

A method for image style transfer based on human-computer interaction gesture recognition

Yanyu Liu^a, Zhihong Zhang^b, Bo Li^{c,*}

School of Electronic Information Engineering, Beihai Vocational College, Beihai, 536000, China
^aliuyanyu@bhzyxy.edu.cn, ^bzhihong_zhang90@163.com, ^c13097794979@163.com

*Corresponding author

Abstract: Image style transfer is a popular research topic in the fields of computer graphics and multimedia, especially artistic stylization. The image style learning algorithm mainly studies the use of computer graphics and machine learning methods for automatic artistic rendering and intelligent processing of real sample data. The current mainstream methods mainly focus on learning static samples of artistic images. However, the information contained in static sample data is isolated and discontinuous; it is difficult to ensure the overall consistency of image style transfer. This article aims to recognize the real gesture movements of painters during the process of image style transfer, and propose a theoretical model and design method for image style transfer based on sequence task learning theory. Mainly completed the following research work: (1) Complete the design path and verification experiment of human-computer interaction gesture recognition in the air; use gesture recognition model, color space, color distribution model, color point probability, and maximum inter class variance method to complete the binary processing of images, and use Babbitt distance and learners to learn the intrinsic patterns of actions, ensuring that the gesture recognition design method in the air can be implemented; (2) Analyze and understand the real gesture process of painters, and design a dynamic decision-making model for image style transfer based on parameter exploration. Implement an image style transfer assistance system based on human-computer interaction gesture recognition.

Keywords: Image style transfer; gesture recognition; natural human-computer interaction; neural network

1. Introduction

Image style transfer has a large application market and has been an important research area in computer vision and image processing for many years. This technology aims to transfer the style of one image to another. In recent years, algorithms based on deep neural networks have shown remarkable performance in various intelligent applications such as facial recognition, pedestrian re identification, and vehicle tracking. Image processing algorithms based on deep neural networks can extract high-level features, which are superior to traditional algorithms based on low-level visual features. For example, the shallow layers of VGGNet^[1] (such as conv1_1 and conv1_2) extract simple features such as edges and brightness. And deep layers (such as conv5_1, conv5_2) extract complex features. VGGNet aims to extract deep features from input images. However, unlike VGGNet, the image style transfer algorithm aims to generate images based on input features. Specifically, image style transfer technology refers to designating the input image as the base image, also known as the content image. At the same time, specify another image or images as the desired image style. As shown in Figure 1, there are many different styles that can be considered as style images. The image style transfer algorithm transforms the image style while ensuring the structure of the content image, so that the final output composite image presents a perfect combination of the input image content and the desired style.

The first step in image style transfer is to analyze an image of a certain style and establish a mathematical or statistical model for that style, then modify the image to be transmitted to better conform to the established model, as shown in Figure 2. The results indicate that the effect is good, but there is still a major drawback: a program can only achieve a certain style or scene. Therefore, the practical application of traditional style transfer research is very limited. As early as the early 2000s, many scholars began to study the issue of image style transfer, but at that time they were focused on the synthesis and transfer of textures. The mathematical methods used are mainly statistical methods for

various image transformations, such as wavelet transform, which have limited effectiveness. Until recent years, inspired by the excellent performance of deep neural networks in large-scale image classification, and relying on their powerful multi-level image feature extraction and representation capabilities, this problem has been well solved, and their performance has been significantly improved.

