# SimplyMime: A Control at Our Fingertips

**Sibi Chakkaravarthy Sethuraman**
Centre of Excellence, Artificial Intelligence & Robotics (AIR),
School of Computer Science and Engineering
VIT-AP University, India
sb.sibi@gmail.com

**Gaurav Reddy Tadkapally**
Centre of Excellence, Artificial Intelligence & Robotics (AIR),
School of Computer Science and Engineering
VIT-AP University, India
gauravreddy008@gmail.com

**Athresh Kiran**
Senior Software Developer
Parallel Reality: AI-based Health Tech
United Kingdom
athresh.kiran@gmail.com

**Saraju P. Mohanty**
Department of Computer Science and Engineering,
University of North Texas,
TX 76207, USA,
saraju.mohanty@unt.edu

**Anitha Subramanian**
Centre of Excellence, Artificial Intelligence & Robotics (AIR),
School of Computer Science and Engineering
VIT-AP University, India
anithachubbu@gmail.com

April 10, 2023

## Abstract

The utilization of consumer electronics, such as televisions, set-top boxes, home theaters, and air conditioners, has become increasingly prevalent in modern society as technology continues to evolve. As new devices enter our homes each year, the accumulation of multiple infrared remote controls to operate them not only results in a waste of energy and resources, but also creates a cumbersome and cluttered environment for the user. This paper presents a novel system, named SimplyMime, which aims to eliminate the need for multiple remote controls for consumer electronics and provide the user with intuitive control without the need for additional devices.

SimplyMime leverages a dynamic hand gesture recognition architecture, incorporating Artificial Intelligence and Human-Computer Interaction, to create a sophisticated system that enables users to interact with a vast majority of consumer electronics with ease. Additionally, SimplyMime has a security aspect where it can verify and authenticate the user utilising the palmprint, which ensures that only authorized users can control the devices. The performance of the proposed method for detecting and recognizing gestures in a stream of motion was thoroughly tested and validated using multiple benchmark datasets, resulting in commendable accuracy levels.

One of the distinct advantages of the proposed method is its minimal computational power requirements, making it highly adaptable and reliable in a wide range of circumstances. The paper proposes incorporating this technology into all consumer electronic devices that currently require a secondary remote for operation, thus promoting a more efficient and sustainable living environment.

## I. Introduction

Smart electronics, such as televisions, air conditioners, speakers, and ceiling fans, have become ubiquitous in modern society. Thus driving the framework of smart cities and smart villages which are intended to operate optimally with limited resources while best utilizing the available resources to improve quality of life [1], [2]. With the proliferation of inexpensive and readily available devices, most households have multiple electronic devices that require remote controls for operation and interaction [3]. As the number of devices increases each year, so too does the number of remote controls needed to operate them. This traditional method of interaction, while widespread and commonly used, is not without its flaws. The use of multiple remote controls not only wastes resources and increases the use of plastic, but also makes it difficult to locate and operate the correct remote for a given device [4].

The first remote-control devices were introduced in the 1950s, but it wasn't until the 1970s that infrared-based remotes began to appear on the market [5]. Despite some advancements in human-computer interaction (HCI) technology, such as keyboard and mouse alternatives and devices that use Bluetooth or WiFi communication, these new approaches have not been able to fully replace traditional remote controls due to a lack of seamless functionality. However, recent advances in HCI technology, such as voice commands, mimics, and gestures used in devices such as tablets, smartphones, and smart homes, have shown that there is still potential for improvement in the field [6].

SimplyMime is a hand gesture-based control system for consumer electronics that addresses the shortcomings of traditional remote controls. Hand tracking is a natural and intuitive mechanism for identifying hand movements, and it has been studied for a long time. Skeleton-based hand tracking, due to its resistance to various background conditions, has proven to be a popular choice for this type of technology [7]. SimplyMime offers a lightweight solution that can track hands and provide real-time output, enabling it to be incorporated into a compact system that can be installed in any electronic appliance. Furthermore, traditional remote controls lack security features. To further enhance security, SimplyMime incorporates facial recognition and detection to authenticate and validate the user's identity. This added security measure ensures that only authorized users are able to control the devices, providing an additional layer of protection for users.
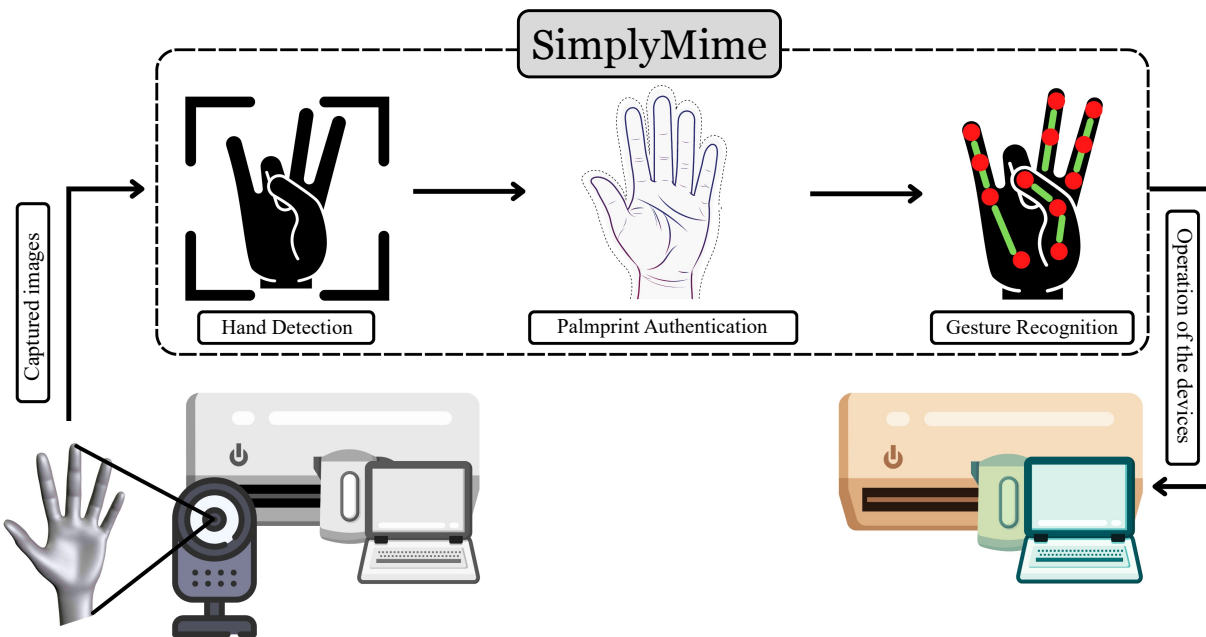


Figure 1. Diagrammatic representation of the architectural working of SimplyMime

SimplyMime harnesses the power of hand detection technology to recognize various user gestures, as outlined in Section IV. This allows for the execution of specific actions, such as controlling embedded devices, such as turning

them ON/OFF, adjusting speeds and volumes, among others. It offers an efficient and intuitive means of interacting with consumer electronics.

The main objectives of SimplyMime are the following:

- **User-centric Design**: As discussed in Section III, various technologies have been proposed to control different devices, but few have given consideration to the user's convenience. SimplyMime prioritizes ease of use for the end user.
- **Seamless Integration**: The SimplyMime system is designed to be integrated seamlessly into any device without altering its internal architecture, making it easy for users to migrate from traditional remote controls to this new technology.
- **Sustainable Solutions**: SimplyMime offers a one-stop solution for controlling multiple devices, reducing the accumulation of e-waste such as remote controls and other controllers.
- **Advanced Security**: SimplyMime includes a security authentication component that verifies and authenticates the user before granting access, providing an added layer of security, unlike traditional remote controls.

Rest of the paper is organized as follows: Section II outlines novel contributions of the current paper. Section III discusses about the existing related research. Section IV provides the system level architecture of the proposed SimplyMime and outlines the novel methods proposed for hand detection, hand landmark detection and Palm print authentication. Section V is used to validate the results of the proposed SimplyMime system and also provides a comparative analysis of SimplyMime against the state of the art solutions. Finally, the paper concludes in Section VI.

## II. Novel Contributions of the Current Paper

### 1) Problem Addressed

The traditional remote controls require line of sight and precise pointing, hindering the user's freedom of movement and reducing the overall comfort level. As discussed earlier, even the recent contributions fail to develop a solution which eliminated the requirement of additional device completely, and provides a better intuitive experience to the user, resulting in unreliable and frustrating user experiences.

### 2) Solution Proposed by SimplyMime

SimplyMime is a state-of-the-art system that employs advanced hand gesture recognition techniques, leveraging the latest advancements in Artificial Intelligence and Human-Computer Interaction, to create a highly efficient and intuitive control framework for consumer electronics. The system utilizes a camera to capture the user's hand gestures, which are then translated into commands for electronic devices. The dynamic architecture of SimplyMime enables it to adapt to various hand gestures and provide a seamless user experience. The architecture is equipped with a palm print authentication module that ensures the device is only accessible to authorized users.

The SimplyMime system is designed to eliminate the need for multiple remote controls, providing users with a streamlined and organized living space. By removing the need for additional devices, SimplyMime offers a more sustainable solution that saves energy and resources. Moreover, the intuitive nature of the system provides users with greater comfort and convenience, as they no longer have to search for and operate multiple remote controls. Additionally, the dynamic architecture of SimplyMime enables it to adapt to new gestures, ensuring that it can keep up with the changing needs of users.

