# Implementation of Interactive System
# with Millimeter Wave and Wavelet Transform

Bo-Heng Chen, Da-Chuan Chen, Wen-Hsiang Yeh,
Hao-Yang Chen, and Yu-Ping Liao[*]

Department of Electrical Engineering, Chung Yuan Christian University,
No. 200, Zhongbei Rd., Zhongli Dist., Taoyuan City 320314, Taiwan (R.O.C.)

In current society, escalating stress in everyday life and work creates significant concern about the potential for both physical and psychological long-term effects if unattended. To address this challenge, in this paper, we introduce a new approach designed to alleviate stress and promote well-being. Our proposed system synchronizes pointer swings with music rhythms to achieve a combination of rhythmic music engagement and interactive physical activity. This solution centers on the intersection of technology and well-being by providing users with an immersive and relaxing audio-visual experience. To achieve higher rhythm extraction precision in different genres of music, we employ wavelet transform and deep learning (DL) to adjust parameters dynamically. The controllable rotating sphere is achieved using a millimeter-wave (mmWave) radar and a thermographic camera to detect preset gestures. This system provides the interaction of rhythm and control over the rotating sphere, which can provide users with an immersive experience and contribute to a stress-relieving experience.

## 1. Introduction

Excessive stress, an acknowledged precursor to physical and psychological disorders, is critical in today's lifestyle. Stress relief is widely acknowledged to be achieved through physical activities and listening to music.[1–4] To tackle this problem, we present an interactive stress relief system.

Our proposed system uses the synchronization movement of a pointer with ongoing music playback to simultaneously engage users through physical activities and listening to music. The system uses an analog-to-digital converter (ADC) for input signals from a microphone and a microcontroller unit (MCU) for rhythm extraction in audio processing. Previous studies have emphasized the importance of parameter selection when performing discrete wavelet transforms to achieve effective music genre extraction.[5,6] Therefore, we adopted a deep learning (DL)

---

model in our embedded system to distinguish music genres, allowing appropriate parameters to be chosen for accurate rhythm extraction.

To further enhance user interaction with our system, we integrate a mmWave radar to detect large gestures, referring to previous research,[7] and recent applications.[8–11] We combined data from a mmWave radar and a thermographic camera to improve the accuracy in gesture recognition and implemented it in an embedded system. Such data then undergoes further processing using a DL model.

In this paper, we build upon our previous research to construct a touchless interactive system that merges music rhythm extraction with gesture recognition. While the audio controls the pointer's movement, users can manipulate the sphere's rotation using gestures, as presented in Fig. 1. This system can successfully reduce stress for users who like music and enjoy engagement in physical activity simultaneously.

The remainder of the paper is organized as follows: In Sect. 2, we examine the two subsystems that control sphere rotation and pointer movement. In Sect. 3, we explain the working principles that support our system's operations. In Sect. 4, we present the results of our experiments and analyze our system's performance.

## 2. System Structure and Communication

In this section, the system structure and interdevice communication are described. The system comprises a pointer control system working with a sphere control system. These systems are further divided into sections and explained separately.

### 2.1 Pointer control system

To control the pointer, we start by extracting the rhythm using wavelet transform and adjusting the parameters to determine the number of current beats per minute (BPM) by music
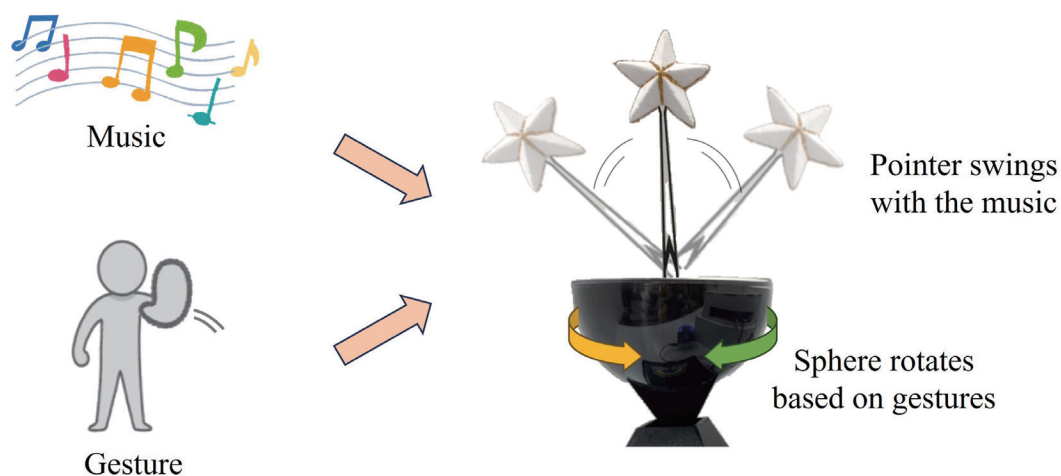


Fig. 1.　(Color online) System interactive scenarios.

genre recognition. Then, the pointer control system uses the BPM value to alternate the speed of pointer swings and synchronize them with the music rhythm, as shown in Fig. 2. We will explain this process in three sequential parts: the genre recognition subsystem, the rhythm extraction subsystem, and the pointer control subsystem.

### 2.1.1   Genre recognition subsystem

The signal is captured and transmitted by a microphone and then converted into log-power spectrograms. A two-dimensional convolutional neural network (CNN) further processes the spectrograms obtained from the audio input. This process generates recognition outcomes that are utilized to identify the genre. The genre identified is then transmitted to the rhythm extraction subsystem, which is responsible for adjusting the parameters used in wavelet transform. The entire process is depicted in Fig. 3 and is crucial for the accurate identification and analysis of the audio input.