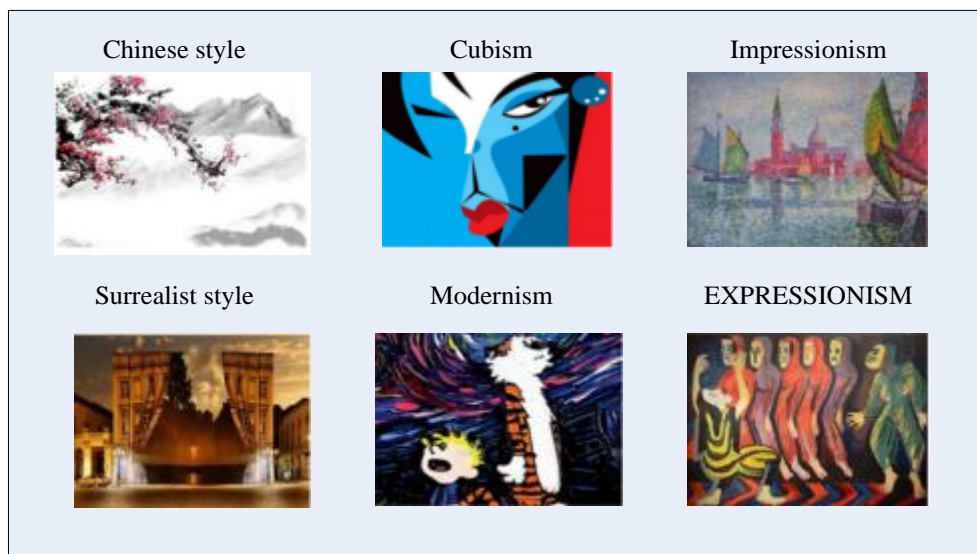


Figure 1 Art Style Image Case

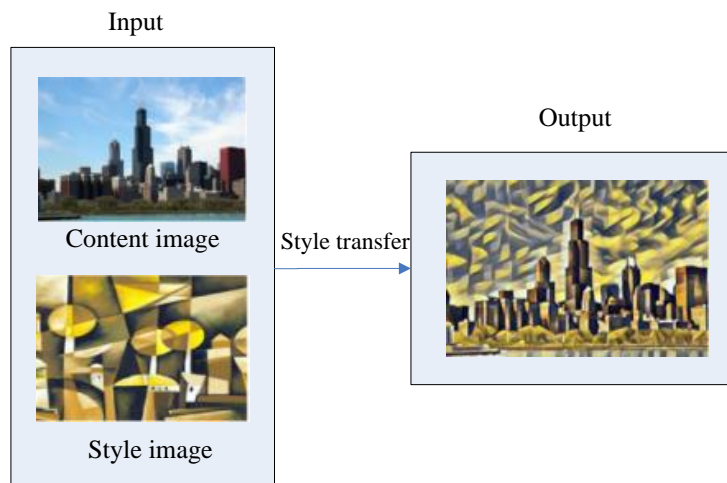


Figure 2 Image Style Transfer Case

Researchers have conducted a series of targeted studies on image style transfer algorithms. To address the issue of semantic mismatch in image style processing, Xie Chuan et al. proposed an image style transfer algorithm based on semantic segmentation, which extracts semantic information of content and style images through a masked R-CNN network. Then, this information will guide the style conversion process. The dataset results from Celeba and Wikiart emphasize that this method is more effective than existing style transfer methods in avoiding image style transfer and rendering issues. In addition, compared with previous style transfer methods, it can successfully overcome the problem of semantic mismatch in the image retrieval process [2]. Kim Minseong et al. developed an end-to-end learning method for improving image style transfer techniques using encoder and decoder networks. This method reduces the computational complexity of the current association perception feature alignment model and minimizes the channel redundancy of encoded features during network training. It is also an ideal choice for multi-scale image style transfer [3]. To address the issue of preserving textures in image style transfer models, Ding H et al. designed an image style transfer model that utilizes wavelet transform and deep neural network techniques to synthesize styles and restore details. This model can match semantic relationships with the help of attention mechanisms and semantic segmentation, while also processing images with small details. The experimental results have verified the superiority of the model in preserving image texture [4]. Li et al. constructed an image segmentation

