

ISSN(O): 2320-9801 ISSN(P): 2320-9798



## International Journal of Innovative Research in Computer and Communication Engineering

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)



Impact Factor: 8.771

Volume 13, Issue 4, April 2025

⊕ www.ijircce.com 🖂 ijircce@gmail.com 🖄 +91-9940572462 🕓 +91 63819 07438

DOI: 10.15680/IJIRCCE.2025.1304042

www.ijircce.com | e-ISSN: 2320-9801, p-ISSN: 2320-9798| Impact Factor: 8.771| ESTD Year: 2013|



International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

## **Enhancing a Speech-Driven AI Chatbot on Raspberry Pi with Remote Connectivity Features**

#### K. Aravinda Shilpa<sup>1</sup>, P. Harshasri<sup>2</sup>, S.L.S. Sanjana<sup>3</sup>, K. Sandhya<sup>4</sup>

Faculty, Department of Electrical Engineering, Andhra University College of Engineering for Women, India<sup>1</sup>

Student, Department of Electrical Engineering, Andhra University College of Engineering for Women, India<sup>2,3,4</sup>

**ABSTRACT:** This work presents the design and development of an AI chatbot on a Raspberry Pi using a range of contemporary technologies to enable interactive and intelligent voice-based dialogue. The system combines Google Web Speech API for speech-to-text, Gemini AI for generating dynamic responses, and Google text-to-speech for speech synthesis. The main features involve the capability to capture images from a camera via OpenCV, listen to audio from YouTube with yt-dlp and VLC, and perform a number of multimedia operations. The chatbot listens to voice commands in order to carry out activities like image capturing, playing or pausing audio, and intelligent response. The primary issues being faced are enhancing the hardware of Raspberry Pi to provide seamless multimedia performance, the handling of real-time voice identification, and optimizing efficient integration across different software parts. Future possibilities might involve enriching the addition of advanced image identification features, the inclusion of more APIs, and enhancing natural language processing towards more sophisticated interaction. This solution provides a customizable platform for any range of AI-based applications ranging from virtual agents to multimedia control.

KEYWORDS: Raspberry Pi, Gemini AI, Google Web Speech API, Google text-to-speech, OpenCV, YouTube.

#### I. INTRODUCTION

Chatbots with artificial intelligence have experienced a tremendous evolution in the past decade, with interactive and smart virtual support being provided through voice interaction. Such systems have transformed user interactions in many areas, ranging from customer service to personal assistants.[1] Prior research has worked towards combining speech recognition and natural language understanding to make chatbots more intelligent. As an example, speech recognition systems have become more effective at translating voice to text, whereas AI algorithms have developed to produce more response-oriented and context-specific answers. Yet, much of this functionality depends on robust hardware or cloud services, which are not practical for use in environments of limited resources, like embedded systems or even low-end devices like the Raspberry Pi.

This paper describes the design and implementation of an AI chatbot running on a Raspberry Pi, bringing together a range of technologies to create an efficient, low-cost, and highly versatile platform for real-time voice-based interaction.[2] Integrating Google Web Speech API for speech-to-text, OpenAI's Gemini AI for intelligent response, and gTTS (Google Text-to-Speech) for text-to-speech, the system provides a low-cost solution for creating an interactive, voice-responsive assistant. The chatbot can respond and understand user queries in a natural manner, allowing for a natural conversational flow.

Beyond basic conversation capabilities, this system also enables the use of multimedia features, broadening the chatbot's scope of interactions. This system builds upon past research on multimedia integration in virtual assistants by using OpenCV-enabled USB cameras for image capture, thereby enabling some level of vision tasks to be performed. The chatbot can also execute commands mid-stream and play audio files from YouTube via yt-dlp and VLC, which allows the voice-activated media tasks to be performed. This integration allows users to interact in a much more engaging way, as responses can be given in text, picture, and sound forms.

This work aims at optimizing hardware of Raspberry Pi so as to support smooth multimedia operations alongside real time voice recognition and response generation. Previous studies have attempted to find solutions to optimization of multimedia applications in embedded systems, but the integration of vision, speech, and real-time audio processing on a single board computer like Raspberry Pi is something that requires more investigation.[5]



This system shows that you can manage resources while still providing good performance in running a low-cost comprehensive AI assistant. Through the combination of speech, text, and multimedia functions, this system illustrates the power of inexpensive machines for enabling the development of real time AI powered virtual assistants which are able to perform intricate tasks. Furthermore, the modular nature of the system enables further developments such as implementing more sophisticated image recognition, additional API integrations for more features, and advancement of the natural language component to permit more sophisticated dialogues.

