

Research Article

Dynamic Gesture Recognition Algorithm Based on 3D Convolutional Neural Network

Yuting Liu ^{1,2}, Du Jiang ^{1,2}, Haojie Duan ^{2,3}, Ying Sun,^{1,2,3} Gongfa Li ^{1,2,3}, Bo Tao ^{1,2}, Juntong Yun ^{2,3}, Ying Liu ^{1,3} and Baojia Chen⁴

¹Key Laboratory of Metallurgical Equipment and Control Technology of Ministry of Education, Wuhan University of Science and Technology, Wuhan 430081, China

²Research Center for Biomimetic Robot and Intelligent Measurement and Control, Wuhan University of Science and Technology, Wuhan 430081, China

³Hubei Key Laboratory of Mechanical Transmission and Manufacturing Engineering, Wuhan University of Science and Technology, Wuhan 430081, China

⁴Hubei Key Laboratory of Hydroelectric Machinery Design & Maintenance, China Three Gorges University, Yichang 443002, China

Correspondence should be addressed to Du Jiang; jiangdu@wust.edu.cn and Gongfa Li; ligongfa@wust.edu.cn

Received 13 July 2021; Accepted 6 August 2021; Published 17 August 2021

Academic Editor: Syed Hassan Ahmed

Copyright © 2021 Yuting Liu et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Gesture recognition is one of the important ways of human-computer interaction, which is mainly detected by visual technology. The temporal and spatial features are extracted by convolution of the video containing gesture. However, compared with the convolution calculation of a single image, multiframe image of dynamic gestures has more computation, more complex feature extraction, and more network parameters, which affects the recognition efficiency and real-time performance of the model. To solve above problems, a dynamic gesture recognition model based on CBAM-C3D is proposed. Key frame extraction technology, multimodal joint training, and network optimization with BN layer are used for making the network performance better. The experiments show that the recognition accuracy of the proposed 3D convolutional neural network combined with attention mechanism reaches 72.4% on EgoGesture dataset, which is improved greatly compared with the current main dynamic gesture recognition methods, and the effectiveness of the proposed algorithm is verified.

1. Introduction

As a common means of communication in people's daily life, gestures also have great application in human-computer interaction. Compared with expressions, actions, and other interactive means, gestures are more intuitive, natural, and comfortable. Therefore, gesture communication is also the most used human-computer interaction means besides language [1, 2]. Gestures have a wide range of applications in human-computer interaction technology, such as intelligent driving, virtual reality, augmented reality [3, 4], medical assistance [5], and so on [6–8]. In intelligent driving, gestures are captured by vehicle intelligent control systems and analyzed by intelligent center. Instructions are sent out to

complete the human control of vehicle navigation and entertainment functions [9, 10]. When talking about virtual reality and argument reality, Microsoft's HoloLens has already realized having entertainment of users in the virtual reality environment through both hands [11]. In medical assistances, gesture recognition can provide assistances for the hearing-impaired groups and realize the normal communication between deaf and dumb people.

Gestures could be mainly divided into dynamic and static. Static gestures focus on the hand posture and shape at a single point in time, such as gesture action "OK." Only the spatial features of gestures are considered in recognition. Dynamic gesture recognition should not only consider hand postures and shapes but also pay attention to the spatial

displacement and spatiotemporal correlation of the whole gesture [12–14]. Compared with static gesture recognition, dynamic gesture is closer to people’s expression habits with more abundant information expression [15–17], which has more practical significance. Nowadays, researchers have proposed a variety of dynamic gesture recognition algorithms, including dynamic gesture feature extraction algorithm such as MEI algorithm, HOG algorithm, and HOF [18] algorithm and classification algorithm such as hidden Markov model [19]. With the development of deep learning, many video classification algorithms, for example, C3D [20] and dual stream convolution network and LSTM [21], have been applied to dynamic gesture recognition [22–24], achieving high recognition rate. However, the amount of network inputs is in large scale because of the need to extract the video spatial information and temporal information, which result in huge number of parameters and calculation. Such networks have complex network structure and low real-time performance. It is possible to increase the effect of dynamic gesture recognition by optimizing input and improving the existing feature extraction methods.

In this manuscript, we propose a dynamic gesture recognition algorithm based on attention mechanism of 3D convolutional neural network, which has several contributions:

- (a) Interframe difference method is optimized. Input video data are processed to improve the problem of data redundancy and format inconsistency.
- (b) 3D convolutional neural network is combined with attention mechanism. CBAM is used to optimize the structure of 3D convolutional neural network to reduce the transmission loss of input information and realize the feature extraction of spatial dimension and time dimension.
- (c) Multimodal joint is applied to train neural network. To improve the effect of gesture recognition, the fusion method of dual-mode feature input is used to realize the feature complementarity of the two modes.

2. Related Work

As deep learning develops by leaps and bounds, computer vision has been promoted. Many excellent image analysis and recognition algorithms have been proposed. For example, the Alexnet which is designed by Srivastava [25] has achieved the best recognition performance in the Imagenet, an image recognition competition [26, 27]. Different from traditional methods of artificial design features, deep learning automatically extracts features through convolution neural network. By training and debugging the feature extraction network, more critical and representative spatiotemporal features can be extracted by deep learning for video classification and action recognition [28, 29]. The dynamic gesture recognition networks based on deep learning are mainly divided into three types: two-stream networks, long short-term memory (LSTM) network, and three-dimensional convolutional neural network (3D-CNN).

The concept of two-stream network was first proposed by Simonyan and Andrew [30] in 2014, and the optimal recognition effect was obtained in the behaviour recognition task of open data sets UCF-101 [31] and HMDB-5 [32]. The two-stream network algorithm also has some defects since the algorithm only gets spatial information through a single image. It is difficult to deal with the large changes in behaviour, and the optical flow image is only suitable for the small changes in motion information capture. Aiming at the problem of two-stream network, Wang et al. [33] designed a temporal segment network (TSN) to sparsely sample long time series images and obtain more robust spatiotemporal features through Inception v2 to improve the effect of action recognition. Based on the work of Wang et al., Feichtenhofer et al. [34] studied the method of fusing spatial and temporal information and found that the feature fusion in higher convolution layer of network has better recognition effect. In two-stream convolution networks, the operation of optical flow information could occupy a lot of memory and affect the recognition rate. Zhu et al. [35] designed a convolution network for optical flow estimation instead of optical flow information operation, cascaded temporal information network and spatial information network, and used multiple images stacked input to complete action recognition. Dynamic gesture recognition [36] is like action recognition. It also uses the algorithm to obtain the spatial and temporal information of the object expression in the video to realize the video action understanding.

