

**Robot Action Recognition Using a Multi-Modal Information Fusion Model**

<sup>1</sup>BURRI NARENDER REDDY, M.Tech Assistant Professor, [narendarburri@gmail.com](mailto:narendarburri@gmail.com)

<sup>2</sup>ACHINA VENKATARAMANA, M.Tech Assistant Professor, [venkatachina5@gmail.com](mailto:venkatachina5@gmail.com)

Department-ECE

Pallavi Engineering College Hyderabad, Telangana 501505.

**ABSTRACT**

*Modern robotics field, driven by cutting-edge IT, is entering a new era of human-machine cooperation. In contrast to older robots, which need isolation rails to keep a safe distance from humans, the latest human-machine cooperation systems are designed to function in close proximity to humans. Maximizing productivity by intelligently allocating jobs to humans and robots, thereby capitalizing on their respective strengths. Approaches that lead to greater efficacy*

**Introduction**

Human-robot interaction systems have emerged as a result of technological development. Achieving human-robot interaction requires two things: (1) the capacity of humans and robots to communicate and (2) the presence of a human or other socially acceptable robot. Is that the robot can interpret human speech, gestures, intentions, and other human cues with remarkable precision [1-3]. At the moment, the intelligence is often used as the medium for human-machine interaction. Robots. The usage of intelligent robots has expanded to various industries; including the three major economic pillars of the country are the service sector, manufacturing, and farming. Artificial intelligence (AI) robots are

cooperative human-robot interaction entering a new phase business strategy steered by the latest wave of info-science and technology. The current generation of robots is dependent on make use of separation fences to keep a safe distance new generation of human-machine cooperation technologies may assist humans in their job rather than taking over removing physical constraints and maximizing the contributions of everybody involved people and machinery by strategically allocating operational activities and enhancing work procedures to get more productivity efficiency. Robots, for instance, will be used in the near future in many types of settings, including manufacturing. Being held accountable for arduous, perhaps harmful, and routine duties, freeing up people to concentrate on more interesting, dynamic, and creative work. Planning, or doing any kind of task that requires adaptability and fortitude. Essential to the development of a human-robot interaction system optimize human-robot communication and collaboration. This is in stark contrast to the keyboard, mouse, and other interface devices often used in traditional human-computer interaction. cumbersome, limits people's freedom of movement, and fails to ease the burden of labour for individuals. People with vocal capabilities In order to direct the robot's behaviour toward a certain goal, one must use gestures or activities, it may significantly

**International Conference Latest Studies In Engineering Research**

reduce the amount of manual labour required from staff. Work. Irrespective of whether the robot is operated by human speech, gesture, crucial criterion is for the robot to be able to mimic human accept instructions from humans. Accurate and fast motion detection is crucial for efficient human-robot interaction in the area of motion-based robotics.

### **Related Knowledge**

The Basic Robotics Operating Principle. Hardware systems on robots often consist of things like cameras, sensors, servos, and the like, whereas software systems are responsible for things like programming and controlling the robot. Components such as circuit boards, communication modules, bodies, and appendages four wheels serve as replacements for the lower legs. Wheels have the abilities to go forward, reverse, turn left or right, and stop. There are four degrees of freedom, which allow for a wide range of motions. An utmost Arms and legs may swing back and forth to lift, move, and position objects. Flexing and gripping. independence The STM32 development board is based on the traditional Cortex CPU, Comprehensive Software Bundle, Extensive large number of supported chip models, a big complement of cost-effective peripherals, affordable in terms of both cost and energy use, as well as having a large number of users, the robot presented in this work was primarily use the STM32 development board for microcontrollers. The STM32 IDE (integrated development environment) A C/C++ development environment called STM32CubeIDE has a peripheral configuration, coding language. \*is functions for code-generation, code-compilation, and code-debugging

Microcontrollers and processors based on the STM32 core. This is the official ST acronym. Libraries for STM32 peripherals, including Using pre-existing device libraries streamlines the process of Construction of a Project. In the meanwhile, ST is used for scripting. Uses a thin layer of abstraction to hide the inner workings of the peripheral library from developers. STM32 internal register structure \*e software architecture figure 1 shows the robot's process flow