### 3) Significance of the Proposed Solution

Our proposed solution, SimplyMime, offers several advantages over traditional remote control systems. First, it eliminates the need for multiple controllers, leading to a cleaner, more organized living space. Second, the system offers a more immersive and intuitive experience for users, making it easier to interact with electronic devices. Moreover, our pipeline incorporates a security feature that is often overlooked in previous solutions, by enabling verification and authentication of the user through the use of palmprint recognition. This added security measure ensures that only authorized users are able to control the devices, providing an additional layer of protection for users. We believe that SimplyMime's innovative approach to remote control systems has the potential to revolutionize the way users interact with their consumer electronics. Its minimal computational requirements make it highly adaptable and reliable in a wide range of circumstances, making it an ideal solution for households and businesses alike.

## III. Prior Related Research

Hand gesture recognition technology has been a subject of study for several years, with developments in different fields, including human-computer interaction, gaming, augmented reality, and assistive technology for those with disabilities. The hand, in particular, is used extensively for gesturing compared to other body parts because it is a natural means of human communication and therefore the most appropriate tool for HCI [7]. And evidently, controlling electronic devices with our hands is a natural and highly intuitive method of interaction. The use of hand gesture recognition to operate consumer electronics such as televisions and home appliances is also gaining traction [8]. This application area has the potential to revolutionise how we interact with electronic devices in our homes by enabling simple hand gestures to control them. However the efforts to develop a reliable and robust hand gesture recognition system for this application ares is still lacking.

Variability in how individuals make gestures is one of the greatest obstacles to hand gesture recognition. People can perform the same gesture differently based on their age, gender, cultural background, and physical capabilities. Researchers have developed various techniques for capturing and analysing hand gestures, including computer vision, machine learning, and sensor-based approaches [9].

### A. Sensor-based Methods

The majority of existing methods can be divided into two categories, with sensor-based systems constituting the first. Utilizing specialised sensors to capture the movement and orientation of the hand, sensor-based methods for hand gesture recognition capture the hand's motion and the orientation. These sensors may include accelerometers, gyroscopes, magnetometers, and other types of sensors that can measure the hand's movement and position [10]. Sensor-based methods have a number of advantages over other methods, including the ability to capture precise information about the hand's movement and position and the capacity to function in environments with poor lighting or visibility. In a recent study, the authors employed the utilization of inertial sensors, specifically accelerometers and gyroscopes, in conjunction with a multilayered perceptron classifier to recognize hand gestures [11]. Adopting a similar methodology, other researchers have developed a finger-worn ring device that captures acceleration data, paired with an LSTM model for the classification of hand gestures [12]–[14]. Likewise, Inviz device incorporates textile-based flexible capacitive sensor arrays for motion detection [15]. This approach has been utilized to develop systems that cater to individuals with paralysis, paresis, weakness, and restricted range of motion. Other notable developments include the WristCam [16], a wrist-worn camera sensor that estimates and recognizes hand trajectories, as well as the EchoFlex [17], a wearable ultrasonic device that tracks hand movements through a novel method. The research makes use of the muscle flexor data collected by the proposed device to develop an algorithm for gesture recognition that combines image processing and neural networks.

Sensor-based systems provide precise information about the hand's movement and position and can function in environments with poor lighting or visibility. However, the requirement for the user to wear additional devices restricts the versatility and scalability of these systems. In many instances, faulty sensor calibration can result in errors in gesture recognition, resulting in abysmal performance.

### B. Vision-based Methods

The integration of computer vision and human-computer interaction has enabled the development of innovative systems aimed at addressing limitations and providing users with a more immersive and efficient experience when interacting with machines. One of the initial approaches in the development of vision-based sensors was the utilization of a colored glove [18]. This method required the user to wear a colored glove, which enabled the system to employ a nearest-neighbor approach for tracking hand movements at interactive rates. Several studies that employed the use of the Kinect sensor utilized color and image depth data to establish a hand model for the analysis of tracked hand gestures. The hand gesture tracking was accomplished through the implementation of a Kalman filter [19]. However, it was observed that the gesture recognition accuracy was relatively low. Similar research utilized RGB-D cameras to extract hand location data through the use of in-depth skeleton-joint information from images [20], [21]. Furthermore, similar to the architecture of SimplyMime, a few works employed neural networks to infer real-time hand landmarks, which were used to establish a skeletal structure of the hand gesture [22], [23].

Computer vision-based systems are non-intrusive and more versatile than sensor-based systems, but they can be affected by poor lighting or visibility. Moreover, recognizing gestures from complex backgrounds is still a challenging task. As a reason, SimplyMime focuses on developing robust and reliable hand gesture recognition systems that can overcome these challenges and enable natural and intuitive interactions with electronic devices.

Table I
COMPARISON OF EXISTING HAND GESTURE RECOGNITION MODELS

| Research | Methodology Employed | Findings and Outcome | Requirement of additional device | Requirement of Calibration | Security Module |
|---|---|---|---|---|---|
| Teachasrisaksakul et al. (2018) [11] | Inertial Sensors | - Achieved an accuracy of 98.33%.<br>- However, the performance will be impacted by external factors, such as the user's body movements or the environmental conditions | Accelerometer and Gyroscope | Yes | None |
| TinyDL (2021) [12] | Inertial Sensors | - Utilized a built-in LSTM model that leveraged data from a finger-worn ring.<br>- The performance will be impacted by factors such as sensor placement, variations among users, and changes in hand gesture execution over time. | Finger-worn device | Yes | None |
| Inviz (2015) [15] | Textile-based capacitive arrays | - Relies on textile-based capacitive arrays built into clothing.<br>- The production of wearable textile capacitive sensor arrays could increase the cost and complexity of the system. | Wearable textile capacitive sensor Arrays | Yes | None |
| EchoFlex (2017) [17] | Ultrasound Imaging | - Utilizes ultrasound sensors to capture hand movement data, which is then analyzed using ML algorithms.<br>- Susceptible to interference from nearby objects or environments with high acoustic noise | Ultrasonographic Device | Yes | None |
| Wang et al. (2009) [18] | Image Processing | - Employs colour segmentation and a tracking algorithm is applied to estimate the hand pose and motion.<br>- Recognition accuracy would be affected by factors such as lighting conditions, colour variations, and occlusions. | Coloured Glove | No | None |
| Feng et al. (2014) [19] | Kinect Sensors | - Demonstrates the feasibility of combining various types of sensors for human motion tracking.<br>- Requires a high-performance computing device to achieve real-time performance | Microsoft Kinect Sensor | Yes | None |
| SHREC (2021) [23] | Skeleton-based hand gesture Recognition | - Achieved high accuracy in recognizing hand gestures in various real-world scenarios | None | No | None |
| **SimplyMime (Current Paper)** | CNN based Skeletal Pose Estimation | - Delivers intuitive control without additional devices<br>- Incorporates palmprint authentication to verify and authenticate users | None | No | Palmprint based Authentication |

## IV. SimplyMime: The Proposed Smart Remote Control

The proliferation of consumer electronics has resulted in the widespread use of traditional remote controls as the primary means of interaction. However, in order to effectively replace such a firmly established technology, an alternative solution must not only be robust, precise, and intuitive, but also possess the added benefits of user-friendliness, compatibility with older devices, and scalability [8]. Our proposed work, SimplyMime, addresses the shortcomings of existing solutions, while also maintaining the immersive experience that traditional remote controls are capable of delivering. The dynamic hand gesture recognition module, represents the most advanced, effective, and ideal replacement for traditional remote controls, providing a unique blend of state-of-the-art performance and efficiency.

The underlying architecture of the system incorporates a hand landmark assignment method, which is utilized to identify and localize key points across the hand, such as finger tips, knuckles, and wrists. These landmarks, assigned by the backend model, form the basis for the gesture recognition algorithm, which is able to identify and classify gestures based on the skeletal structure generated. The system is designed to operate in two distinct modules, which are discussed in greater detail in the following sections. The first module, the hand detection model, is responsible for identifying and isolating the human hand within an image. The second module, the gesture recognition model, processes the detected hand to generate a skeletal structure of the gesture and subsequently recognizes it. The suggested approach constitutes a noteworthy progression in the domain of hand gesture detection, showing potential as a viable substitute for conventional human-computer interaction techniques.

5

## A. Hand Detection Model

The foundation of an effective hand gesture identification model is the accurate detection and isolation of the hand from the given image. To achieve this, SimplyMime employs state-of-the-art neural network technology for localizing the coordinates of the palm. Our model specifically utilizes the Single Shot Detector (SSD) architecture [24]. While conventional RCNN models have been widely used for object localization, they require significant computational power, making them less reliable for real-time applications [25]. In contrast, our SSD algorithm achieves superior real-time performance by eliminating the need for bounding box proposals and the subsequent feature re-sampling stage. This allows for a more efficient and effective approach to hand detection, which is a crucial component of an accurate gesture identification system.



