

HUMAN ACTION RECOGNITION IN DARK VIDEOS

Anonymous authors

Paper under double-blind review

ABSTRACT

In recent years, Human Action Recognition (HAR) in videos has garnered extensive research interest, particularly in applications such as surveillance and autonomous driving. However, the majority of these efforts have been concentrated on well-lit video data, resulting in State of the Art (SOTA) algorithms that often perform poorly in low-light or dark conditions. This paper focuses on addressing the challenge of enhancing action recognition performance in dark or low-light videos. We begin by comparing various frame sampling methods, including uniform sampling, MinMax-Sampling, and random sampling. Furthermore, we investigate image enhancement techniques such as Zero-Reference Deep Curve Estimation (Zero-DCE) and Contrast Limited Adaptive Histogram Equalization (CLAHE) to improve the contrast and brightness of the extracted image frames. Subsequently, the enhanced images are fed into a ResNet-18 (2+1)D neural network for feature extraction, followed by classification using a Support Vector Machine (SVM) and accuracy assessment. The paper also discusses the integration of SOTA methods from each stage into an end-to-end model. By combining these methodologies, this study provides new insights and practical solutions for effective action recognition in low-light conditions.

1 INTRODUCTION

Human action recognition (HAR) has emerged as a pivotal area in the domain of computer vision, boasting significant implications across various applications including human-computer interfaces (Meng et al., 2007), military operations, autonomous vehicles, and intelligent home systems (Zhang et al., 2007). Conventional State of the Art (SOTA) methods, encompassing various techniques, have predominantly concentrated on utilizing high-definition video data for both training and validation phases (Ghadiyaram et al., 2019) (Bello et al., 2021). These methods necessitate an environment with optimal lighting and clear visibility. However, such approaches exhibit suboptimal performance when applied to videos captured in nocturnal or low-light settings, highlighting a critical gap in their applicability under diverse lighting conditions.

In practical scenarios, the efficacy of computer vision transcends human capabilities, particularly in contexts where human visual perception is limited, such as nocturnal or dimly lit environments. This is notably evident in applications like nocturnal security surveillance and military rescue missions conducted under minimal lighting conditions. Achieving precise human action recognition in such low-light settings presents a formidable challenge, demanding extensive research and innovative solutions.

In recent times, the academic community has shown growing interest in refining the quality of low-light video content (Chen et al., 2019) Jiang & Zheng (2019). Certain conventional methodologies have pivoted towards enhancing the visibility in dimly lit footage. However, these methods occasionally compromise the integrity of the data, inadvertently impacting classification accuracy owing to suboptimal enhancement techniques. Concurrently, another research trajectory has been dedicated to devising effective algorithms capable of uniformly segmenting a diverse array of video lengths. This approach, though promising, faces significant challenges in accurately depicting the spatio-temporal attributes inherent in video data.

Tackling the challenge of extracting viable information from dimly-lit, unclear videos demands the selection of suitable enhancement techniques, coupled with the strategic segmentation of extracted

frames. By preprocessing and partitioning extraneous noise, image representation becomes more succinct, thereby streamlining the feature extraction process (Mourya et al., 2020).

Driven by this objective, our study delves into an effective approach for human action recognition in low-light videos. We employ the MinMax-Sampling Singh et al. (2022) method for sampling and the Zero-DCE (Guo et al., 2020) for image enhancement. The enhanced videos are then processed through a 3D-CNN-based ResNet-18 (2+1)D network (Tran et al., 2018), which is pivotal for feature extraction. Subsequently, these features are classified using an Support Vector Machine (SVM), yielding insightful inference results. To further refine our model, we integrate all these methods into a seamless end-to-end system. Here, the training dataset undergoes implicit sampling and enhancement before being inputted into the ResNet-18 (2+1)D network, allowing for a comprehensive training and inference validation.

2 RELATED WORK

The recent landscape of human action recognition algorithms, as delineated in Kong & Fu (2022), categorizes them broadly into two distinct approaches: shallow techniques and deep architectures, as illustrated in Figure 1. The shallow approach encompasses methodologies like action representation, action classification, classifiers for human interaction, and classifiers tailored for RGB-D videos. These techniques primarily employ handcrafted feature engineering to capture local representations. On the other hand, deep architectures involve spatio-temporal networks, multi-stream networks, and hybrid networks, represented chiefly by 3D-CNNs and Two-stream methodologies.

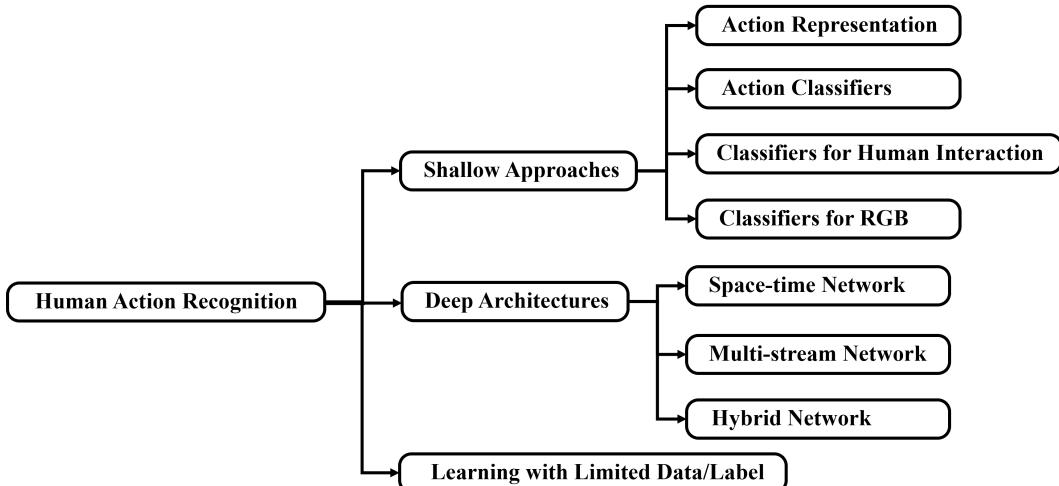


Figure 1: Classical Human Action Recognition Methods.

The mainstream approaches based on shallow architectures primarily consist of two types: holistic representation and local representation. In holistic representation, the method captures the overall motion information of HAR by capturing the global body structure and motion. Motion Energy Images (MEI) and Motion History Images (MHI), proposed by Bobick and Davis, exemplify this approach (Bobick & Davis, 2001). The main challenges of holistic representation are its sensitivity to noise, potential introduction of irrelevant information, the need for a static background, and difficulties in detecting detailed actions. On the other hand, local representation involves extracting local features from motion. Laptev's Spatio-Temporal Interest Points (STIP) method initially detects points of significant change, around which local descriptors are then constructed (Laptev, 2005). Despite its effective local feature extraction, this method sees limited practical application due to its complexity and computational intensity.

The field of human action recognition is advancing towards deeper architectures due to the powerful performance of deep learning, with several robust algorithms for action recognition already in existence. Simonyan and Zisserman (Simonyan & Zisserman, 2014) introduced a two-stream ConvNet architecture comprising spatial and temporal networks for action recognition. However, the

performance of this method was not entirely satisfactory, despite the authors further proposing the Two-Stream Inflated 3D ConvNet (I3D) (Carreira & Zisserman, 2017). To extract spatio-temporal features, attention shifted towards 3D structures, with the 3D-CNN model C3D pioneering this approach (Tran et al., 2015). This was followed by exciting developments like 3D-Resnet Hara et al. (2017) and 3D-Resnext. Moreover, the use of R(2 + 1)D (Tran et al., 2018) has been proposed for spatial and temporal convolution, enhancing the feature extraction for action recognition. The latest state-of-the-art (SOTA) methods revolve around the concept of Transformers, showcasing the prowess of self-attention techniques in this domain.

3 PROPOSED METHODOLOGY

The unique dark video human action recognition architecture proposed in this paper consists of four modules: **image frame sampling**, **image quality enhancement**, **feature extraction network** and **classification network**. The final composition of the end-to-end framework is shown in Figure 2. The framework takes the entire video as input and finally outputs the classification result.

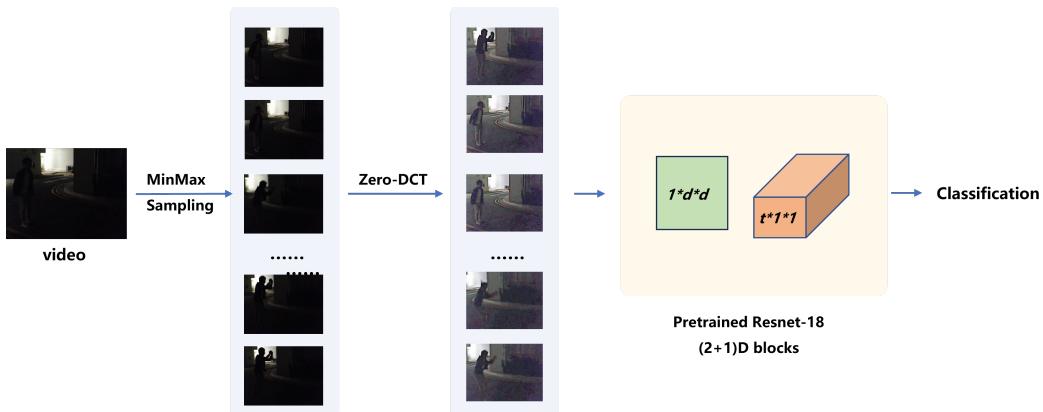


Figure 2: Framework of the proposed methodology.

3.1 IMAGE FRAMES SAMPLING

In human action recognition projects in low-light environments, choosing the correct video frame sampling method is crucial to ensure the effectiveness of model training. In this paper, three commonly used sampling methods are compared: max-min sampling, uniform sampling and random sampling, and the max-min sampling method is finally selected based on the sampling results.

3.1.1 UNIFORM SAMPLING

We first converted the raw input video into a sequence of images, and then applied uniform sampling to acquire one frame in every three image frames until 30 frames were collected (statistically, the number of frames in the training and validation videos varied between 55-100 frames). For less than 90 frames (30×3), a solid black image is used to fill the empty sample. This method is simple, straightforward, and ensures that frames are uniformly sampled from the overall length of the video to represent the action stream at constant intervals, capturing important spatio-temporal information. The frames obtained by uniform sampling are shown in Figure 3. **And for a more straightforward demonstration of the sampling result, the images are selected directly after enhancement.**

3.1.2 RANDOM SAMPLING

The random sampling method is to select frames randomly from the video. This method is simple, captures unexpected and important frames in the video, and is suitable for obtaining diversity and randomness in video content. However, it is very unfriendly to the continuity of actions and may ignore temporally important continuous actions, which is not taken in this paper.

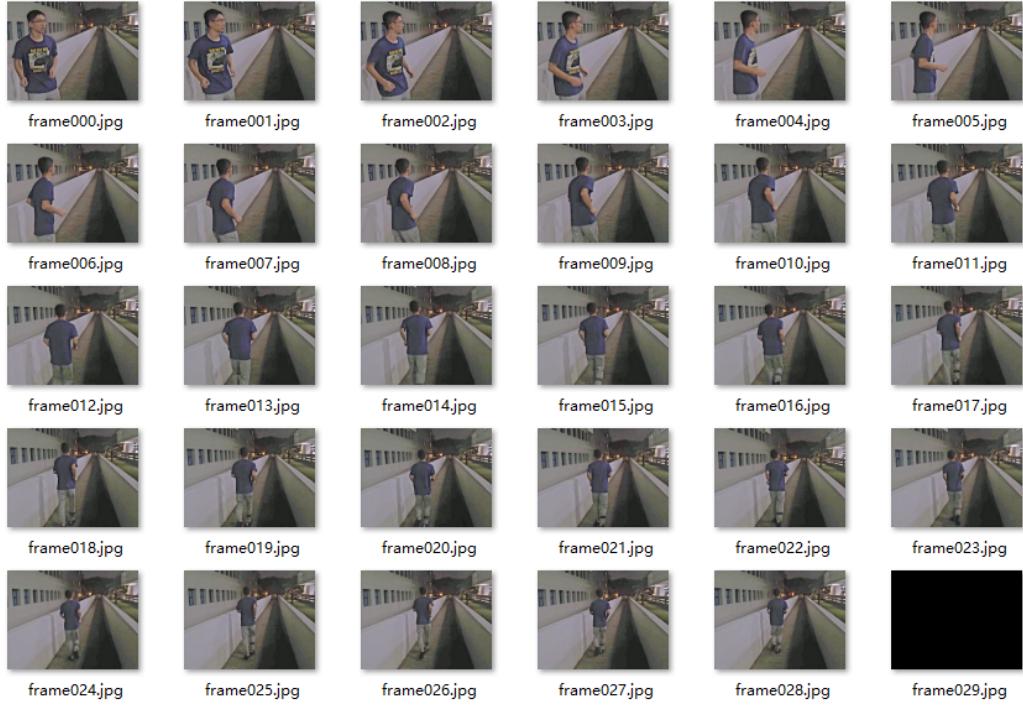


Figure 3: Uniformly sampled image frames from validation video 183.mp4. It can be seen that the last frame is supplemented by a blank frame as the total number of frames is less than 90.

3.1.3 MIN-MAX SAMPLING

Action videos are composed of action preparation, action start and action finish. Different videos have different lengths and distribution of action phases, and general sampling methods may lead to model overfitting. This project used the min-max sampling method to overcome this problem. Consider the input video $F = (F_1, \dots, F_L)$, F_L denotes the F^{th} frame, L denotes the total number of frames. F_r denotes the number of frames required, and the maximum step length is calculated first:

$$S_{\max} = \left\lfloor \frac{F}{F_r} \right\rfloor. \quad (1)$$

Thus the sampling step size S_{req} is as follows:

$$S_{\text{req}} = \text{Max} [\text{Min} [S_{\max}, 4], 1], \quad (2)$$

The principle shows that the max-min sampling method considers the overall length of the video while selecting the frames. **This method ensures that frames are extracted from different parts of the video without over-concentrating on specific parts of the video.** The frames sampled for the training set video *Jump_8_1.mp4* are shown in Figure 4.