### 2.1.2   Rhythm extraction subsystem

In our processing pipeline, the user must establish a Bluetooth connection between a smartphone and our audio receiver board. Once the transmission of audio signal starts from the user's smartphone to the receiver board, the signal passes through the pre-processing circuit. The MCU first digitizes the signal with an ADC and then performs rhythm extraction calculations to compute the current BPM value. When the extraction process is activated, it performs wavelet transformation, full-wave rectification, threshold filtering, dilation, and
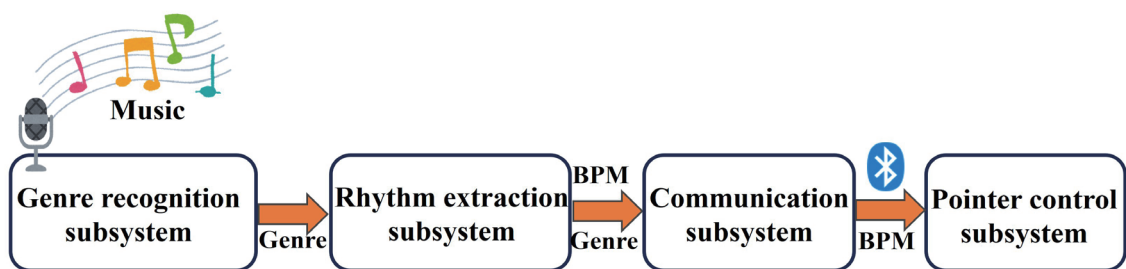
Fig. 2.    (Color online) Flow chart of the pointer control system.
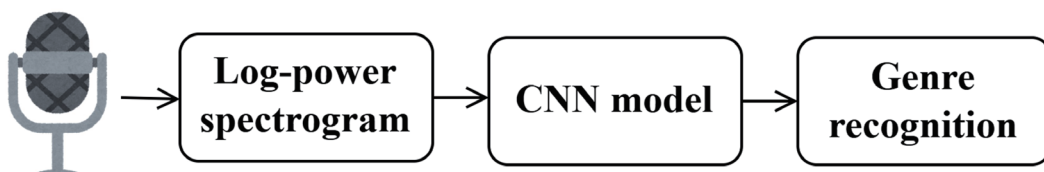
Fig. 3.    (Color online) Genre recognition subsystem flow chart.

erosion, as shown in Fig. 4. Once the current BPM value is computed, the MCU passes it to the communication subsystem with the help of universal asynchronous receiver/transmitter (UART) protocol. The communication subsystem, which utilizes a Raspberry Pi board for implementation, uses the calculated BPM value for pointer control.

### 2.1.3 Pointer control subsystem

Once the BPM value is computed in the rhythm extraction stage, it is also sent to the Raspberry Pi board with Bluetooth low energy (BLE). Next, the pointer control subsystem based on the ESP32 board, as shown in Fig. 5, retrieves the current BPM value. Further control of the stepper motor by modifying the swinging frequency of the pointer will be performed with the ESP32 board with a time-varying BPM value.

## 2.2 Sphere control system

The sphere control subsystem is equipped with a thermographic camera and a mmWave radar to provide various controls for sphere rotation via large gestures, as shown in Fig. 6.
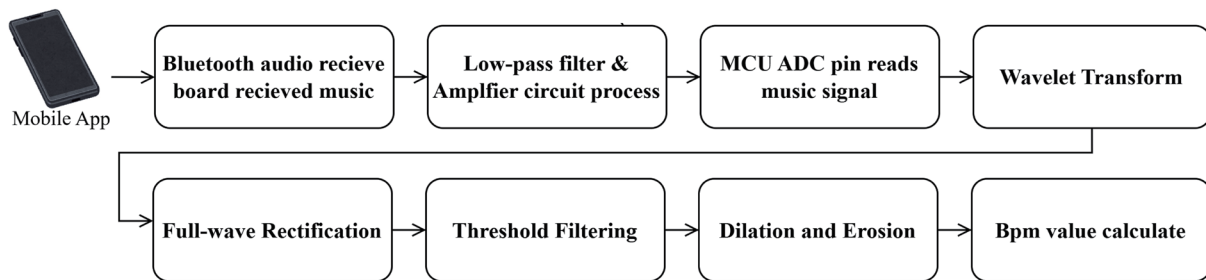


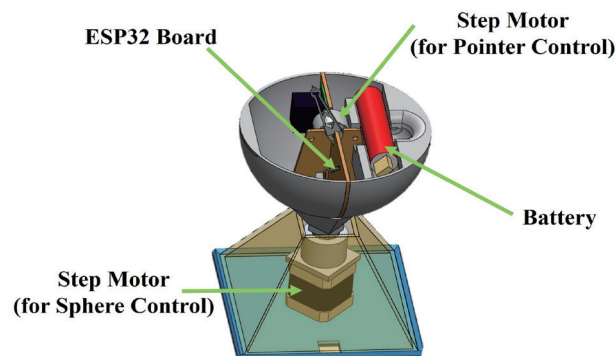Fig. 4.　(Color online) Flow chart of the rhythm extraction subsystem.