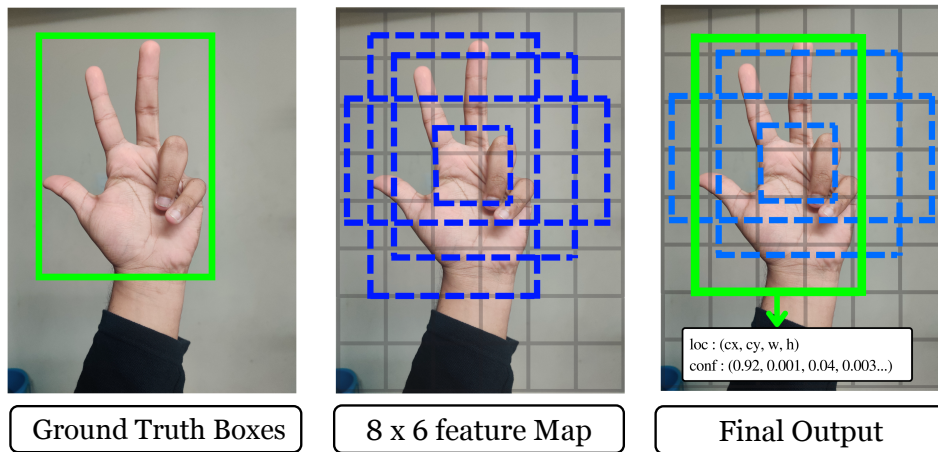| Ground Truth Boxes | 8 x 6 feature Map | Final Output |

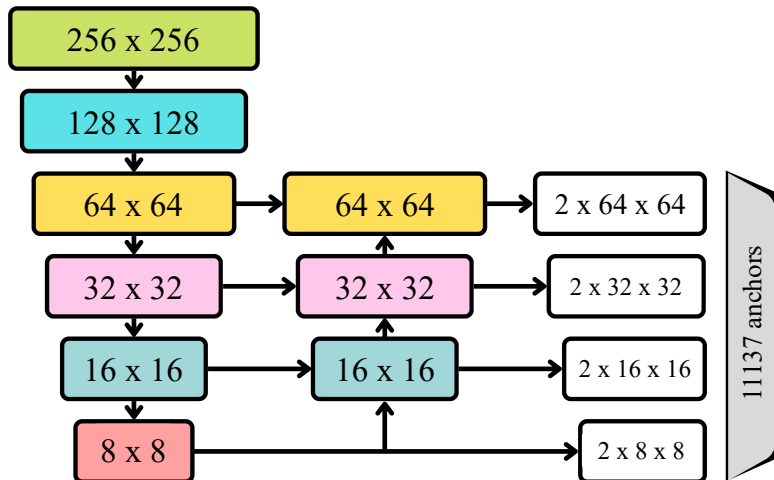Figure 2. Working overview of the Hand detection model



Figure 3. Comprehensive view of the Gesture recognition module's architecture

The hand detection model employed in SimplyMime utilizes a unique approach in which it is initially trained on human palms rather than the entire hand. This approach allows the network to more easily perceive patterns and

generate bounding boxes, as the palms and fists are more rigid objects when compared to a human hand with articulated fingers. In the task of detecting human hands, the last layer of the model creates substantially smaller anchor boxes, as human hands are smaller objects and thus require smaller anchor boxes. The architectural network of SimplyMime's detection model is illustrated in Figure 3. The feature extraction network takes an input RGB image of 224x224px and is followed by a series of convolutional layers, referred to as ConvBlocks, which serve as the bottleneck for the higher abstraction level layers. Essentially, the image is passed through 5 single and 6 double ConvBlocks. The introduced ConvBlocks consist of a series of convolutional layers, followed by a depth-wise and a point-wise convolutional layer. The detailed architecture of our Convblocks is depicted in the Figure 4. Furthermore, our network outperformed a popular light-weight model, MobileNetV2-SSD [26], in terms of both accuracy and inference speed.
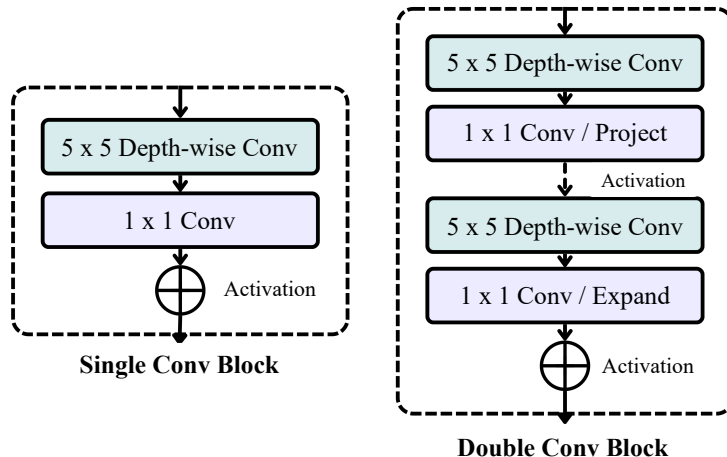


Figure 4. Inner structure of the devised Convolution blocks

The detection architecture of SimplyMime is designed to be highly robust and accurate, and this is achieved through the use of a diverse set of training data. To ensure that the model is able to generalize well to a wide range of scenarios, it is initially trained on three different sources of data. The first of these is an in-the-wild dataset, which comprises a diverse set of images captured in various geographical locations and under different lighting conditions, allowing the model to learn to detect hands under a wide variety of conditions. The second data source is an in-house dataset, which is specifically designed to cover all possible hand angles in a controlled environment. This dataset allows the model to learn to detect hands under consistent conditions, and the combination of these two datasets allows for a well-represented and diversified dataset. Finally, the model is further trained on a synthetic dataset to ensure that it is able to detect hands under a wide range of angles and in a variety of environments, further improving its robustness and accuracy.

## B. Hand Landmark Detection Model

Having successfully localized the hand in the image, the next step in SimplyMime's pipeline is to extract key-points from the isolated hand region. For this purpose, we employ a modified version of the Convolutional Pose Machines (CPM) network, which has been extensively used in the field of human pose estimation [27]. The CPMs provide a confidence map for each keypoint, represented as a Gaussian centered at the true position. These maps are generated based on the size of the input image patch, and the final location of each keypoint is determined by identifying the peak in the corresponding confidence map. By assigning 21 landmarks across different points across the palm, the model achieves to generate a skeletal structure of the palm.

## Dynamic Gesture Detection Algorithm

After extracting the keypoints from the hand region, they are transferred to the gesture detection engine for further analysis. These keypoints, which consist of 21 points, represent distinct areas of the hand and are illustrated in Figure 5. They serve as the basis for identifying the gesture that the user intends to create. This processed information is then utilized to control the SimplyMime-enabled devices. Our algorithm incorporates a sophisticated technique that assigns a FingerState value to each digit. An open finger is designated as 1, while a folded finger is denoted as
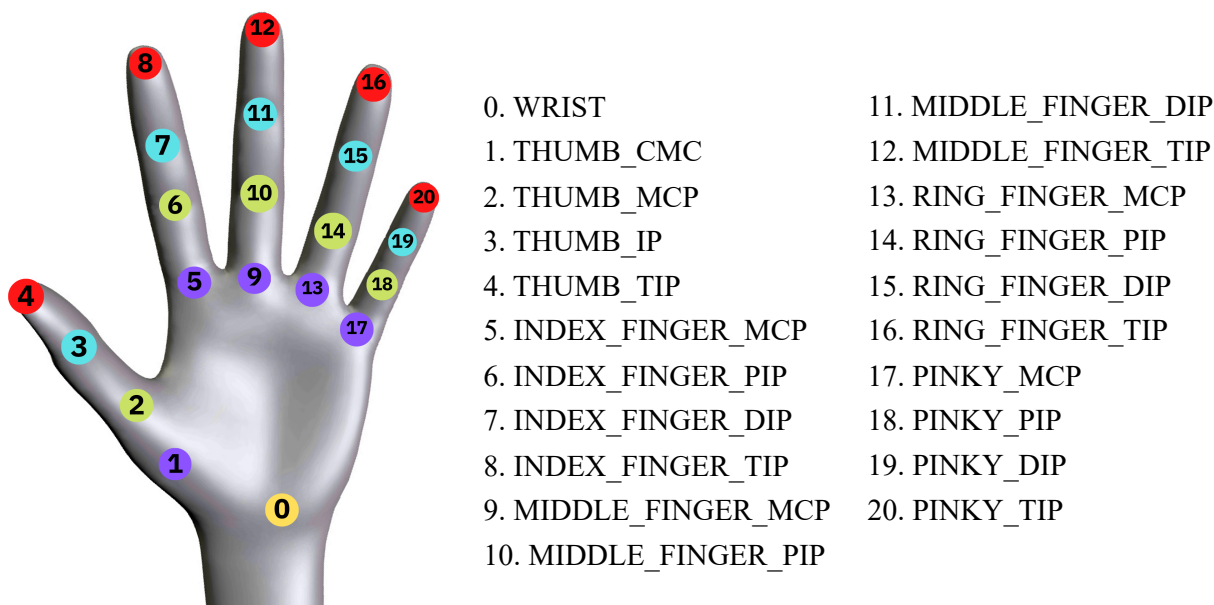
7

0. WRIST
1. THUMB_CMC
2. THUMB_MCP
3. THUMB_IP
4. THUMB_TIP
5. INDEX_FINGER_MCP
6. INDEX_FINGER_PIP
7. INDEX_FINGER_DIP
8. INDEX_FINGER_TIP
9. MIDDLE_FINGER_MCP
10. MIDDLE_FINGER_PIP
11. MIDDLE_FINGER_DIP
12. MIDDLE_FINGER_TIP
13. RING_FINGER_MCP
14. RING_FINGER_PIP
15. RING_FINGER_DIP
16. RING_FINGER_TIP
17. PINKY_MCP
18. PINKY_PIP
19. PINKY_DIP
20. PINKY_TIP

Figure 5. Key-point indexes of all hand landmarks

0. To determine the FingerState, we leveraged the Y-axis coordinates of the metacarpophalangeal (MCP) joint and the fingertip. Nevertheless, the thumb's unique location and alignment necessitated calculating its slope to ascertain its FingerState. This approach results in more precise and accurate gesture recognition and enhances the system's overall performance.

