# Creating a Digital Human Twin: Cloning Voice, Face, and Attitude

1st Karim Ahmed Wahba
*Computer Science(CS)*
*Misr International University (MIU)*
Cairo, Egypt
karim2005494@miuegypt.edu.eg

2nd Khaled Amr Ahmed
*Computer Science(CS)*
*Misr International University (MIU)*
Cairo, Egypt
khaled2002775@miuegypt.edu.eg

3rd Martina Raafat Kamel
*Computer Science(CS)*
*Misr International University (MIU)*
Cairo, Egypt
martina2009403@miuegypt.edu.eg

4th Marwan Fathy
*Computer Science(CS)*
*Misr International University (MIU)*
Cairo, Egypt
marwan2014606@miuegypt.edu.eg

5th Pr.Khaled Hussien Abdelfatah
*Computer Science(CS)*
*Misr International University (MIU)*
Cairo, Egypt
khaled.hussein@miuegypt.edu.eg

6th Sarah Hatem
*Computer Science(CS)*
*Misr International University (MIU)*
Cairo, Egypt
sarah.nabil@miuegypt.edu.eg

*Abstract*—This paper presents and explores the development of digital human twin through mobile app integrated with chat-bot and aimed to deliver dynamic and realistic user interactions. By utilizing advanced technologies such as reinforcement learning, lip-sync animation, and facial expression modeling. This project aims to develop a personal virtual assistant using various technologies. The system integrate various deep learning models such as 3D avatar generation, speech recognition, clone voice from the user, Large Language Model (LLM), 3D avatar after uploading images from the user, and engages in lifelike human conversation. In the evaluation of the system there's a potential to improve user experience. Moreover, the paper also includes a review of related datasets, highlighting the capabilities and challenges of creating a realistic virtual avatar.

## I. INTRODUCTION

The emergence of digital human twins is transforming our live by producing realistic avatar that dynamically responsive and highly personalized with a high quality. There's no doubt that there are a lot of digital avatar have a wide range of applications. It includes serving as customer service agents, virtual companions for the the elderly individuals. Also, there's personalized health monitors. By replacing human behaviour and appearance these realistic avatar offer more immersive and realistic interaction. This technology is to boost the engagement and deliver customized services based on individual user data and interactions. We will discuss the development of a personalized virtual assistant avatar that redefined the boundaries of the computer Interaction. This is by creating a realistic virtual 3D avatar based on some inputs for a specific human photo and creating an admirable experience, by integrating with facial expressions, lip-sync, and full-body motion. Moreover, after generating the virtual 3d avatar with the same characteristics the user provided, the avatar should start to interact with the user from a spoken command processed through Natural Language Processing (NLP) The

model will be trained on the facial expression dataset. The SMPL is used to provide accurate representation of the human body. Our objective is to create a personalized and interactive 3D avatar by integrating various deep learning models. Each model serves a specific purpose in achieving different aspects of the avatar. These aspects include image, voice, and text, which are addressed through the use of speech recognition, a 3D avatar generator, and an LLM ( Large Language Model), respectively. The user has the ability to generate the avatar by uploading an image that is processed to generate the avatar's body and face, capturing or uploading an audio recording for cloning the user's voice, and engaging in a conversation with the LLM ( Large Language Model) model For generating the responses on the user's input. Our project aims to achieve personalized and engaging user experience in communication platforms. Taking advantage of the advanced technologies including lip-sync animations, facial expressions generation, and reinforcement learning, this gives the user a dynamic conversational interface. The avatar component enhances interaction by answering a correct answer according to user's questions, while having the ability to learn more from the user interactions this enables the system to evolve itself into a personalized replica of the user. The implementation details is to integrate all of these features into a mobile application which should be a user-friendly interface and give the user a great experience. After we have finished explaining the system, several subjects will be covered in the upcoming sections. The related work section it reviews the existing literature on large language model (LLM), 3D avatar generation, and voice cloning. it covers these areas and their relevant on developing a realistic avatar to be shown to the users. We will see at datasets section our datasets that we used for this project including audio, text, and images datasets. Alpaca dataset for LLM training, SMPL dataset for 3D body, and VCTK dataset for voice cloning that is coming from the user. At the architecture section we have used layered architecture of the application.

