# Technical Documentation: Image Description Application for the Visually Impaired

The Incredibles

January 21, 2024

# Contents

# 1  Introduction

In an era where visual content dominates our information landscape, ensuring accessibility for all is not just important, it's essential. This technical documentation introduces a pioneering Python application specifically designed to bridge the accessibility gap for visually impaired individuals. Our application harnesses the cutting-edge capabilities of OpenAI's GPT-4 to interpret and describe images, transforming them into detailed, descriptive text.

The core functionality of this application lies in its ability to accurately query GPT-4 with any given image and retrieve a comprehensive textual description. This process involves several intricate steps, from encoding images into a format compatible with the OpenAI API to crafting and sending precise requests that yield meaningful descriptions. The outcome is a powerful tool that can describe an array of images – from simple everyday objects to complex scenes – in a manner that is both understandable and useful for those who rely on non-visual cues to interact with their surroundings.

This documentation will guide you through the application's architecture, its components, and its operation, providing insights into how this technology opens new horizons for inclusivity and accessibility. By leveraging the remarkable capabilities of artificial intelligence, the application stands as a testament to the potential of technology to create a more inclusive world.

# 2  Overview of the Code

The application performs the following steps:

1. Receives the path of an image (any file can be addressed through this program using different models).

2. Encodes the image to base64 format, which is required by the OpenAI API.

3. Sends a request to OpenAI's GPT-4, including the image and a predefined description template.

4. Retrieves and provides the response from GPT-4, which is a descriptive text of the image.

# 3  Installation and Setup Requirements

- Python version: (3.12.0)

- Required libraries: openai[1], requests[2], base64[3], os[4].

- Installation commands or steps for setting up the environment.

# 4 Function Descriptions

## 4.1 encode_image

```
def encode_image(image_path):
    # Function implementation
```

**Purpose:** Encodes an image to a base64 string, preparing it for transmission to the OpenAI API.

**Parameters:**

- image_path: Path to the image file.

**Return Value:** Base64 encoded string of the image.

## 4.2 query_gpt4_with_image

```
def query_gpt4_with_image(image_path, api_key):
    #
```

**Purpose:** Queries GPT-4 with an encoded image and returns a descriptive text, aiding visually impaired users in understanding the image.

**Parameters:**

- image_path: Path to the image file.

- api_key: API key for accessing the GPT-4 service.

**Process Flow:**

1. Encodes the provided image using base64.

2. Prepares a request with the encoded image and a descriptive instruction template.

3. Sends the request to OpenAI's GPT-4 API.

4. Handles the API response and returns the descriptive text.

**Return Value:** Textual description of the image as interpreted by GPT-4.

# 5 Application Usage

## 5.1 Addressing the Challenge

Visual impairment significantly alters how individuals interact with and interpret visual media. Traditional methods like Braille displays, while revolutionary, have their limitations, especially in conveying complex visual information. Existing automated image-to-text descriptions in refreshable Braille displays often fall short in providing detailed and meaningful interpretations of images. This gap in technology presents a unique challenge in educational and informational contexts, particularly for visually impaired (VI) students relying on tactile textbooks and Braille descriptions.

## 5.2 The Need for Advanced Solutions

The current solutions, although invaluable, come with substantial limitations. Manual creation of tactile textbooks and detailed Braille descriptions is not only costly, averaging up to $20,000 per academic year, but also time-consuming. This lag in availability poses a significant barrier to timely access to educational materials for VI students. There's a pressing need for an automated solution that goes beyond basic object classification and short descriptions, providing rich, detailed, and contextually relevant image descriptions.

## 5.3 Our Application's Role

Our application steps into this space with an innovative approach, utilizing generative AI to produce comprehensive image descriptions. By integrating with various file formats, including PowerPoint presentations, PDFs, and image files, our solution aims to greatly enhance the accessibility of visual content for VI individuals. The application is designed to be compatible with advanced Braille displays, enabling users to receive detailed image descriptions in a format they can easily access and understand.

