# AudioBookify: Bridging Text to Sound for the Visually Impaired

Yuraja Kadari, Yash Jivani, Dipen Patel, John Melo, Shrutika Gajbhiye, Vicky Jadhav
*Seidenberg School of CSIS, Pace University, New York City, New York.*
*{yk85418n, yj38882n, dp54214n, jm33812n, sg30238n, vj54818n, }@pace.edu*

*Abstract: This project addresses the pressing need for universal access to written content in the digital age. With over 253 million people worldwide experiencing vision impairment [11] (WHO – World Health Organization), the need for universal access to literature and educational materials is important. "AudioBookify" addresses this challenge by converting images and PDFs into audiobooks. First, Optical Character Recognition (OCR) technology transforms diverse forms of written material into text format. And then using Text to Speech (TTS) technology, it converts text into audio format. This project not only empowers visually impaired individuals but also caters to auditory learners and busy individuals.*

*Index Terms— Universal Access, Optical Character Recognition (OCR) Technology, Text-to-Speech (TTS) Technology, Visual Impairment, Auditory Learners, Literature Accessibility, Audio Books, Advanced Technology*

## I. INTRODUCTION

THE digitally driven world is composed of written content that serves as a fundamental pillar for learning, communication, and entertainment. Yet, accessing written material poses daunting obstacles for individuals with visual impairments or those who favor auditory learning. To bridge this accessibility gap, our project, Audiobookify, endeavors to craft a comprehensive solution that seamlessly converts images and PDFs into audiobooks. By catering to the needs of visually impaired individuals and those who prefer listening over reading, Audiobookify aims to revolutionize access to written content.

Powered by advanced Optical Character Recognition (OCR) technology, Audiobookify conducts intricate analyses of images and PDFs, discerning and extracting text, symbols, tables, diagrams, and coding language content. This sophisticated OCR process facilitates the transformation of diverse forms of written material into accessible text format. Subsequently, the extracted text undergoes audio format synthesis, making an audiobook readily accessible to users. The outcome of the project is transformative: visually impaired individuals gain autonomy in navigating literature and educational resources, while auditory learners and busy individuals benefit from a convenient alternative to traditional reading methods. Audiobookify sets a standard for universal access to written content, fostering inclusivity and equal opportunities for knowledge acquisition in the digital era.

## II. BUSINESS OBJECTIVE

**Accessibility:** Our main goal is to make sure that everyone, including people with vision problems or those who prefer listening to reading, can easily access written content. This helps to include everyone and gives everyone an equal chance to learn.

**Users:** We want all types of people to use our solution. By making it unique and helpful for many kinds of people, like the visually impaired, those who like listening instead of reading, and otherwise busy people, we aim to reach a big audience.

**User Satisfaction:** We want our users to be content with our service. That means we'll make sure it's easy to use and that the audiobooks are always available when needed.

**Partnerships and Collaborations:** We can work with organizations like libraries, schools, and groups that help people with vision problems. By teaming up with them, we can offer more audiobooks and reach more people who need our service.

## III. SOFTWARE REQUIREMENTS

**Optical Character Recognition (OCR) Technology:**

We used advanced OCR technology -Tesseract library which can accurately analyze images and PDFs to extract text, symbols, tables, diagrams, and coding language content.

**Text-to-Audio Synthesis:** We implemented gTTS, robust text-to-audio synthesis technology to convert extracted text into high-quality audiobooks, ensuring clear pronunciation and natural-sounding narration.

## IV. LITERATURE REVIEW

**Origins of Text-to-Speech Technology:** Text-to-speech synthesis was first attempted in 1779 by the Russian professor Christian Kratzenstein. Kratzenstein documented and explained the difference between the sounds of the five long vowels (A/E/I/O/U). To mimic these sounds, he created a machine that mirrored the structure of human vocal system [1]. This machine used vibrating reeds to mimic the sounds of these vowels. It is important to note that the I vowel was not produced with a vibrating reed. It was produced by blowing air into the machine like a flute. Text to speech synthesis was later attempted in Vienna in 1791 by the "Acoustic-Mechanical Speech Machine" invented by Wolfgang von Kempelen. Like Kratzenstein's

machine, Von Kemplen's machine mimicked the human vocal system and was able to produce individual sounds and a few sound combinations [1].