### **Algorithm**

Constructed from Algorithms. Applications such as computer vision and human-computer interaction rely on HAR as their basis. Its primary function is to make it possible for the computer to identify a variety of human gestures, such as the left foot, and making a fist with both hands and elevating them over the head are all examples of HAR gestures. The video's subject matter using the predetermined categories. Which kind of behaviour is represented by the given data? Video and report the verdict. Exhibits by Human Beings a wide range of behavioural patterns throughout the course of their task in addition to everyday existence. A wide variety of physical activities, Actions taken by humans vary depending on the context. Various muscle-upper-body motion combinations Human actions are shown through the upper and lower limbs. In This study suggests a method in order to detect such changes. By combining a deep learning model with multifeature fusion for recognition. Algorithm for learning. Schematic of the structure shown in Figure 3. A technique for recognizing actions has been suggested in this work.

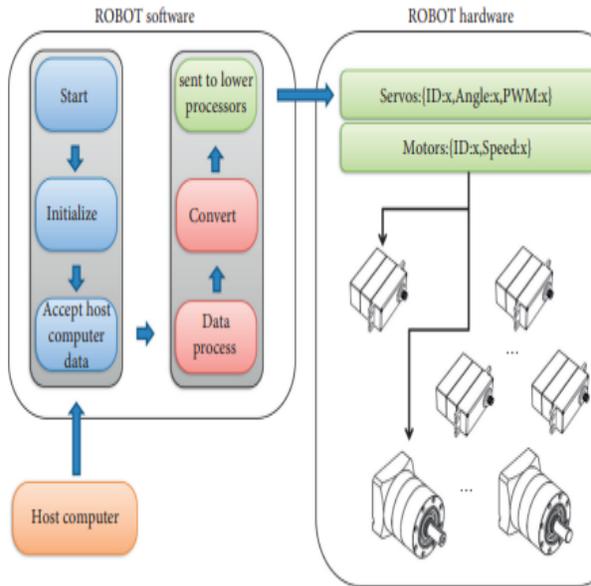


Figure 1: Principle of robot operation

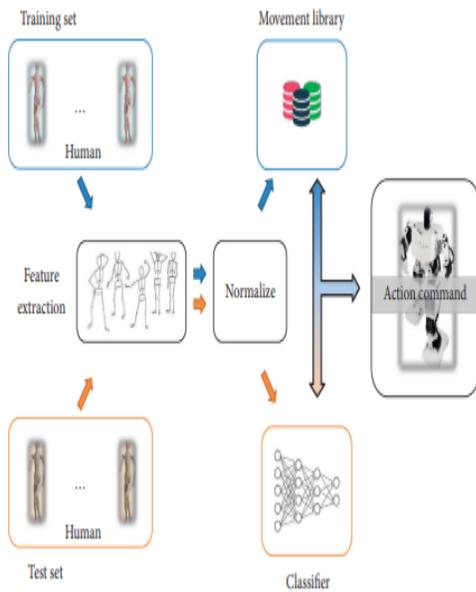


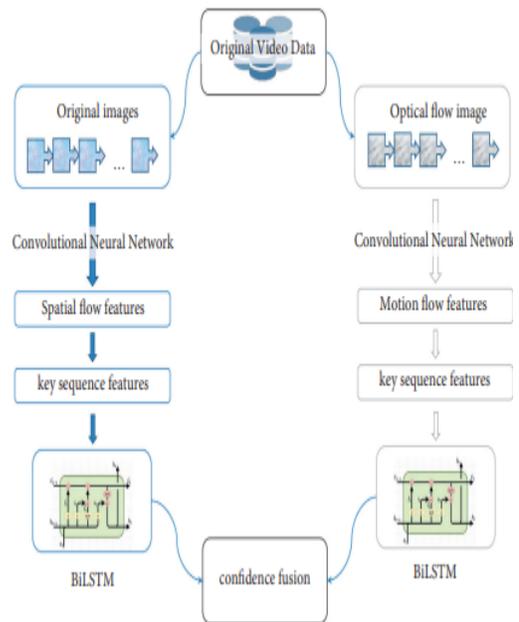
Figure 2: Action recognition process

For the purpose of doing the same operation on each element of the sequence, \*e RNN cells are connected to create a loop. Input is chosen at random. Sequence  $x(x_1, x_2, \dots, x_t)$  is an embedding that contains a secretive cell ( $h$ ) and a result ( $y$ ), with  $T$  denoting the end of the process. At each step. RNN's hidden state  $h_t$  is determined. Depending on the revealed state  $h_{t-1}$  and the current new data coming in  $x_t$ , with the hidden layer and the output layers being computed in this way

