

A.I.dan: ChatGPT Integrated Virtual Assistant

Andrew Scott

Leonardo Garcia

Brahmtegg Minhas

Project #25

TA Hanyin Shao

02/23/23

ECE 445

Table of Contents

1. Introduction	3
1.1 Problem	3
1.2 Solution	3
1.3 Visual Aid	4
1.4 High Level Requirements	4
2. Design	5
2.1 Block Diagram	5
2.2 Physical Design	5
2.3 Remote Block	6
2.3.1 PC	6
2.3.2 ChatGPT API	7
2.4 Onboard Block	7
2.4.1 Microcontroller	8
2.4.2 Input Subsystem	8
2.4.3 Output Subsystem	8
2.4.4 Power	9
2.5 Tolerance Analysis	9
3 Cost and Schedule	11
3.1 Cost Analysis	11
3.2 Schedule	12
4 Ethics and Safety	13
5 Citations	14

1. Introduction

1.1 Problem

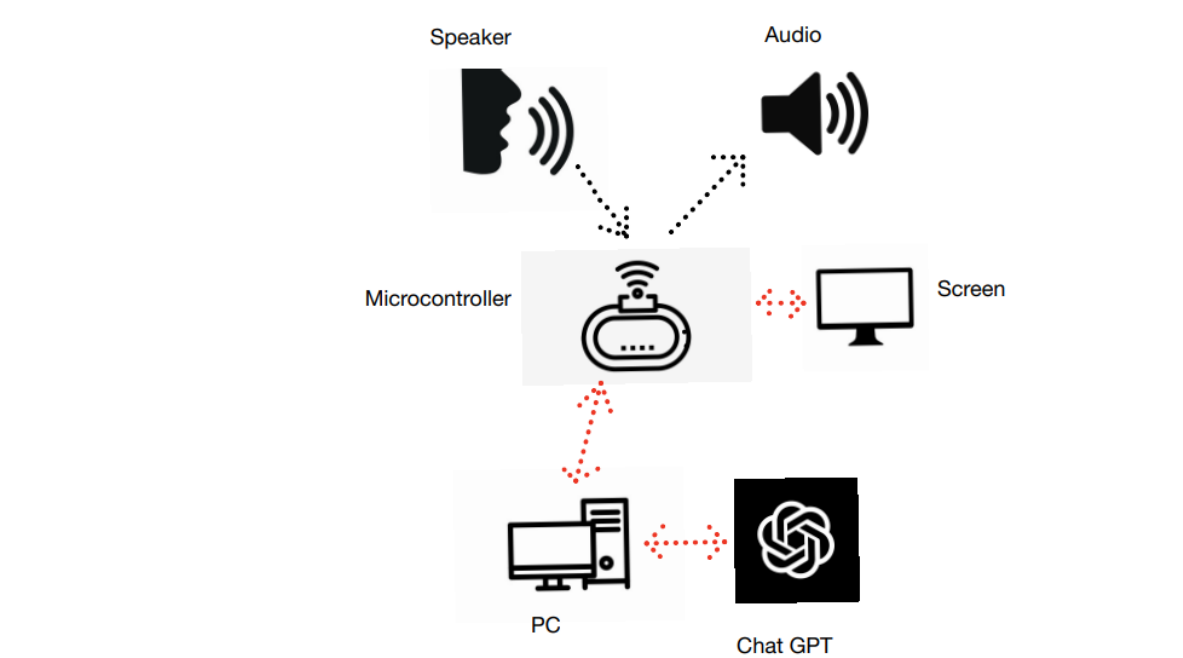
All existing virtual assistants use a search engine (primarily Google Search) as a data retrieval tool for answering questions posed to them by a user. The responses given by a search engine to a virtual assistant is a result or excerpt from the top search result that comes when searching for that question. However, this result is often not useful or relevant. Search engines are built to present multiple information sources from the web when presented with a question, not necessarily output one definitive answer. Furthermore, searches respond poorly to requests, unable to produce an output of a specified form.

OpenAI's ChatGPT is an AI language model that can respond to questions, generate text and have conversations. It is an extremely useful tool to respond to questions or present a singular output in a specified format. However, accessing ChatGPT requires going to the ChatGPT website, logging in and manually typing the question. This process is cumbersome and takes more time than it needs to.