model using an improved semantic segmentation network DeepLab2. This model has the ability to transfer local image styles, and experimental results have confirmed its practicality and transfer efficiency [5]. Chen et al. developed an image style transfer model that utilizes multiple convolutional filters. Meanwhile, the model also includes an autoencoder and a style library learning component. Filter banks are used for multi parameter image smoothing and denoising, and allow the model to produce results comparable to single parameter settings [6]. In order to improve the efficiency of generating style transfer, Huang L et al. developed an algorithm that relies on semantic segmentation and residual networks, as well as Visual Geometry Group Network for feature extraction. According to the experimental results, the model has improved the efficiency of local image style transfer and generation. In addition, this technology can also be developed in the fields of entertainment, film and television, healthcare, and industry [7]. In order to solve the challenging problem of semantic image style transfer, Liao Yongsheng et al. proposed a new method with context awareness and semantics. This model includes a global context network and a local context network, and focuses on programmatic image derivation and semantic context style transfer. According to the experimental results, this method produces stylized results that are more in line with human semantic perception than existing models [8]. Gatys et al. used VGG network for hierarchical image feature extraction, effectively representing images. They used CNN to fuse the semantic content of images with different styles. This new style transfer method has brought good results in the academic community and led to a large number of research results on using deep learning style transfer [9]. Ulyanov et al. trained a compact feed forward convolutional network to generate multiple samples of the same texture of any size and convey the style of art from a given image to any other image. The resulting network is well lit and can produce textures and other features of comparable quality to Gatys, but at a speed hundreds of times faster [10].

In summary, based on the analysis of the current research status, level, and development trends at home and abroad, and our research foundation, although a large number of achievements have been made in the study of image style transfer technology, there are still some shortcomings and deficiencies, mainly reflected in:

①How to accurately capture the real gestures of painters: Currently, mainstream methods mainly focus on studying static samples of artistic images. The processing methods for dynamic behavior sequence data are limited, making it difficult to analyze and understand the dynamic drawing process;

②How to design a dynamic decision model for image style transfer: Currently, there are not many research results on analyzing and understanding the real gesture action process of painters. However, since the decision of real gesture action in the painting process will directly affect the final effect of image style transfer, it is necessary to design a dynamic decision model for image style transfer.

2. Overall frameworks

This article aims to propose a theoretical framework for dynamic decision-making of image stylization transfer for large-scale complex data in the creative process, to learn the gesture behavior of painters in the real creative process, and apply it to automatic art drawing strategies. The basic framework of the research approach for this project is shown in Figure 3.

The research plan in this article is mainly implemented in three steps: firstly, design a high-precision human-machine interaction gesture recognition device, and perform state analysis and feature extraction on the obtained sequence data; Secondly, relying on reverse reinforcement learning methods, feature structured organization and feedback evaluation function construction are carried out; Finally, study the stylized dynamic strategy method based on parameter exploration. Ultimately, an image style transfer assistance system based on dynamic behavior sequence data will be implemented.

Considering the complexity of gesture recognition data in real artistic creativity, this leads to a large variance in the gradient estimation results for such complex sequence data, affecting the stability of model training and slow convergence. To this end, we explore a new dynamic decision-making mechanism for art styles. This article intends to propose a dynamic drawing strategy model and an optimization algorithm based on mixed guided samples (HGS-PGPE) for complex sequence data of creative behavior, based on the parameter exploration strategy gradient (PGPE) theory.

