

EXPRESSION RECOGNITION SYSTEM BASED ON DEEP LEARNING FRAMEWORK

Dan Li ^{1*}, Jinping Sun ^{1*}, Weiwei Liu ² and Likai Wang ²

1. School of Information Engineering(School of Big Data), Xuzhou University of Technology, Xuzhou, Jiangsu, China

2. Traffic police detachment of Xuzhou Public Security Bureau, Xuzhou Jiangsu, China

Corresponding Authors' Email: lidanonline@xzit.edu.cn, 313272361@qq.com

Abstract

Due to the changes of expression, background, position and noise, the automatic recognition of facial expression image is a challenge for computer. The system uses the face detection module in OpenCV and Dlib library, loads 68 key point detection models to detect faces, and annotates the key points on the image. The Fer2013 database is trained to get the position information of 68 key points on the face. The expression set is predicted by the classifier, and the predicted probability is displayed visually..

Keywords: facial expression recognition; convolutional neural network; feature extraction; Tenserflow framework.

1. INTRODUCTION

In recent years, with the continuous development of deep learning[1,2], the research of facial expression recognition has gradually stepped out of the laboratory and into practical application scenarios. At present, deep neural networks are widely used to learn the representation of facial expressions. Facial expressions can best reflect inner activities. With the in-depth research on facial expression recognition, how to obtain better facial expression representation has also become a hot research direction[3]. In addition, facial expression recognition also has many practical application scenarios. For example, after adding the facial expression recognition function, the question answering robot can provide more accurate services according to the owner's expression. In the medical field, doctors can judge the patient's tolerance according to the expression reported by the system. In fatigue driving monitoring, we can judge whether the driver is tired by facial expression.

The traditional method of extracting face features cannot extract deep features. In recent years, deep learning methods have made important contributions in the direction of face detection and registration. Compared with traditional methods, these models are more effective and can deal with face poses from a larger angle. Expression recognition based on deep learning method focuses on the face region, which requires preprocessing steps such as face detection, registration and size normalization. Some scholars have proposed a deep learning method[4,5] based on attention convolution network, which can focus on important parts of the face, and has achieved significant improvements over previous models on multiple data sets such as Fer2013, CK+, Ferg and Jaffe. Wang et al.[6] proposed a new method for static facial expression recognition. Through experiments on Fer2013 data set, a group of static images are divided into seven basic emotions by using CNN model, and effective classification is automatically realized. In the test set, the accuracy of facial expression recognition reached 68.79%. He et al.[7] proposed a multi-channel deep neural network, which learns and integrates space-

time features. The basic idea of this method is to select the optical flow in the transformation process between the emotional face and the neutral face as a kind of temporal information of expression, and at the same time, it can also take an emotional gray-scale image generated by the face as spatial information. The recognition rate of this method is 98.38% in CK + database, 99.17% in Rafd database and 99.59% in MMI database. Inspired by the illumination invariant LBP descriptor, Levi et al.[8] extracted LBP features based on RGB images and mapped them to 3D space, used the LBP 3D mapping and RGB images to train convolutional neural networks respectively, and fused them at the decision level. In addition, Zhang et al.[9] extracted SIFT features based on 68 key points of the face, thereby converting the image into a 68*128 feature matrix as the input of CNN. Unlike ordinary CNN, because each row in the feature matrix represents SIFT features on a key point of the face, one-dimensional convolution and one-dimensional pooling are used to further learn the high-level semantic features related to expression in the feature matrix. Yu et al.[10] used multiple face detectors to detect faces and designed multiple CNN models. Each CNN model was first pre trained by other natural scene expression databases, and then migrated to the target natural scene expression database

Although the feature learning ability of deep learning is very strong, there are still some deficiencies in the application of human expression recognition. First of all, in order to avoid excessive fitting, the deep neural network has a lot of additional needs in the amount of training data. However, by analyzing the existing facial expression database, it cannot train a neural network with a certain depth architecture. In addition, based on the differences of personal attributes, such as actual age, gender, male and female, as well as ethnic background and expressive level, the differences between their subjects can still not be ignored. Of course, in addition to the differences caused by the main body, the changes of posture, illumination intensity and occlusion also have a greater impact on the results. The relationship between

these factors and facial expressions is nonlinear coupling. Therefore, in order to solve the problem of huge intra class variability and learn effective specific expression features, the requirements of deep network must be improved.

