

Energy-based Periodicity Mining with Deep Features for Action Repetition Counting in Unconstrained Videos

Jianqin Yin Yanchun Wu Huaping Liu Yonghao Dang Zhiyi Liu and Jun Liu

Abstract—Action repetition counting is to estimate the occurrence times of the repetitive motion in one action, which is a relatively new, important but challenging measurement problem. To solve this problem, we propose a new method superior to the traditional ways in two aspects, without preprocessing and applicable for arbitrary periodicity actions. Without preprocessing, the proposed model makes our method convenient for real applications; processing the arbitrary periodicity action makes our model more suitable for the actual circumstance. In terms of methodology, firstly, we analyze the movement patterns of the repetitive actions based on the spatial and temporal features of actions extracted by deep ConvNets; Secondly, the Principal Component Analysis algorithm is used to generate the intuitive periodic information from the chaotic high-dimensional deep features; Thirdly, the periodicity is mined based on the high-energy rule using Fourier transform; Finally, the inverse Fourier transform with a multi-stage threshold filter is proposed to improve the quality of the mined periodicity, and peak detection is introduced to finish the repetition counting. Our work features two-fold: 1) An important insight that deep features extracted for action recognition can well model the self-similarity periodicity of the repetitive action is presented. 2) A high-energy based periodicity mining rule using deep features is presented, which can process arbitrary actions without preprocessing. Experimental results show that our method achieves comparable results on the public datasets YT_Segments and QUVA.

Index Terms—Action repetition counting, Deep ConvNets, Fourier transform with multi-stage threshold Filter



1 INTRODUCTION

Visual action repetition in real life appears in many applications, such as sports, music playing and manufacturing assembly. It is important to count the repetition of a specific motion in videos, which can be used for video question answering [1], action classification [2] [3] [4], segmentation [5] [6] [7], 3D reconstruction [8] [9], motion tracking [10] and motion planning of robots [11]. Due to the diversity of motion patterns and the limitations in video capturing (e.g., camera movement), the development of an universal solution for counting the repetitive actions remains under-explored.

To count the action repetition, early methods usually assumed that repetitive motions occurred in fixed scenes with regular periodicity. With this assumption, traditional features were used to analyze the action repetition, including human skeleton obtained by sensor device [8] [9], the wavelength spectrum [12]. However, the actions to be counted are usually captured in complex dynamic scenes and have variable periodicity over different periods, making the traditional features not suitable for the counting tasks. To tackle this complexity involved in real circumstances,

two methods have been proposed in recent years. In [13], multiple repetitive motion modes are simulated to construct the periodicity of the repetitive actions to realize the counting. Because the simulated motion modes are fixed, the algorithm can handle the specified modes well. But the performance significantly decreases for actions with other repetitive modes [14]. In [14], a counting method based on the detection of the moving area is proposed to achieve an improved performance. However, this method relies on additional preprocessing steps to detect the moving area. In sum, the methods based on the simulated action modes are not effective to count the varied types of repetitive motions, and the moving-area methods are highly dependent on the preprocessing performance. This motivates us to find an action repetition counting scheme that can work for unspecified motion modes without relying on extra detection steps. There are many challenges in action repetition counting for unconstrained videos. In the unconstrained videos, besides the interested repetitive action, there exists other motion information such as changes in the background, actions of the false objects and other unrelated movements. How to distinguish the periodic actions from these various irrelevant signals is prerequisite problem to address for counting the repetitive action. Moreover, action repetition modes are very different in different actions. For example, the repetition mode can be rotation, swinging, translation, and other modes. How to discover the relationship between the periodicity and various repetition modes is the another challenge. Additionally, the amplitude and the frequency of the repetition within the same action can also be different. Therefore, even for the same action, how to discover the

- Jianqin Yin, Yanchun Wu, Yonghao Dang, and Zhiyi Liu are with Automation School of Beijing University of Posts and Telecommunications, 100876, Beijing, China.
Huaping Liu is with Department of Computer Science, Tsinghua University, Beijing, China.
Jun Liu is with Department of Mechanical Engineering, City University of Hong Kong, Hong Kong, China.
Jianqin Yin and Yanchun Wu are contributed equally to this work. Jianqin Yin and Jun Liu are the corresponding authors.

Manuscript received Nov 19, 2019

features to count the repetition is an additional challenge. Although the repetition counting is very important, because of the aforementioned difficulties, the research progress on it is relatively slow in recent years. In contrast, great advances have been achieved by using deep learning methods for action recognition [15] [16] [17], which opens up new possibilities to propose new solutions for repetition counting. It has been proven that deep ConvNet can capture versatile and robust action features for action recognition, which illustrates that the deep features contain rich information of the action. For repetition counting, there is an important clue that the self-similarity periodic dynamics is the key. Accordingly, we propose in this paper that the deep features are helpful for mining the self-similarity periodicity of the action and can be used to count the repetition. This insight relieves us from the preprocessing and helps in addressing the challenges of modeling the periodicity of the unconstrained videos.

Although deep features may include periodic rules, the high dimension of the feature space makes it difficult to discover the self-similarity modes. In order to mine the repetitive rules from the deep features, we use Principal Component Analysis (PCA) to reduce the high dimension of the deep features. And we found that the non-stationary repetitive signals frequently appear in the first-dimensional principal component. This is an important insight for modeling the action repetition. Using PCA, the high dimensional deep features can be converted to a one-dimension waveform. In a word, the periodic signal can be generated combining deep features and PCA.

For the video including the repetitive action, most of the energy of the video usually comes from the repetitive motions due to the repetition. From this perspective, once we obtain the waveform that contains the periodic motions, we can locate the periodic motion using the highest energy rules. In detail, frequency analysis is used to locate the signals with highest energy corresponding to the repetitive action automatically, and then we can finish counting based on the located signal.