LSTM is actually a type of recurrent neural network (RNN). Inputs of each layer of RNN consist of the output of the upper layer and the output of the same layer, and outputs of the neuron are the inputs of the same layer. Therefore, RNN could effectively deal with the problem of temporal feature extraction. However, the network can only solve the problem of short time series due to the limitation of structure. LSTM is designed to solve problem of long time series and historical information loss in the iteration. For the problem of information loss in long time sequence, LSTM controls the information processing of neurons with three structures: input gate layer, forget gate layer. and output gate layer. In the dynamic gesture recognition, LSTM uses the common convolutional network to extract the features, serializes the spatial features extracted by the previous network through LSTM, and then classifies them in the full connection layer.

3D-CNN is an improvement of convolution kernel and pooling method on traditional 2D convolution neural network. Continuous motions contain unique temporal information. However, 2D convolution kernel can only extract spatial information from image [37, 38], and 3D convolution kernel is designed to extract features from continuous image to obtain temporal information. Spatial scale pooling and channel scale pooling are also included in pooling process. Many scholars have studied 3D convolutional neural network. For example, Tran et al. [39] proposed C3D network to realize dynamic gesture recognition. The I3D network is proposed by Carreira and Zisserman [40].

Some previous excellent models are used in this algorithm. To solve the problem of excessive computing cost, Qiu et al. [41] proposed the P3D network, which optimized the convolution model.

Moreover, some algorithms use feature fusion to integrate image information and optical flow information, such as dual stream algorithm and MFFs-net algorithm [42]. The accuracy of the MFFs-net algorithm on Jester dataset is 96.28%, but the calculation of optical flow needs a lot of computing resources. Molchanov team has also done a lot of work in dynamic gesture recognition. In paper [25], the team proposed to integrate RGB image, depth image, and radar data to realize dynamic gesture recognition. In paper [43], the team proposed using 3D-CNN to train two different resolution networks and fusing the recognition results to improve the recognition accuracy. In paper [44], the team also used residual neural network to optimize 3D-CNN and verified the effectiveness of its model in SKIG data set and CharLearn2014 data set. Since the sign language data set also contains a large number of dynamic gestures, most dynamic gesture recognition models also use sign language as the data set for training and testing. In paper [45], a dual stream 3D-CNN is designed to realize dynamic gesture recognition based on effective fusion of multimodal data.

Based on the analysis of the above current research situation, it can be found that dynamic gesture recognition is still a research hotspot in the field of computer vision at present. Although it is still in the initial stage of this technology, the recognition of simple gesture has achieved good results in daily application. There are still some problems to be further studied and optimized.

3. 3D Convolutional Neural Network Based on Attention Mechanism

3.1. Structure of Dynamic Gesture Recognition Algorithm. Three aspects should be considered in motion recognition of video with single gesture:

- (a) Appearance and texture feature of gestures
- (b) Changes of gesture features, namely, gesture space features
- (c) Time domain information between images, which is spatiotemporal characteristics of continuous changing gestures

In view of the above three aspects, this paper takes RGB image and depth image as input and designs a dynamic recognition model of three-dimensional convolutional neural network combined with convolutional block attention module (CBAM-C3D). The process (shown in Figure 1) is as follows: firstly, a key frame extraction method based on interframe difference method is designed to process the original input video. In this way, redundant frames of network input are reduced, and scale alignments of input data are realized. Then, the processed depth images and RGB images data are input into CBAM-C3D. Finally, the two data features are fused in series in the feature layer to complete the dynamic gesture recognition.

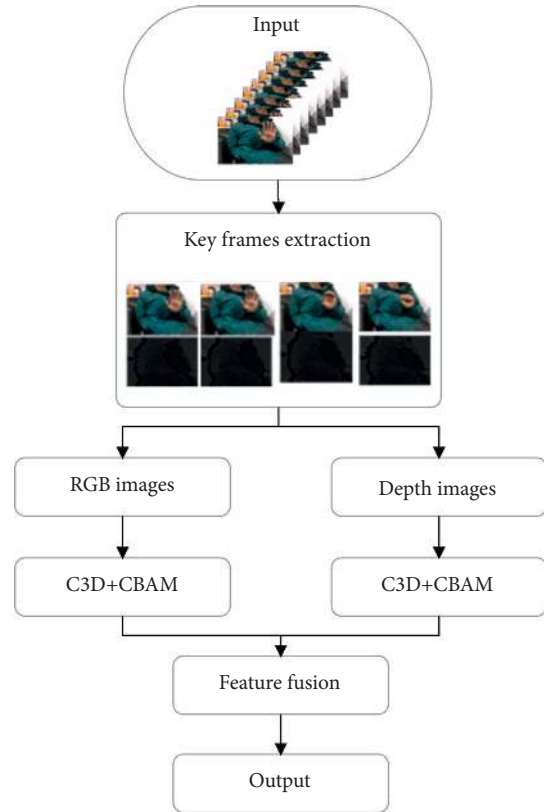


FIGURE 1: The overall framework of the proposed algorithm.

The proposed CBAM-C3D algorithm shown in Figure 2 is optimized according to the C3D network's structure which was proposed by Du et al. Batch normalization and ReLU layer are added into 3D convolution layer. The full connection layer and maximum pooling layer relate to CBAM network to optimize the features. This fusion network can not only reduce transmission loss of input information but also automatically learn important spatiotemporal information contained in images.

3.2. Key Frames Extraction of Dynamic Gesture. Due to the disunity of action standards and personal physical factors, the duration of the same action and gesture changes may be greatly different, which is also one of the difficulties of gesture recognition. There are two cases for the inconsistency of action duration:

- (a) *Redundancy of Frames.* The input video sequences usually have many redundant frames without gesture action. These redundant images could greatly reduce the calculation speed of the network and even affect the recognition effect of the network.
- (b) The time distribution of gesture movement is uneven. For complex gestures, it takes a long time to complete multiple gestures. In general, it cannot guarantee that the actions are carried out at a uniform speed, which makes the recognition process of dynamic gestures more difficult.