In summary, max-min sampling is the most appropriate choice. Firstly, **it provides a more comprehensive coverage of the whole video**, including the onset, execution and end phases of the action, which is important to ensure that the model learns the different phases of the action when it is trained. Secondly, **max-min sampling omits the trial-and-error process of determining the sampling interval for uniform sampling and avoids uneven distribution of key actions**.

3.2 IMAGE QUALITY ENHANCEMENT

Since the overall video of the provided dataset has a mean value of approximately: [0.08149717 0.07725237 0.08080856] and a standard deviation of: [0.17944118 0.17953332 0.1681924], it is evident that the capture environment is very poor, so implementing image enhancement to improve the visual quality is a crucial step.

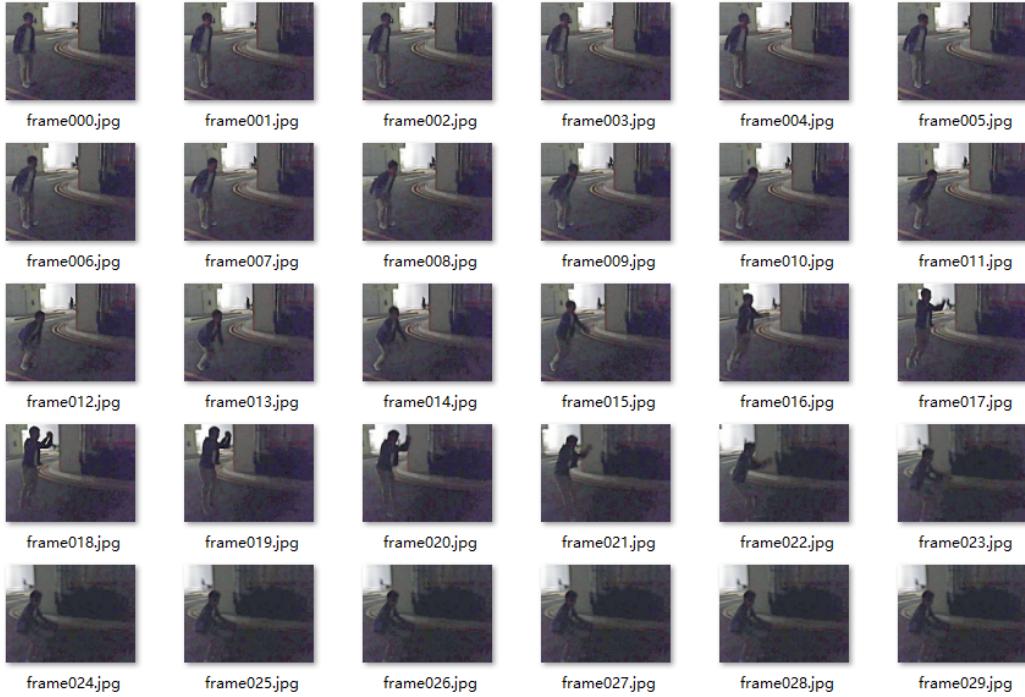


Figure 4: Min-max Sampling method sampled image frames from train video Jump_8_1.mp4.

3.2.1 CONTRAST LIMITED ADAPTIVE HISTOGRAM EQUALIZATION (CLAHE)

The histogram of an image can effectively represent its grayscale information, which is a comprehensive representation of brightness, contrast, color information, etc. Assuming that the input image is $f(x, y)$, the cumulative distribution function is as follows:

$$c(f) = \sum_{t=0}^f p_f(t) = \sum_{t=0}^f \frac{n_t}{n}, \quad f = 0, 1, \dots, L \quad (3)$$

Then after applying histogram equalization (HE) to transform the input image in gray level, the output image is as follows:

$$g = T(f) = \text{round} \left(\frac{c(f) - c_{\min}}{1 - c_{\min}} L \right), \quad c(f) \geq c_{\min} \quad (4)$$

When the process of intensity remapping is conducted across an image's various channels, leveraging its aggregate data, it can inadvertently amplify noise through excessive enhancement. Contrast Limited Adaptive Histogram Equalization (CLAHE) mitigates such noise amplification by partitioning the source image into discrete, non-overlapping contextual regions prior to applying HE. By computing individual histograms for each specified area, CLAHE imposes a cap on contrast enhancement, thereby curbing the prospective intensification of contrast to a permissible threshold. The image after enhancement using CLAHE is shown in Figure 5b.

3.2.2 ZERO-REFERENCE DEEP CURVE ESTIMATION (ZERO-DCE)

As a deep-structured image enhancement network paradigm, Zero-DCE can estimate the best-fit light enhancement curve based on the input image. It consists of three key components: the light enhancement curve, the DCE-Net, and the non-reference loss function. As a deep-structured image enhancement network paradigm, Zero-DCE can estimate the best-fit light enhancement curve based on the input image. It consists of three key components: the light enhancement curve, the DCE-Net, and the non-reference loss function.

- a) **Light-enhancement curve (LE-curve):** With the mapping technique, it automatically adjusts the curve parameters according to the input pixel values, while ensuring that the pixel values are regularized in the range of [0,1] to minimize information loss. What is more wonderful is that the monotonicity of the curve avoids the loss of details.
- b) **DCE-Net:** Deep curve estimation networks are used to learn the mapping relationship from the input image to the optimal curve parameters, outputting a set of pixel-level parameter mappings representing higher-order curves. The network is known for its light weight.
- c) **Non-reference loss functions:** The training of the model is performed by leveraging a linearly combined quintet of reference-based loss functions, specifically: spatial consistency loss (L_{spa}), exposure control loss (L_{exp}), color constancy loss (L_{col}), and illumination smoothness loss (L_{tvA}). This amalgamated loss function is crafted to maintain regional homogeneity and local luminance, while rectifying chromatic aberrations and ensuring the smooth transition of illumination across the image. The mathematical formulation is:

$$L_{\text{total}} = L_{spa} + L_{exp} + W_{col}L_{col} + W_{tvA}L_{tvA}, \quad (5)$$

where W_{col} and W_{tvA} are the weights of the losses.

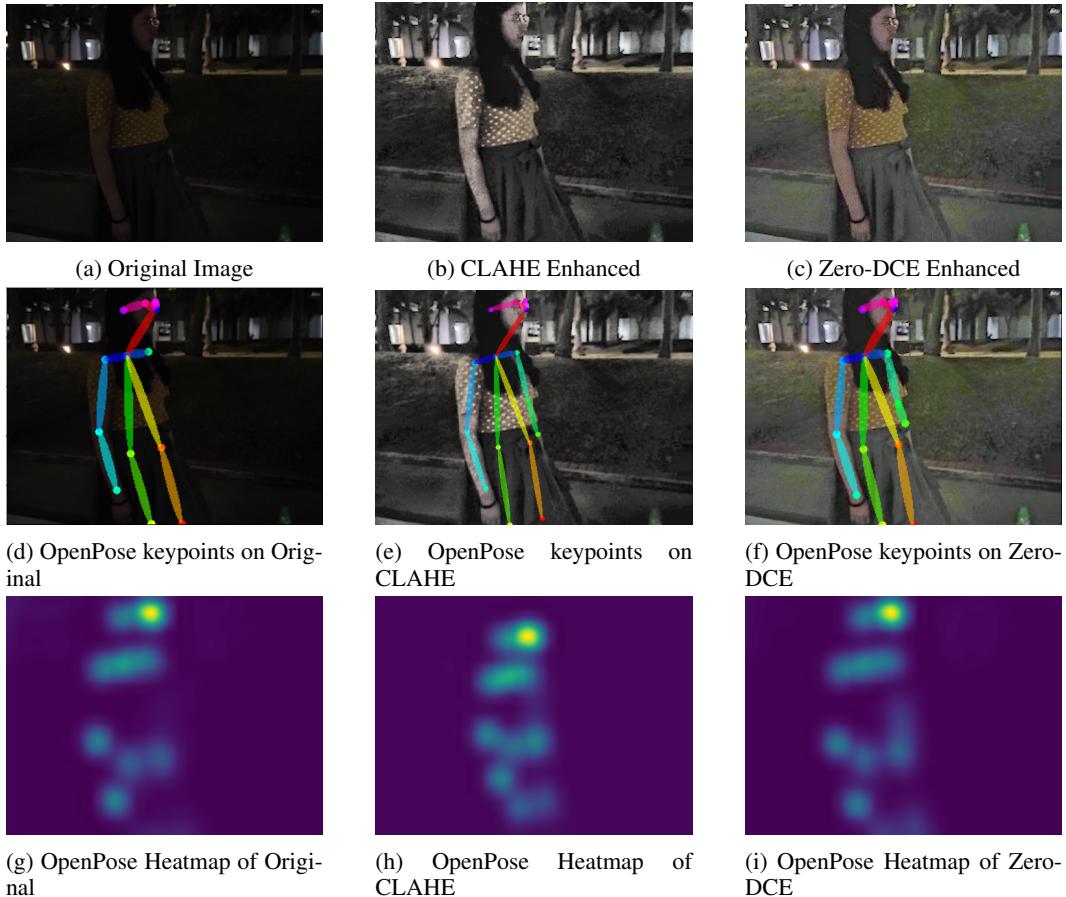


Figure 5: Comparison of Original, CLAHE, and Zero-DCE enhanced images from train video Stand_5_4's frame along with their corresponding Openpose keypoints and heatmaps.

In Figure 5c, the visual enhancement achieved by applying the Zero-DCE network to the original image is exhibited, **demonstrating a notable improvement in color saturation and brightness over the result augmented by CLAHE**. Furthermore, to assess the impact of varying image qualities on automated recognition systems, human skeletal keypoint maps and heatmaps are derived using the Openpose model¹, as depicted in Figure 5. **This comparative analysis highlights that image**

¹<https://github.com/Hzzone/pytorch-openpose>

enhancement substantially influences the accuracy and clarity of the traced human keypoints and associated heatmaps. In summary, the image enhancement technology will significantly improve the contrast of the image while suppressing noise, and at the same time has an excellent correction effect on brightness and color, resulting in a good reproduction of the overall details.

3.3 FEATURE EXTRACTION NETWORK

The *ResNet* model, renowned for its superior residual learning framework, has recently become the preferred paradigm for 2D image classification. This paper employs the *ResNet-18 (2 + 1)D* architecture, a *ResNet*-style architecture that disentangles 3D convolutions into discrete 2D spatial and 1D temporal convolutions to extract clip-level features from videos. Given an input video with L frames of images with 3 channels and spatial dimensions $h \times w$, in a 3D CNN, the tensor dimensionality is four-dimensional, while in the *R (2 + 1)D* model, the original N_i 3D filters of dimensions $N_{(i-1)} \times t \times d \times d$ are substituted with M_i 2D spatial convolutional filters of dimensions $N_{(i-1)} \times 1 \times d \times d$ and N_i temporal convolutional filters of dimensions $M_i \times t \times 1 \times 1$. **Although this decomposition increases the volume of non-linear computations, it notably diminishes the training error.**

Due to the small size of the available dataset, this project directly adopts the use of pre-trained weights from Kinetics 400 for feature extraction from the Tran et al. (2018), where the output of the fully connected layer represents the size of the extracted features according to the network structure. The set of frames extracted from each video is fed into the model and the dimension of the extracted features is 400.

3.4 CLASSIFICATION NETWORK

In order to extensively verify the effect of image sampling method and image enhancement on the accuracy of extracted features, three classification methods are selected in this paper to observe whether the trends are consistent. The three classification methods are: **support vector machine**, **Naive Bayes** and adding a fully connected layer at the end of the R(2+1)D model to **output the classification results directly**.

3.4.1 CLASSIFIERS

In the realm of machine learning, various classifiers exhibit distinctive strengths contingent upon the nature and intricacy of the data. In this study, we employed *Support Vector Machine* (SVM) and *Naive Bayes* classifiers, grounded in their proven efficacy across a multitude of pattern recognition tasks.

Support Vector Machine (SVM): SVM is renowned for its robustness and effectiveness in high-dimensional spaces, which is quintessential for handling the complex feature sets derived from video frames in human action recognition tasks. It operates on the principle of structural risk minimization, striving to find a hyperplane that categorically separates classes while maximizing the margin. This aspect makes it particularly adept at managing datasets with a pronounced distinction between classes, as is often the case in action recognition where specific movements are indicative of particular actions.

Naive Bayes: On the contrary, Naive Bayes classifiers offer a probabilistic perspective, leveraging the Bayes theorem with the "naive" assumption of independence between features. Its allure lies in its simplicity and speed, especially when dealing with large datasets. Moreover, it has demonstrated considerable resilience against irrelevant features and is capable of handling missing data, making it a suitable candidate for preliminary analyses or when computational resources are limited. Despite its simplicity, Naive Bayes can perform remarkably well and is often used as a baseline for classification problems.

3.4.2 END-TO-END CLASSIFICATION

In addressing the constraints posed by the dataset at hand, the conventional methodology of feature extraction followed by independent classification has demonstrated suboptimal performance. In pursuit of a more efficacious approach, this study introduces an alternative strategy. We augment the

Resnet-18 (2+1)D architecture with an appended linear fully-connected layer corresponding to the number of action categories, which is six in this case. This addition facilitates end-to-end training directly utilizing the training dataset while simultaneously outputting the classification results. Concurrently, the validation dataset is assimilated into the model, serving as a benchmark to ascertain the efficacy of each iteration of training. The culmination of this model’s architecture is illustratively presented in Figure 6, encapsulating the integrated approach adopted in this research.

