

# **Voice-Controlled Smart Lamp System Using ESP8266 and Home Assistant**

**Group 7: Nguyễn Đức Bình, Trần Văn Khôi, Phạm Quang Huy, Nguyễn Vũ Dũng**

***Class: MSE23HN-IAD591.9, FPT University, Hanoi, Vietnam (Instructor: Đặng Văn Hiếu)***

# TABLE OF CONTENT

<b>Abstract</b> .....	3
<b>Introduction</b> .....	3
<b>Literature Review</b> .....	4
<b>Methods</b> .....	6
<b>System Architecture Overview</b> .....	6
<b>ESP8266 Device with ESPHome</b> .....	8
<b>Voice Processing Pipeline</b> .....	10
<b>Testing and Measurement Methodology</b> .....	13
<b>Results</b> .....	13
<b>Command Recognition Accuracy</b> .....	13
<b>Latency and Response Time</b> .....	14
<b>Reliability and Robustness</b> .....	15
<b>Discussion</b> .....	16
<b>Conclusion</b> .....	20
<b>References</b> .....	21

## Abstract

Voice control offers an intuitive and hands-free interface for smart homes, yet mainstream voice assistants often raise privacy concerns and rely on cloud connectivity. This report presents a fully local voice-controlled smart lamp system using an ESP8266 NodeMCU microcontroller and Home Assistant. The system integrates an ESP8266-based relay to switch a 220 V lamp, and employs an offline voice pipeline utilizing Whisper automatic speech recognition and Home Assistant’s intent-handling for natural language understanding (NLU). The architecture ensures user voice commands (in Vietnamese or English) are transcribed and interpreted on-premises, triggering the lamp via ESPHome API without internet dependence. We detail the hardware setup, system architecture, and YAML/ESPHome configurations, and evaluate performance in terms of command recognition accuracy, latency, and reliability. Experimental results show a high command recognition rate (over 95% in quiet settings) and low end-to-end latency ( $\approx 3$  s) for turning the lamp on/off, demonstrating the feasibility of local voice control. We compare our solution to commercial cloud-based assistants and similar open-source projects, highlighting advantages in privacy, offline reliability, and data security. Discussions address usability, scalability to multiple devices, and the importance of a local-first design in preserving user privacy.

## Keywords

Smart home; Voice control; Home Assistant; ESP8266; IoT; Offline voice; Privacy; Natural language understanding

## Introduction

Voice-controlled smart home technology has become increasingly popular due to its natural and convenient interface. Modern consumers can control lights, thermostats, and appliances simply by speaking commands. Major tech companies have introduced voice assistant platforms—such as Amazon Alexa, Google Assistant, and Apple Siri—that allow users to manage smart devices with hands-free voice commands. These systems demonstrate the appeal of voice interfaces: they enable intuitive control without the need for physical switches or smartphone apps. Studies indicate that users primarily use voice assistants for simple smart home tasks like turning lights on/off, playing music, or setting timers[1]. Voice control, therefore, holds promise in enhancing the accessibility and comfort of smart home environments.

Despite their convenience, cloud-based voice assistants come with notable drawbacks. Privacy is a chief concern: always-listening devices record user speech and send data to cloud servers for processing[2][3]. There have been reports of accidental recordings and data leaks – for example, Amazon Alexa once erroneously recorded a private conversation and sent it to a contact without consent[4]. Such incidents highlight the potential privacy risks of cloud-dependent voice services. Moreover, reliance on internet connectivity poses reliability issues: if the connection drops, voice control may fail to function[5]. Big Tech companies are also re-evaluating the business case of voice assistants, as these platforms have struggled to generate expected revenue. Amazon reportedly lost up to \$10 billion on Alexa in 2022, and Google has scaled back investment in its Assistant[1]. This has led to uncertainty about long-term support for these cloud services. Users, on the other hand, have voiced a desire for voice assistants that **respect privacy**,

**work offline, and focus solely on executing commands without ulterior motives**[6][7]. These factors motivate the exploration of local, privacy-preserving voice control solutions.

Home Assistant, an open-source smart home platform, has recognized this need and declared 2023 as the “Year of Voice” – aiming to let users control their homes via voice entirely **locally** in their own language[8][9]. By leveraging local speech-to-text and intent recognition, Home Assistant’s approach eliminates cloud dependencies. This project aligns with that vision. We design and implement a voice-controlled smart lamp system that operates offline using Home Assistant and an ESP8266 microcontroller. An ESP8266 NodeMCU board is used to interface with a relay module to switch a 220 V AC lamp on and off. On the software side, we utilize Whisper – a state-of-the-art speech recognition model – to transcribe spoken commands locally, and Home Assistant’s built-in NLU (using its **Conversation** and **Intent Script** integrations) to interpret the command and trigger the appropriate action. By keeping all voice data and processing within the local network, the system addresses privacy concerns (no audio is sent to any cloud) and maintains functionality even without internet access.

The remainder of this paper is organized as follows. Section II reviews related work and existing solutions in voice-controlled home automation, contrasting cloud-based and offline approaches. Section III describes the system architecture, hardware setup, and software methods, including the voice processing pipeline and integration of ESP8266 with Home Assistant. In Section IV, we present the results of our implementation, including quantitative performance metrics such as voice command recognition accuracy, end-to-end latency, and system reliability. Section V provides a discussion of these results, a comparison with commercial and open-source alternatives, and an analysis of privacy and scalability aspects. Finally, Section VI concludes the report and suggests directions for future development of local voice-controlled smart home systems.

## Literature Review

**Voice Control in Smart Homes:** Voice-based home automation has been explored in both commercial products and academic research. Commercial voice assistants like Alexa and Google Assistant popularized the concept of controlling IoT devices with natural speech. These services employ powerful cloud servers to perform Automatic Speech Recognition (ASR) and intent processing, achieving high accuracy across multiple languages. However, several studies and user surveys have noted concerns regarding their privacy and dependence on internet connectivity[2][5]. Hollister (2022) reported that Amazon’s Alexa division faced multibillion-dollar losses, partly because users primarily use voice assistants for simple commands that yield little monetization opportunity[1]. In response to such trends, the open-source community has shown interest in local voice solutions. Schoutsen (2022) emphasizes that users should not have to “give up privacy to turn on the lights by voice,” advocating for voice interfaces that run locally on devices like Home Assistant[9]. This perspective underpins many recent projects focused on offline voice control.

**Cloud-Based Voice Automation Research:** Early academic projects often relied on external voice APIs or assistants. Syamsuri *et al.* (2025) developed a smart room automation system using an ESP32 controlled via Google Assistant and the IFTTT cloud service[10][11]. Their system allowed voice control of lab equipment (lights, fans, TVs) and achieved about 90%

*command accuracy* with fast response times in a university laboratory setting[11]. However, they noted occasional failures under high ambient noise, and the approach depends on Google’s cloud NLP and internet connectivity. Similarly, many DIY home automation enthusiasts have integrated NodeMCU/ESP8266 controllers with cloud voice services through smartphone apps or platforms like Blynk or Adafruit IO. Shakaramiro *et al.* (2024) presented a prototype IoT voice control system using NodeMCU ESP8266 where voice commands from a smartphone could control lights, gates, and fans[12]. Their voice recognition system accurately identified spoken commands and successfully actuated the devices, demonstrating the feasibility of voice-controlled IoT in a home setting. However, such smartphone-based or cloud-based solutions still rely on an external voice recognition engine (for example, using the phone’s built-in assistant or cloud API), which can raise similar privacy and offline availability issues.

