```
1 import yt_dlp
2 import os
3
4 # Step 1: Download Audio from YouTube with Cookies
5 def download_audio(youtube_url, output_path="audio.mp3"):
6     ydl_opts = {
7         'format': 'bestaudio/best',
8         'outtmpl': output_path.replace('.mp3', ''),  # Remove .mp3 for yt-dlp to handle extension
9         'postprocessors': [{
10            'key': 'FFmpegExtractAudio',
11            'preferredcodec': 'mp3',
12            'preferredquality': '192',
13        }],
14        'cookiefile': 'cookies.txt'  # Use cookies to bypass restrictions
15    }
16
17    with yt_dlp.YoutubeDL(ydl_opts) as ydl:
18        ydl.download([youtube_url])
19
20    # Fix the file name if it ends up as audio.mp3.mp3
21    downloaded_file = output_path.replace('.mp3', '.mp3.mp3')
22    if os.path.exists(downloaded_file):
23        os.rename(downloaded_file, output_path)
24
25 # Replace with your YouTube link
26 youtube_url = "https://youtu.be/sK8SILOM37I"
27 download_audio(youtube_url)
```

```
[youtube] Extracting URL: https://youtu.be/sK8SILOM37I
[youtube] sK8SILOM37I: Downloading webpage
[youtube] sK8SILOM37I: Downloading tv client config
[youtube] sK8SILOM37I: Downloading player f6e09c70
[youtube] sK8SILOM37I: Downloading tv player API JSON
[youtube] sK8SILOM37I: Downloading ios player API JSON
[youtube] sK8SILOM37I: Downloading m3u8 information
[info] sK8SILOM37I: Downloading 1 format(s): 251
[download] Destination: audio
[download] 100% of   39.37MiB in 00:00:04 at 9.49MiB/s
[ExtractAudio] Destination: audio.mp3
Deleting original file audio (pass -k to keep)
```

```
1 #!pip install yt_dlp
```

```
1 from IPython.display import Audio
2
3 # Path to the downloaded audio file
4 audio_file = "audio.mp3"
5
6 # Play the audio
7 Audio(audio_file)
```

```
                0:00 / 57:22
```

```
1 #!pip install youtube_transcript_api
```

```
1 #!pip install pydub
```

```
1 #!pip install SpeechRecognition
```

```
1 #!pip install pytube
```

```
1 pip install deepmultilingualpunctuation
```

```
Collecting deepmultilingualpunctuation
  Downloading deepmultilingualpunctuation-1.0.1-py3-none-any.whl.metadata (4.0 kB)
Requirement already satisfied: torch>=1.8.1 in /usr/local/lib/python3.11/dist-packages (from deepmultilingualpunctuation) (2.5.1-
Requirement already satisfied: transformers in /usr/local/lib/python3.11/dist-packages (from deepmultilingualpunctuation) (4.48.3
Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from torch>=1.8.1->deepmultilingualpunctuatio
Requirement already satisfied: typing-extensions>=4.8.0 in /usr/local/lib/python3.11/dist-packages (from torch>=1.8.1->deepmultil
Requirement already satisfied: networkx in /usr/local/lib/python3.11/dist-packages (from torch>=1.8.1->deepmultilingualpunctuatio
Requirement already satisfied: jinja2 in /usr/local/lib/python3.11/dist-packages (from torch>=1.8.1->deepmultilingualpunctuation
Requirement already satisfied: fsspec in /usr/local/lib/python3.11/dist-packages (from torch>=1.8.1->deepmultilingualpunctuation
Collecting nvidia-cuda-nvrtc-cu12==12.4.127 (from torch>=1.8.1->deepmultilingualpunctuation)
  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>=1.8.1->deepmultilingualpunctuation)
  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>=1.8.1->deepmultilingualpunctuation)
  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>=1.8.1->deepmultilingualpunctuation)
  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>=1.8.1->deepmultilingualpunctuation)
  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>=1.8.1->deepmultilingualpunctuation)
  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>=1.8.1->deepmultilingualpunctuation)
  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>=1.8.1->deepmultilingualpunctuation)
  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>=1.8.1->deepmultilingualpunctuation)
  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>=1.8.1->deepmultil
Requirement already satisfied: nvidia-nvtx-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch>=1.8.1->deepmult
Collecting nvidia-nvjitlink-cu12==12.4.127 (from torch>=1.8.1->deepmultilingualpunctuation)
  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>=1.8.1->deepmultilingualpunct
Requirement already satisfied: sympy==1.13.1 in /usr/local/lib/python3.11/dist-packages (from torch>=1.8.1->deepmultilingualpunct
Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from sympy==1.13.1->torch>=1.8.1->c
Requirement already satisfied: huggingface-hub<1.0,>=0.24.0 in /usr/local/lib/python3.11/dist-packages (from transformers->deepmu
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.11/dist-packages (from transformers->deepmultilingualpunctua
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from transformers->deepmultilingualpur
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.11/dist-packages (from transformers->deepmultilingualpunctua
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.11/dist-packages (from transformers->deepmultilingualp
Requirement already satisfied: requests in /usr/local/lib/python3.11/dist-packages (from transformers->deepmultilingualpunctuatic
Requirement already satisfied: tokenizers<0.22,>=0.21 in /usr/local/lib/python3.11/dist-packages (from transformers->deepmultilir
Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.11/dist-packages (from transformers->deepmultilingual
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.11/dist-packages (from transformers->deepmultilingualpunctuat
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from jinja2->torch>=1.8.1->deepmultili
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests->transformers->
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests->transformers->deepmultilir
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests->transformers->deepmu
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests->transformers->deepmu
Downloading deepmultilingualpunctuation-1.0.1-py3-none-any.whl (5.4 kB)
Downloading nvidia_cublas_cu12-12.4.5.8-py3-none-manylinux2014_x86_64.whl (363.4 MB)
   ──────────────────────────────────── 363.4/363.4 MB 3.7 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 22.0 MB/s eta 0:00:00
Downloading nvidia_cuda_nvrtc_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl (24.6 MB)
   ──────────────────────────────────── 24.6/24.6 MB 20.5 MB/s eta 0:00:00
Downloading nvidia_cuda_runtime_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl (883 kB)
```

```python
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  from deepmultilingualpunctuation import PunctuationModel
9  import os
10 # Initialize model once at the top
11 model = PunctuationModel()
12
13 def extract_video_id(video_url):
14     """
15     Extracts the YouTube video ID from various URL formats.
16     """
17     parsed_url = urllib.parse.urlparse(video_url)
18     query_params = urllib.parse.parse_qs(parsed_url.query)
19
20     if "v" in query_params:
21         return query_params["v"][0]
22
23     match = re.search(r"(youtu\.be/|youtube\.com/embed/|youtube\.com/shorts/)([\w-]+)", video_url)
24     if match:
25         return match.group(2)
26
27     return None
28
29 def download_audio(video_url):
30     """
31     Downloads the audio using yt-dlp with cookies and returns the file path.
32     """
33     try:
34         ydl_opts = {
35             'format': 'bestaudio/best',
36             'outtmpl': 'audio.%(ext)s',
37             'cookiefile': 'cookies (1).txt',  # Use the exported cookies
38             'postprocessors': [{
39                 'key': 'FFmpegExtractAudio',
40                 'preferredcodec': 'mp3',
```

```
41              'preferredquality': '192',
42          }],
43      }
44      with yt_dlp.YoutubeDL(ydl_opts) as ydl:
45          info = ydl.extract_info(video_url, download=True)
46          return "audio.mp3"
47  except Exception as e:
48      return f"Error downloading audio: {str(e)}"

50 def convert_audio_to_wav(audio_file):
51     """
52     Converts the downloaded MP3 audio to WAV format using pydub.
53     """
54     wav_file = "audio.wav"
55     try:
56         AudioSegment.from_mp3(audio_file).export(wav_file, format="wav")
57         return wav_file
58     except Exception as e:
59         return f"Error converting to WAV: {str(e)}"

61 def transcribe_audio(audio_path, chunk_length=30):
62     """
63     Splits audio into smaller chunks, transcribes each chunk separately,
64     and adds punctuation using deepmultilingualpunctuation library.
65     """
66     recognizer = sr.Recognizer()
67     audio = AudioSegment.from_wav(audio_path)
68     total_duration = len(audio) / 1000  # Convert to seconds
69     transcribed_text = []

71     # Load punctuation model
72     model = PunctuationModel()

74     print("Transcribing audio in chunks...")

76     # In transcribe_audio()
77     punctuated_chunks = []
78     for chunk_text in transcribed_text:
79         punctuated = model.restore_punctuation(chunk_text)
80         punctuated_chunks.append(punctuated)
81         return " ".join(punctuated_chunks)

83     # Split and transcribe audio in chunks
84     for start in range(0, int(total_duration), chunk_length):
85         end = min(start + chunk_length, int(total_duration))
86         chunk = audio[start * 1000:end * 1000]  # Extract chunk in milliseconds
87         chunk.export("chunk.wav", format="wav")  # Save chunk temporarily

89         with sr.AudioFile("chunk.wav") as source:
90             try:
91                 audio_data = recognizer.record(source)
92                 text = recognizer.recognize_google(audio_data)
93                 transcribed_text.append(text)
94             except sr.UnknownValueError:
95                 transcribed_text.append("[Unintelligible]")
96             except sr.RequestError as e:
97                 return f"Error with the speech recognition service: {str(e)}"

99         os.remove("chunk.wav")  # Clean up temporary chunk file

101     # Combine chunks and add punctuation
102     combined_text = " ".join(transcribed_text)
103     punctuated_text = model.restore_punctuation(combined_text)

105     return punctuated_text

107 def get_transcript_unlisted(video_url):
108     """
109     Tries to fetch the transcript using youtube_transcript_api first,
110     then falls back to downloading and transcribing audio if necessary.
111     """
112     model = PunctuationModel()  # Initialize once
113     video_id = extract_video_id(video_url)

115     if not video_id:
116         return "Invalid YouTube URL."

118     # Try API path with punctuation
119     try:
120         transcript = YouTubeTranscriptApi.get_transcript(video_id)
121         raw_text = " ".join([item['text'] for item in transcript])
122         return model.restore_punctuation(raw_text)  # <-- Critical fix
```

```
123     except:
124         print("Transcript not available via API, attempting audio transcription...")
125
126     # Audio fallback path (existing implementation)
127     # ... rest of audio processing code ...
128     # Download and transcribe audio if no transcript is available
129     audio_file = download_audio(video_url)
130     if "Error" in audio_file:
131         return audio_file
132
133     wav_file = convert_audio_to_wav(audio_file)
134     if "Error" in wav_file:
135         return wav_file
136
137     transcription = transcribe_audio(wav_file)
138
139     # Cleanup temporary files
140     os.remove(audio_file)
141     os.remove(wav_file)
142
143     return transcription
144
145 # Example usage
146 # Example usage
147 if __name__ == "__main__":
148     video_url = input("Enter the YouTube video URL: ")
149     transcript = get_transcript_unlisted(video_url)
150
151     # Save transcript to a text file
152     if "Error" not in transcript and "Invalid YouTube URL." not in transcript:
153         output_file = "transcript.txt"
154         with open(output_file, "w", encoding="utf-8") as file:
155             file.write(transcript)
156         print(f"\nTranscript saved successfully to {output_file}")
157     else:
158         print("\n", transcript)
```

```
/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(
config.json: 100%                                              892/892 [00:00<00:00, 29.9kB/s]

model.safetensors: 100%                                        2.24G/2.24G [00:23<00:00, 129MB/s]

tokenizer_config.json: 100%                                    406/406 [00:00<00:00, 23.3kB/s]

sentencepiece.bpe.model: 100%                                  5.07M/5.07M [00:00<00:00, 77.4MB/s]

tokenizer.json: 100%                                           17.1M/17.1M [00:00<00:00, 154MB/s]

special_tokens_map.json: 100%                                  239/239 [00:00<00:00, 16.0kB/s]
Device set to use cpu
/usr/local/lib/python3.11/dist-packages/transformers/pipelines/token_classification.py:170: UserWarning: `grouped_entities` is depr
  warnings.warn(
Enter the YouTube video URL: https://youtu.be/sK8SILOM37I
Device set to use cpu

Transcript saved successfully to transcript.txt
```

```
1 import re
2 import os
3
4 def format_transcript_sentences(input_file, output_file=None):
5     """
6     Processes a transcript text file to add line breaks after sentences.
7
8     Args:
9         input_file: Path to the original transcript file
10         output_file: Path for formatted file (default: adds '_formatted' suffix)
11
12     Returns:
13         Path to the formatted file or error message
14     """
15     try:
16         # Read input file
17         with open(input_file, 'r', encoding='utf-8') as f:
18             raw_text = f.read().replace('\n', ' ')  # Remove existing newlines
19
20         # Split into sentences using punctuation followed by whitespace
21         sentences = re.split(r'(?<=[.!?]) +', raw_text)
```

```
22
23          # Format with each sentence on new line and proper capitalization
24          formatted_text = []
25          for sentence in sentences:
26              sentence = sentence.strip()
27              if sentence:
28                  # Capitalize first letter of each sentence
29                  formatted_sentence = sentence[0].upper() + sentence[1:]
30                  formatted_text.append(formatted_sentence)
31
32          formatted_text = '\n'.join(formatted_text)
33
34          # Create output filename if not provided
35          if not output_file:
36              base, ext = os.path.splitext(input_file)
37              output_file = f"{base}_formatted{ext}"
38
39          # Write formatted text
40          with open(output_file, 'w', encoding='utf-8') as f:
41              f.write(formatted_text)
42
43          return output_file
44
45      except FileNotFoundError:
46          return f"Error: File '{input_file}' not found"
47      except Exception as e:
48          return f"Error processing file: {str(e)}"
49
50 # Example usage
51 if __name__ == "__main__":
52     input_path = input("Enter path to transcript file: ").strip()
53     result = format_transcript_sentences(input_path)
54
55     if "Error" in result:
56         print(f"\n{result}")
57     else:
58         print(f"\nFormatted transcript saved to: {result}")
59         print("\nFirst 5 lines of formatted text:")
60         with open(result, 'r', encoding='utf-8') as f:
61             print(''.join(f.readlines()[:5]))
```

```
Enter path to transcript file: /content/transcript.txt

Formatted transcript saved to: /content/transcript_formatted.txt

First 5 lines of formatted text:
So, sir, we know that India has seen a huge Revolution with digital payments.
We all thought that India is a place- at least the West thought- that India is a place where, uh, many people do not get a square me
That was the narrative some 30 years ago.
And not many are literate- people cannot read, um.
But then we have now shown that digital payments number one is India, while people thought that it wouldn't even come to top 50, let
```

```
1 import numpy as np
2 from sklearn.metrics.pairwise import cosine_similarity
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sentence_transformers import SentenceTransformer
5
6 def semantic_segmentation(input_file, output_file=None, min_length=3, threshold=0.65):
7     """
8     Segments text into meaningful chunks with semantic coherence and keywords.
9
10    Args:
11        input_file: Path to formatted transcript file
12        output_file: Output path (default: adds '_segmented' suffix)
13        min_length: Minimum sentences per segment
14        threshold: Semantic similarity threshold (0-1)
15
16    Returns:
17        Path to segmented file or error message
18    """
19    try:
20        # Load ML models
21        model = SentenceTransformer('all-MiniLM-L6-v2')
22
23        # Read and split sentences
24        with open(input_file, 'r', encoding='utf-8') as f:
25            sentences = [line.strip() for line in f if line.strip()]
26
27        if len(sentences) < min_length:
28            return f"Need at least {min_length} sentences for segmentation"
```

```
29
30          # Generate sentence embeddings
31          embeddings = model.encode(sentences)
32
33          # Create segments with semantic coherence
34          segments = []
35          current_segment = []
36          current_emb = None
37
38          for sent, emb in zip(sentences, embeddings):
39              emb = emb.reshape(1, -1)
40              if not current_segment:
41                  current_segment.append(sent)
42                  current_emb = emb
43                  continue
44
45              similarity = cosine_similarity(current_emb, emb)[0][0]
46              if similarity >= threshold and len(current_segment) < 5:
47                  current_segment.append(sent)
48                  current_emb = (current_emb * len(current_segment) + emb) / (len(current_segment) + 1)
49              else:
50                  if len(current_segment) >= min_length:
51                      segments.append(current_segment)
52                  current_segment = [sent]
53                  current_emb = emb
54
55          # Finalize remaining sentences
56          if current_segment:
57              if segments and len(current_segment) < min_length:
58                  segments[-1].extend(current_segment)
59              else:
60                  segments.append(current_segment)
61
62          # Extract keywords for each segment
63          results = []
64          for seg in segments:
65              vectorizer = TfidfVectorizer(stop_words='english', ngram_range=(1,2))
66              X = vectorizer.fit_transform([' '.join(seg)])
67              features = vectorizer.get_feature_names_out()
68              keywords = features[np.argsort(X.toarray())[0][-3:]][::-1]
69
70              results.append({
71                  'sentences': seg,
72                  'keywords': keywords,
73                  'count': len(seg)
74              })
75
76          # Create output filename
77          if not output_file:
78              base, ext = os.path.splitext(input_file)
79              output_file = f"{base}_segmented{ext}"
80
81          # Write segmented output
82          with open(output_file, "w") as f:
83              for i, seg in enumerate(results, 1):
84                  f.write(f"Segment {i} ({seg['count']} sentences | Keywords: {', '.join(seg['keywords'])})\n")
85                  f.write('\n'.join(seg['sentences']) + '\n\n')
86
87          return output_file
88
89      except Exception as e:
90          return f"Error during segmentation: {str(e)}"
91
92  # Example usage
93  if __name__ == "__main__":
94      input_path = input("Enter path to formatted transcript file: ").strip()
95      result = semantic_segmentation(input_path)
96
97      if "Error" in result:
98          print(f"\n{result}")
99      else:
100         print(f"\nSegmented transcript saved to: {result}")
```

