



Degree Program in Management and Computer Science

Course of Algorithms

AI-Powered Voice  
Assistant  
Enhancing Hands-Free  
User Interaction  
Vox AI

Prof. Irene Finocchi

---

SUPERVISOR

277001

---

CANDIDATE

Academic Year 2023/2024

# Table of Contents

Introduction .....	1
Background.....	1
Objectives.....	1
Scope .....	1
Chapter 1: System Design and Architecture .....	2
1.1 Purpose and Scope .....	2
1.2 Objective.....	2
1.3 Capabilities.....	2
1.4 Technical Distinctions.....	2
Chapter 2: System Architecture .....	4
2.1 User Interaction Layer.....	4
2.2 Core Assistant Engine.....	5
2.3 Service Integrations .....	5
2.4 Data Management.....	5
2.5 Utility and Support Components.....	5
2.6 Scalability and Customization .....	5
2.7 Client-Side Details (Android Application) .....	6
2.8 Server-Side Configuration.....	6
Chapter 3: Implementation .....	7
Chapter 4: Data Management.....	9
4.1 Local Data Caching.....	9
4.2 Token Management.....	9
4.3 Data Directory Structure .....	9
4.4 Security and Integrity.....	10
Chapter 5: Speech Detection .....	11
5.1 System Initialization and Hotword Detection .....	11
5.2 Voice Activity Detection and Handling Silence .....	11
5.3 Buffer Management and VAD Checks.....	12
5.4 Audio Processing Workflow .....	13
5.5 Transcription of Speech.....	13
5.6 Error Handling and User Feedback .....	13
5.7 Fallback Mechanisms .....	14
Chapter 6: Voice Synthesis.....	15
6.1 Technology and Integration.....	15
6.2 Operational Flow .....	15
Chapter 7: Natural Language Understanding (NLU) .....	17
7.1 Technological Framework.....	17
7.2 Why GPT-3.5 Turbo 16k?.....	18

7.6 Comprehensive Evaluation .....	20
7.7 Functionality .....	20
Chapter 8: Agent Memory .....	21
8.1 Role and Utility of Conversation Token Buffer Memory .....	21
8.2 Technical Implementation .....	21
Chapter 9: Agent Tools .....	23
9.1 Implementation .....	23
9.2 Tool Initialization and Switching .....	23
Chapter 10: Web Search Engine Integration .....	26
10.1 Simple Search .....	26
10.2 Deep Search .....	26
10.3 Academic Search .....	27
10.4 Compliance and Efficiency .....	27
10.5 A Deeper Look at Deep Search .....	27
10.6 Technical Implementation .....	29
10.7 Relevance Ranking System .....	30
10.8 Web Scraping and Content Parsing .....	32
Chapter 11: PDF Handler .....	34
Chapter 12: Integration of Flask and Ngrok for the Android Application .....	38
12.1 Flask Application on the Laptop .....	38
12.2 Ngrok Configuration .....	40
12.3 Python Anywhere as a Middleware .....	41
Chapter 13: Implementation of the Python Anywhere Flask Server for the AI Voice Assistant .....	42
13.1 Server Setup and Configuration .....	43
13.2 Logging and Monitoring .....	43
13.3 Token Management and Secure Communication .....	44
13.4 Database Interactions and Session Management .....	44
Chapter 14: Interface Demonstration of the Android Application .....	45
14.1 Main Tab Overview .....	45
14.2 Search Demonstration .....	45
14.3 Note-saving Feature .....	46
Results .....	47
System Capabilities .....	47
Performance Metrics .....	47
Conclusion .....	49
References .....	51

# Introduction

## Background

Voice assistants have become an integral part of modern technology, enhancing user interaction with devices through hands-free, voice-command-based interfaces. These systems utilize artificial intelligence (AI) to perform various tasks, such as sending emails, managing schedules, playing music, and retrieving information from the web. Their significance lies in the convenience and efficiency they offer, particularly for users who require seamless and intuitive interactions with their digital environments.

## Objectives

This project aims to develop an AI-powered personal voice assistant that enhances user interaction with technology by enabling hands-free operation through voice commands. The primary objectives are to create a user-friendly system capable of performing a wide range of tasks without manual input, leverage AI and machine learning techniques to ensure dynamic and adaptable functionalities, and integrate with various services like Gmail, Google Drive, Google Calendar, Spotify, and Android devices to provide comprehensive support for daily activities. Furthermore, the system is designed to be scalable and operable on minimal hardware setups, from the Raspberry Pi to high-end systems.

## Scope

The scope of the project encompasses the design, development, and implementation of the AI-powered personal voice assistant. It focuses on several key areas including the system architecture which involves the detailed design of the system's components and their interactions. Implementation aspects cover the coding and integration of various functionalities such as speech recognition, natural language processing, and text-to-speech. User interaction is a significant focus, with the development of user interfaces for both desktop and mobile platforms to ensure ease of use and accessibility. Additionally, scalability and customization are crucial, ensuring the system can be customized and scaled to meet different user requirements and hardware capabilities.

# Chapter 1: System Design and Architecture

This chapter provides a comprehensive overview of the AI-powered personal voice assistant, detailing its purpose, objectives, capabilities, and technical distinctions. It aims to outline how the assistant enhances user interaction with technology, emphasizing simplicity, accessibility, and efficiency.

## 1.1 Purpose and Scope

This AI-powered personal voice assistant enhances user interaction with technology through a hands-free, voice-command-based interface, focusing on simplicity and user-friendliness. This project aims to make technology more accessible and efficient, particularly for individuals who are comfortable with advanced tech products.

## 1.2 Objective

The main objective of this voice assistant is to enable a variety of tasks to be executed without manual input, allowing users to manage their digital environments through voice commands effectively. By leveraging function calling features in OpenAI's<sup>[4]</sup> API, the assistant can dynamically perform actions as defined by developers, thus offering broad and adaptable functionalities.

## 1.3 Capabilities

This voice assistant incorporates several key functionalities to support daily activities. It is capable of sending emails, managing files, and backing up data to cloud services like Google Drive. Additionally, the assistant integrates with Google Calendar to assist in managing events and appointments. For entertainment and information access, it connects with Spotify to stream music and performs web searches to deliver news and other information as requested. The assistant also supports interactions with Android devices, allowing users to retrieve files from their computers remotely. Furthermore, the system can execute command-line instructions, incorporating human validation to ensure that commands are processed securely and appropriately.

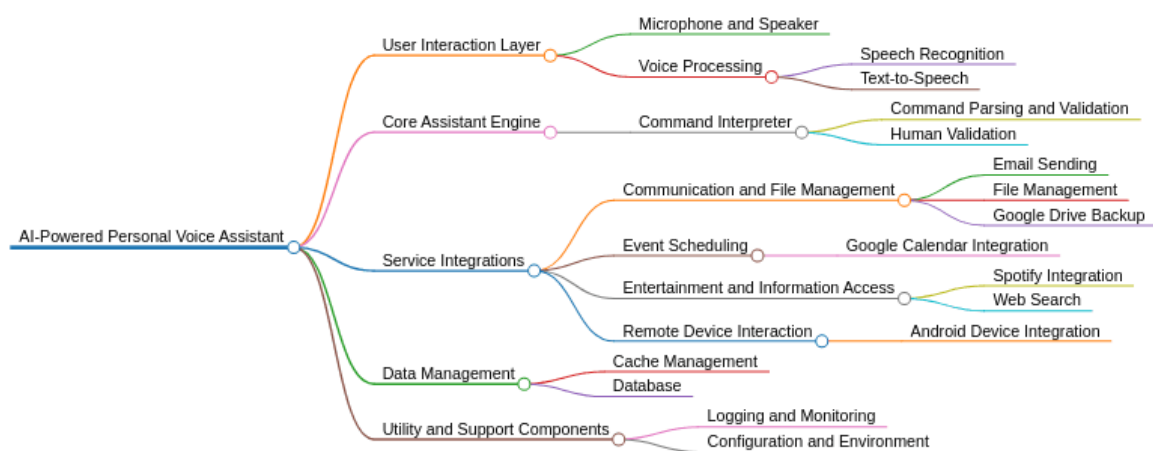
## 1.4 Technical Distinctions

The assistant is designed to operate effectively with minimal hardware, suitable for devices ranging from a Raspberry Pi 4 to mid-range laptops, and can be scaled up to make use of

dedicated GPUs for more resource-intensive tasks. Its architecture supports extensive customization and scalability, allowing it to adapt to various performance needs and user environments. The system's design facilitates significant user-driven customization, making it a versatile tool for a wide range of personal technology applications.

## Chapter 2: System Architecture

This chapter provides an in-depth explanation of the AI personal voice assistant's system architecture. It covers the interaction between the user and digital tasks through voice commands, detailing the components and their functionalities, with a focus on input processing, data management, and service integrations.



*Figure 1 Demonstration of the system architecture*

The architecture of the AI personal voice assistant, as shown in the above diagram, is designed to ensure seamless interaction between the user and digital tasks through voice commands. This section details the components and their functionalities, focusing on how the system processes inputs, manages data, and interacts with external services.

### 2.1 User Interaction Layer

The User Interaction Layer of this voice assistant comprises essential devices such as a microphone and speaker to capture voice commands and deliver audio responses, facilitating real-time auditory interaction. The voice processing aspect includes speech recognition, which converts spoken words into digital text using advanced algorithms to accurately transcribe user

commands into actionable inputs, and text-to-speech, which converts textual information from the system into spoken responses, allowing the assistant to audibly communicate back to the user.

## 2.2 Core Assistant Engine

The Core Assistant Engine features a Command Interpreter that handles the parsing and validation of user commands to understand intent and ensure validity before any action is taken. It also includes human validation to verify commands that could affect system settings or data security, ensuring all actions are intentional and secure.

## 2.3 Service Integrations

Service Integrations are extensive and include Communication and File Management capabilities, such as utilizing the Gmail API for sending emails and integrating with Google Drive for file management including creation, modification, and backup. Event Scheduling through Google Calendar manages scheduling and reminders. For Entertainment and Information Access, Spotify integration provides music streaming, and a web search feature executes internet searches for information and news. Additionally, Android Device Integration facilitates interaction with Android devices, enabling remote file access and control.

## 2.4 Data Management

Data Management includes Cache Management, which handles the temporary storage of data to enhance performance and response times, and a Database for managing long-term data storage dealing with user preferences, command logs, and historical interactions to refine the user experience and system responses.

## 2.5 Utility and Support Components

Utility and Support Components such as Logging and Monitoring are critical for maintaining system health, tracking operations, and logging errors, facilitating troubleshooting and system optimization. The Configuration and Environment component ensures consistent operation across various environments by managing settings and adaptations required for different hardware and software contexts.

## 2.6 Scalability and Customization

Finally, Scalability and Customization aspects include Hardware Support, making the assistant versatile across different hardware configurations, from minimal setups like a Raspberry Pi to more capable systems equipped with GPUs. Customization allows for



extensive user-driven adjustments, enabling the tailoring of functionalities and responses to individual preferences and specific use cases.