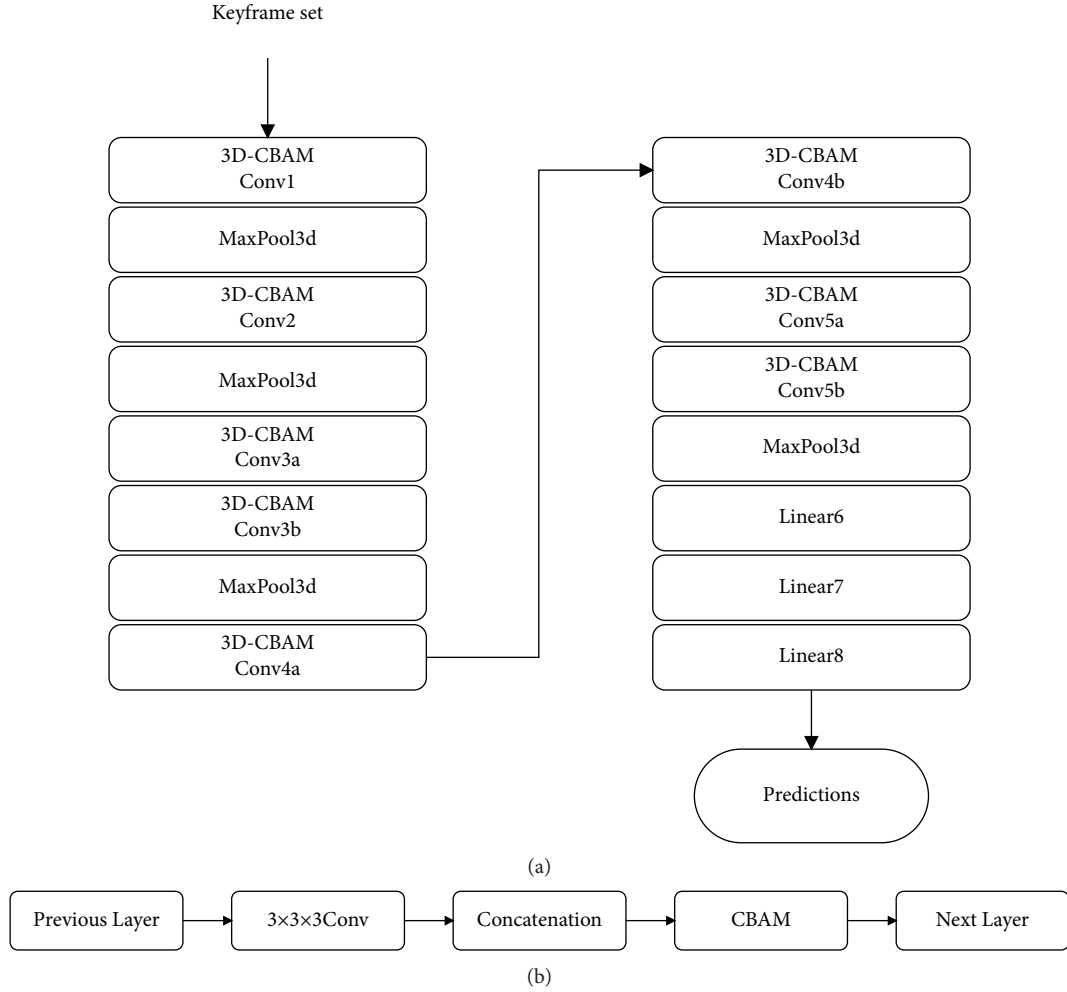


FIGURE 2: Structure of the proposed CBAM-C3D networks. (a) is the connection methods. (b) is the structure of 3D-CBAM in the proposed networks.

For these problems, the method of interframe difference to unify the scale of video data and simplify the data processing is proposed. The whole video is defined by several representative images.

The traditional inferframe difference method is mainly used for moving target monitoring. In this manuscript, it is optimized to obtain more accurate key frame images. Figure 3 shows the calculation process of the inferframe difference method. Firstly, RGB-D image is used to segment the gesture area to obtain the hand image with the background removed. Then, the adjacent image pixel standard deviation algorithm is used to calculate the inferframe difference of the adjacent image in the image sequence. Finally, the size of the inferframe difference is sorted to complete the key frame extraction.

The standard deviation of inferframe difference L_n is the evaluation standard of key frames. For example, the number of key frames K is preset, and the standard deviation of the gray value change of the frame n image is calculated. The continuous images of the input video sequence are supposed to be f_n and f_{n+1} . The pixels on images are (x, y) . Gray values of the corresponding image are $f_n(x, y)$ and $f_{n+1}(x, y)$. According to formula (1), f_n^i represents the gray value of pixel i of image n :

$$L(f_n, f_{n+1}) = \sqrt{\sum_{i=1}^n (f_n^i - f_{n+1}^i)^2}. \quad (1)$$

The maximum and minimum values of the sequence frame difference are counted, and intermediate value $\text{mid}(L)$ is calculated according to formula (2). All local extremums less than $\text{mid}(L)$ are removed and assumed that the number of remaining is m . Finally, the extracted extremum points are sorted, in which the frames corresponding to the first K extremum points are taken as the key frames. If $m \leq K$, the last in the sequence is copied and filled based on m images.

$$\text{mid}(L) = \frac{(\max(L) + \min(L))}{2}. \quad (2)$$

3.3. Optimization of Three-Dimensional Convolutional Neural Network

3.3.1. Optimization of Network Feature Extraction. Compared with SENet which only focuses on channel information extraction, CBAM considers the spatial

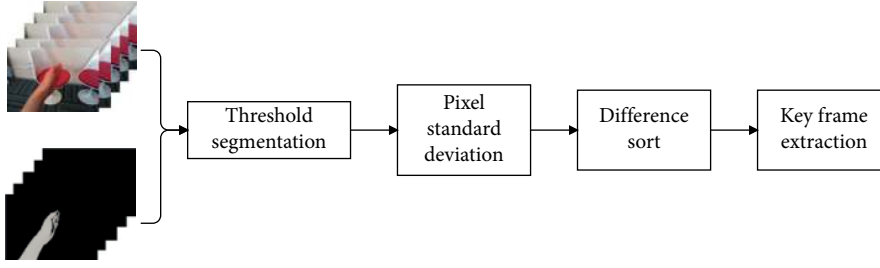


FIGURE 3: Flow chart of the interframe difference method in video sequences.

information and channel information in feature extraction, that is, temporal features. In this paper, CBAM is used to optimize the structure of three-dimensional convolutional neural network, complete the important feature extraction of spatial dimension and time dimension, and strengthen the effect of network feature extraction.

$$\begin{aligned} I' &= A_c(I) \otimes I, \\ I'' &= A_s(I') \otimes I', \end{aligned} \quad (3)$$

where $I \in R^{C \times H \times W}$ is feature map of input and $A_c \in R^{C \times 1 \times 1}$ and $A_s \in R^{1 \times H \times W}$ are one-dimensional and two-dimensional channel attention map, respectively.

Two vectors with only channel dimension are obtained by maximum pooling and average pooling in the spatial dimension. Then, the features are added and sigmoid activated by a two-layer neural network. The input of spatial attention processing can be obtained by multiplying the channel attention vector with the feature graph as follows:

$$\begin{aligned} A_c(I) &= \sigma(\text{MLP}(\text{AvgPool}(I)) + \text{MLP}(\text{MaxPool}(I))) \\ &= \sigma(W_1(W_0(I_{\text{avg}}^c)) + W_1(W_0(I_{\text{max}}^c))), \end{aligned} \quad (4)$$

where σ is the sigmoid activation function. $\text{AvgPool}(\)$ and $\text{MaxPool}(\)$ represent average pooling and maximum, respectively. $W_0 \in R^{C/r \times C}$ and $W_1 \in R^{C \times C/r}$ are weights of MLP. This module is equivalent to a filter, the important channel weight is larger, and the unimportant channel weight is smaller. Therefore, it realizes the attention mechanism in the feature dimension.