The framework of our method is shown in Fig. 1. Our framework includes four steps. At first, deep features are extracted using two-stream deep ConvNets, generating spatial and temporal features separately. Secondly, one-dimensional periodic signal is generated based on the deep features using PCA. Corresponding to the spatial and temporal features, there are two one-dimensional periodic signals. For simplicity, we only show the periodic signal of the spatial features in Fig. 1. Thirdly, the periodicity of the repetitive action is mined based on the highest energy rules using Fourier transform, filtering and inverse Fourier transform. Finally, peak detection is used to finish the repetition counting.

In this paper, a new action repetition counting method is proposed for solving the unconstrained action repetition counting in videos, and main contributions are summarized as follows. (1) We propose an energy-based action repetition counting method without extra preprocessing, which can be used to effectively count the repetition of the action with arbitrary periodic motions for unconstrained videos. (2) We

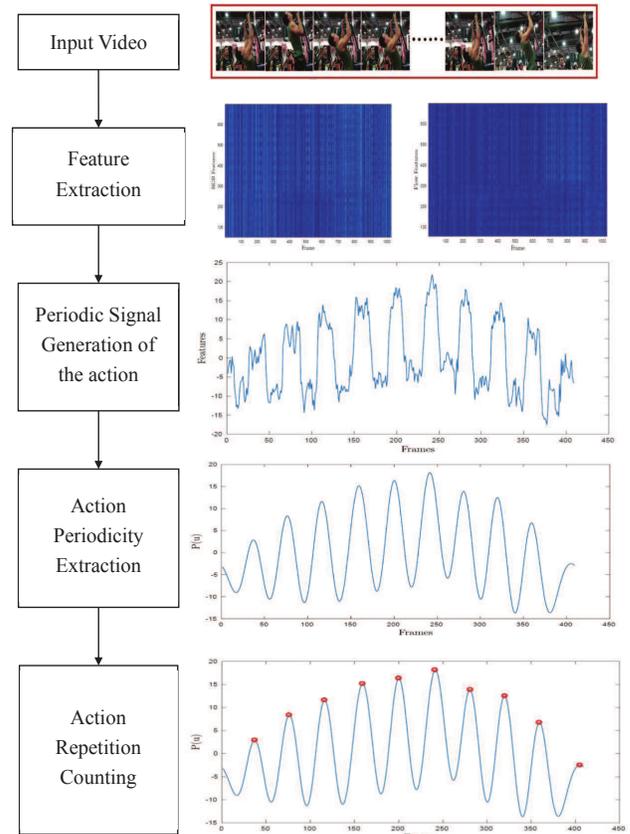


Fig. 1. The framework of action repetition counting in this paper.

give an important insight that the periodic self-similarity movement information can be well modeled by the deep features of the action used in action recognition. (3) We propose the Fourier transform with Multi-stage threshold Filter (FMF) scheme to automatically mine the interesting repetitive action and remove noises to improve the repetition counting performance.

The remainder of the paper is organized as follows. The next section investigates the related work. Section III discusses our algorithm in detail. The datasets and evaluation criteria are described in Section IV. Experimental results and analysis are illustrated in Section V. We conclude our paper in Section VI.

2 RELATED WORK

Action repetition counting is usually realized by converting the video into a one-dimensional waveform with the repetitive motion structures [8] [9] and then analyzing the spectral or frequency component by Fourier transform or wavelet analysis. Waveform analysis is also widely used in the periodic movement analysis [18] [19], which is very related to the repetition counting. At the same time, as discussed in the Introduction, the action feature is another important problem for counting. Therefore, we will review the related work from the following two aspects: periodic movement analysis and deep features for action analysis.

2.1 Periodic movement analysis

The existing methods have achieved remarkable results in video action periodic analysis tasks. Burghouts et al. [12] proposed a spatiotemporal filter bank for estimating of action repetition, which could work online. But it was limited to the motion of stationary scenes, and the filter bank was manually adjusted. Laptev et al. [20] used a matching method for action counting, whose primary work is to detect and segment repetitive motions using the geometrical constraints generated by the same motion repeatedly when the viewpoint changes. Ormoneit et al. [10] used functional analysis to represent cyclic movement. Ribnick et al. [8] [21] found that it is possible to reconstruct accurately periodic movements in 3D from a single camera view. Based on this research, they applied 3D reconstruction to gait recognition. Ren et al. [22] and Li et al. [23] developed two autocorrelation counting systems based on matching visual descriptors. Although both systems completed the repetition counting, they are postprocessing methods, which are only applicable to specific domains of restricted video. Pogalin et al. [18] get the information of a certain part of the body by tracking the object of interest, then performs pPCA and spectral analysis followed by detection and frequency measurement. But its purpose is to estimate the degree of periodic motion but not to count the repetition. Based on the human skeleton points captured by Kinect, Wang et al. [7] proposed an unsupervised repetitive motion segmentation algorithm based on the frequency analysis of the motion parameter, zero-velocity cross detection and adaptive k -means clustering. Although the above methods lay the foundation for the video repetitive counting task, they only realized the simple repetitive motion estimation of a fixed scene and had a poor performance on the diversity and non-stationary motion, which commonly exist in the real applications.

In recent years, Kumdee et al. [3] used the image self-similarity measure as the input of the multi-layer the perceptron neural network to determine whether the input video is a repetitive action. This method is relatively stable to image changes, noise, and low-resolution images. However, they focus on classifying that the video is a repetitive video or not but not counting. Levy et al. [13] proposed a method to count the repetitive action of the videos using the convolutional neural network. They used synthetic data to simulate four motion types for the periodic motion and carried out network training and prediction. In the test, the region of interest was calculated through the motion threshold for the test data. The motion cycle was classified through the classification network to complete the repetitive counting task. The method showed excellent performance on YT_Segments dataset. However, their algorithm decayed a lot when there are actions with different repetitive modes from the trained modes. The wavelet transform was presented in [14] to better deal with more complex and diverse video dynamics. From the flow field and its differentials, they derived 18 totally different repetitive perception. Based on the gradient, curl and divergence, a motion foreground segmentation representation based on flow was realized, and remarkable results were obtained. However their methods needs

the foreground segmentation, which is also a difficult problem. Therefore, we propose a method without extra preprocessing.