## 2.7 Client-Side Details (Android Application)

The Android application mirrors the full capabilities of the laptop version with some adaptations for mobile use. Notably, the app requires manual activation via a microphone icon, aligning with mobile usage patterns where users prefer explicit control over voice interactions. The app is structured into several key sections:

**Main Interaction Box:** Where voice interactions are initiated and responses are displayed.

**Notes Section:** Automatically captures and organizes notes from conversations, with options to sync these with Google Drive.

**News Feed:** Aggregates and summarizes current news, presenting them in an easily accessible format and providing auditory summaries upon request.

**Weather Tab:** Offers real-time weather updates and forecasts using the Open Weather API.

**User Input for File Searching:** Enhances the system's ability to accurately locate files on connected devices, particularly useful when voice recognition struggles with specific file names or formats.

## 2.8 Server-Side Configuration

The server setup, hosted on Python Anywhere, manages all backend processes including data caching, token storage, and communication with the Android app via Firebase. This configuration supports continuous data synchronization and reliable service availability, essential for maintaining the assistant's operational integrity.

## Chapter 3: Implementation

This section elaborates on the methods and mechanisms for configuring the environment and managing dependencies in the AI personal voice assistant system. It explains the role of various configuration files and the significance of the dependencies listed in the requirements.txt file.

### Environment Setup

The system configuration is managed through multiple files, ensuring flexibility and security:

**Environment Variables (.env):** This file secures sensitive information necessary for the application's operation, including API keys for services like OpenAI and Google, and authentication tokens needed for email, calendar integrations, and database connections. Storing these values in environment variables instead of in the codebase enhances security by keeping sensitive data out of version control.

**Configuration Script (config.py):** This Python script defines global settings for the application, including the secret key for secure sessions and the database URI. It utilizes the secrets module to generate a secure, random secret key for session management, and the os module to build paths that are relative to the project's base directory, ensuring the application remains environment agnostic.

**Path Management (config.yaml):** The YAML configuration file specifies critical paths within the project's structure, such as the location of server scripts and other essential directories. This approach decouples path settings from the main application logic, simplifying path management across different stages of development and deployment.

### Dependency Management

The requirements.txt file details all external Python libraries required for the system's functionality, aiding in the creation of a consistent setup across different environments:

**Web Frameworks:** Libraries like Flask provide the web server backbone, facilitating HTTP request handling and response delivery, while Flask-SQLAlchemy offers ORM capabilities for database interactions.

**AI and Machine Learning:** The project leverages AI libraries such as transformers and torch<sup>[7]</sup> for noise filtering and speech detection. Transformers<sup>[13]</sup> is used for leveraging state-of-the-art pre-trained models for text embeddings, while torch for voice related tasks like Voice Activity Detection (VAD).

**Audio Processing:** PyAudio<sup>[8]</sup> and SpeechRecognition<sup>[9]</sup> are critical for handling audio input and output, allowing the system to perform voice-to-text and text-to-voice conversions effectively.

**Data Handling:** Libraries like numpy and pandas are utilized for numerical operations and data manipulation, crucial for handling the data-intensive processes within the assistant.

**Cloud Integration:** Integration with Google's cloud services is managed through libraries such as google-cloud-storage and firebase-admin, which facilitate cloud data storage, real-time data synchronization, and user authentication.

**Communication:** spotipy enables Spotify integration, allowing the assistant to control music playback, and beautifulsoup4 along with requests support web scraping functionalities for fetching and processing online content.

This comprehensive list of dependencies ensures that all necessary functionality is supported and that the assistant can be easily set up in new environments by simply installing these packages using the command `pip install -r requirements.txt`.

## Chapter 4: Data Management

Effective data management is critical for the performance and security of the AI personal voice assistant. The system utilizes both local caching and cloud storage to ensure efficient operation and robust data security.

### 4.1 Local Data Caching

The voice assistant employs a local caching mechanism to enhance performance and reduce latency. Web content that the assistant accesses frequently is cached following compression. This process involves compressing the data using advanced algorithms before storage, which not only saves space but also speeds up subsequent access times. When the assistant encounters previously accessed web content, it retrieves and decompresses the cached data, significantly reducing the time required for data processing and presentation.

### 4.2 Token Management

Authentication tokens, crucial for secure communication with services like Spotify, Gmail, and the Android app, are stored both locally and on the Python Anywhere server. This dual storage approach ensures redundancy and high availability, enabling the assistant to maintain functionality even if one storage point fails. Each Android device connected to the assistant has a unique token, enhancing security and personalized service delivery.

### 4.3 Data Directory Structure

The organization of data within the system is both logical and functional, facilitating easy access and systematic data handling:

**Audio Files:** Stored under the audio directory, these files include pre-recorded messages that the assistant uses during its operation. Examples include notifications for actions like starting up (`start.mp3`), entering sleep mode (`sleep_mode.mp3`), and handling errors (`an_error_occured.mp3`). These audio cues enhance user interaction by providing audible feedback during various operations.

**Google Authentication Files:** Located in the `google_auth` directory, these files (`credentials.json`) manage the authentication process for Google services, ensuring secure API interactions and data transactions.

Logging: The logging directory contains log files (agent\_logs.log, agent\_logs.html) that record system activities and errors. This comprehensive logging mechanism supports system monitoring and troubleshooting by detailing every action, response time, and system error. It plays a critical role in maintaining the reliability and integrity of the system.

#### 4.4 Security and Integrity

The assistant uses robust logging and error-handling mechanisms to ensure the security and integrity of its operations. Errors and system activities are logged meticulously, with details such as token costs and response times monitored continuously. This data helps in analyzing system performance and identifying areas for improvement.

By integrating sophisticated data management practices, the system ensures efficient operation, rapid data access, and high levels of security, all of which are critical for a responsive and reliable personal voice assistant.

## Chapter 5: Speech Detection

This chapter delves into the speech detection mechanisms of the AI personal voice assistant. It covers the initialization process, hotword detection, voice activity management, buffer management, audio processing, speech transcription, error handling, and fallback mechanisms to ensure robust and reliable speech detection and recognition.

### 5.1 System Initialization and Hotword Detection

The speech detection process begins with system initialization, where the AI assistant prepares for operation by setting up necessary components like PyTorch models for voice activity detection and remote server configurations. At this stage, the system is poised to listen for the hotword (“Hey Vox”), which is detected using the `pvporcupine`<sup>[3]</sup> library. This library offers robust hotword detection capabilities, configured to minimize false positives by adjusting sensitivity settings.

```
self.porcupine = pvporcupine.create(  
    access_key=ACCESS_KEY,  
    keyword_paths=[HOTWORD_PATH],  
    sensitivities=[0.5] # Adjusted for minimal false positives  
)
```

Upon detecting the hotword, the system signals readiness through an auditory beep, transitioning into active listening mode, awaiting further user commands.

### 5.2 Voice Activity Detection and Handling Silence

To effectively manage periods of silence and prevent unnecessary processing, the system employs Silero VAD<sup>[2]</sup> (Voice Activity Detector), loaded through `torch.hub`. This tool is critical in determining whether incoming audio contains human speech, which optimizes processing and aids in managing system resources efficiently.

```
self.device = torch.device("cpu") # Can switch to 'cuda' for GPU  
acceleration  
self.torch_model, self.utils = torch.hub.load(  
    repo_or_dir="snakers4/silero-vad", model="silero_vad",  
    force_reload=True  
)  
(self.get_speech_ts, _, self.read_audio, _, _) = self.utils
```

The VAD checks are conducted on a sliding window of audio samples, where the system continuously evaluates whether the detected audio has speech. If prolonged silence is detected,

the system logs this event and may decide to break the listening loop, thus conserving computational resources and improving response efficiency.

```
silence_accumulated = 0.0
frame_duration = self.porcupine.frame_length /
self.porcupine.sample_rate
while True:
    data = audio_stream.read(self.porcupine.frame_length)
    buffer.extend(np.frombuffer(data, dtype=np.int16).tolist())
    vad_buffer = buffer[-WINDOW_SIZE:]
    vad_data = torch.tensor(vad_buffer).float().to(self.device)
    speech_timestamps = self.get_speech_ts(vad_data,
self.torch_model)
    if not speech_timestamps:
        silence_accumulated += frame_duration
        if silence_accumulated >= silence_duration:
            logger.info("Breaking due to silence")
            break
```

### 5.3 Buffer Management and VAD Checks

The management of audio buffers and the execution of voice activity detection checks are crucial for optimizing the speech recognition process. The system maintains a dynamic buffer where audio data is continuously collected and managed. This buffer is crucial for handling raw audio data streamed in real time from the microphone:

```
while True:
    data = audio_stream.read(self.porcupine.frame_length)
    buffer.extend(np.frombuffer(data, dtype=np.int16).tolist())
```

The buffer's content is periodically analyzed to detect active speech using the Silero VAD model. This model is applied to a window of the latest audio samples to determine if speech is present:

```
vad_buffer = buffer[-WINDOW_SIZE:] # Extract the latest audio
samples for analysis
vad_data = torch.tensor(vad_buffer).float().to(self.device) #
Prepare data for VAD
speech_timestamps = self.get_speech_ts(
    vad_data, self.torch_model
) # Determine speech presence
```

If prolonged silence is detected (no speech timestamps identified), the system can decide to reduce resource consumption by breaking out of the listening loop, thus preventing unnecessary processing:

```

if not speech_timestamps:
    silence_accumulated += frame_duration
    if silence_accumulated >= silence_duration:
        logger.info("Breaking due to silence")
        break

```

## 5.4 Audio Processing Workflow

Once speech is detected, the system processes the collected audio data to prepare it for transcription. The raw audio buffer is converted into a format suitable for speech recognition technologies. This conversion involves assembling the audio frames into a contiguous byte string and then creating an `AudioData` object that is compatible with the speech recognition API:

```

audio_data = b''.join(frames) # Combine all audio frames into a
single byte string
audio = sr.AudioData(
    audio_data,
    self.porcupine.sample_rate,
    pa.get_sample_size(pyaudio.paInt16)
)

```

This `AudioData` object is then used to perform actual speech recognition, either through your primary method (`faster-whisper`<sup>[5][6]</sup>) or through fallback mechanisms.

## 5.5 Transcription of Speech

Following successful speech detection and activity confirmation via Silero VAD, the system employs `faster-whisper` for speech transcription. This process involves converting the spoken input into textual data, which the system can then process for understanding and action.

```

audio_array, _ = sf.read(wav_stream)
audio_array = audio_array.astype(np.float32)
segments, _ = self.model.transcribe(audio=audio_array,
language="en", vad_filter=True)
speech = '\n'.join([segment.text.strip() for segment in segments if
segment.text])

```

In instances where `faster-whisper` fails to perform adequately due to audio quality or other issues, the system is equipped to fallback to Google's Speech Recognition API, ensuring that speech recognition is maintained at high reliability levels.