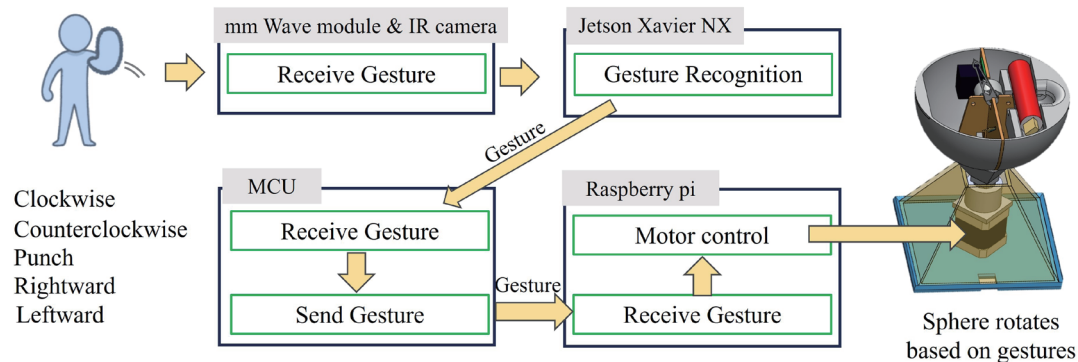Fig. 5.　(Color online) Inner structure of the sphere.

Fig. 6.    (Color online) Flow diagram for the sphere control system.

### 2.2.1   Gesture recognition

Our recognition process uses two different sensors—a thermographic camera and a mmWave radar—together as input to aggregate data. The mmWave radar captures sets of point clouds that undergo several processes, including maximum velocity limiting, 2D data clustering algorithm —density-based spatial clustering application with noise (DBScan), image registration, $k$-means clustering, temporal feature extraction, and data normalization. On the other hand, the images of palms captured by the thermographic camera will undergo palm location recognition, interpolation, and normalization using the YOLOv7 model.[12] After data of both sensors are processed, the results are integrated and collectively computed with a gated recurrent unit (GRU) model for recognizing large control gestures. The process is illustrated in Fig. 7. Recognition results from the GRU model are later sent to the MCU through the UART protocol to aid sphere control.

### 2.2.2   Sphere control

Once the MCU receives the recognition result from the gesture recognition process via the UART protocol, it will pass the result to Raspberry Pi. Raspberry Pi will utilize the received recognition result to control the stepper motor seats below the sphere.

## 3.    Operating Principle

In this section, we describe the algorithms used by each device. The proposed interactive system is able to extract rhythm by wavelet transform, identify music genres using artificial intelligence models, and recognize gestures using mmWave radars and thermal imaging cameras based on DL technology.
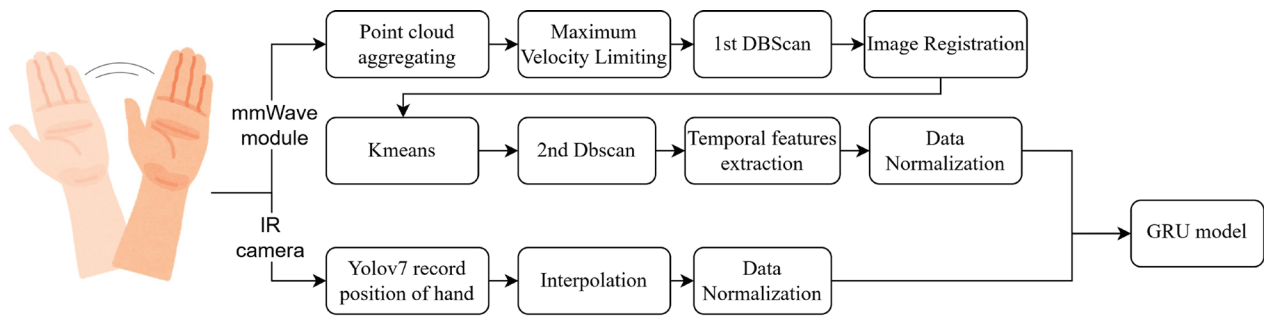
Fig. 7. (Color online) Gesture recognition system flow diagram.

### 3.1 Rhythm extraction using wavelet transform

As part of the process of extracting the rhythm from the music received by the BLE receiver board, the signals undergo initial treatment. This treatment involves a low-pass filter and amplifier circuit that removes all signals with frequencies higher than 50 Hz while enhancing the rest of the frequency components. Then, rhythm extraction from the treated signals is performed with the HT32F52352 MCU, as shown in Fig. 8. The following steps describe how the rhythm extraction proceeds:

Step 1: Daubechies Wavelet Transform—Daubechies wavelets are used to differentiate high- and low-frequency components in the wavelet transform.[5,6] Table 1 provides our parameters for the wavelet transform in accordance with the identified genre of music.

Step 2: Full-wave Rectification—The output data from the previous step is rectified into only positive values for the following step.

Step 3: Threshold Filtering—Any signal with an amplitude lower than a predefined threshold is converted to zero; those with amplitudes higher than the threshold are converted to one.

Step 4: Erosion and Dilation—Fast alternating signals are grouped together to create large spikes.

Step 5: BPM Calculation—With the help of sampling frequency, we use the time between two signal spikes to calculate the current BPM value.

### 3.2 Genre recognition

The Raspberry Pi in our system will perform the following steps to identify the music genre and tweak the wavelet transform parameter settings accordingly after retrieving audio signals from the microphone:

Step 1: Log-power Spectrogram—We use short-time Fourier transform (STFT) to generate power spectrograms with a logarithmic scale on the frequency axis to enhance the variations in lower-frequency regions. The Librosa library is used to aid the computation process.