It explains how these layers interact to provide a seamlessly user experience. The methodology section breaks down the implementation of specific tasks such as voice cloning discuss the use of deep learning and machine learning to clone the user's voice. Also, avatar creation discuss the steps of creating the 3D avatar. The implementation section . Implementation section covers the integration features into mobile app using flutter ensuring real time communication between the user and the digital twin. At the evaluation and result section this cover the integration of the system in order to the user engagement and interaction quality. It also discuss the potential improvements and challenges faced during the development process.

## II. RELATED WORK

### A. Maintaining the Integrity of the Specifications

In "A Survey of Large Language Models" emphasis on models greater than 10B, this offers a thorough study of (LLM). Findings, essential ideas, and methods for comprehending and applying LLM are covered, covering pre-trained, assessment, adaption, and utilization. The difficulties and potential paths for LLMs and their procedure are also covered. All of this is done to protect personal information and secure data. Enhancing the usage of LLMs in traditional NLP tasks, as well as experimental data and conclusions on closed-source models, are covered in the review. Overall, it offers insightful information regarding the capabilities, difficulties, and future direction of LLMs, making it a useful resource for engineers and researchers who are interested in the technology. The article also goes into detail on the dataset that was used to assess how well the Large Language Model (LLM) performed in a range of NLP tasks. The article also describes the zero-shot performance of LLM on these datasets, i.e., models are assessed on tasks for which they have not yet received training. This assessment score reflects LLMs' capacity to generalize to novel tasks. In general, using these datasets provide insightful information about how well LLMs do a range of NLP tasks.

The capabilities of LLMs are covered in The Future of Large Language Models: A Futuristic Dissection on AI and Human Interaction, including their capacity to produce language that matches the language of a human and manage difficult tasks like translation and summarization. It mentions using the 175 billion-parameter GPT-3 model from OpenAI to illustrate these skills. The practical uses of LLMs in healthcare, education, and customer service are discussed, along with any possible advantages and drawbacks. The study also highlights the necessity of continuous investigation, policy formulation, and interdisciplinary cooperation in order to tackle moral quandaries, assess practical implications, and engage the general public in conversations regarding the future of LLMs.

Based on Large Language Models Show Human Behavior The LLMs replicate human behavior in more subtle ways, like confident errors and sensitivity to question changes, in addition to basic language processing for problem solving. Large Language Models (LLMs) have been used for a number of activities, such as summarizing texts, translating languages, and writing creatively. LLMs can be employed as "participants" in psychological studies and can help us understand human behavior better. The similarity to human behavior observed in LLMs creates new avenues for investigation. Comparing studies involving human participants versus those using LLMs has the advantage of allowing for speedier and less constrained research because LLM data collection does not require informed permission or ethical guidelines.

In "Transfer Learning from Speaker Verification to Multispeaker Text-To-Speech Synthesis," a neural network-based system that can produce speech audio in the voices of different speakers—including ones not used during training—is presented. A speaker encoder network, a Tacotron 2–based sequence-to-sequence synthesis network, and a WaveNet–based vocoder network form the system. The model can produce high-quality speech audio with the traits of various speakers by utilizing transfer learning techniques from speaker verification challenges. The VCTK (Voice Cloning Toolkit) and LibriSpeech training datasets were used, and evaluations indicate that the findings are encouraging in terms of speaker similarity and speech naturalness. Still, it might be difficult for the model to adequately capture key subtleties associated with speaker attributes, like prosody and accent. The system may perform better overall for multispeaker text-to-speech synthesis if training data variety and model architecture are improved further. The system's accuracy is mainly evaluated in terms of Speaker Verification Equal Error Rates (SV-EERs) over a number of datasets in this study. These datasets are in use. CTK Dataset: 10.46 for SV-EER on VCTK; Training Set: 98 speakers from the VCTK dataset, SV-EER for LibriSpeech: 29.19 . LibriSpeech Dataset : - On VCTK, SV-EER: 6.26, Training Set: 1.2K speakers from the LibriSpeech dataset and SV-EER for LibriSpeech: 5.08