2.2 Deep features for action analysis

CNNs have been widely used in action recognition. Some of these CNNs use deep architectures with 2D convolutions to extract translation-invariant features in the video frames [15]. Specifically, Karpathy et al. [15] first introduced a CNN based method for action recognition and organized a large-scale sports video dataset (i.e., Sports-1M dataset) for training deep CNNs. To model the temporal information of the action, two-stream based CNN learning framework [16] [17] has been proposed. The two streams mean the spatial stream represented by RGB values and the temporal stream represented by pre-computed optical flow features. Because of the excellent balance between efficiency and effectiveness, the BN-Inception Network [24] is used as the backbone of the framework. The prominent characteristic of BN-Inception network is the Inception module, which carries out multi-scale processing and fusion of image features to extract better feature representation. Moreover, it is well known that the 3×3 convolution kernel has the best performance in VGG [25], and a very effective Batch Normalization(BN) method has been proposed to accelerate the learning of data distribution during training, making the accuracy of the classification improve significantly. In addition, the deep ConvNets can take pictures of any form as input to extract features. The training of deep ConvNets requires a large number of training samples to achieve good performance in action modeling. Nowadays, a large number of publicly available video datasets provide great convenience. Therefore, we extract the deep features of our experimental data using BN-Inception Network in this paper.

3 ACTION REPETITION COUNTING

In this part, we will discuss the proposed algorithm. As shown in Fig. 1, our algorithm includes four steps. Firstly, deep features of the unconstrained videos are extracted using deep ConvNets. Secondly, based on the high-dimensional deep features, the periodic signal is generated using PCA, and can obtain a one-dimensional waveform reflecting the repetitive changes of the videos. Thirdly, the action repetition rules are extracted based on the highest energy rules extracted from the 1D waveform, combining Fourier analysis, a new proposed multi-stage threshold filter, and the inverse Fourier transform. Finally, using the waveform, peak detection is used to count the repetition.

3.1 Deep Feature Extraction based on BN-Inception ConvNets

The Inception v2 based on Batch Normalization network [15] is used to obtain the features of the action, as shown in Fig. 2. It was trained on the large public Kinetics dataset [26], which contains 300,000 clip videos from real scenes,

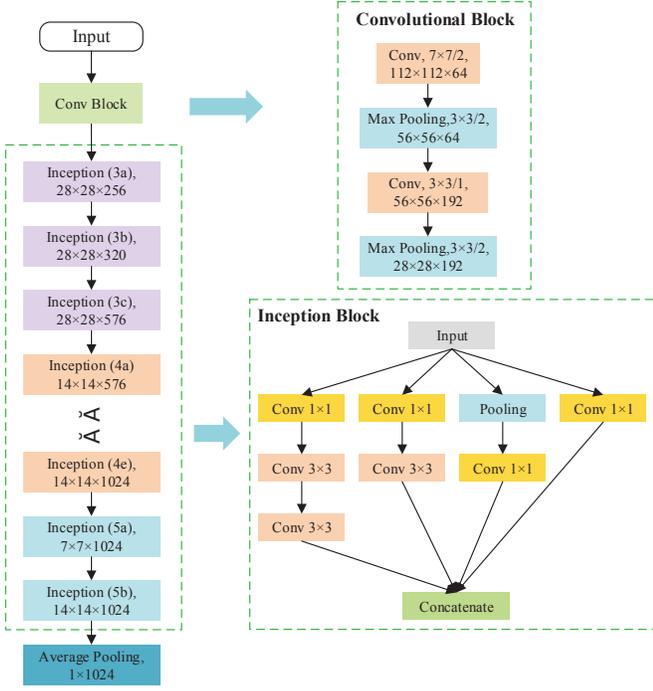


Fig. 2. The framework of BN-Inception Network

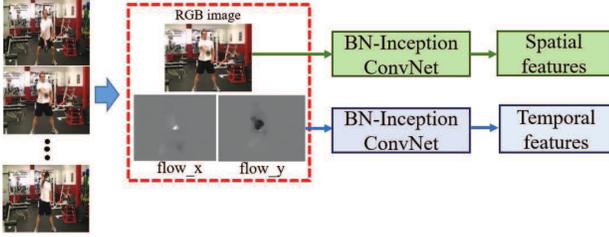
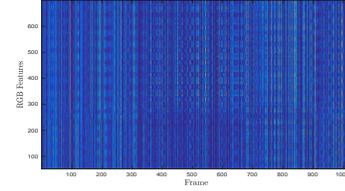


Fig. 3. The structure for spatiotemporal features extraction. The input modalities of BN-Inception are the RGB images and optical flow fields (x, y directions).

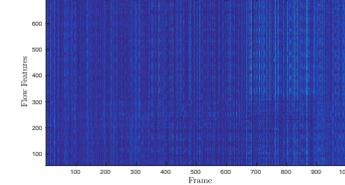
including 400 action categories, which is the largest action recognition dataset by far. The model we used is the best network trained in [25].

Two networks are used to extract robust action features, operating on two components, spatial and optical flow separately. The spatial flow network operates on the RGB image, which extracts spatial features describing the scene and object information. The optical flow network takes the pre-computed optical flow images as the input to extract the temporal features, which describe the motion information of the video. Robust spatiotemporal features are extracted by this method. Fig. 3 is a structure for feature extraction of the action.