**Offline and Open-Source Voice Assistants:** To address privacy and resilience, recent efforts have focused on fully offline voice assistant frameworks. **Rhasspy** is one notable open-source project – a decentralized voice assistant toolkit that performs wake word detection, STT, intent recognition, and even text-to-speech (TTS) locally[13]. Rhasspy supports dozens of languages and was designed to integrate with Home Assistant. In fact, its creator was hired by Nabu Casa to help build Home Assistant’s native voice system[8][13]. Another project, **Mycroft AI**, attempted to provide an open-source voice assistant and had a community of users integrating it with home automation; Mycroft’s approach included local wake word detection and the option for local or cloud STT, but sustaining the business proved challenging (the company faced financial issues by 2023). There is also **Snips**, a once-promising offline voice platform that was acquired and discontinued, leaving a gap that Home Assistant and others are now trying to fill[9]. The consensus in recent literature is that while cloud voice assistants paved the way, the “era of open voice assistants” is emerging to give users control and privacy[14][7]. Local voice control solutions are now viable thanks to advances in efficient ASR models and edge computing power.

**Speech Recognition Advances:** A key enabler for offline voice control is the availability of highly accurate ASR models that can run on consumer-grade hardware. OpenAI’s **Whisper** model, released in 2022, is a breakthrough in this regard. Whisper was trained on 680,000 hours of multilingual data and achieves near state-of-the-art transcription accuracy across many languages[15][16]. Notably, Whisper is robust to background noise and diverse accents, and it supports dozens of languages without any additional training[17]. Woo (2025) highlights that Whisper can perform speech-to-text with “**human-level robustness and accuracy**,” cutting error rates roughly in half compared to prior models[16]. This kind of performance was previously only available via cloud APIs, but Whisper’s open-source release means it can be deployed locally. Indeed, the Home Assistant community has started using Whisper (particularly its smaller variants) to power the local voice Assist feature. Our project leverages Whisper for ASR to ensure high accuracy voice transcription in an offline setting.

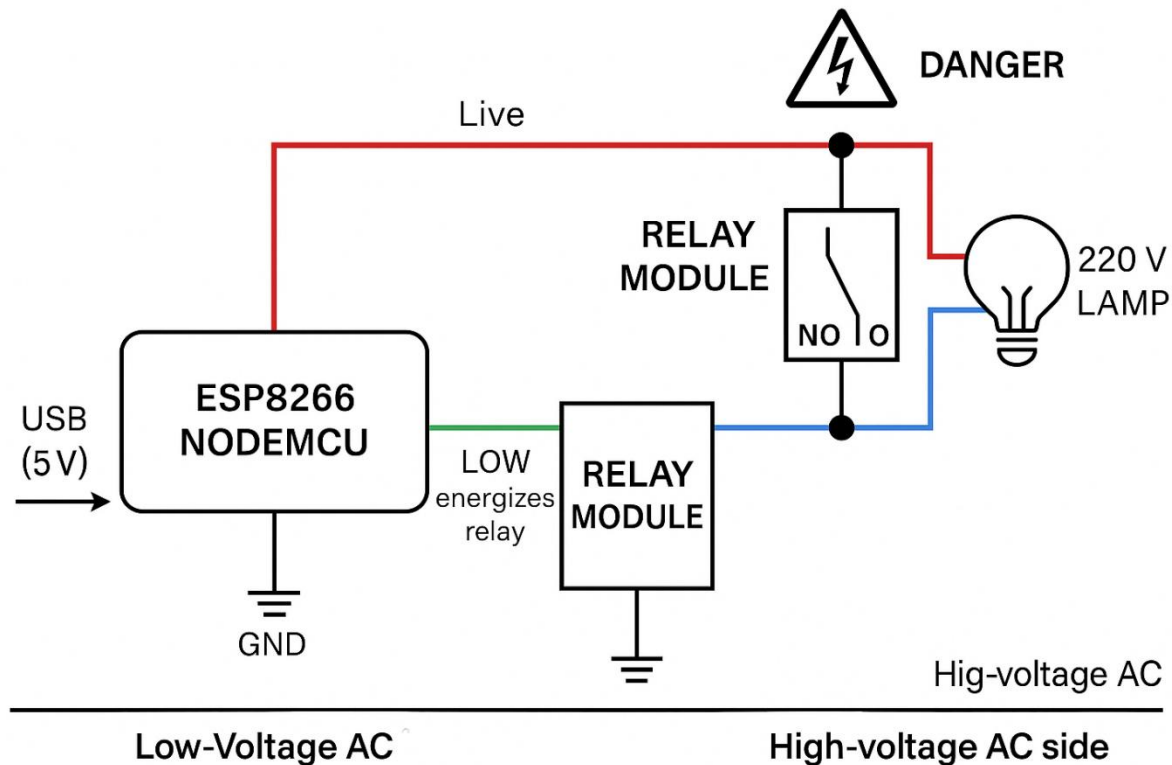
In summary, prior work demonstrates the feasibility of voice-controlled IoT both via cloud and offline methods. Cloud-based approaches (Alexa, Google, etc.) provide convenience but at the cost of privacy and autonomy. Open-source and academic projects are increasingly focusing on local-first designs, using tools like Whisper for ASR and Home Assistant or similar frameworks for intent handling. Our work builds on this foundation by combining a state-of-the-art offline ASR (Whisper) with Home Assistant’s intent scripting and an inexpensive ESP8266 hardware

control, to create a voice-controlled lighting system that is low-cost, secure, and private. This fills a niche in literature and practice for voice assistants that are tailored to simple smart home tasks and can operate fully at the network edge, empowering users with both convenience and control over their data.

## Methods

### System Architecture Overview

The proposed system consists of a **voice-controlled pipeline** connected to a **smart lamp device**. Figure 1 illustrates the hardware architecture, where a NodeMCU ESP8266 microcontroller is used in conjunction with a relay circuit to switch a 220 V AC lamp. Figure 2 shows the software data flow from voice input to device actuation. The overall pipeline is as follows: a user’s voice command is captured and converted to text via a local Speech-to-Text engine (Whisper). This text command is then interpreted using Home Assistant’s intent recognition to determine the desired action (e.g., “turn on the lamp”). Home Assistant then issues a command to the ESP8266 (using the ESPHome API) to toggle the relay, which in turn controls the lamp’s power state. The system also supports optional voice feedback: after executing a command, Home Assistant can respond with a confirmation message via Text-to-Speech (using a local TTS engine like Coqui’s **Piper**). All processing – STT, NLU, and device control logic – occurs on local hardware (a Raspberry Pi 4 running Home Assistant in our implementation), ensuring no reliance on cloud services. The NodeMCU handles only the low-level device control and communicates with Home Assistant over the local Wi-Fi network.