1.2 Solution

Our solution is a device that integrates ChatGPT with a virtual assistant in order to nullify the weaknesses of both tools. Because ChatGPT is trained on human language, it is much better suited for a virtual assistant, which uses human language as input. As such, integrating ChatGPT into a virtual assistant would yield a much better tool for getting information and tailored responses to user questions. Our solution will allow the user to present the virtual assistant with a prompt or question. The device will then return the response that ChatGPT gives to the user input, both through a speaker and through a screen on the assistant.

1.3 Visual Aid

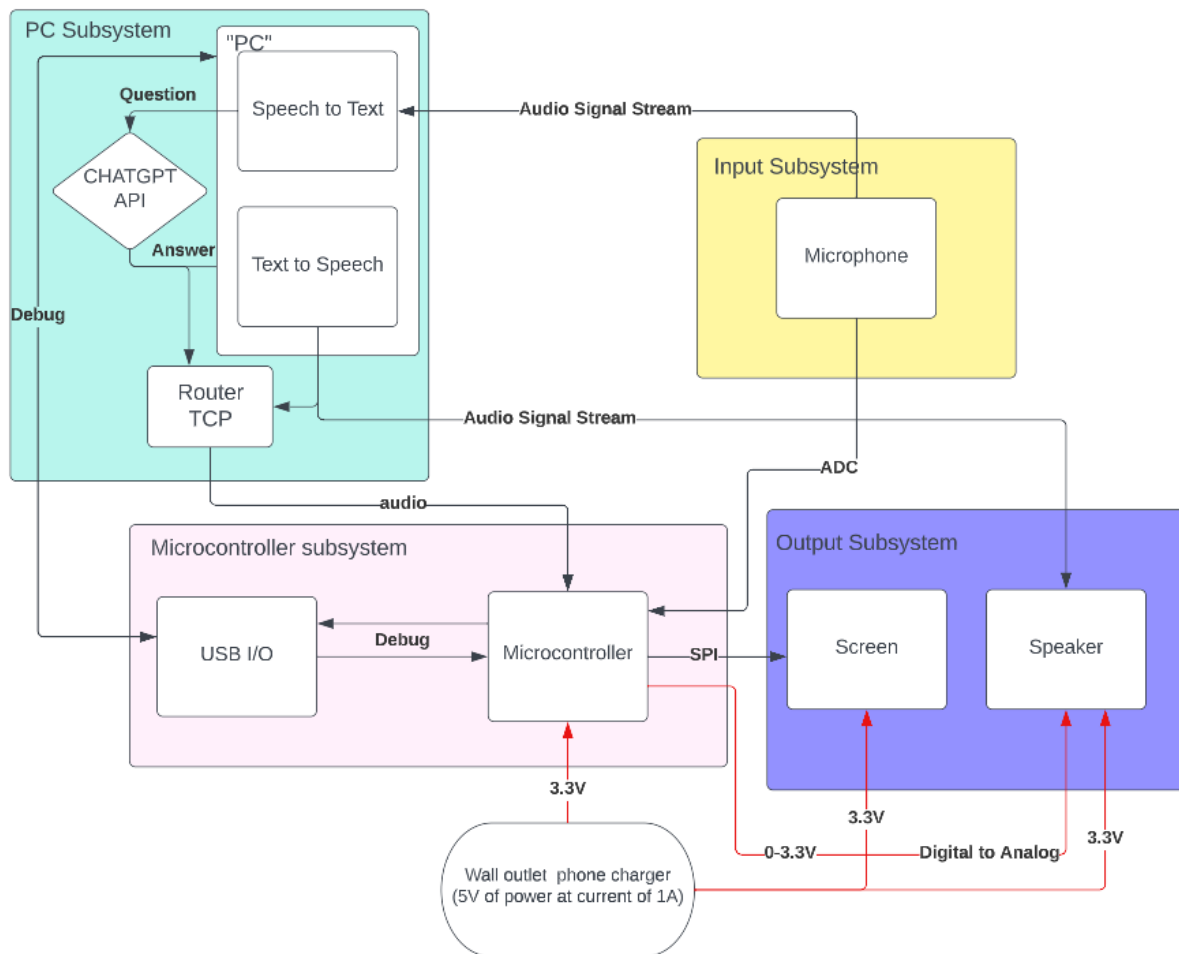


1.4 High Level Requirements

- At least 75% of attempted basic interactions should be successful. Basic interactions are defined as questions based entirely on words included in a pre-trained speech to text model.
- Code (Markdown) or other text must display properly given a successful question. This is verified by inputting the same question to chatGPT on a separate device, or by logging the I/O of the PC's API Requests.
- Given a successful question with an appropriate audio response, the speakers must output an understandable Text To Speech interpretation. While mispronunciation is acceptable, greater than 80% of users must be able to correctly interpret the audio.

2. Design

2.1 Block Diagram



2.2 Physical Design

The four major physical components of our design are a small (2" diagonal) screen, our PCB, a microphone, and a speaker. These components will all be mounted inside a single 3" 3d-printed cube, with the speaker and microphone on top, the screen on the front, and the PCB on bottom. The USB Ports for power and Debug will be cut out of the back side and directly soldered to the pcb.. Outside of the enclosure will be a standard Wall plug DC 5V power supply, wired to the back of our cube through a power-only type c connection.

2.3 Remote Block

The PC block consists of an external computer connected to the microcontroller through WiFi. The PC performs the majority of the heavy compute lifting of the entire project, performing the Speech recognition, ChatGPT API Calls, and Text to speech processing.

2.3.1 PC

The PC accepts audio input from the microcontroller (streamed via WiFi), and constantly translates it to text in chunks, listening for the keyword, “Hey A.I.Dan” (pronounced Aidan). When it hears the phrase, the PC will listen to the next sentence it gets (waiting for a long pause in audio to stop listening), convert the data to text and send it to ChatGPT’s API.¹ The API text output is then transmitted as both text (to be displayed on the screen), and as audio (through a text-to-speech conversion) back to the microcontroller to be output to the user.