To calculate spatial attention which concerns on the position of useful information, the maximum pooling and average pooling on the channel with only two channel dimensions are acquired. Then, the two vectors are input into a two-layer neural network, respectively. After feature addition, it is input into the convolution layer for weight optimization. In this way, a spatial attention filter is generated as follows:

$$\begin{aligned} A_s(I') &= \sigma(f^{7 \times 7}([\text{AvgPool}(I'); \text{MaxPool}(I')])) \\ &= \sigma(f^{7 \times 7}(I_{\text{avg}}^s; I_{\text{max}}^s)), \end{aligned} \quad (5)$$

where f is the convolution operation.

The channel attention feature is input to the spatial attention mechanism network, and the spatial attention mechanism network is used to complete the feature

extraction including space and temporal feature. Temporal attention concerns more on global information, while spatial attention focuses on local information. Therefore, the combination of these two attention modules will effectively extract salient features and enhance the expression of features.

3.3.2. Optimization of Network Training Process. Due to the trend of polarization distribution and uneven distribution of data, gradient may disappear or explode. To solve this problem, BN layer is added to ensure the consistency of input data in each layer. The specific algorithm is as follows:

$$\mu_\beta = \frac{1}{n} \sum_{i=1}^n x_i, \quad (6)$$

$$\sigma_\beta^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_\beta)^2, \quad (7)$$

$$\hat{x}_i = \frac{x_i - \mu_\beta}{\sqrt{\sigma_\beta^2 + \varepsilon}}, \quad (8)$$

$$y_i = \gamma \hat{x}_i + \beta. \quad (9)$$

Formulas (6) and (7), respectively, calculate the mean value and variance of the data, formula (8) standardizes the data, and formula (9) offsets the data. After the data are processed by BN layer, the distribution of data will be more uniform, which is also conducive to improving the generalization performance of the network.

3.3.3. Feature Fusion Strategy for Bimodal Data. In the field of behaviour recognition, many researchers have tried to use multimodal data as input to improve the recognition effect such as RGB image, depth image, and optical flow data. All kinds of modal data contain important information of gesture recognition, but the results of single-modal data recognition are not good enough. Therefore, the recognition efficiency of dynamic gesture is improved in this paper by fusing multiple modal data. Considering that 50 kinds of actions with significant difference in gesture are included, RGB images are used to fuse depth images for gesture recognition.

In this paper, RGB image and depth image are used as input data, and the two kinds of modal information are input





Label	Illustration	Instruction	Label	Illustration	Instruction
1		Scroll hand towards right	9		Zoom out with fists
2		Scroll hand towards left	10		Rotate fists clockwise
3		Scroll hand downward	11		Rotate fists counterclockwise
4		Scroll hand upward	12		Zoom in with fingers
5		Scroll hand forward	13		Zoom out with fingers
6		Scroll hand backward	14		Rotate fingers clockwise
7		Cross index fingers	15		Rotate fingers counterclockwise
8		Zoom in with fists	16		Click with index finger

FIGURE 4: Schematic diagram of some gesture categories in EgoGesture dataset.

into the CBAM-C3D model, respectively, by the way of training, and then the fusion is carried out after the respective features are obtained. Assume that the output feature vectors of the two kinds of modal information are F_{RGB} and F_{Depth} , respectively, after feature extraction, and the final fusion feature vector is F_u .

(a) Average fusion:

$$F_u = \frac{1}{2}(F_{\text{RGB}} + F_{\text{Depth}}). \quad (10)$$

(b) Series fusion:

$$F_u = F_{\text{RGB}} \otimes F_{\text{Depth}}, \quad (11)$$

where \otimes represents a tandem operation.

4. Experiment

4.1. Experimental Dataset. Different countries have different definitions of gesture, so there is no recognized dynamic gesture data set in this field. In order to facilitate multimodal data fusion training, this paper selects the EgoGesture [46] first person gesture database released by the State Key Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences in 2018 (Figure 4).

The dataset contains 2081 RGB-D videos, 24161 gesture samples, and 2953224 frames from six different themes. Each video sample is taken by Intel Realsense sr300 camera, and the data format is RGB-D. Each frame of video is 640×480 pixel resolution and 30 fps recording. There are 33 kinds of static and 50 kinds of dynamic gestures collected by 50 people from six different indoor and outdoor scenes. Figure 5 is a schematic diagram of some gesture categories in EgoGesture database. As this paper mainly focuses on dynamic gesture recognition, we select 50 dynamic gestures and their labels as RGB data and depth data input to filter and adjust the EgoGesture database. Some samples of RGB and depth video data are taken as examples.

4.2. Experimental Environment and Training Parameters. All the experiments are carried out in Window10 system. The graphics card is NVIDIA gtx3060ti 8g. The running software environment is Python 3.6, Python-1.3.0 + torch vision-0.5.0, OpenCV-Python-4.5.0, and other auxiliary Python libraries. The data input is EgoGesture dataset, and the training set, test set, and verification set are divided according to the ratio of 3:1:1. When training the model, the model is verified and adjusted every 20 steps. Before network training, key frames are extracted from RGB images. RGB images and depth images are selected according to