$$h(t) = \text{sig}(W_1 x_t + W_2 h_{t-1}), \quad (1)$$

$$y(t) = g(V h_t),$$

Where  $W_1$  and  $W_2$  are network weight matrices,  $\text{sig}$  and  $\text{sof}$  are sigmoid and softmax activation matrices, and  $a$  and  $b$  are bias and variance vectors the activation and inhibition functions are determined...like so



TM Figure 3: HAR framework

$$\text{sig}(z) = \frac{1}{1 + e^{-z}},$$

$$\text{sof}(z_m) = \frac{e^{z_m}}{\sum_k e^{z_k}}. \tag{2}$$

Long Short-Term Memory (LSTM) uses an input, an output, and a forget gate to update its hidden and storage units depending on the input and the output, respectively. In LSTM, the RNN.\*e formula is defined as follows:

$$\begin{aligned} i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i), \\ f_t &= \sigma(W_f \cdot [h_{t-1}, x_t] + b_f), \\ \tilde{c}_t &= \tanh(W_c \cdot [h_{t-1}, x_t] + b_c), \\ c_t &= f_t \cdot c_{t-1} + i_t \cdot \tilde{c}_t, \\ o_t &= \sigma(W_o \cdot [h_{t-1}, x_t] + b_o), \end{aligned} \tag{3}$$

Where  $i_t$ ,  $f_t$ ,  $o_t$ , and  $c_t$  represent the input, forget, and output gates of the cell and the output at time  $t$ . At time point  $t$ , the  $x_t$  input variables and the  $h_t$  secret variables are as follows. The symbol for the sigmoid activation function is,  $\sigma$ .  $W$  stands for both the weight vector and the bias vector. and  $b$ . A diagram of a long short-term memory (LSTM) cell is shown in Figure 5. Uses essentially three fundamental gates and a c

BiLSTM is a multi-directional LSTM that takes use of both forward and reverse LSTMs. Aiming to model present and to be determined later on based on circumstances. When compared to the LSTM network, Both the input and output layers of a BiLSTM network are path of least resistance, and forward and reverse propagation Front and rear passes are completed for each layer initially. Neural networks do with other sections of the input

sequence, networks function, the BiLSTM model is able to maintain the both the front and rear sequences have unique sets of information. Directions. This system is capable of completely considering contextual information. In a network with two LSTM layers, the vector formulation is modified in the following ways:

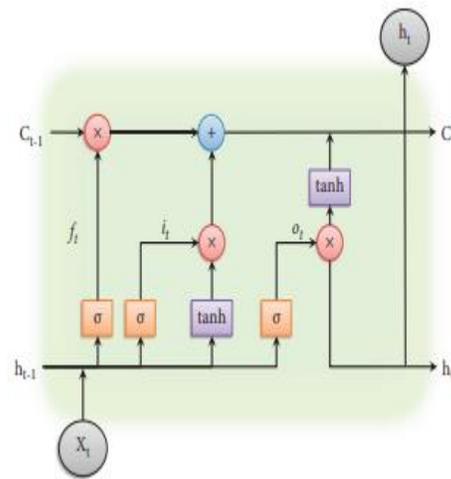


Figure 5: LSTM network structure

$$\begin{aligned} h_{ft} &= H(W_{xh_f}x_t + W_{h_f h_f}h_{f,t-1} + b_{h_f}), \\ h_{bt} &= H(W_{xh_b}x_t + W_{h_b h_b}h_{b,t-1} + b_{h_b}), \end{aligned} \tag{4}$$

Where  $h_b$  Rd and  $h_f$  Rd are the reverse and forward outputs, respectively, and  $e$  is the ultimate result  $y_t$  [ $h_{ft}$ ,  $h_{bt}$ ]. And  $y_t > R2d$  are joined to form. Hybrid of the front and rear layers for use as a just one BiLSTM layer. When using video data for action recognition, for this reason, BiLSTM is used to record the time dependency of sequences, and acquire knowledge of the historical context of various types of actions performed across varying amounts of time. In the study, the data comes from the last layer of pooling.

