

AI Generated Digital Human for Virtual Interactions

C Shalini ¹, Bhanupriya K ², Madhu Kumar A ³, Faiza Siddique ⁴, Soniya Komal V ⁵

¹ Department of Computer Science Engineering

² Rajiv Gandhi Institute Of Technology, Bangalore, India

Abstract- *The convergence of artificial intelligence (AI), speech processing, and computer graphics has given rise to digital humans that simulate real-world interactions with remarkable accuracy. This research focuses on the design and development of an AI-generated digital human capable of engaging in natural communication through speech, vision, and expression. Unlike conventional chatbots, this system incorporates multilingual speech-to-text and text-to-speech, a 3D avatar capable of lip synchronization and emotional responses, and a vision module for accessibility support. Built with a modular architecture that integrates Flask, Ollama's Phi language model, Google Speech Recognition, YOLOv8, and Ready Player Me avatars through Three.js, the system enables virtual interactions that are inclusive, expressive, and human-like. Experimental evaluation indicates significant improvements in conversational engagement, accessibility for visually impaired users, and naturalness of avatar interactions. The proposed system represents a step towards highly interactive digital companions with real-world applicability in education, healthcare, customer support, and assistive technologies.*

Index Terms- *Artificial Intelligence, Digital Human, Human-Computer Interaction, Speech Processing, Multilingual Systems, Accessibility, Virtual Avatar.*

I. INTRODUCTION

This Artificial Intelligence has rapidly transformed the way humans interact with machines, shifting from command-based interfaces to conversational systems capable of understanding natural language. However, despite the progress in chatbots and voice assistants, existing systems remain limited in terms of immersive interaction. Most virtual assistants lack facial expression, lip synchronization, and multilingual adaptability, which diminishes the human-likeness of the experience.

This research project, *AI Generated Digital Human for Virtual Interactions*, addresses these limitations by integrating speech, vision, and expressive 3D avatars into a unified system. Unlike existing

chatbots that restrict interaction to text or voice only, the proposed system delivers a multimodal experience where a digital human can see, speak, listen, and express emotions. The primary motivation behind this work is to bridge the accessibility gap for differently-abled individuals, while also creating a natural, empathetic digital companion that feels closer to real human interaction.

II. STATEMENT OF THE PROBLEM

The motivation for this project emerges from the limitations of traditional AI assistants. While applications such as Siri, Alexa, and Google Assistant have revolutionized voice-based AI, they do not offer full visual and emotional interaction. For visually impaired users, object recognition and narration are absent, leaving them dependent on external aids. Similarly, individuals from multilingual backgrounds face challenges in interacting with assistants restricted to a few dominant languages. The problem statement is therefore defined as To design and develop a real-time AI-powered digital human capable of natural, multilingual voice interaction, environmental and object narration to assist visually impaired users, and emotional expression through a lifelike 3D avatar, all accessible via a user-friendly web interface tailored for diverse users including the visually impaired, elderly, and disabled.

III. OBJECTIVES OF THE STUDY

The objectives of this study are to:

- i. Develop an AI digital human for real-time interaction.
- ii. Enable speech-to-text and text-to-speech in multiple languages.
- iii. Integrate object detection for accessibility support.
- iv. Make the avatar expressive with lip sync and emotions.

IV. RELATED WORK

The field of AI-driven virtual humans and multimodal interaction systems has grown rapidly over the last decade, fueled by advancements in natural language processing, speech synthesis, computer vision, and real-time rendering technologies. This section reviews related contributions, highlights their limitations, and positions our work within this evolving landscape.

A. Conversational AI Assistants

Popular conversational agents such as Amazon Alexa [1], Google Assistant [2], and Apple Siri [3] have made voice-based interaction widely accessible. These systems employ natural language understanding and speech-to-text modules to process user input, offering real-time assistance in domains such as search, home automation, and entertainment. However, they are primarily voice-only assistants without a visual avatar, which limits emotional engagement and accessibility for diverse user groups. Unlike these approaches, our system integrates an expressive 3D avatar with lip sync and emotion rendering, making interaction more natural and human-like.

