# CenterNet-SPP based on multi-feature fusion for basketball posture recognition

Zhouxiang Jin[1,*]

[1]Jiaozuo University, 3066 Renmin Road, Jiaozuo City, Henan Province, 454000 China

## Abstract

Aiming at the problem that the existing posture recognition algorithms can not fully reflect the dynamic characteristics of athletes' posture, this paper proposes a CenterNet-SPP model based on multi-feature fusion algorithm for basketball posture recognition. Firstly, motion posture images are collected by optical image collector, and then gray scale transformation is performed to improve the image quality. Furthermore, body contour and motion posture region are obtained based on shadow elimination technology and inter-frame difference method. Finally, radon transform and discrete wavelet transform are used to extract the motion posture region and body contour, and the two complementary features are fused and then input into the CenterNet-SPP network to realize the final posture recognition. Experimental results show that the recognition accuracy of the proposed method is higher than that of other new methods.

*Corresponding author. Email: snowberry@qq.com

## 1. Introduction

Accurate recognition of athletes' posture movements is very necessary in high-level training and key judgment in major competitions [1-4]. At present, relevant references have been explored in this field. Reference [5] proposed a double-layer background modeling method based on code-book and run-time mean method, and used it to detect moving targets, achieving good results. Reference [6] proposed a method of posture recognition based on motion region, and discussed a two-level posture modeling architecture in detail. Reference [7] proposed a method of diving posture recognition based on visual technology, and achieved accurate recognition results.

In recent years, optical image has gradually become an important branch of digital image processing, and the application of optical image in motion recognition is getting more and more attention [8]. Based on this, a basketball posture recognition algorithm based on CenterNet-SPP model and multi-feature fusion is proposed.

In this method, motion postures such as serve and serve are collected by optical image collector, and the image quality is improved based on gray transform. The body contour and body contour are obtained by shadow elimination technology and inter-frame difference method, and then the body contour and body contour are extracted based on Radon transform and discrete wavelet transform, and the final posture recognition is achieved by combining the two complementary features and network

training. In order to verify the effectiveness of the proposed method, a comparative experiment is designed. The experimental results show that the proposed method achieves higher recognition accuracy than the traditional method.

## 2. Proposed posture recognition method

The existing posture recognition algorithms can not fully reflect the dynamic characteristics of athletes' movement posture. A posture recognition algorithm based on multi-feature fusion and CenterNet-SPP model is proposed. The overall structure of the proposed method is shown in figure 1.
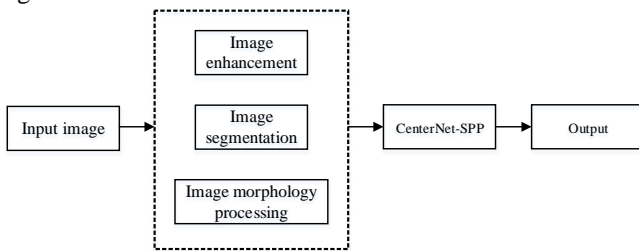


**Figure 1.** Proposed model

Firstly, some basic postures including serving and receiving serve are collected through optical image collector. Then image preprocessing operation, further feature extraction and the final accurate recognition are executed. The process of accurate recognition of athletes' posture movements can be divided into several steps: 1) feature extraction based on body contour and postural region; 2) The recognition and modeling process of motion postures; 3) Multi-feature fusion.

Optical image acquisition and processing technology is a key technology which is gradually mature to assist research.

In order to collect optical images of motion posture, this paper uses CCD as image sensor and CPLD as control core to collect and process optical images of motion posture. Some basic postures including serving, receiving, dribbling, passing and shooting are collected.

### 2.1. Motor posture recognition area

Based on formula (1), the recognition area of athlete's movement posture is defined:

$$Q_g(h) = \int_X h(X)g = \int_0^1 g(h_\varepsilon)d_\varepsilon \qquad (1)$$

Where $X = \{x_1, x_2, \cdots, x_n\}$ represents a finite discrete set of athlete's posture recognition area, so equation (1) can be rearranged as:

$$Q_g(h) = \sum_{i=1}^n g(x_i)[h(x_i) - h(x_i + 1)]$$

$$= \sum_{i=1}^n h(x_i)[g(x_i) - g(X_{i-1})] \qquad (2)$$

Suppose $T = \{t_1, t_2, t_3, t_4\}$ represents the feature of the recognized region, F represents the motion posture to be recognized, and $X = \{x_1, x_2, x_3\}$ represents the posture feature group. The pose recognition model is constructed by the following formula:

$$g_k^h = 0.8 \le k \le 4 \qquad (3)$$

### 2.2. Image preprocessing

Image preprocessing is a very important step in motion recognition. The athlete posture recognition algorithm based on optical image collector divides the image preprocessing process into the following three steps:

(1) Image enhancement. In order to improve the final identification accuracy, the gray scale transformation of the collected image is carried out to improve the image quality. Set the image gray level as L, and the r-level gray level of the original image can be mapped to the s-level gray level of the resulting image through a mapping transformation, namely:

$$s = T(r), r, s \in (0, L-1) \qquad (4)$$

(2) Image segmentation. It is the key step in image processing. In the process of athlete posture recognition, shadow elimination technology and inter-frame difference method are selected to obtain body contour and motion posture region.

$$D(x, y) = \begin{cases} 1, |f_{k+1}(x, y) - f_k(x, y)| > T \\ 0, else \end{cases} \qquad (5)$$

(3) Image morphology processing. Based on morphology, the connectivity analysis of the movement posture region is carried out.

### 2.3. Motion feature extraction

Based on radon transform [9] and discrete wavelet transform [10], motion posture region and body contour features are extracted. Firstly, motion posture region and domain are extracted, and discrete binary images $\xi(x, y)$ are set to carry out standardized processing on the image, and the following formula is used for feature extraction of motion region:

$$\frac{1}{\alpha^2}\int_{-\infty}^{\infty} T_R^2 f(-\rho, \theta \pm \pi) d\rho = -\int_{-\infty}^{\infty} T_R^2 f(v, \theta \pm \pi) \quad (6)$$

Then the body contour features are extracted, and the body contour lines are expanded into one-dimensional features of the Euclidean distance of body contour points by using the following formula:

$$d_i = \sqrt{(x_i - x_c)^2 + (y_i - y_c)^2} \quad (7)$$

## 2.4. CenterNet-SPP model

AlexNet model proposed ReLU function, and the convergence speed of the model using ReLU was much faster than Sigmoid [11], which became one of the advantages of AlexNet model. The model has a total of 8 layers, among which the first 5 layers are convolution layers, the first two convolution layers and the fifth convolution layer have pooling layers, and the last 3 layers are fully connected layers. At the same time, each layer plays different roles. Overlapping pooling layer is to improve accuracy and is not prone to over-fitting, local normalized response is to improve accuracy, and data gain and dropout are to reduce over-fitting.

The essence of VGG16 model is an enhanced version of AlexNet structure, which emphasizes the depth of convolutional neural network design [12]. Each of these convolution layers is followed by a pooling layer. VGG network uses a smaller convolution kernel, which makes the parameters smaller and saves computing resources. Due to the large number of layers, the convolution kernel is relatively small, which makes the whole network have better feature extraction effect.

Inception V3 network is a deep convolutional network developed by Google [13]. The main idea of Inception structure in the model is to find out how to approximate the optimal local sparse structure with dense components. The Inception structure adopts convolution kernels of 3×3 and 5×5 size, and adds convolution kernels of 1×1 size, and proposes BN (Batch Normalization) method [14]. The method of splitting a large two-dimensional convolution into two small one-dimensional convolution is introduced by using the branch structure. This asymmetric convolution structure splitting has better effects on processing more and richer spatial features and increasing feature diversity than symmetric convolution structure splitting, and can reduce the amount of calculation.

ResNet50 model solves the problem that the actual effect becomes worse due to the increase of network depth and width [15]. The deep neural network model sacrifices a lot of computing resources, and the error rate also increases. This phenomenon is mainly caused by the gradient disappearing phenomenon becoming more and more obvious with the increase of neural network layers. In ResNet50 model, residual structure is added, that is, an identity mapping is added to transform the original transformation function H(x) into F(x)+x, which makes the network no longer a simple stack structure and solves the problem of gradient disappearance. Such simple superposition does not add extra parameters and computation to the network, but also improves the effect and efficiency of network training.

