


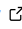

When Content Speaks Volumes: Podcastfy — An Open Source Python Package Bridging Multimodal Data and Conversational Audio with GenAI

Tharsis T. P. Souza ^{1,2}

¹ Columbia University in the City of New York ² Instituto Federal de Educacao, Ciencia e Tecnologia do Sul de Minas (IFSULDEMINAS)

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Abstract

Podcastfy is an open-source Python framework that programmatically transforms multisourced, multimodal content into multilingual, natural-sounding audio conversations using generative AI. By converting various types of digital content - including images, websites, YouTube videos, and PDFs - into conversational audio formats, Podcastfy enhances accessibility, engagement, and usability for a wide range of users. As an open-source project, Podcastfy benefits from continuous community-driven improvements, enhancing its adaptability to evolving user requirements and accessibility standards.

Statement of Need

The rapid expansion of digital content across various formats has intensified the need for tools capable of converting diverse information into accessible and digestible forms ([Chen & Wu, 2023](#); [Johnson & Smith, 2023](#); [McCune & Brown, 2023](#)). Existing solutions often fall short due to their proprietary nature, limited multimodal support, or inadequate accessibility features ([Gupta & Lee, 2023](#); [Marcus & Zhang, 2019](#); [Peterson & Allen, 2023](#)).

Podcastfy addresses this gap with an open-source solution that supports multimodal input processing and generates natural-sounding, summarized conversational content. Leveraging advances in large language models (LLMs) and text-to-speech (TTS) synthesis, Podcastfy aims to benefit a diverse group of users — including content creators, educators, researchers, and accessibility advocates — by providing a customizable solution that transforms digital content into multilingual textual and auditory formats, enhancing accessibility and engagement.

Features

- Generate conversational content from multiple sources and formats (images, websites, YouTube, and PDFs).
- Customize transcript and audio generation (e.g., style, language, structure, length).
- Create podcasts from pre-existing or edited transcripts.
- Leverage cloud-based and local LLMs for transcript generation (increased privacy and control).
- Integrate with advanced text-to-speech models (OpenAI, ElevenLabs, and Microsoft Edge).
- Provide multi-language support for global content creation and enhanced accessibility.
- Integrate seamlessly with CLI and Python packages for automated workflows.

See [audio samples](#).

Use Cases

Podcastfy is designed to serve a wide range of applications, including:

- **Content Creators** can use Podcastfy to convert blog posts, articles, or multimedia content into podcast-style audio, enabling them to reach broader audiences. By transforming content into an audio format, creators can cater to users who prefer listening over reading.
- **Educators** can transform lecture notes, presentations, and visual materials into audio conversations, making educational content more accessible to students with different learning preferences. This is particularly beneficial for students with visual impairments or those who have difficulty processing written information.
- **Researchers** can convert research papers, visual data, and technical content into conversational audio. This makes it easier for a wider audience, including those with disabilities, to consume and understand complex scientific information. Researchers can also create audio summaries of their work to enhance accessibility.
- **Accessibility Advocates** can use Podcastfy to promote digital accessibility by providing a tool that converts multimodal content into auditory formats. This helps individuals with visual impairments, dyslexia, or other disabilities that make it challenging to consume written or visual content.

Implementation and Architecture

Podcastfy implements a modular architecture designed for flexibility and extensibility through five main components, as shown in Figure 1.