Enter path to formatted transcript file: /content/transcript_formatted.txt

modules.json: 100%      349/349 [00:00<00:00, 7.58kB/s]

config_sentence_transformers.json: 100%      116/116 [00:00<00:00, 1.31kB/s]

README.md: 100%      10.5k/10.5k [00:00<00:00, 117kB/s]

sentence_bert_config.json: 100%      53.0/53.0 [00:00<00:00, 750B/s]

config.json: 100%      612/612 [00:00<00:00, 10.9kB/s]

model.safetensors: 100%      90.9M/90.9M [00:00<00:00, 140MB/s]

tokenizer_config.json: 100%      350/350 [00:00<00:00, 15.2kB/s]

vocab.txt: 100%      232k/232k [00:00<00:00, 2.71MB/s]

tokenizer.json: 100%      466k/466k [00:00<00:00, 14.2MB/s]

special_tokens_map.json: 100%      112/112 [00:00<00:00, 6.82kB/s]

config.json: 100%      190/190 [00:00<00:00, 10.9kB/s]

```python
1  import numpy as np
2  from sklearn.metrics.pairwise import cosine_similarity
3  from sklearn.feature_extraction.text import TfidfVectorizer
4  from sentence_transformers import SentenceTransformer
5  import os
6
7  def semantic_segmentation(input_file, output_file=None, min_length=3, threshold=0.65):
8      """
9      Segments text into meaningful chunks with semantic coherence and keywords.
10
11     Args:
12         input_file: Path to formatted transcript file
13         output_file: Output path (default: adds '_segmented' suffix)
14         min_length: Minimum sentences per segment
15         threshold: Semantic similarity threshold (0-1)
16
17     Returns:
18         Path to segmented file or error message
19     """
20     try:
21         # Load ML models
22         model = SentenceTransformer('all-MiniLM-L6-v2')
23
24         # Read and split sentences
25         with open(input_file, 'r', encoding='utf-8', errors='ignore') as f:
26             sentences = [line.strip() for line in f if line.strip()]
27
28         if len(sentences) < min_length:
29             return f"Need at least {min_length} sentences for segmentation"
30
31         # Generate sentence embeddings
32         embeddings = model.encode(sentences)
33
34         # Create segments with semantic coherence
35         segments = []
36         current_segment = []
37         current_emb = None
38
39         for i, (sent, emb) in enumerate(zip(sentences, embeddings)):
40             emb = emb.reshape(1, -1)
41             if not current_segment:
42                 current_segment.append(sent)
43                 current_emb = emb
44                 continue
45
46             similarity = cosine_similarity(current_emb, emb)[0][0]
47             if similarity >= threshold and len(current_segment) < 5:
48                 current_segment.append(sent)
49                 current_emb = (current_emb * len(current_segment) + emb) / (len(current_segment) + 1)
50             else:
51                 # Finalize segment if it meets minimum length
52                 if len(current_segment) >= min_length:
53                     segments.append(current_segment)
54                 else:
55                     # If too short, append to previous segment if possible
56                     if segments:
57                         segments[-1].extend(current_segment)
58                     else:
59                         segments.append(current_segment)
```

```
60                    current_segment = [sent]
61                    current_emb = emb
62
63                # Force finalize segment at the end of the file
64                if i == len(sentences) - 1:
65                    if len(current_segment) >= min_length:
66                        segments.append(current_segment)
67                    else:
68                        # If too short, append to previous segment if possible
69                        if segments:
70                            segments[-1].extend(current_segment)
71                        else:
72                            segments.append(current_segment)
73
74          # Extract keywords for each segment
75          results = []
76          for seg in segments:
77              vectorizer = TfidfVectorizer(stop_words='english', ngram_range=(1,2))
78              X = vectorizer.fit_transform([' '.join(seg)])
79              features = vectorizer.get_feature_names_out()
80              keywords = features[np.argsort(X.toarray())[0][-3:]][::-1]
81
82              results.append({
83                  'sentences': seg,
84                  'keywords': keywords,
85                  'count': len(seg)
86              })
87
88          # Create output filename
89          if not output_file:
90              base, ext = os.path.splitext(input_file)
91              output_file = f"{base}_segmented{ext}"
92
93          # Write segmented output
94          with open(output_file, "w") as f:
95              for i, seg in enumerate(results, 1):
96                  f.write(f"Segment {i} ({seg['count']} sentences | Keywords: {', '.join(seg['keywords'])})\n")
97                  f.write('\n'.join(seg['sentences']) + '\n\n')
98
99          return output_file
100
101     except Exception as e:
102         return f"Error during segmentation: {str(e)}"
103
104 # Example usage
105 if __name__ == "__main__":
106     input_path = input("Enter path to formatted transcript file: ").strip()
107     result = semantic_segmentation(input_path)
108
109     if "Error" in result:
110         print(f"\n{result}")
111     else:
112         print(f"\nSegmented transcript saved to: {result}")
```