## 5.6 Error Handling and User Feedback

Comprehensive error handling mechanisms are embedded within the system to manage and respond to issues during the speech detection phase. When errors occur, the system informs



the user through specific audio feedback signals, and logs these incidents for further analysis and system enhancement.

```
except Exception as e:
    logger.error(f"Error in speech recognition: {str(e)}")
    self.speech_handler.play_audio('data/audio/failed.mp3')

self.speech_handler.play_audio('data/audio/an_error_occured.mp3')
```

## 5.7 Fallback Mechanisms

The system is designed to ensure high reliability in speech recognition by implementing fallback mechanisms when the primary method fails. If faster-whisper fails to accurately transcribe speech, possibly due to audio quality issues or other anomalies, the system switches to Google's Speech Recognition as a secondary option:

```
try:
    speech = self.recognizer.recognize_google(audio_data = audio,
language = "en-US")
except sr.UnknownValueError:
    logger.error("Google Speech Recognition could not understand the
audio")
    self.speech_handler.play_audio('data/audio/failed.mp3')
except sr.RequestError as e:
    logger.error(f"Could not request results from Google Speech
Recognition service; {e}")
    self.speech_handler.play_audio('data/audio/failed.mp3')
```

This dual-layered approach ensures that the system can maintain functionality even under less-than-ideal circumstances, providing robustness and enhancing user experience by reducing the chances of failed speech recognition attempts.

This detailed approach not only ensures robust operation of the speech detection functionality but also enhances user experience by maintaining clear communication regarding the system's status and any issues encountered.

## Chapter 6: Voice Synthesis

Voice synthesis in the AI personal voice assistant serves a critical role in providing audible feedback to the user, enhancing the interactive experience. This feature transforms text responses from the assistant into clear, audible speech, allowing for a seamless conversational interface. The synthesis process is particularly designed to manage and refine text output to ensure that the spoken feedback is user-friendly and appropriate, especially by filtering out elements like code snippets or URLs that are better suited for visual display.

### 6.1 Technology and Integration

The choice of Google Cloud's Text-to-Speech API for implementing voice synthesis was driven by its efficiency, cost-effectiveness, and the high quality of its speech output. Google's API provides a wide range of voice options and extensive language support, making it an ideal solution for applications requiring reliable and dynamic text-to-speech capabilities.

Integration of Google Cloud's Text-to-Speech API within the system is facilitated through HTTP requests made from the Python environment. The API is accessed via a secure key, ensuring all interactions are authenticated and data is transmitted securely.

This method allows the assistant to send text data to Google's servers, where it is processed and converted into audio in MP3 format, which is then played back to the user.

### 6.2 Operational Flow

The operational flow of voice synthesis begins with the preparation of text input, which is cleaned to ensure that only appropriate text is spoken. Special handling is applied to remove or substitute unsuitable content, such as code snippets or complex symbols, which are replaced with more general spoken cues like "See code in terminal."

```
def clean_response(self, text: str) -> str:
    """Replaces paths and code blocks with shorter messages for the
    VA to speak."""
    pattern = r"```.*?```"
    replacement = "See code in terminal."
    text = re.sub(pattern, replacement, text)
    return text
```

Once the text is sanitized and prepared, it is sent to Google's Text-to-Speech API with specific parameters detailing the desired voice characteristics and speaking rate. These parameters are

crucial for ensuring that the voice output matches the intended tone and clarity expected by the users.

```
data = {
    "input": {"text": text},
    "voice": {"languageCode": "en-US", "name": "en-US-Standard-I",
"ssmlGender": "MALE"},
    "audioConfig": {"audioEncoding": "MP3", "speakingRate": 1.15},
}
response = requests.post(url, headers=headers,
data=json.dumps(data))
```

The response from Google's API includes an audio file in MP3 format, which is temporarily stored and then played back using the pydub library. This process ensures that the user receives immediate and clear audio feedback.

```
audio_data = base64.b64decode(response.json()['audioContent'])
fd, mp3_path = tempfile.mkstemp(suffix=".mp3")
with os.fdopen(fd, 'wb') as tmp:
    tmp.write(audio_data)
    sound = AudioSegment.from_file(mp3_path, format="mp3")
    play(sound)
    os.remove(mp3_path)
```

Through this sophisticated integration of Google Cloud's Text-to-Speech API, the AI personal voice assistant effectively communicates with users, providing a dynamic and engaging user experience. The careful management of text input and audio output ensures that the assistant not only understands but also speaks in a way that is clear, contextually appropriate, and user-friendly.

## Chapter 7: Natural Language Understanding (NLU)

Natural Language Understanding (NLU) is central to the functionality of the AI personal voice assistant, enabling the system to interpret and respond to user commands accurately. NLU processes the transcribed text from the user's speech, identifies the intent, and determines the appropriate actions or responses. This capability allows the assistant to interact in a human-like manner, adapting to user requests with context-aware responses.

### 7.1 Technological Framework

The NLU component of the voice assistant leverages OpenAI's GPT-3.5 Turbo model, integrated through the LangChain<sup>[1]</sup> library. This setup utilizes advanced machine learning models to parse, understand, and generate responses based on the user's input. The choice of OpenAI's model is due to its state-of-the-art performance in understanding context and generating coherent, contextually appropriate text.

```
self.chat = ChatOpenAI(openai_api_key=API_KEY,
                       max_tokens=1000,
                       model_name='gpt-3.5-turbo-16k-0613',
                       temperature=0.1)

self.agent_kwargs = {
    "system_message": self.system_message,
    "extra_prompt_messages":
    [MessagesPlaceholder(variable_name="history")],
}

self.basic_tools = self.tool_manager.basic_tools

self.logfile = "data/logging/agent_logs.log"
logger.add(
    self.logfile, rotation="10 MB", retention="10 days",
    colorize=True, enqueue=True
)
self.handler = FileCallbackHandler(self.logfile)

self.agent = initialize_agent(
    tools=self.basic_tools,
    llm=self.chat,
    agent=AgentType.OPENAI_FUNCTIONS,
    max_iterations=3,
    max_execution_time=300,
    verbose=True,
    memory=self.memory,
```

```

agent_kwargs=self.agent_kwargs,
callbacks=[self.handler],
early_stopping_method="force",
handle_parsing_errors=self._handle_error)

```

Additionally, the system uses LangChain to create an AI agent, organizing the tools and managing the interaction flow efficiently, ensuring that the system does not exceed token limits or become overwhelmed by the user’s queries.

## 7.2 Why GPT-3.5 Turbo 16k?

Choosing the right language model is pivotal for optimizing the performance and cost-effectiveness of natural language processing (NLP) applications. This analysis focuses on the GPT-3.5 Turbo 16k model, emphasizing its extended context length and balanced attributes in quality, speed, and price, based on both empirical data and additional insights.

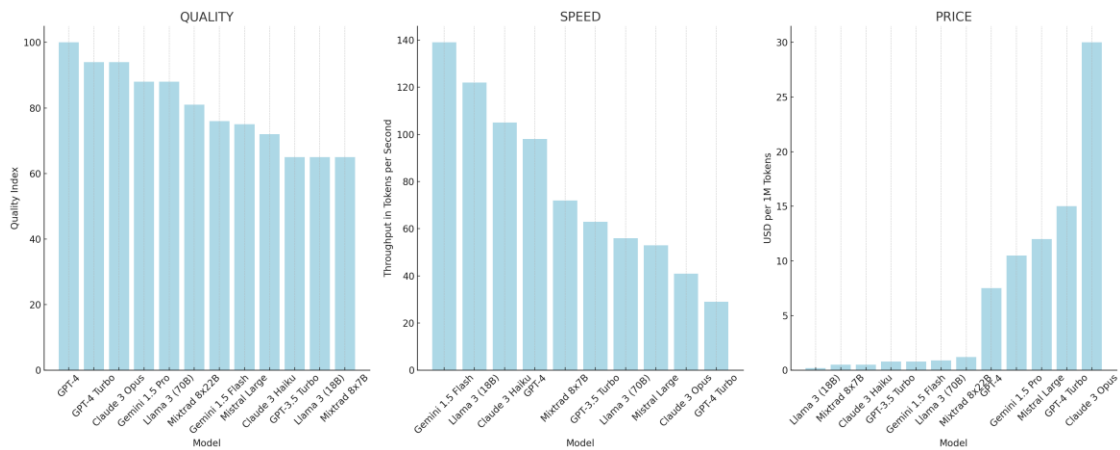


Figure 2 GPT-3.5 Turbo benchmarks against other LLMs [15]

### Quality

Quality is essential in a language model as it directly impacts the accuracy and reliability of generated text. The provided data shows that GPT-3.5 Turbo has a quality index of 65, which is competitive among high-quality models like GPT-4 (100) and Gemini 1.5 Pro (88). The GPT-3.5 Turbo 16k variant builds on this by supporting larger contexts, allowing it to process more extensive text segments. This capability is particularly beneficial for tasks requiring detailed and context-rich responses, such as lengthy documents or complex queries.

## Speed

Speed, measured in throughput (tokens per second), is another critical factor. GPT-3.5 Turbo demonstrates a throughput of 63 tokens per second, placing it in a mid-range position among evaluated models. Although it is outperformed by models such as Gemini 1.5 Flash (139 tokens per second) and Llama 3 (18B) (122 tokens per second), GPT-3.5 Turbo remains efficient for most interactive applications. The ability of the GPT-3.5 Turbo 16k to handle larger chunks of text in fewer API calls further optimizes its speed and efficiency, enhancing user experience and operational effectiveness.

## Price

Cost-effectiveness is crucial, especially for large-scale deployments. The GPT-3.5 Turbo 16k model, priced at \$0.003 per 1,000 input tokens and \$0.004 per 1,000 output tokens, offers a balanced cost structure. While it is more expensive than the base GPT-3.5 Turbo, it is significantly more affordable than the high-end GPT-4 models. This makes it an economically viable choice for applications requiring extensive context handling without incurring prohibitive costs.

## Additional Advantages of GPT-3.5 Turbo 16k

1. **Extended Context Length:** The GPT-3.5 Turbo 16k can handle up to 16,000 tokens per request, roughly equivalent to 20 pages of text. This extended context length is advantageous for tasks requiring comprehensive text processing, such as detailed reports and extensive data analysis.
2. **Fine-Tuning Capabilities:** The model supports fine-tuning, allowing businesses to tailor it to specific needs. Fine-tuning can improve performance by making the model better at following instructions, formatting responses, and adopting a custom tone. This leads to more accurate and contextually appropriate outputs, essential for specialized applications.
3. **Cost Savings Through Fine-Tuning:** Fine-tuning can significantly reduce prompt sizes, lowering the overall cost of API usage. Early tests have shown that fine-tuning can reduce prompt sizes by up to 90%, leading to faster and cheaper API calls.
4. **Versatility and Reliability:** The GPT-3.5 Turbo 16k model is versatile and reliable, capable of handling a wide range of NLP tasks efficiently. Its balanced attributes make it suitable for both real-time applications and extensive text processing tasks, providing a comprehensive solution for various needs.