#### **II. LITERATURE REVIEW**

Jiang et al. (2015) reviewed the automatic online assessment of intelligent assistants, which require continuous advancements in machine learning algorithms to optimize interaction. Their study focused on simply automating virtual assistant assessments in real time with feedback loops and adaptive learning issues that could considerably improve the competence of virtual assistants over time. Apple's Siri, Amazon's Alexa, Google's Assistant, and Microsoft's Cortana have established the benchmark for the industry in AI-powered speech recognition systems. These virtual assistants are powered by natural language processing (NLP), automatic speech recognition (ASR), and deep learning methodology to increase user interaction and task performance. These stimulations are now being further researched so that they can use context and sustain more complex use queries with efficiency.[1]

Piyush et al. (2019) developed a Raspberry Pi-based voice-operated personal assistant that illustrates how low-cost embedded systems can support the development of AI-powered applications. Their research addressed voice recognition, command execution, and the use of cloud-based APIs to facilitate added capability. This report demonstrates the potential of low-cost hardware solutions to allow more people access to AI Assistants.[2]

Research by Leung & Wen (2020) recommended that restaurants integrate AI chatbots to create better consumer experiences and satisfaction levels. The authors clarify, chatbots that are based on context-aware NLP models could assist in analysing user preferences and making personalized food recommendations. This notion further posits that AI implementations could change the customer service experience to benefit users' engagement by providing relevant and customized human interactions.[3]

Karan & Sharma (2023) conducted a study that illustrated support for the notion that chatbot development is reliant on natural language processing (NLP), dialog management, and user interface design. They observed that while traditional rule-based chatbots depended on pre-written scripts, AI-powered chatbots rely on machine learning algorithms to dynamically generate responses. The authors advocated for the implementation of context-aware NLP models which improves the chatbot's ability to understand user intent and provide meaningful replies. The transition to an adaptive AI system provides an improved format for user-to-user conversations and increases the chatbot's efficacy in a plethora of applications.[4]

Renuka and Mulani (2017) investigated voice-controlled robots using Raspberry Pi, verifying that speech recognition commands could be accomplished without necessarily being connected to the cloud all the time. This research points to the significance of edge computing for AI voice assistants; commands can be recognized, and commands constructed, in real time even while offline.[5]

López, Quesada, & Guerrero (2017) advocated a systematic evaluation of chatbot quality. By using various attributes, including response accuracy, flow, and user satisfaction, the authors believed that chatbot evaluation should go beyond to measure only syntactic correctness to assesses semantic understanding and alignment with user expectations. The findings demonstrated the importance of continuing to develop evaluation frameworks to assess chatbot efficacy in a real-world context.[6]

In a study by Labadze et al. (2023), the authors advocate for the use of AI chatbots to improve instruction and student engagement, while also providing the student with personalized tutoring. The authors state that context-aware AI systems and chatbots can enhance student engagement through response tailoring based on interaction history and information about the student's emotional tone. This is in line with helping AI tutors understand the learner's needs and present opportunities that may lead to a more interactive and effective learning environment.[7]



## International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

| e-ISSN: 2320-9801, p-ISSN: 2320-9798| Impact Factor: 8.771| ESTD Year: 2013|

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

In a study by Alsayed et al. (2024), it was indicated that chatbots driven by artificial intelligence have a significant role to play in assisting students' mental health, as in civilizing their psychological disorders, such as anxiety, depression, fear, and stress. The authors stated that AI chatbots can work as conversational agents created to support students in transition to academic pressure, social staged anxiety and distressing emotional responses. The finding of this study shows the possibility for advanced AI for mental health support and a lead into the future for more advanced chatbot systems with emotional intelligence support.[8]

Amruta et al. (2017) explored the capabilities of the Raspberry Pi board, used to control obstacle-avoiding robots through voice recognition. These projects demonstrate the promise of artificial intelligence-based voice assistants in automation and robots. It can be seen how speech recognition may be incorporated into intelligent robotic systems for improved usability and efficiency in robotic systems.[9]

According to Caldarini et al. (2022), utilizing deep learning algorithms (e.g. recurrent neural networks - RNNs, transformer models) can facilitate chatbots in producing more accurate and engaging responses. The authors of this study noted that as AI models advance and grow in complexity, chatbot dialogues are transitioning from pre-scripted responses to dynamic replies, thus enhancing the quality of conversations.[10]

#### **III. OBJECTIVES**

The objectives of this study are as follows:

1. To build a voice-interactive, audio/video enabled, AI chatbot system using a Raspberry Pi that is scientifically interesting for voice interaction methods and chatbot API protocols.