MobileNet model reduces model complexity and improves model speed while maintaining model performance [16]. The basic unit of the model is depth-level separable convolution, and its essence is separable convolution operation. Different from the standard convolution that the convolution kernel acts on all the input channels, the depth separation convolution adopts different convolution kernels each input channel, and at the same time, BN is added and ReLU activation function is used to greatly reduce the amount of computation and the number of model parameters.

CenterNet is a detection network based on the center point, which is simple and fast, and its accuracy is no less than that of the detector based on the anchor frame. In COCO dataset, the AP value of CenterNet model is 20% higher than that of YOLOv2 model, and 6.9% higher than that of Faster-RCNN model [17]. Therefore, CenterNet model is selected for motion posture recognition in this paper. Different from the traditional target detection model, the CenterNet model takes the detection target as a central point instead of the anchor frame, which solves the problems of unbalanced positive and negative samples in the anchor frame and too much calculation. The algorithm firstly obtains the feature image through the feature extraction network, then finds the local feature peak on the feature image as the center point, and obtains the image features such as the target size through the center point regression. In the training process, each target generates a central point without NMS non-maximum suppression, which reduces the amount of computation and training time. Meanwhile, feature maps with higher resolution are used to improve the detection ability of small targets.

The Centernet-SPP network structure adopted in this paper is a feature extraction network based on ResNet50, and the Spatial Pyramid Pooling (SPP) module is added [18-20]. SPP is used to increase the range of receptive fields to improve the reception range of trunk features, and at the same time significantly separate important features to improve the ability of feature extraction.

Based on this network model, the original image $I \in R^{W \times H \times 3}$ is first input into ResNet50 feature extraction network, and then into SPP module. The number of channels is adjusted to 512 through 1×1 convolution, and 12×12 convolution. The maximum pooling of 8X8 and 4×4, the four feature maps were added to obtain the final feature map and the heat map of

key points $\hat{Y} \in [0,1]^{(W/R) \times C}$, that is, the Gaussian distribution map. Where W and H are the width and height of the image respectively, and R is the size scaling ratio. In this test, R=4, and C is the number of detection points. $\hat{Y}_{xyc} = 1$ indicates that a key point has been detected, and $\hat{Y}_{xyc} = 0$ indicates the background.

In central point prediction, the sampled coordinate of key point $p = (\frac{x_1 + x_2}{2}, \frac{y_1 + y_2}{2})$ is $\tilde{p} = [p/R]$. A Gaussian kernel $Y_{xyc} = \exp(-\frac{(y - \tilde{p}_x)^2 + (y - \tilde{p}_y)^2}{2s_p^2})$ is used to distribute the key points of the real frame to the heat map. Focal Loss is used for pixel-level logistic regression, and the loss function is $L_k$. Where $s_p$ is the adaptive standard variance of the target size. $\alpha$ =2 and $\beta$ =4 are focal Loss hyperparameters. $\hat{Y}_{xyc}$ is the predicted value, and $Y_{xyc}$ is the labeled true value.

$$L_k = -\frac{1}{N}\sum_{xyc}\begin{cases}(1-\hat{Y}_{xyc})^{\alpha}\log(\hat{Y}_{xyc}) & if \quad Y_{xyc} = 1 \\ (1-Y_{xyc})^{\beta}(\hat{Y}_{xyc})^{\alpha}\log(1-\hat{Y}_{xyc}) & others\end{cases}$$
(8)

Where N is the number of key points in the image.
In the prediction of center point offset, due to the down-sampling R=4 times of the image, such feature map will bring precision error when remapping to the original image. Therefore, a local offset $\hat{O} \in R^{\frac{W}{R} \times \frac{H}{R} \times 2}$ is used to compensate for each center point. The L1 loss function is used to train the center offset value of the object. The loss function is $L_{off}$, where $\hat{O}_{\tilde{p}}$ is the predicted bias. At this point, $p$ is the coordinate of the center point of the image, and $\tilde{p}$ is the approximate integer coordinate of the center point after scaling.

$$L_{off} = \frac{1}{N}\sum_{p}|\hat{O}_{\tilde{p}} - (p/R - \tilde{p})|$$
(9)

When the size of the target box is predicted, $(x_1^{(k)}, y_1^{(k)}, x_2^{(k)}, y_2^{(k)})$ is taken as the boundary box of the target $k$, and its center point is $p_k = (\frac{x_1^{(k)} + x_2^{(k)}}{2}, \frac{y_1^{(k)} + y_2^{(k)}}{2})$. The key point estimator $\hat{Y}$ is used to generate the center point, and the regression is performed for each target, and the final regression size is $S_k = (x_2^{(k)} - x_1^{(k)}, y_2^{(k)} - y_1^{(k)})$. This value is calculated before the training and is the length and width value after sampling. To reduce the difficulty of regression, $\hat{S}$ is used as the predictive value, $L_{size}$ loss function is used for training.

$$L_{size} = \frac{1}{N}\sum_{k=1}^{N}|\hat{S}_{pk} - S_k|$$
(10)

Where $S_k$ is the real size of the target. $\hat{S}_{pk}$ is the predicted size.