Based on a real-time voice cloning system which uses multiple algorithms to enhance speech quality. Speech synthesis, acoustic parameter creation, and text analysis are all covered by modules in the system. Although the system has benefits like the ability to generate natural sounds and requires less data. To increase voice quality, the study highlights advancements in long sentence processing, vocabulary pronunciation, and noise reduction techniques. For improved voice cloning tasks, the authors also recommend switching from the Tacotron structure to a transformer structure. When compared to other models, the system's overall performance in terms of fluency, naturalness, and clarity is encouraging. Its main goal is to improve the quality of voice synthesis through algorithmic improvements. The real-time voice cloning system outperformed DeepVoice, Tacotron2, and Fastspeech2 models with a mean Opinion

Score of 3.85. It was trained on 1000 audio recordings from the Librispeech ASR corpus.

In Adapting TTS models For New Speakers using Transfer Learning addresses transfer learning strategies to effectively adapt Text-to-Speech (TTS) models to new speakers, the PDF . To help pre-trained TTS models adjust to fresh speaker data, it presents both mixed and direct finetuning techniques. To prevent overfitting, mixed finetuning integrates data from both the original and new speakers. Direct finetuning entails training on new speaker data alone. In terms of naturalness, voice likeness, and style similarity, the document assesses the efficacy of various techniques. As contrasted to starting from scratch with almost 27 hours of data, the results demonstrate that employing just 30 minutes of fresh speaker data can produce competitive outcomes. While offering open-source tools for training high-quality TTS models with little data, the paper suggests more research on the best finetuning iterations. Using the Hi-Fi Multi-Speaker English TTS Dataset, measures such as Equal Error Rate, Gross Pitch Error, Voicing Decision Error, and F0 Frame Error are used to determine the accuracy of the voice cloning system.

The usage of layered surface volumes (LSVs) in the research article "Efficient 3D Articulated Human Generation with Layered Surface Volumes" offers a novel technique for creating high-quality and varied 3D assets of articulated digital humans. While volumetric representations in 3D GANs can be computationally wasteful due to the need for volume rendering, traditional 3D object representations struggle to capture fine features like hair, clothing, and accessories. Volumetric scene representations' representational capacity and template meshes' efficiency are combined by LSVs to address these problems. To realistically display off-surface details and volumetric structures, this is accomplished by introducing volumetric manifolds surrounding the template mesh surface with a non-zero thickness.The results show the importance of hand regularization to achieve natural-looking results and show improvements in face quality through face discrimination. The work highlights the need for more detail in 3D human GAN approaches and also discusses its limits and future research prospects.The AIST++ Dataset, DEEPFASHION Dataset, and SHHQ Dataset are the datasets used for 3D human digitization and modeling.

In TotalSelfScan: Learning Full-body Avatars from Self-Portrait Videos of Faces, Hands, and Bodies they used a powerful technique for creating full-body avatars from multiple monocular self-rotation images that capture the hands, body, and face. It is . Observation frames can be aligned with the canonical space more easily thanks to part deformation fields in this multi-part network approach to representing the human form in a canonical space. Using a part composition technique, TotalSelfScan integrates disparate parts into a unified full-body model to produce a coherent picture of human anatomy. In order to make the reconstructed

avatars more lifelike, the study presents a non-rigid ray transformation that renders photorealistic free-viewpoint films with a variety of human positions. The reconstructed models can be used in a variety of viewing angles and positions thanks to this rendering process, which produces genuine visual outputs. Both quantitative and qualitative reviews highlight the efficacy of TotalSelfScan and highlight its capacity to capture the intricate geometry and look of the human body. The quantitative method presented in the research reconstructs full-body avatars from monocular self-rotation achieves high accuracy.