2. To build a system with direct use of Google Web Speech API for real-time speech-to-text, and OpenAI's Gemini AI for intelligent responses.

3. To investigate multimedia capabilities with streaming audio from any source (e.g. YouTube) and capturing images with OpenCV and a collapsible USB camera.

4. To refine the system to measure its functionality on Raspberry Pi in light of hardware performance limitations, spanning speech recognition, intelligent response generation, and multimedia, with a goal of being as responsive as possible.

5. To build flexible and modular systems that encourage future work towards the thought of more advanced functions, data (text/image), the potential of improved NLP, and new APIs.

#### **IV. SCOPE OF THE STUDY**

This article describes a design and implementation of an economical AI chatbot solution based on the Raspberry Pi platform. The solution tracks several important technologies to enable voice interaction, multisensory experiences, and intelligent responses. The project involves:

- Using the Google Web Speech API for real-time speech recognition.
- Integrating OpenAI's Gemini AI so responses can be generated dynamically
- Integrating gTTS, or Google Text-to-Speech into the system to output speech
- Incorporating multimedia features, such as taking images with OpenCV and playing audio from yt-dlp and VLC.
- Optimizing the project for rapid performance under resource constraints of the Raspberry Pi hardware.

#### V. SIGNIFICANCE OF THE STUDY

This work is significant because it is able to offer a strong, low-cost AI assistant that will run on a low-cost platform like the Raspberry Pi. Unlike other AI assistants that rely heavily on costly hardware or cloud-based services, the proposal uses the cheap, small architecture of the Raspberry Pi but still has the flexibility of a number of functions. This makes it an accessible solution for a wide number of applications including home automation, educational tools, and customized virtual assistants. In addition, by adding multimedia capabilities such as image capture and audio streaming to the chatbot, the functionalities are enhanced to foster more enriched interactions, and allow for more complex use cases. Overall, simply by showing that efficient AI functionalities may be implemented through a cost-effective platform, this work adds to the body of literature on AI systems for embedded devices, and helps encourage opportunities for new cost-effective AI applications in healthcare, entertainment, and IoT.



## International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

| e-ISSN: 2320-9801, p-ISSN: 2320-9798| Impact Factor: 8.771| ESTD Year: 2013|

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

#### VI. METHODOLOGY

#### 6.1 System Overview:

The AI Chatbot system has been developed to support human-computer interaction via voice command. As the main hardware, the project implements a Raspberry Pi (4GB RAM), an I2S microphone for speech input, and a wired speaker to output the voice. The chatbot uses Google Web Speech API for speech-to-text features, OpenAI and Gemini AI for responses, and pyttsx3 & Google TTS for text-to-speech features. OpenCV is also included to capture and store images from a collapsible USB camera. YouTube audio streaming can also be done using yt-dlp and Video LAN Client (VLC). Remote access to the browser for monitoring and debugging is done through RealVNC.





#### 6.2 Speech Recognition:

Speech recognition plays a vital role in the AI Chatbot system, facilitating voice communication between users and the assistant in a natural way. The AI Chatbot uses Google Web Speech API speech recognition technology to transform voice into text using the latest advancements in machine learning. Upon speaking into the I2S microphone, the audio of the user's speech is processed and sent to the API for automatic speech recognition (ASR). ASR involves the study of how speakers produce phonetic patterns in order to segment the audio input into corresponding meaningful units of speech and text. Google Web Speech API understands several languages and dialects. After the speech has been texted into the API, the chatbot calculates the user input using natural language processing (NLP) processes. Google Web Speech API and the AI Chatbot relies on the cloud to assure processing is performed efficiently, although an existing internet connection must be available to enable the process. Speech recognition makes the chatbot capable of recognizing input commands more accurately, which allows voice interaction to feel more ordinary and natural for a user.

#### 6.3 Response Generation:

The AI Chatbot employs the Gemini AI API, which produces intelligent and context-aware responses based on user input. After the speech recognition system converts the spoken command into text, the chatbot sends the text to the Gemini AI API, which processes the text input with advanced natural language processing (NLP) and deep learning models. The API will analyse the text input, identify the user's intent and generate meaningful and relevant responses. Using machine learning algorithms, Gemini AI is capable of understanding conversational flow, generating responses having out potentially actually happened on their prior messaging history, and creating responses which are more natural and engaging.