Table 1: PC Subsystem Requirement and Verification

Requirement	Verification
<i>The PC must maintain an open WiFi port for 1 hour, untouched, to listen for microcontroller input.</i>	<i>Establish a TCP connection between the PC and the ESP32 from the PC. After 1 hour of inactivity, check to see if the connection is still active by checking available devices on the network and searching for the ESP32 IP address.</i>
<i>The PC must be able to convert clear and simple speech into text using a pre-trained model with an 80% or greater accuracy rate</i>	<i>Collect multiple audio recordings that contain clear and simple speech. Run the speech to text model on the audio recording and verify that the text output matches the spoken words with the specified (>80%) accuracy rate.</i>
<i>The PC must be able to analyze text to determine whether a question is valid (contains the wake words “Hey A.I.dan” or similar) with an 80% or greater accuracy rate.</i>	<i>Collect multiple audio recordings that either do or do not contain the wake words “Hey A.I.dan.” Run the audio recordings through the speech to text model. Check the PC console output for a log saying “Wake Word Detected” and ensure this occurs at a frequency >80% for recordings with the wake word included. For recordings that do not contain the wake word, the log should display “Wake Word Detected” less than 20% of the time.</i>

<i>The PC must be able to convert standard text into speech using a pre-trained model, understandable at least 80% of the time.</i>	<i>Input text files containing standard text into the text to speech module. Analyze the output audio files to ensure the output speech is understandable and accurate to the source text, within reason, greater than 80% of the time.</i>
<i>The PC must be able to output both a digital audio signal and text to the ESP32 microcontroller through WiFi, only outputting one or none 20% or less of the time</i>	<i>Collect a digital audio file and a text file. Establish a TCP connection with the microcontroller. Send the text file, followed by the audio file, over the TCP connection to the microcontroller. Ensure that the microcontroller received data of the same size as the digital audio file and the text file for 80% or more of the trials</i>

2.3.2 ChatGPT API

The ChatGPT API generates the ChatGPT response. It receives text input from the PC and outputs the response of chatGPT as text back to the PC.