Clipping and rotation of training images, to decrease the influence of the noise and increase the stability of features, are used to get the image set. In the process of feature extraction, we take the image set as network input and get the features in the *avg_pool* layer. Then the summation and the average features are computed in each dimension. Finally, the spatial features, denoted by $Static_{fea}$, and tem-



(a) RGB features of the action video



(b) Optical features of the action video

Fig. 4. The visualization of deep features

poral features, denoted by $Dynamic_{fea}$ are extracted for the single image, respectively. Due to the designed structure of the network, the dimension of the feature vector is 1024, as shown in equations (1) and (2).

$$Static_{fea} = (s_1, s_2, \dots, s_{1024}) \quad (1)$$

$$Dynamic_{fea} = (d_1, d_2, \dots, d_{1024}) \quad (2)$$

As discussed above, deep features can be obtained using the pre-trained models and we do not need to retrain the model using repetitive actions. The extracted deep features of the video is visualized in Fig. 4, where RGB features are given in Fig. 4(a) and optical flow features are given in Fig. 4(b). From the visualization results, although we can find some specified patterns; it is difficult for us to find the periodic rules using this high-dimension features. Therefore, we need some other methods to mine the periodicity of the repetitive action.

3.2 Periodic Signal Generation

To extract the periodicity information, we mine the hidden periodic action rules from two different features, spatial features ($Static_{fea}$) and temporal features ($Dynamic_{fea}$). The mining method for these two features are the same; therefore, we take spatial features $Static_{fea}$ as an example to explain our extraction method. Although deep features can classify the actions well, the counting of the repetition of the action is totally different from the classification of the action. For classification, one action is considered as a whole, and it focuses on the difference between different actions. On the contrary, action repetition counting focuses on locating the repetition of the same motion pattern. Therefore, the deep features extracted directly for action recognition may not be suitable for counting. We transform the high-dimensional features into an intuitive waveform by extracting the primary component of its covariance

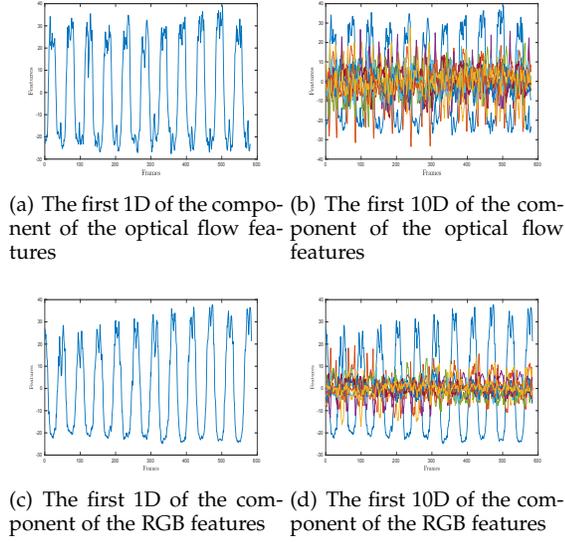


Fig. 5. Periodicity visualization of the repetitive action

matrix.

For simplicity, we represent the static features $Static_{fea}$ of the i th frame as $f_i = Static_{fea}$, where i is the frame index ranging from 1 to N . Therefore, for every action video, its features can be represented by the features of all its frames, denoted by $F = [f_1, f_2, \dots, f_N]$, where F is a 2D matrix with dimensions of $N \times D$, N is the total number of the frames in a video, D is the dimension of the spatial features, 1024. The detailed process is as follows.

Firstly, for the feature matrix F , we preprocess its i th components using $\bar{f}_i = \frac{f_i}{D}$, where, $i = 1, 2, \dots, N$. Then we obtain the mean matrix according to the formula $\bar{F} = [\bar{f}_1, \bar{f}_2 \dots \bar{f}_N]$. Then the transformation matrix \hat{F} is computed using $\hat{F} = F - \bar{F}$. Finally, the covariance matrix is calculated according to equation (3).

$$COV = \frac{1}{D} \hat{F} \hat{F}^T = V \Lambda V^T \quad (3)$$

We also can compute the eigenvalues and eigenvectors of the covariance matrix. The corresponding results are separately denoted by their matrix form as $\Lambda = diag(\lambda_1, \lambda_2 \dots \lambda_D)$ and $V = [\mu_1, \mu_2 \dots \mu_N]$, where each μ_i is a vector with $1 \times D$. We arrange Λ according to the value of its eigenvalue from large to small. And according to the new order of eigenvalues, we rearrange V to V' in columns. Then, we reserve the first eigenvector to get the transformation matrix V'_1 . The size of V'_1 is $D \times 1$. Therefore the mapped matrix $P = (p_1, p_2, p_3 \dots p_N)$ can be computed according to formula (4), where p_i is the principal component of i th frame of the videos, $i = 1, 2, \dots, N$. The size of P_i is $N \times 1$, by which the high-dimensional video features are transformed into the new space constructed by 1D waveform, as shown in Fig. 5.

$$p_i = V'_1 f_i \quad (4)$$

In this paper, the extracted spatial and temporal features are analyzed separately, and the mapping matrix P is obtained after the PCA transformation. To analyze the effect of different principal components, we also compute the first 10-dimensional principal component transformation matrix V'_{10} , the size is $D \times 10$. For each dimension, we get the mapped vector separately, the visualization results are shown in Fig. 5. From this figure, we can see that the first-dimensional feature includes more information on the motion characteristics of repetitive actions. Therefore, in this paper, the first-dimensional principal component is used to count the repetitive action.

3.3 Periodicity Mining and Repetition Counting

Due to the complexity and diversity of the videos captured in the real scene and the non-standardization during the action execution, the principal component contains lots of noises. As shown in Fig. 5(a), although there are some repetitive motion rules in the figure, the lower peak and the noises may lead to poor performance when counting. To locate the repetitive action, we need to distinguish the interesting actions from the unrelated noises. As discussed above, the interesting repetitive actions usually have two features. One is that it usually carries more energy than the other movement. The other is that it usually has a relatively higher frequency compared with the occasional camera motion. Therefore, we propose a high-energy-based repetitive action location method. And the counting is realized using the location results. In order to locate the high-energy repetitive action, we first give the frequency analysis of the repetitive action and then design a multi-stage filtering scheme.