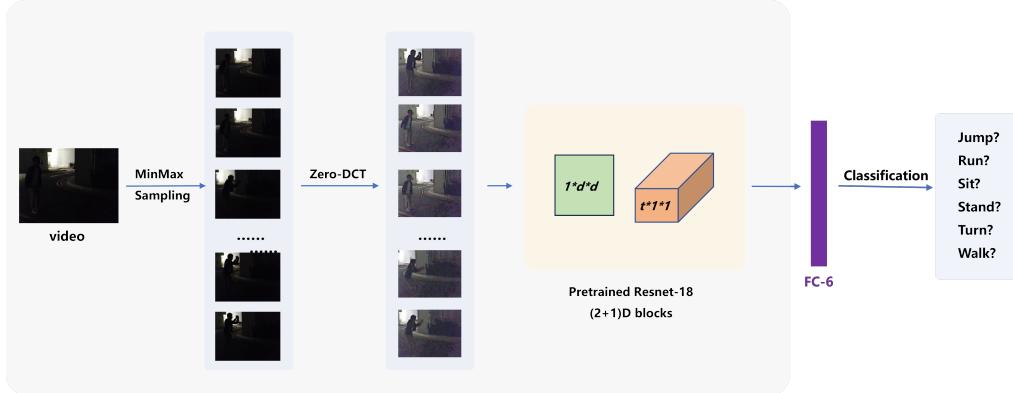


Figure 6: End-to-End Training Classification Architecture.

4 EXPERIMENTAL RESULT

4.1 EXPERIMENTAL DETAILS

All the proposed frameworks are implemented on the the Lenovo Legion y9000p with Intel Core i9-12900k and NVIDIA GEFORCE RTX3060. The training dataset used consists of six action categories, each containing 25 video data, and the corresponding validation dataset consists of 96 videos distributed within the corresponding categories. Experiments were completed using PyTorch. In this paper, we first compare the effects of different image enhancement methods on the classifier accuracy, and the results obtained are shown in Table 1.

Table 1: Comparison of Action Recognition Accuracy

Method	Accuracy	
	SVM	Naive Bayes
Max-min Sampling + R(2+1)D	0.218	0.177
Max-min sampling + CLAHE + R(2+1)D	0.234	0.202
Max-min sampling + Zero-DCE + R(2+1)D	0.243	0.229
Max-min sampling + CLAHE + heatmap + R(2+1)D	0.217	0.199
Max-min sampling + Zero-DCE + heatmap + R(2+1)D	0.245	0.219

The input to the training experiments for the end-to-end model consisted of a series of frames of size $6 \times 3 \times 30 \times 224 \times 224$, 30 frames of RGB three-channel images, with a batch size of 6. The cross-entropy loss function and the AdamW optimizer were used. The training experiments were conducted for 20 epochs and 40 epochs respectively. For the learning rate setting, a dynamic learning rate was adopted, where the learning rate is slowly increased at the beginning of the training to ensure that the error rate is reduced at the beginning of the training, and slowly reduced at about 7 epochs to ensure that the model accuracy is maintained. The training loss and verification loss during the 20 epochs training process are shown in Figure 7a, and the test accuracy is shown in Figure 7b. Meanwhile, Table 2 lists the accuracy of training 20 epochs and 40 epochs.

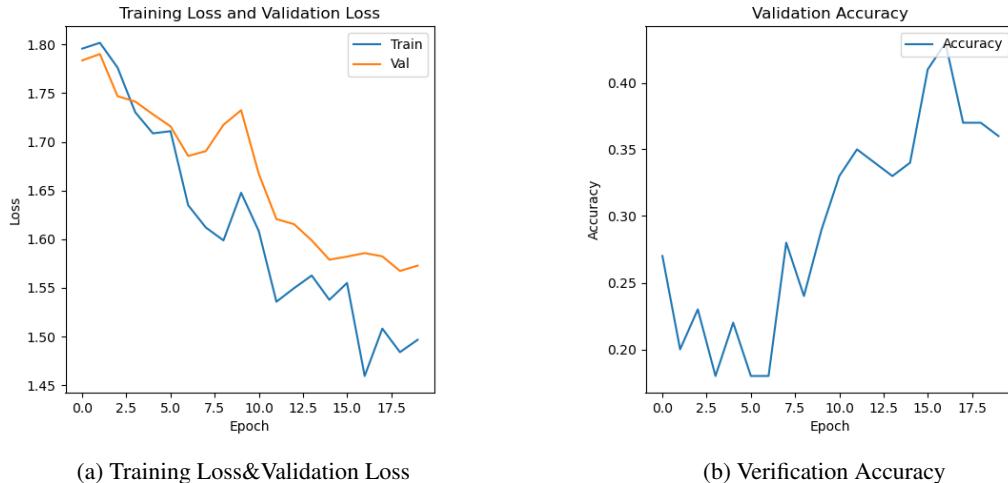


Figure 7: Training loss and verification accuracy for the model trained with 20 epochs. (a) Training Loss&Validation Loss. (b) Verification accuracy.

Table 2: End-to-End Training Accuracy Results

Method	Top-1 (Max)	Top-3
End-to-End Training with 20 epochs	44.38%	89.30%
End-to-End Training with 40 epochs	37.93%	87.52%

4.2 RESULT DISCUSSION

4.2.1 NON-END-TO-END MODEL

The experimental summary reveals several critical insights into the model’s performance. In the non-end-to-end models, the overall classification accuracy was notably low, with individual class accuracies struggling to reach 30%. The potential impacts of image enhancement, dataset quality, and the number of image frames on classification accuracy are discussed as follows:

1. Firstly, the improvement in feature extraction was marginal. Although most dark images were effectively enhanced visually, high-frequency noise persisting within the images, imperceptible to the human eye, may further impact feature extraction and, consequently, classifier performance.
2. Secondly, the dataset size, compared to public datasets with over 700 video materials (such as ARID), was insufficient in both quantity and diversity, leading to a propensity for overfitting and instability during training or testing.
3. Lastly, the number of image frames, as indicated in the appendix, one of the categories with lower accuracy scores, the “Walk” category, is often confused with the “Run” category due to their similarity. The spatiotemporal features of these two categories are very alike, making them challenging to distinguish in models lacking attention mechanisms. One solution is to increase the number of sampled frames per video to ensure the capture of every key action, which, while demanding more computational power, also imposes higher requirements on GPU memory.

4.2.2 END-TO-END MODEL

As depicted in Figure 7, the end-to-end training model exhibits a gradual decrease in training and validation losses, with a concurrent rise in validation accuracy, indicating that the model has not yet entered the overfitting phase. Table 2 further illustrates that both the validation and test sets

achieve a top-1 accuracy exceeding 35%, and the top-3 accuracy is approximately 90%. The superior performance of the end-to-end model over non-end-to-end models can be attributed to the enhanced accuracy of feature extraction achieved through the training of the feature extractor.

5 CONCLUSION

The field of Human Action Recognition (HAR) within video content is rapidly advancing, pushing the boundaries of machine learning and computer vision technologies into new territories such as surveillance and autonomous systems. This paper has taken a deep dive into the realm of HAR, specifically under the challenging conditions of low-light and nocturnal environments where conventional State of the Art (SOTA) methods tend to falter.

Our exploration began with an analytical comparison of various frame sampling techniques, including the traditional uniform sampling, the innovative MinMax-Sampling, and the stochastic random sampling. The objective was to discern the most effective strategy for capturing the quintessence of motion within the constrained luminance of dark video scenes.

In pursuit of refining the visual quality of these frames, we probed into the realms of image enhancement. Techniques such as Zero-Reference Deep Curve Estimation (Zero-DCE) and Contrast Limited Adaptive Histogram Equalization (CLAHE) were employed, each contributing significantly to the amplification of contrast and luminosity in the extracted frames. The result was a set of enriched images that served as a more robust foundation for feature extraction.

The crux of our feature extraction process was the ResNet-18 (2+1)D neural network. The network's adeptness at distilling pertinent features from the enhanced frames was paramount to the success of the subsequent classification phase. For the classification, we utilized a Support Vector Machine (SVM), a decision that was guided by the machine's proven proficiency in handling high-dimensional data. The accuracy of the SVM classifier was meticulously assessed, revealing a marked improvement over traditional methods when applied to low-light conditions.

Crucially, the integration of each discrete state-of-the-art methodology into a cohesive end-to-end model epitomized the innovative spirit of this research. By orchestrating the entire process from sampling to classification within a single streamlined model, we were able to achieve a synergy that not only optimized performance but also unveiled new avenues for HAR in less-than-ideal lighting scenarios.

The outcomes of this study underscore the potent potential of a holistic approach to HAR in dark videos. The top-1 and top-3 accuracies achieved in our experiments affirm the viability and superiority of the proposed methods. This research serves as a testament to the adaptability and resilience of machine learning techniques, even when confronted with the daunting task of deciphering human actions veiled in shadows.

ACKNOWLEDGMENTS

Thanks very much to Professor Xu Yuecong and Professor Jiang Xudong for giving us the treasure of knowledge. Through the exploration of experiments, I have grown a lot, and I hope to do more meaningful and contributing things in the future.

REFERENCES

- Irwan Bello, William Fedus, Xianzhi Du, Ekin Dogus Cubuk, Aravind Srinivas, Tsung-Yi Lin, Jonathon Shlens, and Barret Zoph. Revisiting resnets: Improved training and scaling strategies. *Advances in Neural Information Processing Systems*, 34:22614–22627, 2021.
- Aaron F. Bobick and James W. Davis. The recognition of human movement using temporal templates. *IEEE Transactions on pattern analysis and machine intelligence*, 23(3):257–267, 2001.
- Joao Carreira and Andrew 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, 2017.

- Chen Chen, Qifeng Chen, Minh N Do, and Vladlen Koltun. Seeing motion in the dark. In *Proceedings of the IEEE/CVF International conference on computer vision*, pp. 3185–3194, 2019.
- Deepti Ghadiyaram, Du Tran, and Dhruv Mahajan. Large-scale weakly-supervised pre-training for video action recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12046–12055, 2019.
- Chunle Guo, Chongyi Li, Jichang Guo, Chen Change Loy, Junhui Hou, Sam Kwong, and Runmin Cong. Zero-reference deep curve estimation for low-light image enhancement. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1780–1789, 2020.
- Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Learning spatio-temporal features with 3d residual networks for action recognition. In *Proceedings of the IEEE international conference on computer vision workshops*, pp. 3154–3160, 2017.
- Haiyang Jiang and Yinqiang Zheng. Learning to see moving objects in the dark. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7324–7333, 2019.
- Yu Kong and Yun Fu. Human action recognition and prediction: A survey. *International Journal of Computer Vision*, 130(5):1366–1401, 2022. doi: 10.1007/s11263-022-01594-9.
- Ivan Laptev. On space-time interest points. *International journal of computer vision*, 64:107–123, 2005.
- Hongying Meng, Nick Pears, and Chris Bailey. A human action recognition system for embedded computer vision application. In *2007 IEEE conference on computer vision and pattern recognition*, pp. 1–6. IEEE, 2007.
- Komal Mourya, Sharda Patil, Tabasum Nadaf, Divya Voccaligara, Harsha Chari, and Shailendra Aswale. Techniques for learning to see in the dark: A survey. In *2020 International Conference on Emerging Trends in Information Technology and Engineering (ic-ETITE)*, pp. 1–7. IEEE, 2020.
- Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. *Advances in neural information processing systems*, 27, 2014.
- Himanshu Singh, Saurabh Suman, Badri Narayan Subudhi, Vinit Jakhetiya, and Ashish Ghosh. Action recognition in dark videos using spatio-temporal features and bidirectional encoder representations from transformers. *IEEE Transactions on Artificial Intelligence*, pp. 1–11, 2022. doi: 10.1109/tai.2022.3221912.
- Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pp. 4489–4497, 2015.
- Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 6450–6459, 2018.
- Lun Zhang, Stan Z Li, Xiaotong Yuan, and Shiming Xiang. Real-time object classification in video surveillance based on appearance learning. In *2007 IEEE conference on computer vision and pattern recognition*, pp. 1–8. IEEE, 2007.