This paper designs a human face expression recognition model based on convolutional neural network. The model includes face recognition, face key point description, face expression recognition and face expression classification sub modules. The model is developed in Python language. Based on the deep learning Tensorflow framework, it is developed by using OpenCV, Dlib, Numpy, and Keras libraries.

1.1 data set

A. Fer2013 dataset

In the challenge of feature representation of ICML (International Conference on Machine Learning) in 2013, Google Image Search API proposed Fer2013 dataset. This dataset is a large and less restricted data set collected online. The pictures in the dataset are unified to 48*48. The whole data set includes 28709 images for training, 3589 images for validation and 3589 images for testing. Each image is accompanied by seven expression tags, namely (anger, discrepancy, fear, happiness, safety, surprise and neutral).



Figure 1 Fer2013 database samples

B. CK+ dataset

The extended Cohn-Kanade (CK+) database contains 593 video sequences from 123 subjects, including the label of expression and the label of action units. The database is composed of 123 test objects, including all kinds of expressions and labels of various parts of motor units. The test objects are college students aged 18-30 years old. There are 593 dynamic image sequences, 65% of which are women. The last frame of the dynamic image sequence is used as action unit labels, and 327 sequences containing emotional labels.

C. WIDER FACE dataset

The wide face dataset is used for the evaluation of face detection algorithm. This dataset was proposed by the

Chinese University of Hong Kong to solve the problem of saturation of face detection algorithm in the existing dataset. The whole dataset contains 32203 images and 393703 manually labeled faces. The number and scale of images are much larger than other datasets. These images include face images in different scenes, including face images with fuzzy image quality and different attitudes. An image contains dozens or even hundreds of faces.

2.OVERALL STRUCTURE OF SYSTEM MODEL

Figure 2 shows the flow chart of the whole model, which includes the following five major processes: image preprocessing, feature extraction, model training, expression classification and visualization results.

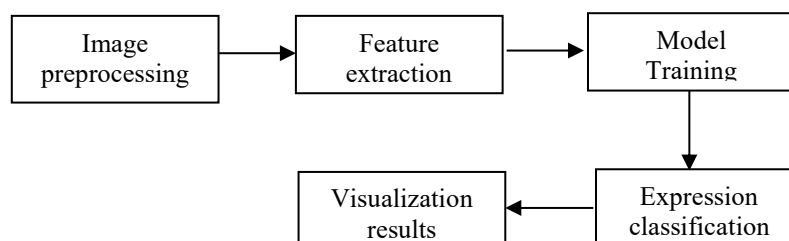


Figure 2 Model processing flow chart

The system mainly includes the following functions: 1) Face recognition sub function: when face collection information appears in the window after running the

program, the background code will use the face detection module in OpenCv and Dlib library to implement face detection. 2) Face key point recognition: use the key point

gradually exposed. When the number of layers of neural network is deepened, the effect of learning is not positively correlated with the change of the number of layers. In other words, the convolutional neural network model appears degradation. Researchers have proposed three possibilities for such problems: over fitting and gradient disappearance or gradient explosion. The increase of network layers cannot improve the classification performance, but reduce the convergence rate and accuracy.

ResNet[11,12] is a feature extraction network model proposed by three Chinese scientists, and it is also one of the most widely used CNN networks at present. From the perspective of information theory, ResNet introduces direct mapping, and the network must contain more information than the previous layer of network, avoiding the problem that the amount of information will decay with the deepening of layers in the forward propagation process caused by DPI (data processing inequality).

2) Residual block

The residual block is the constituent unit of the residual network, which is composed of residual and direct mapping. It is expressed by formula as follows:

$$y_l = h(x_l) + F(x_l, W_l) \quad (1)$$

$$x_{l+1} = f(y_l) \quad (2)$$

For deeper layers, it can be expressed as:

$$x_L = x_l + \sum_{i=1}^{L-1} F(x_i, W_i) \quad (3)$$

The structure of the residual block is shown in Figure 5, where weight refers to convolution:

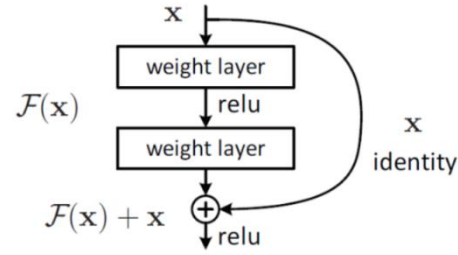


Figure 5 Residual block structure

The loss function loss is expressed as:

$$\frac{\partial loss}{\partial x_l} = \frac{\partial loss}{\partial x_L} \cdot \frac{\partial x_L}{\partial x_l} = \frac{\partial loss}{\partial x_L} \cdot \left[1 + \frac{\partial}{\partial x_L} \sum_{i=1}^{L-1} F(x_i, W_i) \right] \quad (4)$$

3) Bottleneck module

The bottleneck module skillfully controls the dimension of feature mapping by using a 1*1 convolution kernel, so that the classification number of the subsequent 3*3 convolution kernel will not be affected by the previous results.