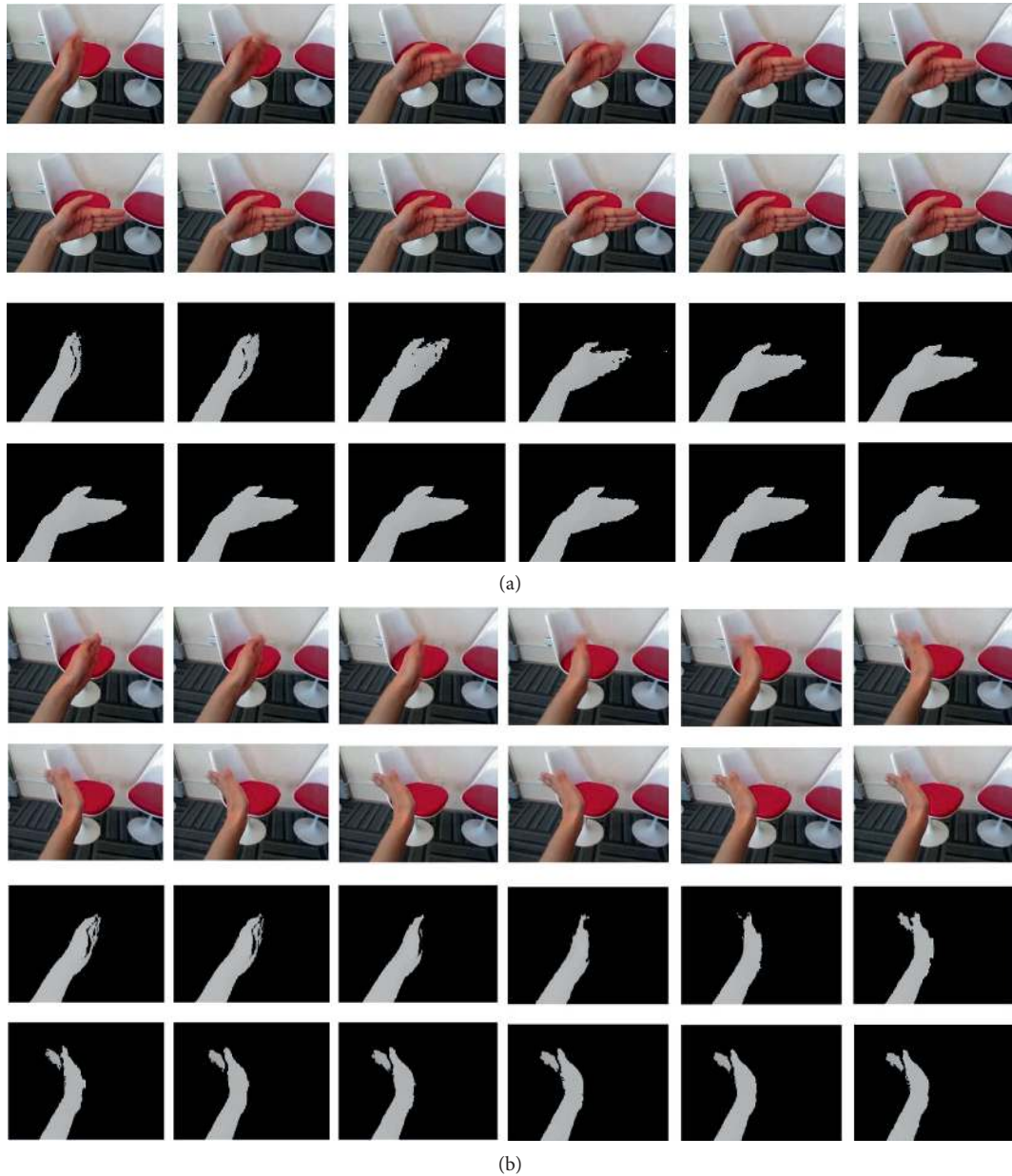


FIGURE 5: Samples of RGB-D images: (a) scroll hand towards right; (b) scroll hand towards left.

the number of frames. To increase the generalization of the model, the image is randomly clipped, and the initial input image of $240 * 240$ is randomly clipped to $112 * 112$. In the model training, the small batch stochastic gradient descent algorithm with momentum is used to optimize the 3D convolutional neural network. The number of training steps is 101, batch size is 16, initial learning rate is 0.01, and learning rate attenuation factor is 0.1 every 3000 iterations.

4.3. Comparison and Analysis

4.3.1. Comparative Experiment of Different Inputs and Fusion Strategies. The single-mode input and dual-mode input are used for comparative experiments. RGB images, depth images, optical flow images, and RGB-depth images are

selected to be input mode. Furthermore, average fusion and series fusion are used for dual-mode input in feature fusion layer. The input is 16-frame image set as input training, and the accuracy of the final dynamic gesture recognition result is shown in Figure 6.

Through the analysis of the experimental results, for the samples in the training set, in the use of the single-mode data input model, the recognition accuracy of RGB image is the highest, with a recognition accuracy of 52.5%. After fusing depth image input, the recognition accuracy is improved by 9.16%. It can be found that the multimodal data input model has better performance than the single-modal data input model. On the other hand, the multimode fusion input based on feature layer series connection has better effect, and the accuracy is 5.07% higher than the average fusion input. What

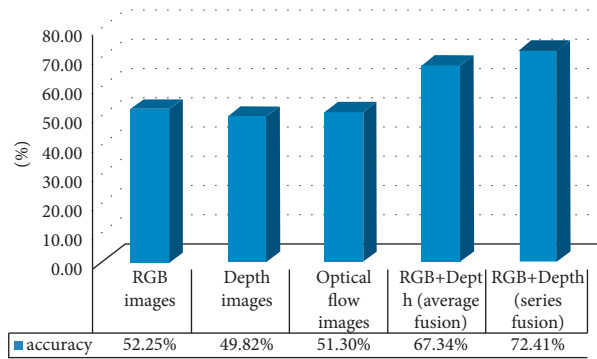


FIGURE 6: Accuracy comparison of input data in different modes.

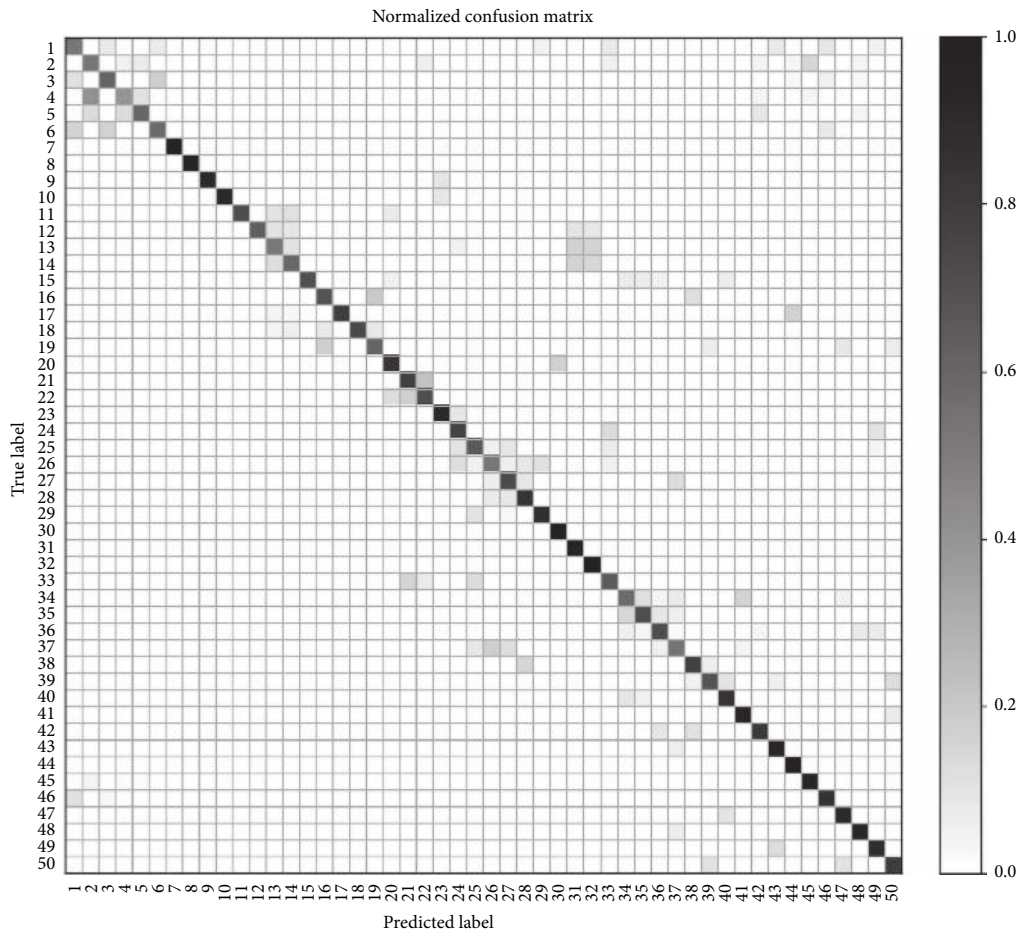


FIGURE 7: Confusion matrix of average fusion model.

can be inferred is that both input modes and fusion method influenced the performance of neural network. When the type of input data mode is fixed, using the appropriate data fusion method can make the characteristics of the object more prominent and train a better network performance.