A APPENDIX

A.1 PYTHON CODE

```

1 import os
2 import cv2
3
4 def min_max_sampling(video_path, output_folder, required_frames):
5     video = cv2.VideoCapture(video_path)
6     total_frames = int(video.get(cv2.CAP_PROP_FRAME_COUNT))
7     S_max = total_frames // required_frames
8

```

```

9     if not os.path.exists(output_folder):
10        os.makedirs(output_folder)
11
12    for i in range(required_frames):
13        frame_id = i * S_max
14        video.set(cv2.CAP_PROP_POS_FRAMES, frame_id)
15        success, frame = video.read()
16        if success:
17            cv2.imwrite(os.path.join(output_folder, f'frame_{i:03d}.jpg'), frame)
18        else:
19            print(f"Failed to read frame {frame_id} from {video_path}")
20
21    video.release()
22
23 def process_videos(src_folder, dst_folder, required_frames):
24     for root, dirs, files in os.walk(src_folder):
25         for file in files:
26             if file.endswith('.mp4'):
27                 video_path = os.path.join(root, file)
28                 relative_path = os.path.relpath(root, src_folder)
29                 output_path = os.path.join(dst_folder, relative_path, os.path.splitext(file)[0])
30                 min_max_sampling(video_path, output_path, required_frames)
31
32 # Paths
33 train_video_folder = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition-\
34 In-The-Dark-master/EE6222_data/train'
35 val_video_folder = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition-In-\
36 The-Dark-master/EE6222_data/validate'
37 output_train_folder = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition-\
38 In-The-Dark-master/EE6222_data/train_img_2'
39 output_val_folder = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition-\
40 In-The-Dark-master/EE6222_data/validate_img_2'
41 required_frames = 30
42
43 # Process training and validation videos
44 process_videos(train_video_folder, output_train_folder, required_frames)
45 process_videos(val_video_folder, output_val_folder, required_frames)

```

```

1 import cv2
2 import numpy as np
3 import math
4 from scipy.ndimage import gaussian_filter
5 import matplotlib.pyplot as plt
6 import torch
7 import sys
8 sys.path.append('..') #
9 from openpose_pytorch.src import util
10 from openpose_pytorch.src.model import bodypose_model
11
12 class Body(object):
13     #
14     def __init__(self, model_path):
15         self.model = bodypose_model()
16         if torch.cuda.is_available():
17             self.model = self.model.cuda()
18         model_dict = util.transfer(self.model, torch.load(model_path))
19         self.model.load_state_dict(model_dict)
20         self.model.eval()
21
22     #
23     def __call__(self, oriImg):
24         # scale_search = [0.5, 1.0, 1.5, 2.0]
25         scale_search = [0.5]
26         boxsize = 368
27         stride = 8
28         padValue = 128
29         threl = 0.1
30         thre2 = 0.05
31         multiplier = [x * boxsize / oriImg.shape[0] for x in scale_search]
32         heatmap_avg = np.zeros((oriImg.shape[0], oriImg.shape[1], 19))
33         paf_avg = np.zeros((oriImg.shape[0], oriImg.shape[1], 38))
34
35         for m in range(len(multiplier)):
36             scale = multiplier[m]
37             imageToTest = cv2.resize(oriImg, (0, 0), fx=scale, fy=scale, interpolation=cv2.INTER_CUBIC)
38             imageToTest_padded, pad = util.padRightDownCorner(imageToTest, stride, padValue)
39             im = np.transpose(np.float32(imageToTest_padded[:, :, :, np.newaxis]), (3, 2, 0, 1)) / 256 - 0.5

```

```

40     im = np.ascontiguousarray(im)
41
42     data = torch.from_numpy(im).float()
43     if torch.cuda.is_available():
44         data = data.cuda()
45     # data = data.permute([2, 0, 1]).unsqueeze(0).float()
46     with torch.no_grad():
47         Mconv7_stage6_L1, Mconv7_stage6_L2 = self.model(data)
48         Mconv7_stage6_L1 = Mconv7_stage6_L1.cpu().numpy()
49         Mconv7_stage6_L2 = Mconv7_stage6_L2.cpu().numpy()
50
51     # extract outputs, resize, and remove padding
52     # heatmap = np.transpose(np.squeeze(net.blobs['output_blobs'].keys()[1].data), (1,
53     #     2, 0)) # output 1 is heatmaps
54     heatmap = np.transpose(np.squeeze(Mconv7_stage6_L2), (1, 2, 0)) # output 1 is
55     # heatmaps
56     heatmap = cv2.resize(heatmap, (0, 0), fx=stride, fy=stride, interpolation=cv2.
57     INTER_CUBIC)
58     heatmap = heatmap[:imageToTest_padded.shape[0] - pad[2], :imageToTest_padded.shape
59     [1] - pad[3], :]
60     heatmap = cv2.resize(heatmap, (oriImg.shape[1], oriImg.shape[0]), interpolation=
61     cv2.INTER_CUBIC)
62
63     # paf = np.transpose(np.squeeze(net.blobs['output_blobs'].keys()[0].data), (1, 2, 0)
64     # ) # output 0 is PAFs
65     paf = np.transpose(np.squeeze(Mconv7_stage6_L1), (1, 2, 0)) # output 0 is PAFs
66     paf = cv2.resize(paf, (0, 0), fx=stride, fy=stride, interpolation=cv2.INTER_CUBIC)
67     paf = paf[:imageToTest_padded.shape[0] - pad[2], :imageToTest_padded.shape[1] -
68     pad[3], :]
69     paf = cv2.resize(paf, (oriImg.shape[1], oriImg.shape[0]), interpolation=cv2.
70     INTER_CUBIC)
71
72     heatmap_avg += heatmap_avg + heatmap / len(multiplier)
73     paf_avg += paf / len(multiplier)
74
75     all_peaks = []
76     peak_counter = 0
77
78     for part in range(18):
79         map_ori = heatmap_avg[:, :, part]
80         one_heatmap = gaussian_filter(map_ori, sigma=3)
81
82         map_left = np.zeros(one_heatmap.shape)
83         map_left[1:, :] = one_heatmap[:-1, :]
84         map_right = np.zeros(one_heatmap.shape)
85         map_right[:-1, :] = one_heatmap[1:, :]
86         map_up = np.zeros(one_heatmap.shape)
87         map_up[:, 1:] = one_heatmap[:, :-1]
88         map_down = np.zeros(one_heatmap.shape)
89         map_down[:, :-1] = one_heatmap[:, 1:]
90
91         peaks_binary = np.logical_and.reduce(
92             (one_heatmap >= map_left, one_heatmap >= map_right, one_heatmap >= map_up,
93             one_heatmap >= map_down, one_heatmap > thre1))
94         peaks = list(zip(np.nonzero(peaks_binary)[1], np.nonzero(peaks_binary)[0])) # #
95         note reverse
96         peaks_with_score = [x + (map_ori[x[1], x[0]],) for x in peaks]
97         peak_id = range(peak_counter, peak_counter + len(peaks))
98         peaks_with_score_and_id = [peaks_with_score[i] + (peak_id[i],) for i in range(len(
99             peak_id))]
100
101         all_peaks.append(peaks_with_score_and_id)
102         peak_counter += len(peaks)
103
104     # find connection in the specified sequence, center 29 is in the position 15
105     limbSeq = [[2, 3], [2, 6], [3, 4], [4, 5], [6, 7], [7, 8], [2, 9], [9, 10], \
106     [10, 11], [2, 12], [12, 13], [13, 14], [2, 1], [1, 15], [15, 17], \
107     [1, 16], [16, 18], [3, 17], [6, 18]]
108     # the middle joints heatmap correpondence
109     mapIdx = [[31, 32], [39, 40], [33, 34], [35, 36], [41, 42], [43, 44], [19, 20], [21,
110     22], \
111     [23, 24], [25, 26], [27, 28], [29, 30], [47, 48], [49, 50], [53, 54], [51,
112     52], \
113     [55, 56], [37, 38], [45, 46]]
114
115     connection_all = []
116     special_k = []
117     mid_num = 10
118
119     for k in range(len(mapIdx)):
120         score_mid = paf_avg[:, :, [x - 19 for x in mapIdx[k]]]

```

```

108 candA = all_peaks[limbSeq[k][0] - 1]
109 candB = all_peaks[limbSeq[k][1] - 1]
110 nA = len(candA)
111 nB = len(candB)
112 indexA, indexB = limbSeq[k]
113 if (nA != 0 and nB != 0):
114     connection_candidate = []
115     for i in range(nA):
116         for j in range(nB):
117             vec = np.subtract(candB[j][:2], candA[i][:2])
118             norm = math.sqrt(vec[0] * vec[0] + vec[1] * vec[1])
119             norm = max(0.001, norm)
120             vec = np.divide(vec, norm)
121
122             startend = list(zip(np.linspace(candA[i][0], candB[j][0], num=mid_num),
123                                 np.linspace(candA[i][1], candB[j][1], num=mid_num)))
124
125             vec_x = np.array([score_mid[int(round(startend[I][1])), int(round(
126                                         startend[I][0])), 0] \
127                             for I in range(len(startend))])
128             vec_y = np.array([score_mid[int(round(startend[I][1])), int(round(
129                                         startend[I][0])), 1] \
130                             for I in range(len(startend))])
131
132             score_midpts = np.multiply(vec_x, vec[0]) + np.multiply(vec_y, vec[1])
133             score_with_dist_prior = sum(score_midpts) / len(score_midpts) + min(
134                 0.5 * oriImg.shape[0] / norm - 1, 0)
135             criterion1 = len(np.nonzero(score_midpts > thre2[0])) > 0.8 * len(
136                 score_midpts)
137             criterion2 = score_with_dist_prior > 0
138             if criterion1 and criterion2:
139                 connection_candidate.append(
140                     [i, j, score_with_dist_prior, score_with_dist_prior + candA[i]
141                         [2] + candB[j][2]])
142
143             connection_candidate = sorted(connection_candidate, key=lambda x: x[2],
144                                         reverse=True)
145             connection = np.zeros((0, 5))
146             for c in range(len(connection_candidate)):
147                 i, j, s = connection_candidate[c][0:3]
148                 if (i not in connection[:, 3] and j not in connection[:, 4]):
149                     connection = np.vstack([connection, [candA[i][3], candB[j][3], s, i, j
150 ]])
151                 if (len(connection) >= min(nA, nB)):
152                     break
153
154             connection_all.append(connection)
155         else:
156             special_k.append(k)
157             connection_all.append([])
158
159 # last number in each row is the total parts number of that person
160 # the second last number in each row is the score of the overall configuration
161 subset = -1 * np.ones((0, 20))
162 candidate = np.array([item for sublist in all_peaks for item in sublist])
163
164 for k in range(len(mapIdx)):
165     if k not in special_k:
166         partAs = connection_all[k][:, 0]
167         partBs = connection_all[k][:, 1]
168         indexA, indexB = np.array(limbSeq[k]) - 1
169
170         for i in range(len(connection_all[k])):
171             found = 0
172             subset_idx = [-1, -1]
173             for j in range(len(subset)):
174                 if subset[j][indexA] == partAs[i] or subset[j][indexB] == partBs[i]:
175                     subset_idx[found] = j
176                     found += 1
177
178                 if found == 1:
179                     j = subset_idx[0]
180                     if subset[j][indexB] != partBs[i]:
181                         subset[j][indexB] = partBs[i]
182                         subset[j][-1] += 1
183                         subset[j][-2] += candidate[partBs[i].astype(int), 2] +
184                                         connection_all[k][i][2]
185             elif found == 2: # if found 2 and disjoint, merge them
186                 j1, j2 = subset_idx

```

```

180     membership = ((subset[j1] >= 0).astype(int) + (subset[j2] >= 0).astype
181             (int))[:-2]
182     if len(np.nonzero(membership == 2)[0]) == 0: # merge
183         subset[j1][:-2] += (subset[j2][:-2] + 1)
184         subset[j1][-2:] += subset[j2][-2:]
185         subset[j1][-2] += connection_all[k][i][2]
186         subset = np.delete(subset, j2, 0)
187     else: # as like found == 1
188         subset[j1][indexB] = partBs[i]
189         subset[j1][-1] += 1
190         subset[j1][-2] += candidate[partBs[i].astype(int), 2] +
191                         connection_all[k][i][2]
192
193     # if find no partA in the subset, create a new subset
194     elif not found and k < 17:
195         row = -1 * np.ones(20)
196         row[indexA] = partAs[i]
197         row[indexB] = partBs[i]
198         row[-1] = 2
199         row[-2] = sum(candidate[connection_all[k][i], :2].astype(int), 2) +
200                     connection_all[k][i][2]
201         subset = np.vstack([subset, row])
202
203     # delete some rows of subset which has few parts occur
204     deleteIdx = []
205     for i in range(len(subset)):
206         if subset[i][-1] < 4 or subset[i][-2] / subset[i][-1] < 0.4:
207             deleteIdx.append(i)
208     subset = np.delete(subset, deleteIdx, axis=0)
209
210     # subset: n*20 array, 0-17 is the index in candidate, 18 is the total score, 19 is the
211     # total parts
212     # candidate: x, y, score, id
213     return candidate, subset, heatmap_avg
214
215 if __name__ == "__main__":
216     body_estimation = Body('..\\model\\body_pose_model.pth')
217     test_image = 'D:\\Human-Action-Recognition-In-The-Dark-master\\Human-Action-Recognition-In
218         -The-Dark-master\\Zero-DCE_code\\data\\result\\LIME\\Zero.jpg'
219     # test_image = 'D:\\Human-Action-Recognition-In-The-Dark-master\\Human-Action-Recognition-
220         -In-The-Dark-master\\EE6222_data\\train_img\\Stand\\Stand_5_4\\frame076.jpg' #
221
222     oriImg = cv2.imread(test_image) # B,G,R order
223     candidate, subset, heatmap_avg = body_estimation(oriImg)
224
225     #
226     canvas = util.draw_bodypose(oriImg, candidate, subset)
227     plt.imshow(canvas[:, :, [2, 1, 0]])
228     plt.title("Keypoints")
229     plt.show()
230
231     #
232     plt.imshow(heatmap_avg[:, :, :-1].sum(axis=2))
233     plt.title("Heatmaps")
234     plt.show()

```

```

1 from PIL import Image
2 import matplotlib.pyplot as plt
3 import cv2
4 import numpy as np
5
6 #
7 image_path = 'D:\\Human-Action-Recognition-In-The-Dark-master\\Human-Action-Recognition-In-The
8 -Dark-master\\EE6222_data\\enhanced_train_img_4\\Jump\\Jump_8_2\\frame000.jpg'
9 # "D:\\Human-Action-Recognition-In-The-Dark-master\\Human-Action-Recognition-In-The-Dark-master\\
10 EE6222_data\\enhanced_train_img_4\\Jump\\Jump_8_2\\frame000.jpg"
11 # image_path = 'E:\\OneDrive - Nanyang Technological University\\Desktop\\Figure_11.png'
12 original_image = cv2.imread(image_path)
13
14 #           BGR           LAB
15 lab_image = cv2.cvtColor(original_image, cv2.COLOR_BGR2LAB)
16
17 #           LAB           L, A,           B
18 l_channel, a_channel, b_channel = cv2.split(lab_image)
19
20 #           CLAHE
21 clahe = cv2.createCLAHE(clipLimit=2.0, tileGridSize=(8, 8))
22 l_channel = clahe.apply(l_channel)
23
24 #
25 updated_lab_image = cv2.merge((l_channel, a_channel, b_channel))

```

```

24
25     #           LAB           BGR
26     clahe_image = cv2.cvtColor(updated_lab_image, cv2.COLOR_LAB2BGR)
27     plt.imshow(cv2.cvtColor(clahe_image, cv2.COLOR_BGR2RGB))
28
29     #
30     plt.figure(figsize=(12, 6))
31     plt.subplot(1, 2, 1)
32     plt.imshow(cv2.cvtColor(original_image, cv2.COLOR_BGR2RGB))
33     plt.title('Original Image')
34     plt.axis('off')
35
36     plt.subplot(1, 2, 2)
37     plt.imshow(cv2.cvtColor(clahe_image, cv2.COLOR_BGR2RGB))
38     plt.title('CLAHE Image')
39     plt.axis('off')
40
41     plt.show()