Unlike traditional rule-based chatbots, which rely only on scripted responses, Gemini AI generates contextaware responses, giving the chatbot more flexibility to respond to a range of questions. The API helps ensure that the chatbot will respond to a diverse range of interactions, from fact-based questions to casual speech. Once the chatbot has developed an intelligent response, it will either display the response in the chat as text or convert the text into speech for voice output - both modes help create a seamless and interactive experience with the user.



### International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

| e-ISSN: 2320-9801, p-ISSN: 2320-9798| Impact Factor: 8.771| ESTD Year: 2013|

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

#### 6.4 Text to speech conversion:

The AI Chatbot uses the Google Text-to-Speech (Google TTS) API to turn text responses into speech that sounds natural. After the chatbot creates a response with the Gemini AI API, it sends that text over to Google TTS. This API then uses deep learning models to produce high-quality speech. Google TTS is versatile, supporting various languages, accents, and voice styles, so users can customize it to their liking.

The process of converting text to speech is quite structured. First, the API conducts a linguistic analysis to ensure that pronunciation, intonation, and rhythm are spot on. Once the speech is synthesized, it's streamed back to the chatbot, which plays it through a wired speaker for real-time interaction. Since Google TTS operates in the cloud, it does need an active internet connection to work. But thanks to its cutting-edge speech synthesis technology, it delivers clear, expressive, and human-like voice output, making the chatbot experience much more enjoyable for users.

#### 6.5 Image Capturing using Open CV:

The AI Chatbot uses OpenCV, a popular open-source computer vision library, to let users process images just by using their voice. With this cool feature, you can snap pictures using a collapsible USB camera that connects to the Raspberry Pi. When you say a specific command like "Capture Image," the chatbot springs into action, activating OpenCV to access the camera, take a shot, and save the image in a specific folder. The whole image capture process is pretty straightforward.

First, OpenCV gets the camera module up and running, tweaking settings like resolution and frame rate to make sure the image looks great. After capturing the image, it processes and saves it in a designated directory with a unique filename, making it easy to find later. The chatbot then gives you a little feedback, either by showing a message or announcing "Image Captured Successfully" through the Google TTS API.

By integrating OpenCV, the chatbot becomes even more powerful, opening the door to basic computer vision tasks. This means it could be used for things like object detection, facial recognition, or even document scanning in future updates. Overall, this setup ensures a smooth and automated image capture experience that works seamlessly within the chatbot's framework.

#### 6.6 YouTube Audio Streaming using Yt-dlp and VLC:

The AI Chatbot comes with a cool feature that lets you stream audio from YouTube just by using your voice. It uses yt-dlp, a handy command-line tool that pulls media from online sources, along with VLC media player to handle the audio playback. So, when you say something like "Play [song name] on YouTube," the chatbot jumps into action, searching for the right YouTube video link. With yt-dlp, it grabs the direct audio stream URL from the video without needing to download the whole thing. Then, it sends that audio to VLC, which plays it right away through the chatbot's wired speaker. This method makes sure you can enjoy audio streaming without hogging too much bandwidth. Since everything happens on the fly, you can listen to songs, podcasts, or any other YouTube audio content without a hitch. This integration really boosts the chatbot's multimedia features, turning it into a more interactive and engaging AI assistant.

#### 6.7 Remote Access using Real VNC viewer:

The AI Chatbot system makes it super easy to access remotely using RealVNC Viewer. This means you can keep an eye on, control, and troubleshoot the chatbot that's running on your Raspberry Pi, all from a device that's miles away. This feature is a game-changer because it lets you manage the chatbot without needing to be right there in front of it. You can easily perform software updates, fix any issues, tweak system settings, and boost performance. With remote control, you can interact with the chatbot from your laptop, desktop, tablet, or smartphone, which really enhances convenience and accessibility.

Setting up remote access is straightforward. You just need to install VNC Server on your Raspberry Pi and configure VNC Viewer on your remote device. Once that's done, you can create a secure, encrypted connection over your local network or the internet, giving you full control over the Raspberry Pi's graphical user interface (GUI) or commandline terminal. This setup allows for smooth interaction with the chatbot, whether you're running scripts, updating AI models, adjusting audio settings, or managing storage. Plus, RealVNC offers both direct network connections and cloudbased remote access, so you can connect from just about anywhere. This is especially handy when the chatbot is in a



place where you can't easily get to the hardware. Being able to access the system remotely also makes maintenance easier, speeds up debugging, and enhances scalability, making the chatbot more versatile for real-world use.