3.3.1 Frequency Analysis of the repetitive action

We use the Fourier transform to analyze the repetition of the action. Fourier transform is usually used to estimate the power spectrum. Specifically, the time-varying signal is decomposed into the superposition of the components in the frequency domain by Fourier transform. The vibration frequency of the signal is separated to get the spectrum by equation (5), where P_u is the first-dimensional principal component of the u th frame, N is the total number of frames for a video.

$$X(k) = \sum_{u=1}^N P_u * exp\left(-j * 2\pi * (k-1) * \frac{u}{N}\right) \quad (5)$$

$$X(k)_{threshold \leq k \leq (L - threshold)} = 0 \quad (6)$$

$$P_u = \frac{1}{N} \sum_{k=1}^N X(k) * exp\left(j * 2\pi * (k-1) * \frac{u-1}{N}\right) \quad (7)$$

Fig. 6 is the visualization of the principal component of a repetitive action video, where Fig. 6(a) is the first dimensional principal component reflecting the repetitive motion, and Fig. 6(b) gives its corresponding Fourier results. From Fig. 6(a), we can see that there are lots of noises in

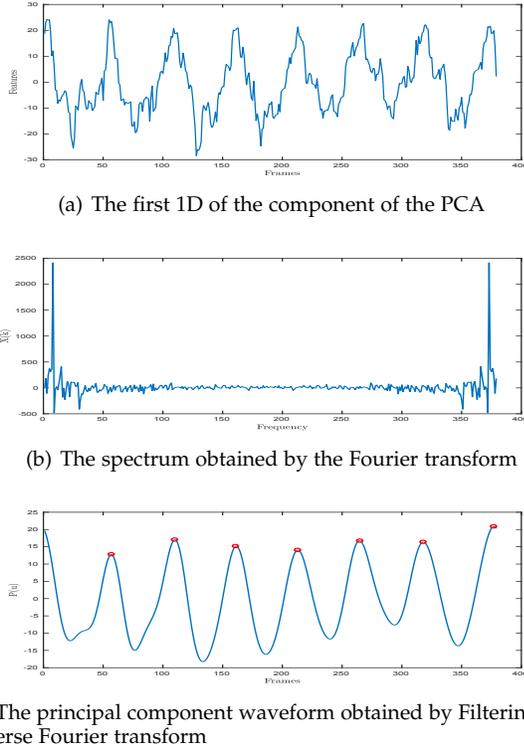


Fig. 6. Repetitive action characteristics over time

the 1D waveform. In other words, there are some other signals besides repetitive actions. Obviously, there are some unrelated actions or other motions in the video. To locate the repetitive action, we use the coefficients of the different frequency components to determine the filtered frequency band, and set the values of their corresponding frequency domain to 0, as in equation (6). Then the original signal P_u is obtained through the inverse Fourier transform by the equation (7). The above method is used to filter the signal noise, so that the frequency graph has a smooth trajectory with repetitive motion rules, as shown in Fig. 6(c). Finally, the counting results can be obtained by simple peak detection.

3.3.2 Fourier Transform with Multi-stage Threshold Filtering Analysis

In order to find the irrelevant frequency band discussed in the above section, we propose a filtering scheme, Fourier transform with Multi-stage threshold Filtering (FMF). To obtain the periodic waveform well reflecting the times of the repetitive action, we need to filter the noises and unrelated actions. The key is to locate the interesting frequencies and delete the unrelated frequencies. The difference between the interested repetitive motion and the unrelated motion lies in its frequency and its corresponding energy. Obviously, the noise often has high frequency and low energy, while the interested repetitive action may often have a relatively low frequency and large energy. However, different repetitive actions may have different frequencies. Therefore, the filtering algorithm must adaptively adjust the filtering frequency threshold according to the motion characteristics of

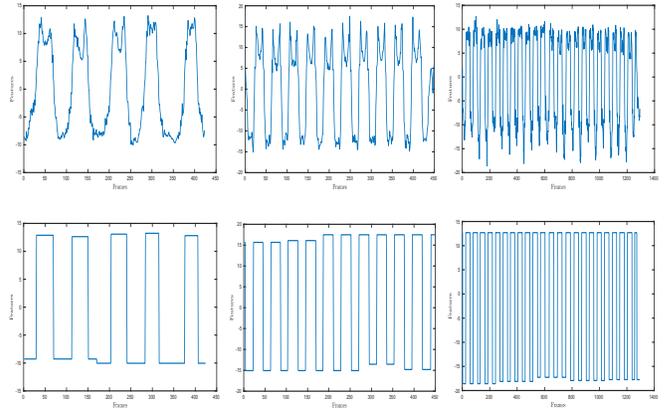


Fig. 7. The number of the high-pass frequency band

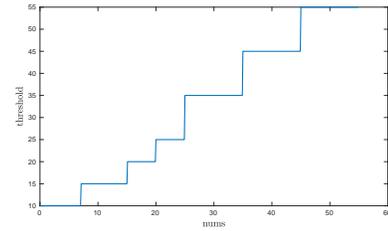


Fig. 8. The multistage thresholds

different frequencies, and we name our filtering algorithm as FMF. In order to do this, the number of the high-pass frequency band is calculated based on the first-dimensional principal component, as shown in Fig. 7. And according to the number, different thresholds are selected using different thresholds shown in Fig. 8, which is empirically obtained by analyzing the effect of different motions on noises and uses different thresholds according to different frequencies of motions. Then the noise is removed by the spectrum obtained by the inverse Fourier transform, like 6(c). By the multi-stage benchmark, we modify the motion waveform and make the repetitive motion have obvious temporal periodic characteristics, which is useful to the counting of repetitive actions.

4 EXPERIMENTS

YT_Segments [13] and QUVA [14] datasets from Youtube are used to evaluate our algorithm. These data are diverse and challenging, and they also include the camera and background movement. The repetitive actions have varying lengths and complex appearance patterns.