```

```

1 import os
2 import pandas as pd
3 import pickle
4 import numpy as np
5 import joblib
6 from sklearn.svm import SVC
7 from sklearn.model_selection import train_test_split
8 from sklearn.naive_bayes import GaussianNB
9 from sklearn.metrics import classification_report, accuracy_score
10 from sklearn.metrics import top_k_accuracy_score
11 #
12 project_root = "D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition-In-The-
13 -Dark-master"
14
15 def load_features_and_labels(features_file, labels_file):
16     #
17     features = joblib.load(features_file)
18     #
19     labels_data = pd.read_csv(labels_file, sep='\t', header=None)
20     labels = labels_data[0].values
21     #
22     features_array = []
23     labels_array = []
24     for path, data in features.items():
25         frame_feature = data['features']
26         frame_label = data['label']
27         features_array.append(frame_feature)
28         labels_array.append(frame_label)
29     return np.array(features_array), np.array(labels_array)
30
31 train_features, train_labels = load_features_and_labels(
32     os.path.join(project_root, 'train_video_features.pkl'),
33     os.path.join(project_root, 'EE6222_data/train.txt')
34 )
35 test_features, test_labels = load_features_and_labels(
36     os.path.join(project_root, 'val_video_features.pkl'),
37     os.path.join(project_root, 'EE6222_data/validate.txt')
38 )
39
40 def reshape_features(features_array):
41     #
42     if features_array.ndim == 3:
43         n_samples = features_array.shape[0]
44         n_features = features_array.shape[2]
45         return features_array.reshape(n_samples, n_features)
46     return features_array
47
48 train_features = reshape_features(train_features)
49 test_features = reshape_features(test_features)
50
51 #           SVM
52 # Ensure your SVM classifier can return probabilities
53 svm_classifier = SVC(probability=True)
54 svm_classifier.fit(train_features, train_labels)
55
56 #
57 svm_probs = svm_classifier.predict_proba(test_features)
58 # The top-1 predictions are simply the class with the highest probability
59 svm_top1_predictions = np.argmax(svm_probs, axis=1)
60 # Calculate Top-5 accuracy using the probabilities

```

```

61 svm_top5_accuracy = top_k_accuracy_score(test_labels, svm_probs, k=3)
62
63 print("SVM Classifier Report:")
64 print(classification_report(test_labels, svm_top1_predictions))
65 print(f"SVM Top-1 Accuracy: {accuracy_score(test_labels, svm_top1_predictions)}")
66 print(f"SVM Top-3 Accuracy: {svm_top5_accuracy}\n")
67
68 # Naive Bayes naturally gives probabilities, so no need to change the classifier
69 # initialization
70 nb_classifier = GaussianNB()
71 nb_classifier.fit(train_features, train_labels)
72
73 # Getting the probabilities for each class
74 nb_probs = nb_classifier.predict_proba(test_features)
75 # The top-1 predictions are simply the class with the highest probability
76 nb_top1_predictions = np.argmax(nb_probs, axis=1)
77 # Calculate Top-5 accuracy using the probabilities
78 nb_top5_accuracy = top_k_accuracy_score(test_labels, nb_probs, k=3)
79
80 print("Naive Bayes Classifier Report:")
81 print(classification_report(test_labels, nb_top1_predictions))
82 print(f"Naive Bayes Top-1 Accuracy: {accuracy_score(test_labels, nb_top1_predictions)}")
83 print(f"Naive Bayes Top-3 Accuracy: {nb_top5_accuracy}")
84 # svm_predictions = svm_classifier.predict(test_features)
85 # print("SVM Classifier Report:")
86 # print(classification_report(test_labels, svm_predictions))
87 # print(f"SVM Accuracy: {accuracy_score(test_labels, svm_predictions)}\n")
88
89 # #
90 # nb_classifier = GaussianNB()
91 # nb_classifier.fit(train_features, train_labels)
92 # #
93 # nb_predictions = nb_classifier.predict(test_features)
94 # print("Naive Bayes Classifier Report:")
95 # print(classification_report(test_labels, nb_predictions))
96 # print(f"Naive Bayes Accuracy: {accuracy_score(test_labels, nb_predictions)}")

```

```

1 from pathlib import Path
2 import os
3
4
5 DATA_DIR = Path('EE6222_data')
6
7 MODEL_NAME = 'resnet18' # For CNN-RNN only
8 MODEL_TYPE = 'r2plus1d_18' # 'cnn-rnn', 'r3d_18', 'r2plus1d_18'
9 IMG_SHAPE = 'BCTHW' #'BCTHW', 'BTCHW'
10 #SAVED_WEIGHTS_DIR = '/kaggle/input/openpose-r21d-model-state-dict-221027/best_val_loss (1).pt'
11
12 PRETRAINED = Path('model/state_dict')
13 CHECKPOINT_DIR = PRETRAINED / 'checkpoint'
14 MODEL_STATE_DIR = PRETRAINED / 'bestloss'
15 OPENPOSE_STATE_DIR = 'model/body_pose_model.pth'
16
17 SUBMISSION_DIR = Path('model/Submission')
18
19 TRAIN_VID_FOLDER = 'train'
20 TRAIN_IMG_FOLDER = 'train_img'
21 TRAIN_VID_DIR = DATA_DIR / TRAIN_VID_FOLDER
22 TRAIN_IMG_DIR = DATA_DIR / TRAIN_IMG_FOLDER
23 TRAIN_CSV = DATA_DIR / 'train.txt'
24
25 VAL_VID_FOLDER = 'validate'
26 VAL_IMG_FOLDER = 'validate_img'
27 VAL_VID_DIR = DATA_DIR / VAL_VID_FOLDER
28 VAL_IMG_DIR = DATA_DIR / VAL_IMG_FOLDER
29 VAL_CSV = DATA_DIR / 'validate.txt'
30
31 INF_VID_FOLDER = 'validate'
32 INF_IMG_FOLDER = 'validate_img'
33 INF_VID_DIR = DATA_DIR / INF_VID_FOLDER
34 INF_IMG_DIR = DATA_DIR / INF_IMG_FOLDER
35 INF_CSV = DATA_DIR / 'validate.txt'
36
37 MAP_TABLE_DIR = DATA_DIR / 'mapping_table_23.txt'
38 NUM_CLASSES = 6
39
40 EXTENSION = ".mp4"
41 EXTENSION1 = ".jpg"
42 N_FRAMES = 12
43 INTERVAL = 5

```

```

43 SPLITS = 5
44
45
46 IMG_DIM = 224
47
48 IMAGENET_MEAN = [0.485, 0.456, 0.406] # RGB
49 IMAGENET_STD = [0.229, 0.224, 0.225] # RGB
50
51 DROPOUT = 0.3
52 RNN_HIDDEN_SIZE = 128
53 RNN_NUM_LAYERS = 1
54 ACCUM_ITER = 8
55
56 TRAIN_BATCH_SIZE = 6
57 VAL_BATCH_SIZE = 1 * TRAIN_BATCH_SIZE
58 TEST_BATCH_SIZE = VAL_BATCH_SIZE
59
60 LR = 1e-4 # 1e-3
61 NUM_EPOCHS = 40
62 NUM_WARMUP_EPOCHS = 5
63 NUM_COS_CYCLE = 0.5 # 0.5
64
65 TOP_K = 5

```

```

1 import os
2 import utils
3 from config.constants import (TRAIN_VID_FOLDER, TRAIN_IMG_FOLDER, TRAIN_VID_DIR,
4                               VAL_VID_FOLDER, VAL_IMG_FOLDER, VAL_VID_DIR,
5                               EXTENSION, N_FRAMES)
6
7 listOfCategories = os.listdir(TRAIN_VID_DIR)
8
9 def list_cat(train_path):
10     for cat in listOfCategories:
11         print("category:", cat)
12         path2acat = os.path.join(train_path, cat)
13         listOfSubs = os.listdir(path2acat)
14         print("number of sub-folders:", len(listOfSubs))
15         print("-" * 50)
16     return
17
18 def sv_frame(path, subfolder, subfolder_jpg):
19     for root, dirs, files in os.walk(path, topdown=False):
20         #print(f'root: {root}')
21         for name in files:
22             if EXTENSION not in name:
23                 continue
24             path2vid = os.path.join(root, name)
25             #print(path2vid)
26             frames, vlen = utils.get_frames(path2vid, n_frames=N_FRAMES)
27             path2store = path2vid.replace(subfolder, subfolder_jpg)
28             path2store = path2store.replace(EXTENSION, "")
29             print(path2store)
30             os.makedirs(path2store, exist_ok=True)
31             utils.store_frames(frames, path2store)
32             print("-" * 50)
33
34 #sv_frame(TRAIN_VID_DIR, TRAIN_VID_FOLDER, TRAIN_IMG_FOLDER)
35 #sv_frame(VAL_VID_DIR, VAL_VID_FOLDER, VAL_IMG_FOLDER)

```

```

1 import cv2
2 import numpy as np
3 import torch
4 import glob
5 import os
6 import albumentations as a
7 from pathlib import Path
8 from config.constants import N_FRAMES, INTERVAL, IMG_DIM, IMG_SHAPE, OPENPOSE_STATE_DIR
9 from torch.utils.data import Dataset
10 from .openpose_pytorch.src.body import Body_HM
11
12 # Dataset class with shape (T, C, H, W)
13
14 class VideoDataset(Dataset):
15     def __init__(self, df, transforms, img_path_col="path", labelAvailable=True):
16         self.df = df
17         self.transforms = transforms #
18         self.img_path_col = img_path_col #
19         self.labelAvailable = labelAvailable #

```

```

21     def __len__(self):
22         print("here 1") #
23         return len(self.df) #
24
25     def __getitem__(self, idx): #obtain single sample
26
27         # Video path
28         img_folder = self.df[self.img_path_col].iloc[idx] + "/*.jpg"
29         #print(img_folder)
30         # Append path of all frames in a video
31         all_img_path = glob.glob(img_folder)
32         all_img_path = sorted(all_img_path)
33         v_len = len(all_img_path)
34         #print(v_len)
35
36         # Uniformly samples N_FRAMES number of frames
37         if v_len > N_FRAMES*INTERVAL:
38             frame_list = np.arange(np.int64((v_len - (N_FRAMES*INTERVAL))*0.5), np.int64((v_len - (N_FRAMES*INTERVAL))+0.5) + N_FRAMES*INTERVAL, INTERVAL)
39         else:
40             frame_list = np.arange(0, N_FRAMES*INTERVAL, INTERVAL)
41         #print("here 2") p r o c e s s
42         img_path = []
43         for fn in range(v_len):
44             if (fn in frame_list):
45                 img_path.append(all_img_path[fn])
46         #print(img_path)
47         #print(img_path) output 12 path
48         images = []
49         for p2i in img_path:
50             p2i = Path(p2i)
51             img = cv2.imread(p2i.as_posix())
52             img = cv2.cvtColor(img, cv2.COLOR_BGR2RGB)
53             images.append(img)
54         #print("here 3") #after RGB
55         while (len(images) < N_FRAMES):
56             #if not enough ,the black screen is added for the rest of the frames
57             images.append(np.zeros((IMG_DIM,IMG_DIM,3), np.uint8))
58         images_tr = []
59
60         body_estimation_hm = Body_HM(OPENPOSE_STATE_DIR)
61
62         for image in images:
63             #print("here 4") # 12
64             if len(images_tr) == 0:
65                 augmented = self.transforms(image=image)
66                 image = augmented['image']
67                 body_heatmap = body_estimation_hm(image)
68                 hm = torch.Tensor((1 - body_heatmap[:, :, -1]))
69                 hm = torch.stack((hm,hm,hm), dim = 0)
70                 images_tr.append(torch.Tensor(hm))
71                 data_replay = augmented['replay']
72
73             else:
74                 image = a.ReplayCompose.replay(data_replay, image=image)
75                 image = image['image']
76                 body_heatmap = body_estimation_hm(image)
77                 hm = torch.Tensor((1 - body_heatmap[:, :, -1]))
78                 hm = torch.stack((hm,hm,hm), dim = 0)
79                 images_tr.append(torch.Tensor(hm))
80
81         if len(images_tr)>0:
82             images_tr = torch.stack(images_tr)
83         #print("here 5") # 12
84         if self.labelAvailable == True:
85             # Label
86             label = self.df["label"].iloc[idx]
87             return img_folder, images_tr, label
88         else:
89             return img_folder, images_tr
90
91     # Train image transformation function
92
93     def get_train_transforms(img_dim):
94         trans = a.ReplayCompose([
95             a.CLAHE(clip_limit=(10, 10), tile_grid_size=(8, 8), always_apply=True),
96             a.RandomGamma((150, 150), always_apply=True),
97             a.PadIfNeeded(img_dim, img_dim),