Text-to-speech technology (TTS) took a significant leap in 1968 with the development of the first full text-to-speech English system by Noriko Umeda in the Electrotechnical Laboratory in Japan. This system used an articulatory model to articulate sounds along with a syntactic analysis module. This allowed the system to generate intelligible English speech. However, the speech that was generated did not include prosody and was monotonous [1]. Finally, in 1939 the first speech synthesizer known as VODER (Voice Operating Demonstrator) was introduced by Homer Dudley. The VODER synthesizer worked by analyzing speech into acoustic parameters which would synthesize a speech signal. The signal would then pass through ten filters where the output was manipulated by fingers. Although the output of the VODER synthesizer was mediocre, this device set the foundation for modern-day TTS technologies [1]. Currently, many institutions like universities have limited resources for visually impaired students. To digest information, visually impaired individuals rely on braille or on verbal translation from their peers. To help visually impaired individuals digest information more efficiently, a text-to-audio converter system will need to be designed. To design this system, we will need to identify factors such as the appropriate API, the type of texts (books, articles, PDFs) the system will convert, and design the database for storing the audible material.

**Advancements in Text-to-Speech technology:** Audiobooks have been proven to be just as effective as traditional reading methods in helping people retain information. The University of Oregon conducted an experiment that showed comparable competencies between listening to and reading materials. The experiment showed that the test subjects remembered and forgot the same portions of the material regardless of the medium. The reason for this is because the part of the brain that is responsible for language comprehension works the same regardless of whether the individual is listening to or reading the material [2]. Since audiobooks are comparable to traditional reading for information retention, audiobooks are especially beneficial to individuals with vision impairments. Advancements in text-to-speech technology (TTS) have single-handedly led to the recent success of audiobooks. Recent advancements in TTS deep learning have improved the ability to recreate natural – sounding audio that mimics human speech. For example, it was recorded that the United States experienced annual audiobook sales revenue of about $1.3 billion in the year 2020 alone [3]. Additionally, by using advanced text-to-speech technology, the National Library Service for the blind has expanded their audiobook library by recording tens of thousands of books. Despite these recent TTS advancements, TTS systems still have drawbacks when it comes to prosody and intonation variation. For example, machine – generated audio stays within a small constant range of pitch and volume, making the speech monotone. To develop a more natural prosody mechanism, the system must contain text-character/prosody analysis prediction capabilities. Text-

character/prosody analysis will determine how text-characters are read and the pitch at which they are read [3].

Most text-to-speech synthesizers contain certain modules that perform separate important functions. One module is called the Natural Language Processor (NLP). The NLP creates a phonetic transcription of the chosen text with prosody. The second module is called the Digital Signal Processor (DSP). The DSP module converts symbolic text from the NLP module into an intelligible audio medium [4]. In addition, NLP has two main functions. The first main function is text analysis. Text analysis works by segmenting text into tokens. The second main function is implementation of pronunciation rules. Once text analysis is successfully completed, pronunciation rules can be implemented. Finally, the product of the NLP module is sent to the Digital Signal Processor (DSP) module for processing [4]. Advancements in TTS technology have made audiobooks a staple resource for people with visual impairments. Experiments have shown that there is no difference between information retention rates from reading text versus listening to text. Having an adequate TTS system prevents visually impaired individuals from being at a disadvantage. To develop an efficient TTS system, we will need to design our natural language processor module along with the digital processor module.

**OCR Technology:** Machine-learning is a vital form of technology in the development of Audiobook systems. Machine-learning can be used to convert speech to text, text to speech, and even identify objects in an image [5]. Optical Character Recognition (OCR) technology is an example of machine-learning that is fundamental to audiobook systems. Optical Character Recognition is a method that converts any form of text into a modifiable digital format. Optical character recognition works by prompting a device to detect text in documents and/or images. Optical character recognition has some dependents, such as font, image quality, and input document quality.

The process of optical character recognition is made up of six main phases: preprocessing, segmentation, normalization, feature extraction, classification, and postprocessing [6]. Preprocessing is the first OCR phase that prepares the image and/or text for the segmentation phase by converting the image/text into a convertible format. The preprocessing phase will employ methods such as compression, filtering, and slant removal. Segmentation is the second OCR phase in which text is identified and isolated in an image. This is otherwise known as document segmentation. Normalization is the third OCR phase in which an image is converted into a m*n matrix. This matrix will assist in eliminating redundant data from document. Feature extraction is the fourth OCR phase in which feature vectors are built by extracting relevant data from alphabets and/or objects by using a classifier [6]. Classification is the fifth OCR phase in which inputs are classified and group into homogenous classes based on the qualities of the input. This is done to facilitate pattern recognition methodologies necessary for the functionality of OCR technology. Lastly, postprocessing is the sixth OCR phase. The postprocessing phase seeks to

identify and correct errors in the OCR's transcription. An example would be spell-check.