Step 2: CNN Model—The power spectrogram generated in the previous step is used to perform genre recognition with a CNN model consisting of conv2d, max-pooling, and dense layers, as shown in Fig. 9.
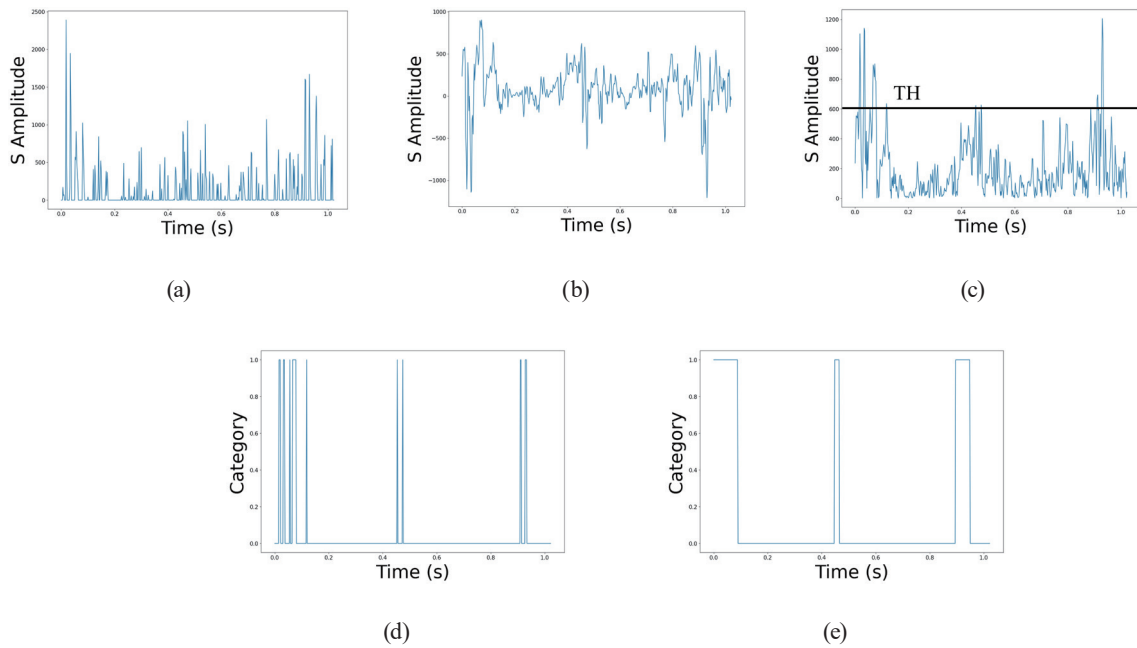
Fig. 8. (Color online) Sample result of digital signal processing for music extraction. (a) ADC receives input signal. (b) Wavelet transform. (c) Full-wave rectification. (d) Threshold filtering. (e) Erosion and dilation.

Table 1
Lookup table of wavelet parameters and music genres.

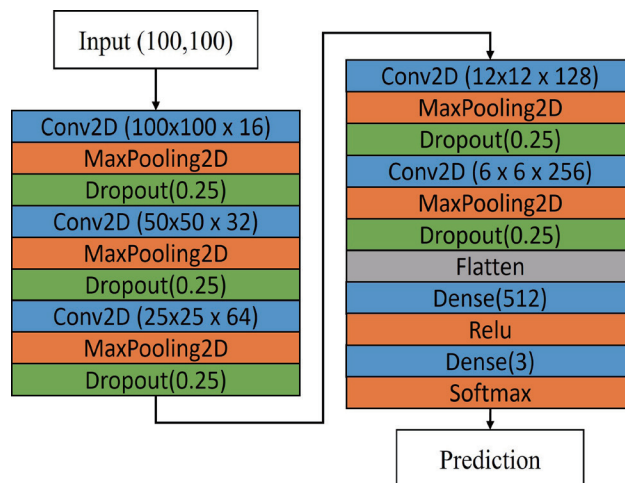| Genre | dB | Level |
|---|---|---|
| Classical | 2 | 4 |
| POP | 6 | 8 |
| Rock | 4 | 8 |



Fig. 9. (Color online) CNN model structure.

### 3.3　Gesture recognition

We use gesture data aggregated from both mmWave radar and thermal imaging camera sensors, perform preprocessing, and use DL techniques to recognize different genres with Jetson Xavier NX. In this section, we outline the preprocessing pipeline and the DL model deployed for recognition.

### 3.3.1　mmWave radar signal processing

The mmWave radar we adopted detects moving objects with the phase difference between two different frequency-modulated continuous waves (FMCWs) with a given time interval. The primary goal of preprocessing data from a mmWave radar is to capture only the data related to palm movements and to eliminate noise. All gestures can be performed within this time frame of 2.5 s, according to the results of our experiments, and 20 frames of point cloud data are generated. The following noise reduction steps are implemented to combine these 20 frames:

Step 1: Maximum Speed Limitation—Any data in the gathered point cloud exceeding two m/s or extending beyond 1.5 m from the sensor is discarded to reduce noise.

Step 2: Outlier Rejection Using DBScan—A density-based clustering algorithm such as DBScan can work effectively with our data with two important parameters: eps and min_samples. The algorithm will group data within a radius of eps if there are more than min_samples data points. When the radius contains fewer data points than min_samples, they will be labeled as noise data and removed. We can successfully perform outlier rejection operations by setting eps and min_samples to 0.1 and 20.

Step 3: Image Registration—We use data rotation and repositioning techniques to increase model recognition accuracy and reduce data complexity for each data point.

Step 4: Body and Palm Separation with $K$-means—To separate the palm and body presented in the data, we use the $K$-means algorithm, which splits the data into two groups. The separation mechanism is achieved using the distance between the center of each group and data points. Since there are only two groups of data to process, we use two for the number of data groups ($K$).

Step 5: Additional Outlier Removal with DBScan—The same process as in Step 2 is performed again to further remove outliers generated in the previous step with eps and min_samples set to 0.1 and 5.

Step 6: Temporal Feature Extraction—All point cloud data across all frames are centered by $x$, $y$, and $z$ coordinates to generate temporal data for our DL model.