The concluding stage of the process involves the recognition of gestures by analyzing the extracted keypoints from the hand region. Through the use of a posture array, each gesture can be accurately identified and differentiated from one another. These posture arrays, each unique to a particular gesture, can be designated as trigger functions to regulate a diverse range of consumer electronics. To add a dynamic element to the algorithm, the metacarpophalangeal joint (MCP) of the middle finger is employed as the focal point to align and center the camera, ensuring that the palm remains in view. Moreover, the midpoint between the tips of the thumb and index finger is used to locate the cursor, if necessary, for a specific device. Figure 6 demonstrates several images and their corresponding posture arrays.

## C. Palm Print Authentication Module

The protection of security is of paramount importance in hand gesture-based systems utilized for controlling consumer electronics, as these systems are susceptible to unauthorized access and control [28]. Without robust security mechanisms, these systems can be easily compromised, leading to data breaches, unauthorized access to sensitive information, and potential damage to devices. Therefore, it is essential to incorporate advanced security measures during the design and development of hand gesture-based systems to mitigate these risks. The integration of strong security measures is critical for ensuring the reliability and integrity of hand gesture-based systems for consumer electronics.

SimplyMime employs advanced biometric authentication methods, specifically, PalmPrint identification. To achieve this, we utilize a Siamese Network architecture, a common approach used in similarity recognition tasks like facial recognition [29]. Our network includes two parallel feature extraction blocks that process two distinct palm images, producing a 4096-dimensional tensor for each image. By measuring the dissimilarity between the two sets of embeddings utilizing the euclidean distance function, the similarity between the two palms is verified. A pre-determined threshold is established to establish the authenticity of the user. Access to the system is restricted if the dissimilarity measure exceeds the threshold, ensuring that only authorized individuals have access to the system.
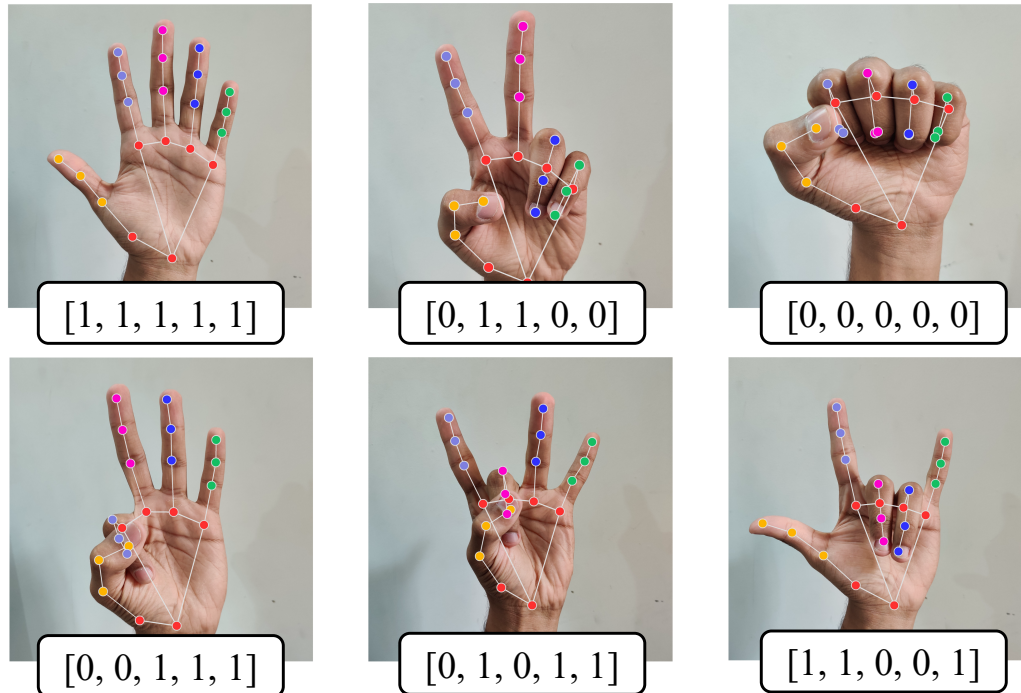
Figure 6. Working results and subsequent arrays generated by the landmark detection model

This multi-factor approach ensures that only authorized individuals have access to the system, protecting against unauthorized attempts to control the consumer electronics. Furthermore, implementing these security measures guarantees the protection of personal and sensitive information, enhancing the system's reliability and integrity.

During the training process, the network is presented with three input images, including an anchor image, a positive image, and a negative image. The anchor and positive images correspond to the palm print of the same individual, while the negative image represents the palm print of a different individual. The network is trained for 100 epochs using the Triplet Loss function, which is optimized using Adam optimizer [30]. The XceptionNet [31] serves as the base model, acting as the encoder. The output of the training process is a set of two distances, which are input into the Triplet Loss function [32], as shown in Equation 1:

$$TripleLoss = \sum_{N}^{i}[\|f(x_i^a) - f(x_i^p)\|^2$$
$$-\|f(x_i^a) - f(x_i^n)\|^2 + \alpha] \tag{1}$$

Figure 7 illustrates the architectural functioning of the SimplyMime's authentication module.

In the Equation 1, the function $f(x)$ maps each input image to a 4096-dimensional embedding, represented by a tensor. The input images, denoted as $a$, $p$, and $n$, respectively correspond to the anchor, positive, and negative samples used in the triplet loss function. The margin parameter $\alpha$ is used to control the relative distance between the positive and negative pairs, with the goal of maximizing the separation between them.

## D. Hardware Implementation and Setup

The hardware setup of SimplyMime is a crucial aspect of its overall design and functionality. The camera is mounted on a cuboidal cardboard structure that is equipped with a motor on its right side, which enables control over the Y-axis movement. The cardboard structure is further mounted on a CD disk which is connected to another motor, allowing for control over the X-axis movement. These motors are connected to an On-board microcontroller, which facilitates communication between the hardware components and the software algorithms. Figure 9 provide a visual representation of the hardware setup, while the circuit diagram in the Figure 8 illustrates the connections between the various components. The hardware setup is designed to be compact, portable, and easy to set up, enabling users to easily control consumer electronics with hand gestures.
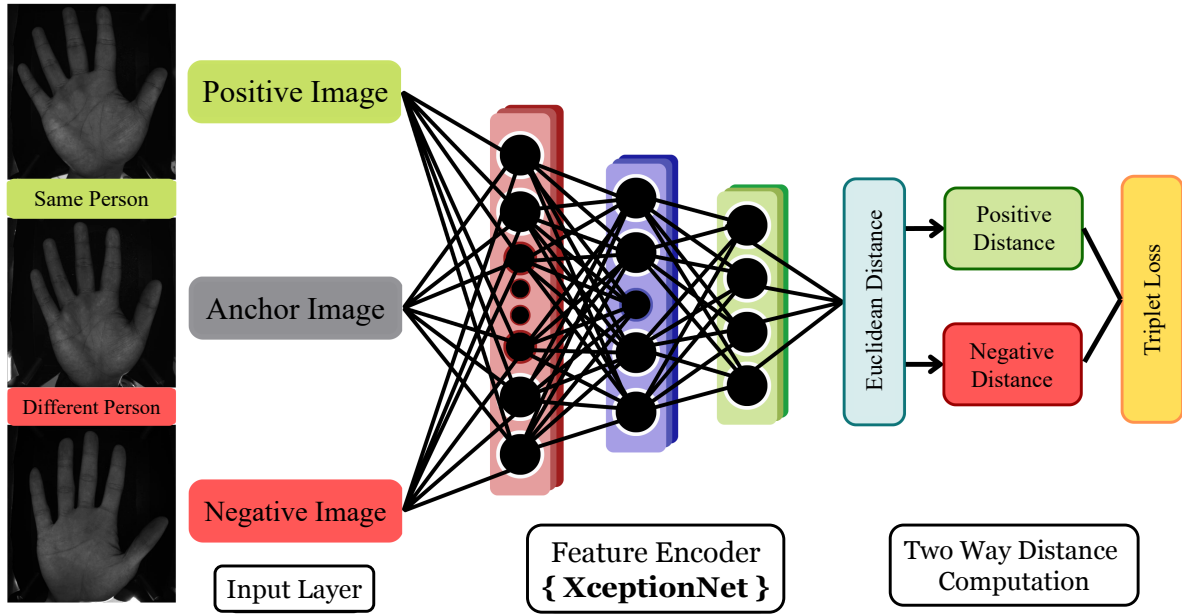
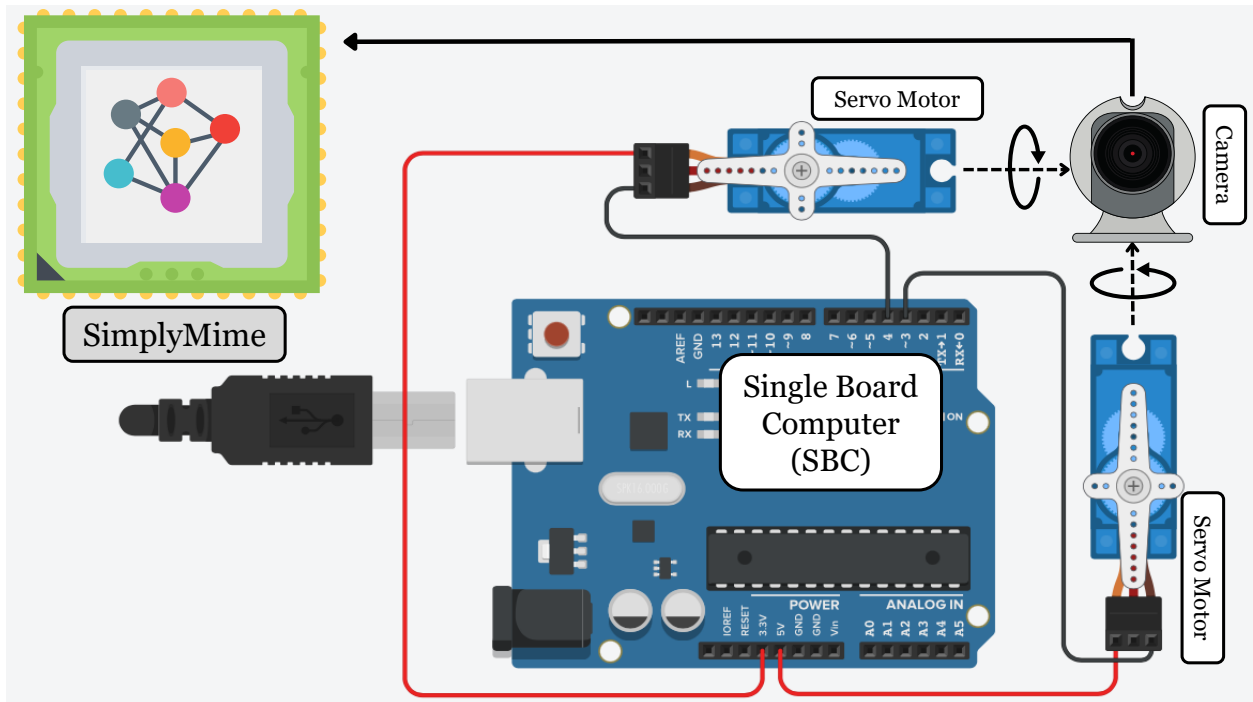Figure 7. Working pipeline of the Siamese network based authentication model



