformers import pipeline
2
3 def segment_transcript_local(text):
4     """Uses a local LLM (Hugging Face) to segment the transcript into meaningful topics."""
5
6     summarizer = pipeline("summarization", model="facebook/bart-large-cnn")
7
8     chunks = [text[i:i+1024] for i in range(0, len(text), 1024)]

```