Step 7: Data Normalization—We use the MinMaxScaler technique to normalize the $x$, $y$, and $z$ coordinates to speed up our DL model converging process. This process limits previously processed point cloud data values to a range of 0 to 1. Furthermore, MaxAbsScaler is applied to velocity data to scale data into a range of $-1$ to 1.

### 3.3.2 Thermal imaging camera signal processing

Thermal imaging cameras are effective in detecting long-wave infrared (LWIR) emissions and their energy. They convert this energy into temperature values, which are displayed as colors in the captured images. Our main goal is to identify the location of the palm in an image by the following preprocessing steps:

Step 1: YOLOv7 Recognition—We process thermal images continuously for 2.5 s with the YOLOv7 model to recognize the palms of users in images.

Step 2: Interpolation—As the thermographic camera can only capture 12 frames of images, we use interpolation to expand the data set to 20 frames to match the volume of mmWave radar data.

Step 3: Data Normalization—We divide all data by 400 and restrict their values to the range of 0 to 1. This process can increase the convergence speed of our model.

### 3.3.3 GRU model for gesture recognition

Research on recurrent neural networks has shown that the GRU model has higher performance per watt and can recognize temporal data more rapidly than the long short-term memory (LSTM) network. [13] Hence, we have opted for a GRU model that comprises two GRU layers and a dense layer, as illustrated in Fig. 10, instead of LSTM for gesture recognition.

## 4. Results

This interactive system not only allows users to interact with it using gestures, but the pointer in the sphere will also swing to the beat of the music.

### 4.1 Gesture design and sphere control

To control the rotation of the sphere, we designed five gestures: clockwise, counterclockwise, leftward, rightward, and punch motions, as shown in Fig. 11. The movement of these gestures is
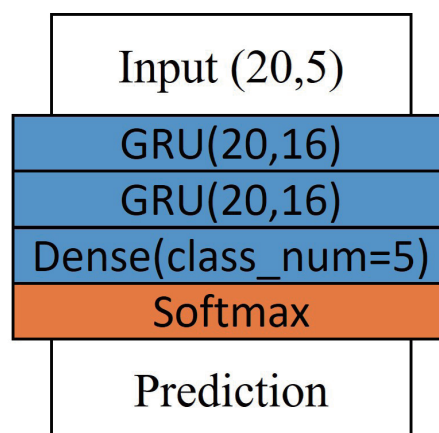


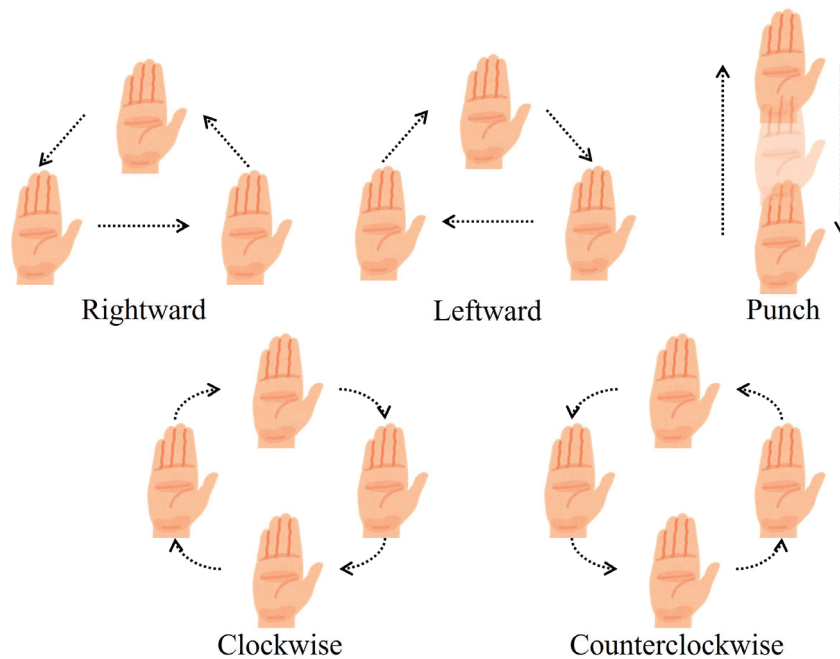Fig. 10.   (Color online) GRU model structure.

Fig. 11. (Color online) Schematic of designed gesture motions.

limited to the *x-y* plane and has minimal variation on the *z*-axis coordinate. Each gesture has a corresponding predefined rotation mode, which involves right turn, left turn, deceleration, acceleration, and stop, Table 2 shows the relationship between the gestures and the sphere rotation that we designed. Users can control the sphere interactively with these five different gestures and switch the rotation mode.

## 4.2 Thermal imaging camera processing

Thermal imaging cameras generate colored images of object surface temperatures. Our initial process assigns color only to areas with temperatures higher than the set single threshold of 30 °C. However, this process is ineffective when the ambient temperature increases, making it difficult to distinguish the user's palm from the rest of the body, as shown in Fig. 12(a). This results in inaccurate palm recognition. To improve this process, we introduced a multiple threshold color assignment technique, as depicted in Fig. 12(b). Our later experiment with this technique reinforces our decision to adopt multiple thresholds, which significantly enhances recognition accuracy. Still, if the temperature of the surrounding environment is close to that of the user's palm, our experiments show that recognition by the YOLOv7 model can still be difficult.