## 7.6 Comprehensive Evaluation

When evaluating the trade-offs between quality, speed, and price, GPT-3.5 Turbo 16k stands out as a balanced and practical choice. Its extended context window enhances its ability to manage large texts, making it suitable for detailed and complex tasks. The model's competitive pricing and fine-tuning capabilities further strengthen its position as a cost-effective solution for a wide range of applications, from academic research to commercial deployment.

## 7.7 Functionality

The NLU process begins once speech has been detected and transcribed. The transcribed text undergoes a cleaning process where unnecessary elements such as code snippets or complex file paths are simplified or removed. This ensures that the text passed to the NLU system is clear and manageable.

Once cleaned, the text is passed to the GPT-3.5 Turbo model, where it is analyzed for intent. The system distinguishes between different types of requests—whether the user is initiating a conversation, asking a question, or commanding a specific action like retrieving a file or conducting a web search.

The NLU system is adept at parsing the user's intent and determining if any specific tools are required to fulfill the request. If a particular action is needed, the system utilizes OpenAI's function calling feature, which allows the LLM to interact dynamically with external tools or databases coded specifically for tasks like file management or web searches.

This process is managed efficiently to maintain low token usage and clear memory of irrelevant data, preventing the system from becoming overloaded with information and ensuring cost-effective operations.

```
self.memory = ConversationTokenBufferMemory(llm=self.chat,  
                                             max_token_limit=16000,  
                                             memory_key="history",  
                                             return_messages=True)
```

The system's NLU capabilities are crucial for providing a responsive and intuitive user experience, enabling the voice assistant to handle a variety of tasks and interact with users in a meaningful way.

## Chapter 8: Agent Memory

The *Conversation Token Buffer Memory* used in the AI personal voice assistant plays a critical role in managing and optimizing the natural language understanding process. This mechanism is essential for handling the conversation history and ensuring that the interaction between the user and the assistant is smooth, contextually aware, and efficient in resource usage.

### 8.1 Role and Utility of Conversation Token Buffer Memory

The Conversation Token Buffer Memory is designed to store and manage the history of interactions between the user and the AI system. This memory buffer allows the AI to maintain context over the course of a conversation, which is crucial for providing coherent and contextually appropriate responses. Here's how this system enhances the functionality:

**Context Management:** By storing the history of the conversation, the AI can refer back to previous exchanges to understand the context of the current interaction. This ability is vital for tasks that require continuity, such as follow-up questions or multi-step processes.

**Efficiency and Performance:** The memory buffer helps manage the token economy of the system, ensuring that the number of tokens<sup>[12]</sup> processed by the AI at any given time does not exceed the model's capacity. This management is crucial for maintaining the efficiency of the system, avoiding overloading the AI, and keeping computational costs in check.

**Dynamic Content Pruning:** The system includes mechanisms to dynamically prune older or less relevant content from the memory as new inputs are added. This pruning helps maintain the relevance and manageability of the information that the AI needs to keep in active memory.

```
if curr_buffer_length > self.max_token_limit:
    pruned_memory = []
    while curr_buffer_length > self.max_token_limit:
        pruned_memory.append(buffer.pop(0))
        curr_buffer_length =
self.llm.get_num_tokens_from_messages(buffer)
```

### 8.2 Technical Implementation

The implementation of the Conversation Token Buffer Memory involves several key components and processes:



Memory Initialization and Management: The memory system is instantiated with a predefined token limit, which represents the maximum number of tokens that can be stored. This limit helps ensure that the memory usage stays within the capabilities of the underlying language model.

```
class ConversationTokenBufferMemory (BaseChatMemory) :  
    max_token_limit = 2000 # Maximum tokens that can be stored in  
    memory (modifiable)
```

Message Handling and Storage: When new user inputs or AI outputs are generated, they are processed and stored in the memory. The system uses a method to convert these inputs and outputs into a format that is efficient for storage and retrieval, managing them as discrete messages.

Retrieval and Utilization: When generating responses, the AI can retrieve the stored conversation history to maintain context. This retrieval is managed through functions that compile the history into a single string or a set of structured messages, depending on the needs of the model at that moment.

```
def load_memory_variables (self, inputs: Dict[str, Any]) -> Dict[str,  
Any]:  
    return {self.memory_key: self.buffer} # Load the conversation  
    history for the AI model
```

Adaptive Pruning: The memory system is designed to adaptively prune its contents to stay within the token limit. This adaptive pruning ensures that the AI has access to the most relevant information without exceeding resource limits.

```
# Method to prune memory if it exceeds the token limit  
def prune_memory_if_needed (self):  
    while (  
  
self.llm.get_num_tokens_from_messages (self.chat_memory.messages)  
    > self.max_token_limit  
    ):  
        self.chat_memory.messages.pop (0)
```

The Conversation Token Buffer Memory is a crucial element of the AI's architecture, enabling sophisticated and context-aware interactions while managing computational resources effectively. This system not only enhances the user experience but also optimizes the operational efficiency of the AI, making it capable of handling extensive and complex dialogues with ease.

## Chapter 9: Agent Tools

Agent tools within the context of our system refer to a suite of integrated services and utilities that the AI agent can access and utilize to perform a variety of tasks. These tools encompass a wide range of functionalities, including but not limited to file management, email sending, music playing, web searching, academic research, and scheduling events. The primary purpose of these tools is to enhance the agent's capability to interact with external services and perform complex operations autonomously.

### 9.1 Implementation

The implementation of agent tools is facilitated by LangChain, which serves as the LLM (Large Language Model) architecture connecting the agent to these tools. The tools are organized into distinct groups based on their functionality. Each tool is initialized with specific credentials and configurations required to interact with the corresponding service. The integration of these tools into the agent's workflow allows for seamless task execution.

The tools are grouped as follows:

1. Google Drive Tools: Used for backing up and restoring files.
2. Basic Tools: Includes command terminal access, session management, and other fundamental utilities.
3. File Management Tools: Enables creating, modifying, deleting, and searching for files.
4. Search Tools: Facilitates various types of searches, including academic, simple, and deep searches.
5. Gmail Tools: Provides functionalities for sending emails.
6. Music Tools: Allows playing music through Spotify.
7. Question Answering Tools: Utilizes vector databases and retrieval-augmented generation to extract answers from PDFs.
8. Google Calendar Tools: Manages scheduling events and reminders.

### 9.2 Tool Initialization and Switching

Tools are accessed and managed through a central class, Tools, which initializes and provides access to the various tool groups. The tools can be switched based on the task

requirements using the `change_tool` method. This method allows dynamic switching of the toolset the agent utilizes, enabling it to adapt to different tasks efficiently.

```
class Tools:
    """Class for agent tools"""

    def __init__(
        self,
        change_tool_func,
        clear_memory,
        terminate_session,
        stop_session,
        audio_event,
    ):
        self.basic_tools = [...]
        self.file_management_tools = [...]
        self.gmail_tools = [...]
        # other tools initialization

class SpeechRecognition:
    """Voice recognition system"""

    def __init__(self):
        self.tool_manager = Tools(self.change_tool, ...)
        self.agent =
initialize_agent(tools=sel.tool_manager.basic_tools)

    def change_tool(self, tool_name_id: int) -> str:
        tools_dict = {
            1: self.tool_manager.gdrive_tools,
            2: self.tool_manager.basic_tools,
            3: self.tool_manager.file_management_tools,
            4: self.tool_manager.search_tools,
            5: self.tool_manager.gmail_tools,
            6: self.tool_manager.spotify_tools,
            7: self.tool_manager.pdf_tools,
            8: self.tool_manager.google_calendar_tools,
        }
        if tool_name_id in tools_dict:
            self.agent = initialize_agent(tools_dict[tool_name_id],
            ...)

            return f"Switched to {tool_names[tool_name_id]}"
        else:
            return "Invalid tool ID."
```

By organizing tools into groups and allowing dynamic switching, the agent can perform a wide range of tasks efficiently. This modular approach ensures that the system remains scalable and adaptable to new functionalities as needed.

## Chapter 10: Web Search Engine Integration

The AI personal voice assistant's web search capability is a sophisticated system designed to retrieve information efficiently and accurately from the internet. This system employs multiple search strategies to cater to the complexity and variety of user queries, enabling it to provide comprehensive and relevant information swiftly.

### 10.1 Simple Search

The simple search function is the assistant's first line of information retrieval for straightforward queries such as factual questions ("Who won the FIFA World Cup 2022?"). This method utilizes APIs from Google and Bing to fetch the top ten results from each, providing a total of twenty snippets. These snippets are directly analyzed by the assistant to extract and deliver the answer. This approach is designed for speed and efficiency, handling common queries by sifting through readily available data without delving into the deeper content of the web pages.

### 10.2 Deep Search

For more complex queries that require nuanced understanding or multiple perspectives ("What are the pros and cons of drinking milk in the morning?"), the assistant employs a deep search strategy. This method begins with the generation of multiple refined queries by an AI model specifically trained to expand and explore various facets of the original question. Typically, three distinct queries are crafted to cover different aspects related to the topic.

Once these queries are generated, they are sent to the Google Search API, which returns results for each. The assistant then reads and parses the full content of each returned article using sophisticated libraries designed to extract the main content from web pages, eliminating any irrelevant information or boilerplate text. This ensures that the AI processes only the essential content, thereby reducing the risk of misinformation and enhancing the quality of the response.

Each query's results are individually processed, and the most relevant article from each is chosen to be presented to the user. This meticulous process involves multithreading and parallel processing, enabling the assistant to handle multiple data streams simultaneously, significantly speeding up the search without compromising on detail.

### 10.3 Academic Search

The academic search function targets scholarly articles and research papers, particularly useful for queries that require authoritative and scientific evidence. Utilizing specialized APIs such as PubMed and arXiv, this search method retrieves research papers related to the query. Depending on the user's preference or the query's nature, the assistant can either download and analyze the full documents using its question-answering tools or summarize the findings based on abstracts and snippets provided by the APIs.

### 10.4 Compliance and Efficiency

The assistant is designed to respect internet protocols and ethical guidelines, adhering strictly to robots.txt rules to prevent unauthorized scraping of content. This respect for digital rights ensures that the assistant's web search functionality remains sustainable and legally compliant.

Additionally, the system's capability to perform parallel processing allows it to manage extensive data from multiple sources quickly, making the search process not only comprehensive but also remarkably swift. This efficiency is critical in maintaining user engagement and delivering a seamless information retrieval experience.

### 10.5 A Deeper Look at Deep Search

The *Deep Search* function of the voice assistant is one of its most powerful features, designed to handle complex queries by engaging multiple advanced techniques to ensure the relevance and accuracy of the information provided to the user.

#### Generating Diversified Queries

The process starts when a user poses a complex question. The AI does not simply relay this question directly to a search engine. Instead, it engages an advanced language model to reframe and diversify the query into multiple distinct versions. This is crucial because a single query may not capture all relevant facets or might bias the search results towards a particular interpretation.