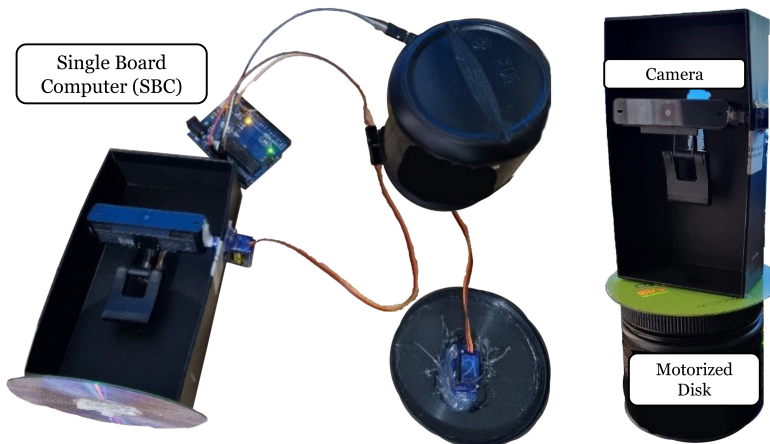Figure 8. Circuit implementation of SimplyMime

Figure 9.  Hardware setup of SimplyMime

## V. Experimental Results

In order to thoroughly evaluate the performance and accuracy of our hand gesture recognition model, we conducted extensive testing on two benchmark datasets. The first dataset was utilized to assess the model's ability to accurately detect gestures in images, resulting in an impressive 96.16% accuracy [33]. The second dataset [34] was utilized to evaluate the model's detection rate, yielding a accuracy of 87.37%. Additionally, we also evaluated our Siamese network-based palmprint authentication system using the CASIA dataset [35], resulting in an accuracy rate of over 90%. These results demonstrate the effectiveness and robustness of our proposed model in identifying and authenticating hand gestures for controlling consumer electronics. Figures 10(a), 10(b), and 10(c) set our few samples from all the utilised benchmark datasets.
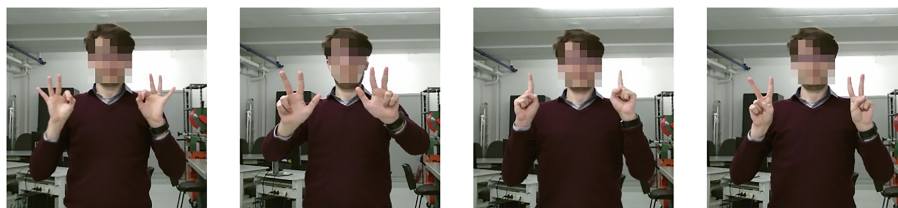
### A. HANDS Dataset

The HANDS dataset comprises a substantial corpus of samples collected from 5 diverse participants, consisting of 3 male and 2 female individuals [33]. The participants were instructed to perform 16 pre-defined gestures, comprising of 12 single-handed gestures performed with both arms, and 4 double-handed gestures. Each gesture was captured in 150 RGB frames, resulting in a total of 11,250 images across all participants. Additionally, the gestures were captured from varying distances, resulting in a diverse range of depths and variations in the hand poses.

Upon evaluation of the dataset, our proposed model achieved an overall accuracy of 96.16%. The evaluation metric, the recall in particular, is extremely crucial for evaluating SimplyMime as they provide an understanding of how well the model is able to correctly identify hand gestures among the input data, and balance the trade-off between false positives and false negatives. These metrics aid in evaluating the overall performance of the model and the ability of the model to generalize to unseen data. The Recall is used compute the proportion of true positive predictions among all actual positive instances in the data. In the context of SimplyMime, a high recall score indicates that the model is able to identify a high proportion of true hand gestures among all gestures present in the input data. These performances are further highlighted in Figure 11 and Table II, where the gesture-specific results are also illustrated.
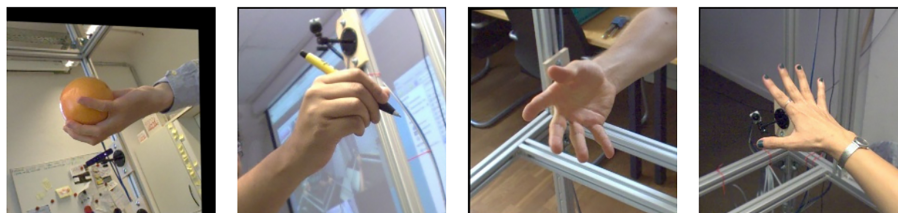
Furthermore, in order to evaluate the effectiveness of our proposed system, we conducted a comparative analysis against other notable models in the field. One such study utilized Generative Adversarial Networks for hand gesture detection [36], while another employed a similar pipeline to ours and leveraged the popular MobileNetV2 architecture as the baseline model [37]. Despite the high accuracy results achieved by these models, they required significant computational power. In contrast, our proposed system, SimplyMime, achieved comparable performance and precision while requiring significantly less computational resources as shown in Table III.
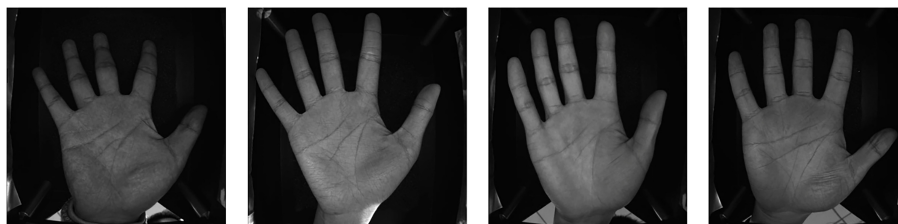
### B. FreiHand Dataset

A robust hand gesture-based control system requires a strong hand detector as its backbone. The hand detector's primary function is to accurately localize the hand region within an image, which is crucial for the subsequent gesture

(a) HANDS Dataset [33]



(b) FreiHand Dataset [34]



(c) CASIA Dataset [35]