```

```

98         a.CenterCrop(img_dim, img_dim, always_apply=True),
99         a.HorizontalFlip(p=0.4),
100        a.VerticalFlip(p=0.5),
101        a.Rotate(limit=120, p=0.8),
102    ])
103    #print("transferred")
104    return trans
105
106
107 # Validation image transformation function
108 def get_val_transforms(img_dim):
109     trans = a.ReplayCompose([
110         a.CLAHE(clip_limit=(10, 10), tile_grid_size=(8, 8), always_apply=True),
111         a.RandomGamma((150, 150), always_apply=True),
112         a.PadIfNeeded(img_dim, img_dim),
113         a.CenterCrop(img_dim, img_dim, always_apply=True),
114     ])
115    #print("transformed-val")
116    return trans
117
118 #diy COLLATE function
119 def collate_fn(batch):
120     if len(batch[0]) == 3:
121         #print("collate-1")
122         img_folder_batch, imgs_batch, label_batch = list(zip(*batch))
123         label_batch = [torch.tensor(l) for l, imgs in zip(label_batch, imgs_batch) if len(imgs)
124                         >0]
125         labels_tensor = torch.stack(label_batch)
126     else:
127         #print("collate-2")
128         img_folder_batch, imgs_batch = list(zip(*batch))
129     print("collate-3")
130     img_folder_batch = [folders for folders in img_folder_batch if len(folders)>0]
131     imgs_batch = [imgs for imgs in imgs_batch if len(imgs)>0]
132
133     #imgs_tensor = torch.stack(torch.Tensor(imgs_batch))
134     if IMG_SHAPE == 'BCTHW':
135         #print("if")           BCTHW
136         imgs_batch = [torch.transpose(imgs, 1, 0) for imgs in imgs_batch if len(imgs)>0]
137     elif IMG_SHAPE == 'BTCHW':
138         imgs_batch = [imgs for imgs in imgs_batch if len(imgs)>0]
139     #print("collect-5")
140     imgs_tensor = torch.stack(imgs_batch)
141
142     if len(batch[0]) == 3:
143         #print("collate-6")
144
145         return img_folder_batch, imgs_tensor, labels_tensor
146     else:
147         return img_folder_batch, imgs_tensor

```

```

1 import os
2 import cv2
3 import numpy as np
4 from tqdm import tqdm
5 from scipy.ndimage import gaussian_filter
6 import matplotlib.pyplot as plt
7 import torch
8 import sys
9 import math
10 sys.path.append('..') #
11 from openpose_pytorch.src import util
12 from openpose_pytorch.src.model import bodypose_model
13
14 #          CLAHE
15 def apply_clahe(image):
16     #          BGRLAB
17     lab = cv2.cvtColor(image, cv2.COLOR_BGR2LAB)
18
19     #          LAB
20     lab_planes = list(cv2.split(lab)) #
21
22     #          CLAHE
23     clahe = cv2.createCLAHE(clipLimit=10.0, tileSize=(8, 8))
24     lab_planes[0] = clahe.apply(lab_planes[0]) #          LCLAHE
25
26     #          lab = cv2.merge(lab_planes)
27
28

```

```

29     #           LABBGR
30     enhanced_img = cv2.cvtColor(lab, cv2.COLOR_LAB2BGR)
31     return enhanced_img
32
33     #           OpenPose
34 class BodyPoseEstimator(object):
35     def __init__(self, model_path):
36         self.model = bodypose_model()
37         if torch.cuda.is_available():
38             self.model = self.model.cuda()
39         model_dict = util.transfer(self.model, torch.load(model_path))
40         self.model.load_state_dict(model_dict)
41         self.model.eval()
42
43     def __call__(self, oriImg):
44         # scale_search = [0.5, 1.0, 1.5, 2.0]
45         scale_search = [0.5]
46         boxsize = 368
47         stride = 8
48         padValue = 128
49         thre1 = 0.1
50         thre2 = 0.05
51         multiplier = [x * boxsize / oriImg.shape[0] for x in scale_search]
52         heatmap_avg = np.zeros((oriImg.shape[0], oriImg.shape[1], 19))
53         paf_avg = np.zeros((oriImg.shape[0], oriImg.shape[1], 38))
54
55         for m in range(len(multiplier)):
56             scale = multiplier[m]
57             imageToTest = cv2.resize(oriImg, (0, 0), fx=scale, fy=scale, interpolation=cv2.
58                                     INTER_CUBIC)
59             imageToTest_padded, pad = util.padRightDownCorner(imageToTest, stride, padValue)
60             im = np.transpose(np.float32(imageToTest_padded[:, :, :, np.newaxis]), (3, 2, 0,
61                                   1)) / 256 - 0.5
62             im = np.ascontiguousarray(im)
63
64             data = torch.from_numpy(im).float()
65             if torch.cuda.is_available():
66                 data = data.cuda()
67                 # data = data.permute([2, 0, 1]).unsqueeze(0).float()
68                 with torch.no_grad():
69                     Mconv7_stage6_L1, Mconv7_stage6_L2 = self.model(data)
70                     Mconv7_stage6_L1 = Mconv7_stage6_L1.cpu().numpy()
71                     Mconv7_stage6_L2 = Mconv7_stage6_L2.cpu().numpy()
72
73                 # extract outputs, resize, and remove padding
74                 # heatmap = np.transpose(np.squeeze(net.blobs['output_blobs'].keys()[1].data), (1,
75                 #                                2, 0)) # output 1 is heatmaps
76                 heatmap = np.transpose(np.squeeze(Mconv7_stage6_L2), (1, 2, 0)) # output 1 is
77                 # heatmaps
78                 heatmap = cv2.resize(heatmap, (0, 0), fx=stride, fy=stride, interpolation=cv2.
79                                     INTER_CUBIC)
80                 heatmap = heatmap[:imageToTest_padded.shape[0] - pad[2], :imageToTest_padded.shape
81                               [1] - pad[3], :]
82                 heatmap = cv2.resize(heatmap, (oriImg.shape[1], oriImg.shape[0]), interpolation=
83                                     cv2.INTER_CUBIC)
84
85                 # paf = np.transpose(np.squeeze(net.blobs['output_blobs'].keys()[0].data), (1, 2, 0)
86                 # ) # output 0 is PAFs
87                 paf = np.transpose(np.squeeze(Mconv7_stage6_L1), (1, 2, 0)) # output 0 is PAFs
88                 paf = cv2.resize(paf, (0, 0), fx=stride, fy=stride, interpolation=cv2.INTER_CUBIC)
89                 paf = paf[:imageToTest_padded.shape[0] - pad[2], :imageToTest_padded.shape[1] -
90                           pad[3], :]
91                 paf = cv2.resize(paf, (oriImg.shape[1], oriImg.shape[0]), interpolation=cv2.
92                                     INTER_CUBIC)
93
94                 heatmap_avg += heatmap_avg + heatmap / len(multiplier)
95                 paf_avg += paf / len(multiplier)
96
97                 all_peaks = []
98                 peak_counter = 0
99
100                for part in range(18):
101                    map_ori = heatmap_avg[:, :, part]
102                    one_heatmap = gaussian_filter(map_ori, sigma=3)
103
104                    map_left = np.zeros(one_heatmap.shape)
105                    map_left[1:, :] = one_heatmap[:-1, :]
106                    map_right = np.zeros(one_heatmap.shape)
107                    map_right[:-1, :] = one_heatmap[1:, :]
108                    map_up = np.zeros(one_heatmap.shape)
109                    map_up[:, 1:] = one_heatmap[:, :-1]

```

```

100     map_down = np.zeros(one_heatmap.shape)
101     map_down[:, :-1] = one_heatmap[:, 1:]
102
103     peaks_binary = np.logical_and.reduce(
104         (one_heatmap >= map_left, one_heatmap >= map_right, one_heatmap >= map_up,
105          one_heatmap >= map_down, one_heatmap > thre1))
106     peaks = list(zip(np.nonzero(peaks_binary)[1], np.nonzero(peaks_binary)[0])) # note reverse
107     peaks_with_score = [x + (map_ori[x[1], x[0]],) for x in peaks]
108     peak_id = range(peak_counter, peak_counter + len(peaks))
109     peaks_with_score_and_id = [peaks_with_score[i] + (peak_id[i],) for i in range(len(peak_id))]
110
111     all_peaks.append(peaks_with_score_and_id)
112     peak_counter += len(peaks)
113
114     # find connection in the specified sequence, center 29 is in the position 15
115     limbSeq = [[2, 3], [2, 6], [3, 4], [4, 5], [6, 7], [7, 8], [2, 9], [9, 10], \
116                 [10, 11], [2, 12], [12, 13], [13, 14], [2, 1], [1, 15], [15, 17], \
117                 [1, 16], [16, 18], [3, 17], [6, 18]]
118     # the middle joints heatmap correpondence
119     mapIdx = [[31, 32], [39, 40], [33, 34], [35, 36], [41, 42], [43, 44], [19, 20], [21,
120                 22], \
121                     [23, 24], [25, 26], [27, 28], [29, 30], [47, 48], [49, 50], [53, 54], [51,
122                         52], \
123                         [55, 56], [37, 38], [45, 46]]
124
125     connection_all = []
126     special_k = []
127     mid_num = 10
128
129     for k in range(len(mapIdx)):
130         score_mid = paf_avg[:, :, [x - 19 for x in mapIdx[k]]]
131         candA = all_peaks[limbSeq[k][0] - 1]
132         candB = all_peaks[limbSeq[k][1] - 1]
133         nA = len(candA)
134         nB = len(candB)
135         indexA, indexB = limbSeq[k]
136         if (nA != 0 and nB != 0):
137             connection_candidate = []
138             for i in range(nA):
139                 for j in range(nB):
140                     vec = np.subtract(candB[j][2:], candA[i][2:])
141                     norm = math.sqrt(vec[0] * vec[0] + vec[1] * vec[1])
142                     norm = max(0.001, norm)
143                     vec = np.divide(vec, norm)
144
145                     startend = list(zip(np.linspace(candA[i][0], candB[j][0], num=mid_num),
146                                         \
147                                         np.linspace(candA[i][1], candB[j][1], num=mid_num)))
148
149                     vec_x = np.array([score_mid[int(round(startend[I][1])), int(round(
150                         startend[I][0])), 0] \
151                             for I in range(len(startend))])
152                     vec_y = np.array([score_mid[int(round(startend[I][1])), int(round(
153                         startend[I][0])), 1] \
154                             for I in range(len(startend))])
155
156                     score_midpts = np.multiply(vec_x, vec[0]) + np.multiply(vec_y, vec[1])
157                     score_with_dist_prior = sum(score_midpts) / len(score_midpts) + min(
158                         0.5 * oriImg.shape[0] / norm - 1, 0)
159                     criterion1 = len(np.nonzero(score_midpts > thre2)[0]) > 0.8 * len(
160                         score_midpts)
161                     criterion2 = score_with_dist_prior > 0
162                     if criterion1 and criterion2:
163                         connection_candidate.append(
164                             [i, j, score_with_dist_prior, score_with_dist_prior + candA[i][2] + candB[j][2]])
165
166                     connection_candidate = sorted(connection_candidate, key=lambda x: x[2],
167                                         reverse=True)
168                     connection = np.zeros((0, 5))
169                     for c in range(len(connection_candidate)):
170                         i, j, s = connection_candidate[c][0:3]
171                         if (i not in connection[:, 3] and j not in connection[:, 4]):
172                             connection = np.vstack([connection, [candA[i][3], candB[j][3], s, i, j
173                                         ]])
174                         if (len(connection) >= min(nA, nB)):
175                             break
176
177

```

```

connection_all.append(connection)
else:
    special_k.append(k)
    connection_all.append([])

# last number in each row is the total parts number of that person
# the second last number in each row is the score of the overall configuration
subset = -1 * np.ones((0, 20))
candidate = np.array([item for sublist in all_peaks for item in sublist])

for k in range(len(mapIdx)):
    if k not in special_k:
        partAs = connection_all[k][:, 0]
        partBs = connection_all[k][:, 1]
        indexA, indexB = np.array(limbSeq[k]) - 1

        for i in range(len(connection_all[k])):
            found = 0
            subset_idx = [-1, -1]
            for j in range(len(subset)): # 1:size(subset,1):
                if subset[j][indexA] == partAs[i] or subset[j][indexB] == partBs[i]:
                    subset_idx[found] = j
                    found += 1

            if found == 1:
                j = subset_idx[0]
                if subset[j][indexB] != partBs[i]:
                    subset[j][indexB] = partBs[i]
                    subset[j][-1] += 1
                    subset[j][-2] += candidate[partBs[i].astype(int), 2] +
                        connection_all[k][i][2]
            elif found == 2: # if found 2 and disjoint, merge them
                j1, j2 = subset_idx
                membership = ((subset[j1] >= 0).astype(int) + (subset[j2] >= 0).astype(int))[:-2]
                if len(np.nonzero(membership == 2)[0]) == 0: # merge
                    subset[j1][-2:] += (subset[j2][-2:] + 1)
                    subset[j1][-2:] += subset[j2][-2:]
                    subset[j1][-2] += connection_all[k][i][2]
                    subset = np.delete(subset, j2, 0)
                else: # as like found == 1
                    subset[j1][indexB] = partBs[i]
                    subset[j1][-1] += 1
                    subset[j1][-2] += candidate[partBs[i].astype(int), 2] +
                        connection_all[k][i][2]

            # if find no partA in the subset, create a new subset
        elif not found and k < 17:
            row = -1 * np.ones(20)
            row[indexA] = partAs[i]
            row[indexB] = partBs[i]
            row[-1] = 2
            row[-2] = sum(candidate[connection_all[k][i, :2].astype(int), 2]) +
                connection_all[k][i][2]
            subset = np.vstack([subset, row])
    # delete some rows of subset which has few parts occur
deleteIdx = []
for i in range(len(subset)):
    if subset[i][-1] < 4 or subset[i][-2] / subset[i][-1] < 0.4:
        deleteIdx.append(i)
subset = np.delete(subset, deleteIdx, axis=0)

# subset: n*20 array, 0-17 is the index in candidate, 18 is the total score, 19 is the
# total parts
# candidate: x, y, score, id
return candidate, subset, heatmap_avg

# OpenPose
body_pose_estimator = BodyPoseEstimator('../model/body_pose_model.pth')

#
src_train_folder = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition-In-The-Dark-master/EE6222_data/enhanced_train_img_4'
src_val_folder = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition-In-The-Dark-master/EE6222_data/enhanced_validate_img_4'
dst_train_folder = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition-In-The-Dark-master/EE6222_data/train_img_6'
dst_val_folder = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition-In-The-Dark-master/EE6222_data/validate_img_6'