Here's how the assistant handles query diversification:

**Initial Query Reception:** The original user query is taken as input. This is the question that needs exploration, such as "What are the pros and cons of drinking milk in the morning?"

Query Transformation: The AI assistant uses the OpenAIAssistant class, specifically designed to interact with OpenAI's language models. This interaction involves setting up the model with specific instructions to generate output that is directly usable as search queries.

```
class OpenAIAssistant:
    def __init__(self, model_name="gpt-3.5-turbo"):
        # Load environment variables
        load_dotenv()
        self.model_name = model_name
        self.output_parser = CommaSeparatedListOutputParser()

    def query(self, question, temperature=0.2, max_tokens=500):
        import openai

        response = openai.ChatCompletion.create(
            model=self.model_name,
            temperature=temperature,
            max_tokens=max_tokens,
            messages=[
                {
                    "role": "system",
                    "content": "You are a specialized AI assistant
...",
                },
                {
                    "role": "user",
                    "content": f'Provide three comprehensive and
distinct ...',
                },
            ],
        )
        model_answer = response["choices"][0]["message"]["content"]
        return self.output_parser.parse(model_answer), model_answer
```

Model Configuration: The model is prompted to generate outputs that are straightforward lists of comma-separated queries, ensuring that the output is in a format that can be immediately processed. The system's configuration emphasizes generating responses devoid of any extraneous text or formatting, focusing solely on delivering three distinct and comprehensive queries related to the user's question.

Handling Model Responses: The language model's response is parsed to extract a list of queries. These queries are designed to be diverse, covering different angles or aspects of the original question. This ensures a broad and multi-dimensional search process, enhancing the likelihood of retrieving comprehensive information.

If the parsing fails or no diversified queries are generated, the system reverts to using the original query to prevent the search process from stalling.

```
if not diversified_queries:
    logger.warning(
        "Parsing failed or generated no queries. Using original
        query instead."
    )
    diversified_queries = [query]
```

**Logging and Monitoring:** Throughout this process, the system logs information about the generated queries and any issues encountered. This not only helps in debugging but also in refining the query generation process over time.

```
logger.info(f"Model's Generated Queries: {model_answer}")
logger.info(f"Parsed Generated Queries: {diversified_queries}")
logger.info(f"Using cached content for link: {link}")
logger.error(f"Unable to fetch robots.txt from {robots_url} due to
{e}")
```

## Integration with Search APIs

Once the diversified queries are ready, they are passed to search APIs like Google Custom Search. Each query is independently used to fetch results, which are then evaluated for relevance before further processing. This approach not only maximizes the breadth of the search but also enhances the depth, as different queries can lead to different types of information being uncovered.

## 10.6 Technical Implementation

The AI assistant's ability to generate and utilize diversified queries relies heavily on the interaction with OpenAI's language models. This interaction is carefully crafted to ensure that the model understands its task precisely:

**System Instructions:** Before sending the user query, the system sets up the model with instructions that specify the expected format and content of the response. This helps in minimizing the need for post-processing and ensures that the outputs are in a ready-to-use format for web searches.

**Response Processing:** The system uses a specialized output parser (CommaSeparatedListOutputParser) to handle the model's response. This parser is adept at converting the model's text output into a structured list of queries by dealing with various formatting challenges, such as handling quoted commas and newline-separated items.



```
class CommaSeparatedListOutputParser(ListOutputParser):
    def parse(self, text: str) -> List[str]:
        reader = csv.reader([text], skipinitialspace=True)
        return [item.strip('"') for item in next(reader)]
```

This meticulous process of generating diversified queries significantly enhances the depth and relevance of the search results, thereby providing users with rich and varied information. This capability is central to the AI assistant's effectiveness in handling complex inquiries where multiple perspectives or detailed explanations are required. If further details are needed or specific parts of the process need more explanation, please let me know, and we can continue to refine this description.

## 10.7 Relevance Ranking System

After the AI generates diversified queries and retrieves multiple search results, the next crucial step is to determine which of these results are most relevant. This is where the Relevance Ranking System comes into play. It evaluates each link's relevance to the query before any extensive content processing occurs, thus optimizing both the efficiency and effectiveness of the search process.

### Operation of the Relevance Ranking System

The operation of the Relevance Ranking System is broken down into several steps:

**Collection of Search Results:** As the system gathers links from search engines using the diversified queries, it prepares to assess each link's potential relevance. This is crucial because not all links returned by a search query will be equally useful or relevant to the user's needs.

**Relevance Evaluation:** Each link, along with its title, is passed to a specialized AI model trained to assess relevance based on the content's context and the specific nuances of the query. This model operates under strict guidelines to focus solely on the relevance of the titles to the query, excluding any extraneous considerations.

**AI Model Configuration:** The `MostRelevantQuerySelector` class encapsulates the functionality of this AI model. It is configured to interact with OpenAI's language models, specifically tailored to evaluate lists of URLs and titles for relevance.

**Model Interaction:** The model is queried with a structured prompt that asks it to choose the most relevant article based solely on the titles provided. This prompt ensures that the model's response is focused and directly applicable to the task at hand.

```
response = openai.ChatCompletion.create(
    model=self.model_name,
    messages=[
        {
            "role": "system",
            "content": "Evaluate the relevance of these links based
on their ...",
        },
        {
            "role": "user",
            "content": f'Pick the most relevant link related to
"{question}" ...',
        },
    ],
)
```

**Selection of the Most Relevant Link:** The model processes the information and returns the URL of the article it assesses as most relevant. If no title matches the relevance criteria closely enough, the model is capable of returning None, indicating that none of the provided links are sufficiently relevant.

**Processing the Selected Link:** Once the most relevant link is identified, the system proceeds to fetch and process the content from this link only. This targeted approach not only saves resources but also ensures that the information presented to the user is the most pertinent and valuable.

## Benefits

This relevance ranking system provides several benefits:

**Efficiency:** By filtering out less relevant links early in the process, the system conserves resources and reduces the time needed to deliver answers to the user.

**Accuracy:** Focusing on the most relevant sources increases the likelihood that the information provided is accurate and directly addresses the user's query.

**User Satisfaction:** By consistently providing highly relevant information, the system enhances user trust and satisfaction.

## 10.8 Web Scraping and Content Parsing

Following the relevance assessment, the selected links undergo a scraping process where the full content of the pages is retrieved. The system employs scraping techniques that respect the site's robots.txt rules, ensuring ethical data gathering practices.

During parsing, unnecessary elements such as ads, navigation bars, and any non-informative text are filtered out to focus solely on the main content. This step is supported by libraries designed to identify and extract the core textual content of a webpage effectively.

### Intelligent Caching and Blacklisting

To enhance efficiency, the system incorporates an intelligent caching mechanism. Web pages that have been previously scraped and processed are stored in a local database, allowing the assistant to quickly retrieve cached content without re-scraping the same pages.

```
cached_content = self._load_from_cache(link)
if cached_content:
    return self._process_cached_content(cached_content)
```

Moreover, a blacklist system is maintained where links that consistently fail to provide valuable content or that breach scraping policies are noted and avoided in future searches, optimizing the search process and resource allocation.

### Handling Special Content Types

The deep search is also equipped to handle special content types like PDFs and Reddit<sup>[14]</sup> threads. For PDFs, the system can extract text directly from the document, analyzing it for relevance to the query. For Reddit, it aggregates top comments and discussions relevant to the query, providing a rich, user-generated context to the answers.

```
if link.endswith(".pdf"):
    return self._handle_pdf_link(link, title, search_query)
elif "reddit.com" in link:
    return self._handle_reddit_link(link, title)
```

### Multithreading and Parallel Processing

To manage the complexity and volume of data processed during deep search, the system leverages multithreading and parallel processing techniques. This approach allows it to handle multiple links and queries simultaneously, drastically reducing the time required to gather and process information.

```

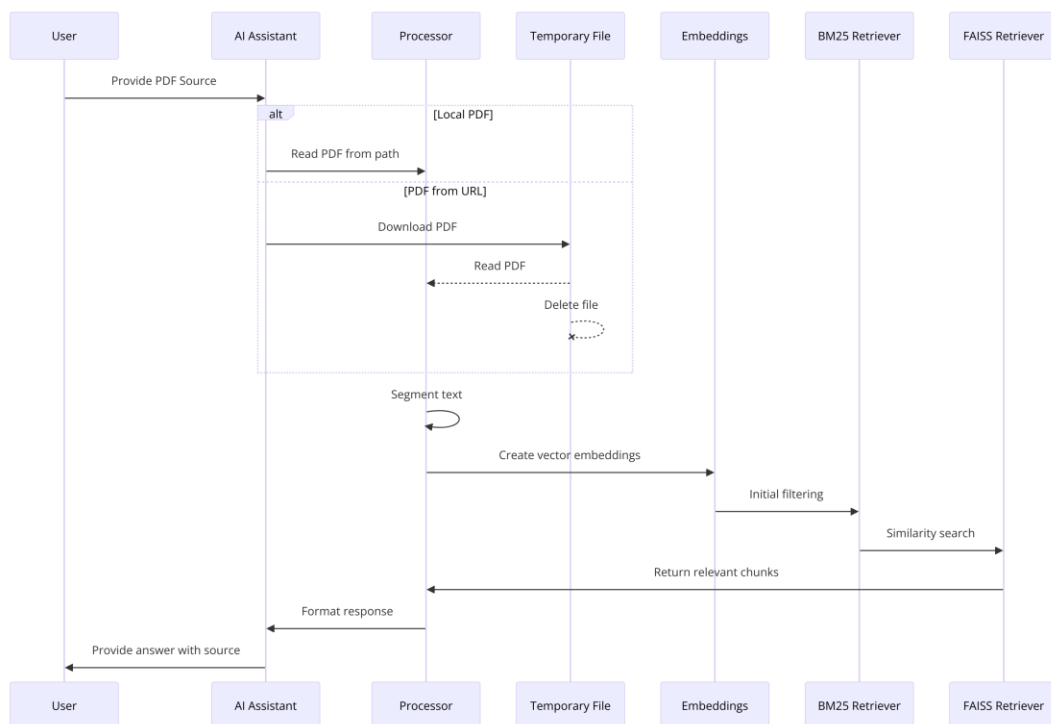
def _process_search_results(self, results, query):
    """Handles multithreading for page content extraction."""
    summaries = []
    with ThreadPoolExecutor(max_workers=5) as executor:
        future_to_item = {
            executor.submit(self._fetch_and_process_link, item,
query): item
            for item in results
        }
        for future in
concurrent.futures.as_completed(future_to_item):
            item = future_to_item[future]
            try:
                result = future.result()
                if result:
                    summaries.append(result)
            except Exception as exc:
                logger.error(f"Link {item['link']} generated an
exception: {exc}")
    return summaries

```

This sophisticated orchestration of tasks ensures that the deep search function is not only thorough but also exceptionally fast, providing users with detailed and accurate responses in a fraction of the time typically required for such extensive searches.