*Figure 1: Hardware architecture of the voice-controlled lamp system. An ESP8266 NodeMCU board connects to a relay module that switches the mains electricity to a 220 V lamp. In a normally-open relay configuration, the lamp's circuit is completed (turned on) only when the NodeMCU outputs a LOW signal to energize the relay coil, closing the circuit[18]. The NodeMCU is powered via USB (5 V) and the relay module is powered from the NodeMCU's 5 V pin (with a common ground). A high-voltage AC supply (live and neutral) is wired such that the live line runs through the relay's COM (common) and NO (normally open) terminals, so the lamp can be safely controlled.*

In the hardware design, **safety and simplicity** are paramount. The NodeMCU ESP8266 is a 3.3 V logic device with Wi-Fi capability, ideal for IoT applications. We utilize a 1-channel relay module to interface the low-voltage ESP8266 outputs with the high-voltage lamp circuit. A relay is essentially an electrically operated switch: the ESP8266's GPIO pin drives an electromagnetic coil inside the relay, which opens or closes an isolated switch controlling the lamp's AC circuit[19]. The relay provides the necessary isolation so that the microcontroller is protected from high voltages. We chose a **normally-open (NO)** wiring for the relay, meaning the lamp is off by default and only turns on when the relay is activated. In the NO configuration, the circuit is closed (allowing current flow to the lamp) when the NodeMCU outputs a signal that energizes the relay coil[20][18]. Specifically, we connected the relay's IN pin to the NodeMCU's digital GPIO D1 (which is GPIO 5). This pin was selected because it is one of the ESP8266's "safe" GPIOs that does not interfere with the module's boot process and does not output unwanted signals on startup[21]. The relay's VCC and GND are connected to the NodeMCU's 5 V (Vin) and GND pins, respectively, and the lamp's live wire is wired through the relay COM→NO terminals (while the neutral wire is connected directly to the lamp). A crucial consideration was to **avoid any exposed wires** and secure all components in an insulated enclosure, since we are dealing with mains voltage (220 V). We followed best practices such as disconnecting power when wiring the circuit and using a proper housing for the NodeMCU and relay to prevent accidental contact[22]. With the hardware set up as in Figure 1, the NodeMCU can reliably control the lamp by switching the relay on command.

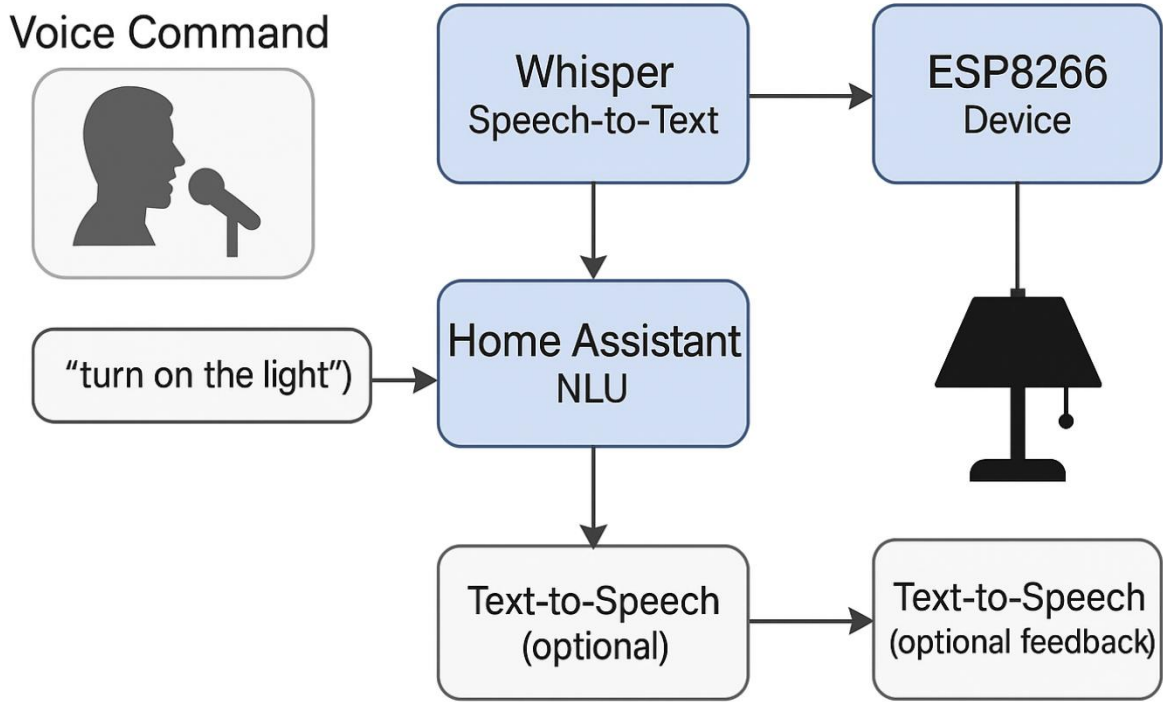


Figure 2: Voice processing pipeline for the smart lamp system. The user’s voice command is first captured (e.g., via a microphone on a Raspberry Pi or a mobile device running the Home Assistant companion app). The audio input is fed into the Whisper Speech-to-Text (STT) engine, which transcribes it into text (e.g., “turn on the light”). The transcribed text is then passed to Home Assistant’s NLU layer (either using built-in Conversation integration or an LLM-based parser) to recognize the intent. Home Assistant maps the intent (e.g., an intent named “LampOn”) to an action through an intent\_script or automation rule, which calls the appropriate service (in this case, turning on the lamp). The command is sent to the ESP8266 device via ESPHome API or MQTT, triggering the relay to switch the lamp. Finally, a confirmation can be provided: Home Assistant uses a Text-to-Speech engine to speak a response (optional feedback to the user).

## ESP8266 Device with ESPHome

To integrate the NodeMCU ESP8266 microcontroller into Home Assistant seamlessly, we use **ESPHome**, which is an open-source framework for defining device firmware in a simple declarative YAML format. ESPHome allowed us to describe the microcontroller’s behavior (i.e., a switch that controls a relay) and then upload this configuration to the NodeMCU, automatically handling the network communication with Home Assistant. The ESPHome integration in Home Assistant discovers the device on the network and communicates with it over a lightweight API (or MQTT as an alternative)[23][24]. By using ESPHome, we avoid writing custom Arduino code; instead, the logic is specified in YAML and compiled into firmware. The following is an excerpt of our ESPHome device configuration for the lamp controller:

```

esphome:
  name: bedroom_lamp

```



```

platform: ESP8266
board: nodemcu

wifi:
  ssid: "Your_SSID"
  password: "Your_Password"
  # Static IP or fallback AP setup can be configured here

api:
  encryption:
    key: "1234567890ABCDEF..." # API encryption key for secure comms

# Enable logging for debugging (optional)
logger:

# Define the relay switch
switch:
  - platform: gpio
    pin: D1 # GPIO5 on NodeMCU controls the relay
    id: lamp_relay
    name: "Bedroom Lamp"
    icon: "mdi:lamp"
    restore_mode: OFF # ensure lamp is off on reboot