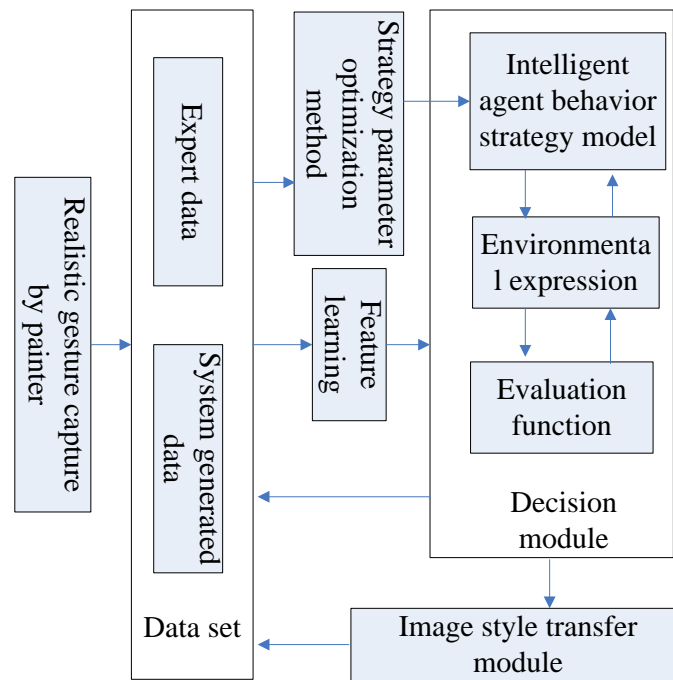


Figure 3 Overall frameworks for dynamic decision-making of image style transfer

3. Design path and verification experiment for human-machine interaction gesture recognition in the air

(1) Overall framework

With the increasingly mature application of multi-channel interaction means, it can evolve into a new service supply platform under the Internet environment, and such network platforms can break the space constraints and have a more natural interaction mode. Due to the increasing frequency of software and hardware updates in the information age, traditional human-computer interaction gesture recognition methods have been optimized and upgraded. It is necessary to establish gesture recognition models by introducing semantic feedback and information integration mechanisms from artificial intelligence technology to further ensure the accuracy of gesture recognition.

The specific gesture recognition framework includes: user module, language agent, vibrator agent, presentation module, dialogue control module, and application module.

The user module is mainly responsible for collecting interactive tools and devices in the interactive object library, completing information input, and obtaining corresponding information through feedback devices.

Language agent and vibrator agent are both part of the human-computer interaction object module, and their performance lies in controlling input and output devices. They can convert user gesture actions into interaction primitives and store them in the interaction primitive queue. This module can not only receive the presentation module driver, but also transmit the processing results to the user. Because abstracting gesture actions can obtain interaction primitives, each level primitive needs to maintain a certain degree of association with the device, and the purpose of creating each module is to be responsible for adding, deleting, and changing devices to ensure channel integrity.

The performance of the presentation module lies in abstracting the interaction primitives within the interaction object, obtaining the interaction concept primitives, and receiving the response information returned by the dialogue control module. The concept primitive refers to the information used to process human-computer interaction object modules, while the main object for object display processing is the dialogue control module.

The performance of the dialogue control module is demonstrated by integrating user interaction information to parse interaction intentions, creating interaction task primitives, and adding them to the primitive queue. This module can receive application feedback formed by the application and then pass it on to the corresponding users through the presentation module.

The function of the application module is to process tasks submitted by users and pass the processing results to the dialogue control module. The structure of this module includes a non-interactive calculation part, and the calculation results are stored in the abstract part of the user interface before being transmitted to the user after passing through the user interface.

(2) Gesture segmentation

When designing, it is necessary to use the MCG database as the basis, select components Ar and Ab, and express the architecture color distribution in the color space as follows:

$$Pc(x)=\alpha Arg(x, \mu Ar, \sigma Ar)+\alpha Abg(x, \mu Ab, \sigma Ab)$$

Among them, α represents the color component coefficient, g represents the color filtering value, and μ refers to the uniform quantization coefficient. According to the above color distribution expression, it is necessary to use background elimination to merge and build Gaussian models. The specific formula for describing the Gaussian mixture model is:

$$Py(x)=\alpha_h g(x, \mu_h, \sigma_h)+\alpha_f g(x, \mu_f, \sigma_f)$$

Among them, $(\alpha_h, \alpha_f, \mu_h, \mu_f, \sigma_h, \sigma_f)$ refers to the parameters of the probability density function, α_h and α_f have the following relationship:

$$\alpha_h+\alpha_f=1$$

Since gesture behavior is always in front of the body during human-computer interaction, μ_h in the Gaussian mixture model is set as the gesture depth value, while μ_f is used to represent the body depth value. The expression formula for the depth threshold is:

$$Td=(\mu_h+\mu_f)/2$$

Input image pixels into the color distribution model, then obtain the probability of pixel color points, and use the method of presenting them in the grayscale range [0,255] to achieve the goal of reconstructing skin similarity images, ensuring that the skin similarity images can contain 256 grayscale levels, and the similarity can increase with the increase of pixel grayscale values. At the same time, it is necessary to combine the maximum inter class variance method to achieve binary processing of images, using gray pixels to represent skin color points and black pixels to represent remaining points, in order to segment gesture behavior.

(3) Gesture tracking

When designing gesture tracking, it is necessary to create a depth histogram and select a series of discrete functions to represent the image histogram in the grayscale range of [0, Q-1]. The specific calculation formula is:

$$h(r_k)=n_k$$

Among them, r_k represents k-level grayscale, while n_k represents the number of pixels containing r_k grayscale in the image. In the process of designing gesture tracking, attention should also be paid to normalizing the image data and transforming it into a fixed standard form. The specific histogram calculation formula is:

$$P(r_k)=n_k/n$$

Among them, $P(r_k)$ represents the estimated probability of r_k grayscale formation, and the above formula can achieve the superposition of each region of the histogram, with a total of 1. Afterwards, the absolute value of pixels should be used to calculate the absolute depth in the minimum depth value, in order to set the relative depth range, that is, [0, I-1]. If the above formula $h(r_k)=n_k$ is used to represent the relative depth histogram, then the k-level relative depth can be represented as r_k . As for the number of pixels containing r_k depth level, it should be n_k , and the value range of k should be within [0, I-1]. The method of creating the relative depth histogram should be consistent with that of the same grayscale histogram. Generally speaking, probability distribution evaluation methods are used in designing gesture tracking to accurately determine histogram similarity. Therefore, the Babbitt distance based on discrete probability distribution can be used to define the domain boundary. The formula is:

$$Mb(p, q)=-\ln(Bs(p, q))$$

Among them, $Bs(p, q)$ represents the Babbitt coefficient, and after calculating the Babbitt distance, gesture tracking can be calculated using the following formula:

$$G_x = P(r_k) \sum_{x \in X} \sqrt{p(x)} q(x)$$

Among them, x represents the speed of gesture movement, and the above is the implementation path of isolated gesture tracking, which can provide support for precise gesture recognition in the future [11].

(4) Gesture recognition

In order to further optimize the efficiency of gesture recognition and upgrade traditional gesture recognition methods, it is necessary to introduce learners and utilize the learning of posture patterns in each stage of learners to complete gesture recognition more quickly. The specific cascading gesture recognition process is: gesture action to be recognized \rightarrow ①T \rightarrow ②T \rightarrow ③T \rightarrow non controlled posture \rightarrow post-processing.

In the design process, H_i needs to be used to represent the i -level learner. However, when the posture sample reaches level i , if the learner H_i is difficult to recognize with high confidence, the recognition task needs to be conveyed to the next level learner until the recognition result is obtained before the task can be completed. Due to the cascading design concept, learners at all levels can complete the learning of various posture patterns. Therefore, the first two levels of learners can be used to recognize non control postures with lower difficulty, while the relatively difficult parts can be identified by the later learners, thus avoiding the influence of non-control posture recognition. At the same time, when using a base table structure for hierarchical recognition of non-recognition control poses, the increase in the number of levels will to some extent reduce the imbalance of categories between control poses and non-control poses. By focusing the attention of subsequent learning devices on the differential patterns of different poses, it can ensure that even if there is sample imbalance, it will not affect the recognition effect [12].