To the BiLSTM network from the Inception-V3 network. Acquire the ability to decode the information contained in various video clips. Last key subsequence layer's output is described as

$$\mathbf{x}' = [\mathbf{x}^{i1}, \mathbf{x}^{i2}, \dots, \mathbf{x}^{in}]. \quad (5)$$

The number of frames in the it key subsequence is denoted by in in the following equation: \*e sequence characteristics acquired by the BiLSTM network are defined as

$$r = \text{BiLSTM}(\mathbf{x}'). \quad (6)$$

$$s = \text{soft max}(r^i). \quad (7)$$

Figure 6 shows a graphical representation of the BiLSTM model's structural details. Model

Complementary effects of independently collected characteristics from the data are the basis of the fusion strategy used to get these results. Spatial flow characteristics, time-series features and motion-flow features are

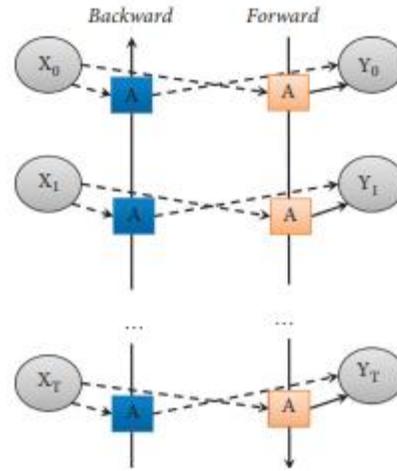


Figure 6: BiLSTM network structure

Complimentary. Choosing a suitable fusion approach for\she classification results derived from diverse characteristics may considerably enhance the accuracy of action identification. Confidence level of a classifier is a crucial metric for assessing the efficacy of the categorization job, as it establishes the refuse-recognition threshold and is critical in the integration of many classifiers. Confidence level utilized in this study I

$$z_v(x) = (1 - \alpha)(P_v^{\max}(x) - P_v^{\text{sub max}}(x)) + \alpha \left( P_v^{\max}(x) - \frac{1}{n-1} \sum_{j=1}^{c-1} p_{v,j}(x) \right) \quad (8)$$

$$s.t. \quad p_{v,j}(x) \neq P_v^{\max}(x),$$

$$y(x) = y_s(x) * z_s(x) + y_m(x) * z_m(x). \quad (9)$$

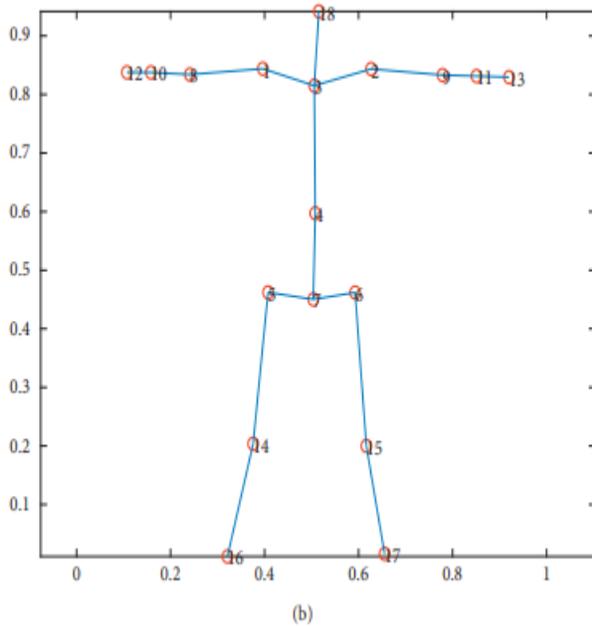
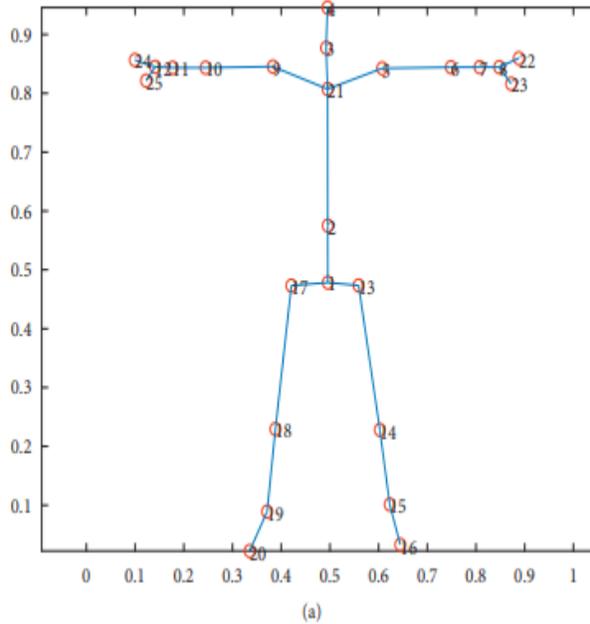