→ Enter path to formatted transcript file: /content/transcript_formatted.txt

    Segmented transcript saved to: /content/transcript_formatted_segmented.txt

```
1 #Now doing the transcription with timestamps
```

```
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 from deepmultilingualpunctuation import PunctuationModel
9 import os
10 # Initialize model once at the top
11 model = PunctuationModel()
12
13 def extract_video_id(video_url):
14     """
15     Extracts the YouTube video ID from various URL formats.
16     """
17     parsed_url = urllib.parse.urlparse(video_url)
18     query_params = urllib.parse.parse_qs(parsed_url.query)
19
20     if "v" in query_params:
21         return query_params["v"][0]
```

```
22
23    match = re.search(r"(youtu\.be/|youtube\.com/embed/|youtube\.com/shorts/)([\w-]+)", video_url)
24    if match:
25        return match.group(2)
26
27    return None
28
29 def download_audio(video_url):
30     """
31     Downloads the audio using yt-dlp with cookies and returns the file path.
32     """
33     try:
34         ydl_opts = {
35             'format': 'bestaudio/best',
36             'outtmpl': 'audio.%(ext)s',
37             'cookiefile': 'cookies (1).txt',  # Use the exported cookies
38             'postprocessors': [{
39                 'key': 'FFmpegExtractAudio',
40                 'preferredcodec': 'mp3',
41                 'preferredquality': '192',
42             }],
43         }
44         with yt_dlp.YoutubeDL(ydl_opts) as ydl:
45             info = ydl.extract_info(video_url, download=True)
46             return "audio.mp3"
47     except Exception as e:
48         return f"Error downloading audio: {str(e)}"
49
50 def convert_audio_to_wav(audio_file):
51     """
52     Converts the downloaded MP3 audio to WAV format using pydub.
53     """
54     wav_file = "audio.wav"
55     try:
56         AudioSegment.from_mp3(audio_file).export(wav_file, format="wav")
57         return wav_file
58     except Exception as e:
59         return f"Error converting to WAV: {str(e)}"
60
61 def transcribe_audio(audio_path, chunk_length=30):
62     """
63     Splits audio into smaller chunks, transcribes each chunk separately,
64     and adds punctuation using deepmultilingualpunctuation library.
65     """
66     recognizer = sr.Recognizer()
67     audio = AudioSegment.from_wav(audio_path)
68     total_duration = len(audio) / 1000  # Convert to seconds
69     transcribed_text = []
70
71     # Load punctuation model
72     model = PunctuationModel()
73
74     print("Transcribing audio in chunks...")
75
76     # In transcribe_audio()
77     punctuated_chunks = []
78     for chunk_text in transcribed_text:
79         punctuated = model.restore_punctuation(chunk_text)
80         punctuated_chunks.append(punctuated)
81         return " ".join(punctuated_chunks)
82
83     # Split and transcribe audio in chunks
84     for start in range(0, int(total_duration), chunk_length):
85         end = min(start + chunk_length, int(total_duration))
86         chunk = audio[start * 1000:end * 1000]  # Extract chunk in milliseconds
87         chunk.export("chunk.wav", format="wav")  # Save chunk temporarily
88
89         with sr.AudioFile("chunk.wav") as source:
90             try:
91                 audio_data = recognizer.record(source)
92                 text = recognizer.recognize_google(audio_data)
93                 transcribed_text.append(text)
94             except sr.UnknownValueError:
95                 transcribed_text.append("[Unintelligible]")
96             except sr.RequestError as e:
97                 return f"Error with the speech recognition service: {str(e)}"
98
99     os.remove("chunk.wav")  # Clean up temporary chunk file
100
101     # Combine chunks and add punctuation
102     combined_text = " ".join(transcribed_text)
103     punctuated_text = model.restore_punctuation(combined_text)
```

```python
104
105      return punctuated_text
106
107 def get_transcript_unlisted(video_url):
108      """
109      Tries to fetch the transcript using youtube_transcript_api first,
110      then falls back to downloading and transcribing audio if necessary.
111      """
112      model = PunctuationModel()  # Initialize once
113      video_id = extract_video_id(video_url)
114
115      if not video_id:
116          return "Invalid YouTube URL."
117
118      # Try API path with punctuation and timestamps
119      try:
120          transcript = YouTubeTranscriptApi.get_transcript(video_id)
121          formatted_transcript = []
122          for item in transcript:
123              start_time = convert_time(item['start'])
124              end_time = convert_time(item['start'] + item['duration'])
125              formatted_transcript.append(f"[{start_time}-{end_time}] {item['text']}")
126          return model.restore_punctuation(" ".join(formatted_transcript))
127      except:
128          print("Transcript not available via API, attempting audio transcription...")
129
130      # Audio fallback path (existing implementation)
131      # ... rest of audio processing code ...
132      # Download and transcribe audio if no transcript is available
133      audio_file = download_audio(video_url)
134      if "Error" in audio_file:
135          return audio_file
136
137      wav_file = convert_audio_to_wav(audio_file)
138      if "Error" in wav_file:
139          return wav_file
140
141      transcription = transcribe_audio(wav_file)
142
143      # Cleanup temporary files
144      os.remove(audio_file)
145      os.remove(wav_file)
146
147      # For audio fallback, timestamps are not directly available
148      # You might need to manually add timestamps or use a different approach
149      return transcription
150
151 def convert_time(seconds):
152      """Converts seconds to [hrs:mins:seconds] format."""
153      hrs = int(seconds // 3600)
154      mins = int((seconds % 3600) // 60)
155      secs = round(seconds % 60, 2)
156      return f"{hrs:02d}:{mins:02d}:{secs:05.2f}"
157
158 # Example usage
159 # Example usage
160 if __name__ == "__main__":
161      video_url = input("Enter the YouTube video URL: ")
162      transcript = get_transcript_unlisted(video_url)
163
164      # Save transcript to a text file
165      if "Error" not in transcript and "Invalid YouTube URL." not in transcript:
166          output_file = "transcript.txt"
167          with open(output_file, "w", encoding="utf-8") as file:
168              file.write(transcript)
169          print(f"\nTranscript saved successfully to {output_file}")
170      else:
171          print("\n", transcript)
172
173 # Example usage
174 # Example usage
175 if __name__ == "__main__":
176      video_url = input("Enter the YouTube video URL: ")
177      transcript = get_transcript_unlisted(video_url)
178
179      # Save transcript to a text file
180      if "Error" not in transcript and "Invalid YouTube URL." not in transcript:
181          output_file = "transcript.txt"
182          with open(output_file, "w", encoding="utf-8") as file:
183              file.write(transcript)
184          print(f"\nTranscript saved successfully to {output_file}")
```

```
185        else:
186            print("\n", transcript)
```

```
Device set to use cpu
/usr/local/lib/python3.11/dist-packages/transformers/pipelines/token_classification.py:170: UserWarning: `grouped_entities` is depre
  warnings.warn(
Enter the YouTube video URL: https://youtu.be/sK8SILOM37I
Device set to use cpu
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 f6e09c70
[youtube] sK8SILOM37I: Downloading tv player API JSON
[youtube] sK8SILOM37I: Downloading ios player API JSON
[youtube] sK8SILOM37I: Downloading m3u8 information
[info] sK8SILOM37I: Downloading 1 format(s): 251
[download] Destination: audio.webm
[download] 100% of   39.37MiB in 00:00:01 at 30.62MiB/s
[ExtractAudio] Destination: audio.mp3
Deleting original file audio.webm (pass -k to keep)
Device set to use cpu
/usr/local/lib/python3.11/dist-packages/transformers/pipelines/token_classification.py:170: UserWarning: `grouped_entities` is depre
  warnings.warn(
Transcribing audio in chunks...

Transcript saved successfully to transcript.txt
Enter the YouTube video URL: https://youtu.be/sK8SILOM37I
Device set to use cpu
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 f6e09c70
[youtube] sK8SILOM37I: Downloading tv player API JSON
[youtube] sK8SILOM37I: Downloading ios player API JSON
[youtube] sK8SILOM37I: Downloading m3u8 information
[info] sK8SILOM37I: Downloading 1 format(s): 251
[download] Destination: audio.webm
[download] 100% of   39.37MiB in 00:00:01 at 33.05MiB/s
[ExtractAudio] Destination: audio.mp3
Deleting original file audio.webm (pass -k to keep)
Device set to use cpu
/usr/local/lib/python3.11/dist-packages/transformers/pipelines/token_classification.py:170: UserWarning: `grouped_entities` is depre
  warnings.warn(
Transcribing audio in chunks...

Transcript saved successfully to transcript.txt
```

```
1 pip install nltk scikit-learn
```

```
Requirement already satisfied: nltk in /usr/local/lib/python3.11/dist-packages (3.9.1)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (1.6.1)
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)
Requirement already satisfied: numpy>=1.19.5 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.26.4)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.13.1)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.5.0)
```

```
 1 import nltk
 2 from sklearn.feature_extraction.text import TfidfVectorizer
 3 from sklearn.metrics.pairwise import cosine_similarity
 4 from itertools import combinations
 5 import numpy as np
 6 nltk.download('punkt')
 7 nltk.download('stopwords')
 8 nltk.download('wordnet')
 9 nltk.download('omw-1.4')
10 nltk.download('punkt_tab')
11 # Ensure required resources are downloaded
12 nltk.download('punkt')
13
14 # Function to read transcript from a .txt file
15 def read_transcript(file_path):
16     with open(file_path, 'r', encoding='utf-8') as file:
17         transcript = file.read()
18     return transcript
19
20 # Function to split transcript into individual sentences
21 def split_into_sentences(transcript):
22     sentences = nltk.sent_tokenize(transcript)
23     return sentences
```

```
24
25 # Function to compute cosine similarity between sentence pairs
26 def compute_cosine_similarity(sentences):
27     vectorizer = TfidfVectorizer()
28     tfidf_matrix = vectorizer.fit_transform(sentences)
29     similarity_matrix = cosine_similarity(tfidf_matrix)
30     return similarity_matrix
31
32 # Function to find similar sentence triplets based on cosine similarity
33 def find_similar_triplets(sentences, similarity_matrix, threshold=0.5):
34     triplets = []
35     n = len(sentences)
36
37     # Generate all combinations of triplets
38     for comb in combinations(range(n), 3):
39         i, j, k = comb
40         # Check if all pairs within the triplet are similar
41         if (similarity_matrix[i][j] > threshold and
42             similarity_matrix[j][k] > threshold and
43             similarity_matrix[i][k] > threshold):
44             triplets.append([sentences[i], sentences[j], sentences[k]])
45
46     return triplets
47
48 # Function to write the segmented sentences to a new .txt file
49 def write_segments_to_file(triplets, output_file):
50     with open(output_file, 'w', encoding='utf-8') as file:
51         for idx, triplet in enumerate(triplets, 1):
52             file.write(f"Segment {idx}:\n")
53             for sentence in triplet:
54                 file.write(sentence + "\n")
55             file.write("\n")
56
57 def main():
58     input_file = '/content/transcript (5).txt'  # Input .txt file path
59     output_file = 'segmented_transcript.txt'  # Output .txt file path
60
61     # Reading and processing transcript
62     transcript = read_transcript(input_file)
63     sentences = split_into_sentences(transcript)
64     similarity_matrix = compute_cosine_similarity(sentences)
65
66     # Finding similar sentence triplets
67     triplets = find_similar_triplets(sentences, similarity_matrix, threshold=0.5)
68
69     # Writing segments to output file
70     write_segments_to_file(triplets, output_file)
71
72     print(f"Segmented transcript saved to {output_file}")
73
74 if __name__ == "__main__":
75     main()
```

```
[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!
[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!
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]    Package punkt_tab is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]    Package punkt is already up-to-date!
Segmented transcript saved to segmented_transcript.txt
```

```
1 pip install nltk scikit-learn numpy
```

```
Requirement already satisfied: nltk in /usr/local/lib/python3.11/dist-packages (3.9.1)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (1.6.1)
Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (1.26.4)
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)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.13.1)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.5.0)
```

```
1 import nltk
2 from nltk.tokenize import sent_tokenize
3 from sklearn.feature_extraction.text import TfidfVectorizer
```

```python
4 from sklearn.metrics.pairwise import cosine_similarity
5 import numpy as np
```

```python
1 nltk.download('punkt')
2 nltk.download('punkt_tab')
3 nltk.download('stopwords')
4 nltk.download('wordnet')
5 nltk.download('omw-1.4')
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
[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]   Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
True
```

```python
 1 def read_transcript(file_path):
 2     with open(file_path, 'r') as file:
 3         transcript = file.read()
 4     return transcript
 5
 6 def split_into_sentences(transcript):
 7     return sent_tokenize(transcript)
 8
 9 def calculate_cosine_similarity(sentences):
10     vectorizer = TfidfVectorizer()
11     tfidf_matrix = vectorizer.fit_transform(sentences)
12     cosine_sim = cosine_similarity(tfidf_matrix)
13     return cosine_sim
14
15 def group_sentences(sentences, cosine_sim, min_sentences=3, max_sentences=10):
16     grouped_sentences = []
17     used_indices = set()
18
19     for i in range(len(sentences)):
20         if i in used_indices:
21             continue
22         group = [sentences[i]]
23         used_indices.add(i)
24         for j in range(i + 1, len(sentences)):
25             if j in used_indices:
26                 continue
27             if cosine_sim[i][j] > 0.5:  # Adjust the threshold as needed
28                 group.append(sentences[j])
29                 used_indices.add(j)
30                 if len(group) >= max_sentences:
31                     break
32         if len(group) >= min_sentences:
33             grouped_sentences.append(group)
34     return grouped_sentences
35
36 def process_transcript(file_path):
37     transcript = read_transcript(file_path)
38     sentences = split_into_sentences(transcript)
39     cosine_sim = calculate_cosine_similarity(sentences)
40     grouped_sentences = group_sentences(sentences, cosine_sim)
41     return grouped_sentences
```

```python
1 file_path = '/content/transcript (5).txt'
2 grouped_sentences = process_transcript(file_path)
3
4 for i, group in enumerate(grouped_sentences):
5     print(f"Segment {i + 1}:")
6     for sentence in group:
7         print(sentence)
8     print("\n")
```

```
Segment 1:
how do I do the transition?
how do we do this?
how do we start?


Segment 2:
all kinds of learning are being critized, including learning of soft skills, employability skills, life skills, your hand skills, y
and all these are skills.
okay, number one: the curricular structure now specifically provides for employability skills, soft skills and life skills.
```

Segment 3:
absolutely there.
absolutely, yes, absolutely.
absolutely, we can.


Segment 4:
okay.
okay.
physics, right, okay.
okay.


Segment 5:
why?
why, I see, wonderful.
why not?
so for why see why this was not not even thought about?
why?


Segment 6:
can I can.
can?
we can.