The overall loss function is the sum of the central point prediction loss function, the central point offset loss function and the target frame size loss function. Each loss has a weight.

$$L_{det} = L_k + \lambda_{size}L_{size} + \lambda_{off}L_{off}$$
(11)

Where $\lambda_{size} = 0.1$, $\lambda_{off} = 1$.

## 2.5. Feature fusion

The final gesture recognition is achieved by fusing two complementary features. First, it builds the training sequence of movement posture:

$$X_g = \{X_{g1}, X_{g2}, \cdots, X_{gm}\}$$
(12)

Then it builds the motion posture test sequence:

$$X_p = \{X_{p1}, X_{p2}, \cdots, X_{p+n}\}$$
(13)

Where $m$ and $n$ represent the frames of the motion posture sequence of two athletes respectively. $X_{ij}$ represents the j-th feature vector in the i-th athlete's movement sequence. The distance between the movement period of the posture training set and the k-th sub-sequence of the test set is calculated based on the following formula:

$$dis_{(x_p(k),x_g)}(l) = \sum_{J=1}^{N}\| x_{p,k+j} - x_{g,l+j} \|$$
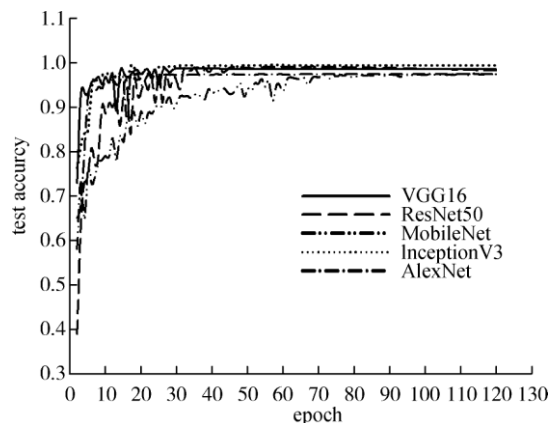(14)

To calculate the similarity of all motion sequences, the following formula is adopted:

$$sim(x_p, x_g) = 1 - \frac{1}{k}\sum_{k=1}^{k}\min_{l}(dis_{x_p(k)}, x_g(l))$$
(15)

## 3. Experiments and analysis

The experimental conditions of this paper are: Windows10, 64-bit operating system, Cuda version 10.0, Tensorflow and Keras deep learning framework based on Python programming language are adopted. The PC is configured with a GeForceGTX 1060 video card, 6G video memory, Intel (R) Core (TM) i59400F processor, 2.90ghz.
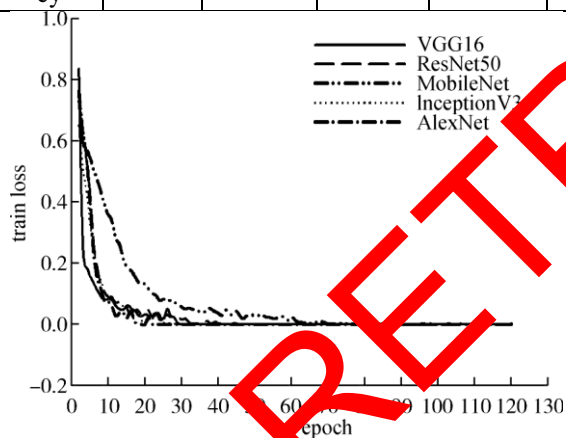
In this paper, the same data set is used to train basketball locomotion posture in different situations under five different convolutional neural network models (VGG16, ResNet50, Inception V3, Mobilenet, AlexNet). For each model, the method of transfer learning [21,22] is used to initialize the parameters of the pre-trained ImageNet classification model, and the model is iterated 120 times. Cross entropy loss function is used. Adam optimizer is used at the same time, and the initial learning rate is 0.000 1 and the momentum factor was 0.1. After 5 epochs, if the model performance is not improved, the learning rate will be reduced. The final loss values of the five models tended to be stable, while the accuracy of the test set stabilized at a relatively high value. Table 1 shows the accuracy of the test sets of the five models. Figure 2 shows the accuracy of training set loss, test set loss and test set.
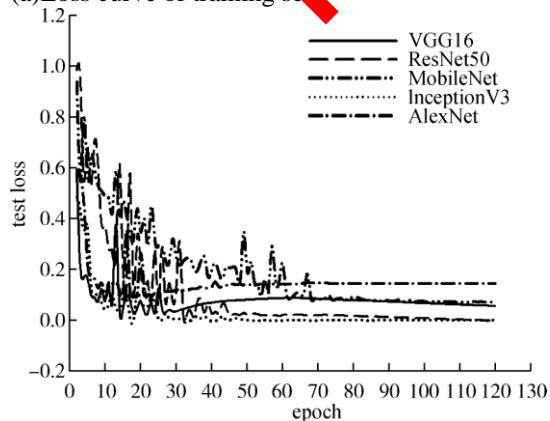
Table 1. Accuracy comparison of the model/%