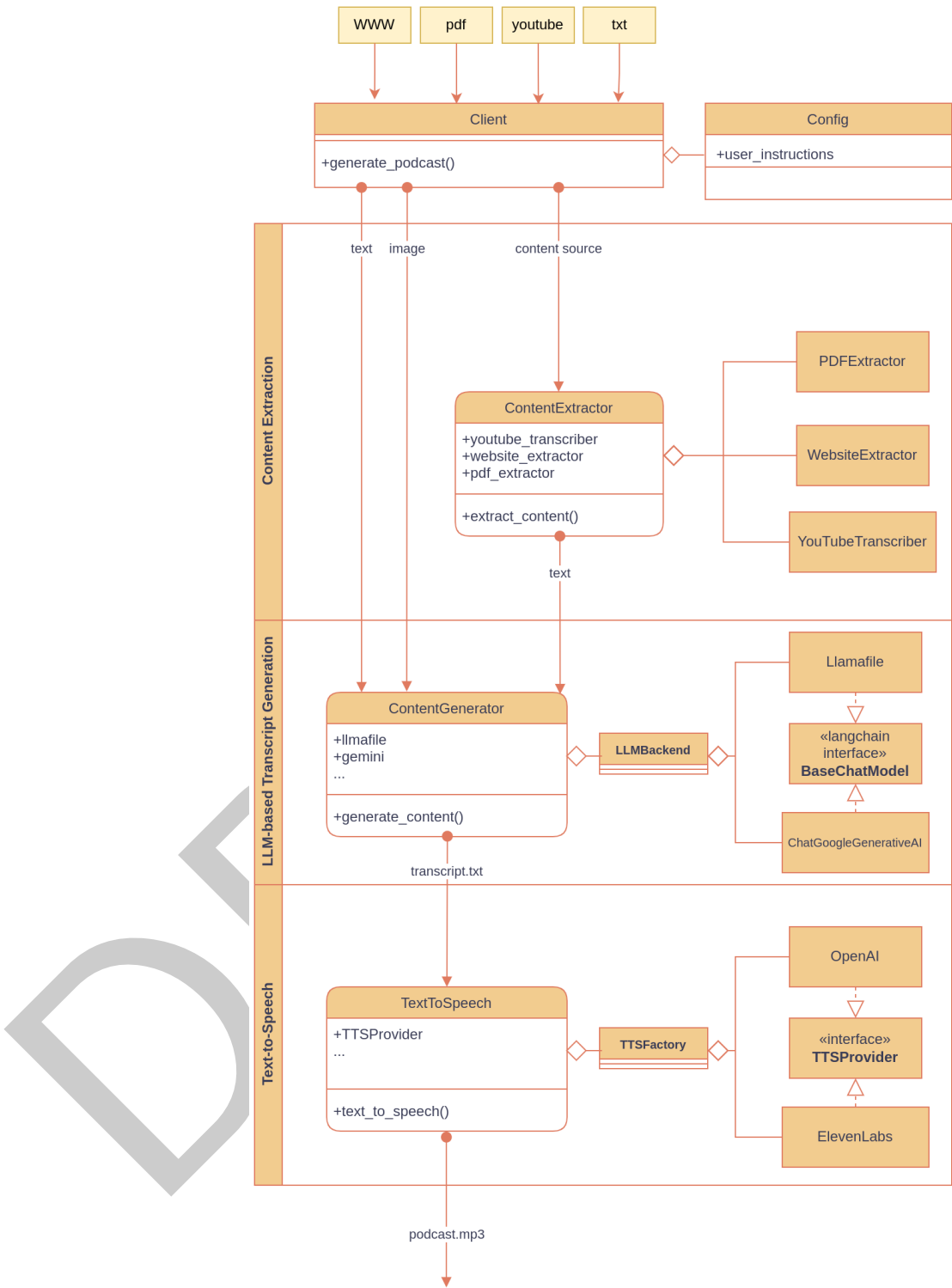


Figure 1: Podcastfy’s simplified architecture and workflow diagram showing the main components and their interactions.

- 60
- 61
- 62
1. Client Interface

 - Provides both CLI (Command-Line Interface) and API interfaces.
 - Coordinates the workflow between processing layers.

63 ▪ Implements a unified interface for podcast generation through the `generate_podcast()`
64 method.

65 **2. Configuration Management**

66 ▪ Offers extensive customization options through a dedicated module.

67 ▪ Manages system settings and user preferences, such as podcast name, language,
68 style, and structure.

69 ▪ Controls the behavior of all processing layers.

70 **3. Content Extraction Layer**

71 ▪ Extracts content from various sources, including websites, PDFs, and YouTube
72 videos.

73 ▪ The `ContentExtractor` class coordinates three specialized extractors:

74 – `PDFExtractor`: Handles PDF document processing.

75 – `WebsiteExtractor`: Manages website content extraction.

76 – `YouTubeTranscriber`: Processes YouTube video content.

77 ▪ Serves as the entry point for all input types, providing standardized text output to
78 the transcript generator.