Segment 7:
school C, us mhm.
the way we teach them in school, mhm.
mhm.

```python
 1 import nltk
 2 from nltk.tokenize import sent_tokenize
 3 from sklearn.feature_extraction.text import TfidfVectorizer
 4 from sklearn.metrics.pairwise import cosine_similarity
 5 import numpy as np
 6
 7 # Download NLTK data
 8 nltk.download('punkt')
 9
10 def read_transcript(file_path):
11     try:
12         with open(file_path, 'r', encoding='utf-8') as file:
13             transcript = file.read()
14         return transcript
15     except Exception as e:
16         print(f"Error reading file: {e}")
17         return None
18
19 def split_into_sentences(transcript):
20     return sent_tokenize(transcript)
21
22 def calculate_cosine_similarity(sentences):
23     vectorizer = TfidfVectorizer()
24     tfidf_matrix = vectorizer.fit_transform(sentences)
25     cosine_sim = cosine_similarity(tfidf_matrix)
26     return cosine_sim
27
28 def group_sentences(sentences, cosine_sim, min_sentences=3, max_sentences=10, similarity_threshold=0.5):
29     grouped_sentences = []
30     used_indices = set()
31
32     for i in range(len(sentences)):
33         if i in used_indices:
34             continue
35         group = [sentences[i]]
36         used_indices.add(i)
37         for j in range(i + 1, len(sentences)):
38             if j in used_indices:
39                 continue
40             if cosine_sim[i][j] > similarity_threshold:  # Adjust the threshold as needed
41                 group.append(sentences[j])
42                 used_indices.add(j)
43                 if len(group) >= max_sentences:
44                     break
45         if len(group) >= min_sentences:
46             grouped_sentences.append(group)
47     return grouped_sentences
48
49 def process_transcript(file_path):
```

```
50     transcript = read_transcript(file_path)
51     if transcript is None:
52         return []  # Return an empty list if the file couldn't be read
53
54     sentences = split_into_sentences(transcript)
55     cosine_sim = calculate_cosine_similarity(sentences)
56     grouped_sentences = group_sentences(sentences, cosine_sim, similarity_threshold=0.6)  # Adjusted threshold
57     return grouped_sentences
58
59 # Replace with your actual file path
60 file_path = '/content/transcript (5).txt'
61 grouped_sentences = process_transcript(file_path)
62
63 if grouped_sentences:
64     for i, group in enumerate(grouped_sentences):
65         print(f"Segment {i + 1}:")
66         for sentence in group:
67             print(sentence)
68         print("\n")
69 else:
70     print("No sentences were processed. Check the file path and file content.")
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
Segment 1:
okay.
okay.
okay.


Segment 2:
why?
why not?
why?


Segment 3:
can I can.
can?
we can.
```

```
1 import nltk
2 from nltk.tokenize import sent_tokenize
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.metrics.pairwise import cosine_similarity
5 import numpy as np
6
7 # Download NLTK data
8 nltk.download('punkt')
9
10 def read_transcript(file_path):
11     try:
12         with open(file_path, 'r', encoding='utf-8') as file:
13             transcript = file.read()
14         return transcript
15     except Exception as e:
16         print(f"Error reading file: {e}")
17         return None
18
19 def split_into_sentences(transcript):
20     return sent_tokenize(transcript)
21
22 def calculate_cosine_similarity(sentences):
23     vectorizer = TfidfVectorizer()
24     tfidf_matrix = vectorizer.fit_transform(sentences)
25     cosine_sim = cosine_similarity(tfidf_matrix)
26     return cosine_sim
27
28 def group_sentences(sentences, cosine_sim, min_sentences=3, max_sentences=10, similarity_threshold=0.5):
29     grouped_sentences = []
30     used_indices = set()
31
32     for i in range(len(sentences)):
33         if i in used_indices:
34             continue
35         group = [sentences[i]]
36         used_indices.add(i)
37         for j in range(i + 1, len(sentences)):
38             if j in used_indices:
39                 continue
40             if cosine_sim[i][j] > similarity_threshold:  # Adjust the threshold as needed
```

```
41                group.append(sentences[j])
42                used_indices.add(j)
43                if len(group) >= max_sentences:
44                    break
45        if len(group) >= min_sentences:
46            grouped_sentences.append(group)
47    return grouped_sentences
48
49 def process_transcript(file_path):
50    transcript = read_transcript(file_path)
51    if transcript is None:
52        return []  # Return an empty list if the file couldn't be read
53
54    sentences = split_into_sentences(transcript)
55    cosine_sim = calculate_cosine_similarity(sentences)
56    grouped_sentences = group_sentences(sentences, cosine_sim, similarity_threshold=0.6)  # Adjusted threshold
57    return grouped_sentences
58
59 # Replace with your actual file path
60 file_path = '/content/transcript (5).txt'
61 grouped_sentences = process_transcript(file_path)
62
63 if grouped_sentences:
64    for i, group in enumerate(grouped_sentences):
65        print(f"Segment {i + 1}:")
66        for sentence in group:
67            print(sentence)
68        print("\n")
69 else:
70    print("No sentences were processed. Check the file path and file content.")
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
Segment 1:
okay.
okay.
okay.


Segment 2:
why?
why not?
why?


Segment 3:
can I can.
can?
we can.
```

```
1 #Cosine Similarity approach
```

```
1 import nltk
2 from nltk.tokenize import sent_tokenize
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.metrics.pairwise import cosine_similarity
5 import pandas as pd
6
7 # Download NLTK data
8 nltk.download('punkt')
9
10 def read_transcript(file_path):
11    try:
12        with open(file_path, 'r', encoding='utf-8') as file:
13            transcript = file.read()
14        return transcript
15    except Exception as e:
16        print(f"Error reading file: {e}")
17        return None
18
19 def split_into_sentences(transcript):
20    return sent_tokenize(transcript)
21
22 def calculate_cosine_similarity(sentences):
23    vectorizer = TfidfVectorizer()
24    tfidf_matrix = vectorizer.fit_transform(sentences)
25    cosine_sim = cosine_similarity(tfidf_matrix)
26    return cosine_sim
27
28 def save_cosine_similarity_to_csv(cosine_sim, sentences, output_file):
29    # Create a DataFrame for the cosine similarity matrix
```

```python
30     df = pd.DataFrame(cosine_sim, index=sentences, columns=sentences)
31
32     # Save the DataFrame to a CSV file
33     df.to_csv(output_file)
34     print(f"Cosine similarity matrix saved to {output_file}")
35
36 def process_transcript(file_path, output_file):
37     # Read the transcript
38     transcript = read_transcript(file_path)
39     if transcript is None:
40         return
41
42     # Split the transcript into sentences
43     sentences = split_into_sentences(transcript)
44     print("Sentences extracted:")
45     for i, sentence in enumerate(sentences):
46         print(f"{i + 1}: {sentence}")
47
48     # Calculate cosine similarity between sentences
49     cosine_sim = calculate_cosine_similarity(sentences)
50
51     # Save the cosine similarity matrix to a CSV file
52     save_cosine_similarity_to_csv(cosine_sim, sentences, output_file)
53
54 # File paths (adjust as needed)
55 file_path = '/content/transcript (5).txt'  # Path to your uploaded transcript file
56 output_file = '/content/cosine_similarity_matrix.csv'  # Output CSV file name
57
58 # Process the transcript and generate the CSV file
59 process_transcript(file_path, output_file)
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]    Package punkt is already up-to-date!
Sentences extracted:
1: so, sir, we know that India has seen a huge Revolution with digital payments.
2: we all thought that India is a place- at least the West thought- that India is a place where, uh, many people do not get a squ
3: that was the narrative some 30 years ago.
4: and not many are literate- people cannot read, um.
5: but then we have now shown that digital payments number one is India, while people thought that it wouldn't even come to top 5
6: I think, immediately after UPA, the next big revolution, personally, I think, is in education, and the complete homework for t
7: sir, my question is: do you think ncrf plus NEP put together will be the next big revolution after UPA in India?
8: absolutely, and why I think so is because in education, the last policy came up many, many years ago- that was in 1986, which
9: so so many changes have happened in the real world.
10: so many changes have happened in the requirement of the industry, requirement of Manpower for the industry.
11: however, there were no corresponding changes which happened?
12: education system.
13: so, therefore, I feel that this was the right time when we brought in the education policy 2020, honorable prime minister, de
14: 2020 is a policy, the NP 2020 is a policy, and for implementing a policy, you need a framework.
15: now, why we call it a framework?
16: we call it a framework because this is very flexible.
17: this allows you all the Innovation, the way you educate your kids, you educate your students.
18: still, it provides you the basic guidelines, the, the framework, the outer layer it provides, and that layer is mostly the en
19: that is such an enabling layer that it has broken the shackles which were there in the education sector, right?
20: so, yes, it is a big Revolution and this is going to change the way we have been educating our kids and this will be Game Cha
21: so, sir, I think, um, uh, let's go with this example of, let's say, I started off living in a small 2bhk apartment M and I sl
22: right, and one fine day, you came and you changed my kitchen completely.
23: I was using a bicycle, I moved to a scooter and a car and now you're asking me to fly and you're giving me an aircraft.
24: NEP sounds more like that for me.
25: how do I do the transition?
26: I fear that I will.
27: I will crash if I use a aeroplane without training.
28: I'm talking about all the teachers in the in the country, all the schools in the country.
29: we have been driving buses at Max.
30: now we should fly.
31: how do we do this?
32: okay, look at the requirement of the industry.
33: requirement of the industry has been moving very fast.
34: the technology is emerging every day and the industry is moving with that speed.
35: so when a student is coming out of your Institute and is going out in the market, he finds that whatever he has been taught h
36: when industry is moving that fast, when the requirement is moving that fast, don't you think it is important for us to change
37: yes, how long can we wait?
38: yes, it will take a lot of effort for every one of us to adapt to this change.
39: but this change is going to be not only beneficial, mhm, but also very facilitative for all of us, very liberating for all of
40: this is going to be highly liberative and choice-based system.
41: there are number of choices which are available to you, which are available to every student.
42: yes, when we introduce a new system, we have to really create, create new things, create new ways of doing things, learn some
43: but once we learn it, there is no limit to Innovation and creativity which this will bring.
44: there's going to be a learning curve here.
45: you mean, it's going to be difficult to fly a fighter, not really difficult to apply, because all of us we have that kind of
46: only thing is that we were not translating our vision and Innovative Minds into the education sector, whereas we are applying
47: know we are, we are moving very fast.
48: it's only education and I think the education sector is the smartest sector, which is very good at learning, and I'm I'm very
49: I can tell you that already a number of Institutions have adopted the NEP and ncrf to varying degrees and whereever they have
50: so then, um, um, I have spoken to many academics.
51: most of us do not understand NEP.
52: let on ncrf, which is the next step on how to implement NEP.
53: if I can request you to give me an elevator pitch for NEP and then an elevator pitch for ncrf, for the exclusive reason that
```

```python
1  import nltk
2  from nltk.tokenize import sent_tokenize
3  from sklearn.feature_extraction.text import TfidfVectorizer
4  from sklearn.metrics.pairwise import cosine_similarity
5  import pandas as pd
6
7  # Download NLTK data
8  nltk.download('punkt')
9
10 def read_transcript(file_path):
11     try:
12         with open(file_path, 'r', encoding='utf-8') as file:
13             transcript = file.read()
14         return transcript
15     except Exception as e:
16         print(f"Error reading file: {e}")
17         return None
18
19 def split_into_sentences(transcript):
20     return sent_tokenize(transcript)
21
22 def calculate_cosine_similarity(sentences):
23     vectorizer = TfidfVectorizer()
24     tfidf_matrix = vectorizer.fit_transform(sentences)
25     cosine_sim = cosine_similarity(tfidf_matrix)
26     return cosine_sim
27
28 def save_cosine_similarity_to_csv(cosine_sim, sentences, output_file):
29     # Create sentence names (S1, S2, S3, ...)
30     sentence_names = [f"S{i+1}" for i in range(len(sentences))]
31
32     # Create a DataFrame for the cosine similarity matrix
33     df = pd.DataFrame(cosine_sim, index=sentence_names, columns=sentence_names)
34
35     # Save the DataFrame to a CSV file
36     df.to_csv(output_file)
37     print(f"Cosine similarity matrix saved to {output_file}")
38
39 def process_transcript(file_path, output_file):
40     # Read the transcript
41     transcript = read_transcript(file_path)
42     if transcript is None:
43         return
44
45     # Split the transcript into sentences
46     sentences = split_into_sentences(transcript)
47     print("Sentences extracted:")
48     for i, sentence in enumerate(sentences):
49         print(f"S{i + 1}: {sentence}")
50
51     # Calculate cosine similarity between sentences
52     cosine_sim = calculate_cosine_similarity(sentences)
53
54     # Save the cosine similarity matrix to a CSV file
55     save_cosine_similarity_to_csv(cosine_sim, sentences, output_file)
56
57 # File paths (adjust as needed)
58 file_path = '/content/transcript (5).txt'  # Path to your uploaded transcript file
59 output_file = '/content/cosine_similarity_matrix1.csv'  # Output CSV file name
60
61 # Process the transcript and generate the CSV file
62 process_transcript(file_path, output_file)
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]    Package punkt is already up-to-date!
Sentences extracted:
S1: so, sir, we know that India has seen a huge Revolution with digital payments.
S2: we all thought that India is a place- at least the West thought- that India is a place where, uh, many people do not get a so
S3: that was the narrative some 30 years ago.
S4: and not many are literate- people cannot read, um.
S5: but then we have now shown that digital payments number one is India, while people thought that it wouldn't even come to top
S6: I think, immediately after UPA, the next big revolution, personally, I think, is in education, and the complete homework for
S7: sir, my question is: do you think ncrf plus NEP put together will be the next big revolution after UPA in India?
S8: absolutely, and why I think so is because in education, the last policy came up many, many years ago- that was in 1986, which
S9: so so many changes have happened in the real world.
S10: so many changes have happened in the requirement of the industry, requirement of Manpower for the industry.
S11: however, there were no corresponding changes which happened?
S12: education system.
S13: so, therefore, I feel that this was the right time when we brought in the education policy 2020, honorable prime minister, c
S14: 2020 is a policy, the NP 2020 is a policy, and for implementing a policy, you need a framework.
S15: now, why we call it a framework?
```

S16: we call it a framework because this is very flexible.
S17: this allows you all the Innovation, the way you educate your kids, you educate your students.
S18: still, it provides you the basic guidelines, the, the framework, the outer layer it provides, and that layer is mostly the e
S19: that is such an enabling layer that it has broken the shackles which were there in the education sector, right?
S20: so, yes, it is a big Revolution and this is going to change the way we have been educating our kids and this will be Game Ch
S21: so, sir, I think, um, uh, let's go with this example of, let's say, I started off living in a small 2bhk apartment M and I s
S22: right, and one fine day, you came and you changed my kitchen completely.
S23: I was using a bicycle, I moved to a scooter and a car and now you're asking me to fly and you're giving me an aircraft.
S24: NEP sounds more like that for me.
S25: how do I do the transition?
S26: I fear that I will.
S27: I will crash if I use a aeroplane without training.
S28: I'm talking about all the teachers in the in the country, all the schools in the country.
S29: we have been driving buses at Max.
S30: now we should fly.
S31: how do we do this?
S32: okay, look at the requirement of the industry.
S33: requirement of the industry has been moving very fast.
S34: the technology is emerging every day and the industry is moving with that speed.
S35: so when a student is coming out of your Institute and is going out in the market, he finds that whatever he has been taught
S36: when industry is moving that fast, when the requirement is moving that fast, don't you think it is important for us to chang
S37: yes, how long can we wait?
S38: yes, it will take a lot of effort for every one of us to adapt to this change.
S39: but this change is going to be not only beneficial, mhm, but also very facilitative for all of us, very liberating for all o
S40: this is going to be highly liberative and choice-based system.
S41: there are number of choices which are available to you, which are available to every student.
S42: yes, when we introduce a new system, we have to really create, create new things, create new ways of doing things, learn son
S43: but once we learn it, there is no limit to Innovation and creativity which this will bring.
S44: there's going to be a learning curve here.
S45: you mean, it's going to be difficult to fly a fighter, not really difficult to apply, because all of us we have that kind of
S46: only thing is that we were not translating our vision and Innovative Minds into the education sector, whereas we are applyin
S47: know we are, we are moving very fast.
S48: it's only education and I think the education sector is the smartest sector, which is very good at learning, and I'm I'm ver
S49: I can tell you that already a number of Institutions have adopted the NEP and ncrf to varying degrees and whereever they hav
S50: so then, um, um, I have spoken to many academics.
S51: most of us do not understand NEP.
S52: let on ncrf, which is the next step on how to implement NEP.
S53: if I can request you to give me an elevator pitch for NEP and then an elevator pitch for ncrf, for the exclusive reason that

```
1 df = pd.read_csv('/content/cosine_similarity_matrix1.csv')
2 df
```