```

```

240 #
241 action_classes = os.listdir(src_train_folder)
242 for action_class in action_classes:
243     action_class_path = os.path.join(src_train_folder, action_class)
244     video_folders = os.listdir(action_class_path)
245     for video_folder in video_folders:
246         video_folder_path = os.path.join(action_class_path, video_folder)
247         frame_files = os.listdir(video_folder_path)
248         for frame_file in frame_files:
249             frame_path = os.path.join(video_folder_path, frame_file)
250             oriImg = cv2.imread(frame_path)
251             #apply CLAHE
252             #enhanced_img = apply_clahe(oriImg)
253             # assecc heatmap
254             _, heatmap_avg = body_pose_estimator(oriImg)
255
256             #
257             dst_video_folder = os.path.join(dst_train_folder, action_class, video_folder)
258             if not os.path.exists(dst_video_folder):
259                 os.makedirs(dst_video_folder)
260             #
261             # heatmap_img = np.uint8(255 * heatmap_avg / np.max(heatmap_avg)) #
262
263             heatmap_path = os.path.join(dst_video_folder, frame_file)
264             if len(heatmap_avg.shape) == 3 and heatmap_avg.shape[2] != 1:
265                 #
266                 heatmap_img = np.mean(heatmap_avg, axis=2)
267             else:
268                 heatmap_img = heatmap_avg
269             #
270             heatmap_img_normalized = cv2.normalize(heatmap_img, None, 0, 255, cv2.NORM_MINMAX)
271             #
272             heatmap_img_normalized = np.uint8(heatmap_img_normalized)
273             #
274             heatmap_img_color = cv2.applyColorMap(heatmap_img_normalized, cv2.COLORMAP_JET)
275             #
276             cv2.imwrite(heatmap_path, heatmap_img_color)
277             print("train data available")
278
279             video_folders = os.listdir(src_val_folder)
280             for video_folder in video_folders:
281                 video_folder_path = os.path.join(src_val_folder, video_folder)
282                 if os.path.isdir(video_folder_path):
283                     frame_files = os.listdir(video_folder_path)
284                     for frame_file in frame_files:
285                         frame_path = os.path.join(video_folder_path, frame_file)
286                         oriImg = cv2.imread(frame_path)
287
288                         # CLAHEOpenPose
289                         # enhanced_img = apply_clahe(oriImg)
290                         _, heatmap_avg = body_pose_estimator(oriImg)
291
292                         #
293                         dst_video_folder = os.path.join(dst_val_folder, video_folder)
294                         if not os.path.exists(dst_video_folder):
295                             os.makedirs(dst_video_folder)
296
297                         # heatmap_img = np.uint8(255 * heatmap_avg / np.max(heatmap_avg)) #
298
299                         heatmap_path = os.path.join(dst_video_folder, frame_file)
300                         if len(heatmap_avg.shape) == 3 and heatmap_avg.shape[2] != 1:
301                             #
302                             heatmap_img = np.mean(heatmap_avg, axis=2)
303                         else:
304                             heatmap_img = heatmap_avg
305                         #
306                         heatmap_img_normalized = cv2.normalize(heatmap_img, None, 0, 255, cv2.NORM_MINMAX)
307                         #
308                         heatmap_img_normalized = np.uint8(heatmap_img_normalized)
309                         #
310                         heatmap_img_color = cv2.applyColorMap(heatmap_img_normalized, cv2.COLORMAP_JET)
311                         #
312                         cv2.imwrite(heatmap_path, heatmap_img_color)
313             print("Congratulations All Done!!")

```

```

1 import torch
2 import torchvision.transforms as transforms
3 from torchvision.models.video import r3d_18, r2plus1d_18, R3D_18_Weights, R2Plus1D_18_Weights

```

```

4 import torch.nn as nn
5 import cv2
6 import os
7 import numpy as np
8 import pandas as pd
9 import pickle
10 from pathlib import Path
11 import os
12
13 # from config.constants import (TRAIN_VID_DIR, TRAIN_VID_FOLDER, TRAIN_IMG_FOLDER, VAL_VID_DIR,
14 #     , VAL_VID_FOLDER, VAL_IMG_FOLDER,
15 #         MAP_TABLE_DIR, TRAIN_CSV, TRAIN_IMG_DIR, EXTENSION, EXTENSION1,
16 #             VAL_CSV, VAL_IMG_DIR,
17 #                 IMG_DIM, TRAIN_BATCH_SIZE, VAL_BATCH_SIZE, MODEL_TYPE,
18 #                     MODEL_NAME, DROPOUT, RNN_HIDDEN_SIZE, RNN_NUM_LAYERS, LR,
19 #                         ACCUM_ITER,
20 #                             NUM_WARMUP_EPOCHS, NUM_EPOCHS, NUM_COS_CYCLE, MODEL_STATE_DIR,
21 #                                 TOP_K, NUM_CLASSES, CHECKPOINT_DIR, SUBMISSION_DIR)
22 #
23 model = r2plus1d_18(weights=R2Plus1D_18_Weights.DEFAULT)
24 model.eval() #
25
26 DATA_DIR = Path('EE6222_data')
27 MAP_TABLE_DIR = DATA_DIR / 'mapping_table_23.txt'
28
29 labels = pd.read_csv(MAP_TABLE_DIR, sep="\t", header=None, index_col=0)
30 label_dict = {}
31 for i in range(len(labels)):
32     label_dict[i] = labels.loc[i, 1]
33
34 label_list = labels.index.tolist()
35
36 transform = transforms.Compose([
37     transforms.ToTensor(),
38     transforms.Normalize(mean=[0.298, 0.276, 0.257], std=[0.209, 0.198, 0.180])
39 ])
40
41 def extract_features_from_frames(frames, size=(240, 320)):
42     #                                     4D
43     frames_resized = [cv2.resize(frame, size, interpolation=cv2.INTER_AREA) for frame in
44         frames]
45     frames_tensor = torch.stack([transform(cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)) for frame
46         in frames_resized])
47     #                                     (batch_size, channels, depth, height, width)
48     #                                     [depth, channels, height, width]
49     frames_tensor = frames_tensor.permute(1, 0, 2, 3).unsqueeze(0) #                                     [batch_size,
50     channels, depth, height, width]
51
52     with torch.no_grad():
53         features = model(frames_tensor) #                                     unsqueeze(0)
54     if features.dim() > 2:
55         pooled_features = torch.mean(features, dim=2)
56     else:
57         pooled_features = features
58     return pooled_features
59
60 def process_video_frames(video_frames_folder):
61     frames = []
62     for frame_name in sorted(os.listdir(video_frames_folder)):
63         frame_path = os.path.join(video_frames_folder, frame_name)
64         frame = cv2.imread(frame_path)
65         if frame is not None:
66             frames.append(frame)
67     return extract_features_from_frames(frames)
68
69 project_root = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition-In-The
70 -Dark-master/'
71 TRAIN_CSV = os.path.join(project_root, 'EE6222_data/train.txt')
72 VAL_CSV = os.path.join(project_root, 'EE6222_data/validate.txt')
73 TRAIN_IMG_DIR = os.path.join(project_root, 'EE6222_data/train_img_6')
74 VAL_IMG_DIR = os.path.join(project_root, 'EE6222_data/validate_img_6')
75 EXTENSION = '.mp4'
76
77 df_train = pd.read_csv(TRAIN_CSV, sep="\t", header=None, index_col=0)
78 df_train.columns = ['label', 'path']
79 df_train['path'] = str(TRAIN_IMG_DIR) + '/' + df_train['path'].str.replace(EXTENSION, "", regex
80 =True)

```

```

76 df_val = pd.read_csv(VAL_CSV, sep="\t", header=None, index_col=0)
77 df_val.columns = ['label', 'path']
78 df_val['path'] = str(VAL_IMG_DIR) + '/' + df_val['path'].str.replace(EXTENSION, "", regex=True)
79
80 #
81 #def extract_and_save_features(df, filename):
82 #    features_dict = {}
83 #    for _, row in df.iterrows():
84 #        video_folder_name = row['path'].replace(".mp4", "") # . m p 4
85 #        video_folder_path = os.path.join(TRAIN_IMG_DIR if 'train' in filename else VAL_IMG_DIR
86 #                                         , video_folder_name)
87 #        if not os.path.exists(video_folder_path):
88 #            print(f"          : {video_folder_path}")
89 #            continue
90 #        label = row['label']
91 #        features = process_video_frames(video_folder_path)
92 #        features_dict[video_folder_path] = {'features': features.cpu().numpy(), 'label': label}
93 #    print(f"Processed {video_folder_path}, Features Shape: {features.shape}")
94
95 with open(filename, 'wb') as f:
96     pickle.dump(features_dict, f)
97     print(f"          : {filename}")
98
99 #
100 extract_and_save_features(df_train, 'train_video_features2.pkl')
101 #
102 #
103 extract_and_save_features(df_val, 'val_video_features2.pkl')
104
105 print("          ")

```

```

1 import cv2
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import random
5
6 def extract_frames(video_path, sample_interval=5, num_frames=16, random_sampling=False):
7     cap = cv2.VideoCapture(video_path)
8     if not cap.isOpened():
9         print(f"Error opening video file: {video_path}")
10        return []
11
12     total_frames = int(cap.get(cv2.CAP_PROP_FRAME_COUNT))
13     if total_frames == 0:
14         print(f"No frames found in video: {video_path}")
15         return []
16
17     if random_sampling:
18         frame_indices = random.sample(range(total_frames), min(num_frames, total_frames))
19     else:
20         frame_indices = [i * sample_interval for i in range(num_frames)]
21
22     frames = []
23     for i in range(max(frame_indices) + 1):
24         ret, frame = cap.read()
25         if i in frame_indices and ret:
26             frames.append(cv2.cvtColor(frame, cv2.COLOR_BGR2RGB))
27         elif i in frame_indices:
28             frame_height = int(cap.get(cv2.CAP_PROP_FRAME_HEIGHT))
29             frame_width = int(cap.get(cv2.CAP_PROP_FRAME_WIDTH))
30             frames.append(np.zeros((frame_height, frame_width, 3), dtype=np.uint8))
31
32     cap.release()
33     return frames
34
35 def plot_frames(frames, title):
36     fig, axs = plt.subplots(4, 4, figsize=(10, 10))
37     axs = axs.flatten()
38     for i in range(16):
39         if i < len(frames) and frames[i].size != 0:
40             axs[i].imshow(frames[i])
41         else:
42             axs[i].imshow(np.zeros((100, 100, 3), dtype=np.uint8)) #
43
44         axs[i].axis('off')
45         axs[i].set_title(f'frame_{i+1}') #
46     plt.suptitle(title)
47     plt.tight_layout()

```

```

47     plt.show()
48
49 #
50 video_path = 'D:\\Human-Action-Recognition-In-The-Dark-master\\Human-Action-Recognition-In-The
-Dark-master\\EE6222_data\\train\\Jump\\Jump_8_5.mp4' #
51
52 #
53 uniform_frames = extract_frames(video_path, sample_interval=5, num_frames=16, random_sampling=
    False)
54 plot_frames(uniform_frames, "Uniform Sampling")
55
56 #
57 random_frames = extract_frames(video_path, num_frames=16, random_sampling=True)
58 plot_frames(random_frames, "Random Sampling")

```

```

1 import torch
2 import torchvision
3 from torchvision import transforms
4 from PIL import Image
5 import os
6 import numpy as np
7 import model
8 import glob
9 import time
10
11 def lowlight(image_path, model_path, result_path):
12     data_lowlight = Image.open(image_path)
13     data_lowlight = (np.asarray(data_lowlight)/255.0)
14     data_lowlight = torch.from_numpy(data_lowlight).float()
15     data_lowlight = data_lowlight.permute(2, 0, 1)
16     data_lowlight = data_lowlight.cuda().unsqueeze(0)
17
18     DCE_net = model.enhance_net_nopool().cuda()
19     DCE_net.load_state_dict(torch.load(model_path))
20     _, enhanced_image, _ = DCE_net(data_lowlight)
21
22     if not os.path.exists(os.path.dirname(result_path)):
23         os.makedirs(os.path.dirname(result_path))
24     torchvision.utils.save_image(enhanced_image, result_path)
25
26 def process_dataset(dataset_path, result_root, model_path):
27     os.environ['CUDA_VISIBLE_DEVICES']='0'
28     for root, dirs, files in os.walk(dataset_path):
29         for file in files:
30             if file.endswith('.jpg') or file.endswith('.png'):
31                 image_path = os.path.join(root, file)
32                 result_path = image_path.replace(dataset_path, result_root).replace('\\', '/')
33                 #print(f'Processing {image_path}')
34                 lowlight(image_path, model_path, result_path)
35
36 if __name__ == '__main__':
37     model_path = 'snapshots/Epoch99.pth'
38     train_data_path = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition
- In-The-Dark-master/EE6222_data/train_img_5'
39     val_data_path = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition-
In-The-Dark-master/EE6222_data/validate_img_5'
40     result_train_root = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-
Recognition-In-The-Dark-master/EE6222_data/enhanced_train_img_7'
41     result_val_root = 'D:/Human-Action-Recognition-In-The-Dark-master/Human-Action-Recognition
- In-The-Dark-master/EE6222_data/enhanced_validate_img_7'
42
43     with torch.no_grad():
44         process_dataset(train_data_path, result_train_root, model_path)
45         process_dataset(val_data_path, result_val_root, model_path)

```

```

1 import pandas as pd
2 import torch
3 from config.constants import (INF_CSV, INF_IMG_DIR, EXTENSION, IMG_DIM, TEST_BATCH_SIZE,
    NUM_CLASSES,
        MODEL_TYPE, MODEL_NAME, DROPOUT, RNN_HIDDEN_SIZE,
        RNN_NUM_LAYERS, MODEL_STATE_DIR, SUBMISSION_DIR,
        VAL_CSV, VAL_IMG_DIR, TOP_K, INF_VID_DIR, INF_VID_FOLDER,
        INF_IMG_FOLDER)
4
5 from model.data import get_val_transforms, collate_fn, VideoDataset
6 from model.model import HARModel
7 from model.inference_gndtruth import inference_loop
8 from torch.utils.data import DataLoader
9 from torchvision.models.video import r3d_18, r2plus1d_18, R3D_18_Weights, R2Plus1D_18_Weights
10