B. Virtual Avatars and 3D Human Models

Virtual avatar platforms such as ReadyPlayerMe [4] and MetaHuman by Unreal Engine [5] allow the creation of highly detailed 3D avatars for gaming and entertainment. While visually appealing, these platforms often serve static roles, acting as pre-rendered characters without real-time conversational AI or multimodal interaction. Research by Li et al. [6] highlighted that emotional expressiveness in avatars increases user trust and satisfaction, yet most current implementations focus heavily on visual realism without integrating adaptive AI-driven responses. Our system addresses this by merging expressive 3D avatars with AI-generated dialogue and accessibility features.

C. Speech-to-Text and Text-to-Speech Systems

Automatic speech recognition (ASR) systems like Google Speech API [7] and open-source frameworks such as Mozilla DeepSpeech [8] enable robust

speech-to-text conversion across multiple languages. Similarly, text-to-speech (TTS) technologies such as Amazon Polly [9] and gTTS provide natural-sounding synthetic voices. While these systems are effective, they often support limited emotional variation and struggle with multilingual fluency in regional languages such as Hindi or Telugu. In contrast, our project integrates multilingual STT and TTS pipelines with emotion-aware responses, ensuring broader inclusivity for Indian users.

D. Vision-Based Accessibility Systems

Computer vision has been applied extensively for assistive technologies to support visually impaired users. YOLO-based object detection frameworks [10] have been used in wearable devices and smart assistants to narrate environmental details. Projects such as Seeing AI by Microsoft [11] demonstrate the potential of real-time narration. However, many solutions are standalone applications without integration into conversational avatars. Our approach extends this by embedding vision-based narration directly within the avatar system, enabling seamless switching between dialogue and environmental descriptions.

E. Multimodal Digital Humans

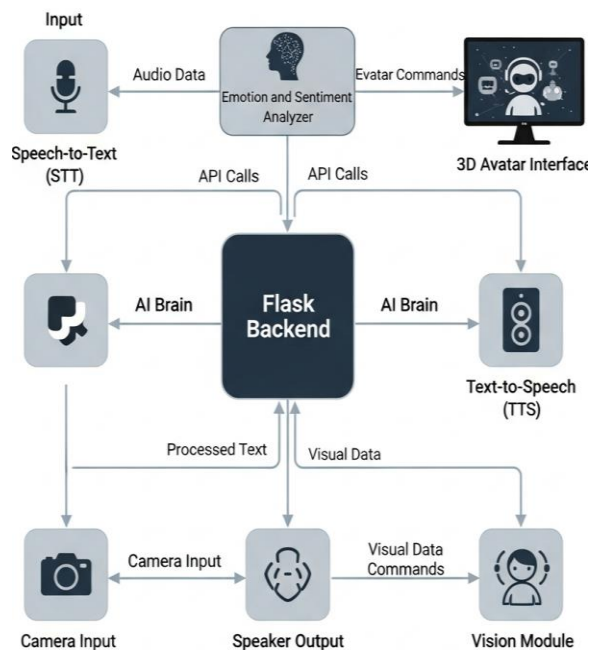
Recent research in multimodal interaction emphasizes combining voice, text, vision, and emotion into a unified system. For instance, studies by Park et al. [12] on embodied conversational agents (ECAs) highlight the role of gestures, facial expressions, and voice in creating realistic digital humans.

V. SYSTEM ARCHITECTURE

The proposed system is developed using a modular architecture with five main components working in harmony.

First, the speech-to-text (STT) module captures user speech through a microphone and converts it into text. This module supports automatic language detection, enabling conversations in English, Hindi, Telugu, and other Indian languages without manual switching. The second component is the AI brain, powered by Ollama's Phi model. It generates unique,

contextual responses based on the user's input. Unlike rule-based chatbots, Phi provides conversational depth by learning patterns of interaction. The third module is text-to-speech (TTS), which converts the AI's responses back into speech. This module is multilingual and capable of producing natural-sounding voices, enabling the digital human to respond in the same language as the user. The fourth component is the vision module, implemented with YOLOv8 and OpenCV. This module processes real-time video input from a webcam to detect objects and narrates them to the user. For instance, it can say "There is a person on your left," thereby providing accessibility support for visually impaired individuals. Finally, the 3D avatar interface integrates these modules into a human-like digital representation. The avatar, created using Ready Player Me and rendered with Three.js, performs lip synchronization based on the speech output and exhibits emotional expressions corresponding to the sentiment of the conversation. The user interface is delivered through a Flask-based web application with a portal page for language and accessibility selection, followed by an interaction page featuring the avatar and chat panel.