4) Residual network structure

ResNet network is an improvement of VGG19 network, which integrates residual units into the network through short circuit mechanism. The important change is that the convolution with step size of 2 is used for down sampling, and GPA (global average pooling) is used to replace the full connection layer. When designing ResNet network, in order to ensure the complexity of network layer, the product of feature map size and number is set as a fixed value. When the size is halved, the number will double. Figure 6 shows the comparison between ResNet structure and other network structures:

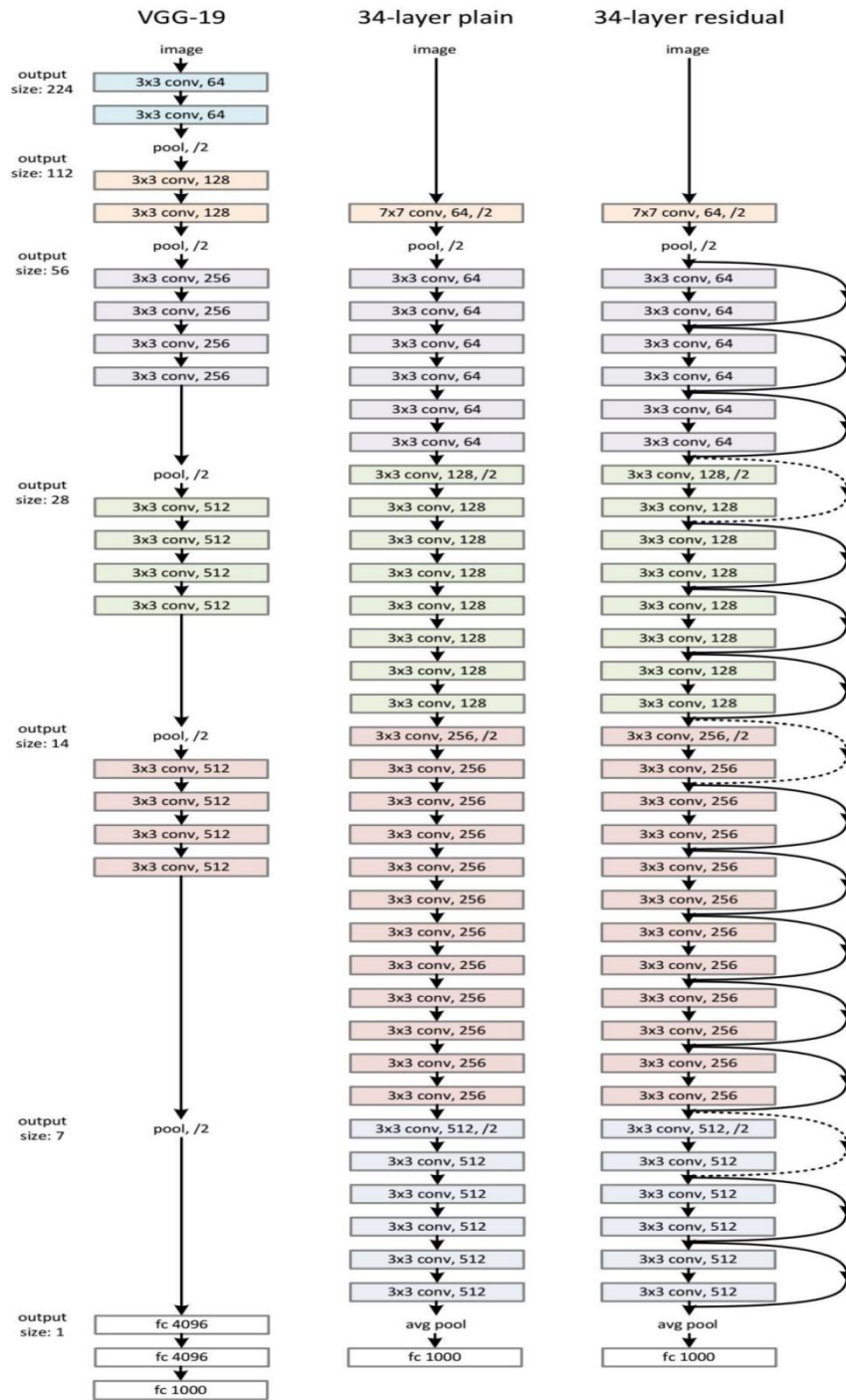


Figure 6 Comparison between ResNet network and other networks

F. Facial expression classification

The cascade classifier of OpenCV is used to load Haar face detection model and torch is used to load CNN or ResNet model respectively. Turn up the camera to enter the shooting state, continue to read the captured video, analyze the video frame by frame, and obtain the gray-scale image of each frame to judge whether there is a face in the image. If one or more human faces are detected, locate and extract the face, and do gray-scale and binary processing on the detected image, and input the processed image into the classifier for analysis. All faces in the camera shooting area are locked by the square box, and

the probability of showing the expression of the current facial state is displayed. For the computer, when the face is recognized, the face data is expanded and normalized immediately, and the neural network model is used to recognize the expression, predict and classify it, and return it to the user.