```

```

11 from sklearn.metrics import accuracy_score, top_k_accuracy_score, confusion_matrix
12 from data_prep import sv_frame
13 import torch.nn as nn
14
15 # sv_frame(INF_VID_DIR, INF_VID_FOLDER, INF_IMG_FOLDER)
16
17 df_test = pd.read_csv(VAL_CSV, sep="\t", header=None, index_col=0)
18 df_test.columns = ['label', 'path']
19 df_test['path'] = str(INF_IMG_DIR) + '/' + df_test['path'].str.replace(EXTENSION, "", regex=True)
20 print(df_test.head())
21
22 #df_test = pd.read_csv(INF_CSV, sep="\t", header=None, index_col=0)
23 #df_test.columns = ['path']
24 #df_test['path'] = str(INF_IMG_DIR) + '/' + df_test['path'].str.replace(EXTENSION, "")
25
26 inf_transforms = get_val_transforms(IMG_DIM)
27 print("here-1,inf")
28 inf_data = VideoDataset(df=df_test,
29                         transforms=inf_transforms,
30                         labelAvailable=True)
31 print("here-2,inf")
32 inf_loader = DataLoader(inf_data,
33                         batch_size=TEST_BATCH_SIZE,
34                         shuffle=False,
35                         collate_fn=collate_fn)
36 print("here-3,inf")
37 device = 'cuda' if torch.cuda.is_available() else 'cpu'
38
39 if MODEL_TYPE == 'cnn-rnn':
40     model = HARModel(
41         model_name=MODEL_NAME,
42         dropout=DROPOUT,
43         rnn_hidden_size=RNN_HIDDEN_SIZE,
44         rnn_num_layers=RNN_NUM_LAYERS,
45         num_classes=NUM_CLASSES,
46         pretrained=True)
47 elif MODEL_TYPE == 'r3d_18':
48     model = r3d_18(weights=None, progress=False)
49     num_features = model.fc.in_features
50     model.fc = nn.Linear(num_features, NUM_CLASSES)
51 elif MODEL_TYPE == 'r2plus1d_18':
52     model = r2plus1d_18(weights=None, progress=False)
53     num_features = model.fc.in_features
54     model.fc = nn.Linear(num_features, NUM_CLASSES)
55 print(MODEL_STATE_DIR,'MODEL_STATE_DIR')
56 model_paths = sorted(list(MODEL_STATE_DIR.glob('best*')))
57 print(model_paths,'model_paths')
58 overall_results = dict()
59 print(overall_results,'overall_results')
60 for m in range(len(model_paths)):
61
62     model_path = model_paths[m]
63
64     overall_results[m] = inference_loop(model, model_path, inf_loader, device)
65
66
67 for i, model_result in overall_results.items():
68     model_folder, model_logits, model_pred, model_target, model_loss = model_result
69     model_score = accuracy_score(model_target, model_pred)
70     model_con_mat = confusion_matrix(model_target, model_pred, labels=[0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
71     model_top_k_score = top_k_accuracy_score(model_target, model_logits, k=TOP_K, labels=[0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
72     print(f'Model accuracy score: {model_score}')
73     print(f'Model Top-{TOP_K} score: {model_top_k_score}')
74     print(f'Model loss: {model_loss}')
75     print(f'Confusion matrix: \n{model_con_mat}')
76     predictions = pd.Series(model_pred, name='prediction')
77     predictions.to_csv(SUBMISSION_DIR / 'vr.txt', sep='\t', header=False)

```

```

1 import pandas as pd
2 import numpy as np
3 from torch.utils.data import DataLoader
4 from config.constants import (TRAIN_VID_DIR, TRAIN_VID_FOLDER, TRAIN_IMG_FOLDER, VAL_VID_DIR,
5                               VAL_VID_FOLDER, VAL_IMG_FOLDER,
6                               MAP_TABLE_DIR, TRAIN_CSV, TRAIN_IMG_DIR, EXTENSION, VAL_CSV,
6                               VAL_IMG_DIR,
6                               IMG_DIM, TRAIN_BATCH_SIZE, VAL_BATCH_SIZE, MODEL_TYPE,

```

```

7             MODEL_NAME, DROPOUT, RNN_HIDDEN_SIZE, RNN_NUM_LAYERS, LR,
8             ACCUM_ITER,
9             NUM_WARMUP_EPOCHS, NUM_EPOCHS, NUM_COS_CYCLE, MODEL_STATE_DIR,
10            TOP_K, NUM_CLASSES, CHECKPOINT_DIR, SUBMISSION_DIR)
11
12 from model.data import get_train_transforms, get_val_transforms, collate_fn, VideoDataset
13 from transformers import get_cosine_schedule_with_warmup
14 from model.model import HARModel
15 from model.train import train_loop, val_loop
16 from data_prep import sv_frame
17 from torchvision.models.video import r3d_18, r2plus1d_18, R3D_18_Weights, R2Plus1D_18_Weights
18 import torch.nn as nn
19 import torch
20 import matplotlib.pyplot as plt
21
22 # sv_frame(TRAIN_VID_DIR, TRAIN_VID_FOLDER, TRAIN_IMG_FOLDER)
23 # sv_frame(VAL_VID_DIR, VAL_VID_FOLDER, VAL_IMG_FOLDER)
24
25 # Label dictionary
26 # {0: 'Drink', 1: 'Jump', 2: 'Pick', 3: 'Pour', 4: 'Push', 5: 'Run', 6: 'Sit', 7: 'Stand', 8:
27 #   'Turn', 9: 'Walk'}
28 labels = pd.read_csv(MAP_TABLE_DIR, sep="\t", header=None, index_col=0)
29 label_dict = {}
30 for i in range(len(labels)):
31     label_dict[i] = labels.loc[i, 1]
32
33 label_list = labels.index.tolist()
34
35 # Load training data
36 df_train = pd.read_csv(TRAIN_CSV, sep="\t", header=None, index_col=0)
37 df_train.columns = ['label', 'path']
38 df_train['path'] = str(TRAIN_IMG_DIR) + '/' + df_train['path'].str.replace(EXTENSION, "", regex=True)
39
40 df_val = pd.read_csv(VAL_CSV, sep="\t", header=None, index_col=0)
41 df_val.columns = ['label', 'path']
42 df_val['path'] = str(VAL_IMG_DIR) + '/' + df_val['path'].str.replace(EXTENSION, "", regex=True)
43
44 # for testing, if training ,remove this two rows
45 # df_train = df_train.iloc[:10]
46 # df_val = df_val.iloc[:10]
47
48 train_transforms = get_train_transforms(IMG_DIM)
49 val_transforms = get_val_transforms(IMG_DIM)
50
51 train_data = VideoDataset(df=df_train,
52                           transforms=train_transforms,
53                           labelAvailable=True)
54
55 val_data = VideoDataset(df=df_val,
56                         transforms=val_transforms,
57                         labelAvailable=True)
58
59 train_loader = DataLoader(train_data, batch_size= TRAIN_BATCH_SIZE,
60                           shuffle=True, collate_fn= collate_fn)
61 val_loader = DataLoader(val_data, batch_size= VAL_BATCH_SIZE,
62                         shuffle=False, collate_fn= collate_fn)
63
64 fold_results = dict()
65 device = 'cuda' if torch.cuda.is_available() else 'cpu'
66
67 if MODEL_TYPE == 'cnn-rnn':
68     model = HARModel(
69         model_name=MODEL_NAME,
70         dropout=DROPOUT,
71         rnn_hidden_size=RNN_HIDDEN_SIZE,
72         rnn_num_layers=RNN_NUM_LAYERS,
73         num_classes=NUM_CLASSES,
74         pretrained=True)
75 elif MODEL_TYPE == 'r3d_18':
76     model = r3d_18(weights=R3D_18_Weights.DEFAULT, progress=False)
77     num_features = model.fc.in_features
78     model.fc = nn.Linear(num_features, NUM_CLASSES)
79 elif MODEL_TYPE == 'r2plus1d_18':
80     model = r2plus1d_18(weights=R2Plus1D_18_Weights.DEFAULT, progress=False)
81     num_features = model.fc.in_features
82     model.fc = nn.Linear(num_features, NUM_CLASSES) #
83
84 model = model.to(device)
85 print("CUDA")
86 optimizer = torch.optim.AdamW(model.parameters(), lr=LR)
87 # Configs for scheduler

```

```

84 #print(len(train_loader), 'len(train_loader)')
85 #print(ACCUM_ITER, 'ACCUM_ITER') 8
86 num_steps_per_epoch = np.round_(len(train_loader)/ACCUM_ITER)
87 #print(num_steps_per_epoch, 'num_steps_per_epoch') 5
88 num_warmup_steps = NUM_WARMUP_EPOCHS * num_steps_per_epoch
89 #print(num_warmup_steps, 'num_warmup_steps') 5
90 num_training_steps = NUM_EPOCHS * num_steps_per_epoch
91 #print(num_training_steps, 'num_training_steps') 10
92
93 scheduler = get_cosine_schedule_with_warmup(
94     optimizer=optimizer,
95     num_warmup_steps=num_warmup_steps,
96     num_training_steps=num_training_steps,
97     num_cycles=NUM_COS_CYCLE
98 )
99
100 best_accuracy = 0
101 best_val_loss = 999_999
102 fold_train_loss = []
103 fold_val_loss = []
104 fold_accuracy = []
105 fold_k_score = []
106
107 for epoch_num in range(NUM_EPOCHS):
108     lr_list = []
109     epoch_info = f'Epoch {epoch_num+1}/{NUM_EPOCHS}'
110     print(epoch_info, 'epoch_info')
111     lr_list, train_losses = train_loop(
112         model, train_loader, device, optimizer, ACCUM_ITER,
113         scheduler=scheduler, epoch_info=epoch_info
114     )
115
116     val_losses, score, top_k_score = val_loop(
117         model, val_loader, device, label_list=label_list, k=TOP_K, epoch_info=epoch_info
118     )
119
120     avg_train_loss = np.mean(train_losses)
121     avg_val_loss = np.mean(val_losses)
122     avg_val_accuracy_score = np.mean(score)
123     avg_val_top_k_accuracy = np.mean(top_k_score)
124
125     fold_train_loss.append(avg_train_loss)
126     fold_val_loss.append(avg_val_loss)
127     fold_accuracy.append(avg_val_accuracy_score)
128     fold_k_score.append(avg_val_top_k_accuracy)
129
130     # Save model if val loss improves
131     if avg_val_loss < best_val_loss:
132         model_state_dict = model.state_dict()
133
134         model_name = f'best_val_loss.pt'
135         torch.save(
136             model.state_dict(),
137             MODEL_STATE_DIR / model_name
138         )
139
140         best_val_loss = avg_val_loss
141
142     # Save checkpoint at the last
143     if epoch_num == NUM_EPOCHS - 1:
144         model_name = f'model_checkpoint.pt'
145         torch.save({'epoch': epoch_num,
146                     'model_state_dict': model.state_dict(),
147                     'optimizer_state_dict': optimizer.state_dict()},
148                     CHECKPOINT_DIR / model_name)
149
150 fold_results[0] = (fold_train_loss, fold_val_loss, fold_accuracy, fold_k_score)
151
152 (train_loss, val_loss, accuracy, k_score) = fold_results[0]
153
154 plt.figure(figsize=(5, 5))
155 plt.plot(train_loss, label='Train')
156 plt.plot(val_loss, label='Val')
157 plt.legend(loc='upper right')
158 plt.xlabel('Epoch')
159 plt.ylabel('Loss')
160 plt.title(f'Training Loss and Validation Loss')
161 plt.show()
162
163 plt.figure(figsize=(5, 5))
164 plt.plot(fold_accuracy, label='Accuracy')

```

```

165 plt.legend(loc='upper right')
166 plt.xlabel('Epoch')
167 plt.ylabel('Accuracy')
168 plt.title(f'Validation Accuracy')
169 plt.show()

```

A.2 RESULT

Table 3: Classification Performance

Classifier	Precision	Recall	F1-Score	Support	Accuracy
SVM Classifier					
0	0.21	0.35	0.27	17	
1	0.19	0.20	0.19	15	
2	0.12	0.07	0.09	15	
3	0.20	0.25	0.22	16	
4	0.44	0.24	0.31	17	
5	0.13	0.12	0.13	16	
Average	0.22	0.20	0.20	96	0.2083
Top-5 Accuracy					0.90625
Naive Bayes Classifier					
0	0.18	0.17	0.18	17	
1	0.07	0.07	0.07	15	
2	0.19	0.20	0.19	15	
3	0.20	0.44	0.27	16	
4	0.57	0.24	0.33	17	
5	0.10	0.06	0.08	16	
Average	0.22	0.20	0.19	96	0.1979
Top-5 Accuracy					0.875

Table 4: Classification Performance

Classifier	Precision	Recall	F1-Score	Support	Accuracy
SVM Classifier					
0	0.17	0.12	0.14	17	
1	0.14	0.27	0.19	15	
2	0.7	0.07	0.07	15	
3	0.14	0.38	0.21	16	
4	0.38	0.29	0.33	17	
5	0.03	0.02	0.01	16	
Average	0.14	0.18	0.15	96	0.17708
Top-3 Accuracy					0.6046
Naive Bayes Classifier					
0	0.07	0.07	0.07	17	
1	0.17	0.17	0.18	15	
3	0.14	0.44	0.23	16	
2	0.25	0.33	0.29	15	
4	0.47	0.41	0.44	16	
5	0.01	0.03	0.02	17	
Average	0.18	0.22	0.19	96	0.2179
Top-3 Accuracy					0.4479