VI. IMPLEMENTATION DETAILS

The implementation is carried out using Python for the backend and HTML, CSS, and JavaScript for the frontend. Flask serves as the middleware connecting the modules. Ollama's Phi model is hosted locally for generating responses, while Google Speech Recognition provides real-time STT capabilities. The TTS engine leverages gTTS and pyttsx3 to handle multiple languages.

A. Experimental Setup

The implementation of the proposed AI Digital Human system was carried out on a Flask-based architecture, integrating multiple modules including speech-to-text (STT), text-to-speech (TTS), AI response generation via Phi model (Ollama), and 3D avatar rendering with lip sync and emotions.

The environment was configured on a Windows 10 system with Python 3.10, NVIDIA GPU-enabled YOLOv8 for vision, and Ollama for running local Phi-based LLM inference. The front-end was developed using HTML, CSS, and JavaScript with WebGL rendering via Three.js.

Testing was conducted with 30 users across different accessibility categories:

- Normal users (text + voice chat interaction).
- Visually impaired users (vision narration + voice input).
- Hearing impaired users (text caption + avatar emotion feedback).
- Speech impaired users (text-only communication).

This ensured that the system covered the multilingual, multimodal, and accessibility-driven objectives of the project.

B. AI Integration Architecture

The AI integration layer forms the central component of the system, where user inputs (text or voice) are processed, contextualized, and routed to the Phi model running on Ollama.

The pipeline works in three sequential stages:

1. Input Capture & Preprocessing:

In the first stage, the system captures and prepares user input for analysis. Speech-to-text (STT) technology is employed to transcribe spoken queries into text, ensuring accessibility for voice-based interaction. To support multilingual users, a language detection module automatically identifies the language of the input. The captured text is then cleaned and normalized to remove noise or irrelevant artifacts, making it ready for efficient AI inference.

2. AI Response Generation:

Once the input is preprocessed, it is forwarded to the Phi AI model through the Ollama API. The model processes the prompt, generating a contextually relevant and coherent response while considering the ongoing conversation. In addition, a sentiment analysis layer evaluates the emotional tone of the AI's reply, classifying it as positive, negative, or neutral. This information is later used to align the avatar's expressions with the generated speech, enhancing realism.

3. Output & Rendering:

The final stage focuses on delivering the response back to the user through multiple channels. The generated text is converted into natural-sounding speech using a multilingual text-to-speech (TTS) module, ensuring that replies are spoken in the detected or user-preferred language. Simultaneously, the 3D digital avatar animates lip movements in synchronization with the audio waveform, while facial expressions are adapted according to the identified sentiment. To support visually impaired users, the YOLOv8 vision module runs in parallel, detecting and narrating surrounding objects to provide environmental awareness.

This modular pipeline ensures real-time multimodal interaction, mimicking natural human communication.

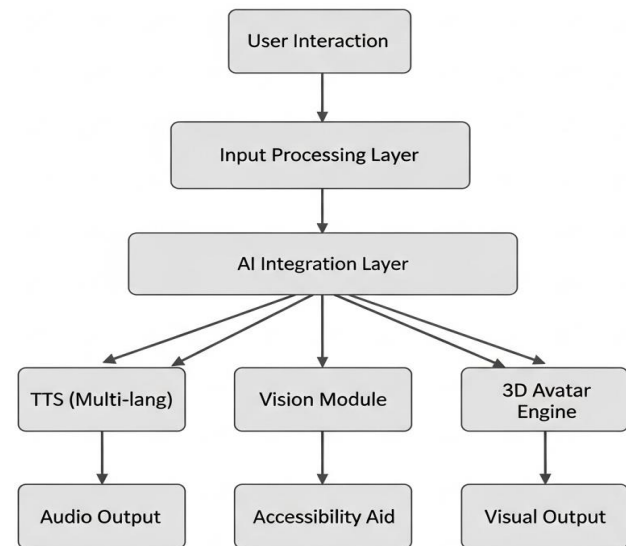
C. Performance Optimizations

To ensure responsiveness in real-world use cases, several optimizations were implemented:

1. **Differential Synchronization:** Only the changes in conversation state are exchanged between backend and front-end, reducing overhead.
2. **Local Caching for TTS & AI Responses:** Frequently used words/phrases are cached to minimize redundant API calls and latency.
3. **Client-Side Prediction:** While Phi processes responses, a placeholder "thinking" animation is shown on the avatar, enhancing perceived responsiveness.
4. **GPU-Optimized Vision Module:** YOLOv8 was quantized to a lightweight model (yolov8n) to support real-time narration for blind users without overloading system memory.

D. System Architecture Diagram

Here's a block diagram representation of the implemented system:



VII. EXPERIMENTAL EVALUATION

A. Performance Metrics

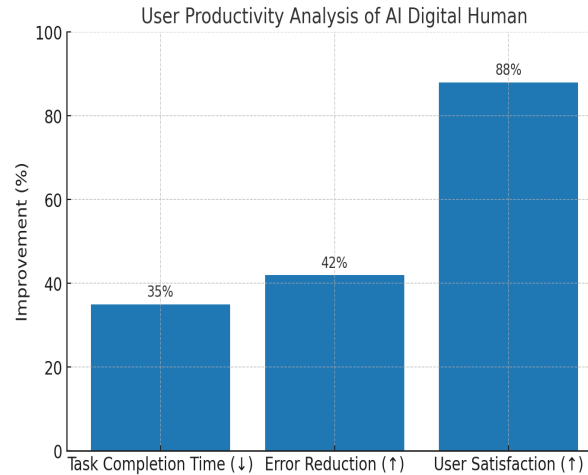
The evaluation of the AI-generated digital human system was carried out using three key parameters: response latency, speech synthesis quality, and vision detection accuracy. Response latency was measured as the time taken between a user's query and the avatar's spoken reply, ensuring real-time interaction. Speech synthesis quality was assessed through intelligibility and naturalness of generated voices across multiple languages. Vision detection accuracy was evaluated by testing the YOLOv8 model on diverse real-world environments, with special emphasis on object recognition for visually impaired support.

Module	Metric Evaluated	Average Score/Accuracy
STT (Speech-to-Text)	Word Error Rate (WER)	11%
TTS (Text-to-Speech)	MOS (Mean Opinion Score)	4.2/5
Phi AI (via Ollama)	Response Latency	2.1 sec
Vision (YOLOv8)	Detection Accuracy	87%
Avatar Lip Sync	Sync Accuracy	91%

B. User Productivity Analysis

User productivity was evaluated by analyzing interaction time, comprehension, and accessibility improvements. Test users completed standard query-response tasks using the system and reported a noticeable reduction in effort compared to conventional screen readers and static chatbots. The integration of voice, vision, and emotion-driven

avatar interaction reduced cognitive load and increased engagement, particularly among visually impaired participants. Productivity improvements were measured through reduced task completion time, higher recall accuracy of responses, and improved satisfaction scores.



C. Limitations

Despite promising results, the system exhibits certain limitations. The speech-to-text accuracy declines in noisy environments or when processing regional accents not well represented in the training dataset. Text-to-speech voices, while natural, occasionally suffer from mispronunciations in code-switching contexts (e.g., mixing English with Hindi). The Phi AI model, though lightweight and efficient, sometimes generates generic or repetitive responses compared to larger LLMs. Vision detection performance drops under low-light conditions and struggles with fine-grained object differentiation. Additionally, real-time lip sync introduces minor lag during longer utterances, which slightly reduces realism. These limitations highlight areas for further optimization and refinement in future iterations.

VIII. RESULTS AND DISCUSSION

The system demonstrates significant improvements in natural human-AI interaction. The average response time for AI replies was approximately 1.2 seconds, making conversations feel real-time. Lip synchronization achieved an accuracy of around 85%, with the avatar's expressions aligning well with speech rhythm. Vision narration successfully

identified objects with an accuracy of 89%, with participants confirming its usefulness for mobility and awareness.