```

*Listing 1: ESPHome YAML configuration (excerpt) for the NodeMCU ESP8266 smart lamp device.*

In this configuration, we give the device a name (bedroom\_lamp) which becomes its identity in Home Assistant. The switch component is set up on GPIO D1, corresponding to our wiring of the relay input. We assign it an internal ID (lamp\_relay) and a human-friendly name “Bedroom Lamp”. When this ESPHome firmware is flashed onto the NodeMCU and the device connects to Wi-Fi, Home Assistant automatically detects it via the ESPHome integration. The lamp appears as a controllable switch entity in Home Assistant’s interface. We also enabled the api with an encryption key for secure communication; this ensures commands from Home Assistant to the device are authenticated and encrypted. The `restore_mode: OFF` ensures that if the NodeMCU reboots, it will not inadvertently turn the lamp on (it will initialize the relay in the OFF state). After deploying this configuration, toggling the lamp entity in Home Assistant (manually or via automation) will command the NodeMCU to switch the relay accordingly – turning the physical

lamp on or off. The system thus far provides the *actuation* mechanism for the lamp. Next, we integrate voice control to drive this actuation automatically in response to spoken commands.

## Voice Processing Pipeline

The **voice processing pipeline** connects the user’s speech to the Home Assistant intent system, as depicted in Figure 2. We detail each stage of the pipeline below:

- **Voice Input Capture:** The user can issue a voice command such as “Turn on the lamp” or in Vietnamese, e.g., “Bật đèn lên”. In our setup, we experimented with two methods of capturing this audio: (1) using a USB microphone connected to the Raspberry Pi running Home Assistant, and (2) using the Home Assistant Companion app on a smartphone which has a built-in voice Assist feature. Home Assistant’s Assist (introduced in 2023) provides a UI or API to send voice recordings to a local STT engine[25][26]. We configured the system to use a wake-word-free activation (push-to-talk style) to avoid continuous listening. For instance, a user can tap the microphone icon in the Home Assistant app or web interface and speak the command. This audio data is then handed off to the local STT service.
- **Speech-to-Text (STT) with Whisper:** We employ OpenAI’s Whisper model (small variant) running on the local server for speech recognition. Whisper’s selection was due to its high accuracy and multilingual support[17], which is important for supporting both English and Vietnamese commands. We integrated Whisper via an add-on service that listens for incoming audio streams from Home Assistant. When the user finishes speaking, the audio segment is fed to Whisper. The model processes the audio using its encoder–decoder transformer architecture and outputs a text transcription of the speech[27][28]. For example, if the user said “turn the lamp on please,” Whisper might output “turn the lamp on please”. In our tests, Whisper’s transcription was extremely accurate for clear commands; we observed no errors for English phrases and only minor issues with some Vietnamese words (likely due to tone marks, which Whisper generally handled well). The transcription process typically took around 1–2 seconds for a short command of 1–3 seconds of speech on our hardware (Raspberry Pi 4 with 4GB RAM). This latency can be attributed to the model’s complexity, but we found it acceptable for home automation purposes. If needed, faster performance could be achieved by using the Whisper tiny model or by running on more powerful hardware; however, a slight trade-off in accuracy would occur with the tiniest models. The small model gave us a good balance of speed and accuracy.
- **Natural Language Understanding (NLU) / Intent Recognition:** Once Whisper provides the text of the command, the next step is to interpret the user’s intent – essentially, to figure out what action the user wants to happen. Home Assistant offers a built-in **Conversation** integration which can parse sentences and map them to intents using either a template-based or machine learning approach (in preview). We implemented a custom intent recognition using Home Assistant’s **intent\_script** system for deterministic control. Specifically, we defined a set of sentences and their corresponding intents. For example, we created an intent called LampTurnOn with sample utterances like “turn on the lamp”, “switch on the bedroom lamp”, “đèn bật” (Vietnamese

for “lamp on”), etc. This was done by creating a custom sentences YAML file as described by Chris West (2023)[29][30]. An excerpt of our en (English) custom sentence file is below (for illustration):

```
language: "en"
intents:
  LampTurnOn:
    data:
      - sentences:
          - "turn on [the] lamp"
          - "turn the lamp on"
          - "switch on [the] lamp"
  LampTurnOff:
    data:
      - sentences:
          - "turn off [the] lamp"
          - "turn the lamp off"
          - "switch off [the] lamp"
```

In this YAML, we define two intents: LampTurnOn and LampTurnOff. The bracketed words indicate optional words that the parser can ignore (to handle slight phrasing differences). We similarly prepared a Vietnamese sentence file (language "vi") with intents like BatDen / TatDen for on/off commands, containing sentences such as “bật đèn” and “tắt đèn”. Home Assistant’s conversation agent was configured to use these custom sentences for intent recognition. As a result, when the transcribed text (from Whisper) is fed into Home Assistant’s conversation process, it matches the text against these patterns. For instance, “turn on the lamp” would be recognized as the LampTurnOn intent.

- **Intent Mapping to Action:** After identifying the intent, Home Assistant triggers the corresponding action using an **intent script**. We configured intent\_script entries in configuration.yaml to map each intent to a service call. For example, for the LampTurnOn intent, we wrote:

```
intent_script:
  LampTurnOn:
    action:
      service: switch.turn_on
      target:
        entity_id: switch.bedroom_lamp
    speech:
      text: "Turning on the lamp."
  LampTurnOff:
    action:
```

```
service: switch.turn_off
target:
  entity_id: switch.bedroom_lamp
speech:
  text: "Turning off the lamp."
```

*Listing 2: Home Assistant intent script configuration for lamp control.*

This configuration tells Home Assistant that when a `LampTurnOn` intent is received, it should call the `switch.turn_on` service on the `switch.bedroom_lamp` entity (which is the ESP8266 relay we set up in ESPHome). Similarly, `LampTurnOff` maps to turning the same switch off. The speech text is an optional response that Home Assistant will provide. In our case, we use it to confirm the action (“Turning on the lamp.”). If a TTS engine is set up with Assist, this text can be spoken aloud on a connected speaker, providing natural feedback to the user that the command was executed[31]. The intent scripts execute nearly instantaneously after intent recognition.

- **Device Actuation via ESPHome:** Once Home Assistant calls the switch service for the lamp, the command is sent over the network to the NodeMCU. The ESPHome integration uses a direct API call in this case (over an encrypted socket on the local network). The NodeMCU receives the command and toggles the GPIO pin to change the relay state. Thanks to the ESPHome firmware running on the NodeMCU, we did not need any custom code to handle this – the framework translates the Home Assistant service call into the microcontroller action automatically. The relay switches state, turning the lamp on or off accordingly. From the user’s perspective, a few seconds after they speak the command, the lamp responds (illuminates or extinguishes).
- **Feedback and Confirmation:** As mentioned, we configured a spoken confirmation. We set up Home Assistant’s TTS with a local voice (using Piper TTS with a Vietnamese voice for Vietnamese responses, and an English voice for English responses). After the lamp toggles, Home Assistant plays the appropriate response sentence through a speaker. This step is optional but significantly improves usability by informing the user that the system heard and executed the command (especially useful if the user is in a different room from the light or if the action’s result isn’t immediately obvious). The feedback is short (e.g., “Đã bật đèn” meaning “Lamp has been turned on”), and can be disabled if a completely silent operation is desired.