On top of that, you can boost your remote access by integrating SSH (Secure Shell). This allows advanced users to manage the system using terminal commands without needing a graphical interface. So even if you're dealing with a slow network connection, developers can still control and update the chatbot effectively. By combining VNC for visual control and SSH for command-line management, the system strikes a great balance between user-friendliness, security, and flexibility, ensuring a strong and scalable remote experience.

#### 6.8 Implementation Workflow:

The below Flowchart (Fig 2) shows us the implementation workflow of our chatbot.



The AI Chatbot operates through a well-organized workflow that brings together various technologies like speech recognition, natural language processing (NLP), response generation, text-to-speech conversion, image processing, and streaming audio from YouTube. When a user gives a voice command, the chatbot picks up the audio through a microphone and uses the Google Web Speech API to convert it from speech to text. This transcribed text is then analysed by the Gemini AI API, which crafts a smart response based on the context.

That response is transformed into speech using the Google TTS API and played back through a wired speaker, allowing for real-time interaction. If the command requires image processing, OpenCV kicks in to capture an image with a USB camera, process it, and save it in a specific folder. For playing multimedia, the chatbot utilizes yt-dlp and VLC media player to stream YouTube audio when requested by the user. Plus, the whole system can be accessed remotely via RealVNC Viewer, letting users monitor, control, and troubleshoot the chatbot from any device. This modular approach guarantees smooth and efficient operation of all the chatbot's features.



### International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

| e-ISSN: 2320-9801, p-ISSN: 2320-9798| Impact Factor: 8.771| ESTD Year: 2013|

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

#### 6.9 Conclusion:

The AI Chatbot brings together speech recognition, natural language processing, image processing, and media streaming into a neat, interactive system powered by a Raspberry Pi. It's designed for remote access making it a versatile and scalable option for voice-driven interactions between humans and computers.

#### **VII. RESULTS AND DISCUSSIONS**

This section will dive into a thorough discussion about how well the chatbot is performing. 7.1 Response Generation using Gemini AI:

The chatbot did a great job generating responses through the Gemini AI API, offering replies that were both relevant and aware of the context. It handled user queries efficiently and provided answers quickly, especially when the network was stable. The chatbot kept the conversation flowing smoothly and responded accurately to a range of prompts. That said, there was a bit of a delay when it came to more complex queries, as the responses relied on cloud-based API calls. Looking ahead, incorporating local AI models could help speed things up and enhance offline capabilities.



fig -2: RESPONSE GENERATION USING GEMINI AI

#### 7.2 YouTube Audio Streaming using yt-dlp and VLC:

The combination of yt-dlp and VLC media player made it possible to stream YouTube audio in real-time using just voice commands. The setup effectively pulled audio from YouTube and played it through the connected speaker. When we tested it, the streaming quality was crisp and smooth as long as the internet connection was strong. However, we did notice some occasional buffering, which was likely due to network fluctuations and delays in API processing. To enhance playback and minimize interruptions, we could look into using caching mechanisms or adaptive streaming techniques.