Based on XAGen: 3D Expressive Human Avatars Generation For creating photo-realistic human avatars with expressive and disentangled controllability over a variety of variables, including shape, body position, jaw pose, and hand pose, XAGen is a state-of-the-art 3D generating model. In order to improve the reliability of the created avatars, the framework makes use of a multi-scale and multi-part 3D representation to capture features in small-scale regions such as hands and faces. XAGen disentangles the synthesis of the body, face, and hands by utilizing multi-part discriminators and a multi-part rendering approach. This enhances geometric quality and makes model training more efficient.Due to the method's novel rendering technique, various body components can be rendered independently, giving the user precise control over facial expressions, jaw positions, torso poses, and hand poses. In the paper, the XAGen model was tested on four datasets: DeepFashion, MPV, UBC, and SHHQ.

## III. DATASETS DESCRIPTION

The dataset gathered for this project is from a number of research papers. Three categories of datasets apply. Audio, Text and Images Because we used audios to create the voice clone, Text to generate Large Language Model and photos to generate the 3D avatar.

### A. SMPL Dataset

The SMPL models use linear positions to determine the 3D positions of the body joints, creating a 3D model of the human body. This will make it possible to accurately show the human body and its movements. 10 form attributes that will record variations in body properties define this model. Additionally, it has 69 posture characteristics that specify the rotation of 23 distinct body joints in addition to one global joint.

### B. LLM Dataset

Alpaca is a dataset of 52,000 guidelines and examples generated by OpenAI's text-davinci-003 engine. By using this instruction data, language models can go through instruction-tuning, enhancing the language model's ability to follow instructions. The authors built on the data generation pipeline from Self-Instruct framework and made some modifications such as using text-davinci-003 engine instead of DaVinci,

more aggressive batch decoding was used, single instance was generated for each instruction, instead of 2 to 3 instances as in Self-Instruct.

| input<br>string · lengths | output<br>string · lengths | text<br>string · lengths |
|---|---|---|
| 0          2.47k | 0          4.18k | 154          4.5k |
| quickly the brown fox jumped | The quick brown fox jumped quickly. | Below is an instruction that describes a task,… |
| The world has been greatly impacted by the… | The tone of the text is one of concern and… | Below is an instruction that describes a task,… |
| [2, 3, 7, 8, 10] | The median of the given data is 7. | Below is an instruction that describes a task,… |
| Although it is generally accepted that the… | The internet has allowed us to connect globally,… | Below is an instruction that describes a task,… |
| | The logo should feature a green motif that is… | Below is an instruction that describes a task… |
| | Joy flows through me like a river clear, Bringing… | Below is an instruction that describes a task… |

Fig. 1.  A sample of Alpaca dataset



Fig. 2.  layered architecture

### C. Voice Cloning Dataset

The VCTK dataset, which stands for the Voice Cloning Toolkit dataset, contains clean speech data from multiple speakers. The dataset was used by the system to train multiple vocoder and synthesis networks. In the context of voice cloning and speaker adaption tasks, this dataset is particularly useful for training and testing speech synthesis models.

## IV.  ARCHITECTURE

### A. layered architecture

The architecture that has been used in DIGI application is layred architecture. The layer known as the Presentation Layer is in charge of displaying the user interface. This covers components such as the home page, avatar page, login page, and sign-up form. It also has components, such an interaction interface and a customization page, that enable modification and interaction. The business layer, is in charge of the web application's main features. Voice recording, questionnaire processing, LLM (Large Language Model) processing, and voice cloning. Avatar creation, user identification, and digital twin customization are also included. The persistence layer, is in charge of storing user data. It is made up of a server for local data storage and server-side communication. The data is stored on the database layer, also known . The databases for SMPL, Alpaca, and Deepfashion are displayed in the picture. Essentially, the business layer processes the data and connects with the persistence layer to store or retrieve data as needed when a user interacts with the web application through the presentation layer.
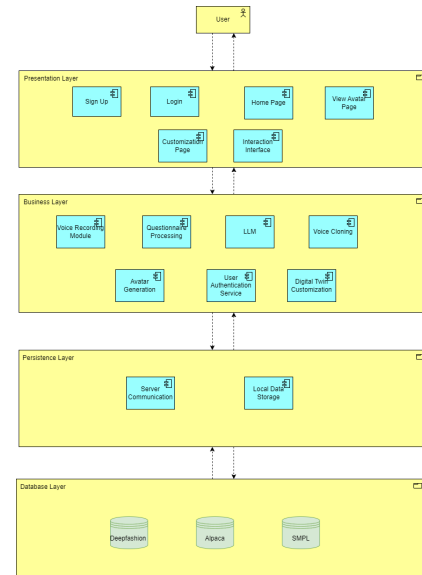
### B. Mobile Application with Flutter

