

Scientext: Streamlining Access to Educational Content and Research Insights through Advanced Speech and Text Processing

github: <https://github.com/ayesha-mohsin/Scientext>

Shaik Ayesha Mohsin || Areeba Hassan

Abstract:

"Scientext" embarks on a pioneering journey to redefine the landscape of educational and scientific engagement. At its core, Scientext is dedicated to two primary endeavors: the precise transcription of TED Talks through advanced speech recognition, and the meticulous summarization of scientific papers utilizing the latest in natural language processing (NLP) technologies. This ambitious project aims to dismantle the barriers that restrict access to the vast reservoirs of knowledge contained within these mediums. By converting spoken words into text and distilling complex scientific research into digestible summaries, Scientext aspires to serve a dual purpose: to democratize information access and to enhance the comprehension of diverse audiences. In doing so, Scientext is not just facilitating knowledge dissemination but is fostering a culture of inclusivity and innovation across the global educational and scientific communities.

1. Motivation:

In the digital era, the proliferation of educational content and scientific research has been exponential, offering boundless opportunities for learning, discovery, and innovation. Yet, this wealth of information often remains ensconced within realms that are not easily navigable by all, hindered by linguistic complexities, technical jargon, and the sheer volume of available material. Scientext is conceived from the recognition of this paradox—the simultaneous abundance and inaccessibility of knowledge. It aims to bridge this divide by leveraging cutting-edge speech recognition and NLP technologies to automate the process of making educational talks and scientific findings not just accessible but also comprehensible. The project is driven by a vision where everyone, regardless of their academic background or expertise, can tap into the full potential of these invaluable resources for personal growth, academic advancement, and societal contribution.

2. Significance:

The implications of Scientext are profound and far-reaching. By streamlining the process of accessing and understanding educational and scientific content, Scientext positions itself as a catalyst for intellectual empowerment and interdisciplinary collaboration. Its capacity to translate spoken discourse into text and to condense intricate research into concise summaries has the potential to revolutionize how knowledge is consumed and applied. For educators and students, it promises a more interactive and engaging learning experience. For researchers and professionals, it offers a tool for swiftly grasping new developments outside their fields of expertise. Furthermore, Scientext stands to play a crucial role in leveling the educational playing field, providing underserved communities with the means to partake in the global conversation on science and education. Ultimately, Scientext is about more than just facilitating access; it's about igniting curiosity, inspiring innovation, and fostering a more informed and interconnected world.

By embodying the principles of accessibility, comprehension, and collaboration, Scientext not only anticipates the needs of its users but also endeavors to anticipate the evolving landscape of education and research. As such, it represents not just a step forward in the application of technology to learning and discovery but a leap towards a future where the gates of knowledge are open to all.

3. Objectives:

- TED Talk Speech Recognition and Text Generation:
 - Develop a robust speech recognition system capable of transcribing TED Talks with high accuracy, leveraging the TED-LIUM 3 dataset for model training and validation.
 - Employ advanced models such as DeepSpeech 2 or wav2vec 2.0, fine-tuned on the TED-LIUM dataset, to ensure the system can handle diverse accents and specialized terminology found in TED Talks.
 - Generate searchable, editable text from transcriptions that can be used for further analysis and summarization.
- Scientific Paper Summarization:
 - Implement and train state-of-the-art NLP models for summarizing scientific papers, with a focus on transformer-based architectures like BERT for extractive summarization and GPT-3 for abstractive summarization.
 - Utilize datasets from arXiv and PubMed, ensuring models are well-versed in a variety of scientific domains and capable of identifying key findings and concepts.
 - Create summaries that are informative, concise, and accessible to non-specialists, facilitating a broader understanding of complex research findings.
- Integration and Usability:
 - Seamlessly integrate the speech recognition and text summarization components into a cohesive, user-friendly platform.
 - Ensure the platform supports easy submission of TED Talk videos and scientific papers, and provides an intuitive way for users to access the generated transcriptions and summaries.
 - Conduct rigorous testing and user studies to validate the system's effectiveness, usability, and accessibility, adjusting based on feedback to meet user needs.

4. Features:

- Advanced Speech Recognition:
 - Utilization of state-of-the-art models such as DeepSpeech 2 or wav2vec 2.0, optimized for the unique challenges of TED Talk transcriptions, including diverse topics and diction.
 - Integration of noise-reduction and audio-enhancement techniques to improve transcription accuracy in varied acoustic environments.
- Sophisticated Summarization:
 - Employment of transformer-based models like BERT for extracting critical information and GPT-3 for generating coherent and engaging summaries of scientific papers.