4.1 Datasets

YT_Segments. It contains 100 videos with repetitive actions, including exercise, cooking, architecture, biology, and so on. To create a clean benchmark from very diverse videos, which are pre-split and only contain repetitive actions. The number of repetitive actions is pre-labeled per video. The

smallest and largest numbers of the repetition are 4 and 50, respectively. The average duration of one video is 14.96s. Meanwhile, there are 30 videos with varying degrees of camera movement.

QUVA. It’s also made up of 100 videos and shows various kinds of repetitive video dynamics, including swimming, stirring, cutting, carding, and music-making. Compared with the YT_Segments dataset, it has more challenges in cycle length, motion appearance, camera motion, and background complexity. Therefore, the dataset is a more realistic and challenging benchmark for estimating repetitive action.

4.2 Evaluation Metrics and Baselines

Metrics. We use the same evaluation criteria [13] as those that used in the baselines as the metric for this task. For every video, the percentage of the absolute difference between ground truth G and the predicted value R is used: $\frac{|G-R|}{G} \times 100$. For N videos, we calculate the Mean Absolute Error(MAE) \pm standard deviation (σ) [13] as the evaluation metrics, where

$$MAE = \frac{1}{N} \sum_{i=1}^N \frac{|G_i - R_i|}{G_i} \quad (8)$$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (G - R)^2$$

Baselines. We compare our method with one classical method [18] and two recent methods in [13] [14]. When reporting the results, we directly make use of the results provided in the paper [14].

4.3 Results

Different thresholds and FMF. Using spatial features, we test the repetition counting results with different thresholds and FMF. At first, we test the results of different thresholds, as shown in Table 1. We set the different thresholds empirically as $\alpha = 10, 15, 20, 25, 30, 35, \dots$. On the YT_Segments dataset, the action repetition counting experimental results are shown in Table 1. From Table 1, we can see that the results of the single threshold are much worse than the multi-stage thresholds. Because the repetitive actions have different frequencies, the interested frequency in one video may become noise frequency in the other video. And we can also find that the relatively higher accuracy can be obtained by threshold 15. We think the possible reason is that lots of actions may have this frequency. However, for some high-frequency actions, this threshold will filter the interested frequency out, which makes that the multi-stage threshold performs better.

Using spatial features, temporal features and the fusion features of spatial and temporal features, we further compare the repetition counting performance using FMF module and without FMF, as shown in Table 2. And the influence of FMF module of our method is also analyzed.. From Table 2, we can see that the accuracy is improved significantly by adding FMF module on the YT_Segments dataset. It also shows that the spatial features based on RGB achieve the highest accuracy on the YT_Segments dataset.

TABLE 1
Comparative Results of different thresholds

α	MAE
10	23.3 \pm 63.4
15	19.9 \pm 42.8
20	30.9 \pm 38.14
25	44.6 \pm 42.8
30	56.0 \pm 62.7
35	67.2 \pm 91.6
multi-stage α	8.7 \pm 3.9

TABLE 2
Experimental Results Comparison about Whether Remove FMF on YT_Segments Dataset

MAE Without FMF			MAE With FMF		
RGB	FLOW	RGB FLOW	RGB	FLOW	RGB FLOW
13.7 \pm 6.2	29.2 \pm 21.3	18.2 \pm 9.4	8.7\pm3.9	21.9 \pm 12.7	15.8 \pm 6.6

Comparisons with state-of-the-art baselines. We compare the performance of our method with state-of-the-art baselines. For YT_Segments and QUVA, the spatial features and temporal features are used separately. At the same time, for convenience of comparison, we give the average of the methods on YT_Segments and QUVA, and we name it Overall. The quantitative comparisons are presented in Table 3. Compared with the existing baselines, our method achieves comparable performance without extra preprocessing. For the YT_Segments dataset, the method of [13] performs best with the MAE of 6.5. The MAE in [14] is 10.3, which is better than the article [18]. Our method is superior to articles [18] [14] with the MAE of 8.7, but the standard error achieves the best performance compared with the above methods, which illustrates that the worst performance of our method is the best in all the methods. Moreover, the results also show that our method can achieve good performance under the static background. In the more challenging QUVA dataset, our experimental results also achieved decent performance. The method [13] performed the worst with the MAE of 48.2, because their network considered only four types of action during training. The method of [18] was 38.5. In [14], the MAE was 23.2. The performance of our method also comes in the second place, which is very close to the best performance. And the standard error is also the least. These results show that our method can also adapt well to the dynamic background. In summary, compared to the above methods, we get the best standard error on the two public datasets. At the same time, we get the MAE very close to the state-of-the-art baseline. The results show that our method can achieve comparable results in counting action repetition for unconstrained videos with a decent framework, not relying on preprocessing.

Detailed results and analysis. To further validate our contribution, on YT_Segments and QUVA, we give the counting results in detail, as shown in Fig. 9. From Fig. 9(a), we can see that the counting results of most of the actions are very close to their groundtruth and the differences cluster around positive or negative 1. Our peak-

TABLE 3
Comparisons with the state-of-the-art baselines

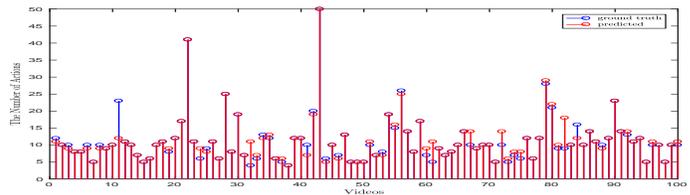
Methods	Datasets	YT_Segments	QUVA	Overall
Pogalin et al. [18]		21.9±30.1	38.5±37.6	30.2±33.9
Levy & Wolf [13]		6.5±9.2	48.2±61.5	27.4±35.4
Runia & Snoek [14]		10.3±19.8	23.2±34.4	16.8±27.1
Ours		8.7±3.9	25.1±26.3	16.9±10.4

detection counting scheme can well explain this. Because we use the number of the peaks as the counting results, it is slightly different from the repetitive case, where the repetition is, in fact, a cycle. Therefore, using more detailed cycle detection may solve this problem. In addition, there is only 1 sample (video 11) whose error is relatively large, and the corresponding video is a repetitive action with the sub-movements of the left-arm and right-arm. Our method regards one repetitive action, which includes the repetitive motion of the left-arm and right-arm as once. However, the groundtruth regards that the motion of the left-arm and right-arm are twice. We think our result is interpretable because these two movements are, in fact, two sub-movements of one movement. From this point, although our performance is not good enough, but the actual performance of our method is largely better than we gave.