## 4.3 Gesture recognition accuracy

We have gathered a total of 14480 pieces of temporal data for gesture recognition, 2896 pieces per gesture. This dataset was segregated into three subsets: 60% for training, 20% for

Table 2
Lookup table of gestures and sphere rotation.

| Gesture | Sphere rotation |
|---|---|
| Clockwise | Turn right |
| Counterclockwise | Turn left |
| Leftward | Decelerate |
| Rightward | Accelerate |
| Punching | Stop |



(a)                                                                                               (b)
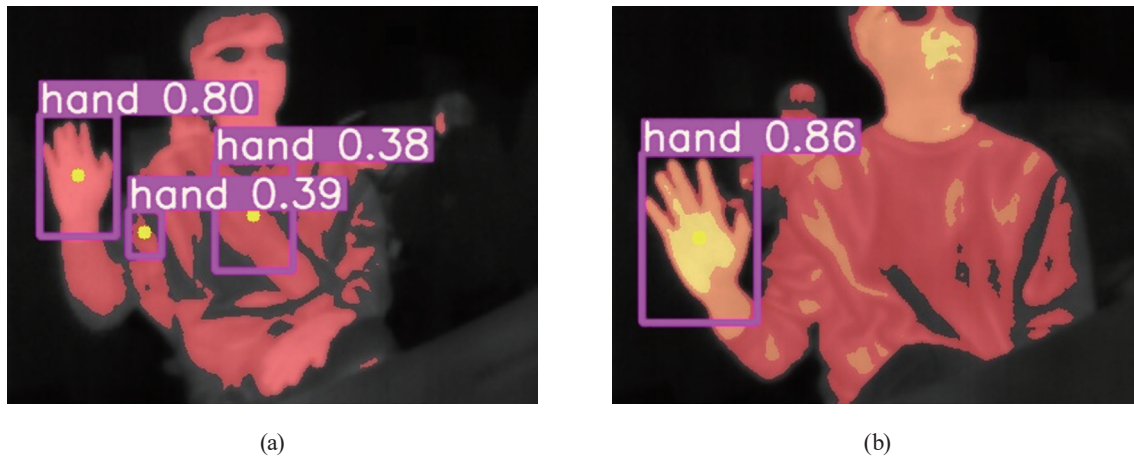
Fig. 12.   (Color online) Thermal camera images: (a) single and (b) multiple thresholds.

validation, and the last 20% for testing. Upon the completion of the model training phase, we analyzed the outcomes by employing a confusion matrix, which is depicted in Fig. 13, which shows that our model can recognize all of the designed gestures with minimal misclassification.

To further evaluate our system's real-world performance, we also collected gesture data from participants not presented in our dataset. However, the recognition accuracy was significantly lower than the training results owing to individual differences in gesture execution. Our model struggled to differentiate between leftward and counterclockwise gestures and rightward and clockwise gestures. In addition, the recognition of the punch gesture showed unsatisfactory performance because of the limited recognition capabilities of the thermographic camera. The full results are presented in Table 3.

### 4.4   Genre recognition accuracy

Three music genres—classical, pop, and rock—from the GTZAN dataset were selected, and we employed a 3 s window with a 0.1 s shift for continuous recognition when reading audio data. Our experimental results, presented in Table 4, show that the rock genre has a much lower accuracy, likely because of the similarity of instruments used in both the pop and rock genres.

### 4.5   BPM calculation and pointer control

We calculated BPM with Eq. (1), where $f_s$ is the sampling frequency of the ADC and is set to 512 Hz. $Data_{num}$ is the number of data points between spikes, as shown in Fig. 14. We limited the
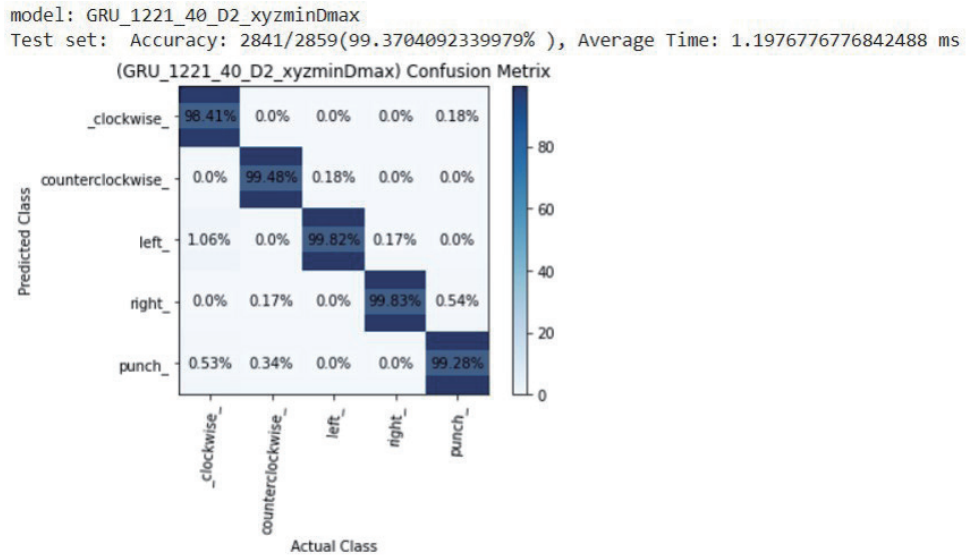
```
model: GRU_1221_40_D2_xyzminDmax
Test set:  Accuracy: 2841/2859(99.3704092339979% ), Average Time: 1.1976776776842488 ms
```

Fig. 13.   (Color online) Confusion matrix of GRU model.

Table 3
Accuracy of gesture recognition result.

| Gesture | Accuracy (%) |
|---|---|
| Clockwise | 70 |
| Counterclockwise | 40 |
| Leftward | 80 |
| Rightward | 60 |
| Punching | 20 |

Table 4
Accuracy of gesture recognition result.