4.3.2. Comparative Experiment of Dual-Mode Data Fusion Strategy. For better reflecting the effectiveness of the multimodal feature fusion strategy, the confusion matrix is used

to show effect of two different fusion methods. Figures 7 and 8 describe confusion matrices of two fusion methods of 50 kinds of gestures in EgoGesture dataset. The horizontal axis represents the categories predicted by the model for dynamic gestures, the vertical axis represents the real labels of dynamic gestures, and the right colour graph represents the prediction accuracy value and corresponding colour performance.

Comparing the two figures, it can be found that the average fusion model is easy to confuse some gesture

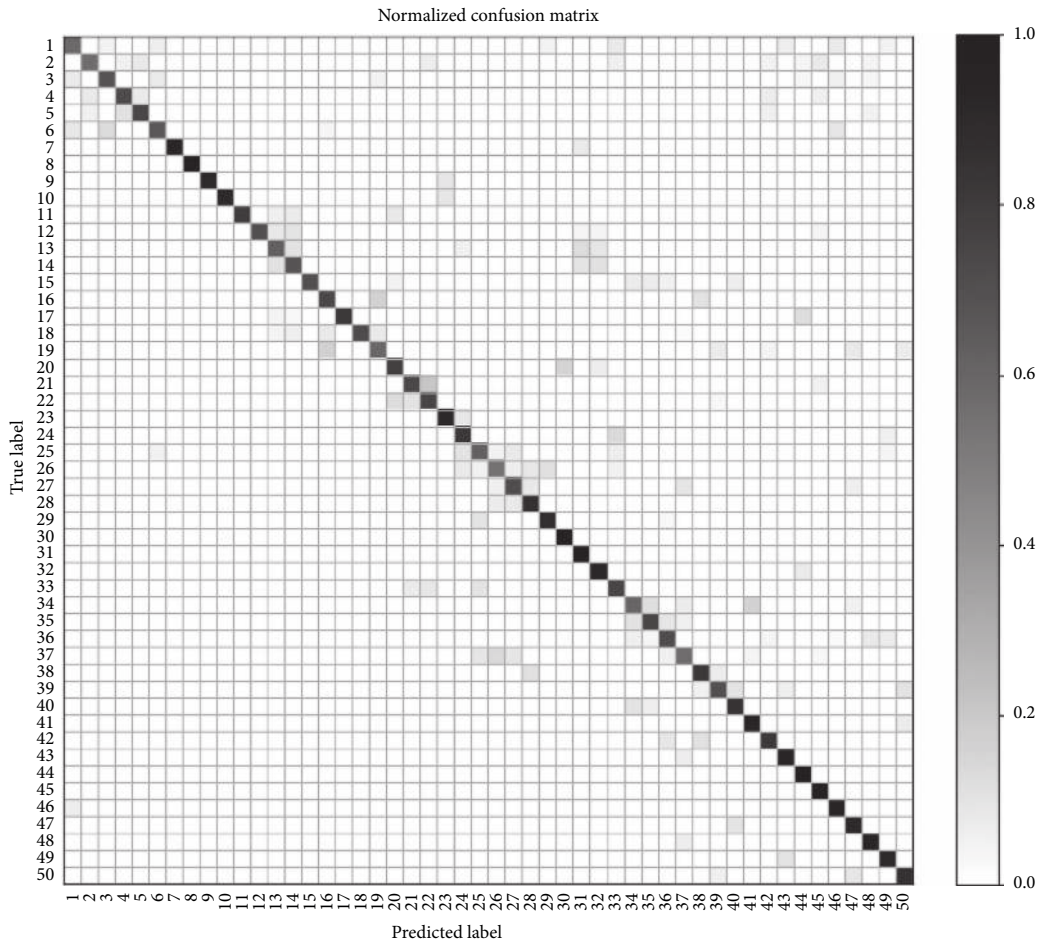


FIGURE 8: Confusion matrix of series fusion model.

recognition. For example, some confusions exist in gesture 1, 3, and 6 or gesture 2, 4, 5. These gestures do have similar features in motion trajectory and hand posture, which leads to misjudgment between them. After using the series fusion model, the false detection probability between these similar gestures decreases to a certain extent and the recognition accuracy is improved to a certain extent. Take gesture 2 as an example. In the average fusion confusion matrix, it can be found that the probability that gesture 4 is mistaken for gesture 2 is greater than 0.4 (the grid colour in the matrix can be compared with the colour bar on the right). However, in the series fusion mode, the probability that gesture 4 is mistaken for gesture 2 is less than 0.2. In general, the recognition effect and average recognition rate of the series fusion model are better than those of the average fusion model.

The above experiments show that in dynamic gesture recognition, the recognition model based on series fusion features can achieve better results mainly because it can save the features of each part when fusing. This method can avoid the loss of feature masking caused by direct fusion and can provide more complete feature information of the classifier so as to improve the performance of the whole model.

4.3.3. Comparative Experiment of Key Frame Extraction. Most of the action frames range from 20 to 50, so 8 frames, 12 frames, 16 frames, and 20 frames are selected for experimental comparison. The interframe difference method optimized in this paper is compared with the traditional interframe difference method, and the advantages and disadvantages of the method are judged by accuracy, which is shown in Figure 9.

As can be seen from Figure 8, the accuracy of the optimization method in extracting 8, 12, 16, and 20 key frames is significantly improved compared with the traditional method, and the maximum recognition rate is increased by 2.82% from 60.45% to 63.27%. With the increase in the number of frames, the recognition effect is also improving. However, when the number of frames reaches 20, the recognition rate decreases. It is speculated that the number of frames of some gestures is less than 20. The key extraction process needs to expand the number, resulting in redundant frames, which will have a negative impact on the effect of gesture recognition.