3.Experimental results and analysis

The model shows the final running results in visual form. The purple box is used to locate the face, and the key points of the face are depicted at the same time, and the corresponding expression probability is displayed on the

interface. The effect of the system is tested under different conditions.

1) Light intensity is different: as shown in Figure 7. There are certain differences in the accuracy of the same expression under different illumination intensity. When

the illumination intensity is good, the accuracy of expression recognition will be improved. When the illumination intensity is weak, although the expression can still be recognized, the accuracy will be reduced.

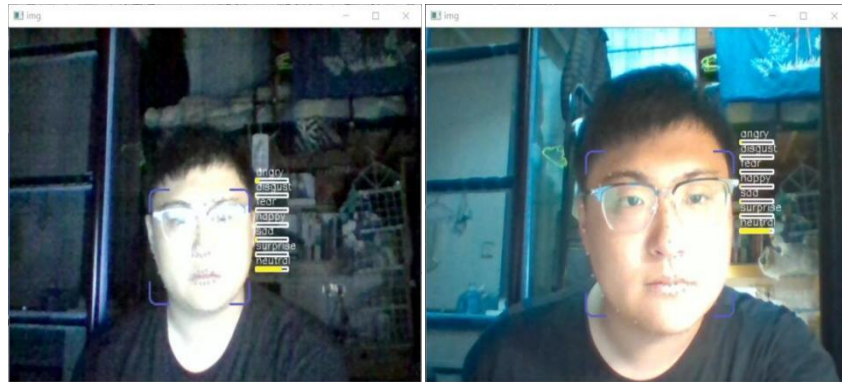


Figure 7 Results of different light intensities

2) Different scenes: the detection accuracy will also affect face recognition in different scenes, as shown in Figure 8. In the indoor test, the accuracy will decline, which is not

accurate enough, but in the outdoor test, the accuracy will be more accurate.

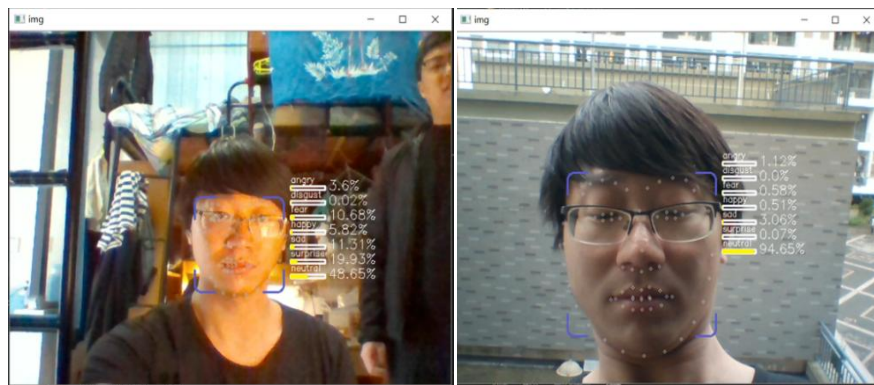


Figure 8 Indoor and outdoor results

3) Different shooting angles: in the same scene and under the same light, test the face from different angles, as shown in Figure 9. It is found that when the face is not in

the lens, the face cannot be detected, because the characteristic value of the face cannot be extracted at this time, so it cannot be detected.