## Chapter 11: PDF Handler

The *PDF Query Processor* within the AI personal voice assistant is a component designed to handle complex question-answering tasks specifically related to content stored within PDF documents. This functionality not only enhances the assistant's versatility but also its ability to provide detailed and accurate answers directly from a wide array of textual resources, including academic papers, reports, and other document formats stored as PDFs.



*Figure 3 Demonstration of how the PDF Handler works*

The PDF Query Processor integrates several advanced techniques for text extraction, text segmentation, vector embeddings, and similarity searches to retrieve relevant information from PDF documents. This comprehensive approach ensures that the assistant can effectively parse and interpret the contents of PDFs to answer user queries accurately.

### Loading and Preparing PDF Text

The first step in processing a PDF document involves loading the text. This can happen in two scenarios: the PDF is available as a local file or needs to be fetched from a URL. Here's how the processor handles each case:

**Local PDFs:** The processor directly opens and reads the PDF from the specified path.

**PDFs from URLs:** The PDF is downloaded to a temporary file, which is then read into the processor. This temporary file is deleted after the text extraction to clean up resources.

```

response = requests.get(pdf_file_path_or_url, stream=True,
timeout=5)
# Writing to a temporary file
with NamedTemporaryFile(suffix=".pdf", delete=False) as temp_file:
    for chunk in response.iter_content(chunk_size=8192):
        temp_file.write(chunk)
    pdf_file_path = temp_file.name

```

The text extraction uses libraries such as PdfReader to pull text from each page of the PDF, which is then concatenated into a single string representing the document's textual content.

## Text Segmentation

Given the potentially large amount of text within a PDF, the processor segments the text into manageable chunks. This is done to facilitate more effective embedding and similarity searches later in the process. The segmentation respects a predefined chunk size and may overlap segments to ensure continuity of context across chunks.

```

text_splitter = RecursiveCharacterTextSplitter(
    chunk_size=self.chunk_size, chunk_overlap=self.chunk_overlap)
return text_splitter.split_text(text=text)

```

## Vector Embedding and Similarity Search Setup

Once the text is segmented, the next step involves creating vector embeddings for each chunk. These embeddings are crucial for performing efficient similarity searches to find the most relevant text segments in response to a user query.

**Vector Store:** If embeddings for a particular PDF have already been stored, they are loaded directly; otherwise, new embeddings are generated using models like sentence-transformers/all-MiniLM-L6-v2.

```

if os.path.exists(file_path):
    with open(file_path, "rb") as f:
        return pickle.load(f)
else:
    embeddings = HuggingFaceEmbeddings(
        model_name="sentence-transformers/all-MiniLM-L6-v2"
    )
    vector_store = FAISS.from_texts(chunks,
embedding=embeddings).as_retriever()

```

The use of the *all-MiniLM-L6-v2 model* from Sentence Transformers for generating embeddings is strategic, focusing on the balance between performance and efficiency:

**Quality of Embeddings:** This model generates high-quality embeddings that capture the semantic richness of the text, which is vital for the subsequent similarity search tasks. The fine-grained understanding enabled by these embeddings ensures that the responses generated by the AI are not only relevant but also contextually accurate.

**Operational Efficiency:** The choice of all-MiniLM-L6-v2 reflects a need for operational efficiency. Its design optimizes for quicker inference times without a substantial sacrifice in the quality of embeddings, making it suitable for real-time applications where response time is critical.

**Retrieval Models:** The system employs both BM25<sup>[11]</sup> and Faiss<sup>[10]</sup> retrievers, possibly in an ensemble, to leverage the strengths of both text-based and embedding-based retrieval mechanisms. This hybrid approach increases the accuracy and relevance of the search results.

The *BM25 Retriever* is employed within the system primarily for its effectiveness in initial filtering of text segments from PDFs based on their relevance to the query. This is particularly useful in environments where text data is extensive and varied:

**Relevance Scoring:** BM25 provides a robust mechanism for scoring segments of text based on their relevance to user queries. This scoring helps in quickly narrowing down the potential candidates from large volumes of text that need more detailed analysis, thus optimizing the processing pipeline.

**Efficient Preliminary Filtering:** By applying BM25 before engaging more computationally intensive processes like vector-based similarity checks, the system efficiently focuses resources on the most promising text excerpts, improving response times and reducing computational load.

The *Faiss Retriever* is integrated for its ability to perform rapid similarity searches among high-dimensional vectors representing text segments:

**Enhanced Accuracy:** After BM25 selects potentially relevant text chunks, Faiss assists in the fine-grained comparison of these chunks through vector similarity. This step significantly enhances the accuracy of the final document selection by ensuring that the semantic nuances of the text are considered.

Scalability for High-Dimensional Data: Given the high-dimensional nature of text embeddings, Faiss's efficient handling of large vector spaces is crucial. It allows the system to scale up to handle large datasets without a linear increase in search time or resource consumption.

## Query Processing and Source Formatting

When a user query is received, it is processed against the prepared vector store to retrieve the most relevant document chunks. The results are then formatted to provide a clean and readable response to the user, highlighting where in the PDF the answers were sourced.

```
def get_sources(self, query, pdf_file_path_or_url):
    """Returns sources based on the query to answer questions about
    pdf files.

    Args:
        query (str): the question to get sources for from the pdf
        file
        pdf_file_path_or_url (str): absolute path or url of the pdf
        file

    Returns:
        str: sources
    """
    logger.info(
        "Initializing PDFQueryProcessor with chunk_size: {} and
        chunk_overlap: {}".format(
            self.chunk_size,
            self.chunk_overlap,
        )
    )
    logger.info("Getting sources for query: {}".format(query))
    docs = self.process_query(query, pdf_file_path_or_url)
    response_text = self.format_sources(docs)
    logger.success("Done.")
    return response_text

def format_sources(self, docs):
    logger.info("Formatting sources...")
    return f"\n{'-' * 100}\n".join(
        [f"Source {i+1}:\n\n" + d.page_content for i, d in
         enumerate(docs)]
    )
```

The formatting includes clearly marking each source, making it easy for users to understand the context and origins of the information provided.



## Chapter 12: Integration of Flask and Ngrok for the Android Application

This chapter describes the integration of Flask and Ngrok to facilitate communication between the desktop environment and the Android application of the AI Voice Assistant. It details the setup of the Flask application, the configuration of Ngrok, and the role of Python Anywhere as a middleware, highlighting their roles within the system's architecture.

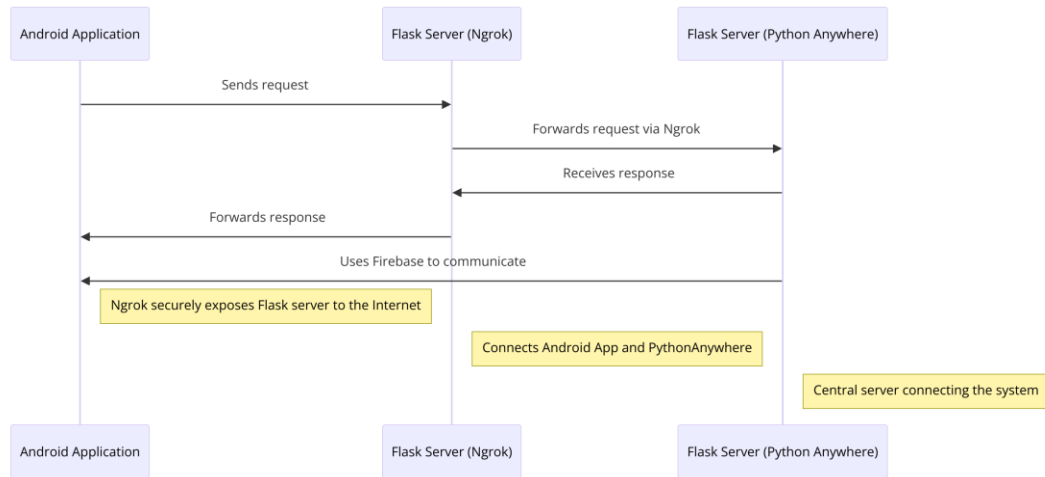


Figure 4 Demonstration of how the system setup work across the web

I employed a multifaceted server setup to ensure robust communication between the desktop environment and the Android application. This setup uses Flask as the primary web framework, enhanced by Ngrok for secure exposure to the internet, and Python Anywhere to bridge the interaction between the desktop and mobile platforms. Below, I detail the rationale and implementation of these components, emphasizing their roles within the system's architecture.

### 12.1 Flask Application on the Laptop

The Flask application acts as the central nerve of the AI Voice Assistant on the laptop, handling requests and processing them through the AI models. This server is designed to be lightweight and responsive, capable of managing multiple requests simultaneously due to its threaded configuration. The core functionalities of this Flask application include:

Endpoint for Questions: The Flask server defines an endpoint `/ask` that receives questions from the Android application. It processes these inquiries through the AI agent, which uses language

models to generate responses. The robust error handling ensures that any issues during request processing are caught and logged, maintaining system integrity.

```
@app.route("/ask", methods=["POST"])
def ask():
    try:
        # Get the question from the request data
        data = request.get_json()
        question = data.get("question")

        if not question:
            return jsonify({"response": "No question provided"}),
400

        response = agent_interface.get_response_from_agent(question)

        # Make sure each line in the response ends with a full stop
        response = ensure_full_stop_at_line_end(response)

        if not response:
            return jsonify({"response": "Failed to get a response
from the agent"}), 500

        # Return the response as 'response' to match Android app's
        expected response key
        return jsonify({"response": response})

    except Exception as e:
        logger.error(f"An error occurred: {str(e)}")
        return jsonify({"error": "An error occurred processing your
request"}), 500
```

**Real-time Logging:** A vital feature of this setup is the real-time log tracking. Logs are streamed to the web interface, allowing remote monitoring of the assistant's operations. This is crucial when I am not physically near the laptop, as it enables me to view the assistant's status and troubleshoot any potential issues from my Android device. Logs include detailed operational messages and are maintained in a rolling fashion to prevent overflow and ensure that only recent activities are displayed.