Assuming that the expression for the set of gesture postures in the air is:

$$G = \{G_1, G_2, G_3, \dots, G_n\}$$

Where n represents the number of posture types. G represents the control action. To accurately identify G , it is necessary to ensure that the non control action \bar{G} satisfies the following set expression:

$$\bar{G} = \{G_2, G_3, G_4, \dots, G_n\}$$

Selecting any set of samples from various posture types as training information, the expression is:

$$G_1 = \{X_{i1}, X_{i2}, X_{i3}, \dots, X_{im}\}$$

In the formula, X_{im} represents the i -th pose with a sample size of m . Afterwards, randomly select a set of cascaded classification devices and describe them using the following formula.

$$T = \{(M_1, F_1), (M_2, F_2), (M_3, F_3), \dots, (M_i, F_i)\}$$

In the formula, M_i represents the i -layer classifier, and F_i represents the corresponding feature of M_i . Therefore, the following feature set formula can be obtained:

$$F = \{F_1, F_2, F_3, \dots, F_i\}$$

The specific gesture recognition process description process is as follows: complete the initialization operation of $C = G, \bar{C} = \Phi, F_i \in F, T = \Phi$; If $C = \{G_1\}$, then it is necessary to extract gesture actions from the gesture set C that have significant differences from G_1 , and then add the relevant gesture actions to the set \bar{C} :

$$C = G_1 - \bar{C}$$

For feature set F , in the process of feature extraction, it is necessary to select F_i , which takes the least amount of time, as well as the training recognition binary classifier, and then make F_i infinitely close to 1, and add M_i misclassified samples \bar{G} to set G ; Add (M_i, F_i) to the cascaded classifier, where the value of i needs to be incremented by 1; Return the cascaded classifier [13].

Assisted by artificial intelligence, aerial gesture recognition in human-computer interaction is a more convenient and natural mode of human-computer interaction, which can further improve the recognition efficiency and accuracy of aerial gestures, approaching the ideal state of human-computer interaction.

(5) Experimental verification

In order to verify the feasibility of the above design, it is necessary to conduct an experiment on human-machine interaction gesture recognition with the assistance of artificial intelligence. During the experiment, KINECT software is used to collect gesture actions, and then multiple images of 50×50 pixels in the dataset are randomly selected as experimental samples. The computer equipment used should be equipped with a 28GHz dual core processor and have a running memory of at least 6GB to ensure computational efficiency. At the same time, during the experimental process, self-learning sparse representation methods and the gesture recognition model mentioned in the article should be combined to recognize gesture actions such as grasping and releasing of the experimental object. The performance evaluation index used is mainly F1, and the larger the value of this parameter, the higher the accuracy of gesture recognition. The specific expression formula is:

$$F1=2 \times \text{PRECISION} \times \text{RECALL} / (\text{PRECISION} + \text{RECALL})$$

Among them, PRECISION represents accuracy, RECALL represents recall rate, and the corresponding calculation formulas for the two are:

$$\text{PRECISION} = \text{TP} / (\text{TP} + \text{FP})$$

$$\text{RECALL} = \text{TP} / (\text{TP} + \text{FN})$$

TP represents the accurately estimated sample size, FP represents the incorrectly evaluated sample size, and FN refers to the estimated sample size for the current category. The specific number of gesture recognition FI indicators is shown in Table 1.

Table 1 Information on the number of FI indicators for gesture recognition

Image Number	Get			Put		
	LEAP Method(%)	Self-learning Sparse Method(%)	Research model(%)	LEAP Method(%)	Self-learning Sparse Method(%)	Research model(%)
5	84.12	85.41	98.65	89.01	88.14	97.48
12	87.56	87.14	99.35	91.04	91.04	97.77
29	85.12	89.34	98.14	88.62	87.36	98.63
37	86.12	85.64	96.34	88.04	84.66	99.14
48	89.00	87.45	97.47	84.71	85.14	98.14

According to Table 1, the recognition performance of the research model is much higher than that of the self-learning sparse method and LEAP method, and the evaluation index is basically above 95%, which is a more ideal gesture recognition method. However, the evaluation index of the self-learning sparse method is below 90%, and the evaluation index of the LEAP method is mostly below 89%, which is difficult to meet practical needs.