Figure 10. Benchmark datasets utilized in SimplyMime Experimentation
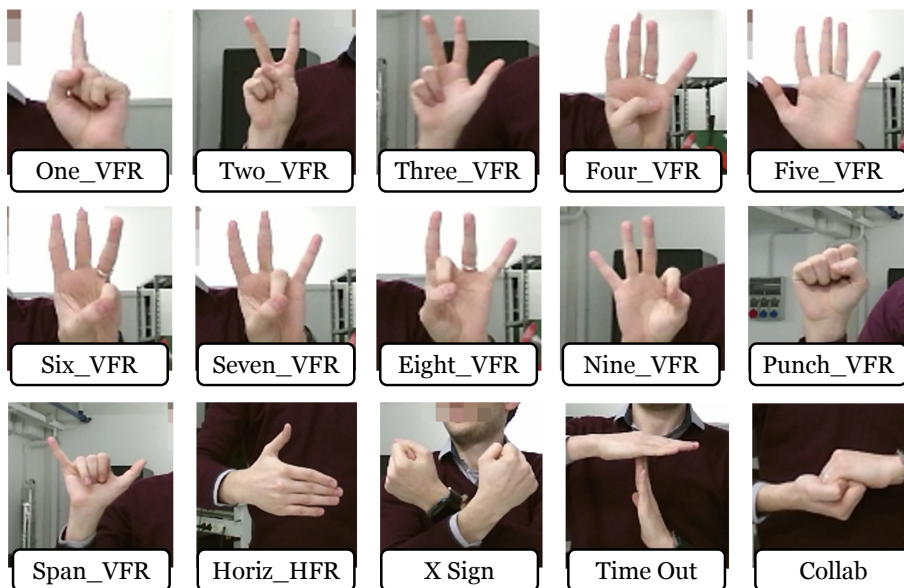


Figure 11. Results of various Gestures classified and labeled

Table II
RESULTS OF OUR HAND GESTURE RECOGNITION MODEL ON THE HANDS DATASET [33]

| Gesture name | Total frames | Accurately predicted frames | Falsely predicted frames | Accuracy % | Error % | Recall |
|---|---|---|---|---|---|---|
| Collab | 750 | 670 | 80 | 89.33 | 10.67 | 0.89 |
| TimeOut | 750 | 686 | 64 | 91.46 | 8.53 | 0.91 |
| XSign | 750 | 704 | 46 | 93.86 | 6.13 | 0.94 |
| Eight_VRF | 750 | 708 | 42 | 94.40 | 5.60 | 0.94 |
| Seven_VRF | 750 | 714 | 36 | 95.20 | 4.80 | 0.95 |
| Eight_VRF | 750 | 718 | 32 | 95.73 | 4.26 | 0.96 |
| Horiz_HRF | 750 | 727 | 23 | 96.93 | 3.06 | 0.97 |
| Span_VRF | 750 | 728 | 22 | 97.06 | 2.93 | 0.97 |
| Six_VRF | 750 | 732 | 18 | 97.60 | 2.40 | 0.98 |
| Five_VRF | 750 | 733 | 17 | 97.73 | 2.26 | 0.98 |
| Four_VRF | 750 | 736 | 14 | 98.13 | 1.86 | 0.98 |
| Three_VRF | 750 | 738 | 12 | 98.40 | 1.60 | 0.98 |
| Two_VRF | 750 | 739 | 11 | 98.53 | 1.46 | 0.99 |
| One_VRF | 750 | 741 | 9 | 98.80 | 1.20 | 0.99 |
| Punch_VRF | 750 | 742 | 8 | 98.93 | 1.06 | 0.99 |
| **Total** | **11250** | **10816** | **434** | **96.16** | **3.85** | **0.9613** |

Table III
COMPARISON OF OUR HAND GESTURE RECOGNITION MODEL AGAINST EXISTING SOLUTION

| Research | Architecture Used | Accuracy |
|---|---|---|
| Feng et al. (2022) [36] | GANs | 96% |
| Dang et al. (2022) [37] | MobileNetV2 | 94% |
| **SimplyMime** | CNN based skeletal pose estimation | 96% |

recognition stage. Therefore, to evaluate the performance of our hand gesture recognition model, SimplyMime, we also conducted evaluations to measure the capacity of our hand detector. The results of these evaluations provide insight into the model's ability to accurately detect and localize the hand region, which is crucial for the overall performance of the system. Additionally, by comparing the results of our hand detector with those of other models, we can gain a better understanding of the performance of our system in relation to the state-of-the-art.

The FreiHand dataset is a benchmark dataset specifically designed to evaluate the performance of hand detection and hand pose estimation models [34]. The dataset contains over 32,560 frames of synchronized RGB and depth data, captured using Microsoft Kinect v2 sensors, collected by 32 people. This dataset is considered to be one of the most challenging datasets for hand detection and pose estimation, as it contains a wide range of hand poses and motion, captured under various lighting conditions and backgrounds. To evaluate the performance of SimplyMime, we used the Freihand dataset to test our model's ability to detect hands and estimate hand poses. The dataset was particularly useful in evaluating the model's performance under challenging conditions such as low resolution, low contrast, and occlusions [34]. Our model was able to achieve a high level of accuracy in detecting hands and estimating hand poses, even under these challenging conditions. Our system's results on the dataset are depicted in Table IV.

Table IV
MODEL'S PERFORMANCE ON FRIEHAND DATASET [34]

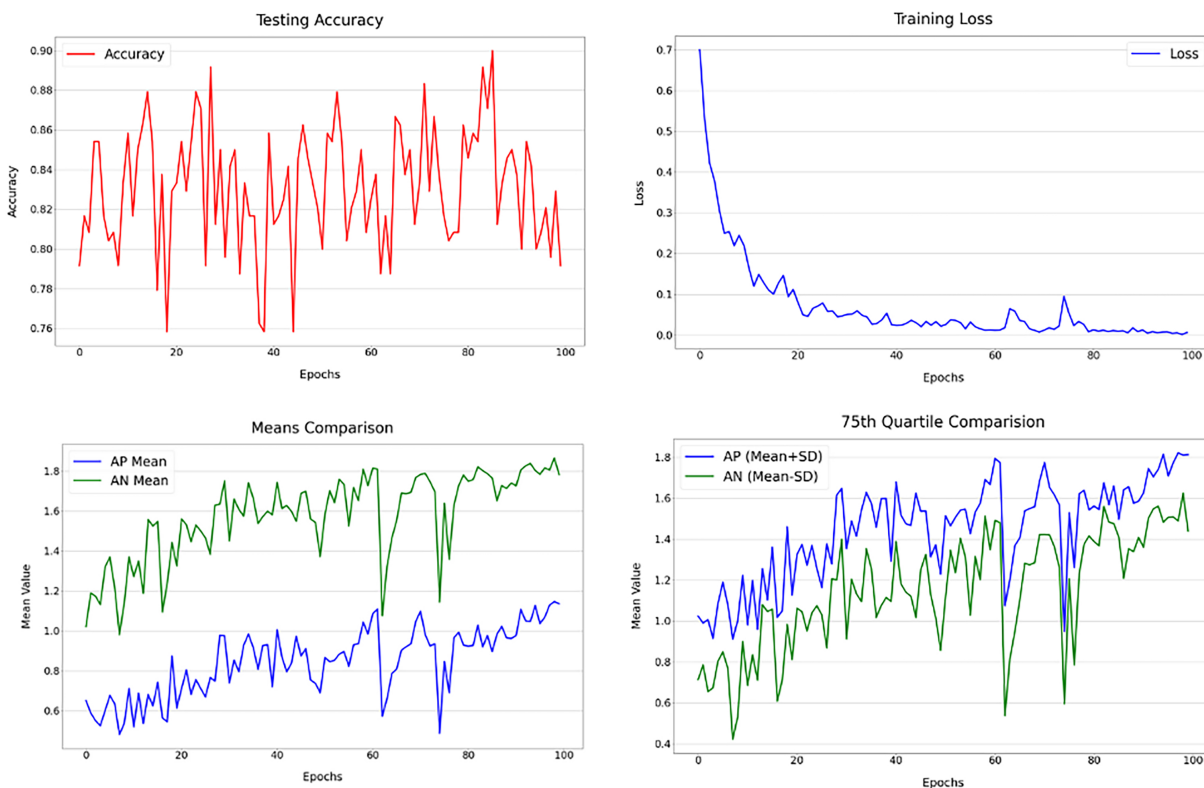| Evaluation Metric | Value |
|---|---|
| **Total images** | 32560 |
| **Truly detected images** | 28448 |
| **Falsely detected images** | 4112 |
| **Accuracy** | 87.37% |
| **Error** | 12.62% |

Figure 12. Performance measures of the proposed palm print authentication model
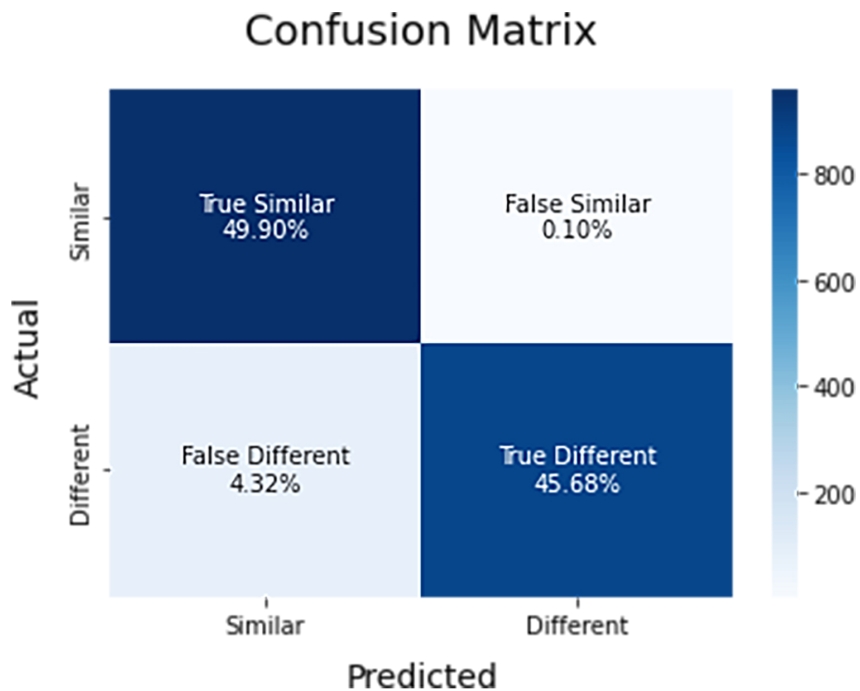