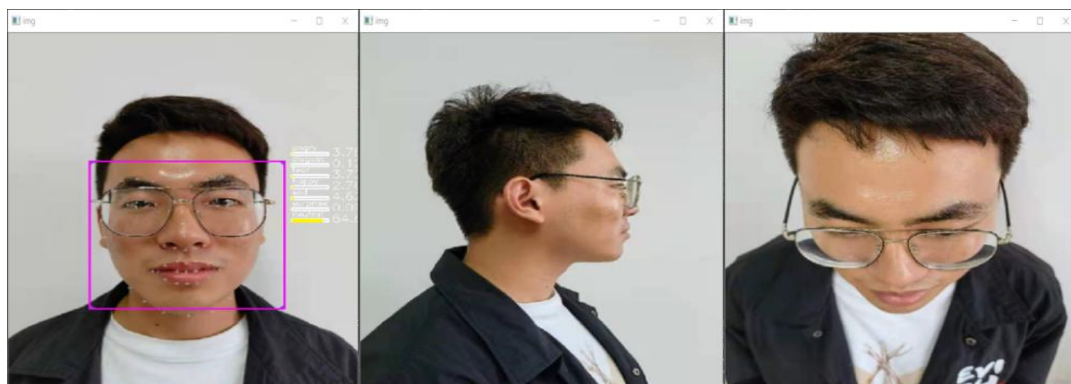


Figure 9 Results at different angles

(4) Gender difference: when using the test set to train the model, the gender of the user is not distinguished, but the gap between the male and female ratio is not large, so the

impact of gender on the accuracy is not particularly obvious in the final test. Under the same expression, the

accuracy of male and female testers is almost the same, as shown in Figure 10.

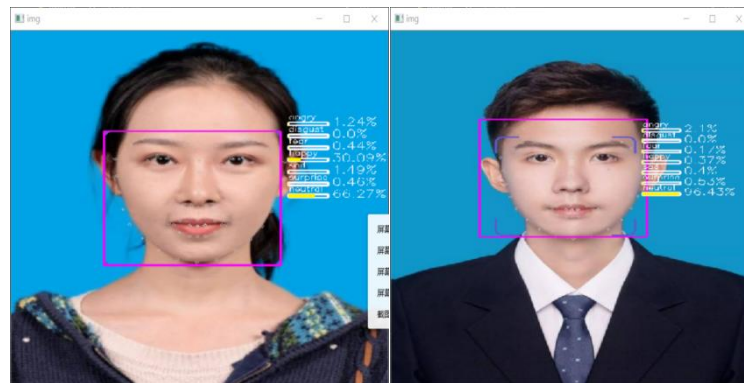


Figure 10 Female and male test results

5) Different ethnic backgrounds: because most of the training sets use pictures of white and black people, and the facial features of different people are also different,

the accuracy of the system in detecting white and black people will be higher during the test. As shown in Figure 11.

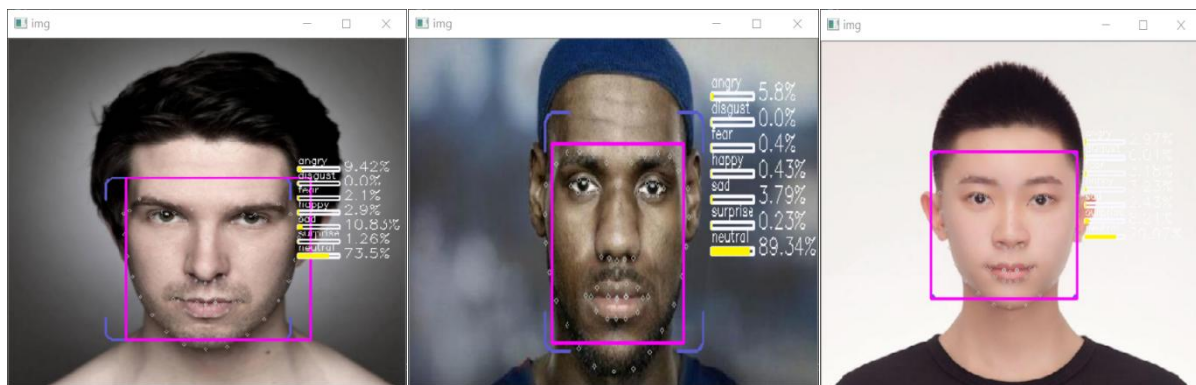


Figure 11 schematic diagram of test results of white, black and yellow people

6) Face occlusion: during the test, it is found that whether to wear glasses has a certain impact on the results. As shown in Figure 12, because facial occlusions such as glasses and masks will hinder the extraction of facial

feature values, some key points cannot be detected, and although the pictures in the training set used have pictures of people wearing glasses, the number is not large. These factors together lead to the deviation of the test results.