## 5.4 Current Capabilities and Future Prospects

In its current state, showcased in this hackathon, the application demonstrates its capability with images, particularly focusing on PowerPoint files. The long-term vision is to expand this functionality to a wide array of file formats, making the solution increasingly versatile and useful in various settings, including university education and beyond.

## 5.5 Implementation and Scale-up

While the initial implementation showcases the core functionality, scaling up to accommodate a broader range of file types is anticipated to be a linear process in terms of software engineering effort. Our team, consisting of five members, is committed to progressively enhancing the application, making it a fully-fledged solution that can seamlessly integrate with different mediums, thereby addressing the pressing needs of VI individuals in accessing and comprehending visual information comprehensively.

# 6 Error Handling

## 6.1 Current Troubleshooting

As of now, the application's troubleshooting process is relatively straightforward. In its current prototype stage, it primarily addresses basic issues such as incorrect image paths, invalid API keys, and common connectivity problems.

Users encountering these issues are advised to check their file paths, verify their OpenAI API credentials, and ensure a stable internet connection.

## 6.2   Future Troubleshooting and Error Handling

Looking ahead, our vision for the application encompasses a far more comprehensive and adaptive troubleshooting system. We aim to extend the application's capabilities to process a wide range of document types, including books, academic papers, PowerPoint presentations, and other varied formats. This expansion will undoubtedly introduce new layers of complexity in error handling and troubleshooting.

In anticipation of these challenges, we plan to implement a more dynamic troubleshooting mechanism. This will include:

- **Advanced Error Detection:** Automated systems to identify specific types of errors based on the file format and user input, enabling more targeted troubleshooting advice.

- **Response Time Optimization:** Addressing delays in processing and response times, especially for larger or more complex files, to ensure a seamless user experience.

- **Customized Problem-Solving Approaches:** Depending on the nature of the file and the issue encountered, the system will suggest various methods to rectify the problem, such as file format conversion, image resolution adjustments, or alternative querying strategies.

- **User Feedback Integration:** Actively incorporating user feedback to continuously improve error handling protocols and enhance the overall effectiveness of the application.

Through these advancements, we are committed to evolving our application into a robust, versatile tool that not only addresses a wide array of user needs but also maintains a high standard of reliability and user-friendliness. The ultimate goal is to ensure that regardless of the document type or the complexity of the request, users will

# 7   Conclusion

This documentation has outlined the technical aspects of a prototype application designed to revolutionize the way visually impaired individuals interact with visual content. While in its early stages, this application promises a transformative leap in accessibility, particularly in educational and informational contexts. By converting images into descriptive text through the innovative use of OpenAI's GPT-4, it bridges a crucial gap in the way visually impaired individuals access and understand visual information.

This prototype serves as a foundational step towards a more inclusive approach to information consumption. Its ability to describe images accurately and contextually holds immense potential, not only in daily life but also in academic and professional settings. It can significantly enhance the experience of reading books, browsing the internet, and engaging with various forms of visual media, making these activities more inclusive.

As we continue to develop and refine this technology, we foresee a future where the barriers faced by blind and visually impaired people in accessing visual information are significantly reduced, if not completely removed. This application is not just a tool; it is a step towards an inclusive world where everyone has equal access to information and knowledge, regardless of visual ability.

# References

1. OpenAI. *OpenAI API Reference (Python)* https://platform.openai.com/docs/api-reference?lang=python. Accessed: 2024-01-20.

2. Python Requests. *Requests: HTTP for Humans* https://requests.readthedocs.io/en/latest/. Accessed: 2024-01-20. 2023.

3. Python Software Foundation. *base64 — Base16, Base32, Base64, Base85 Data Encodings* https://docs.python.org/3/library/base64.html. Accessed: 2024-01-20. 2023.

4. Python Software Foundation. *os — Miscellaneous operating system interfaces* https://docs.python.org/3/library/os.html. Accessed: 2024-01-20. 2023.