79 **4. LLM-based Transcript Generation Layer**

80 ▪ Uses large language models to generate natural-sounding conversations from ex-
81 tracted content.

82 ▪ The `ContentGenerator` class manages conversation generation using different LLM
83 backends:

84 – Integrates with `LangChain` to implement prompt management and common
85 LLM access through the `BaseChatModel` interface.

86 – Supports both local (`LlamaFile`) and cloud-based models.

87 – Uses `ChatGoogleGenerativeAI` for cloud-based LLM services.

88 ▪ Allows customization of conversation style, roles, and dialogue structure.

89 ▪ Outputs structured conversations in text format.

90 **5. Text-to-Speech (TTS) Layer**

91 ▪ Converts input transcripts into audio using various TTS models.

92 ▪ The `TextToSpeech` class implements a factory pattern:

93 – The `TTSFactory` creates appropriate providers based on configuration.

94 – Supports multiple backends (`OpenAI`, `ElevenLabs`, and `Microsoft Edge`) through
95 the `TTSProvider` interface.

96 ▪ Produces the final podcast audio output.

97 The modular architecture enables independent development and maintenance of each compo-
98 nent. This pipeline design ensures a clean separation of concerns while maintaining seamless
99 data transformation between stages. This modular approach also facilitates easy updates and
100 extensions to individual components without affecting the rest of the system.

101 The framework is offered as a Python package, with a command-line interface as well as a
102 REST API, making it accessible to users with different technical backgrounds and requirements.

103 Quick Start

104 Prerequisites

- 105 ▪ Python 3.11 or higher
- 106 ▪ `$ pip install ffmpeg` (for audio processing)

107 Setup

- 108 1. Install from PyPI `$ pip install podcastfy`
- 109 2. Set up [API keys](#)

110 Python

```
110 from podcastfy.client import generate_podcast

audio_file = generate_podcast(urls=["<url1>", "<url2>"])
```

111 CLI

```
112 python -m podcastfy.client --url <url1> --url <url2>
```

113 Customization Examples

114 Podcastfy offers various customization options that make it versatile for different types of
115 content transformation. To accomplish that, we leverage LangChain's (LangChain, 2024)
116 prompt management capabilities to dynamically construct prompts for the LLM, adjusting
117 conversation characteristics such as style, roles, and dialogue structure. Below are some
118 examples that demonstrate its capabilities.

119 Academic Debate

120 The following Python code demonstrates how to configure Podcastfy for an academic debate:

```
120 from podcastfy import generate_podcast

debate_config = {
    "conversation_style": ["formal", "debate"],
    "roles_person1": "main presenter",
    "roles_person2": "opposing viewpoint",
    "dialogue_structure": ["Introduction", "Argument Presentation", "Counterarguments",
}

generate_podcast(
    urls=["PATH/T0/academic-article.pdf"],
    conversation_config=debate_config
)
```

121 In this example, the roles are set to “main presenter” and “opposing viewpoint” to simulate an
122 academic debate between two speakers on a chosen topic. This approach is especially useful
123 for educational content that aims to present multiple perspectives on a topic. The output is
124 structured with clear sections such as introduction, argument presentation, counterarguments,
125 and conclusion, allowing listeners to follow complex ideas easily.

126 Storytelling Adventure

127 The following Python code demonstrates how to generate a storytelling podcast:

```
127 from podcastfy import generate_podcast

story_config = {
    "conversation_style": ["adventurous", "narrative"],
    "creativity": 1.0,
    "roles_person1": "narrator",
    "roles_person2": "character",
    "dialogue_structure": ["Introduction", "Adventure Begins", "Challenges", "Resolution"]
}
```

```
generate_podcast(  
    urls=["SAMPLE/WWW.URL.COM"],  
    conversation_config=story_config  
)
```

128 In this example, Podcastfy creates an engaging story by assigning roles like “narrator” and
129 “character” and adjusting the creativity parameter for richer descriptions. Using this con-
130 figuration, Podcastfy can generate engaging narrative content. By adjusting the creativity
131 parameter, Podcastfy can create a story involving multiple characters, unexpected plot twists,
132 and rich descriptions.