| | Unnamed: 0 | S1 | S2 | S3 | S4 | S5 | S6 | S7 | S8 | S9 | ... | S458 | S459 | S460 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | S1 | 1.000000 | 0.144704 | 0.018483 | 0.000000 | 0.246278 | 0.091548 | 0.230970 | 0.024449 | 0.060579 | ... | 0.000000 | 0.066609 | 0.000000 |
| 1 | S2 | 0.144704 | 1.000000 | 0.034696 | 0.150481 | 0.230813 | 0.084158 | 0.141569 | 0.094473 | 0.073300 | ... | 0.000000 | 0.104808 | 0.000000 |
| 2 | S3 | 0.018483 | 0.034696 | 1.000000 | 0.000000 | 0.025220 | 0.038285 | 0.011759 | 0.295992 | 0.018366 | ... | 0.000000 | 0.031778 | 0.000000 |
| 3 | S4 | 0.000000 | 0.150481 | 0.000000 | 1.000000 | 0.056117 | 0.116833 | 0.000000 | 0.145069 | 0.120146 | ... | 0.000000 | 0.000000 | 0.127874 |
| 4 | S5 | 0.246278 | 0.230813 | 0.025220 | 0.056117 | 1.000000 | 0.055102 | 0.053087 | 0.019148 | 0.028013 | ... | 0.088463 | 0.096468 | 0.000000 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 462 | S463 | 0.019424 | 0.026095 | 0.024720 | 0.000000 | 0.056461 | 0.026382 | 0.027060 | 0.065532 | 0.054914 | ... | 0.093098 | 0.033396 | 0.000000 |
| 463 | S464 | 0.061773 | 0.042473 | 0.034434 | 0.000000 | 0.109107 | 0.113300 | 0.036147 | 0.028488 | 0.033839 | ... | 0.051788 | 0.062190 | 0.000000 |
| 464 | S465 | 0.058840 | 0.064283 | 0.012052 | 0.071909 | 0.063745 | 0.085637 | 0.083713 | 0.093791 | 0.049924 | ... | 0.089716 | 0.091426 | 0.000000 |
| 465 | S466 | 0.087732 | 0.005907 | 0.011192 | 0.000000 | 0.011617 | 0.120295 | 0.104724 | 0.026829 | 0.010998 | ... | 0.046377 | 0.029276 | 0.000000 |
| 466 | S467 | 0.000000 | 0.000000 | 0.000000 | 0.027857 | 0.060456 | 0.018291 | 0.043638 | 0.027710 | 0.000000 | ... | 0.000000 | 0.000000 | 0.000000 |

467 rows × 468 columns

```
1 # Import necessary libraries
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5
6 # Load the uploaded CSV file
7 file_path = '/content/cosine_similarity_matrix1.csv'
8 data = pd.read_csv(file_path)
9
10 # Set the first column as the index and remove the 'Unnamed: 0' column
11 data.set_index('Unnamed: 0', inplace=True)
12
13 # Generate the correlation matrix
14 correlation_matrix = data.corr()
15
16 # Print the correlation matrix summary
17 print("Correlation Matrix Summary:")
```

```
18 print(correlation_matrix.describe())
19
20 # Visualize the correlation matrix using a heatmap
21 plt.figure(figsize=(12, 10))
22 sns.heatmap(correlation_matrix, annot=False, cmap='coolwarm')
23 plt.title('Correlation Matrix Heatmap')
24 plt.show()
```
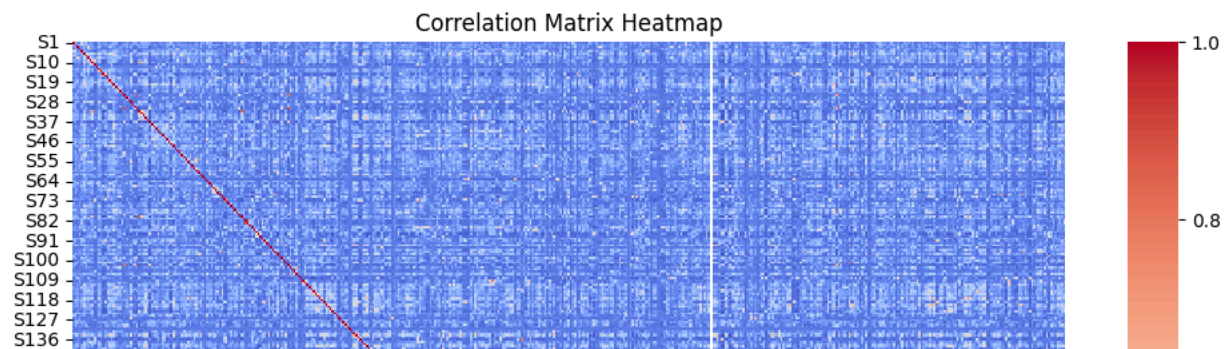
```
Correlation Matrix Summary:
              S1          S2          S3          S4          S5          S6    \
count  466.000000  466.000000  466.000000  466.000000  466.000000  466.000000
mean     0.032835    0.075637    0.045360    0.033332    0.078733    0.150554
std      0.084350    0.111592    0.089276    0.093679    0.101235    0.143893
min     -0.063055   -0.099254   -0.074572   -0.080767   -0.080913   -0.126219
25%     -0.024712   -0.006972   -0.014654   -0.031607    0.000692    0.039624
50%      0.014187    0.062486    0.030391    0.004558    0.070433    0.141213
75%      0.071348    0.135054    0.079533    0.075962    0.139049    0.249313
max      1.000000    1.000000    1.000000    1.000000    1.000000    1.000000

              S7          S8          S9         S10   ...       S458    \
count  466.000000  466.000000  466.000000  466.000000  ...  466.000000
mean     0.069673    0.077276    0.061499    0.083850  ...    0.062465
std      0.103796    0.104836    0.106447    0.127883  ...    0.125686
min     -0.097341   -0.090755   -0.090614   -0.106429  ...   -0.111943
25%     -0.003994    0.003120   -0.018873   -0.015841  ...   -0.035533
50%      0.056906    0.061219    0.046579    0.064468  ...    0.035272
75%      0.128066    0.135750    0.115556    0.147639  ...    0.135771
max      1.000000    1.000000    1.000000    1.000000  ...    1.000000

             S459        S460        S461        S462        S463        S464  \
count  466.000000  466.000000  466.000000  466.000000  466.000000  466.000000
mean     0.079398   -0.015808    0.099405    0.091793    0.031906    0.099529
std      0.143378    0.103117    0.129924    0.116551    0.111166    0.129182
min     -0.123665   -0.103533   -0.095591   -0.089121   -0.084055   -0.117284
25%     -0.037101   -0.059438   -0.003425    0.004502   -0.036968   -0.008316
50%      0.043437   -0.038165    0.083142    0.079135   -0.001110    0.087770
75%      0.177067   -0.016561    0.164672    0.156014    0.063807    0.182374
max      1.000000    1.000000    1.000000    1.000000    1.000000    1.000000

             S465        S466        S467
count  466.000000  466.000000  466.000000
mean     0.108244    0.083739    0.034295
std      0.124239    0.110004    0.104721
min     -0.119478   -0.085893   -0.084523
25%      0.012929    0.003796   -0.029172
50%      0.101166    0.064838    0.000970
75%      0.179246    0.148362    0.063946
max      1.000000    1.000000    1.000000

[8 rows x 467 columns]
```