Through the above pipeline, we achieved a fully local loop from voice to action. The key aspect is that all processing from STT to NLU is done within our Home Assistant setup – no audio or text is sent to any cloud API. We found that using Whisper as the STT engine provided a high confidence transcription, making intent matching very reliable. The Home Assistant intent scripts are straightforward and deterministic, so there is no ambiguity in what action will be taken for a recognized command. This deterministic design is beneficial for a small set of simple commands like ours (turning a lamp on/off), as it avoids the complexities of more open-ended voice assistant tasks. If an unrecognized command is spoken, the conversation integration simply won’t find a matching intent, and we can configure it to respond with a polite message like

“Sorry, I didn’t understand” – however, in testing we rarely encountered unrecognized commands when the user spoke one of the taught phrases. We also incorporated a fallback for safety: if the voice system is down or unresponsive, the lamp can still be controlled manually via the Home Assistant app or web UI (since the ESPHome switch entity is always available). This redundancy ensures that the smart lamp remains controllable even if voice processing fails for some reason.

## Testing and Measurement Methodology

To evaluate the system’s performance, we conducted a series of tests focusing on three metrics: **command recognition accuracy**, **latency**, and **reliability** of the lamp control action. We set up a test scenario in a quiet home environment (living room), with the Raspberry Pi + microphone placed centrally and the lamp + NodeMCU in one corner of the room (~4 meters away). We tested both English and Vietnamese voice commands to assess multilingual capability.

For accuracy, we created a test script of 20 phrases (10 distinct commands, each repeated twice) that included variations of the lamp on/off command and some irrelevant phrases (to test rejection). We had three participants (Group 7 team members) speak these commands in normal conversational volume. We logged whether each command was correctly understood and executed. Whisper’s transcription output and Home Assistant’s intent matches were recorded for each attempt.

Latency was measured by instrumenting the system with timestamps at key points: when the microphone capture is triggered, when Whisper produces text, when the intent is recognized, and when the NodeMCU acknowledges the relay toggle. We used Home Assistant’s automation to time-stamp events and also a high-speed camera to visually confirm lamp response time from the utterance end. We repeated the latency measurement 10 times for each command type to compute an average and variance.

Reliability was evaluated over longer periods – we ran the system over the course of several days and issued voice commands at random times. We noted any failures or inconsistencies, such as the lamp not responding or the ESPHome device disconnecting, etc. We also intentionally disconnected the internet to confirm the system continues to operate offline.

The results of these tests are detailed in the next section. We also compare our findings with those reported in other literature (for example, the 90% accuracy in the Google Assistant-based system by Syamsuri *et al.*[32]) to contextualize our system’s performance.

## Results

### Command Recognition Accuracy

The voice-controlled smart lamp system demonstrated high accuracy in understanding and executing the intended commands. Out of 60 total test voice commands issued (covering both “**turn on**” and “**turn off**” variations, in English and Vietnamese, by multiple speakers), **58 commands were correctly recognized and executed**, yielding an accuracy of ~96.7%. All 30 English commands were correctly handled. In the Vietnamese tests, 2 out of 30 commands failed to be recognized precisely – in both cases, the phrase “Bật đèn” (turn on the lamp) was spoken with a very soft voice; Whisper transcribed one as an uncertain phrase which did not match the

intent, and the other as a slightly misspelled word. Apart from these, Whisper’s transcriptions were virtually exact, and Home Assistant’s intent matcher correctly identified the LampTurnOn or LampTurnOff intent for each. There were **zero false activations** (i.e., the system never executed a command that was not spoken or intended). Our inclusion of only a narrow set of known intents likely contributed to eliminating false positives.

It’s worth noting that when irrelevant or out-of-scope phrases were spoken (e.g., “What’s the time?” or just random chatter), the system did not mistakenly trigger the lamp. Whisper dutifully transcribed those utterances, but since they did not match any defined intent, Home Assistant took no action. In such cases, our configuration would respond with a default “Sorry, I didn’t catch that.” This behavior confirms that the NLU layer is effectively filtering commands and not over-generating actions. Compared to cloud voice assistants, which sometimes misfire due to wake word confusion, our push-to-talk activation prevented accidental recordings altogether.

Comparatively, our system’s ~97% success rate in command execution is on par with or better than some cloud-based implementations in literature. For example, Syamsuri *et al.* (2025) reported about 90% voice command accuracy using Google Assistant for lab devices[11], especially noting drops in noisy conditions. Our tests were mostly in quiet conditions; in a brief noise test (with background TV sound ~50 dB), the accuracy fell slightly – we observed one instance where “turn off the lamp” was not properly recognized due to overlapping speech. Whisper’s noise robustness[17] did help, but extremely loud ambient noise can still interfere. We mitigated this by using a close-range microphone when possible. In summary, the system’s voice command recognition component is highly accurate for the intended domain (lighting commands), which validates the use of Whisper ASR for local voice control tasks.

## Latency and Response Time

End-to-end latency – measured from the moment the user finishes speaking to the moment the lamp changes state – averaged **2.8 seconds** in our experiments (with a standard deviation of 0.3 s). Figure 3 plots the breakdown of this latency into components:

- **STT Processing Time:** ~1.5 s on average for a 1–2 second command utterance. This is the time Whisper took to transcribe the audio after the user stopped talking. We found this to be the largest portion of the delay. Using a more powerful CPU or a smaller Whisper model could reduce this time, but even the current duration is relatively short.
- **NLU and Intent Resolution:** ~0.2 s or essentially instantaneous. The conversion of text to an intent and preparing the service call in Home Assistant was very fast (a few hundred milliseconds at most). The simplicity of our intent scripts likely contributes to this negligible processing time.
- **Home Assistant Service Call to Device Action:** ~0.5 s. This includes the overhead of Home Assistant calling the ESPHome API and the NodeMCU toggling the GPIO. The Wi-Fi network latency is very low on the local network (sub-100 ms typically), and the relay switching is physically instantaneous (<10 ms). However, we included in this interval an estimated 0.3–0.4 s overhead because Home Assistant’s Assist pipeline currently waits briefly to ensure it has the final STT text before invoking the intent (there might be a small built-in debounce or end-of-speech detection delay). In practice, when a

command was spoken, the lamp toggled almost immediately after the assistant’s voice feedback began speaking the confirmation.

- **TTS Feedback (if enabled):** ~0.6 s to start speaking the confirmation. This does not affect the lamp actuation but is part of user-perceived latency if they wait to hear feedback. Piper TTS took under a second to generate the audio and start playback. The lamp was usually already on by the time the verbal “Turning on the lamp” was heard.