4. Dynamic Decision Model for Image Style Transfer

(1) Research on State Analysis and Feature Extraction for Large Scale Complex Dynamic Behavior Sequence Data

In traditional reinforcement learning methods for sequential tasks, the representation of states in the environment is often manually designed and given by experts based on relevant domain knowledge. The quality of state description directly affects the performance of learning algorithms. In this article, the author utilizes the powerful description and abstraction capabilities of deep neural networks for high-dimensional data to assist in the automatic representation of states. However, due to the continuous sequence decision-making characteristics of reinforcement learning, it differs from existing mature deep learning techniques, such as the absence of a large amount of annotated data, temporal correlation of sequence data, and non-independent and identically distributed data. Therefore, classical deep learning methods cannot be directly applied to reinforcement learning frameworks. To this end, firstly, this article designs and develops a creative behavior gesture recognition device; Secondly, a correlation analysis method is proposed to match the sequence data of creative behavior with the

sequence data of the painting generation process; Finally, design a deep network structure for state representation that conforms to the characteristics of the reinforcement learning framework itself.

(2) Research on Feature Structured Organization and Feedback Evaluation Method Based on Reverse Reinforcement Learning

Traditional image style transfer mainly involves static feature analysis and stylized description of artistic works. However, because the information contained in static sample data is isolated and discontinuous, it is difficult to ensure the overall consistency of image style transfer. In terms of local optimization of complex dynamic data, if the behavior cloning method based on supervised learning is not ideal, after experimental testing, we believe that the feature structured organization and feedback evaluation method based on reverse reinforcement learning is relatively effective in achieving consistency in image style transfer.

(3) Research on Stylistic Dynamic Strategy Method Based on Parameter Exploration

The result of gradient estimation for complex sequence data of art works has a large variance, which affects the stability of model training and converges slowly. To this end, we explore a new dynamic decision-making mechanism for image styles. This project intends to introduce the policy gradient theory based on parameter exploration (PGPE) into the framework of artist drawing behavior assisted decision-making, combining the deterministic policy model with the probability distribution of real data based on prior knowledge of real artists. At the same time, it is proposed to introduce the optimal baseline method to further ensure the stability of policy learning performance.

The architecture of the image style conversion assistance system based on human-computer interaction gesture recognition is shown in Figure 4.

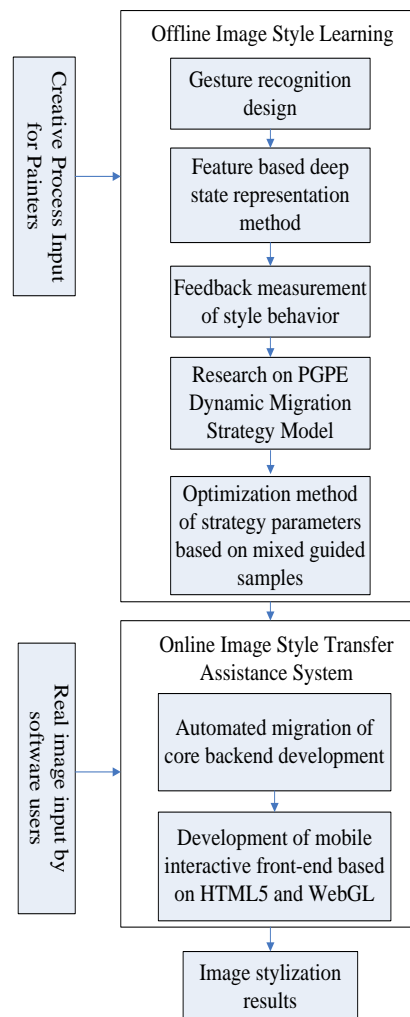


Figure 4 Architecture of Image Style Transfer Assistance System Based on Human Computer Interaction Gesture Recognition