The mobile application, developed using Flutter, serves as the user interface for interacting with the digital human twin.

- **Cross-Platform Development**: Flutter's single codebase allows deployment on both Android and iOS, ensuring a consistent user experience across platforms.
- **UI/UX Design**: The application features a user-friendly interface that supports voice and text interactions. The UI components are designed to be intuitive, enhancing user engagement.
- **Real-Time Communication**: Flutter integrates with backend services using HTTP and WebSocket protocols to facilitate real-time communication between the user and the digital twin.

## V.  METHODOLOGY

The following sections make up the methodology of our work as it was put into practice:

### A. Voice Cloning

*1) Deep Learning:-:* Deep neural networks are used to derive speaker embeddings from reference speech, producing fixed-dimension vectors for use in speech recognition. Tacotron 2, a deep learning network, aligns text inputs with speaker embeddings to produce mel spectrograms. WaveNet, a deep generative model, generates time-domain waveform samples for superior speech synthesis.

*2) Machine Learning:-:* The model uses open datasets like LibriSpeech and VCTK, with clear audio recordings of British-accented speakers. Preprocessing involves dividing data into train, validation, and test subsets. Performance evaluation includes speech naturalness, speaker resemblance, and generalization to unknown speakers.

### B. Avatar Creation

*1) Data acquisition:-:* In this stage, 2D photos of people in various positions and viewpoints are gathered. Moreover, 3D geometry may be captured by using depth data from sensors such as structured light scanners or LiDAR.

*2) Preprocessing:-:* In order to extract pertinent features, align poses, and normalize the input for model training, the gathered data is preprocessed. This stage could involve data augmentation, image scaling, and background removal in order to improve the dataset's quality and diversity.

*3) Model Training:-:* To discover the underlying patterns and connections between 2D images and 3D human shapes, deep learning models like PIFuHD, PyTorch SMPL, and GANStyle are trained on the preprocessed data. These models gather complicated spatial information and produce realistic 3D representations using techniques such as generative adversarial networks (GANs), recurrent neural networks (RNNs), and convolutional neural networks (CNNs).

*4) Evaluation:-:* The performance of the trained models in recreating 3D human shapes, producing realistic textures, and maintaining minute features like hair and accessories is measured using metrics including accuracy, precision, recall, and F1 score. The efficacy of the models is verified in part by quantitative comparisons with ground truth data and qualitative evaluations performed through eye testing.

### C. Large language Model

*1) Machine Learning:-:* The process of aligning large language models (LLMs) with human values and preferences is analyzed in the machine learning sector. In order to train process-supervised reward models (PRMs) and improve base policies via expert iteration, it involves building and using fine-grained datasets. The study evaluates how difficult it is for Reinforcement Learning from Human Feedback (RLHF) to match LLM actions with human desires and highlights how important it is to use process-supervision signals from PRMs to inform policy change. Moreover, it describes the two primary phases of expert iteration—policy improvement and distillation—as crucial elements in enhancing base policies via supervised fine-tuning grounded in superior alignment datasets.