Provide 10 mil top me top me top me top	RESTART: /home/cohoronylac/Dackbar/ani/cation
The transmission of the type and community. In the is shown on the same of the decision of the same of	9.8
Ar Chaibot is scribe. By something Ar Chaibot is scribe. By something by something.	am the pygame community, https://www.pygame.org/contribute.html
<ul> <li>Litter and men eine eine</li></ul>	t is active. Say something
Processing speech." With the weak help you tousy? Processing speech Processing spe	
Nor wand is milding         Nor wand is mil	g speech
Ar i now can I neip you Eoday? Ar i now can I neip you Eoday? Processing spectrum: Processing spectrum: Processing spectrum: Processing intervention of the III Processing intervention of the IIII Processing intervention of the IIIII Processing intervention of the IIIII Processing intervention of the IIIII Processing intervention of the IIIIII Processing intervention of the IIIIIIII Processing intervention of the IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII	hello
<ul> <li>Intendage</li> <li>Inte</li></ul>	an I help you today?
Listendag Hording	
Processing speech Here and the second se	
Listening Togesting speech Togesting speech	g speech
Processing speech The construction of the	
<ul> <li>You said: piay jingte bells</li> <li>You said: piay jingte be</li></ul>	g speech
Seerching YouTube for: jingle bells Playing: https://fisupple.tells Playing: https://fisupple.tells Playing: https://fisupple.tells Playing: https://fisupple.tells State: a constraint of the state of t	play jingle bells
<ul> <li>Pinying: https://rfs-sn-ghpa.fidey.good/wiles.cov/lide/eex/fixessaadadasexades ts286a5a50x201668Hsg1wp6x384A481p2200890x1001000000000000000000000000000000</li></ul>	YouTube for: jingle bells
Veryesät Ise=TVHTML5&sefc=1&Exp=4532434ackprev2cvprv%2cvpucK2cmmk2cmk%2cmk%2cmk%2cmk%2cmk%2cmk%2cmk%	10.722017/GCHS9TWDOKS9KAABJD=2409%SAA0709KCA100CKA100CKAA0607KA71000KBA430805A0101 10.0-ALSADDKE04YEMBEDIRICS500009Xx31003V2005124421101200-313262298009-sni- essLyes&Ryc-Egy02205NQX3DX3DA080L471715%p1-4154785780400252004011510446933080 FX2C3n_0Mp3-h556880=3134505W94Cub752291851210459124032 FX2C3n_0Mp3-h556880=3134505W94Cub75229185120459124032 FX2C3n_0Mp3-h556880=3134505W94Cub75229185120459124032 FX2C3n_0Mp3-h556880=3134505W94Cub75229185120459124032 FX2C3n_0Mp3-h556880=3134505W94Cub75229185120459124032 FX2C3n_0Mp3-h556880=3134505W94Cub7522918512040048045130480 FX2C3n_0Mp3-h556880=3134505W94Cub75229 FX2C3n_0Mp3-h556880=3134505W94Cub75229 FX2C3n_0Mp3-h556880 FX2C3n_0Mp3-h558880 FX2C3n_0Mp3-h558880 FX2C3
Ve joči nazči ta na kačeno ve kači redu i resis Uzici k Učina zemini zemini zemini zemini zemini za ve kači k Učina za ve kači	18c=TVHTML5&sefc=1&txp=4532434&n=0%2Cbu1%2Cvprv%2Csvpuc%2Cmime%2Cns%2Cr
CIDR/2	acitsad%2Cspurce%2Crequiress1%2Cxpc%2Cmu%2Cmu%2Cmu%2Cmu%2Cmu%2Cmu%2Cmu%2Cmu
qh%2Cg1r%2Ct timacCost Sig=ACuhMUOwRgThAWGL3SVUC2phrb2 ms%2Ch1tCwndpp&ls1g=ACuhMUOwRgThAWGL3SVUC2phrb2 ms%2Ch1tCwndpp&ls1g=AJrQdSwkAIg2zPKD6t65J1gAAZ t EAsHLe11AgruCDorKOPDQBVcJWf1XOCt -4%VUS5B051g=AJrQdSwkAIg2zPKD6t65J1gAAZ W9FLcbj6HnEjggWam2GeNIU-79gACIHGRH7alroZtm1JM6StNZBNJ-656-51PuBZDqVL4-fcz W9FLcbj6HnEjggWam2GeNIU-79gACIHGRH7alroZtm1JM6StNZBNJ-656-51PuBZDqVL4-fcz	2011ag%2Cdur%2Clmt&lsparams=met%2Cmt%2chpo2l0d8fFIgNlbtBX2zHfZQ0klZMkAA1
m ms%2CinitcwndDp5stDoDgBvcIwrIXOCi-4yXID5I0fo60x30s20,00-580-51puBZDqVI4-fcz EASHLeiJAgruCDoFKOPDQBvcIwrIXOCi-4yXID5I0fo60x30s20,00-580-51puBZDqVI4-fcz W9FLcbj6HnEjggWamzGeNIU-79gACIHGRH7alroZtmIJMoSTNZBNJ-680-51puBZDqVI4-fcz	2CC LEMS2 Cd La ACubMUOWRgIhAMv6LJSV9CK2phPlata AlfodSswRAIgZzPXD6tc6JIgS4x
a EASHLe11AgruCDoTKOPJQOVC79gACIHGRH7alroZtm1JMoStNZBN3 000 W9FLcbj6HnEjggWamzGeNIU-79gACIHGRH7alroZtm1JMoStNZBN3 000	cwndbpsa tsig
(1) 10 10 10 10 10 10 10 10 10 10 10 10 10	grucDorkOFDQUVCI nEjggWamzGeNIU-79gACIHGRH7alroZtmlJMeStN2BN3-000 12
33°C Haze	0 10 10 10 10 10 10 10 10 10 10 10 10 10
	9 33*C Haze ~ 0 + -