Figure 7: Schematic diagram of the joints of the two data sets. (a) NTU-RGB + D. (b) MSR Action 3D.

The vector of the category scores is recorded as  $y(x) = (q_1, q_2, \dots, q_c)$ , in addition, Sample x has the label Action Category:

$$Y(x) = \arg \max_j q_j \tag{10}$$

**Conclusion**

Robots are already used in numerous sectors, including manufacturing, farming, and the service sector because to advancements in AI technology. Common applications of robots to carry out a task that is time-consuming, physically demanding, or potentially hazardous jobs that people aren't good at, so that we may concentrate on what we are on activities that call for adaptability or improvisation toughness. In terms of helping, conventional robots aren't as workers since they can't help humans with their jobs without the constraints of location or time. To accomplish a smarter robot-human communication, better suited techniques for people to influence machines need to be rapidly researched. Therefore, the study of robotics based on \sHAR has evolved. Accurate, in-the-moment machine identification of human behaviours is becoming more important to enhance the effectiveness of human-robot interaction. Considering that numerous characteristics may completely define data, to enhance action recognition precision, in this analysis, we disentangle the spatial flow characteristics of the dataset and motion flow features, inputs each feature independently into To get the best results, BiLSTM use the confidence fusion technique to build a model that can predict labelled data.

**References**

- [1] M. Kraus, N. Wagner, Z. Colleges, and W. Minker, "The role of trust in proactive conversational assistants," *IEEE Access*, vol. 9, pp. 112821–112836, 2021.
- [2] J. Pustejovsky and N. Krishnaswamy, "Situating meaning in multimodal dialogue: human-robot and human-computer interactions," *Traitement Automatique des Langues*, vol. 61, no. 3, pp. 17–41, 2020.
- [3] M. Jarosz, P. Nawrocki, B. Snieżyński, and B. Indurkha, "MULTI-PLATFORM intelligent system for multimodal human-computer interaction," *Computing and Informatics*, vol. 40, no. 1, pp. 83–103, 2021.
- [4] G. L. Sravanthi, M. V. Devi, K. S. Sandeep, A. Naresh, and A. P. Gopi, "An efficient classifier using machine learning technique for individual action identification," *International Journal of Advanced Computer Science and Applications*, vol. 11, no. 6, pp. 513–520, 2020.
- [5] N. A. Othman and I. Aydin, "Challenges and limitations in human action recognition on unmanned aerial vehicles: a comprehensive survey," *Traitement du Signal*, vol. 38, no. 5, pp. 1403–1411, 2021.
- [6] P. Gupta, A. Atipelli, A. Aggarwal et al., "Quo vadis, skeleton action recognition?" *International Journal of Computer Vision*, vol. 129, no. 7, pp. 2097–2112, 2021.
- [7] M. La Cascia, "3D skeleton-based human action classification: a survey," *Pattern Recognition*, vol. 53, no. 53, pp. 130–147, 2016.