## V. SURVEY

Before embarking on our project's development phase, we recognized the necessity of understanding the prevailing trends and preferences in audiobook consumption during the literature review research. We designed a Google Form questionnaire to gather firsthand insights and distributed it among a diverse sample population. The data collected from this survey provided valuable insights into contemporary reading habits and preferences. One striking revelation was that a significant majority, comprising 88 percent of respondents, expressed a preference for audiobooks over physical books.



*fig 1: Data Visualization of Reading vs Listening.*

This overwhelming preference for audiobooks served as a catalyst, instilling confidence in us to delve deeper into this burgeoning field. With this initial validation in hand, we conducted further analysis of the dataset, uncovering additional insights that informed our project direction and strategy. These insights not only affirmed the relevance and significance of audiobooks in modern literature consumption but also shed light on nuanced preferences across different demographics and occupations, providing invaluable guidance for our project development journey.

**Audiobook Preference Across Gender:** We observed that both males and females exhibit a notable interest in audiobooks, although preferences may vary slightly between the genders. We also discovered that males are more likely to prefer physical books than females.
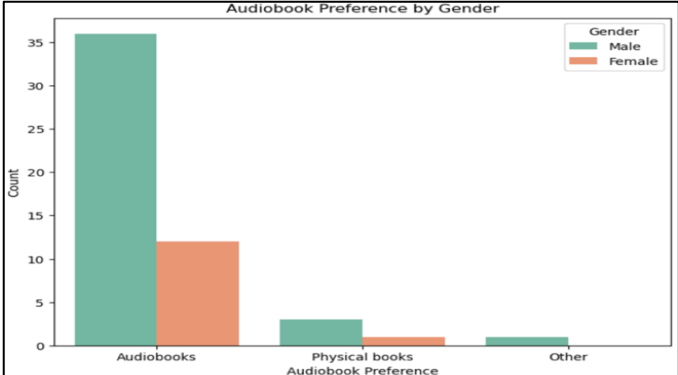


*fig 2: Bar Chart Audiobook Preference by Gender*

**Age-Genre Preference Correlation:** Our analysis uncovered interesting correlations between age groups and preferred book genres. Younger individuals often show a preference for contemporary genres, while older age groups may lean towards classic or literary genres.
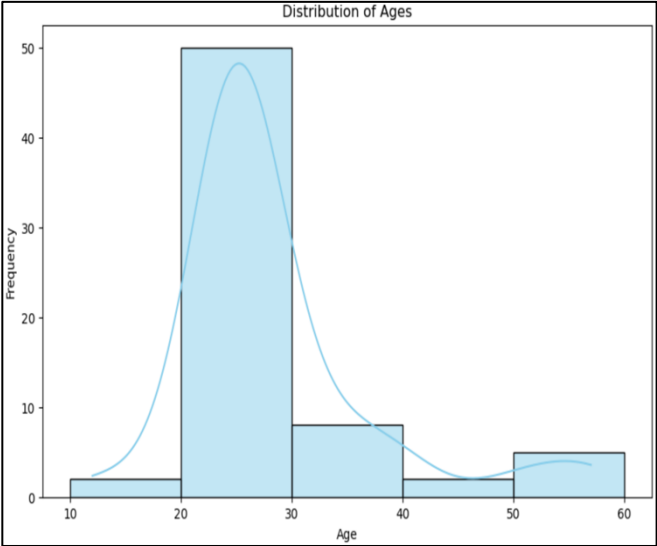


*fig 3: Bar Chart Distribution of Age*

**Occupation and Reading/Listening Frequency**: The frequency of reading or listening to books varies across different occupations. Professionals in certain fields may find audiobooks more convenient due to time constraints, while others may prefer the traditional reading experience.
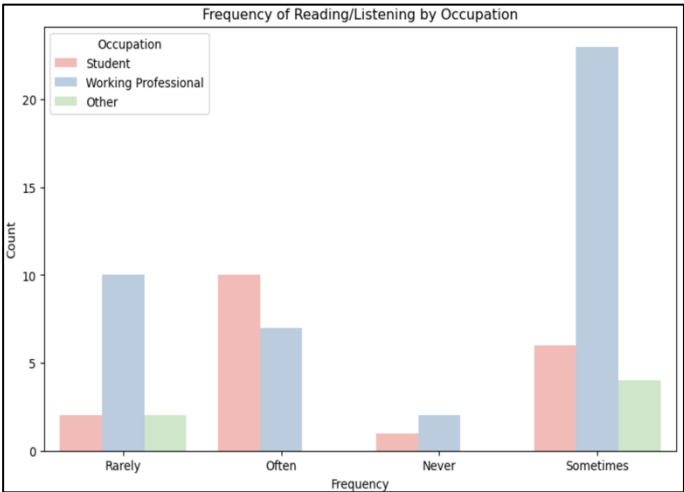