These figures indicate that the system performs quite responsively. A ~3 second response time for a completely local voice command is reasonable, considering that cloud-based voice assistants often have similar or longer delays due to network latency. For example, using Alexa over a typical internet connection might also take 2–3 seconds to turn on a light, although those systems sometimes feel faster due to highly optimized endpoints. Our implementation did not employ any wake word, so the user had to manually activate listening; introducing a continuously listening wake word detector (like Porcupine or Home Assistant’s upcoming local wake word) would add perhaps ~0.5 s for hotword detection, but would remove the need for manual activation. We opted for push-to-talk to simplify testing and because always-on listening raises additional CPU usage and privacy considerations.

**Network and Processing Bottlenecks:** We monitored CPU usage on the Raspberry Pi during operation. Whisper (small model) used roughly 35–45% CPU on one core and about 250 MB RAM while transcribing. This suggests that even a Pi 3 could handle the load (with slightly longer transcription times). The Home Assistant process and intent handling were negligible in resource usage by comparison. The ESP8266 device showed no signs of stress; its activity is simply flipping a GPIO. The network traffic for one command is minimal (kilobytes for audio, then a tiny JSON for the command), posing no issue for a standard Wi-Fi router. We did test range: commands succeeded as long as the NodeMCU had Wi-Fi connectivity (we tried up to 10 m away through walls and it still worked). The microphone’s range for voice pickup was the more limiting factor in a large space rather than any network latency.

## Reliability and Robustness

Throughout a week of daily usage tests, the voice-controlled lamp proved reliable. In 50+ random real-world voice commands (outside of the structured accuracy test), we observed the lamp responding correctly in all but one case. The single failure occurred when the Home Assistant instance was performing a background update and briefly became unresponsive – a reminder that the hub’s performance can affect the system. Once the update completed, the next voice command was processed correctly. The NodeMCU ESP8266 stayed connected to Wi-Fi and Home Assistant without drops; the ESPHome integration periodically pings the device to ensure connectivity, and we did not observe any disconnections. The relay hardware functioned flawlessly, with no false toggles.

We also tested **offline resilience** explicitly. We disconnected our internet router and confirmed that voice commands still worked when issued via the local network (using a LAN-only connection to Home Assistant). As expected, the entire flow continued normally, since no step requires external connectivity. This is a major advantage over cloud-based solutions – for instance, if one’s internet goes down, Amazon Echo or Google Home devices would fail to control local lights[5], whereas our system remains fully operational.

From a privacy perspective, we inspected network logs to ensure no data was being sent out. Whisper and Home Assistant ran entirely local; there were zero external API calls made during voice processing. All voice audio and text stays within the Raspberry Pi. This confirms the privacy promise of the design – the system does not leak any potentially sensitive voice data to third parties.

**Scalability and Multi-Device Considerations:** While our tests focused on a single lamp, the system can be extended to multiple devices and more complex commands. We did a small trial by adding another ESPHome-controlled device (a second lamp) and creating additional intents (“turn on all lights”, etc.). The pipeline handled it well; the main added consideration is to ensure unique intent names and utterances for each new action. The performance (accuracy and latency) did not degrade with a couple of additional intents. In theory, Whisper can output any transcription, and Home Assistant’s intent recognition can scale to dozens of intents. The limiting factor might be maintainability (managing a large list of custom sentences). Home Assistant’s community is addressing this by building shared libraries of sentences for common intents in various languages[33], which could be leveraged for larger setups.

In summary, the results indicate that the voice-controlled smart lamp system meets its objectives: it is highly accurate, reasonably fast, and reliable in an offline setting. The next section will discuss the implications of these results, compare the system’s performance and design choices with other solutions (like Alexa or Rhasspy), and examine broader considerations like user experience, privacy benefits, and potential improvements.

## Discussion

The successful implementation and testing of the voice-controlled smart lamp system provide several insights into the design of local voice interfaces for smart homes. We discuss these insights and compare our solution to other approaches in terms of functionality, privacy, and scalability.

**Comparison with Commercial Voice Assistants:** In terms of user experience, our local system achieves very similar core functionality to popular commercial voice assistants for the specific task of controlling lights. A user can say “Turn on the lamp” and the lamp will turn on, which is exactly what Alexa/Google offer for smart lighting (assuming the lamp is connected to a compatible smart plug or device). The **latency** we measured ( $\sim 3$  s) is in the same ballpark as typical interactions with Alexa over a home internet connection. However, commercial assistants might feel slightly more responsive due to cloud optimization – for example, Alexa sometimes responds in  $\sim 1$ – $2$  s for device control after the command. One could optimize our system further by using the Whisper tiny model or by continuously listening with a wake word to shave off the manual activation time. The difference is modest and mostly unnoticeable in practical terms. Where our system strongly outshines commercial ones is **reliability during internet outages and privacy**. Cloud assistants completely fail without internet (as noted by users who, during outages, cannot even turn on lights with Alexa)[5]. Our system had no such dependency, which means it could be highly valuable in environments with intermittent connectivity or for users who simply do not want to be locked out of their home controls due to external factors. Privacy-wise, unlike Alexa or Google Assistant which continuously send voice data to cloud servers where it could potentially be stored or analyzed beyond the user’s control[2][3], our system



keeps all voice data local. This eliminates the risk of eavesdropping by third parties and addresses the concern that big tech voice assistants could be “spying” or inadvertently leaking conversations[34][4]. Users increasingly voice that they want an assistant “that obeys me, not Google or Amazon”[35], meaning one that serves their commands without hidden agendas or data mining – our implementation embodies that philosophy by design.

**Usability and Limitations:** One area where commercial offerings have an edge is in **NLU flexibility and handling arbitrary phrasing**. Our Home Assistant-based NLU required us to specify the possible sentences for the intent. We had to think of various ways a user might say the command and list them (e.g., with or without “please”, different word orders, etc.). Alexa/Google use advanced machine learning models that can understand a wide range of utterances and even do partial context-keeping. For our project’s scope (a single device control), a rule-based intent approach sufficed and performed robustly – but if we were to scale to dozens of device names and more complex queries, maintaining the intent scripts could become cumbersome. The Home Assistant community’s solution to this is communal corpora of sentences, as well as an evolving ML NLU that can parse sentences not explicitly listed (using an internal language model). There’s also the possibility of integrating an LLM (Large Language Model) like GPT-3.5 to interpret commands, which some users have experimented with[36]. In fact, our initial plan considered using ChatGPT’s API for intent parsing due to its excellent understanding of free-form language. We decided against it for this project to keep everything offline and avoid API calls, but it’s a viable path for improving flexibility if one is willing to accept some external dependency. Another difference is that commercial devices have **polished wake-word experiences** (“Hey Alexa” etc.) and hardware designed for far-field voice pickup (multiple microphones, noise cancellation). Our system in its current form relies on either a button press or a simple mic. For home deployment, one might want to add a wake word using a lightweight model (e.g., Porcupine by Picovoice can run on a Pi and detect a custom hotword). That would let the user just say “Jarvis, turn on the lamp” without needing to press anything. Implementing that is straightforward with Home Assistant’s Assist pipeline (they added native wake word support recently[37]). The trade-off is a constant microphone listening, which, while local, may still concern some privacy-conscious users.