- Customization of summarization techniques to adjust the length and detail of summaries based on user preferences and the specific requirements of different scientific domains.
- Seamless User Interface:
 - Development of a web-based platform that is simple to navigate, allowing users to easily upload TED Talk videos and scientific papers.
 - Implementation of features such as real-time progress tracking, options to customize summary length and style, and the ability to download or share summaries directly from the platform.
- Integration and Workflow Optimization:
 - Creation of a backend pipeline that efficiently processes input videos and documents, ensuring timely transcription and summarization.
 - Use of cloud-based services and APIs to scale the processing capabilities of the platform, accommodating a high volume of concurrent requests.

5. Dataset

- TED-LIUM 3: This dataset provides a vast resource of over 450 hours of TED Talk audio specifically designed for speech recognition tasks, making it an invaluable asset for developing robust transcription systems. The transcripts for these TED Talks are sourced directly from TensorFlow's datasets catalog at (<https://www.tensorflow.org/datasets/catalog/tedlium>), ensuring easy access and integration into speech recognition models.
- Scientific Papers: For the summarization component of our project, we will utilize a curated collection of scientific papers. This collection has been meticulously compiled to span various domains of scientific inquiry, offering a broad spectrum of research findings and discussions. The dataset is readily accessible for download from a dedicated Google Drive link: (<https://drive.google.com/file/d/1RPsgg32xLBUwbMr3eV0uTqGuVgUQSGzu/view>). This dataset focuses on the abstracts and content of scientific papers, providing a rich source of material for training and evaluating our text summarization models.

6. Workflow

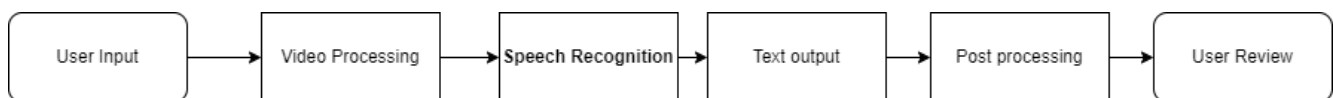


Fig1.Workflow Diagram

Workflow Diagram for Speech Recognition

- User Input:
 - Users access the Scientext platform and submit a TED Talk video link for transcription.
- Video Processing:
 - The system extracts the audio track from the submitted TED Talk video.
- Speech Recognition:
 - The extracted audio is processed by the speech recognition system.
 - Dataset Used: TED-LIUM 3 dataset serves as the training and validation source for the speech recognition model, enhancing its ability to accurately transcribe TED Talks.

- Text Output:
 - The speech recognition model transcribes the audio into text, producing a written version of the TED Talk.
- Post-Processing:
 - The transcribed text undergoes post-processing, including punctuation correction and paragraph formatting, to improve readability.
- User Review:
 - The transcribed TED Talk is displayed to the user on the Scientext platform for review, download, or further action.

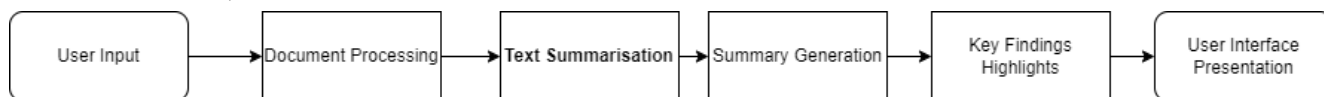


Fig2.Workflow Diagram

Workflow Diagram for Scientific Paper Summarization

- User Submission:
 - Users navigate to the Scientext platform and upload or link to a scientific paper they wish to summarize.
- Document Processing:
 - The system processes the submitted document, extracting the text content from the paper.
- Text Summarization:
 - The extracted text is fed into the text summarization model.
 - Dataset Used: The Scientific Paper Summary Corpus, accessible via the provided Google Drive link, is employed to train and fine-tune the summarization model. This ensures the model is adept at identifying and condensing key findings and concepts from a wide array of scientific disciplines.
- Summary Generation:
 - The model generates a concise summary of the scientific paper, highlighting essential findings, methodologies, and conclusions.
- Key Findings Highlighting:
 - The system enhances the generated summary by emphasizing key findings and significant insights, making them easily accessible to the user.
- User Interface Presentation:
 - The summary, along with highlighted key findings, is presented to the user on the Scientext platform. Users can review, save, or share the summary as needed.

7. Conclusion

Scientext aims to be a groundbreaking tool in making educational content and scientific research more accessible and understandable. By leveraging advanced technologies in speech recognition and text summarization, Scientext will empower individuals from various backgrounds to engage with and benefit from the wealth of knowledge available in TED Talks and scientific literature.