Table 2: ChatGPT API Subsystem Requirement and Verification

Requirement	Verification
<i>The API must be able to generate an answer to a text input with a text output within 4 seconds of the API call</i>	<i>Collect a text file containing a prompt for chatGPT. Run the text file through the ChatGPT API model, using python's timeit library to track how long it took Analyze the output text file of the ChatGPT API model to ensure that it contained a reasonable response to the prompt and received it within the timeframe.</i>

2.4 Onboard Block

The onboard block takes in audio data through a microphone and transmits it through WiFi to the PC. It also receives data through WiFi to be output to a screen and through a speaker

2.4.1 Microcontroller

Consists of an ESP32 microcontroller with built-in WiFi. It serves as the liaison between input and output data. Input data is received from the microphone and transmitted to the PC through WiFi. From the PC, it receives data back both in the form of text and audio. The microcontroller sends the audio directly to the speaker to be output to the user and converts the text input to visual and sends it to the screen through the SPI protocol.

Table 3: Microcontroller Subsystem Requirement and Verification

Requirement	Verification
<i>The microcontroller must be able to host a TCP server using Wifi to interface with the PC.</i>	<i>Establish a TCP connection between the PC and the microcontroller. Ensure the TCP connection is valid by checking the remote client (the PC) for a message from the microcontroller which states: "ESP32 connected".</i>

2.4.2 Input Subsystem

The output subsystem contains a microphone, which is constantly streaming input audio to the microcontroller's built-in analog-to-digital converter.

2.4.3 Output Subsystem

The output subsystem displays data received from the microcontroller through a screen through the SPI protocol and through the speaker following a digital-to-analog conversion.

Table 4: Output Subsystem Requirement and Verification

Requirement	Verification
<i>The screen must be able to properly display code snippets, which requires at least 16-bit color.</i>	<i>Connect the microcontroller to the PC via USB. Download an SPI Code Display Test program which contains a formatted code snippet with text of 16 different colors. Ensure that the code displays legibly and 16 colors are distinguishable and present on the display.</i>

2.4.4 Power

The power subsystem powers all onboard components.

Table 5: Power Estimates for Selected Components

Component	Power Draw
Screen	1 W ₂
Microcontroller	2.50 W ₃
Speaker	0.5 W ₄
Digital-to-Analog Converter	< 0.1 W ₅
Total	4.0 - 4.1 W

Based on the above estimates, a standard wall plug DC 5V power supply operating at 1A should be sufficient.

2.5 Tolerance Analysis

Data Throughput (PC-Microcontroller data transfer)

The largest identified issue associated with this design is the data management on the microcontroller, and its interaction with the PC. Even the most robust microcontrollers only contain around 512 KiB of memory, which corresponds to roughly a second of audio information. Data will be streamed to the PC as fast as possible, and taking and running Text-To-Speech will act similarly quickly, but both streams could prove too much for the microcontroller.

To demonstrate feasibility, assume 100 KiB of memory allocated for each individual stream, up and down. Wifi speed is quite high relative to these numbers, so the actual communication shouldn't be an issue (ESP32-C3 supports 150mbps, and any PC we would be using will support at least 300mbps) Assume standard 44khz audio at 12 bits of digital resolution.

$$44000\text{hz} \cdot 12b = 528000 \text{ b/s} = 528 \text{ Kib/s} \quad (1)$$

A single audio stream comes out to 528 Kib/S. With the initial 100 KiB memory stream assumption, we get:

$$\frac{528 \text{ Kib/s}}{100 \text{ Kib stream}} = 5.28 \text{ refreshes/s} < 160 \text{ mhz} \quad (2)$$

That means we need to fully refresh and clear the memory 5 times over each second. This should be acceptable, given the SRAM's posted speed of 160mhz.

3 Cost and Schedule

3.1 Cost Analysis

UIUC Graduates in Electrical Engineering make \$80,000 per year on average. Prorating to a 2000 hour work year, that comes out to \$40 an hour. Based on this, and assuming that we work roughly 8 hours a week on this project the whole 16 week semester,

$$\text{Labor Cost} = \$40 * 8 \text{ hrs} * 16 \text{ weeks} * 3 \text{ people} * 2.5 = \$38,400 \quad (3)$$