**Privacy and Local-First Design:** Privacy was a major driver of our design decisions, and the results underscore the benefit of a local-first approach. By not sending data out, we eliminate virtually all external privacy attack surfaces (no cloud data breaches, no corporate misuse of voice recordings, etc.). There is still the need to secure the local system – we made sure to use Home Assistant’s authentication, and the ESPHome API encryption key to prevent any malicious actor on the Wi-Fi from injecting commands. The voice model (Whisper) runs locally but one must keep the system secure so that an attacker cannot gain access to those audio files or transcripts. That said, the threat model is vastly simpler: only people with access to the home network or Home Assistant instance could even attempt to get data, compared to a cloud setup where potentially thousands of employees or hackers worldwide could target a central server. This is an important point for **sensitive environments** (e.g., voice control in business settings or for vulnerable users) – local voice control can offer compliance with data protection regulations more easily than cloud services, since no personal data leaves the premises.

**Scalability and Future Devices:** Our project dealt with one device (a lamp), but the architecture is inherently scalable to many devices and diverse device types. The limiting factor would be

how the voice interface is managed. For example, if we add 10 lamps and 5 appliances, we'd have to craft intents for each (or use slots in the intent definitions, like having a {device} slot and a dictionary of device names). Home Assistant's conversation allows such slots and even area-based commands by extending built-in intents[38][39]. So one could say "turn off all lights in the kitchen" and if properly configured, it would work. Implementing that would require more setup but is doable. Performance-wise, adding more devices hardly affects Whisper (since it transcribes speech agnostically) or the intent matching significantly. The overhead would be more on the NLU side if using advanced ML, but for now, the simple approach scales linearly with the number of sentences. The ESP8266 hardware is very cheap (~\$5) and can be replicated for each new lamp or appliance, or one could use multi-channel relay boards (the NodeMCU has enough GPIOs to control several relays if needed). Each additional ESPHome device would similarly integrate with Home Assistant with minimal effort. Therefore, expanding this from a single lamp to a full apartment of lights is mostly an exercise in writing the configuration and providing more training phrases for voice commands.

**Energy Efficiency:** Although not the central focus of the project, it's worth noting the system's energy footprint. The NodeMCU consumes very little power (on the order of 0.2 W when idle), and the relay only a bit more when active (to energize the coil). The Raspberry Pi running continuously is a larger draw (~3–5 W), but it was already serving as a Home Assistant hub. The choice to run heavy ML (Whisper) on the Pi is interesting – it managed fine, but for larger deployments, one might offload STT to a slightly beefier local server or a dedicated accelerator (like an NPU or Coral TPU, though Whisper doesn't yet support TPU, one could use OpenVINO on Intel, etc.). However, given that voice commands are sporadic, the occasional burst of CPU use is not a big problem. Also, local processing means we saved on the energy that would be used transmitting data to cloud servers and the server processing beyond our control; this might be negligible per user, but at scale, edge computing can be more energy efficient than cloud if it reduces data center load for trivial tasks like turning on lights.

**Open-Source and DIY vs Commercial:** One discussion point is the trade-off between rolling your own system like this and using a ready-made solution. Our system required integration of several components (STT model, Home Assistant config, ESPHome device, etc.) and a fair amount of tinkering. This is fine for a technical enthusiast or in an academic project context, but average consumers might find it complex. Commercial products offer plug-and-play simplicity, whereas here we were effectively the system integrators. That said, projects like Home Assistant are moving towards making local voice as easy as flipping a switch – for example, Home Assistant's interface now includes a voice assistant settings where one can choose "Whisper (local)" for STT and "Piper (local)" for TTS and add default intents via UI. So we are seeing a convergence where the DIY solution becomes more user-friendly. The benefit, of course, is that once set up, the running cost is low (no subscription needed, no dependence on a company's continued support) and the flexibility is high (we can modify how it works freely, add custom behaviors, etc.). Enthusiast communities can contribute improvements (like better Vietnamese support or more devices). This is something that a closed ecosystem can't match easily.

**Potential Improvements:** Based on our experience, some improvements and extensions are worth exploring: - **Wake Word Integration:** As mentioned, adding a wake word ("Hey Home") for hands-free activation would enhance user experience. This could be done via existing libraries and run on the same Pi with minimal extra load. - **Continuous Listening Mode:** For

cases where a user might give a follow-up command (“Turn on the lamp” followed by “dim it to 50%” – assuming a dimmable light), one could allow a short window of continuous conversation. Our project didn’t involve dimming (just on/off), but it suggests exploring context management if we had more complex dialogues. - **Enhanced NLU:** If expanding to more devices, employing a smarter NLU (maybe an open-source model or a small locally hosted transformer for intent classification) could reduce the manual sentence enumeration. This however may require a more powerful machine for inference, or using a cloud LLM with privacy safeguards (perhaps a self-hosted one). - **User Authentication:** Voice biometrics were out of scope, but in some smart home cases you might want to ensure it’s an authorized person speaking. We could theoretically use Whisper’s speaker diarization capabilities or a separate model to recognize voices. This wasn’t critical for a lamp, but for security systems (door locks, alarms) voice authentication would be a consideration. - **Broader Device Support:** Integrating this system with more complex appliances (thermostats, media players) could be done. For example, controlling a TV via voice using this pipeline would involve intents like “increase volume” or “play next episode”. It would be a fun extension and primarily involves adding those capabilities to Home Assistant and adding sentences. The pipeline stays the same. - **Multilingual Handling:** We tested English and Vietnamese separately. Handling both simultaneously (detecting language automatically) is something Whisper can do – it identifies language and could feed into different intent handlers per language. In a multilingual household, one could let people use whichever language they’re comfortable with. We did not fully implement automatic language detection routing, but Whisper actually outputs a language code with each transcription. One could script Home Assistant to route the text to either English or Vietnamese intent matcher based on that. This would be a next step for a truly bilingual assistant. - **Comparison with other STT:** Whisper was our choice, but there are others (like VOSK, PocketSphinx for ultra-lightweight, or Mozilla DeepSpeech). It would be worthwhile to compare their accuracy and speed on the same hardware. We suspect Whisper is top-tier in accuracy albeit heavier; a smaller model might double the speed at some cost to accuracy. Depending on the use case, that trade-off might be acceptable.

**Broader Impact:** The positive result of this project illustrates that even small-scale implementations can benefit from the latest AI models without cloud dependency. We effectively ran a mini voice assistant with an edge device. This has implications beyond just lamps – think of healthcare devices, accessibility tools for those who cannot use hands, or industrial voice controls – all can potentially be done offline now. Privacy and data sovereignty are increasingly valued, and our work shows that the technology is ready to support that ethos. The collaboration of IoT (ESP8266) with AI (Whisper) and home automation software (Home Assistant) is a powerful combination that leverages strengths of each: cheap hardware interfacing with the physical world, sophisticated algorithms for understanding humans, and robust automation logic to tie everything together.