Correlation Matrix Heatmap

```
 1 import nltk
 2 from nltk.tokenize import sent_tokenize
 3 from sklearn.feature_extraction.text import TfidfVectorizer
 4 from sklearn.metrics.pairwise import cosine_similarity
 5 import pandas as pd
 6
 7 # Download NLTK data
 8 nltk.download('punkt')
 9 nltk.download('punkt_tab')
10 nltk.download('stopwords')
11 nltk.download('wordnet')
12 nltk.download('omw-1.4')
13 nltk.download('punkt')
14 def read_transcript(file_path):
15     try:
16         with open(file_path, 'r', encoding='utf-8') as file:
17             transcript = file.read()
18         return transcript
19     except Exception as e:
20         print(f"Error reading file: {e}")
21         return None
```

```python
22
23 def split_into_sentences(transcript):
24     return sent_tokenize(transcript)
25
26 def calculate_cosine_similarity(sentences):
27     vectorizer = TfidfVectorizer()
28     tfidf_matrix = vectorizer.fit_transform(sentences)
29     cosine_sim = cosine_similarity(tfidf_matrix)
30     return cosine_sim
31
32 def segment_sentences(sentences, cosine_sim, threshold=0.5):
33     visited = [False] * len(sentences)
34     segments = []
35
36     for i in range(len(sentences)):
37         if not visited[i]:
38             segment = [sentences[i]]
39             visited[i] = True
40
41             # Check for similar sentences
42             for j in range(i + 1, len(sentences)):
43                 if not visited[j] and cosine_sim[i][j] >= threshold:
44                     segment.append(sentences[j])
45                     visited[j] = True
46
47             segments.append(segment)
48
49     return segments
50
51 def print_segments(segments):
52     for idx, segment in enumerate(segments, start=1):
53         print(f"\nSegment {idx}:")
54         for sentence in segment:
55             print(f" - {sentence}")
56
57 def process_transcript(file_path, threshold=0.5):
58     # Read the transcript
59     transcript = read_transcript(file_path)
60     if transcript is None:
61         return
62
63     # Split the transcript into sentences
64     sentences = split_into_sentences(transcript)
65     print("Sentences extracted:")
66     for i, sentence in enumerate(sentences):
67         print(f"S{i + 1}: {sentence}")
68
69     # Calculate cosine similarity between sentences
70     cosine_sim = calculate_cosine_similarity(sentences)
71
72     # Segment sentences based on similarity
73     segments = segment_sentences(sentences, cosine_sim, threshold)
74
75     # Print segmented sentences
76     print("\nSegmented Sentences:")
77     print_segments(segments)
78
79 # File path to the transcript
80 file_path = '/content/transcript (8).txt'  # Replace with your file path
81
82 # Process the transcript and print segments with a similarity threshold of 0.5
83 process_transcript(file_path, threshold=0.1)
```

```
[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.
Sentences extracted:
S1: 0 - 30: so sorry we know that India has seen a huge Revolution with digital payments we all thought that India is a place at
30 - 60: I think immediately after you play the next big revolution personally I think is an education and the complete homework
60 - 90: do you think ncrf plus any people together will be the next big revolution after up in India absolutely and why I think
90 - 120: how many changes have happened in the real world so many changes have happened in the requirement of the industry requi
120 - 150: July 2020 and we recently celebrated the 4th anniversary of NP 2020 ncrf has been brought to implement the intent of 2
150 - 180: this allows you all the Innovation the way you educate your kids you educate your students still it provides you the k
180 - 210: broken the shackles of which were there in the education sector right so yes it is a big Revolution and this is going
210 - 240: 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 char
240 - 270: we have been driving buses at Max now we should fly how do we do this look at the requirement of the industry requirem
270 - 300: link that whatever he has been taught has no relevance to the real life world when industry is moving that fast when i
300 - 330: will you be not only beneficial but also very facilitated for all of us very liberating for all of us this is going to
330 - 360: new things create new ways of doing things learn something new but once we learn it there is no limit to Innovation an
360 - 390: creating our vision and Innovative Minds into the education sector where is we are applying it elsewhere everywhere el
390 - 420: to be fun for everyone and I can I can tell you that already a number of Institutions have adopted the any pain and CF
420 - 450: 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 p
450 - 480: it allows for creditor of all learnings weather in academics orange killing or an experiential learning and all these
```

```
480 - 510: and people who are already skilled or already in the professional area not there the experiential learning would play
510 - 540: increase of technology which has been created by single Department good question so therefore all of us we have to wor
540 - 570: all kinds of learning are being contractors including learning of soft skills employability skills life skills your ha
570 - 600: 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
600 - 630: 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
630 - 660: turn off infrastructure India has today we created in last 1775 years we are going to double that infrastructure in ne
660 - 690: more number of other branches even liberal arts social sciences if I want to be an award Society I would need a proper
690 - 720: have you already seen the Fallout of this know you can see the photo you can see how many Engineers are there for Desi
720 - 750: [Unintelligible]
750 - 780: any of the new technology machines current any laser-based machines any automated operating machines robotic operatior
780 - 810: robotic process and that is killing is equally important and this is important in multiple areas and therefore multidi
810 - 840: me and you want to know design a VTech in CSC syllabus or be taking AI syllabus that is Nip complaint under the ncrf i
840 - 870: are you teaching Teddy teaching Terry is not sufficient if you want the student to really understand and reply that cc
870 - 900: those horses are skill bass courses you divide every subject into 30 and its application how do you apply the theory a
900 - 930: who is learning which is happening which which you are going through so then this looks like let's say my student stay
930 - 960: extension of BSC physics or BSC chemistry and you give him the actual knowledge of computer science right in the first
960 - 990: that's unbelievable so that 50% of the time which means two full years and a btech program a person can stay outside t
990 - 1020: are you learning outcomes and their alignment with the overall curricular structure and then once it comes back we ha
1020 - 1050: he gets the credits okay so here is where I have talked with some inhibitions about the entire setup where you are s
1050 - 1080: write although we can keep that a check it is not easy for us to keep some zones green some zones red operational cc
1080 - 1110: what is the guarantee that the student is learning in the campus is there a is there an accident on some kind of a g
1110 - 1140: either online or with some time stamps it is it is being documented know what time is done and CVT has videographed
1140 - 1170: 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
1170 - 1200: Julie appointed by the awarding body which Awards the certificate and then that video is kept forever so you imagine
1200 - 1230: 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 c
1230 - 1260: experimenting something like this I don't think we'll be damaging I think all the Institute should come out of their
1260 - 1290: [Unintelligible]
1290 - 1320: play nursery and she says it lightly and the answer lies in this new framework and know how many times as you right]
1320 - 1350: it can be adopted but colleges that are approved by a city let's say a state Technical University may not know these
1350 - 1380: of course a price them go and talk to them and then educate them but over and Beyond this there are some subtle prot
1380 - 1410: play by 50% or I'll engage my faculty for the betterment of my students at a level which need not necessarily be tea
1410 - 1440: picture Affiliated to these 1200 University leaving all the ionis that is Institute of national importance but they
S2: 3 4 4
1440 - 1470: do you think this ratio is sufficient to teach a technical subject know once we are sending our students out again t
1470 - 1500: it's not easy to create a project which is outcome based creating that project itself is going to take a lot of time
```

```python
1  import nltk
2  from nltk.tokenize import sent_tokenize
3  from sklearn.feature_extraction.text import TfidfVectorizer
4  from sklearn.metrics.pairwise import cosine_similarity
5  import pandas as pd
6
7  # Download NLTK data
8  nltk.download('punkt')
9
10 def read_transcript(file_path):
11     try:
12         with open(file_path, 'r', encoding='utf-8') as file:
13             transcript = file.read()
14         return transcript
15     except Exception as e:
16         print(f"Error reading file: {e}")
17         return None
18
19 def split_into_sentences(transcript):
20     return sent_tokenize(transcript)
21
22 def calculate_cosine_similarity(sentences):
23     vectorizer = TfidfVectorizer()
24     tfidf_matrix = vectorizer.fit_transform(sentences)
25     cosine_sim = cosine_similarity(tfidf_matrix)
26     return cosine_sim
27
28 def segment_sentences(sentences, cosine_sim, threshold=0.5):
29     visited = [False] * len(sentences)
30     segments = []
31
32     for i in range(len(sentences)):
33         if not visited[i]:
34             segment = [sentences[i]]
35             visited[i] = True
36
37             # Check for similar sentences
38             for j in range(i + 1, len(sentences)):
39                 if not visited[j] and cosine_sim[i][j] >= threshold:
40                     segment.append(sentences[j])
41                     visited[j] = True
42
43             segments.append(segment)
44
45     return segments
46
47 def print_segments(segments):
48     for idx, segment in enumerate(segments, start=1):
```

```
49                print(f"\nSegment {idx}:")
50                for sentence in segment:
51                    print(f" - {sentence}")
52
53  # 🆕 New Function to Save Segments to a Text File
54  def save_segments_to_file(segments, output_file):
55      try:
56          with open(output_file, 'w', encoding='utf-8') as file:
57              for idx, segment in enumerate(segments, start=1):
58                  file.write(f"Segment {idx}:\n")
59                  for sentence in segment:
60                      file.write(f" - {sentence}\n")
61                  file.write("\n")
62          print(f"Segmented output saved to {output_file}")
63      except Exception as e:
64          print(f"Error saving file: {e}")
65
66  def process_transcript(file_path, threshold=0.5, output_file='segmented_output.txt'):
67      # Read the transcript
68      transcript = read_transcript(file_path)
69      if transcript is None:
70          return
71
72      # Split the transcript into sentences
73      sentences = split_into_sentences(transcript)
74      print("Sentences extracted:")
75      for i, sentence in enumerate(sentences):
76          print(f"S{i + 1}: {sentence}")
77
78      # Calculate cosine similarity between sentences
79      cosine_sim = calculate_cosine_similarity(sentences)
80
81      # Segment sentences based on similarity
82      segments = segment_sentences(sentences, cosine_sim, threshold)
83
84      # Print segmented sentences
85      print("\nSegmented Sentences:")
86      print_segments(segments)
87
88      # Save segments to a text file
89      save_segments_to_file(segments, output_file)
90
91  # File paths
92  file_path = '/content/transcript (5).txt'  # Replace with your file path
93  output_file = '/content/segmented_output.txt'  # Output text file path
94
95  # Process the transcript, segment it, and save to a text file
96  process_transcript(file_path, threshold=0.15, output_file=output_file)
```

```
Sentences extracted:
S1: so, sir, we know that India has seen a huge Revolution with digital payments.
S2: we all thought that India is a place- at least the West thought- that India is a place where, uh, many people do not get a so
S3: that was the narrative some 30 years ago.
S4: and not many are literate- people cannot read, um.
S5: but then we have now shown that digital payments number one is India, while people thought that it wouldn't even come to top
S6: I think, immediately after UPA, the next big revolution, personally, I think, is in education, and the complete homework for
S7: sir, my question is: do you think ncrf plus NEP put together will be the next big revolution after UPA in India?
S8: absolutely, and why I think so is because in education, the last policy came up many, many years ago- that was in 1986, which
S9: so so many changes have happened in the real world.
S10: so many changes have happened in the requirement of the industry, requirement of Manpower for the industry.
S11: however, there were no corresponding changes which happened?
S12: education system.
S13: so, therefore, I feel that this was the right time when we brought in the education policy 2020, honorable prime minister, c
S14: 2020 is a policy, the NP 2020 is a policy, and for implementing a policy, you need a framework.
S15: now, why we call it a framework?
S16: we call it a framework because this is very flexible.
S17: this allows you all the Innovation, the way you educate your kids, you educate your students.
S18: still, it provides you the basic guidelines, the, the framework, the outer layer it provides, and that layer is mostly the e
S19: that is such an enabling layer that it has broken the shackles which were there in the education sector, right?
S20: so, yes, it is a big Revolution and this is going to change the way we have been educating our kids and this will be Game Ch
S21: so, sir, I think, um, uh, let's go with this example of, let's say, I started off living in a small 2bhk apartment M and I s
S22: right, and one fine day, you came and you changed my kitchen completely.
S23: I was using a bicycle, I moved to a scooter and a car and now you're asking me to fly and you're giving me an aircraft.
S24: NEP sounds more like that for me.
S25: how do I do the transition?
S26: I fear that I will.
S27: I will crash if I use a aeroplane without training.
S28: I'm talking about all the teachers in the in the country, all the schools in the country.
S29: we have been driving buses at Max.
S30: now we should fly.
S31: how do we do this?
S32: okay, look at the requirement of the industry.
S33: requirement of the industry has been moving very fast.
S34: the technology is emerging every day and the industry is moving with that speed.
S35: so when a student is coming out of your Institute and is going out in the market, he finds that whatever he has been taught
```