fig -3: YOUTUBE AUDIO STREAMING USING yt-dlp and VLC



#### 7.3 Image Capturing using OpenCV:

The chatbot did a great job of snapping pictures using OpenCV whenever it heard a voice command. The USB camera worked well, capturing images and saving them in a specific folder with filenames that made it easy to find them later. The image quality was pretty good in bright lighting, but things got a bit fuzzy in low-light situations. It would be awesome to add some automatic brightness adjustments or real-time enhancements to boost performance. Plus, incorporating features like edge detection or object recognition could really take the chatbot's image-processing skills to the next level in future updates.

ocoupler/L	HDLE Shell 2 9/21
ow Help	File Edit Shejl Debug Options Window Help
L	Hython 3.9.2 (default, Dec 1 2024, 12:12:57) [GCC 10.2.1 20210110] on linux Type "help", "copyright", "credits" or "license()" for more information. >>>
	RESTART: /home/robocoupler/Desktop/vol/volce.py Hello from the pygame community. https://www.pygame.org/contribute.html AI Chatbot is active. Say something
	Listening Processing speech
xPBB9JNr	Processing speech You said: hello
	AI: How can I help you today?
rack if	Listening Processing speech You said: capture image
HUCH II	Listening
tures/"	Processing speech Image saved at: /home/robocoupler/Pictures/pictures/captured_image.jpg AI: How can I help you today?
0	Listening Processing speech
	Listening
t_noise	📴 📫 🚱 🏹

fig -4: IMAGE CAPTURING USING OpenCV

#### VIII. LIMITATIONS AND FUTURE EXPANSIONS

#### 8.1 Limitations:

While the AI Chatbot does a great job of integrating features like speech recognition, response generation, YouTube audio streaming, and image processing, we did notice a few limitations during testing:

- *Network Dependency* The chatbot relies quite a bit on a stable internet connection for its various functions, including the Gemini AI API for generating responses, Google TTS for text-to-speech, YouTube streaming via yt-dlp, and the Google Web Speech API for speech recognition. If the network is weak or unstable, it can lead to delays or even failures in generating responses and streaming media.
- *Latency in Response Generation* Because the chatbot processes user queries through a cloud-based API, we observed some minor delays in getting responses. This lag could be a bit of a hiccup for real-time conversations, especially in interactive settings.
- *YouTube Streaming Delays* The yt-dlp and VLC-based streaming system performed well when the internet connection was strong. However, we did encounter some buffering issues and occasional delays in pulling video URLs when the network speed varied.
- **Speech Recognition Challenges** The chatbot did well in quiet settings but had a tough time in noisy environments, which led to some misinterpretations or incomplete transcriptions. This can affect how accurately it executes commands.
- *Image Quality Variations* The OpenCV-based image capturing worked nicely in good lighting, but images taken in low-light conditions didn't come out as clear. Plus, the system doesn't have advanced image processing features like automatic brightness adjustments.
- *Hardware Limitations* The chatbot operates on a Raspberry Pi with 4GB of RAM, which can limit its multitasking abilities. Running multiple AI processes at the same time can slow down performance, particularly during speech processing and streaming.



International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

| e-ISSN: 2320-9801, p-ISSN: 2320-9798| Impact Factor: 8.771| ESTD Year: 2013|

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

#### **8.2 Future Expansions:**

To tackle these challenges and boost the chatbot's performance, we can consider a few exciting improvements for the future:

- **Offline** AI Models By incorporating a local AI model for generating responses (think on-device NLP), we can cut down on our dependence on cloud-based APIs. This would not only speed up response times but also allow for offline functionality.
- *Optimized Speech Recognition* Introducing noise reduction algorithms or offline speech recognition models can significantly improve command accuracy, especially in noisy settings.
- *Enhanced YouTube Streaming* Implementing caching techniques or preloading content that's frequently watched could help minimize buffering times and make playback much smoother.
- *Improved Image Processing* By adding features like automatic brightness and contrast adjustments in OpenCV, we can enhance image quality even in low-light situations. Future upgrades might also include object detection or facial recognition to broaden the range of applications.
- *Alternative Remote Access Methods* Rather than relying solely on RealVNC, we could explore a web-based dashboard or SSH-based command control, offering a more lightweight and user-friendly remote management option.
- *Hardware Upgrades* Upgrading to a Raspberry Pi model with more RAM or adding an external processor could really boost multitasking capabilities and overall system performance.