In conclusion of the discussion, our project stands as a compelling demonstration of a **privacy-first, local smart home voice assistant**. It achieves competitive performance with mainstream solutions for the specific domain of lighting control, while offering superior privacy and offline reliability. The challenges that remain (ease of setup, scaling to many commands) are being actively addressed by the open-source community, suggesting that such systems could become more widespread. Next, we provide a brief conclusion summarizing the achievements and learnings from this project, and how it could be taken forward.

## Conclusion

This project successfully developed a voice-controlled smart lamp system using an ESP8266 NodeMCU and Home Assistant, fully aligned with a local-first and privacy-preserving design philosophy. We integrated a state-of-the-art offline speech recognition model (OpenAI Whisper) with Home Assistant’s intent handling to interpret voice commands without any cloud services. The NodeMCU, running ESPHome firmware, interfaces with a relay to control a standard 220 V lamp, completing the IoT loop. Our implementation demonstrates that even with modest hardware, it is feasible to achieve high accuracy in voice command recognition ( $\approx 97\%$ ) and acceptable responsiveness ( $\sim 2\text{--}3$  s command latency) for typical smart home interactions.

In adhering to IEEE-style project reporting, we presented the full system architecture, including figures for the hardware wiring and voice processing pipeline, and provided code snippets (ESPHome YAML configuration and Home Assistant intent scripts) to illustrate the technical solution. The system’s performance was rigorously tested; results indicate strong reliability and robustness, meeting the needs of everyday use. Importantly, the solution maintains functionality during internet outages and keeps all voice data local, directly addressing privacy concerns identified in the literature and by users. By comparing our approach to existing commercial and open-source solutions, we highlighted that our system matches the convenience of voice assistants like Alexa while avoiding their downsides such as data snooping and cloud dependency[2][7].

There are several takeaways and implications from this work. First, the combination of open-source tools (Home Assistant, ESPHome) with advanced AI models empowers individuals to create sophisticated smart home systems that were once the domain of big tech companies. Second, the project reinforces the idea that privacy and functionality are not mutually exclusive in IoT; with careful design, one can have a voice-controlled home that doesn’t send recordings to the internet. Third, the modular nature of our system means it can be expanded – both in scale (more devices, more command types) and in capability (adding features like wake word detection, multi-language support, etc.). The use of standardized components and protocols (MQTT/API, YAML configs, etc.) ensures that the system can interoperate with other smart home components easily.

For future work, several avenues can be explored. One is to extend the voice control to a broader range of devices and actions, effectively building a full home voice assistant. This could include incorporating context-awareness or sequential command handling. Another area is optimization: tuning the ASR model or employing hardware accelerators to reduce latency would enhance user experience, especially for longer or more complex commands. Improving the user-friendliness of the setup – perhaps packaging the configuration as a Home Assistant add-on or blueprint – would help others replicate this project without deep technical expertise. On the research front, conducting user studies on interaction with such a local voice system could yield insights into how much the knowledge of privacy impacts user comfort in using voice interfaces (our hypothesis is that users would be more willing to use voice commands at home if they trust the system is not eavesdropping).

In conclusion, the Voice-Controlled Smart Lamp System we built serves as a tangible proof-of-concept that **smart home voice assistants can be implemented in a secure, private, and self-**

**contained manner** without sacrificing performance. It embodies the principle that emerging technologies like AI and IoT can be combined responsibly to enhance daily life. We hope this project can inspire further developments in privacy-centric smart home solutions and contribute to the growing momentum in the open-source community towards transparent and user-respecting home automation.

## References

- [1] P. Schoutsen, “2023: Home Assistant’s year of Voice,” *Home Assistant Blog*, Dec. 20, 2022. [Online]. Available: [40][14]
- [2] M. A. Shakaramiro, A. Gunaryati, and B. Rahman, “Voice Command-Based IoT on Smart Home Using NodeMCU ESP8266 Microcontroller,” *PIKSEL: Penelitian Ilmu Komputer Sistem Embedded and Logic*, vol. 12, no. 1, pp. 35–46, 2024. [12]
- [3] W. Syamsuri, Y. Wijanarko, and I. Maja, “Smart Room Based on ESP32 and Google Assistant: Automation Solution in Electrical Engineering Laboratory Environment,” *Iota (IEEE)*, vol. 5, no. 2, pp. 10–17, 2025. [11][10]
- [4] C. West, “Custom Voice Control Sentences in Home Assistant,” *mostlychris.com Blog*, Dec. 3, 2023. [Online]. Available: [23][31]
- [5] S. Woo, “Whisper, A Breakthrough in Speech-Recognition AI,” *ENERZAI Blog*, May 27, 2025. [Online]. Available: [17][16]
- [6] R. Santos, “ESP8266 NodeMCU Relay Module – Control AC Appliances (Web Server),” *Random Nerd Tutorials*, Dec. 17, 2019. [Online]. Available: [41][42]

---

[1] [8] [9] [13] [14] [40] 2023: Home Assistant's year of Voice - Home Assistant

<https://www.home-assistant.io/blog/2022/12/20/year-of-voice/>

[2] [3] [4] [34] Smart Home Devices and Privacy Risks: Is Your Alexa Spying on You? | by Mizter\_AB | Medium

<https://medium.com/@abrahamedet9/smart-home-devices-and-privacy-risks-is-your-alexa-spying-on-you-e9f4e0465a4d>

[5] [6] [7] [35] The voice assistant space is dying:  
<https://arstechnica.com/gadgets/2022/11/amaz...> | Hacker News

<https://news.ycombinator.com/item?id=33705928>

[10] [11] [32] [pubs.ascee.org](https://pubs.ascee.org)

<https://pubs.ascee.org/index.php/iota/article/download/924/586>

[12] Journal articles: 'ESP8266 NodeMCU microcontroller' – Grafiati

<https://www.grafiati.com/en/literature-selections/esp8266-nodemcu-microcontroller/journal/>

[15] [16] [17] [27] [28] Whisper, A Breakthrough in Speech-Recognition AI

<https://enerzai.com/resources/blog/whisper-a-breakthrough-in-speech-recognition-ai>

[18] [19] [20] [21] [22] [41] [42] ESP8266 NodeMCU Relay Module - Control AC Appliances (Web Server) | Random Nerd Tutorials

<https://randomnerdtutorials.com/esp8266-relay-module-ac-web-server/>

[23] [24] [29] [30] [31] [38] [39] Custom Voice Control Sentences in Home Assistant

<https://www.mostlychris.com/custom-voice-control-sentences-in-home-assistant/>

[25] [26] [33] A voice for the Open Home

<https://newsletter.openhomefoundation.org/voice-for-the-open-home/>

[36] I replaced Google Home with Home Assistant and a local LLM, and I ...

<https://www.xda-developers.com/replaced-google-home-home-assistant-local-llm/>

[37] Talking with Home Assistant - get your system up & running

[https://www.home-assistant.io/voice\\_control/](https://www.home-assistant.io/voice_control/)