S36: when industry is moving that fast, when the requirement is moving that fast, don't you think it is important for us to chang
S37: yes, how long can we wait?
S38: yes, it will take a lot of effort for every one of us to adapt to this change.
S39: but this change is going to be not only beneficial, mhm, but also very facilitative for all of us, very liberating for all o
S40: this is going to be highly liberative and choice-based system.
S41: there are number of choices which are available to you, which are available to every student.
S42: yes, when we introduce a new system, we have to really create, create new things, create new ways of doing things, learn som
S43: but once we learn it, there is no limit to Innovation and creativity which this will bring.
S44: there's going to be a learning curve here.
S45: you mean, it's going to be difficult to fly a fighter, not really difficult to apply, because all of us we have that kind of
S46: only thing is that we were not translating our vision and Innovative Minds into the education sector, whereas we are applyir
S47: know we are, we are moving very fast.
S48: it's only education and I think the education sector is the smartest sector, which is very good at learning, and I'm I'm ver
S49: I can tell you that already a number of Institutions have adopted the NEP and ncrf to varying degrees and whereever they hav
S50: so then, um, um, I have spoken to many academics.
S51: most of us do not understand NEP.
S52: let on ncrf, which is the next step on how to implement NEP.
S53: if I can request you to give me an elevator pitch for NEP and then an elevator pitch for ncrf, for the exclusive reason that
S54: okay, see, if I talk about one simple line on NP and ncf, that is, it is crediti of all learning.
S55: it allows for crediti of all learnings, whether in academics or in Skilling or in experiential learning, and all these three

```python
1  import nltk
2  from nltk.tokenize import sent_tokenize, word_tokenize
3  from sklearn.feature_extraction.text import TfidfVectorizer
4  from sklearn.metrics.pairwise import cosine_similarity
5  import pandas as pd
6
7  # Download NLTK data
8  nltk.download('punkt')
9  nltk.download('stopwords')
10 from nltk.corpus import stopwords
11
12 def read_transcript(file_path):
13     try:
14         with open(file_path, 'r', encoding='utf-8') as file:
15             transcript = file.read()
16         return transcript
17     except Exception as e:
18         print(f"Error reading file: {e}")
19         return None
20
21 def split_into_sentences(transcript):
22     return sent_tokenize(transcript)
23
24 def calculate_cosine_similarity(sentences):
25     vectorizer = TfidfVectorizer()
26     tfidf_matrix = vectorizer.fit_transform(sentences)
27     cosine_sim = cosine_similarity(tfidf_matrix)
28     return cosine_sim
29
30 def segment_sentences(sentences, cosine_sim, threshold=0.5):
31     visited = [False] * len(sentences)
32     segments = []
33     for i in range(len(sentences)):
34         if not visited[i]:
35             segment = [sentences[i]]
36             visited[i] = True
37             for j in range(i + 1, len(sentences)):
38                 if not visited[j] and cosine_sim[i][j] >= threshold:
39                     segment.append(sentences[j])
40                     visited[j] = True
41             segments.append(segment)
42     return segments
43
44 def remove_stopwords(segment):
45     stop_words = set(stopwords.words('english'))
46     filtered_words = []
47     for sentence in segment:
48         words = word_tokenize(sentence)
49         filtered_words.extend([word.lower() for word in words if word.isalnum() and word.lower() not in stop_words])
50     return filtered_words
51
52 def find_keywords(filtered_words):
53     if len(filtered_words) < 2:
54         return None, None, 0  # Not enough words for comparison
55
56     vectorizer = TfidfVectorizer()
57     tfidf_matrix = vectorizer.fit_transform(filtered_words)
58     cosine_sim = cosine_similarity(tfidf_matrix)
59
60     max_sim = 0
61     keyword_pair = (None, None)
62     for i in range(len(filtered_words)):
```

```
63              for j in range(i + 1, len(filtered_words)):
64                  if cosine_sim[i][j] > max_sim:
65                      max_sim = cosine_sim[i][j]
66                      keyword_pair = (filtered_words[i], filtered_words[j])
67      return keyword_pair[0], keyword_pair[1], max_sim
68
69  def save_segments_to_csv(segments, output_file):
70      with pd.ExcelWriter(output_file) as writer:
71          for idx, segment in enumerate(segments, start=1):
72              filtered_words = remove_stopwords(segment)
73              word1, word2, max_sim = find_keywords(filtered_words)
74
75              # Create a DataFrame for the segment
76              df = pd.DataFrame({
77                  'Word 1': [word1] if word1 else [],
78                  'Word 2': [word2] if word2 else [],
79                  'Cosine Similarity': [max_sim] if word1 and word2 else [],
80                  'Keyword': [f"{word1}, {word2}"] if word1 and word2 else []
81              })
82
83              # Save each segment as a separate sheet
84              df.to_excel(writer, sheet_name=f'Segment {idx}', index=False)
85
86      print(f"Segmented keywords saved to {output_file}")
87
88  def process_transcript(file_path, threshold=0.5, output_file='segmented_keywords.xlsx'):
89      transcript = read_transcript(file_path)
90      if transcript is None:
91          return
92
93      sentences = split_into_sentences(transcript)
94      cosine_sim = calculate_cosine_similarity(sentences)
95      segments = segment_sentences(sentences, cosine_sim, threshold)
96
97      save_segments_to_csv(segments, output_file)
98
99      print("\nChosen Keywords for Each Segment:")
100     for idx, segment in enumerate(segments, start=1):
101         filtered_words = remove_stopwords(segment)
102         word1, word2, max_sim = find_keywords(filtered_words)
103         if word1 and word2:
104             print(f"Segment {idx}: Keywords = {word1}, {word2} (Similarity: {max_sim:.2f})")
105         else:
106             print(f"Segment {idx}: Not enough data for keywords")
107
108 # File paths
109 file_path = '/content/transcript (5).txt'  # Replace with your file path
110 output_file = '/content/segmented_keywords.xlsx'  # Output Excel file path
111
112 # Run the process
113 process_transcript(file_path, threshold=0.1, output_file=output_file)
114
```

```
[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 keywords saved to /content/segmented_keywords.xlsx

Chosen Keywords for Each Segment:
Segment 1: Keywords = sir, sir (Similarity: 1.00)
Segment 2: Keywords = 30, 30 (Similarity: 1.00)
Segment 3: Keywords = many, many (Similarity: 1.00)
Segment 4: Keywords = many, many (Similarity: 1.00)
Segment 5: Keywords = education, education (Similarity: 1.00)
Segment 6: Keywords = brought, brought (Similarity: 1.00)
Segment 7: Keywords = call, call (Similarity: 1.00)
Segment 8: Keywords = allows, allows (Similarity: 1.00)
Segment 9: Keywords = um, um (Similarity: 1.00)
Segment 10: Keywords = car, car (Similarity: 1.00)
Segment 11: Keywords = yes, yes (Similarity: 1.00)
Segment 12: Keywords = know, know (Similarity: 1.00)
Segment 13: Keywords = use, use (Similarity: 1.00)
Segment 14: Keywords = level, level (Similarity: 1.00)
Segment 15: Keywords = yes, yes (Similarity: 1.00)
Segment 16: Keywords = going, going (Similarity: 1.00)
Segment 17: Keywords = thing, thing (Similarity: 1.00)
Segment 18: Keywords = number, number (Similarity: 1.00)
Segment 19: Keywords = okay, okay (Similarity: 1.00)
Segment 20: Keywords = students, students (Similarity: 1.00)
Segment 21: Not enough data for keywords
Segment 22: Keywords = single, single (Similarity: 1.00)
Segment 23: Keywords = good, good (Similarity: 1.00)
Segment 24: Keywords = curricular, curricular (Similarity: 1.00)
Segment 25: Keywords = catch, catch (Similarity: 1.00)
```

```
Segment 26: Keywords = short, short (Similarity: 1.00)
Segment 27: Keywords = need, need (Similarity: 1.00)
Segment 28: Keywords = important, important (Similarity: 1.00)
Segment 29: Keywords = unable, unable (Similarity: 1.00)
Segment 30: Keywords = syllabus, syllabus (Similarity: 1.00)
Segment 31: Keywords = already, already (Similarity: 1.00)
Segment 32: Keywords = theory, theory (Similarity: 1.00)
Segment 33: Keywords = student, student (Similarity: 1.00)
Segment 34: Keywords = three, three (Similarity: 1.00)
Segment 35: Keywords = talking, talking (Similarity: 1.00)
Segment 36: Not enough data for keywords
Segment 37: Not enough data for keywords
Segment 38: Not enough data for keywords
Segment 39: Not enough data for keywords
Segment 40: Not enough data for keywords
Segment 41: Keywords = project, project (Similarity: 1.00)
Segment 42: Keywords = inside, inside (Similarity: 1.00)
Segment 43: Keywords = think, think (Similarity: 1.00)
Segment 44: Keywords = contain, contain (Similarity: 1.00)
Segment 45: Keywords = really, really (Similarity: 1.00)
Segment 46: Not enough data for keywords
Segment 47: Not enough data for keywords
Segment 48: Keywords = mhm, mhm (Similarity: 1.00)
Segment 49: Not enough data for keywords
Segment 50: Not enough data for keywords
Segment 51: Not enough data for keywords
```

```python
1  import nltk
2  from nltk.tokenize import sent_tokenize, word_tokenize
3  from sklearn.feature_extraction.text import TfidfVectorizer
4  from sklearn.metrics.pairwise import cosine_similarity
5  import pandas as pd
6
7  # Download NLTK data
8  nltk.download('punkt')
9  nltk.download('stopwords')
10 from nltk.corpus import stopwords
11
12 def read_segmented_file(file_path):
13     try:
14         with open(file_path, 'r', encoding='utf-8') as file:
15             content = file.read()
16             # Split the file content into segments based on a delimiter (e.g., "\n\n")
17             segments = [segment.strip() for segment in content.split("\n\n") if segment.strip()]
18         return segments
19     except Exception as e:
20         print(f"Error reading file: {e}")
21         return None
22
23 def remove_stopwords(segment):
24     stop_words = set(stopwords.words('english'))
25     words = word_tokenize(segment)
26     return [word.lower() for word in words if word.isalnum() and word.lower() not in stop_words]
27
28 def calculate_word_cosine_similarity(filtered_words):
29     if len(filtered_words) < 2:
30         return None, None, 0, None  # Not enough words for comparison
31
32     vectorizer = TfidfVectorizer()
33     tfidf_matrix = vectorizer.fit_transform(filtered_words)
34     cosine_sim = cosine_similarity(tfidf_matrix)
35
36     max_sim = 0
37     keyword_pair = (None, None)
38     for i in range(len(filtered_words)):
39         for j in range(i + 1, len(filtered_words)):
40             if cosine_sim[i][j] > max_sim:
41                 max_sim = cosine_sim[i][j]
42                 keyword_pair = (filtered_words[i], filtered_words[j])
43
44     return keyword_pair[0], keyword_pair[1], max_sim, cosine_sim
45
46 def save_word_similarity_to_excel(cosine_sim, filtered_words, sheet_name, writer):
47     if cosine_sim is not None:
48         df = pd.DataFrame(cosine_sim, index=filtered_words, columns=filtered_words)
49         df.to_excel(writer, sheet_name=sheet_name)
50     else:
51         df = pd.DataFrame(columns=['Info'])
52         df.loc[0] = ["Not enough words for comparison"]
53         df.to_excel(writer, sheet_name=sheet_name, index=False)
54
55 def save_keywords_summary_to_csv(keywords, output_file):
56     summary_df = pd.DataFrame(keywords, columns=['Segment', 'Word 1', 'Word 2', 'Cosine Similarity', 'Keyword'])
57     summary_df.to_csv(output_file, index=False)
58     print(f"Keywords summary saved to {output_file}")
```

```
 59
 60 def process_segments(file_path, output_excel='segment_word_similarity.xlsx', summary_csv='segment_keywords_summary.csv'):
 61     segments = read_segmented_file(file_path)
 62     if segments is None:
 63         return
 64
 65     keywords_summary = []
 66
 67     with pd.ExcelWriter(output_excel) as writer:
 68         for idx, segment in enumerate(segments, start=1):
 69             # Remove stopwords
 70             filtered_words = remove_stopwords(segment)
 71
 72             # Calculate cosine similarity between words
 73             word1, word2, max_sim, cosine_sim = calculate_word_cosine_similarity(filtered_words)
 74
 75             # Save word similarity matrix for each segment to Excel
 76             save_word_similarity_to_excel(cosine_sim, filtered_words, f'Segment {idx}', writer)
 77
 78             # Prepare summary data
 79             if word1 and word2:
 80                 keyword = f"{word1}, {word2}"
 81                 keywords_summary.append([f"Segment {idx}", word1, word2, max_sim, keyword])
 82             else:
 83                 keywords_summary.append([f"Segment {idx}", "N/A", "N/A", 0, "N/A"])
 84
 85     # Save summary of keywords to CSV
 86     save_keywords_summary_to_csv(keywords_summary, summary_csv)
 87
 88     # Print chosen keywords for each segment
 89     print("\nChosen Keywords for Each Segment:")
 90     for row in keywords_summary:
 91         print(f"{row[0]}: Keywords = {row[4]} (Similarity: {row[3]:.2f})")
 92
 93 # File paths
 94 segmented_file_path = '/content/segmented_output.txt'  # Replace with your segmented text file path
 95 output_excel = '/content/segment_word_similarity.xlsx'  # Excel file for word similarity matrices
 96 summary_csv = '/content/segment_keywords_summary.csv'  # CSV file for keywords summary
 97
 98 # Run the process
 99 process_segments(segmented_file_path, output_excel=output_excel, summary_csv=summary_csv)
100
```

```
[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!
Keywords summary saved to /content/segment_keywords_summary.csv

Chosen Keywords for Each Segment:
Segment 1: Keywords = sir, sir (Similarity: 1.00)
Segment 2: Keywords = thought, thought (Similarity: 1.00)
Segment 3: Keywords = 30, 30 (Similarity: 1.00)
Segment 4: Keywords = think, think (Similarity: 1.00)
Segment 5: Keywords = however, however (Similarity: 1.00)
Segment 6: Keywords = education, education (Similarity: 1.00)
Segment 7: Keywords = 2020, 2020 (Similarity: 1.00)
Segment 8: Keywords = call, call (Similarity: 1.00)
Segment 9: Keywords = way, way (Similarity: 1.00)
Segment 10: Keywords = provides, provides (Similarity: 1.00)
Segment 11: Keywords = uh, uh (Similarity: 1.00)
Segment 12: Keywords = N/A (Similarity: 0.00)
Segment 13: Keywords = saying, saying (Similarity: 1.00)
Segment 14: Keywords = nep, nep (Similarity: 1.00)
Segment 15: Keywords = say, say (Similarity: 1.00)
Segment 16: Keywords = fact, fact (Similarity: 1.00)
Segment 17: Keywords = without, without (Similarity: 1.00)
Segment 18: Keywords = talking, talking (Similarity: 1.00)
Segment 19: Keywords = N/A (Similarity: 0.00)
Segment 20: Keywords = okay, okay (Similarity: 1.00)
Segment 21: Keywords = student, student (Similarity: 1.00)
Segment 22: Keywords = yes, yes (Similarity: 1.00)
Segment 23: Keywords = yes, yes (Similarity: 1.00)
Segment 24: Keywords = new, new (Similarity: 1.00)
Segment 25: Keywords = going, going (Similarity: 1.00)
Segment 26: Keywords = minds, minds (Similarity: 1.00)
Segment 27: Keywords = already, already (Similarity: 1.00)
Segment 28: Keywords = um, um (Similarity: 1.00)
Segment 29: Keywords = okay, okay (Similarity: 1.00)
Segment 30: Keywords = already, already (Similarity: 1.00)
Segment 31: Keywords = interdisciplinarity, interdisciplinarity (Similarity: 1.00)
Segment 32: Keywords = single, single (Similarity: 1.00)
Segment 33: Keywords = good, good (Similarity: 1.00)
Segment 34: Keywords = curricular, curricular (Similarity: 1.00)
Segment 35: Keywords = learning, learning (Similarity: 1.00)
Segment 36: Keywords = multiple, multiple (Similarity: 1.00)
```

```
Segment 37: Keywords = three, three (Similarity: 1.00)
Segment 38: Keywords = short, short (Similarity: 1.00)
Segment 39: Keywords = think, think (Similarity: 1.00)
Segment 40: Keywords = infrastructure, infrastructure (Similarity: 1.00)
Segment 41: Keywords = N/A (Similarity: 0.00)
Segment 42: Keywords = feel, feel (Similarity: 1.00)
Segment 43: Keywords = important, important (Similarity: 1.00)
Segment 44: Keywords = imagine, imagine (Similarity: 1.00)
Segment 45: Keywords = unable, unable (Similarity: 1.00)
Segment 46: Keywords = design, design (Similarity: 1.00)
Segment 47: Keywords = computer, computer (Similarity: 1.00)
Segment 48: Keywords = courses, courses (Similarity: 1.00)
Segment 49: Keywords = 50, 50 (Similarity: 1.00)
Segment 50: Keywords = keep, keep (Similarity: 1.00)
Segment 51: Keywords = N/A (Similarity: 0.00)
```

```python
1  import nltk
2  from nltk.tokenize import word_tokenize
3  import re
4  import pandas as pd
5
6  # Download NLTK data if not already available
7  nltk.download('punkt')
8  nltk.download('stopwords')
9  from nltk.corpus import stopwords
10
11 # Define filler words to remove
12 FILLER_WORDS = {'umm', 'uh', 'oh', 'okay', 'like', 'you know', 'actually', 'basically', 'literally', 'well', 'so', 'just', 'i mean',
13
14 def read_segmented_file(file_path):
15     """Read segmented text file and split into segments."""
16     try:
17         with open(file_path, 'r', encoding='utf-8') as file:
18             content = file.read()
19             # Split the file content into segments based on double newlines
20             segments = [segment.strip() for segment in content.split("\n\n") if segment.strip()]
21         return segments
22     except Exception as e:
23         print(f"Error reading file: {e}")
24         return None
25
26 def clean_text(segment):
27     """Remove stopwords, fillers, and dates from a segment."""
28     stop_words = set(stopwords.words('english'))
29
30     # Remove dates (e.g., 12/03/2023, March 12, 2023)
31     segment = re.sub(r'\b(\d{1,2}[/-]\d{1,2}[/-]\d{2,4}|\b(?:January|February|March|April|May|June|July|August|September|October|Nov
32
33     # Tokenize and clean words
34     words = word_tokenize(segment)
35     filtered_words = [
36         word.lower() for word in words
37         if word.isalnum() and word.lower() not in stop_words and word.lower() not in FILLER_WORDS
38     ]
39     return ' '.join(filtered_words)
40
41 def save_to_csv(segments, cleaned_segments, output_file):
42     """Save original and cleaned segments to a CSV file."""
43     df = pd.DataFrame({
44         'Original Segment': segments,
45         'Cleaned Segment (No Stopwords or Fillers)': cleaned_segments
46     })
47     df.to_csv(output_file, index=False)
48     print(f"Cleaned segments saved to {output_file}")
49
50 def process_segments(file_path, output_csv='cleaned_segments.csv'):
51     """Process the segmented text file."""
52     segments = read_segmented_file(file_path)
53     if segments is None:
54         return
55
56     cleaned_segments = [clean_text(segment) for segment in segments]
57     save_to_csv(segments, cleaned_segments, output_csv)
58
59     # Print a preview of cleaned segments
60     print("\nPreview of Cleaned Segments:")
61     for idx, (orig, clean) in enumerate(zip(segments, cleaned_segments), start=1):
62         print(f"\nSegment {idx} (Original): {orig}")
63         print(f"Segment {idx} (Cleaned): {clean}")
64
65 # File paths
66 segmented_file_path = '/content/segmented_output.txt'  # Replace with your segmented text file path
67 output_csv = '/content/cleaned_segments.csv'  # CSV file for cleaned segments
68
```

```
69 # Run the process
70 process_segments(segmented_file_path, output_csv=output_csv)
71
```

Cleaned segments saved to /content/cleaned_segments.csv

Preview of Cleaned Segments:

Segment 1 (Original): Segment 1:
  - so, sir, we know that India has seen a huge Revolution with digital payments.
  - but then we have now shown that digital payments number one is India, while people thought that it wouldn't even come to top 5
  - sir, my question is: do you think ncrf plus NEP put together will be the next big revolution after UPA in India?
  - know we are, we are moving very fast.
  - wonderful, sir.
  - we all know a lot of scandals.
  - so that has been provided here.
  - not only that India is grappling with this problem.
Segment 1 (Cleaned): segment 1 sir know india seen huge revolution digital payments shown digital payments number one india peopl