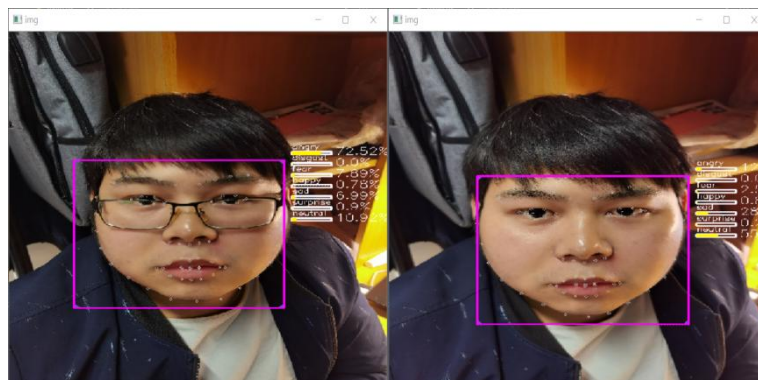


Figure 12 schematic diagram of results with and without glasses

4.CONCLUSION

A facial expression recognition model based on convolutional neural network is developed. The model designs convolutional neural network and carries out

facial expression recognition training on Fer2013 and other datasets. Then, the trained model is used to complete the required detection for the real-time face. The experimental results show that the facial expression recognition system developed based on Tensorflow deep

learning framework is feasible and reliable, and can better meet the expression recognition needs of a certain environment.

ACKNOWLEDGMENT

This work was supported in part by the Xuzhou Science and Technology Plan Project (Grant: KC21303), Jiangsu Industry University Research Cooperation Project (Grant:BY2021159), Jiangsu Educational Science "14th five year plan" Project(Grant:C-c/2021/01/65), the Sixth "333 project" of Jiangsu Province, Natural Science Research Projects of Colleges and Universities in Jiangsu Province(Grant: 22KJA520012).

REFERENCES

- [1] Bo, C. , et al. "The Hierarchical Beta Process for Convolutional Factor Analysis and Deep Learning." ICML Omnipress, 2020.
- [2] Memon, F. A. , et al. "Predicting Actions in Videos and Action-Based Segmentation Using Deep Learning." IEEE Access PP.99(2021):1-1.
- [3] Shahid, A. R. , S. Khan , and H. Yan . "Contour and region harmonic features for sub-local facial expression recognition." Journal of Visual Communication and Image Representation 73.2(2020):102949.
- [4] Zhang, F. , et al. "A Unified Deep Model for Joint Facial Expression Recognition, Face Synthesis, and Face Alignment." IEEE Transactions on Image Processing 29(2020):6574-6589.
- [5] Trimech, I. H. , A. Maalej , and N. Amara . "Facial Expression Recognition Using 3D Points Aware Deep Neural Network." Traitement du Signal 38.2(2021):321-330.
- [6] Wang X, Huang J, Zhu J, Yang M, Yang F. "Facial expression recognition with deep learning. " Proceedings of the 10th International Conference on Internet Multimedia Computing and Service. New York, NY, USA: Association for Computing Machinery, 2018
- [7] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in Proceedings of the IEEE conference on computer vision and pattern recognition, 2016.
- [8] Levi G and Hassner T. "Emotion recognition in the wild via convolutional neural networks and mapped binary patterns. " in Proceedings of the 2015 ACM on international conference on multimodal interaction. 2015. ACM: 503-510
- [9] Zhang T, Zheng W, Cui Z, et al., "A deep neural network-driven feature learning method for multi-view facial expression recognition. " IEEE Transactions on Multimedia, 2016.18(12):2528-2536
- [10] Yu Z and Zhang C. "Image based static facial expression recognition with multiple deep network learning. " in Proceedings of the 2015 ACM on International Conference on Multimodal Interaction. 2015. ACM: 435-442
- [11] Tachibana, R. , et al. "Comparative performance of self-supervised 3D-ResNet-GAN for electronic cleansing in single- and dual-energy CT colonography." Imaging Informatics for Healthcare, Research, and Applications 2021.
- [12] Ikechukwu, A. V. , et al. "ResNet-50 vs VGG-19 vs training from scratch: A comparative analysis of the segmentation and classification of Pneumonia from chest X-ray images." Global Transitions Proceedings 2. 2(2021):375-381.