| Model | VGG 16 | ResNet 50 | Inception V3 | Mobile net | AlexNet |
|-------|--------|-----------|--------------|------------|---------|
| accuracy | 98.81 | 99.31 | 99.61 | 97.91 | 97.61 |



(a)Loss curve of training set



(b)Testing set Loss curve



(c)Accuracy curve of test set

**Figure 2.** Curve of five models

In basketball motion images, 50 images are randomly selected to verify each model respectively, and the confusion matrix obtained is shown in figure 3. InceptionV3 model has the highest average accuracy of 98.11% (Table 2), and it takes 0.12s on average to identify an image. The results show that all the five models can recognize basketball motion image accurately.

Table 2. Average accuracy rate of model

| Model | Accuracy/% |
|-------|-----------|
| VGG16 | 96.11 |
| ResNet50 | 96.67 |
| InceptionV3 | 98.11 |
| MobileNet | 96.11 |
| AlexNet | 90.11 |

| | Serve | Pass | Dribble | Shoot |
|---|-------|------|---------|-------|
| Serve | 62.1 | 0.5 | 2.3 | 32.4 |
| Pass | 25.4 | 73.5 | 4.8 | 12.1 |
| Dribble | 5.8 | 8.9 | 69.2 | 11.5 |
| Shoot | 21.3 | 14.7 | 6.8 | 70.8 |

(a)confusion matrix of VGG16

| | Serve | Pass | Dribble | Shoot |
|---|-------|------|---------|-------|
| Serve | 73.5 | 14.8 | 5.6 | 22.4 |

| | | | |
|---|---|---|---|
| Pass | 18.9 | 66.8 | 32.1 | 10.6 |
| Dribble | 6.7 | 12.5 | 70.8 | 9.6 |
| Shoot | 10.5 | 11.7 | 24.3 | 79.6 |

(b)confusion matrix of ResNet50

| | Serve | Pass | Dribble | Shoot |
|---|---|---|---|---|
| Serve | 86.5 | 14.7 | 0.8 | 26.3 |
| Pass | 15.7 | 77.4 | 5.2 | 24.7 |
| Dribble | 16.8 | 11.2 | 82.3 | 3.6 |
| Shoot | 25.8 | 14.6 | 12.3 | 81.1 |

(c)confusion matrix of InceptionV3

| | Serve | Pass | Dribble | Shoot |
|---|---|---|---|---|
| Serve | 80.4 | 1.5 | 7.8 | 14.7 |
| Pass | 15.7 | 72.6 | 17.7 | 10.8 |
| Dribble | 3.7 | 12.8 | 80.1 | 21.4 |
| Shoot | 12.8 | 10.3 | 2.9 | 74.4 |