Figure 13. Confusion matrix generated by the Siamese Networks

14

**C. Evaluation of Palm Print Identification Module**

In addition to evaluating the performance of our hand gesture recognition model, we also sought to assess the accuracy of our palmprint identification component. To do so, we utilized the CASIA Dataset [35], which comprises a large and diverse collection of palmprint images. The dataset includes images from over 100 individuals, captured under various lighting conditions. However, we utilised a subset of the data, particularly the images taken in lowest wavelength sample and white light specifically. As a result, we acquired a dataset with 12000 sample in total. Further on, we pre-processed the dataset leveraging a custom data loader that generated test triplets from the CASIA dataset, where a triplet consists of an anchor image, a positive image, and a negative image. By training our model on these triplets, it is able to learn to differentiate between the palmprints of different individuals and accurately identify a user based on their palmprint. We trained our model using a triplet loss function [29] and Adam optimizer [30], which allows the model to optimize its error by comparing the similarity between the anchor and positive images to that of the anchor and negative images. This approach allows the model to learn the underlying features of a palmprint that are unique to an individual, which enables it to accurately identify a user based on their palmprint. The proposed model's training metrics are depicted in the Figure 12.

One of the key factors for optimized performance of our model is the Triplet Loss [29]. We utilised Triplet loss to compare the relative similarity of three inputs: an anchor image, a positive image, and a negative image. The anchor image is taken to be the "neutral" image, while the positive image is a image from the same subject, in our case, and the negative image is from a different subject. The functions was utilised to minimize the distance between the anchor and positive images, while maximizing the distance between the anchor and negative images. This is done by adjusting the model's weights and biases to better discriminate between the three inputs. As a result, the model becomes better at recognizing the similarities and differences between the anchor and positive images, and can more effectively differentiate between the anchor and negative images. The Figure 13 illustrates the confusion matrix acquired from the test set of the data.

Finally, the evaluation results of our proposed SimplyMime model on multiple benchmark datasets indicate its effectiveness in both hand gesture recognition and palmprint identification. Compared to existing solutions, our model has demonstrated superior performance. Further, by incorporating palmprint identification as an additional security measure, the model's practicality and usability in real-world applications have been further enhanced, positioning it as a viable alternative to conventional remote control devices. The novel combination of hand gesture recognition and palmprint identification in one system presents an exciting advancement in the field of human-computer interaction.

## VI. Conclusion

In conclusion, this paper presents SimplyMime, a novel hand gesture-based control system that aims to provide an immersive, efficient, and secure user experience while eliminating the need for multiple remote controls for consumer electronics. The system leverages advanced hand gesture recognition techniques, incorporating the latest developments in Artificial Intelligence and Human-Computer Interaction, to create a sophisticated architecture that can recognize a wide range of hand gestures with exceptional accuracy. Additionally, SimplyMime incorporates a palm print authentication module, which enhances the security of the system by ensuring that only authorized users can access the device. Through thorough testing and evaluation, SimplyMime demonstrated remarkable performance, achieving high accuracy levels of 96.16%, 87.37%, and 90% in hand detection, recognition, and palm print authentication, respectively. These results served as a testament to the effectiveness and efficiency of SimplyMime. Overall, SimplyMime offers significant advantages over traditional remote control systems, making it an excellent alternative for users looking for a more intuitive and efficient way of controlling their consumer electronics.

Despite the impressive performance of SimplyMime, there is still scope for improvement and further research. In the future, we plan to enhance the robustness of the system by incorporating additional sensors, such as proximity and depth sensors, to improve the accuracy and reliability of the system. Moreover, we aim to reduce the computational power required while improving the accuracy of the model. Another area of future research is the potential for SimplyMime to be integrated into other applications, such as virtual and augmented reality, to expand its capabilities and utility. By incorporating these enhancements, SimplyMime has the potential to become a highly versatile and widely adopted hand gesture-based control system that can revolutionize the way we interact with consumer electronics.

# References

[1] A. M. Joshi, P. Jain, S. P. Mohanty, and N. Agrawal, "iGLU 2.0: A New Wearable for Accurate Non-Invasive Continuous Serum Glucose Measurement in IoMT Framework," *IEEE Trans. Consumer Electron.*, vol. 66, no. 4, pp. 327–335, 2020. [Online]. Available: https://doi.org/10.1109/TCE.2020.3011966

[2] L. Rachakonda, S. P. Mohanty, E. Kougianos, and P. Sundaravadivel, "Stress-Lysis: A DNN-Integrated Edge Device for Stress Level Detection in the IoMT," *IEEE Trans. Consumer Electron.*, vol. 65, no. 4, pp. 474–483, 2019. [Online]. Available: https://doi.org/10.1109/TCE.2019.2940472

[3] A. S. G. Andrae and O. Andersen, "Life cycle assessments of consumer electronics — are they consistent?" *The International Journal of Life Cycle Assessment*, vol. 15, no. 8, pp. 827–836, Sep. 2010. [Online]. Available: https://doi.org/10.1007/s11367-010-0206-1

[4] T. Issa and P. Isaias, "Usability and Human–Computer Interaction (HCI)," in *Sustainable Design: HCI, Usability and Environmental Concerns*, T. Issa and P. Isaias, Eds. London: Springer, 2022, pp. 23–40. [Online]. Available: https://doi.org/10.1007/978-1-4471-7513-1_2

[5] M. Oudah, A. Al-Naji, and J. Chahl, "Hand Gesture Recognition Based on Computer Vision: A Review of Techniques," *Journal of Imaging*, vol. 6, no. 8, p. 73, Aug. 2020. [Online]. Available: https://www.mdpi.com/2313-433X/6/8/73

[6] E. Seiter, H. Borchers, G. Kreutzner, and E.-M. Warth, Eds., *: Television, Audiences, and Cultural Power*. London: Routledge, May 2013.

[7] M. S. Sodhi and S. Lee, "An analysis of sources of risk in the consumer electronics industry," *Journal of the Operational Research Society*, vol. 58, no. 11, pp. 1430–1439, Nov. 2007. [Online]. Available: https://doi.org/10.1057/palgrave.jors.2602410

[8] S. S. Rautaray and A. Agrawal, "Vision based hand gesture recognition for human computer interaction: a survey," *Artificial Intelligence Review*, vol. 43, no. 1, pp. 1–54, Jan. 2015. [Online]. Available: http://link.springer.com/10.1007/s10462-012-9356-9

[9] H. S. Hasan and S. A. Kareem, "Human Computer Interaction for Vision Based Hand Gesture Recognition: A Survey," in *2012 International Conference on Advanced Computer Science Applications and Technologies (ACSAT)*, Nov. 2012, pp. 55–60.

[10] L. Guo, Z. Lu, and L. Yao, "Human-Machine Interaction Sensing Technology Based on Hand Gesture Recognition: A Review," *IEEE Transactions on Human-Machine Systems*, vol. 51, no. 4, pp. 300–309, Aug. 2021.

[11] K. Teachasrisaksakul, L. Wu, G.-Z. Yang, and B. Lo, "Hand Gesture Recognition with Inertial Sensors," in *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, Jul. 2018, pp. 3517–3520, iSSN: 1558-4615.

[12] B. Coffen and M. Mahmud, "TinyDL: Edge Computing and Deep Learning Based Real-time Hand Gesture Recognition Using Wearable Sensor," in *2020 IEEE International Conference on E-health Networking, Application & Services (HEALTHCOM)*, Mar. 2021, pp. 1–6.

[13] K. Nguyen-Trong, H. N. Vu, N. N. Trung, and C. Pham, "Gesture Recognition Using Wearable Sensors With Bi-Long Short-Term Memory Convolutional Neural Networks," *IEEE Sensors Journal*, vol. 21, no. 13, pp. 15 065–15 079, Jul. 2021.

[14] G. Yuan, X. Liu, Q. Yan, S. Qiao, Z. Wang, and L. Yuan, "Hand Gesture Recognition Using Deep Feature Fusion Network Based on Wearable Sensors," *IEEE Sensors Journal*, vol. 21, no. 1, pp. 539–547, Jan. 2021.

[15] G. Singh, A. Nelson, R. Robucci, C. Patel, and N. Banerjee, "Inviz: Low-power personalized gesture recognition using wearable textile capacitive sensor arrays," in *2015 IEEE International Conference on Pervasive Computing and Communications (PerCom)*, Mar. 2015, pp. 198–206.

[16] F. Chen, H. Lv, Z. Pang, J. Zhang, Y. Hou, Y. Gu, H. Yang, and G. Yang, "WristCam: A Wearable Sensor for Hand Trajectory Gesture Recognition and Intelligent Human–Robot Interaction," *IEEE Sensors Journal*, vol. 19, no. 19, pp. 8441–8451, Oct. 2019.