*fig 4: Distribution of Occupation*

**Importance of Audiobooks:** Our findings suggest that audiobooks play a significant role in modern information consumption habits, with a substantial portion of the population showing interest in this format. The pie chart visualization highlights the proportion of individuals who prefer audiobooks over traditional reading methods.
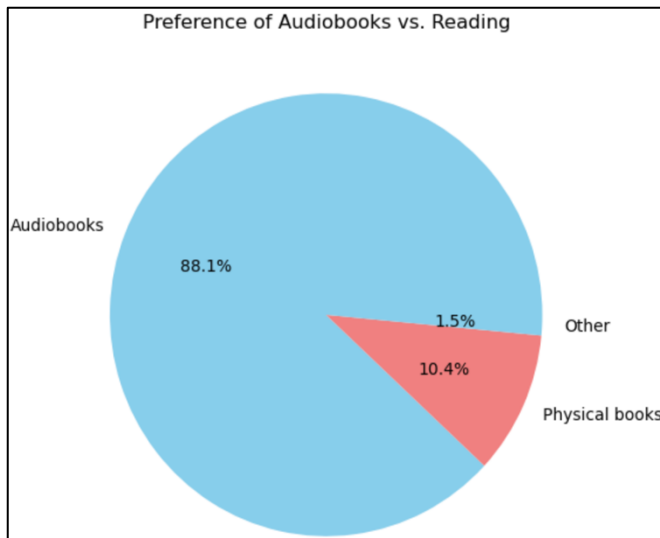
*fig 5: Pie Chart Preference of Audiobook vs Books*

**Consumption:** The pair plot visualization revealed potential relationships between audiobook preference and other numerical variables, suggesting that audiobook preference may influence how individuals consume information.

## VI. SOFTWARE METHODOLOGY

The decision to utilize Agile methodology for developing an Audiobook Management System is driven by several compelling factors. Firstly, Agile's adaptability and flexibility are paramount. Its iterative approach allows seamless adjustments to evolving requirements, ensuring that the final product aligns closely with stakeholders' expectations. Furthermore, Agile's incremental development process is particularly well-suited to the intricate nature of audiobook management systems. This method enables the gradual incorporation of features, maximizing the value delivered with each iteration while maintaining flexibility to respond to changing needs.

Additionally, Agile methodology places a strong emphasis on active stakeholder involvement, a crucial element in projects where user experience is paramount. By facilitating rapid prototyping and feedback loops, Agile ensures continuous alignment with user needs, ultimately resulting in a product that resonates deeply with its intended audience. Moreover, Agile's proactive approach to risk management enables the early identification and mitigation of potential issues. This not only minimizes project disruptions but also enhances overall project resilience, ensuring a smoother development journey and a higher likelihood of success.

Agile methodology is a strong choice for developing an audiobook management system due to its emphasis on flexibility, incremental development, stakeholder involvement, rapid prototyping, and risk management. Its iterative approach allows for seamless adaptation to evolving requirements, ensuring that the final product meets stakeholder expectations effectively. Moreover, Agile's incremental development process enables the gradual addition of features, maximizing

value with each iteration while accommodating changes smoothly.

On the other hand, methodologies like Waterfall and the V-Model may not be as suitable for audiobook management system development. These linear, sequential approaches lack the flexibility needed to address evolving requirements and user needs. Their rigid structures emphasize upfront planning and documentation, which may not align well with audiobook management system projects' iterative and dynamic nature.

While Rapid Application Development (RAD) shares Agile's focus on rapid prototyping and iterative development, it may lack the necessary structure and project management framework. This could potentially lead to challenges in managing project scope, timelines, and stakeholder expectations. Overall, Agile methodology's adaptive and collaborative approach makes it the preferred choice for developing audiobook management systems, ensuring efficient development processes and successful project outcomes.