(d)confusion matrix of MobileNet

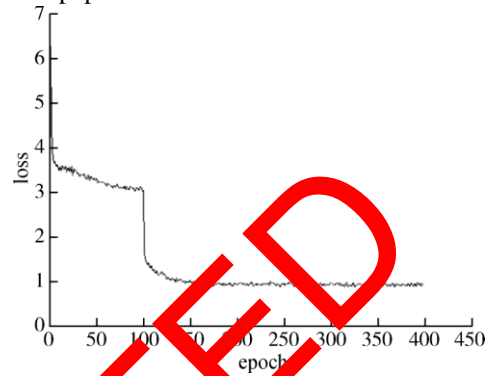| | Serve | Pass | Dribble | Shoot |
|---|---|---|---|---|
| Serve | 61.8 | 12.5 | 6.4 | 11.9 |
| Pass | 12.5 | 77.7 | 12.3 | 9.9 |
| Dribble | 25.1 | 14.5 | 71.3 | 9.2 |
| Shoot | 6.4 | 18.9 | 3.9 | 77.9 |

(e)confusion matrix of AlexNet

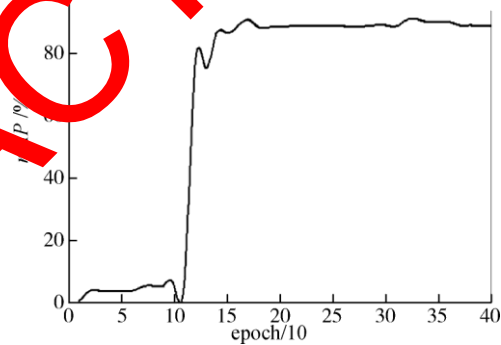**Figure 3.** Confusion matrix of the five models

In the experiment, LabelImg (image annotation tool) is used to label basketball posture manually according to PASCAL VOC2007 format. Centernet-SPP structure is used and the pre-training model of VOC dataset is used to set initialization parameters. A total of 400 epochs are trained in the model. The loss precision of the model decreased rapidly in the first 100 times. Since the model is thawed after 100 times (the first 100 epochs trained the model after the backbone feature network and trained all the networks in the last 300 times), the loss value of the model decreases rapidly and then gradually tends to be stable. This indicates that the training effect of the model is good, and its training loss curve is shown in figure 4(a).

In order to select a model with high enough overall performance, the posture targets with confidence greater than 0.5 are retained first, and the weight files with the highest mAP values are found. Then the model is carried out 400 iterations, and one model is output every 10 iterations. Therefore, a total of 40 models are obtained, and a model with the highest mAP value should be found among the 40 models as shown in figure4(b). When the mAP value tends to be stable at the end of iteration, the maximum value is 90.03%, which is the selected model in this paper.

(a)Curve of Loss value changing with the number of iterations

(b)The mean of average accuracy changes with the number of iterations

**Figure 4.** Training results of CenterNet-SPP model

The above experimental data show that the method proposed in this paper can effectively extract various postures of basketball.

## 4. Conclusion

An algorithm for athlete posture recognition based on multi-feature fusion is proposed. In this method, motion posture images are collected by optical image collector, and the image quality is improved by gray scale transformation. Furthermore, body contour and motion posture region are obtained based on shadow elimination technology and inter-frame difference method. Finally, radon transform and discrete wavelet transform are used

to extract the motion posture region and body contour, and the final posture recognition is achieved through the fusion of the two complementary features and CenterNet-SPP network training. Experiments show that the recognition accuracy of the proposed method is higher than other new methods. In the future, it is hoped that the proposed method can be extended to other types of athletes' movement recognition, so as to better improve the overall training level of athletes.

## References

[1] J. E. Kiriazi, S. M. M. Islam, O. Borić-Lubecke and V. M. Lubecke, "Sleep Posture Recognition With a Dual-Frequency Cardiopulmonary Doppler Radar," in IEEE Access, vol. 9, pp. 36181-36194, 2021, doi: 10.1109/ACCESS.2021.3062385.

[2] Shoulin Yin, Jie Liu, Ye Zhang, Lin Teng. Cuckoo search algorithm based on mobile cloud model[J]. International Journal of Innovative Computing, Information and Control. Volume 12, Number 6. pp.1809-1819. 2016.

[3] W. Ren, O. Ma, H. Ji and X. Liu, "Human Posture Recognition Using a Hybrid of Fuzzy Logic and Machine Learning Approaches," in IEEE Access, vol. 8, pp. 135628-135639, 2020, doi: 10.1109/ACCESS.2020.3011697.

[4] Liu T, Yin S. An improved particle swarm optimization algorithm used for BP neural network and multimedia course-ware evaluation[J]. Multimedia Tools & Applications, 76(9):11961-11974, 2017.