[17] J. McIntosh, A. Marzo, M. Fraser, and C. Phillips, "EchoFlex: Hand Gesture Recognition using Ultrasound Imaging," in *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, ser. CHI '17. New York, NY, USA: Association for Computing Machinery, May 2017, pp. 1923–1934. [Online]. Available: https://doi.org/10.1145/3025453.3025807

[18] R. Y. Wang and J. Popović, "Real-time hand-tracking with a color glove," *ACM Transactions on Graphics*, vol. 28, no. 3, pp. 63:1–63:8, Jul. 2009. [Online]. Available: https://doi.org/10.1145/1531326.1531369

[19] S. Feng and R. Murray-Smith, "Fusing Kinect sensor and inertial sensors with multi-rate Kalman filter," in *IET Conference on Data Fusion & Target Tracking 2014: Algorithms and Applications (DF&TT 2014)*, Apr. 2014, pp. 1–8.

[20] J. Xu, H. Wang, J. Zhang, and L. Cai, "Robust Hand Gesture Recognition Based on RGB-D Data for Natural Human–Computer Interaction," *IEEE Access*, vol. 10, pp. 54 549–54 562, 2022.

[21] D.-S. Tran, N.-H. Ho, H.-J. Yang, S.-H. Kim, and G. S. Lee, "Real-time virtual mouse system using RGB-D images and fingertip detection," *Multimedia Tools and Applications*, vol. 80, no. 7, pp. 10 473–10 490, Mar. 2021. [Online]. Available: https://doi.org/10.1007/s11042-020-10156-5

[22] S. Shin and W.-Y. Kim, "Skeleton-Based Dynamic Hand Gesture Recognition Using a Part-Based GRU-RNN for Gesture-Based Interface," *IEEE Access*, vol. 8, pp. 50 236–50 243, 2020.

16

[23] A. Caputo, A. Giachetti, S. Soso, D. Pintani, A. D'Eusanio, S. Pini, G. Borghi, A. Simoni, R. Vezzani, R. Cucchiara, A. Ranieri, F. Giannini, K. Lupinetti, M. Monti, M. Maghoumi, J. J. LaViola Jr, M.-Q. Le, H.-D. Nguyen, and M.-T. Tran, "SHREC 2021: Skeleton-based hand gesture recognition in the wild," *Computers & Graphics*, vol. 99, pp. 201–211, Oct. 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0097849321001382

[24] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single Shot MultiBox Detector," 2016, vol. 9905, pp. 21–37, arXiv:1512.02325 [cs]. [Online]. Available: http://arxiv.org/abs/1512.02325

[25] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," Oct. 2014, arXiv:1311.2524 [cs] version: 5. [Online]. Available: http://arxiv.org/abs/1311.2524

[26] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications," Apr. 2017, arXiv:1704.04861 [cs]. [Online]. Available: http://arxiv.org/abs/1704.04861

[27] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh, "Convolutional Pose Machines," Apr. 2016, arXiv:1602.00134 [cs]. [Online]. Available: http://arxiv.org/abs/1602.00134

[28] A. Cyril Jose and R. Malekian, "Smart Home Automation Security: A Literature Review," *The Smart Computing Review*, Aug. 2015. [Online]. Available: http://smartcr.org/view/download.php?filename=smartcr_vol5no4p004.pdf

[29] T. Müller, G. Pérez-Torró, and M. Franco-Salvador, "Few-Shot Learning with Siamese Networks and Label Tuning," Apr. 2022, arXiv:2203.14655 [cs]. [Online]. Available: http://arxiv.org/abs/2203.14655

[30] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," Jan. 2017, arXiv:1412.6980 [cs]. [Online]. Available: http://arxiv.org/abs/1412.6980

[31] F. Chollet, "Xception: Deep Learning with Depthwise Separable Convolutions," Apr. 2017, arXiv:1610.02357 [cs] version: 3. [Online]. Available: http://arxiv.org/abs/1610.02357

[32] E. Hoffer and N. Ailon, "Deep metric learning using Triplet network," Dec. 2018, arXiv:1412.6622 [cs, stat]. [Online]. Available: http://arxiv.org/abs/1412.6622

[33] C. Nuzzi, S. Pasinetti, R. Pagani, G. Coffetti, and G. Sansoni, "HANDS: an RGB-D dataset of static hand-gestures for human-robot interaction," *Data in Brief*, vol. 35, p. 106791, Apr. 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2352340921000755

[34] C. Zimmermann, D. Ceylan, J. Yang, B. Russell, M. Argus, and T. Brox, "FreiHAND: A Dataset for Markerless Capture of Hand Pose and Shape from Single RGB Images," Sep. 2019, arXiv:1909.04349 [cs]. [Online]. Available: http://arxiv.org/abs/1909.04349

[35] Z. Sun, T. Tan, Y. Wang, and S. Li, "Ordinal palmprint representation for personal identification [representation read representation]," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, vol. 1, Jun. 2005, pp. 279–284 vol. 1, iSSN: 1063-6919.

[36] J. Feng, P. Ji, and F. Ma, "Gesture Position Detection Based on Generative Adversarial Networks," in *2022 2nd International Conference on Robotics and Control Engineering*, ser. RobCE 2022. New York, NY, USA: Association for Computing Machinery, Apr. 2022, pp. 39–44. [Online]. Available: https://doi.org/10.1145/3529261.3529268

[37] T. L. Dang, S. D. Tran, T. H. Nguyen, S. Kim, and N. Monet, "An improved hand gesture recognition system using keypoints and hand bounding boxes," *Array*, vol. 16, p. 100251, Dec. 2022. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2590005622000844

**Sibi C. Sethuraman** (M'18) received his Ph.D from Anna University in the year 2018. He is an Associate Professor in the School of Computer Science and Engineering at Vellore Institute of Technology – Andhra Pradesh (VIT-AP) University. Further, he is the coordinator for Artificial Intelligence and Robotics (AIR) Center at VIT-AP. He is an active reviewer in many reputed journals of IEEE, Springer, and Elsevier. He is a recipient of DST fellowship.

**Gaurav Reddy Tadkapally** (S'19) is a Bachelor of Technology student at Vellore Institute of Technology, Amaravati (VIT-AP). He specializes in Deep Learning and the applications of Artificial Intelligence in Consumer Electronics and Embedded Hardware. As a member of the Artificial Intelligence and Robotics center at VIT-AP, Gaurav has developed several cutting-edge technologies, including iDrone, an IoT-powered drone for detecting wildfires, and MagicEye, an intelligent wearable designed towards independent living of Visually impaired.



**Athresh Kiran** (M'22) received the bachelor's degree in Computer Science & Engg. from the VIT-AP University, India, in 2021. He is currently working in Parallel Reality. His research interest includes reality platforms, Metaverse, IT, full stack development etc.



**Saraju P. Mohanty** (Senior Member, IEEE) received the bachelor's degree (Honors) in electrical engineering from the Orissa University of Agriculture and Technology, Bhubaneswar, in 1995, the master's degree in Systems Science and Automation from the Indian Institute of Science, Bengaluru, in 1999, and the Ph.D. degree in Computer Science and Engineering from the University of South Florida, Tampa, in 2003. He is a Professor with the University of North Texas. His research is in "Smart Electronic Systems" which has been funded by National Science Foundations (NSF), Semiconductor Research Corporation (SRC), U.S. Air Force, IUSSTF, and Mission Innovation. He has authored 450 research articles, 5 books, and 9 granted and pending patents. His Google Scholar h-index is 49 and i10-index is 211 with 10,800 citations. He is regarded as a visionary researcher on Smart Cities technology in which his research deals with security and energy aware, and AI/ML-integrated smart components. He introduced the Secure Digital Camera (SDC) in 2004 with built-in security features designed using Hardware Assisted Security (HAS) or Security by Design (SbD) principle. He is widely credited as the designer for the first digital watermarking chip in 2004 and first the low-power digital watermarking chip in 2006. He is a recipient of 16 best paper awards, Fulbright Specialist Award in 2020, IEEE Consumer Electronics Society Outstanding Service Award in 2020, the IEEE-CS-TCVLSI Distinguished Leadership Award in 2018, and the PROSE Award for Best Textbook in Physical Sciences and Mathematics category in 2016. He has delivered 15 keynotes and served on 14 panels at various International Conferences. He has been serving on the editorial board of several peer-reviewed international transactions/journals, including IEEE Transactions on Big Data (TBD), IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD), IEEE Transactions on Consumer Electronics (TCE), and ACM Journal on Emerging Technologies in Computing Systems (JETC). He has been the Editor-in-Chief (EiC) of the IEEE Consumer Electronics Magazine (MCE) during 2016-2021. He served as the Chair of Technical Committee on Very Large Scale Integration (TCVLSI), IEEE Computer Society (IEEE-CS) during 2014-2018 and on the Board of Governors of the IEEE Consumer Electronics Society during 2019-2021. He serves on the steering, organizing, and program committees of several international conferences. He is the steering committee chair/vice-chair for the IEEE International Symposium on Smart Electronic Systems (IEEE-iSES), the IEEE-CS Symposium on VLSI (ISVLSI), and the OITS International Conference on Information Technology (OCIT). He has mentored 2 post-doctoral researchers, and supervised 14 Ph.D. dissertations, 26 M.S. theses, and 18 undergraduate projects.

**Anitha S** (M'22) received the bachelor's degree in Electronics & Communication Engg. from the Anna University, India, in 2017, the master's degree in Power Electronics from the Anna University University in 2019, and she is currently pursuing her Ph.D. degree in Electronics Engineering from VIT-AP University.