**Interactive Web Interface:** The server hosts a simple yet functional web interface, accessible through the root endpoint /. This interface utilizes server-sent events (SSE) to display real-time logs, which are crucial for remote diagnostics and monitoring the assistant's interactions and performance.

```

@app.route("/logs")
def sse_logs():
    """
    Stream Server-Sent Events (SSE) logs to the client.

    This endpoint streams log messages to the client using Server-
    Sent Events (SSE).
    It first sends recent logs, followed by continuously sending new
    logs as they arrive.

    Returns:
        Response: A Flask Response object with the mimetype set to
        'text/event-stream',
        which streams log messages to the client.
    """

    def generate():
        """
        Generator function that yields log messages in SSE format.

        This function first yields the recent logs and then
        continuously
        yields new logs from the log queue as they become available.

        Yields:
            str: Log messages formatted as SSE data.
        """

        # First, send the recent logs
        for log in recent_logs:
            yield f"data: {log}\n\n"

        # Now get new logs from the queue
        while True:
            log = log_queue.get()
            yield f"data: {log}\n\n"

    return Response(stream_with_context(generate()),
                    mimetype="text/event-stream")

```

## 12.2 Ngrok Configuration

Ngrok plays a critical role in this architecture by securely exposing the local Flask server to the internet. It creates a tunnel to the Flask application, allowing it to send and receive data over the internet without exposing the entire device to potential security risks. This is

particularly important for maintaining the confidentiality and integrity of the data processed by the AI Voice Assistant. Here's how Ngrok integrates:

**Secure Tunneling:** Ngrok starts before the Flask server, setting up a secure tunnel specified by a custom domain name. This tunnel is configured to forward requests to the Flask server running locally on a predefined port. The use of Ngrok ensures that the communication between the Flask server on the laptop and the external server on Python Anywhere is encrypted and secure.

**Domain Configuration:** The domain provided by Ngrok (e.g., regularly-alive-roughly.ngrok-free.app) is used to configure the external Flask server on Python Anywhere, allowing it to route requests correctly between the Android application and the laptop.

### 12.3 Python Anywhere as a Middleware

Python Anywhere serves as the middleware connecting the Ngrok-enabled Flask server on the laptop and the Android application. It hosts another Flask application that acts as a relay, managing communication and data synchronization between the mobile interface and the laptop. This setup ensures that requests from the Android app are correctly forwarded to the laptop and that responses travel back to the user's device.

**Firestore Integration:** The server uses Firestore for real-time data exchange with the Android application. This integration allows for efficient and reliable communication, ensuring that messages are delivered promptly, enhancing the responsiveness of the AI Voice Assistant.

**Central Management:** Python Anywhere centralizes the management of requests and responses, handling potential discrepancies and errors that might occur during data transmission. It ensures that the system remains robust and that the user experience is seamless across devices.

The combination of Flask, Ngrok, and Python Anywhere provides a comprehensive and secure system architecture for the AI Voice Assistant. This setup not only facilitates smooth communication between different components of the system but also ensures that the assistant is accessible, efficient, and secure, regardless of the user's location. The next section will explore the implementation details of the Pythonanywhere server, demonstrating how it connects the Android application with the Ngrok server.

## Chapter 13: Implementation of the Python Anywhere Flask Server for the AI Voice Assistant

This chapter explores the implementation of the Python Anywhere Flask server, which functions as a critical intermediary facilitating robust communication between the AI Voice Assistant on the laptop and the Android application. The setup is designed to handle complex data flows, manage sessions, logs, and database interactions efficiently. Here, I detail the server's configuration and its pivotal role in sustaining the assistant's operations, emphasizing the server's logging and monitoring mechanisms, database management, and token handling.

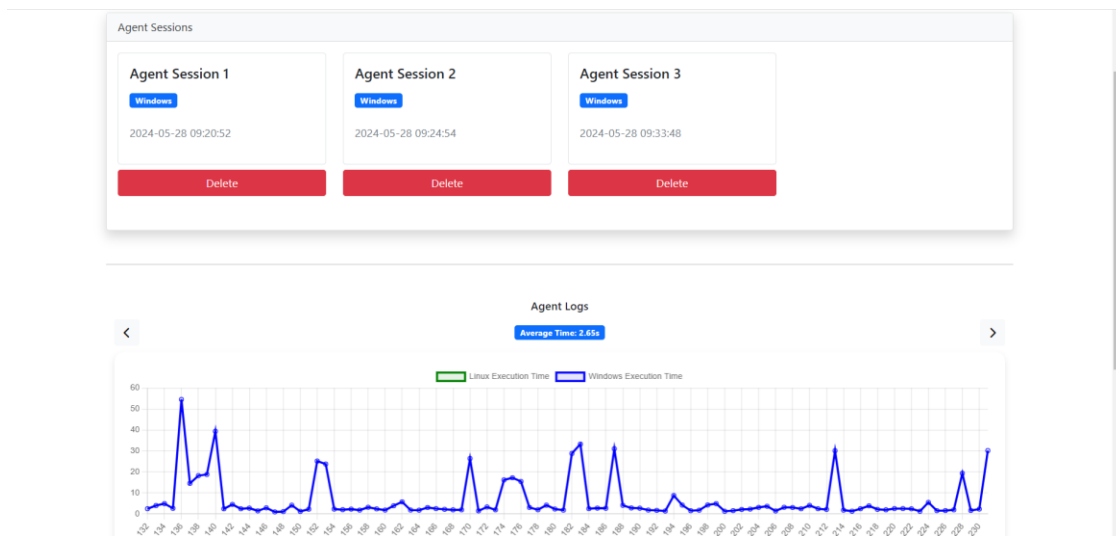


Figure 5 Demonstration of how the conversations are stored on the web server

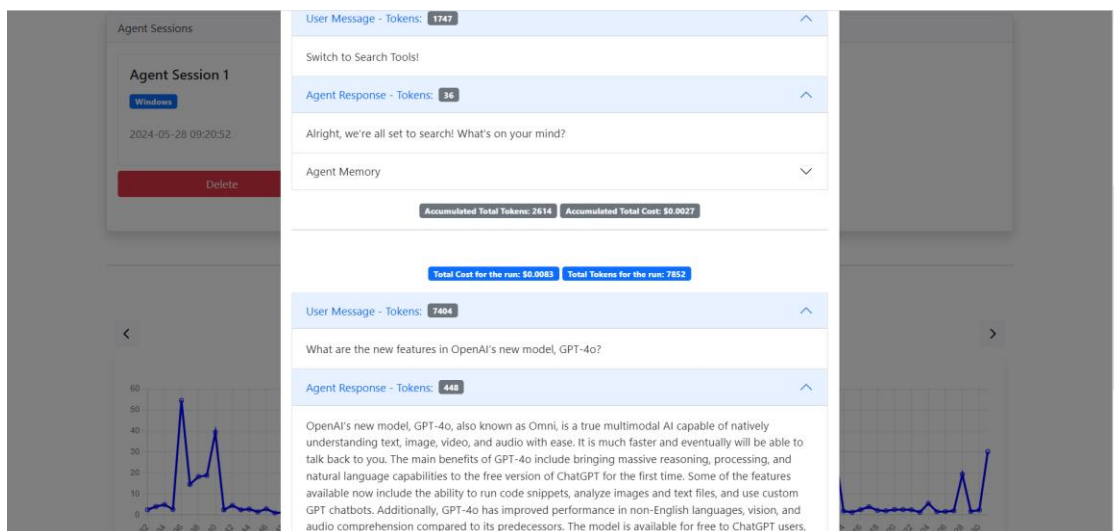


Figure 6 Demonstration of how the conversation history look like

The Python Anywhere Flask server functions as a critical intermediary, facilitating robust communication between the AI Voice Assistant on the laptop and the Android

application. This setup is designed to handle complex data flows and manage sessions, logs, and database interactions efficiently. Here, I detail the server's configuration and its pivotal role in sustaining the assistant's operations, emphasizing the server's logging and monitoring mechanisms, database management, and token handling.

## 13.1 Server Setup and Configuration

The Python Anywhere Flask server is configured to manage interactions seamlessly across multiple platforms:

**Flask Configuration:** The server is set up with Flask, a lightweight web framework suitable for handling HTTP requests efficiently. Flask's built-in capabilities are extended with SQLAlchemy for database operations, integrating models that represent tokens, session logs, and agent interactions.

**Database Initialization:** Using SQLAlchemy, the server interacts with a SQLite database (`vox_data.db`), which stores session data, token information, and logs. This relational database setup facilitates structured data storage and efficient query handling, crucial for tracking agent interactions and system statuses.

## 13.2 Logging and Monitoring

Comprehensive logging and monitoring are vital for maintaining the AI Voice Assistant's reliability and performance:

**Agent Session Management:** Each interaction with the AI agent is logged in `AgentSessionLog`, detailing user messages, agent responses, and token costs. This granularity allows for detailed analysis of the assistant's performance and resource usage.

**Error and Execution Logs:** The server maintains logs for errors (`ErrorLog`) and agent execution metrics (`AgentLog`), which are essential for debugging and optimizing the system. These logs help identify patterns and potential issues that could impact the assistant's functionality.

**Real-Time Log Streaming:** Incorporating real-time feedback mechanisms, the server streams logs to an administrative interface, enabling immediate visibility into the system's operations. This feature is crucial for remote monitoring, especially when direct access to the server or the assistant's local environment is not feasible.

### 13.3 Token Management and Secure Communication

Token management is a cornerstone of secure communication between the server and the Android application:

**Firestore Integration:** Firestore Admin SDK is used to manage notifications and data synchronization between the server and the Android app. This setup ensures that messages are delivered reliably and securely.

**Dynamic Token Handling:** The server manages authentication tokens meticulously, refreshing them as needed to maintain secure connections. This process is vital for protecting data exchanges from unauthorized access and ensuring that communications are encrypted.

### 13.4 Database Interactions and Session Management

The server's database interactions are structured to support robust data management and retrieval:

**Session and Log Deletion:** Old sessions and logs are purged regularly to maintain database performance and prevent data overflow. This automated cleanup helps manage storage efficiently and keeps the system responsive.

**Interactive Response Handling:** Through endpoints like `/ask`, the server processes user queries and fetches corresponding responses from the AI model. This interaction is logged meticulously, ensuring each step of the conversation is tracked for analysis and review.

The Python Anywhere Flask server is an integral component of the AI Voice Assistant's architecture, bridging the laptop application and the Android interface. It handles complex data flows, manages secure communications, and maintains extensive logs for monitoring and analysis. This server setup not only enhances the functionality and accessibility of the AI Voice Assistant but also ensures it operates securely and efficiently, adapting to various operational environments and user interactions. The subsequent sections will delve into the Android application and its integration with this server setup, highlighting how mobile interactions are managed and synchronized with the core AI functionalities.

## Chapter 14: Interface Demonstration of the Android Application

In the following section, I present a visual demonstration of the Android application interface, which serves as a crucial touchpoint for users interacting with the AI Voice Assistant. This interface is designed to be intuitive, facilitating seamless interaction between the user and the AI, enhancing the overall user experience.