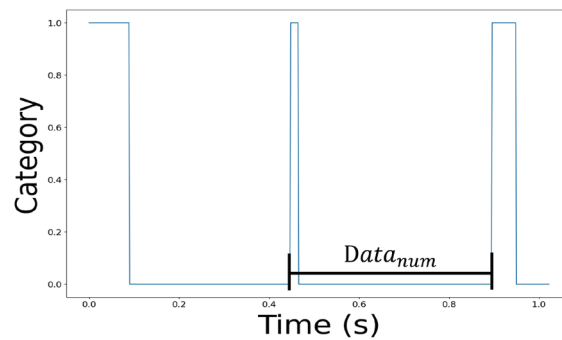| Music genre | Accuracy (%) |
|---|---|
| Classical music | 100 |
| POP music | 97.8 |
| Rock music | 64.8351 |

Fig. 14.   (Color online) $Data_{num}$ between two peaks.

maximum BPM to 100 and halved it if it exceeded this value to ensure that the stepper motor can operate even with fast music.

When designing the pointer, we set the oscillation period per cycle to 1.6 s. A delay time based on Eq. (2) will be applied if the current BPM falls below 100 to synchronize the pointer's movement with the music's tempo.

$$BPM = 60 / \left( Data_{num} * \frac{1}{f_s} \right), \; f_s = 512 \tag{1}$$

$$Delay\,time\,(\text{ms}) = 1000\,(\text{ms}) - BPM * 10 \tag{2}$$

### 4.6  Implemented system

We have implemented a demonstration system of the proposed interactive system, as depicted in Fig. 15. The device on the left-hand side is the pointer in a sphere, which can be controlled to swing and rotate. A thermal imaging infrared camera (IR camera) and mmWave radar are located at the top of the structure on the right-hand side. In addition, Jetson Xavier NX is located at the middle level, and it receives data transmitted from the sensors above and performs gesture recognition. The motor controller and signal processing circuits, that is, the Bluetooth receiver board, low-pass filter, amplifier circuit, HT32F52352 MCU, and Raspberry Pi, are located at the bottom level and are responsible for controlling the pointer in the sphere and rhythm extraction. The block diagram of the entire proposed system is illustrated in Fig. 16.

## 5.  Discussion

Stress relief using visual effects or music has been explored in numerous studies.[14–17] In this paper, we suggested integrating music, visual effects, and physical activity for stress relief. Employing a mmWave radar for gesture recognition and wavelet transform for rhythm extraction, motors were driven to rotate the pointer and sphere. The mmWave radar and thermographic camera data were utilized for detecting gesture coordinates, forming the basis for training the GRU model. This approach offers user stress relief through multiple modalities simultaneously.
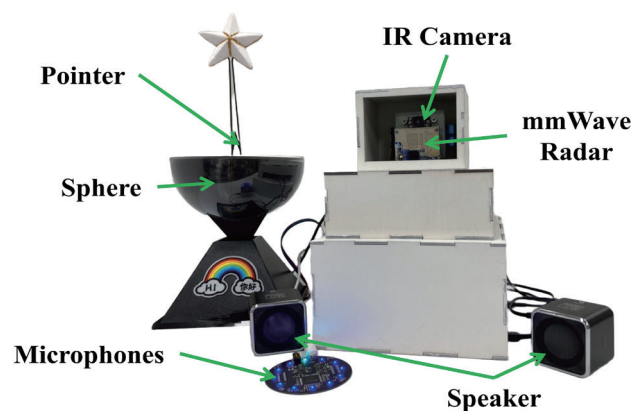


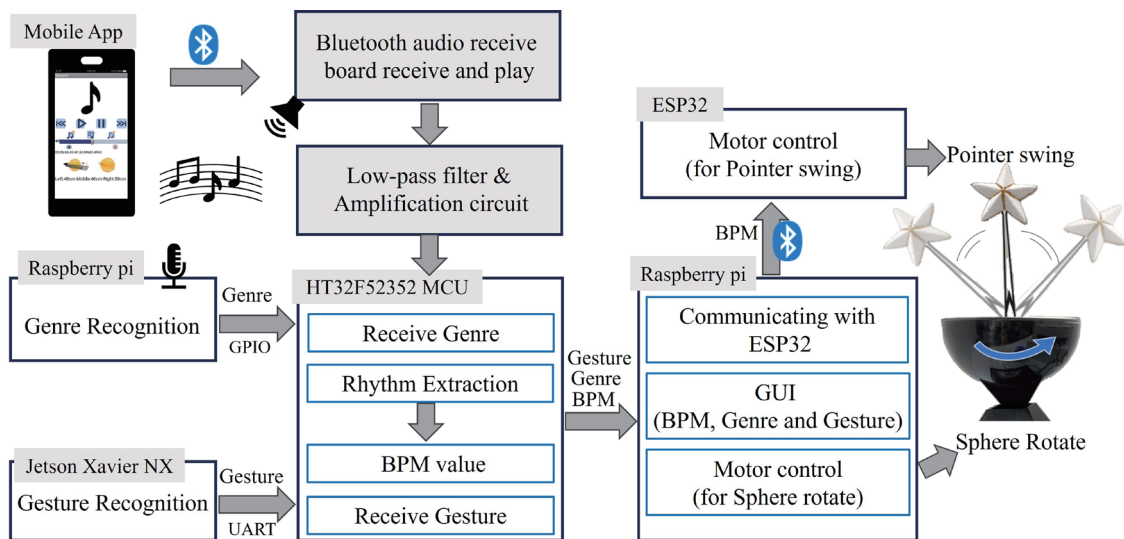Fig. 15.   (Color online) Implemented system.