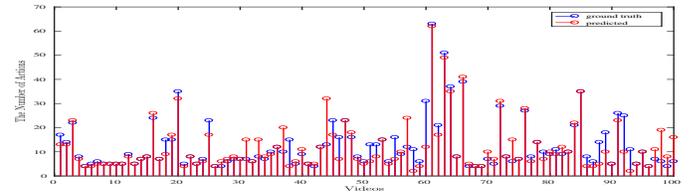
From Fig. 9(b), besides the above mentioned negative and positive 1 difference problem, there are also other problems. Some groundtruth results are almost twice as our results, for example, video 34. We found this error occurs because those actions often include two coupled actions; for example, in video 34, the whole action includes the left leg and right leg sub-movements. This will contribute 2 times of 1 action in our algorithm. The other error is due to the changing of the view angles of the action, like video 31. For these two problems, we think they can be solved using more advanced signal processing methods. And we will focus on it in our future work.

In summary, the main error of our method lies in two aspects. One is the difference between the peak and the cycle. The other is the coupled submovements in action. Both kinds of errors become from the processing of the deep features. Therefore, this also validates our insight that the 1D principle component of the deep features can well model the periodicity of the repetitive action.

Feature analysis. Fig. 10 shows the visualization results of different features by our method on YT_Segments and QUVA datasets. Only three parts are shown here, and the first part is the input data to get the deep features by BN-Inception Network, the second part is the sequence waveform of action obtained by PCA when $k = 1$, and the third part is the final result for counting. The experimental results indicate that static features achieve better results in a relatively clean background, and results are poorer with serious background interference. However, the method based on optical flow features achieved prominent results. They demonstrate the effectiveness of the proposed method in this paper. Moreover, we also can find that



(a) The counting results on YT_Segments



(b) The counting results on QUVA

Fig. 9. The detailed counting results of our method

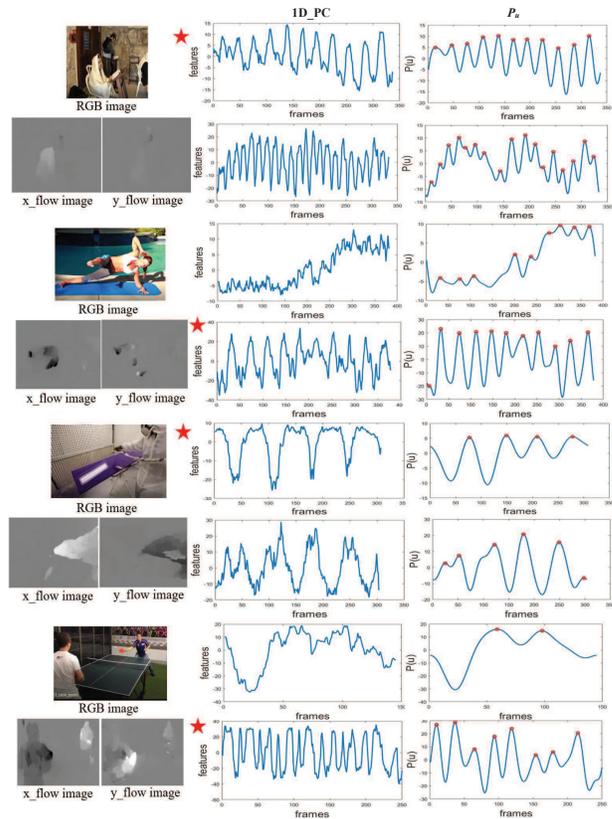


Fig. 10. The comparison of counting results based on RGB features and optical flow features, where 1D PC is the first dimensional principal component by PCA and the corresponding result by FMF based on the 1D PC over the time sequence. The position of Red star indicates the optimal detection result.

the 1D principal component of the deep features includes good periodic information of the action. If we can further improve the signal processing after the periodicity mining, the performance can improve further.

5 CONCLUSIONS

We propose an important insight that the periodicity of the action can be well modeled by the deep features extracted from the action recognition task. We think this insight is very important for repetition counting due to two reasons. On the one hand, the repetition counting method can borrow the state-of-the-art results or experiences from action recognition, which decreases the gap between the development of action recognition and the repetition counting. On the other hand, this insight can simplify the repetition counting task from the trivial preprocessing or synthetic mode generation. Based on this insight, we propose a new counting method using high energy rules for unconstrained video. Due to the introduction of energy, our method can solve unconstrained action modes in unconstrained videos. In detail, by using the training model based on Kinetics, we extract reliable deep features, including the temporal evolution characteristics of video actions and the unique appearance and spatiotemporal characteristics of motion patterns by deep ConvNets. Then, the periodic movement information can be obtained by the PCA based on the deep features of the action. Besides, we compute the frequency spectrum by Fourier transform to remove noise information by multi-stage threshold filtering. Then, the time sequence waveform is smoothed, and the action repetition counting task is completed according to peak detection. Extensive experimental results show the effectiveness of our method.

Compared with the existing methods, our method is simple and flexible without preprocessing. However, it still has poor performance when there is interference or chaotic background in the motion, especially when there are coupled repetitive actions and continuous changing of the viewpoints. These non-periodic interference motions make it impossible to analyze the motion characteristics of the target object accurately. We will focus on these problems in the future.