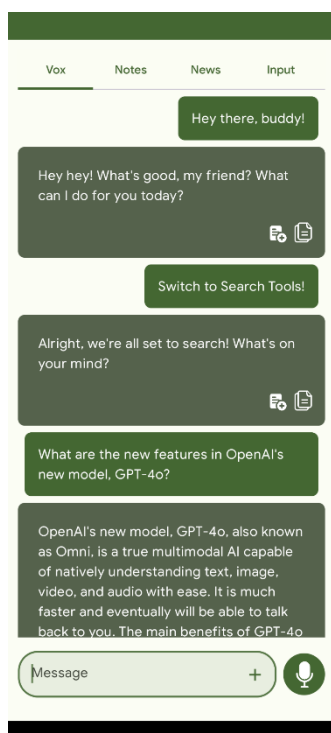


Figure 7 Main App Tab

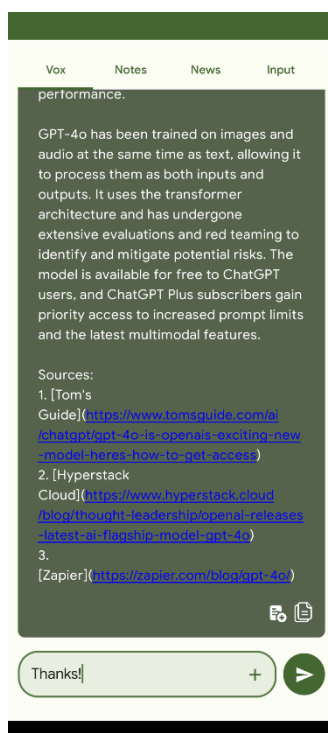


Figure 8 Demonstration of how the search works

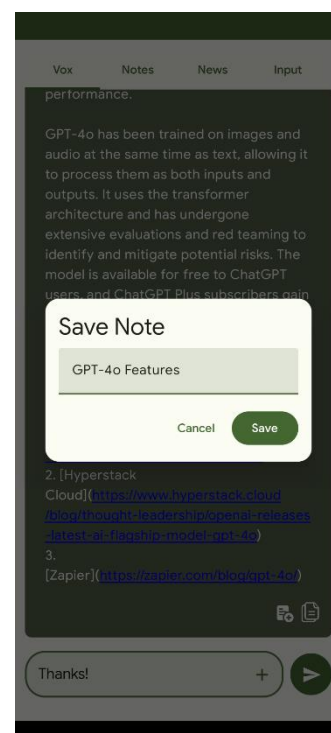


Figure 9 Demonstration of how answers can be saved as notes

### 14.1 Main Tab Overview

The main tab (first screenshot) of the Android application is where users primarily interact with the AI Voice Assistant. It features a simple and clean layout with a prominent area for entering voice commands. This design ensures that users can easily initiate conversations with the AI without navigating through complex menus.

### 14.2 Search Demonstration

Included in the demonstration (second screenshot) showcasing a typical search operation. This example illustrates how a user can inquire about general knowledge or specific queries. The AI processes these inquiries and displays answers directly within the interface,



leveraging the robust capabilities of the backend server to fetch and present information efficiently.

### 14.3 Note-saving Feature

An essential feature (third screenshot) of the application is the ability to save responses from the AI. After receiving an answer, users have the option to save this information as a note within the app. This functionality is demonstrated in the screenshots, highlighting how users can easily store and retrieve important information from their interactions with the AI. This feature not only adds convenience but also enhances the app's utility by allowing users to keep a record of critical data.

These screenshots not only illustrate the application's functionality but also emphasize its user-friendly design and efficient interaction flow. The integration between the Android application and the backend server ensures that the user experience is smooth and the responses from the AI are timely and relevant.

## Results

The AI-powered personal voice assistant developed in this project demonstrates a range of capabilities and robust performance across multiple tasks. Key functionalities include communication, file management, event scheduling, entertainment, information retrieval, and remote device interaction. Below is a summary of the main findings, highlighting the system's capabilities and performance metrics.

### System Capabilities

The system capabilities of this voice assistant include several key functions. In Communication and File Management, the assistant successfully integrates with Gmail, enabling users to send and manage emails through voice commands. It also interacts with Google Drive, allowing for the creation, modification, deletion, and backup of files, ensuring efficient file management without manual input. In Event Scheduling, Google Calendar integration lets the assistant schedule, modify, and delete events, providing users with seamless event management capabilities.

For Entertainment and Information Access, users can control their music on Spotify through voice commands. The assistant also performs web searches to fetch news, weather updates, and general information, enhancing its utility as an information retrieval tool. Remote Device Interaction includes Android Device Support, which allows remote file access and management, enabling users to manage files on their computers from their mobile devices.

### Performance Metrics

Performance metrics of the assistant highlight its compatibility and efficiency. It operates effectively on minimal hardware setups like Raspberry Pi 5 and mid-range laptops, and scales up for more resource-intensive tasks using dedicated GPUs. The system's speech recognition capabilities deliver quick and accurate transcriptions of user commands, while its use of advanced NLP models like OpenAI's GPT-3.5 Turbo ensures swift and coherent responses.

The assistant demonstrates high accuracy in task execution, credited to robust API integration and AI models. Its deep search functionality ensures high relevance and accuracy in web-searched information. Voice interaction is enhanced by Google Cloud's Text-to-Speech API, providing clear and natural voice responses. Real-time logging and error handling mechanisms provide immediate user feedback, improving the user experience.

Scalability and customization are achieved through the system's modular design, allowing for extensive customization and scalability to suit different user requirements and environments. Users can also tailor the assistant's functionalities and responses to their preferences, ensuring a personalized experience.

The AI-powered personal voice assistant developed in this project successfully demonstrates the integration of advanced AI technologies to provide a versatile, efficient, and user-friendly tool for hands-free operation. Its robust performance across various functionalities, coupled with its scalability and customization capabilities, highlights its potential as a valuable addition to personal technology solutions.

## Conclusion

This bachelor thesis on an AI-powered personal voice assistant highlights the ways AI can enhance user interaction with technology. The key contributions include versatile functionality, where the assistant supports various services like Gmail, Google Drive, Google Calendar, Spotify, and Android devices, aiding daily tasks through voice commands. It handles tasks such as sending emails, managing files, scheduling events, playing music, performing web searches, and executing command-line instructions securely.

Technical robustness is evident as the system works well on minimal hardware, from Raspberry Pi 5 to mid-range laptops, and can scale up for more demanding tasks using dedicated GPUs. It uses advanced AI models for speech recognition, natural language processing, and text-to-speech, ensuring accurate and efficient performance. The user-centric design of the assistant offers an easy-to-use interface, improving accessibility and efficiency, especially for users with physical limitations or those preferring hands-free operation. Its modular architecture allows for extensive customization and scalability, letting users adjust the system to their specific needs and hardware configurations.

However, the project has limitations such as dependency on external APIs and third-party services for critical functions like authentication, search, AI operations, voice detection, and voice transcription. If these services experience outages or changes, the assistant's functionality could be compromised. Limited testing resources have restricted the evaluation of the assistant's full performance and scalability since it was primarily tested on a Raspberry Pi 4, a mid-range laptop, and a desktop with an AMD GPU. Resource constraints as a student project hindered the ability to implement more advanced AI features and technologies that require substantial computational power and investment.

Future work to address these limitations and enhance the system includes reducing dependency on external APIs by developing in-house capabilities for critical functions. Enhancing testing and development resources by securing access to high-end hardware and more diverse testing environments would enable thorough performance evaluation and optimization. Expanding

computational resources would allow the implementation of more advanced AI features and technologies. Adapting to technological advances by staying updated with the rapid advancements in AI and integrating new features and improvements continuously is also proposed. Additionally, building a collaborative development team with diverse expertise in AI, software development, and user experience design would accelerate the development process, bring in new perspectives, and enhance the overall quality and functionality of the assistant.

In conclusion, this bachelor's thesis demonstrates the significant potential of AI in developing advanced personal voice assistants. While the project has notable limitations, addressing these through future work could result in a highly capable and reliable assistant that surpasses existing solutions like Google Assistant and Siri. The current implementation, despite its constraints, offers more features and flexibility, highlighting the promise of further development. With more resources, collaborative efforts, and continuous adaptation to technological advances, the AI personal voice assistant can evolve into an indispensable tool for everyday use. Here is a [quick demo](#) showcasing some of the capabilities of the voice assistant.

## References

- [1] Langchain-Ai. (n.d.). *GitHub - langchain-ai/langchain: 🦜🗨️ Build context-aware reasoning applications*. GitHub. <https://github.com/langchain-ai/langchain>
- [2] Snakers. (n.d.). *GitHub - snakers4/silero-vad: Silero VAD: pre-trained enterprise-grade Voice Activity Detector*. GitHub. <https://github.com/snakers4/silero-vad>
- [3] Picovoice. (n.d.). *GitHub - Picovoice/porcupine: On-device wake word detection powered by deep learning*. GitHub. <https://github.com/Picovoice/porcupine>
- [4] Openai. (n.d.). *GitHub - openai/openai-python: The official Python library for the OpenAI API*. GitHub. <https://github.com/openai/openai-python>
- [5] Openai. (n.d.-b). *GitHub - openai/whisper: Robust Speech Recognition via Large-Scale Weak Supervision*. GitHub. <https://github.com/openai/whisper>
- [6] Systran. (n.d.). *GitHub - SYSTRAN/faster-whisper: Faster Whisper transcription with CTranslate2*. GitHub. <https://github.com/SYSTRAN/faster-whisper>
- [7] Pytorch. (n.d.). *GitHub - pytorch/pytorch: Tensors and Dynamic neural networks in Python with strong GPU acceleration*. GitHub. <https://github.com/pytorch/pytorch>
- [8] PyAudio: Cross-platform audio I/O for Python, with PortAudio. (n.d.). <https://people.csail.mit.edu/hubert/pyaudio/>
- [9] Uberi. (n.d.). *GitHub - Uberi/speech\_recognition: Speech recognition module for Python, supporting several engines and APIs, online and offline*. GitHub. [https://github.com/Uberi/speech\\_recognition](https://github.com/Uberi/speech_recognition)
- [10] Facebookresearch. (n.d.). *GitHub - facebookresearch/faiss: A library for efficient similarity search and clustering of dense vectors*. GitHub. <https://github.com/facebookresearch/faiss>
- [11] Xianchen. (n.d.). *GitHub - xianchen2/Text\_Retrieval\_BM25: Python implementation of the BM25 for file retrieval*. GitHub. [https://github.com/xianchen2/Text\\_Retrieval\\_BM25](https://github.com/xianchen2/Text_Retrieval_BM25)
- [12] Openai. (n.d.-b). *GitHub - openai/tiktoken: tiktoken is a fast BPE tokeniser for use with OpenAI's models*. GitHub. <https://github.com/openai/tiktoken>
- [13] Huggingface. (n.d.). *GitHub - huggingface/transformers: 🤗 Transformers: State-of-the-art Machine Learning for Pytorch, TensorFlow, and JAX*. GitHub. <https://github.com/huggingface/transformers>
- [14] Praw-Dev. (n.d.). *GitHub - praw-dev/praw: PRAW, an acronym for "Python Reddit API Wrapper", is a python package that allows for simple access to Reddit's API*. GitHub. <https://github.com/praw-dev/praw>
- [15] GPT-3.5 Turbo - Quality, Performance & Price Analysis | Artificial Analysis. (n.d.). <https://artificialanalysis.ai/models/gpt-35-turbo>