#### **IX. CONCLUSION**

The AI Chatbot has been successfully developed and tested, bringing together speech recognition, response generation, YouTube audio streaming, and image processing into a compact and efficient system powered by a Raspberry Pi. It showcased its ability to process voice commands in real time, generate meaningful responses using the Gemini AI API, convert text to speech with Google TTS, capture images through OpenCV, and stream audio from YouTube using yt-dlp and VLC. Plus, with remote access via RealVNC, monitoring and control became a breeze, enhancing the system's overall usability.

The results show that the chatbot offers accurate voice-based interaction, making it a practical solution for a variety of real-world applications. The response generation system performed admirably, keeping conversations contextually accurate and engaging. Speech recognition worked well in quiet settings, although it faced some challenges in noisy environments. YouTube streaming was reliable with a strong internet connection, and image capturing was effective, though improving low-light image processing could boost overall performance.

Even with its reliance on a network and some minor latency issues, the chatbot proved to be effective under stable conditions, demonstrating its potential as an AI-driven voice assistant. Its modular design opens the door for future enhancements, like offline speech processing, local AI models for generating responses, better image enhancement techniques, and optimized multimedia streaming. Additionally, incorporating edge AI processing could lessen the dependence on cloud services, speeding up response times and making the chatbot even more versatile.

With ongoing improvements, this chatbot could be expanded for smart home automation, virtual assistants, educational tools, and AI-powered customer service applications. Its knack for interpreting voice commands, generating intelligent responses, and handling multimedia tasks makes it a scalable and adaptable solution for voice-based human-computer interaction.

#### REFERENCES

- Jiang, J., Awadallah, H. A., Jones, R., Ozertem, U., Zitouni, I., Kulkarni, R. G., & Khan, O. Z. (2015). Automatic online evaluation of intelligent assistants. Proceedings of the 24th International Conference on World Wide Web, 506–516.
- 2. Piyush, V., Singh, J. P., Jain, P., & Kumar, J. (2019). Raspberry PI based voice-operated personal assistant. International Conference on Electronics and Communication and Aerospace Technology (ICECA).
- 3. Leung, X. Y., & Wen, H. (2020). Chatbot usage in restaurant takeout orders: A comparison study of three ordering methods. Journal of Hospitality and Tourism Management, 45, 377-386.

© 2025 IJIRCCE | Volume 13, Issue 4, April 2025|

DOI: 10.15680/IJIRCCE.2025.1304042

www.ijircce.com



International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

| e-ISSN: 2320-9801, p-ISSN: 2320-9798| Impact Factor: 8.771| ESTD Year: 2013|

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

- 4. Karan, A., & Sharma, V. (2023). Design & Development of Chatbots. International Journal of Engineering Technology and Management Sciences, 7(3), 96-105.
- Renuka, P. K., & Mulani, A. O. (2017). Raspberry Pi-based voice-operated Robot. International Research Journal of Engineering and Technology (IRJET), 2(12).
- 6. López, G., Quesada, L., & Guerrero, L. A. (2017). *Evaluating Quality of Chatbots and Intelligent Conversational Agents.* arXiv preprint arXiv:1704.04579.
- 7. Labadze, L., Grigolia, M., & Machaidze, L. (2023). *Role of AI chatbots in education: Systematic literature review.* International Journal of Educational Technology in Higher Education, 20(56).
- 8. Alsayed, S., Assayed, S. K., Alkhatib, M., & Shaalan, K. (2024). *Impact of Artificial Intelligence Chatbots on Student Well-being and Mental Health: A Systematic Review.* People and Behavior Analysis, 2(2).
- 9. Amruta, N., Doddamani, A., Deshpande, D., & Manjramkar, S. (2017). *Raspberry Pi-based obstacle avoiding robot*. International Research Journal of Engineering and Technology (IRJET), 4(2).
- 10. Caldarrini, G., Jaf, S., & McGarry, K. (2022). *A literature survey of recent advances in chatbots*. Information, 13(1), 41.



INTERNATIONAL STANDARD SERIAL NUMBER INDIA







# **INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH**

IN COMPUTER & COMMUNICATION ENGINEERING

🚺 9940 572 462 应 6381 907 438 🖂 ijircce@gmail.com



www.ijircce.com