ACKNOWLEDGMENTS

This work was supported in part by the National Natural Science Foundation of China under Grant 61673192, U1613212, and in part by the Basic Scientific Research Project of the Beijing University of Posts and Telecommunications under Grant 2018RC31.

REFERENCES

- [1] Yunseok Jang, Yale Song, Chris Dongjoo Kim, Youngjae Yu, Youngjin Kim, Gunhee Kim. "Video question answering with spatio-temporal reasoning," *International Journal of Computer Vision*, vol. 127, pp.1385-1412, 2019.
- [2] S. M. Seitz and C. R. Dyer, "View-invariant analysis of cyclic motion," *International Journal of Computer Vision (IJCV)*, vol. 25, on. 3, pp. 231-251, 1997.
- [3] O. Kumdee, P. Ritthipravit, "Repetitive motion detection for human behavior understanding from video images," in 2015 IEEE International Symposium on Signal Processing and Information Technology (ISSPIT), IEEE, 2015, pp. 484-489.
- [4] O. Kihl, D. Picard, P. H. Gosselin, "Local polynomial space-time descriptors for action classification," *Machine Vision and Applications (MVA)*, vol. 27, on. 3, pp. 351-361, 2016.
- [5] C. Lu and N. J. Ferrier, "Repetitive motion analysis: Segmentation and event classification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 26, on. 2, pp. 258-263, 2004.
- [6] A. B. Albu, R. Bergevin, S. Quirion, "Generic temporal segmentation of cyclic human motion," *Pattern Recognition*, vol. 41, no. 1, pp. 6-21, 2008.
- [7] Q. Wang, G. Kurillo, F. Ofli, et al., "Unsupervised temporal segmentation of repetitive human actions based on kinematic modeling and frequency analysis," in 2015 international conference on 3D vision (IEEE), 2015, pp. 562-570.
- [8] E. Ribnick, R. Sivalingam, N. Papanikolopoulos, and K. Daniilidis, "Reconstructing and analyzing periodic human motion from stationary monocular views," *Computer Vision and Image Understanding (CVIU)*, vol. 116, on. 7, pp. 815-826, 2012.
- [9] B. Wandt, H. Ackermann, B. Rosenhahn, "3d reconstruction of human motion from monocular image sequences," *IEEE transactions on Pattern Analysis and Machine Intelligence (PAMI)*, vol. 38, on. 8, pp. 1505-1516, 2016.
- [10] D. Ormoneit, M. J. Black, T. Hastie, et al., "Representing cyclic human motion using functional analysis," *Image and Vision Computing*, vol. 23, no. 14, pp.1264-1276, 2005.
- [11] T. F. Iversen, L. P. Ellekilde, "Kernel density estimation based self-learning sampling strategy for motion planning of repetitive tasks," in 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2016, pp. 1380-1387.
- [12] G. J. Burghouts and J.-M. Geusebroek, "Quasi-periodic spatiotemporal filtering," *IEEE Transactions on Image Processing (TIP)*, vol. 15, on. 6, pp. 1572-1582, 2006.
- [13] O. Levy and L. Wolf, "Live Repetition Counting," in Proceedings of the IEEE International Conference on Computer Vision (CVPR). 2015, pp. 3020-3028.
- [14] T. F. H. Runia, C. G. H. Snoek, A. W. M. Smeulders, "Real-World Repetition Estimation by Div, Grad and Curl," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). 2018, pp. 9009-9017.
- [15] A. Krizhevsky, I. Sutskever, and G. Hinton, "ImageNet classification with deep convolutional neural networks," in Advances in neural information processing systems (NIPS), 2012, pp. 1097-1105.
- [16] K. Simonyan and A. Zisserman "Two-stream convolutional networks for action recognition in videos," in Advances in neural information processing systems (NIPS), 2014, pp. 568-576.
- [17] L. Wang, Y. Xiong, Z. Wang, "Temporal segment networks: Towards good practices for deep action recognition," in European Conference on Computer Vision. Springer, Cham, 2016, pp. 20-36.
- [18] E. Pogalin, A. W. M. Smeulders, and A. H. Thean. "Visual quasi-periodicity," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2008.
- [19] R. Cutler, L. S. Davis, "Robust real-time periodic motion detection, analysis, and applications," *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, vol. 22, on. 8, pp. 781-796, 2000.
- [20] I. Laptev, S. Belongie, P. Perez, and J. Wills, "Periodic motion detection and segmentation via approximate sequence alignment," in Tenth IEEE International Conference on Computer Vision (ICCV), 2005, pp. 816-823.
- [21] E. Ribnick and N. Papanikolopoulos, "3D reconstruction of periodic motion from a single view," *International Journal of Computer Vision (IJCV)*, vol. 90, on. 1, pp. 28-44, 2010.
- [22] Y. Ren, B. Fan, W. Lin, X. Yang, H. Li, W. Li, and D. Liu, "An efficient framework for analyzing periodical activities in sports videos," in 2011 4th International Congress on Image and Signal Processing. IEEE, 2011, pp. 502-506.
- [23] G. Li, X. Han, W. Lin, and H. Wei, "Periodic motion detection with ro-based similarity measure and extrema-based reference selection," *IEEE Transactions on Consumer Electronics*, vol. 58, on. 3, pp.947-95, 2012.
- [24] S. Ioffe, C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.
- [25] K. Simonyan, A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [26] W. Kay, J. Carreira, K. Simonyan, "The kinetics human action video dataset," *arXiv preprint arXiv:1705.06950*, 2017.
- [27] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway: IEEE, 2016, 770-778.
- [28] N. Xiao, P. Yang, Y. Yan, "From Communication to Sensing: Recognizing and Counting Repetitive Motions with Wireless Backscattering," *arXiv preprint arXiv:1810.11707*, 2018.

- [29] R. Polana, R. C. Nelson, "Detection and recognition of periodic, nonrigid motion," *International Journal of Computer Vision (IJCV)*, vol. 23, no. 3, pp. 261-282, 1997.