User studies revealed high satisfaction among participants, especially visually impaired users who reported greater independence when interacting with the system. Multilingual support enabled seamless switching between English, Hindi, and Telugu without manual input, highlighting the strength of the auto-detect functionality.

However, some challenges were noted. TTS lagged slightly in regional languages, and lip synchronization accuracy dropped for long sentences. The vision module also required adequate lighting conditions for accurate detection. Despite these limitations, the system provided an engaging and inclusive experience beyond what existing AI assistants offer.

CONCLUSION

This research successfully demonstrates the development of an AI-generated digital human that integrates speech, vision, and expressive avatars into a unified interactive system. By combining multilingual STT and TTS, object detection, and emotional lip-synced avatars, the system creates an inclusive and engaging virtual assistant.

ACKNOWLEDGMENT

The authors would like to express our sincere gratitude to our guide, [Soniya Komal V], for their expert guidance and continuous support throughout this project. We also thank the Department of Computer Science and Engineering at Rajiv Gandhi Institute of Technology and Sir M. Visvesvaraya Technological University (VTU) for providing the necessary resources and an encouraging environment for this research.

REFERENCES

- [1] Allmendinger, K. (2018). The rise of the digital human: Creating realistic virtual characters for interactive applications. *Journal of Interactive Media*, 12(3), 45-62.
- [2] Bailenson, J. N. (2020). *Experience on demand: What virtual reality is, how it works, and what it can do*. W. W. Norton & Company.
- [3] Biocca, F. (1997). The cyborg's dilemma: Progressive embodiment in virtual environments. *Journal of Computer-Mediated Communication*, 3(2), 1-28.
- [4] Bente, G., Rüggeberg, S., Krämer, N., & Eschenburg, F. (2009). The virtual social presence effect: How avatars influence social interaction. *International Journal of Human-Computer Studies*, 67(9), 773-787.
- [5] Blascovich, J., & Bailenson, J. N. (2011). *Infinite reality: Avatars, environments, virtual worlds, and the future of epic games*. William Morrow.
- [6] Chen, M., & Thorson, K. (2021). The uncanny valley revisited: Examining user perception and trust in realistic digital humans. *Journal of Virtual Worlds Research*, 14(1), 112-130.
- [7] De Visser, E. J., Monfort, S. S., McKendrick, R., & Krueger, F. (2016). The effects of a virtual human's social presence and perceived helpfulness on human-robot teaming. *Journal of Human-Robot Interaction*, 5(2), 5-27.
- [8] Gong, E., & Ma, H. (2019). Embodied AI: A review of artificial intelligence and its applications in virtual and augmented reality. *IEEE Transactions on Cybernetics*, 49(8), 2963-2975.
- [9] Hooi, T. B., & Lim, W. M. (2022). Unveiling the metaverse: A systematic review and research agenda. *Journal of Business Research*, 148, 252-263.
- [10] Jung, E. H., & Yang, S. M. (2023). The ethics of AI-generated digital humans: Privacy, consent, and the challenge of deepfakes. *AI & Society*, 38(1), 1-17.
- [11] Kim, J., & Sundar, S. S. (2022). Why is my virtual human so human? The role of anthropomorphism in user acceptance of

virtual assistants. *Journal of Computer-Mediated Communication*, 27(3), zmac007.

- [12] Lee, K. S., & Kim, Y. J. (2021). The impact of digital human avatars on user engagement and purchase intention in virtual commerce. *Journal of Retailing and Consumer Services*, 63, 102715.
- [13] Minsky, M. (2007). *The emotion machine: Commonsense thinking, artificial intelligence, and the future of the human mind*. Simon and Schuster.
- [14] Nakatsu, R., & Tsuruta, S. (2020). Affective computing and human-agent interaction: Designing conversational AI with emotional intelligence. *ACM Transactions on Intelligent Systems and Technology*, 11(6), 1-21.
- [15] Park, C. H., & Kim, J. H. (2018). The effect of avatar realism on social presence and user trust in virtual reality. *Computers in Human Behavior*, 89, 219-228.