4.3.4. Overall Performance Evaluation of CBAM-C3D Model. The proposed dynamic gesture recognition method and other representative algorithms are tested on the EgoGesture

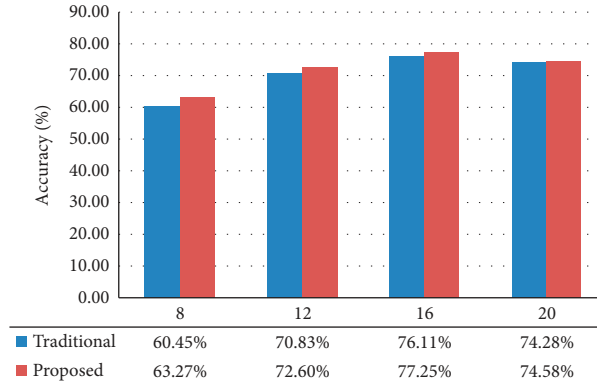


FIGURE 9: Comparison of accuracy of two key frame extraction methods.

TABLE 1: Results of different network models in EgoGesture.

Recognition methods	Input categories	Validation set of EgoGseture (%)	Test set of EgoGseture (%)
Harris3.5D	RGB + depth	35.8	36.1
HOG3D	RGB + depth	43.8	44.6
HON4D	Depth	57.2	58.3
C3D	RGB + depth	56.1	57.4
T3D	RGB + depth	62.4	63.8
R3D	RGB + depth	64.8	66.1
Proposed algorithm	RGB + depth	71.5	72.4

dataset. The feasibility and effectiveness of this method are verified by comparative experiments. In Table 1, the accuracy of the proposed method is 72.4%, which is greatly improved compared with the traditional dynamic gesture recognition method.

5. Conclusions

A dynamic gesture recognition algorithm based on attention mechanism of 3D convolutional neural network is proposed in this manuscript. The optimized interframe difference method is used to deal with problems about video data redundancy and format disunity. Combined with CBAM network, the important spatiotemporal features are enhanced and invalid features are suppressed to realize the prominent expression of features. Finally, the dual-mode feature input fusion method is adopted to realize the complementary features of the two modes and improve the effect of gesture recognition. Meanwhile, the model designed in this paper is compared with other mainstream methods on EgoGesture dataset to verify the effectiveness of the proposed method. The recognition accuracy of the designed method is 72.4%, which is better than other networks.

This method also has some defects, such as large amount of network parameters and slow network prediction, resulting in poor real-time performance. It only recognized gesture for video containing a single action. Moreover, the dual-modal data fusion strategy proposed in this paper has mentioned that the input modal RGB image and depth image are related to each other, but this paper only fuses

them on the feature layer. In the future work, it can be considered to fuse the dual-modal data in the preprocessing stage and increase the input modal information, such as the addition of optical flow data, to improve the recognition accuracy.

Data Availability

The data used to support the findings of this study have not been made available because the relevant data involve legal issues and related confidentiality.

Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this article.

Acknowledgments

This work was supported by grants of the National Natural Science Foundation of China (52075530, 51575407, 51505349, 61733011, and 41906177), the grants of Hubei Provincial Department of Education (D20191105), the grants of National Defense Preresearch Foundation of Wuhan University of Science and Technology (GF201705), the Open Fund of the Key Laboratory for Metallurgical Equipment and Control of Ministry of Education in Wuhan University of Science and Technology (2018B07 and 2019B13), and the Open Fund of Hubei Key Laboratory of Hydroelectric Machinery Design & Maintenance in Three Gorges University (2020KJX02).