5. Conclusion

This article proposes a theoretical model and design method for intelligent art style transfer based on sequence task learning theory, targeting the real gesture movements of painters in the process of image style transfer. We have completed the design path and verification experiment of human-computer interaction gesture recognition, analyzed and understood the real gesture action process of painters, and designed a dynamic decision model for image style transfer based on parameter exploration. We have designed and implemented an image style transfer assistance system based on human-computer interaction gesture recognition by integrating image style transfer theory and algorithms into the open-source "Online Reinforcement Learning Intelligent Agents (REINFORCEjs)" framework at Stanford University in the United States. This study has certain theoretical innovation and significant practical application value.

Acknowledgment

The research is supported by the 2024 Research Project of Beihai Vocational College in Guangxi, China, titled "Research on Image Style migration Method Based on Dynamic Behavior Sequence Data", No. 2024YKY05.

References

- [1] K. Simonyan, A. Zisserman, *Very deep convolutional networks for large-scale image recognition*, *Comput. Sci.* (2014).
- [2] C. Xie, Z.Z. Wang, H.B. Chen, X.L. Ma, W. Xing, L. Zhao, W. Song, Z.J. Lin, *Image style transfer algorithm based on semantic segmentation*, *IEEE Access* 9 (1) (2021) 274–287.
- [3] M. Kim, H.C. Choi, *Uncorrelated feature encoding for faster image style transfer*, *Neural Network*. 140 (1) (2021) 148–157.
- [4] H. Ding, G. Fu, Q. Yan, C. Jiang, T. Cao, W. Li, S. Hu, C. Xiao, *Deep attentive style transfer for images with wavelet decomposition*, *Inf. Sci.: Int. J.* 587 (587) (2022) 63–81.
- [5] Y. Li, Y. Wang, H.W. Tseng, H.K. Hang, C. Chen, *Optimized and improved methods of image style transfer for local reinforcement*, *Sensor. Mater.: An Int. J. Sensor Technol.* 33 (9) (2021) 3325–3332, pt.4.
- [6] D. Chen, L. Yuan, J. Liao, N. Yu, G. Hua, *Explicit filterbank learning for neural image style transfer and image processing*, *IEEE Trans. Pattern Anal. Mach. Intell.* 43 (7) (2021) 2373–2387.
- [7] L. Huang, P. Wang, C.F. Yang, T. Hsien-Wei, *Rapid local image style transfer method based on residual convolutional neural network*, *Sensor. Mater.: An Int. J. Sensor Technol.* 33 (4) (2021) 1343–1352. Pt.2.
- [8] Y.S. Liao, C.R. Huang, *Semantic context-aware image style transfer*, *IEEE Trans. Image Process.* 31 (1) (2022) 1911–1923.
- [9] L.A. Gatys, A.S. Ecker, M. Bethge, *Image style transfer using convolutional neural networks*, in: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE Computer Society, 2016.
- [10] D. Ulyanov, V. Lebedev, A. Vedaldi, V.S. Lempitsky, *Texture networks: Feedforward synthesis of textures and stylized images*, in: *ICML, 2017*, pp. 1349–1357.
- [11] Zhang Menghuan, Wang Yagang, *Gesture recognition and assisted identity authentication based on CNN and ultrasound sensing*, *Sensors and Microsystems.* 41 (05)(2022) 110-113+117
- [12] Hu Zongcheng, Zhou Yatong, Shi Baojun. *Static gesture recognition algorithm combining attention mechanism and feature fusion*, *Computer Engineering.* 48 (04)(2022)240-246
- [13] Zhang Yunfeng, Zhang Chao, Lv Zhao. *Dynamic gesture recognition method based on keypoint residual fully connected network*, *Journal of Anhui University (Natural Science Edition).* 46 (02)(2022) 30-38