The diagram below is the Gantt chart. It includes phases such as research, development, testing, and deployment. Key tasks include selecting libraries, developing OCR functionality, testing text-to-speech synthesis, and finalizing the user interface. Milestones include completing OCR integration and conducting user acceptance testing. The chart ensures efficient project management and timely completion of Audiobookify.
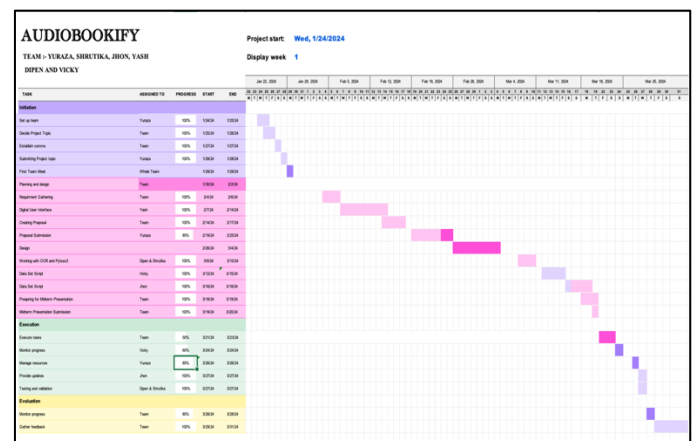


*fig 6: Gannt Chart Diagram*

## VII. PROJECT PLANNING

The Agile software development methodology has been selected for its iterative and incremental approach, offering flexibility and adaptability throughout the project lifecycle. The Agile plan unfolds over several sprints, each spanning a two-week timeframe. During the Communication and Planning phase in the first two weeks, the team establishes communication channels, defines project goals and objectives, and conducts initial sprint planning to lay the groundwork for subsequent development phases. As the project progresses, Sprint 1 focuses on requirement gathering and analysis, alongside the design of the initial architecture and user

interface, setting the foundation for subsequent development cycles. Subsequent sprints, ranging from Sprint 2 to Sprint 6, each tackle specific development tasks, including the development of core functions, integration of functionalities like Optical Character Recognition (OCR) and text-to-speech synthesis, as well as testing, deployment, and user feedback gathering.

Finally, in Sprint 7, user feedback is incorporated into updates, leading to the finalization of deployment, thus completing the Agile development cycle. This structured approach enables efficient progress tracking, continuous feedback incorporation, and iterative improvement throughout the project, ensuring the delivery of a high-quality Audiobook Management System that meets stakeholders' expectations.
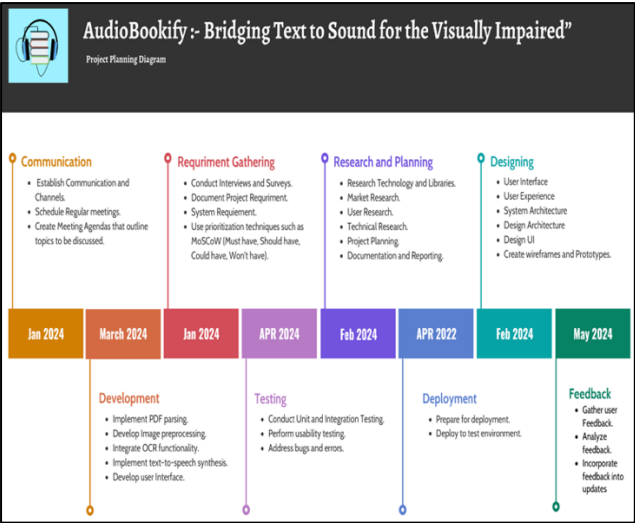


*Fig 7: Project Schedule*

The above diagram shows the timeline of the project. It shows the software development lifecycle steps that we follow while developing the software.

## VIII. DEVELOPMENT

Using a multi-phase approach, the project development phase aims at converting text from PDFs into audiobooks. This methodical procedure encompasses several steps, each leveraging various Python libraries. Specifically, we have utilized PyPDF2 for PDF parsing, Pillow for image processing, Pytesseract for OCR (image-to-text) functionality, and potentially pyttsx3 for text-to-speech synthesis. The project commences with the selection of an input file, which can either be a PDF or an image. Based on this file format, appropriate preprocessing, and OCR text extraction techniques such as converting it to grayscale or enhancing its quality are applied to the input file.

A pivotal aspect of our approach is text recognition, which extends beyond conventional text to include recognition of handwritten text, tables, diagrams, and programming language code. These additional features are intended to be seamlessly integrated into our system. Following the text extractor, text-to-speech synthesis will be employed to convert the extracted text

into audio format.