133 Additional Examples

134 Daily News Briefing

```
news_config = {  
    "word_count": 1500,  
    "conversation_style": ["concise", "informative"],  
    "podcast_name": "Morning Briefing",  
    "dialogue_structure": [  
        "Headlines",  
        "Key Stories",  
        "Market Update",  
        "Weather"  
    ],  
    "roles_person1": "news anchor",  
    "roles_person2": "field reporter",  
    "creativity": 0.3  
}  
  
generate_podcast(  
    urls=[  
        "https://news-source.com/headlines",  
        "https://market-updates.com/today"  
    ],  
    conversation_config=news_config  
)
```

135 Language Learning Content

```
language_config = {  
    "output_language": "Spanish",  
    "word_count": 1000,  
    "conversation_style": ["educational", "casual"],  
    "engagement_techniques": [  
        "vocabulary explanations",  
        "cultural context",  
        "pronunciation tips"  
    ],  
    "roles_person1": "language teacher",  
    "roles_person2": "curious student",  
    "creativity": 0.6  
}  
  
generate_podcast(  
    urls=["https://spanish-content.com/article"],
```

```
        conversation_config=language_config
    )
```

136 Technical Tutorial

```
tutorial_config = {
    "word_count": 2500,
    "conversation_style": ["instructional", "step-by-step"],
    "roles_person1": "expert developer",
    "roles_person2": "learning developer",
    "dialogue_structure": [
        "Concept Introduction",
        "Technical Background",
        "Implementation Steps",
        "Common Pitfalls",
        "Best Practices"
    ],
    "engagement_techniques": [
        "code examples",
        "real-world applications",
        "troubleshooting tips"
    ],
    "creativity": 0.4
}

generate_podcast(
    urls=["https://tech-blog.com/tutorial"],
    conversation_config=tutorial_config
)
```

137 Working with Podcastfy Modules

138 Podcastfy's components are designed to work independently, allowing flexibility in updating
139 or extending each module. The data flows from the ContentExtractor module to
140 ContentGenerator and finally to the TexttoSpeech converter, ensuring a seamless transformation
141 of multimodal content into audio. In this section, we provide some examples of how
142 to use each module.

143 Content Extraction

144 Podcastfy's content_extractor.py module allows users to extract content from a given URL,
145 which can be processed further to generate a podcast. Below is an example of how to use the
146 content extraction component:

```
from podcastfy.content_extractor import ContentExtractor

# Initialize the content extractor
extractor = ContentExtractor()

# Extract content from a URL
url = "https://example.com/article"
extracted_content = extractor.extract_content(url)

print("Extracted Content:")
print(extracted_content)
```

147 This example demonstrates how to extract text from a given URL. The extracted content is
148 then passed to the next stages of processing.

149 Content Generation

150 The `content_generator.py` module is responsible for generating conversational content based
151 on textual input. Below is an example of how to use the content generation component:

```
from podcastfy.content_generator import ContentGenerator

# Initialize the content generator
generator = ContentGenerator(api_key="<GEMINI_API_KEY>")

# Generate conversational content
input_text = "This is a sample input text about artificial intelligence."
generated_conversation = generator.generate_conversation(input_text)

print("Generated Conversation:")
print(generated_conversation)
```

152 Users can opt to run a cloud-based LLM (Gemini) or run a local (potentially Open Source)
153 LLM model ([see local llm configuration](#)).

154 Text-to-Speech Conversion

155 The `text_to_speech.py` module allows the generated transcript to be converted into audio.
156 Below is an example of how to use the text-to-speech component:

```
from podcastfy.text_to_speech import TextToSpeech

# Initialize the text-to-speech converter
tts = TextToSpeech(model='elevenlabs', api_key="<ELEVENLABS_API_KEY>")

# Convert the generated conversation to speech
input_text = "<Person1>This is a sample conversation generated by Podcastfy.</Person1><Person2>This is a sample conversation generated by Podcastfy.</Person2>"
output_audio_file = "output_podcast.mp3"
tts.convert_to_speech(input_text, output_audio_file)

print(f"Audio saved to {output_audio_file}")
```

157 This example demonstrates how to use the `TextToSpeech` class to convert generated text into
158 an audio file. Users can specify different models for TTS, such as `elevenlabs`, `openai`, or
159 `edge` (free to use).