*2) Deep Learning:-:* The process-supervised reward models (PRMs) are trained using large language models (LLMs) and fine-grained datasets in order to improve base policies via expert iteration. It highlights how process-supervision signals from PRMs can be used to direct policy reform by bringing LLM behaviors into line with human values and preferences. Through supervised fine-tuning based on high-quality alignment datasets, the work iteratively refines base policies, highlighting the complexities of Reinforcement Learning from Human Feedback (RLHF) in this context.

## VI. ALGORITHM VIEWPOINTS

### A. Overview

Our project focuses on creating a digital human twin, encapsulating the user's voice, shape, and attitude. The system leverages advanced AI models and modern development frameworks to achieve a seamless and realistic interaction experience. The core components include a Large Language Model (LLM) for text generation, Real-Time Voice Cloning (RTVC) for voice synthesis, a Flutter-based mobile application, and a Flask server for backend operations. This section details the algorithms and methodologies employed in each component.

### B. Text Generation with Large Language Models

The text generation module is powered by a state-of-the-art Large Language Model (LLM). We utilize a pre-trained transformer-based model (e.g., GPT-4), which has been fine-tuned for generating contextually relevant and coherent responses based on user input.

- **Pre-training and Fine-tuning**: The LLM is pre-trained on diverse datasets and fine-tuned with domain-specific data to align with the desired personality and attitude of the digital twin.
- **Inference**: During interaction, user input is processed, and the LLM generates text that mimics the user's style and tone. The model employs attention mechanisms to maintain context and coherence over long conversations.

### C. Voice Cloning with Real-Time Voice Cloning (RTVC)

The voice cloning component utilizes Real-Time Voice Cloning (RTVC) to synthesize the user's voice.

- **Speaker Embedding**: The RTVC system begins by extracting speaker embeddings from a sample of the user's voice. This is achieved using a pre-trained speaker encoder that generates a fixed-size vector representing the user's unique vocal characteristics.
- **Synthesis**: The embedding is then fed into a synthesizer model, which generates mel-spectrograms corresponding to the generated text. This step ensures the voice output matches the user's vocal attributes.
- **Vocoder**: Finally, a vocoder (e.g., WaveGlow or Griffin-Lim) converts the mel-spectrograms into waveform audio, producing natural-sounding speech that closely resembles the user's voice.

### D. Backend Server with Flask

The backend server, implemented using Flask, handles the core logic and integration of various components.

- **API Endpoints**: Flask provides RESTful API endpoints for processing user requests, managing sessions, and orchestrating the interactions between the LLM, RTVC, and the mobile app.

- **Session Management**: The server maintains user sessions, ensuring continuity and context preservation across interactions.
- **Scalability and Performance**: The server is optimized for low latency and high throughput, ensuring smooth real-time communication. Load balancing and caching mechanisms are employed to enhance performance under high user loads.

### E. Integration and Workflow

The system's workflow is as follows:

1) **User Interaction**: The user interacts with the mobile application through voice or text.
2) **Request Processing**: The input is sent to the Flask server, which processes the request and forwards it to the LLM for text generation.
3) **Response Generation**: The LLM generates the response text, which is sent back to the server.
4) **Voice Synthesis**: The text is then passed to the RTVC system to generate the corresponding voice output.
5) **Response Delivery**: The synthesized voice is sent back to the mobile app, providing a seamless conversational experience to the user.

This integrated approach ensures that the digital human twin accurately reflects the user's voice, style, and personality, creating a highly personalized interaction experience.

## VII. RESULTS

### A. Voice Cloning Performance

Several criteria were used to evaluate the voice cloning module in order to determine the accuracy and quality of the synthesized speech. The following are the outcomes:

- **Speaker Similarity**: On a scale of 0 to 1, the average similarity score between the original speaker's voice and the cloned voice was 0.87, suggesting a high degree of likeness.
- **Naturalness**: The cloned voice seemed natural to the majority of listeners, as evidenced by the Mean Opinion Score (MOS) of 3.8 out of 5.
- **Generalization**: The model maintained a similarity score of 0.82 on the test set, demonstrating strong performance with speakers that had not yet been encountered.