The final output will be stored as an MP3 audiobook file. The project's workflow ensures user-friendly access to audiobooks, promoting inclusivity and diversity by enhancing accessibility to various types of textual content. Fig x shows a flowchart that illustrates a detailed streamlined process of converting text from PDFs or images into audiobooks, enhancing accessibility and diversity in accessing textual content.
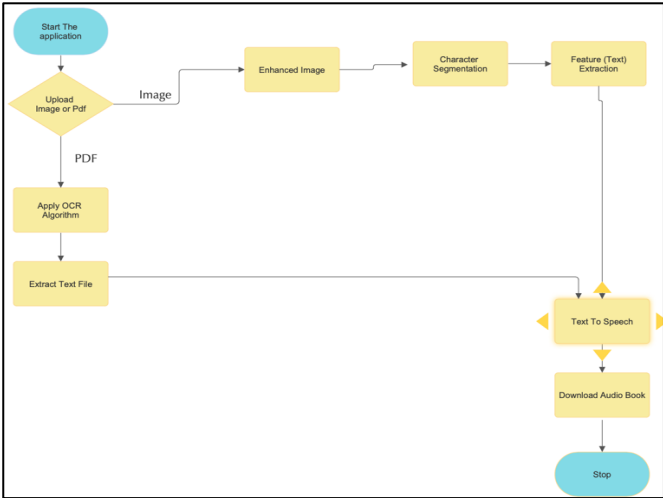


*fig 8: System Design*

The below is use case diagram about Audiobookify.
The project's use-case diagram illustrates how users interact with the Audiobookify system to convert text from PDFs or images into audiobooks. Users select an input file, which the system preprocesses, extracts text using OCR, synthesizes it into audio, and saves the audiobook. The system aims to enhance accessibility to various textual content types.
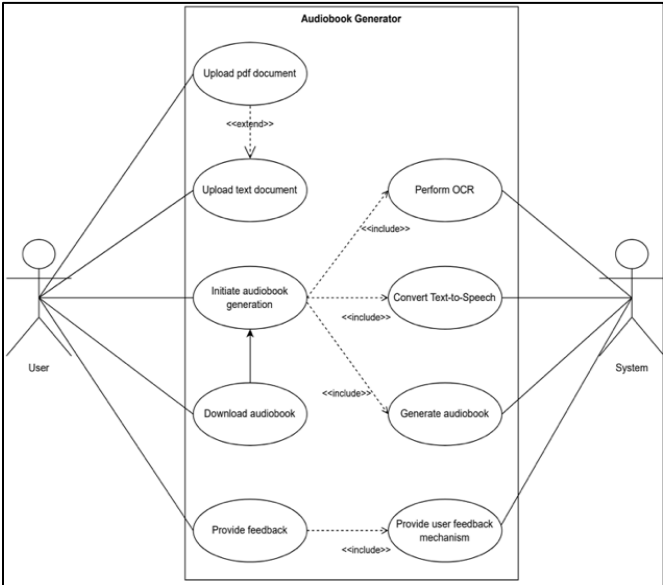


*fig 9: Use Case Diagram*

## IX. TESTING

In addition to running scripts to identify potential errors, comprehensive testing protocols have been implemented to ensure the robustness and reliability of the Audiobookify system. These integrity tests encompass a diverse array of input files in various formats, allowing us to anticipate and address any anomalies that may arise when users execute the code.

Our testing framework is designed to meticulously examine both input/output (I/O) errors and processing errors that could potentially occur due to the integration of Optical Character Recognition (OCR) and text-to-speech (TTS) APIs. By subjecting the system to rigorous tests, we have endeavored to fortify its resilience against unexpected contingencies, enhancing the overall user experience. Furthermore, all files processed during testing procedures are systematically stored in a temporary directory, serving as a contingency resource in the event of system disruptions or failures. This proactive measure facilitates troubleshooting and ensures seamless operations continuity in case of unexpected errors.

Looking ahead, we recognize the importance of optimizing system efficiency and resilience against API failures or memory constraints, particularly on devices with limited resources. As part of future enhancements, caching mechanisms and strategies for easy information retrieval will be explored to bolster system performance and mitigate potential downtimes. Additionally, the implementation of try-except blocks to handle API failures or memory errors represents a pivotal aspect of our ongoing development efforts, as we remain committed to delivering a robust and dependable solution to our users.