Segment 2 (Original): Segment 2:
  - we all thought that India is a place- at least the West thought- that India is a place where, uh, many people do not get a squ
  - and not many are literate- people cannot read, um.
  - how do we do this?
  - you have place for everything.
  - okay, you look at the kind of the, the requirement of infrastructure India has today.
  - so here is where I have um, um, um thought, with some inhibitions, about the entire setup, where you are simply assuming a lot
  - so that I know the student is learning, the student knows that here is a genuine place and the person there also knows that we
  - the university is not giving me what I thought it would.
  - we have, uh, which is not appearing at the surface level, which is the student teacher ratio.
  - it is uh, it is really uh, against the students, right, right, already, it's already so, but then, uh, the colleges should not
  - so here's where I think there's a foot for Thought, where I will, uh, try to slightly disagree.
  - I I see that as as this is a wonderful thought, but I think the students wants to see those short videos not in these alignmer
  - that is where the problem lies.
  - that that is where the problem lies.
  - that is, the hecaton is is all about.
  - that's a very good thought.
  - why is it that we have not thought about all these things worldwide?
  - so for why see why this was not not even thought about?
  - we did that, mhm, and that has been done not only once but at least twice, multiple times.
  - that is the framework.
  - it gives me some non-trivial money, at least enough for me to eat and stay in a place.
Segment 2 (Cleaned): segment 2 thought india least west india place many people get square meal right many people read um place e

Segment 3 (Original): Segment 3:
  - that was the narrative some 30 years ago.
  - absolutely, and why I think so is because in education, the last policy came up many, many years ago- that was in 1986, which
  - I was trying to answer.
  - that was the second part which I wanted to answer.
  - it was.
  - Manan is going to be 30 seconds only.
  - and even in that 7 and a half minutes, that concept is broken into three parts of 2 minutes each, followed by 30 Seconds of so
  - earlier this was not there.
Segment 3 (Cleaned): segment 3 narrative 30 years ago absolutely think education last policy came many many years 1986 slightly 1

Segment 4 (Original): Segment 4:
  - I think, immediately after UPA, the next big revolution, personally, I think, is in education, and the complete homework for 1
  - so so many changes have happened in the real world.
  - so many changes have happened in the requirement of the industry, requirement of Manpower for the industry.
  - so, therefore, I feel that this was the right time when we brought in the education policy 2020, honorable prime minister, dec
  - we call it a framework because this is very flexible.
  - that is such an enabling layer that it has broken the shackles which were there in the education sector, right?

```
 1 import json
 2 import re
 3
 4 def load_transcript_from_file(file_path):
 5     """
 6     Loads the transcript from a file.
 7     Args:
 8         file_path (str): Path to the transcript file.
 9     Returns:
10         list: List of dictionaries with 'start', 'end', and 'text' keys.
11     """
12     try:
13         with open(file_path, "r", encoding="utf-8") as file:
14             if file_path.endswith(".json"):
15                 transcript = json.load(file)
16             else:
17                 # For plain text files, assume each line is in the format: [start] - [end]: [text]
18                 transcript = []
19                 for line in file:
20                     match = re.match(r"(\d+\.\d+) - (\d+\.\d+): (.+)", line.strip())
21                     if match:
22                         start, end, text = match.groups()
23                         transcript.append({
24                             "start": float(start),
```

```python
25                    "end": float(end),
26                    "text": text
27                })
28        return transcript
29    except Exception as e:
30        print(f"Error loading transcript: {e}")
31        return None
32
33 def save_transcript_with_timestamps(transcript, output_path="transcript_with_timestamps.txt"):
34    """
35    Saves the transcript with timestamps to a text file.
36    Args:
37        transcript (list): List of dictionaries with 'start', 'end', and 'text'.
38        output_path (str): The path to save the output file.
39    """
40    with open(output_path, "w", encoding="utf-8") as file:
41        for segment in transcript:
42            file.write(f"{segment['start']} - {segment['end']}: {segment['text']}\n")
43    print(f"Transcript with timestamps saved to {output_path}")
44
45 # Example usage
46 if __name__ == "__main__":
47    # Ask the user for the transcript file path
48    transcript_path = input("Enter the path to the transcript file: ")
49    output_path = input("Enter the path to save the transcript with timestamps (default: transcript_with_timestamps.txt): ") or "tra
50
51    # Load the transcript from the file
52    transcript = load_transcript_from_file(transcript_path)
53    if not transcript:
54        print("Failed to load transcript. Exiting.")
55        exit()
56
57    # Print the transcript with timestamps
58    print("\nTranscript with Timestamps:")
59    for segment in transcript:
60        print(f"{segment['start']} - {segment['end']}: {segment['text']}")
61
62    # Save the transcript with timestamps to a file
63    save_transcript_with_timestamps(transcript, output_path)
```

```
Enter the path to the transcript file: https://youtu.be/sK8SILOM37I
Enter the path to save the transcript with timestamps (default: transcript_with_timestamps.txt):
Error loading transcript: [Errno 2] No such file or directory: 'https://youtu.be/sK8SILOM37I'
Failed to load transcript. Exiting.


Transcript with Timestamps:
--------------------------------------------------------------------
TypeError                                  Traceback (most recent call last)
<ipython-input-31-41113dfd0e0a> in <cell line: 0>()
     57        # Print the transcript with timestamps
     58        print("\nTranscript with Timestamps:")
---> 59        for segment in transcript:
     60            print(f"{segment['start']} - {segment['end']}: {segment['text']}")
     61

TypeError: 'NoneType' object is not iterable
```

```python
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 import yt_dlp
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
```

```python
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': '/content/cookies (2).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:
```

```python
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 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 f6e09c70
[youtube] sK8SILOM37I: Downloading tv player API JSON
[info] sK8SILOM37I: Downloading 1 format(s): 251
[download] Destination: audio.webm
[download] 100% of    39.37MiB in 00:00:01 at 31.86MiB/s
[ExtractAudio] Destination: audio.mp3
Deleting original file audio.webm (pass -k to keep)
Transcribing audio in chunks...

Transcript with Timestamps:
0 - 30: so sorry we know that India has seen a huge Revolution with digital payments we all thought that India is a place at leas
30 - 60: I think immediately after you play the next big revolution personally I think is an education and the complete homework
60 - 90: do you think ncrf plus any people together will be the next big revolution after up in India absolutely and why I think
90 - 120: how many changes have happened in the real world so many changes have happened in the requirement of the industry requi
120 - 150: July 2020 and we recently celebrated the 4th anniversary of NP 2020 ncrf has been brought to implement the intent of 2
150 - 180: this allows you all the Innovation the way you educate your kids you educate your students still it provides you the b
180 - 210: broken the shackles of which were there in the education sector right so yes it is a big Revolution and this is going
210 - 240: 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 char
240 - 270: we have been driving buses at Max now we should fly how do we do this look at the requirement of the industry requiren
270 - 300: link that whatever he has been taught has no relevance to the real life world when industry is moving that fast when t

```
300 - 330: will you be not only beneficial but also very facilitated for all of us very liberating for all of us this is going to
330 - 360: new things create new ways of doing things learn something new but once we learn it there is no limit to Innovation ar
360 - 390: creating our vision and Innovative Minds into the education sector where is we are applying it elsewhere everywhere el
390 - 420: to be fun for everyone and I can I can tell you that already a number of Institutions have adopted the any pain and CF
420 - 450: 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 r
450 - 480: it allows for creditor of all learnings weather in academics orange killing or an experiential learning and all these
480 - 510: and people who are already skilled or already in the professional area not there the experiential learning would play
510 - 540: increase of technology which has been created by single Department good question so therefore all of us we have to wor
540 - 570: all kinds of learning are being contractors including learning of soft skills employability skills life skills your ha
570 - 600: 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
600 - 630: 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
630 - 660: turn off infrastructure India has today we created in last 1775 years we are going to double that infrastructure in ne
660 - 690: more number of other branches even liberal arts social sciences if I want to be an award Society I would need a proper
690 - 720: have you already seen the Fallout of this know you can see the photo you can see how many Engineers are there for Desi
720 - 750: [Unintelligible]
750 - 780: any of the new technology machines current any laser-based machines any automated operating machines robotic operatior
780 - 810: robotic process and that is killing is equally important and this is important in multiple areas and therefore multidi
810 - 840: me and you want to know design a VTech in CSC syllabus or be taking AI syllabus that is Nip complaint under the ncrf i
840 - 870: are you teaching Teddy teaching Terry is not sufficient if you want the student to really understand and reply that cc
870 - 900: those horses are skill bass courses you divide every subject into 30 and its application how do you apply the theory a
900 - 930: who is learning which is happening which which you are going through so then this looks like let's say my student stay
930 - 960: extension of BSC physics or BSC chemistry and you give him the actual knowledge of computer science right in the first
960 - 990: that's unbelievable so that 50% of the time which means two full years and a btech program a person can stay outside 1
990 - 1020: are you learning outcomes and their alignment with the overall curricular structure and then once it comes back we ha
1020 - 1050: he gets the credits okay so here is where I have talked with some inhibitions about the entire setup where you are s
1050 - 1080: write although we can keep that a check it is not easy for us to keep some zones green some zones red operational cc
1080 - 1110: what is the guarantee that the student is learning in the campus is there a is there an accident on some kind of a g
1110 - 1140: either online or with some time stamps it is it is being documented know what time is done and CVT has videographed
1140 - 1170: 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
1170 - 1200: Julie appointed by the awarding body which Awards the certificate and then that video is kept forever so you imagine
1200 - 1230: 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 c
```

```
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 3.1 MB/s eta 0:00:00
Downloading yt_dlp-2025.2.19-py3-none-any.whl (3.2 MB)
                                      3.2/3.2 MB 36.3 MB/s eta 0:00:00
Installing collected packages: yt-dlp
Successfully installed yt-dlp-2025.2.19
```

```
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
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-transcrip
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests->youtube-transcript-api) (3.16
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests->youtube-transcript-api
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests->youtube-transcript-api
Downloading youtube_transcript_api-0.6.3-py3-none-any.whl (622 kB)
                                      622.3/622.3 kB 9.6 MB/s eta 0:00:00
Installing collected packages: youtube-transcript-api
Successfully installed youtube-transcript-api-0.6.3
```

```
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 34.0 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 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 3.2 MB/s eta 0:00:00
Installing collected packages: pytube
Successfully installed pytube-15.0.0
```

```python
1  import nltk
2  from nltk.tokenize import sent_tokenize
3  from sklearn.feature_extraction.text import TfidfVectorizer
4  from sklearn.metrics.pairwise import cosine_similarity
5  import pandas as pd
6
7  # Download NLTK data
8  nltk.download('punkt')
9
10 def read_transcript(file_path):
11     try:
12         with open(file_path, 'r', encoding='utf-8') as file:
13             transcript = file.read()
14         return transcript
15     except Exception as e:
16         print(f"Error reading file: {e}")
17         return None
18
19 def split_into_sentences(transcript):
20     return sent_tokenize(transcript)
21
22 def calculate_cosine_similarity(sentences):
23     vectorizer = TfidfVectorizer()
24     tfidf_matrix = vectorizer.fit_transform(sentences)
25     cosine_sim = cosine_similarity(tfidf_matrix)
26     return cosine_sim
27
28 def segment_sentences(sentences, cosine_sim, threshold=0.5, min_sentences=5):
29     visited = [False] * len(sentences)
30     segments = []
31
32     for i in range(len(sentences)):
33         if not visited[i]:
34             segment = [sentences[i]]
35             visited[i] = True
36
37             for j in range(i + 1, len(sentences)):
38                 if not visited[j] and cosine_sim[i][j] >= threshold:
39                     segment.append(sentences[j])
40                     visited[j] = True
41
42             segments.append(segment)
43
44     # Merge smaller segments
45     merged_segments = []
46     temp_segment = []
47
48     for segment in segments:
49         temp_segment.extend(segment)
50         if len(temp_segment) >= min_sentences:
51             merged_segments.append(temp_segment)
52             temp_segment = []
53
54     if temp_segment:
55         if merged_segments:
56             merged_segments[-1].extend(temp_segment)
57         else:
58             merged_segments.append(temp_segment)
59
60     return merged_segments
61
62 def print_segments(segments):
63     for idx, segment in enumerate(segments, start=1):
64         print(f"\nSegment {idx}:")
65         for sentence in segment:
66             print(f"  - {sentence}")
67
68 def save_segments_to_file(segments, output_file):
69     try:
70         with open(output_file, 'w', encoding='utf-8') as file:
71             for idx, segment in enumerate(segments, start=1):
72                 file.write(f"Segment {idx}:\n")
73                 for sentence in segment:
74                     file.write(f"  - {sentence}\n")
75                 file.write("\n")
76         print(f"Segmented output saved to {output_file}")
77     except Exception as e:
78         print(f"Error saving file: {e}")
79
80 def process_transcript(file_path, threshold=0.5, min_sentences=5, output_file='segmented_output.txt'):
81     transcript = read_transcript(file_path)
```

```
 82     if transcript is None:
 83         return
 84
 85     sentences = split_into_sentences(transcript)
 86     print("Sentences extracted:")
 87     for i, sentence in enumerate(sentences):
 88         print(f"S{i + 1}: {sentence}")
 89
 90     cosine_sim = calculate_cosine_similarity(sentences)
 91
 92     segments = segment_sentences(sentences, cosine_sim, threshold, min_sentences)
 93
 94     print("\nSegmented Sentences:")
 95     print_segments(segments)
 96
 97     save_segments_to_file(segments, output_file)
 98
 99 file_path = '/content/transcript.txt'  # Replace with your file path
100 output_file = '/content/segmented_output.txt'
101 process_transcript(file_path, threshold=0.15, min_sentences=5, output_file=output_file)
```