### B. Avatar Generation Quality

Both qualitative and quantitative indicators were used to evaluate the quality of the 3D avatars that were generated:

- **Visual Fidelity**: A panel of specialists assessed the created avatars visually, giving them an average score of 4.3 out of 5 for realism and intricacy.
- **Pose Accuracy**: There was a 2.5 mm margin of error in the 3D models' body joint locations when compared to ground truth data.
- **Facial Expression Matching**: Compared to a dataset of labeled expressions, the system recreated facial expressions with a 92



Fig. 3. Generated 3D avatar

### C. Large Language Model (LLM) Effectiveness

The capacity of the LLM to produce coherent and contextually relevant responses served as the benchmark for evaluating its performance:

- **Response Relevance**: Human reviewers assessed 91
- **Coherence over Long Conversations**: The LLM received a coherence score of 4.5 out of 5 for maintaining cohesive context over an average of 10 exchanges.
- **User Satisfaction**: The quality of interactions with the digital twin was rated an average of 4.2 out of 5 in user satisfaction surveys.

### D. Mobile Application Usability

Through user testing and feedback, the mobile application's usability was assessed:

- **User Interface (UI) Design**: For intuitiveness and usability, the UI received a rating of 4.6 out of 5.
- **Performance and Responsiveness**: The program offered a flawless user experience with an average response time of 250 ms.
- **Feature Integration**: With a functioning rating of 4.4 out of 5, users praised the incorporation of text and voice interactions.

## VIII. CONCLUSION

In this research paper, the development of the digital twin of human personalized virtual assistant avatar, has significant implications for the field of Human-Computer Interaction HCI. By creating a realistic 3D avatar that integrates facial expressions, lip sync, and full-body motion, we have redefined the boundaries of engagement and interaction.
Through the integration of various deep learning models, including speech recognition, facial expressions dataset training

and testing, Text-To-Speech model, reinforcement learning model, and accurate representation of the human body, we have successfully achieved a personalized and interactive 3D avatar. This avatar has the ability to generate responses based on the user input, creating a dynamic conversational interface that enhances communication platform AI agents.

Our project not only demonstrates the practical application of advanced technologies such as deep learning and natural language processing but also highlights the potential of large language models in understanding user input. The LLM tasks such as summarizing text, generating human-like expressions, and development in fields like customer service, and education. Furthermore, our exploration of transfer training techniques for Text-to-Speech TTS models has improved voice quality and synthesis, achieving a more natural and engaging user experience. By advancing algorithmic improvements in areas such as long sentence processing, vocabulary pronunciation, and noise reduction, we have enhanced the quality of voice synthesis.

Our application has the potential to revolutionize Human-Computer Interaction and user experiences across various domains. The user-friendly mobile application interface and the ability of the avatar to learn and evolve to be more replica of the user further enhance the Ai-companion.

## REFERENCES

[1] G. Eason, B. Noble, and I. N. Sneddon, "On certain integrals of Lipschitz-Hankel type involving products of Bessel functions," Phil. Trans. Roy. Soc. London, vol. A247, pp. 529–551, April 1955.

[2] J. Clerk Maxwell, A Treatise on Electricity and Magnetism, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.

[3] I. S. Jacobs and C. P. Bean, "Fine particles, thin films and exchange anisotropy," in Magnetism, vol. III, G. T. Rado and H. Suhl, Eds. New York: Academic, 1963, pp. 271–350.

[4] K. Elissa, "Title of paper if known," unpublished.

[5] R. Nicole, "Title of paper with only first word capitalized," J. Name Stand. Abbrev., in press.

[6] Y. Yorozu, M. Hirano, K. Oka, and Y. Tagawa, "Electron spectroscopy studies on magneto-optical media and plastic substrate interface," IEEE Transl. J. Magn. Japan, vol. 2, pp. 740–741, August 1987 [Digests 9th Annual Conf. Magnetics Japan, p. 301, 1982].

[7] M. Young, The Technical Writer's Handbook. Mill Valley, CA: University Science, 1989.