Table 6: Price Estimates for Selected Components

Description	Manufacturer	Part No.	Quantity	Cost
ESP32 microcontroller with Wifi	Espressif Systems	ESP32-U4WDH	1	\$2.07
Power only USB-C connector	Kycon inc.	KUSBX-SMT-CS-6-BTR	1	\$0.91
Wire terminated low power speaker	Soberton inc.	SP-1304	1	\$1.75
Analog microphone	PUI Audio inc	AOM-5024P-HD-MB-R	1	\$3.83
320x240 RGB SPI Screen	Adafruit industries LLC	4311	1	\$19.50
USB-A Connector	Switchcraft inc.	RAHPCUA30E	1	\$1.50

Parts list ignores minor components such as resistors, capacitors, etc.

$$\text{Part cost} = \sum \text{ind. parts} = \$29.56 \quad (4)$$

$$\text{Total Costs} = \text{Labor Cost} + \text{Part Cost} = \$38,400 + \$29.56 = \$38,429.56 \quad (5)$$

3.2 Schedule

Table 7: Planned Timetable for Project Design Cycle

Week	Deadlines	Tasks		
		Brahmteq	Leo	Andrew
2/13		Prepare Design Document and Team Contract		
2/20	Design Document, Team Contract	Text-to-Speech Model	Finish Design Document	PCB Design, Solidworks CAD, Finish Design Document
2/27	Design Review	Speech-To-Text Model	Finalize PCB	Finalize PCB to be ordered
3/06	First Round PCB Orders, Teamwork Evaluation 1	Train Speech-To-Text Model on Wake Word ChatGPT API	Data Stream	ChatGPT API Call, Order components for PCB
3/13		Program text-to-SPI and screen functionality	Data Stream	Soldering/ testing PCB
3/20		Establish TCP routing between PC and ESP32	Configure Microphone/ Speaker	Glue code week-connect different working parts
3/27	Second Round PCB Orders	Assemble full system		
4/03		Testing/Verification		
4/10	Team Contract Fulfillment	Prepare Final Demo, Final Presentation and Final Paper		

4 Ethics and Safety

When AI solutions are being developed, ethical and safety concerns abound. Based on IEEE's Code of Ethics, there are three major ethical concerns that should be addressed with this design. The first, violating point I.2, is that AIs like ChatGPT are capable of misinformation, despite most users considering them infallible. To address this, we will be adding a disclaimer to the attached screen that displays on startup, informing the user of the limitations that an AI presents. That disclaimer will be something to the effect of "The responses provided by this machine are created automatically using an AI engine. These responses may be misleading or untrue. Do not rely on this device for medical or vital information." This disclaimer also addresses the second ethical issue, which is that some users may rely on the bot for life saving information, which it cannot guarantee is factual or helpful (which violates IEEE Code of Ethics point I.1). Hopefully by informing users of the limitations beforehand, these issues will be minimized.

One final ethical concern presented by this project is the issue of collecting user data, which also touches on IEEE's Code of Ethics point I.2. In order for this project to work, a microphone must be constantly on and recording the room around it, presenting a safety risk to users were that data to be leaked or tapped. While many companies use this data to reinforce their own learning models, we will be discarding it as soon as humanly possible to protect our users.

5 Citations

1. *OpenAI API*. [Online]. Available: <https://platform.openai.com/docs/models/gpt-3>. [Accessed: 18-Feb-2023].
2. "Specification for LCD module module no: AFK320240A0-3.5n12nth revision ..." [Online]. Available: <https://www.orientdisplay.com/wp-content/uploads/2021/10/AFK320240A0-3.5N12NTH.pdf>. [Accessed: 18-Feb-2023].
3. "Esp32-C6 series - espressif.com." [Online]. Available: https://www.espressif.com/sites/default/files/documentation/esp32-c6_datasheet_en.pdf. [Accessed: 18-Feb-2023].
4. Overestimate based on: "30mm PCB mount 8 ohm 82DB 700-5khz speaker - jameco electronics." [Online]. Available: https://www.jameco.com/z/36MS30008LF-PN-Jameco-Valuepro-30mm-PCB-Mount-8-Ohm-82dB-700-5kHz-Speaker_135722.html. [Accessed: 18-Feb-2023].
5. "Low-power, quad, 10-bit voltage-output DAC with serial interface." [Online]. Available: <https://www.mouser.com/datasheet/2/256/MAX5250-1292624.pdf>. [Accessed: 18-Feb-2023].