## X. RESULTS

The Audiobookify system successfully processed the storybook "Lost in the Fog," demonstrating the system's capacity to transform written content into an audible format. Both a text file and an audio (MP3) file were produced because of this processing, showcasing the system's capacity to convert text into a format that people can easily understand. Through an exhaustive verification procedure, the precision of the extracted text was carefully verified, guaranteeing that every piece of material from the storybook was accurately identified by the OCR technology that has been integrated into the system.

This thorough verification highlights Audiobookify's dependability and accuracy in effectively translating textual information, improving literature and educational resources' accessibility for a variety of users, including those who are visually impaired.

Fig 10 shows the PDF of the book which we processed and got text and audio file:
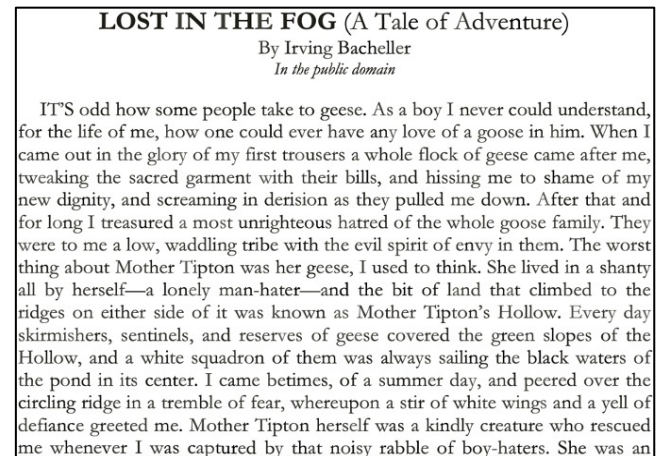


fig 10: Storybook Pdf

The figure below shows the text file and audio file we get after running the code and passing the storybook pdf.
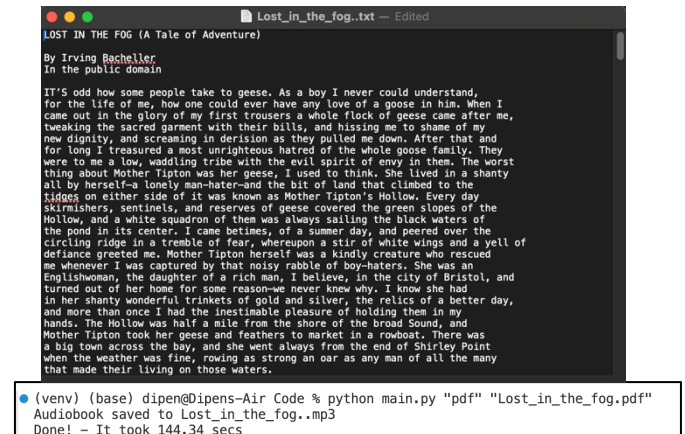


fig 11: Storybook Text file and Audio file.

## XI. FUTURE SCOPE

**AI-Powered Narration as A Customized and Immersive Experience:** AI can completely change the narration of audiobooks. Imagine listening to audiobooks where the narrator changes their tone of voice to suit your tastes or where various characters have distinct voices created by AI. This field of study might examine how AI can improve and customize audiobooks' emotional impact.

**AI Translation:** Overcoming Language Barriers for Worldwide Accessibility AI translation could let audiobooks cross language boundaries and become widely available. To ensure that audiobooks maintain their impact and emotional depth across all languages, this future scope explores the idea of AI that not only translates text effectively but also captures the soul of the narration.

**Intelligent Suggestions:** Revealing Your Upcoming Audio Obsession Automated recommendation systems for

audiobooks may be used in the future, going beyond simple suggestions. Imagine a system that determines which audiobooks exactly suit your tastes based on an analysis of your ratings, reviews, and listening history. This field of study investigates how AI might produce a more personalized and interesting audiobook experience.

**Collaboration with Libraries and Educational Institutions:** Working with libraries, educational institutions, and organizations that assist the blind and visually impaired can help Audiobookify reach a wider audience and have a greater impact. Audiobookify can be incorporated into curricular, digital libraries, and assistive technology initiatives by collaborating with educational organizations.

**Context-Aware Audiobooks:** Adapting to Your Environment Context-aware feature development could lead to an even more seamless audiobook experience. This future scope investigates the possibility of audiobooks being environment adaptive. Consider audiobooks that change the volume or pace of the narrator according to your lifestyle—at home, when working out, or while commuting.