160 Limitations

161 Podcastfy has several limitations, including:

- 162 ▪ **Content Accuracy and Quality**
 - 163 – The accuracy of generated conversations depends heavily on the capabilities of the
 - 164 underlying LLMs.
 - 165 – Complex technical or domain-specific content may not always be accurately inter-
 - 166 preted or summarized.
 - 167 – The framework cannot guarantee the factual correctness of generated content,
 - 168 requiring human verification for critical applications.
- 169 ▪ **Language Support Constraints**

- While multilingual support is available, performance may vary significantly across different languages.
- Less common languages may have limited TTS voice options and lower-quality speech synthesis.
- Nuanced cultural contexts and idioms may not translate effectively across languages.
- **Technical Dependencies**
 - Reliance on third-party APIs (OpenAI, ElevenLabs, Google) introduces potential service availability risks.
 - Local LLM options, while providing independence, require significant computational resources.
 - Network connectivity is required for cloud-based services, limiting offline usage.
- **Content Extraction Challenges**
 - Complex webpage layouts or dynamic content may not be accurately extracted.
 - PDF extraction quality depends on document formatting and structure.
 - YouTube video processing depends on the availability of transcripts.
- **Accessibility Considerations**
 - Generated audio may not fully meet all accessibility standards.
 - Limited support for real-time content processing.
 - May require additional processing for users with specific accessibility needs.

These limitations highlight areas for future development and improvement of the framework. Users should carefully consider these constraints when implementing Podcastfy for their specific use cases and requirements.

Conclusion

Podcastfy contributes to multimodal content accessibility by enabling the programmatic transformation of digital content into conversational audio. The framework addresses accessibility needs through automated content summarization and natural-sounding speech synthesis. Its modular design and configurable options allow for flexible content processing and audio generation workflows that can be adapted for different use cases and requirements.

As an open-source project, Podcastfy benefits from continuous community-driven improvements and adaptations, helping support its long-term value and relevance in meeting evolving user requirements and accessibility standards.

We invite contributions from the community to further enhance the capabilities of Podcastfy. Whether it's by adding support for new input modalities, improving the quality of conversation generation, or optimizing the TTS synthesis, we welcome collaboration to make Podcastfy more powerful and versatile.

Acknowledgements

We acknowledge the open-source community and the developers of the various libraries and tools that make Podcastfy possible. Special thanks to the developers of LangChain, Llamafire and HuggingFace. We are particularly grateful to all our [contributors](#) who have helped improve this project.

References

- Chen, R., & Wu, Y. (2023). Digital accessibility tools for multimodal content processing: A systematic review. *Digital Transformation Review*, 5(3), 91–109.
- Gupta, S., & Lee, A. (2023). Advances in adaptive user interfaces for enhanced accessibility. *Journal of Accessibility and User Experience*, 14(3), 203–215.

- 215 Johnson, L., & Smith, K. (2023). Adaptive user interfaces for accessibility across digital content
216 modalities. *Journal of Multimodal Accessibility*, 17(2), 43–58. [https://link.springer.com/
217 article/10.1007/s12193-023-00427-4](https://link.springer.com/article/10.1007/s12193-023-00427-4)
- 218 LangChain. (2024). *LangChain: Building applications with LLMs through composability*.
219 <https://www.langchain.com/>
- 220 Marcus, A., & Zhang, T. (2019). Design for multimodal human-computer interaction. In
221 *Lecture notes in computer science* (Vol. 1157, pp. 160–176). [https://link.springer.com/
222 chapter/10.1007/978-3-319-52162-6_18](https://link.springer.com/chapter/10.1007/978-3-319-52162-6_18)
- 223 McCune, B., & Brown, L. (2023). Accessibility of digital information: Standards, frameworks,
224 and tools related to information literacy and information technology. *Journal of Information
225 Literacy and Accessibility*, 9(4), 112–129.
- 226 Peterson, L., & Allen, J. (2023). Web accessibility and multimodal digital engagement.
227 *Technology Accessibility Journal*, 12(1), 23–34.

DRAFT