[5] Zhang S, Callaghan V. Real-time human posture recognition using an adaptive hybrid classifier[J]. International Journal of Machine Learning and Cybernetics, 2020:1-11.

[6] Wang J, Yuan R, Shi H. Emotional state representation and detection method of users in library space based on body posture recognition[J]. Digital Library Perspectives, 2020, ahead-of-print(ahead-of-print).

[7] Zhen H, Miao H C, Ning W D, et al. A wearable-based posture recognition system with AI-assisted approach for healthcare IoT[J]. Future Generation Computer Systems, 2021.

[8] H. Diao et al., "Deep Residual Networks for Sleep Posture Recognition With Unobtrusive Miniature Scale Smart Mat System," in IEEE Transactions on Biomedical Circuits and Systems, vol. 15, no. 1, pp. 111-121, Feb. 2021, doi: 10.1109/TBCAS.2021.3053602.

[9] Li H, Yin S, Liu J, et al. Novel gaussian approximate filter method for stochastic non-linear system[J]. International Journal of Innovative Computing, Information and Control. 13(1): 201-218, 2017.

[10] Shoulin Yin, Ye Zhang, Shahid Karim. Large Scale Remote Sensing Image Segmentation Based on Fuzzy Region Competition and Gaussian Mixture Model[J]. IEEE Access. volume 6, pp: 26069 - 26080, 2018.

[11] Lin L, Zhang G, Wang J, et al. Utilizing transfer learning of pre-trained AlexNet and relevance vector machine for regression for predicting healthy older adult's brain age from structural MRI[J]. Multimedia Tools and Applications, 2021(1).

[12] Yin Shoulin, Liu Jie, Li Hang. A Self-Supervised Learning Method for Shadow Detection in Remote Sensing Imagery[J]. 3D Research, vol. 9, no. 4, December 1, 2018.

[13] K. Liu, S. Yu and S. Liu, "An Improved InceptionV3 Network for Obscured Ship Classification in Remote Sensing Images," in IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, vol. 13, pp. 4738-4747, 2020, doi: 10.1109/JSTARS.2020.3017676.

[14] Shoulin Yin, Hang Li, Asif Ali Laghari, et al. A Bagging Strategy-Based Kernel Extreme Learning Machine for Complex Network Intrusion Detection[J]. EAI Endorsed Transactions on Scalable Information Systems. 21(33), e8, 2021. http://dx.doi.org/10.4108/eai.6-10-2021.171247

[15] Elpeltagy M, Sallam H. Automatic prediction of COVID19 from chest images using modified ResNet50[J]. Multimedia Tools and Applications, 2021, 80(17):26451-26463.

[16] Liu J, Wang X. Correction to: Early recognition oftomato gray leaf spot disease based on MobileNetv2-YOLOv3 model[J]. Plant Methods, 2021, 17(1):1-1.

[17] Yin, S., Li, H. & Teng, L. Airport Detection Based on Improved Faster RCNN in Large Scale Remote Sensing Images [J]. Sensing and Imaging, vol. 21, 2020. https://doi.org/10.1007/s11220-020-00314-2

[18] Li M, Wang Y, Wang C. Recursive residual atrous spatial pyramid pooling network for single image deraining[J]. Signal Processing Image Communication, 2021, 99(4):116430.

[19] Sun Y, Wang L, Chen Y, et al. Accurate Lane Detection with Atrous Convolution and Spatial Pyramid Pooling for Autonomous Driving[C]// 2019 IEEE International Conference on Robotics and Biomimetics (ROBIO). IEEE, 2020.

[20] Yong S T, Lim K M, Tee C, et al. Convolutional neural network with spatial pyramid pooling for hand gesture recognition[J]. Neural Computing and Applications, 2020:1-13.

[21] Shoulin Yin, Hang Li*, Desheng Liu and Shahid Karim. Active Contour Modal Based on Density-oriented BIRCH Clustering Method for Medical Image Segmentation [J]. Multimedia Tools and Applications. Vol. 79, pp. 31049-31068, 2020.

[22] S. Yin and H. Li. Hot Region Selection Based on Selective Search and Modified Fuzzy C-Means in Remote Sensing Images[J]. IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, vol. 13, pp. 5862-5871, 2020, doi: 10.1109/JSTARS.2020.3025582.