**Creating the Framework for the Growth of Audiobooks**: The increasing demand for audiobooks calls for a strong infrastructure. This study aims to investigate the need for creative ways to store, distribute, and stream audiobooks effectively. It looks at ways to ensure that audiobook services can meet the constantly rising demand from listeners worldwide.

**Iterative development and continuous user feedback:** Audiobookify prioritizes executing iterative development cycles and obtaining user feedback to guarantee ongoing relevance and usability. By gathering feedback from individuals with visual impairments, educators, and advocates for accessibility, Audiobookify can pinpoint areas in need of development, adjust to user requirements, and introduce feature updates that are in line with the changing needs of its intended user population.

## XII. CONCLUSION

Nowadays, audiobooks are a strong and adaptable format that appeals to a variety of consumers. They provide a practical and easy way to interact with books, especially for those who have busy schedules or visual impairments. Thanks to text-to-speech technology, audiobooks are becoming a common resource for those who are blind or visually impaired, and more studies into advanced audio augmentation could make listening better for everyone.

Optical Character Recognition (OCR) systems selection and research is essential processes in creating effective text-to-speech software. By exploring OCR technologies and implementing them into Audiobookify, we have emphasized how important and accurate text extraction and synthesis are to providing a flawless user experience. We have improved the usefulness and dependability of Audiobookify by giving priority to the selection of strong OCR technologies, guaranteeing that visually impaired people may access content with accuracy and clarity.

Artificial intelligence (AI) has a bright future for audiobooks and is expected to have a revolutionary impact. The audiobook experience might be greatly improved by AI-powered features like context-aware playback, intelligent recommendations, and personalized narration. AI translation also holds the potential to reduce language barriers and increase the accessibility of audiobooks for a worldwide audience.

But in addition to the fascinating opportunities, there are challenges to consider. Because audiobooks are being widely used, a scalable infrastructure must be created to store, distribute, and stream them effectively. It will be essential to ensure that audiobooks continue to be a welcoming and educational experience for everyone as they develop.

## XIII. REFERENCES

[1] Lemmetty, Sami. Review of Speech Synthesis Technology - SPA, Helsinki University of Technology, Mar. 1999, research.spa.aalto.fi/publications/theses/lemmetty_mst/thesis.pdf.

[2] Best, Emily. Audiobooks and Literacy - A Rapid Review of the Literature. National Literacy Trust, Feb. 2020. Accessed Feb. 2024.

[3] Pethe, Chart, et al. "Prosody Analysis of Audiobooks." Cornell University, Stony Brook University, 10 Oct. 2023, https://arxiv.org/abs/2310.06930. Accessed Feb. 2024.

[4] Isewon, Itunuoluwa, et al. Design and Implementation of Text to Speech Conversion for Visually Impaired People. International Journal of Applied Information Systems, Apr. 2014. Accessed Feb. 2024.

[5] LeCun, Yann, et al. "Deep learning." Nature, vol. 521, no. 7553, 27 May 2015, pp. 436444

[6] Hamad, Karez, and Mehmet Kaya. "A detailed analysis of Optical Character Recognition Technology." International Journal of Applied Mathematics, Electronics and Computers, vol. 4, no. Special Issue-1, 22 Dec. 2016, pp. 244–244, https://doi.org/10.18100/ijamec.270374.

[7] J. O. Wobbrock, S. K. Kane, K. Z. Gajos, and S. Harada, "Ability-based design: Concept, principles and examples," ACM Transactions on Accessible Computing (TACCESS), vol. 11, no. 1, pp. 1-27, 2018.

[8] R. Deahl, "The Benefits of Audiobooks for All Readers," Reading Rockets, 2017. [Online]. Available: https://www.readingrockets.org/topics/educationaltechnolo

gy/articles/benefits-audiobooks-all-readers. March 1, 2024]. [Accessed:]

[9] M. C. Sanchez and J. D. Rodriguez, "The Impact of Audiobooks on Reading Comprehension and Literacy Development in Diverse Learners," Journal of Adolescent & Adult Literacy, vol. 60, no. 4, pp. 421-431, 2016.

[10] M. Wessel and M. Jones, "Accessible Learning Materials for Postsecondary Education: A Review of the Literature," Journal of Postsecondary Education and Disability, vol. 28, no. 2, pp. 159-176, 2015.

[11] WHO | Visual impairment and blindness. WHO, April 7 1948 .http://www.who.int/mediacentre/factsheets/fs28 2/en/.