References

- [1] B. K. Chakraborty, D. Sarma, M. K. Bhuyan, and K. F. MacDorman, "Review of constraints on vision-based gesture recognition for human-computer interaction," *IET Computer Vision*, vol. 12, no. 1, pp. 3–15, 2018.
- [2] D. Jiang, G. Li, Y. Sun, J. Kong, B. Tao, and D. Chen, "Grip strength forecast and rehabilitative guidance based on adaptive neural fuzzy inference system using sEMG," *Personal and Ubiquitous Computing*, 2019.
- [3] Y. Liu, M. Peng, M. R. Swash, T. Chen, R. Qin, and H. Meng, "Holoscopic 3D microgesture recognition by deep neural network model based on viewpoint images and decision fusion," *IEEE Transactions on Human-Machine Systems*, vol. 51, no. 2, pp. 162–171, 2021.
- [4] J. Li, B. Yang, D. Chen, N. Wang, G. Zhang, and H. Bao, "Survey and evaluation of monocular visual-inertial SLAM algorithms for augmented reality," *Virtual Reality & Intelligent Hardware*, vol. 1, no. 4, pp. 386–410, 2019.
- [5] B. Luo, Y. Sun, G. Li, D. Chen, and Z. Ju, "Decomposition algorithm for depth image of human health posture based on brain health," *Neural Computing and Applications*, vol. 32, no. 10, pp. 6327–6342, 2020.
- [6] D. Jiang, G. Li, Y. Sun, J. Hu, J. Yun, and Y. Liu, "Manipulator grabbing position detection with information fusion of color image and depth image using deep learning," *Journal of Ambient Intelligence and Humanized Computing*, 2021.
- [7] S. Liao, G. Li, H. Wu et al., "Occlusion gesture recognition based on improved SSD," *Concurrency and Computation: Practice and Experience*, vol. 33, no. 6, p. e6063, 2021.
- [8] T. Chen, L. Peng, J. Yang, and G. Cong, "Analysis of user needs on downloading behavior of English vocabulary APPs based on data mining for online comments," *Mathematics*, vol. 9, no. 12, p. 1341, 2021.
- [9] F. Demim, A. Nemra, A. Boucheloukh, E. Kobzili, M. Hamerlain, and A. Bazoula, "SLAM based on adaptive SVSF for cooperative unmanned vehicles in dynamic environment," *IFAC-PapersOnLine*, vol. 52, no. 8, pp. 73–80, 2019.
- [10] W. Shu, K. Cai, and N. N. Xiong, "Research on strong agile response task scheduling optimization enhancement with optimal resource usage in green cloud computing," *Future Generation Computer Systems*, vol. 124, pp. 12–20, 2021.
- [11] H. P. Gupta, H. S. Chudgar, S. Mukherjee, T. Dutta, and K. Sharma, "A continuous hand gestures recognition technique for human-machine interaction using accelerometer and gyroscope sensors," *IEEE Sensors Journal*, vol. 16, no. 16, pp. 6425–6432, 2016.
- [12] K. A. Bhaskaran, A. G. Nair, K. D. Ram, K. Ananthanarayanan, and H. N. Vardhan, "Smart gloves for hand gesture recognition: sign language to speech conversion system," in *Proceedings of the 2016 International Conference on Robotics and Automation for Humanitarian Applications (RAHA)*, pp. 1–6, IEEE, Amritapuri, India, December 2016.
- [13] H. Duan, Y. Sun, W. Cheng et al., "Gesture recognition based on multi-modal feature weight," *Concurrency and Computation: Practice and Experience*, vol. 32, no. 5, p. e5991, 2021.
- [14] T. Chen, J. Rong, L. Peng, J. Yang, G. Cong, and J. Fang, "Analysis of social effects on employment promotion policies for college graduates based on data mining for online user review in China during the COVID-19 pandemic," *Healthcare*, vol. 9, no. 7, p. 846, 2021.
- [15] Y. Peng, H. Tao, W. Li, H. Yuan, and T. Li, "Dynamic gesture recognition based on feature fusion network and variant ConvLSTM," *IET Image Processing*, vol. 14, no. 11, pp. 2480–2486, 2020.
- [16] C. Tan, Y. Sun, G. Li, G. Jiang, D. Chen, and H. Liu, "Research on gesture recognition of smart data fusion features in the IoT," *Neural Computing and Applications*, vol. 32, no. 22, pp. 16917–16929, 2020.
- [17] T. Chen, X. Yin, L. Peng, J. Rong, J. Yang, and G. Cong, "Monitoring and recognizing enterprise public opinion from high-risk users based on user portrait and random forest algorithm," *Axioms*, vol. 10, no. 2, p. 106, 2021.
- [18] R. H. Huan, C. J. Xie, F. Guo et al., "Human action recognition based on HOIRM feature fusion and AP clustering BOW," *PLoS One*, vol. 14, no. 7, Article ID e0219910, 2019.
- [19] D. Jiang, G. Li, Y. Sun, J. Kong, and B. Tao, "Gesture recognition based on skeletonization algorithm and CNN with ASL database," *Multimedia Tools and Applications*, vol. 78, no. 21, pp. 29953–29970, 2019.
- [20] Y. Zhang, L. Shi, Y. Wu, K. Cheng, J. Cheng, and H. Lu, "Gesture recognition based on deep deformable 3D convolutional neural networks," *Pattern Recognition*, vol. 107, Article ID 107416, 2020.
- [21] E. Tsironi, P. Barros, C. Weber, and S. Wermter, "An analysis of convolutional long short-term memory recurrent neural networks for gesture recognition," *Neurocomputing*, vol. 268, pp. 76–86, 2017.
- [22] P. Bao, A. I. Maqueda, C. R. del-Blanco, and N. García, "Tiny hand gesture recognition without localization via a deep convolutional network," *IEEE Transactions on Consumer Electronics*, vol. 63, no. 3, pp. 251–257, 2017.
- [23] D. Jiang, Z. Zheng, G. Li et al., "Gesture recognition based on binocular vision," *Cluster Computing*, vol. 22, no. 3, pp. 13261–13271, 2019.
- [24] F. Xiao, G. Li, D. Jiang et al., "An effective and unified method to derive the inverse kinematics formulas of general six-DOF manipulator with simple geometry," *Mechanism and Machine Theory*, vol. 159, Article ID 104265, 2021.
- [25] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [26] D. Jiang, G. Li, C. Tan, L. Huang, Y. Sun, and J. Kong, "Semantic segmentation for multiscale target based on object recognition using the improved faster-RCNN model," *Future Generation Computer Systems*, vol. 123, pp. 94–104, 2021.
- [27] Y. Cheng, G. Li, M. Yu et al., "Gesture recognition based on surface electromyography-feature image," *Concurrency and Computation: Practice and Experience*, vol. 33, no. 6, p. e6051, 2021.
- [28] W. Wei, Q. Dai, Y. Wong, Y. Hu, M. Kankanhalli, and W. Geng, "Surface-electromyography-based gesture recognition by multi-view deep learning," *IEEE Transactions on Biomedical Engineering*, vol. 66, no. 10, pp. 2964–2973, 2019.
- [29] L. Dong, M. N. Satpute, W. Wu, and D. Du, "Two-phase multidocument summarization through content attention-based subtopic detection," *IEEE Transactions on Computational Social Systems*, pp. 1–14, 2021.
- [30] K. Simonyan and Z. Andrew, "Two-stream convolutional networks for action recognition in videos," 2014, <https://arxiv.org/abs/1406.2199>.
- [31] S. Hare, S. Golodetz, A. Saffari et al., "Struck: structured output tracking with kernels," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 10, pp. 2096–2109, 2015.

- [32] E. Rosten and T. Drummond, "Machine learning for high-speed corner detection," in *European Conference on Computer Vision*, pp. 430–443, Springer, Berlin, Germany, 2006.
- [33] L. Wang, Y. Xiong, Z. Wang et al., "Temporal segment networks: towards good practices for deep action recognition," in *European Conference on Computer Vision*, pp. 20–36, Springer, Cham, Germany, 2016.
- [34] C. Feichtenhofer, A. Pinz, and A. Zisserman, "Convolutional two-stream network fusion for video action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1933–1941, Las Vegas, NV, USA, June 2016.
- [35] Y. Zhu, Z. Lan, S. Newsam, and A. Hauptmann, "Hidden two-stream convolutional networks for action recognition," in *Asian Conference on Computer Vision*, pp. 363–378, Springer, Cham, Germany, 2018.
- [36] G. Li, H. Wu, G. Jiang, S. Xu, and H. Liu, "Dynamic gesture recognition in the internet of things," *IEEE Access*, vol. 7, pp. 23713–23724, 2019.
- [37] Y. Weng, Y. Sun, D. Jiang et al., "Enhancement of real-time grasp detection by cascaded deep convolutional neural networks," *Concurrency and Computation: Practice and Experience*, vol. 33, no. 5, p. e5976, 2021.
- [38] Y. He, G. Li, Y. Liao et al., "Gesture recognition based on an improved local sparse representation classification algorithm," *Cluster Computing*, vol. 22, pp. 10935–10946, 2019.
- [39] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3D convolutional networks," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4489–4497, Santiago, Chile, December 2015.
- [40] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6299–6308, Seattle, WA, USA, July 2017.
- [41] Z. Qiu, T. Yao, and T. Mei, "Learning spatio-temporal representation with pseudo-3d residual networks," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5533–5541, Venice, Italy, October 2017.
- [42] G. Hinton, Y. LeCun, and Y. Bengio, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [43] P. Molchanov, X. Yang, S. Gupta, K. Kim, S. Tyree, and J. Kautz, "Online detection and classification of dynamic hand gestures with recurrent 3D convolutional neural network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4207–4215, Las Vegas, NV, USA, June 2016.
- [44] P. Molchanov, S. Gupta, K. Kim, and J. Kautz, "Hand gesture recognition with 3D convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1–7, Boston, MA, USA, June 2015.
- [45] P. Molchanov, X. Yang, S. Gupta, K. Kim, S. Tyree, and J. Kautz, "Online detection and classification of dynamic hand gestures with recurrent 3D convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, June 2016.
- [46] Y. Zhang, C. Cao, J. Cheng, and H. Lu, "EgoGesture: a new dataset and benchmark for egocentric hand gesture recognition," *IEEE Transactions on Multimedia*, vol. 20, no. 5, pp. 1038–1050, 2018.