Fig. 16.    (Color online) Block diagram of the proposed system.

## 6.    Conclusions

In this paper, we demonstrated the seamless integration of genre recognition, rhythm extraction, and gesture recognition, resulting in the successful development of a contactless interactive system. The system utilizes the wavelet transform to extract the music rhythm. In addition, the parameters of the wavelet transform are adjusted through the CNN model trained to accurately estimate the BPM value, synchronizing the pointer swings with the rhythm of the music. The system also utilizes both a mmWave radar and a thermal imaging camera to detect gesture coordinates. Then, the GRU model is used to recognize the gesture and control the rotation of the sphere. However, temperature and individual differences in gesture execution can disturb the accuracy of gesture recognition. In future work, we aim to improve accuracy by expanding the gesture dataset to include gesture motion data from multiple individuals. Additionally, we plan to enhance the multi-threshold layering to display a wider range of colors in the thermal images, which will help to mitigate environmental disturbances.

## References

1   M. Chennafi, M. A. Khan, G. Li, Y. Lian, and G. Wang: Proc. 2018 IEEE Asia Pacific Conf. Circuits and Systems (IEEE,2018) 131.
2   J. Alagha and A. Ipradjian: GJHSS-A. **17** (2017) 1. https://socialscienceresearch.org/index.php/GJHSS/article/view/2336
3   R. Hansmann, S. M. Hug and K. Seeland: Urban For. Urban Greening **6** (2007) 213. https://doi.org/10.1016/j.ufug.2007.08.004
4   S. T. Gura: Work **19** (2002) 3. https://api.semanticscholar.org/CorpusID:25955243
5   Y. T. Zhang, Y. L. Cheng, H. X. Wu, and Y. P. Liao: J. Phys. Conf. Ser. **2345** (2022) 1. https://doi.org/10.1088/1742-6596/2345/1/012021
6   B. H. Chen, Y. T. Zhang, H. X. Wu, R. C. Lu, and Y. P. Liao: 2023 The 8th Int. Conf. Precision Machinery and Manufacturing Technology, Kenting, Taiwan (2003).

7  Y. P. Liao, F. K. Huang, Y. J. Xia and H. Cheng: 2022 IEEE Int. Conf. Consumer Electronics (IEEE, 2022) 439.
8  J. T. Yu, L. Yen, and P. H. Tseng: 2020 IEEE 91st Vehicular Technology Conf. (IEEE, 2020) 1.
9  S. Skaria, A. Al-Hourani, and D. Huang: 2021 IEEE Sensors (IEEE, 2021) 1.
10  Q. Li, L. Liy, S. Hao, and G. Wan: 2022 5th Int. Conf. Pattern Recognition and Artificial Intelligence (IEEE, 2022) 63.
11  J. Liu, K. Furusawa, T. Tateyama, Y. Iwamoto and Y. W. Chen: 2019 IEEE Int. Conf. Image Processing (IEEE, 2019) 375.
12  C. Y. Wang, A. Bochkovskiy, and H. Y. M. Liao: 2023 IEEE/CVF Conf. Computer Vision and Pattern Recognition (IEEE, 2023) 7464.
13  S. Yang, X. Yu, and Y. Zhou: 2020 Int. Workshop on Electronic Communication and Artificial Intelligence (IEEE, 2020) 98.
14  M. Chennafi, M. A. Khan, G. Li, Y. Lian, and G. Wang: 2018 IEEE Asia Pacific Conf. Circuits and Systems (IEEE, 2018) 131.
15  M. M. Hasan, S. Cirstea, and M. N. Shraboni: 2023 15th Int. Conf. Software, Knowledge, Information Management and Applications (IEEE, 2023) 231.
16  Y. Wu, Y. Jin, A. Yokoyama, C. Wang, Y. Li, Y. Liu, and G. Wang: 2019 IEEE Biomedical Circuits and Systems Conf. (IEEE, 2019) 1.
17  A. K. F. Lui, K. F. Wong, S. C. Ng, and K. H. Law: 2012 6th Int. Conf. Distributed Smart Cameras (IEEE, 2012) 1.

## About the Authors

**Bo-Heng Chen** received his B.S. degree in electrical engineering from Chung Yuan Christian University, Taoyuan, Taiwan, in 2022 and is currently pursuing his M.S. degree in electrical engineering from the same university. His research interests include signal processing, machine learning, and AIoT. (rich11319@cycu.org.tw)

**Da-Chuan Chen** received his B.S. degree in electrical engineering from Chung Yuan Christian University, Taoyuan, Taiwan, in 2023 and is currently pursuing his M.S. degree in electrical engineering in the same university. His research interests include signal processing, machine learning, and AIoT. (dachuan516@gmail.com)

**Wen-Hsiang Yeh** received his B.S. and M.S. degrees in electrical engineering, from Chung Yuan University, Taiwan, in 2021 and 2023, respectively. His current research interests include artificial intelligence applications, embedded systems, and deep learning. (peter880127@gmail.com)

**Hao-Yang Chen** is currently pursuing his B.S. degree in electrical engineering from Chung Yuan Christian University, Taoyuan, Taiwan. His research interests include machine learning and AIoT. (chenhaoyang921@gmail.com)

**Yu-Ping Liao** received her B.S. degree in physics and M.S. and Ph.D. degrees in electrical engineering from the National Taiwan University, Taipei, Taiwan, in 1981, 1983, and 1985, respectively. She is currently a professor in the Department of Electrical Engineering at Chung Yuan Christian University, Taiwan. She is the author/coauthor of more than 10 textbooks. Her research interests are in artificial intelligence over Internet of Things systems and FPGA applications. (lyp@cycu.org.tw)