

Nrityantar: Pose oblivious Indian classical dance sequence classification system

Vinay Kaushik*
Dept. of Electrical Engineering
IIT Delhi
eez158117@ee.iitd.ac.in

Prerana Mukherjee*
Dept. of Electrical Engineering
IIT Delhi
eez138300@ee.iitd.ac.in

Brejesh Lall
Dept. of Electrical Engineering
IIT Delhi
brejesh@ee.iitd.ac.in

ABSTRACT

In this paper, we attempt to advance the research work done in human action recognition to a rather specialized application namely Indian Classical Dance (ICD) classification. The variation in such dance forms in terms of hand and body postures, facial expressions or emotions and head orientation makes pose estimation an extremely challenging task. To circumvent this problem, we construct a pose-oblivious shape signature which is fed to a sequence learning framework. The pose signature representation is done in two-fold process. First, we represent person-pose in first frame of a dance video using symmetric Spatial Transformer Networks (STN) to extract good person object proposals and CNN-based parallel single person pose estimator (SPPE). Next, the pose basis are converted to pose flows by assigning a similarity score between successive poses followed by non-maximal suppression. Instead of feeding a simple chain of joints in the sequence learner which generally hinders the network performance we constitute a feature vector of the normalized distance vectors, flow, angles between anchor joints which captures the adjacency configuration in the skeletal pattern. Thus, the kinematic relationship amongst the body joints across the frames using pose estimation helps in better establishing the spatio-temporal dependencies. We present an exhaustive empirical evaluation of state-of-the-art deep network based methods for dance classification on ICD dataset.

CCS Concepts

•Computing methodologies → Supervised Learning; Feature selection; *Computer Vision; Image Processing;*

Keywords

Pose Signature; Dance classification; Action Recognition; Motion and Video Analysis; Deep Learning; Supervised Learning; LSTM

*Equal contribution

1. INTRODUCTION

Research in action recognition from video sequences has expedited in the recent years with the advent of large-scale data sources like ActivityNet[4] and ImageNet [13] and availability of high computing resources. Along with other computer vision areas, Convolutional neural networks (ConvNets) has been extended for this task as well by utilizing the spatio-temporal filtering approaches [10], multi-channel input streams [26, 11], dense trajectory estimation with optical flow [5]. However, the extent of success is not that overwhelming as is in the case of image classification and recognition tasks. Long Short Term Memory networks (LSTMs) capture the temporal context such as pose and motion information effectively in action recognition [22]. We attempt to advance the research work done in human action recognition to a rather specialized domain of Indian Classical Dance (ICD) classification.

There is a rich cultural heritage prevalent in the Indian Classical Dance or Shashtriyā Nritya¹ forms. The seven Indian classical dance forms include: Bharatanatyam, Kathak, Odissi, Manipuri, Kuchipuri, Kathakali and Mohiniattam. Each category varies in the hand (or hasta-mudras) and body postures (or anga bhavas), facial expressions or emotions depicting the nava-rasas², head orientation, as well as the rhythmic musical patterns accompanied in these. Even the dressing attires and make-up hugely differs across them. All these put together constitute the complex rule engines which govern the gesture and movement patterns in these dance forms. In this work, however we do not address the gesture aspect of these complex dance forms. In this paper, we classify these categories based on the human body-pose estimation utilizing sequential learning. Instead of feeding a simple chain of joints in the sequence learner which generally hinders the network performance, we constitute a feature vector of the normalized distance vectors which captures the adjacency configuration in the skeletal pattern. Thus, the kinematic relationship amongst the body joints across the frames helps in better establishing the spatio-temporal dependencies.

During an Indian classical dance performance, the dancer may be often required to sit in a squat position or half-sitting

¹“Shastriya Nritya” is the Sanskrit equivalent of “Classical dance”.

²“Rasa” is an emotion experienced by the audience created by the facial expression or the feeling of the actor. “Nava” refers to nine of such emotions namely, erotic, humor, pathetic, terrible, heroic, fearful, odious, wondrous and peaceful.



Figure 1: Indian Classical Dance (ICD) Forms: (a) Bharatnatyam (b) Kuchipudi (c) Manipuri (d) Kathak (e) Mohiniattam (f) Odissi

pose, cross-legged pose, take circular movements or turns, thus leading to severe body-occlusion which renders the human pose estimation as a highly challenging task. The addition of certain dress attires further complicates the pose estimation. As in the case of Manipuri dance, women dancers wear an elaborately decorated barrel shaped long skirt which is stiffened at the bottom and closed near the top due to which the leg positions are obscured. In some dance forms such as Kathakali, only the facial expressions are highlighted and the dancers wear heavy masks. In this work, we leverage the feature strength by constructing dance pose shape signature which is fed to a sequence learning framework. It encodes the motion information along with the geometric constraints to handle the aforementioned problems. In this paper, we address six oldest Indian dance classes namely Bharatnatyam, Kathak, Odissi, Kuchipudi, Manipuri and Mohiniattam. Bharatnatyam is one of most popular ICD belonging to the southern belt of India, particularly originated in Tamil Nadu. Kuchipudi emanated from Andhra Pradesh and Telangana regions. Mohiniattam belongs to Kerala. Kathak originated in the northern part of India. Manipuri and Odissi are from eastern part of India, Manipur and Orissa respectively. We have not considered Kathakali ICD as it involves mainly facial expressions and does not contain enough visual dance stances to be catered in the adopted classification pipeline. The different dance forms have been depicted in Fig. 1.

Remaining sections in the paper are organized as follows. In Sec. 2 we discuss the related work in dance and action recognition classification. In Sec. 3, we outline the methodology we propose to detect the person dancing, track the dancer’s pose and movement as well as characterize the dancer’s trajectory to understand their behavior. In Sec. 4, we discuss experimental results and conclude the paper in Sec. 5.

2. RELATED WORKS

2.1 Handcrafted feature based approaches

Most of the traditional works in action recognition domain constitute of constructing handcrafted features to discriminate action labels [15, 21, 14, 2, 6]. Next, these features were encoded into Bag of Visual Words (BoVW) or Fisher Vector (FV) encodings. The classifier models are then trained on these encodings and classify the videos with suitable labels. However, these feature sets were not opti-

mized and failed to capture the temporal dependencies in the video streams to encode the high-level contextual information. Histogram of Oriented Gradients [6] have remained the de-facto choice in capturing the shape information in action recognition task. Histogram of Optical Flow magnitude and orientation (HOF) [7] models the motion. They obtain better performance since they represent the dynamic content of the cuboid volume. Space time interest point (STIP) detectors [14] are extensions of 2D interest point detectors that incorporate temporal information. Similarly, as in the 2D case, STIP features are stable under rotation, viewpoint, scale and illumination changes. Spatio-temporal corners are located in region that exhibits a high variation of image intensity in all three directions(x,y,t). This requires that such corner points are located at spatial corners such that they invert motion in two consecutive frames (high temporal gradient variation). They are identified from local maxima of a cornerness function computed for all pixels across spatial and temporal scales. Holistic representation of each action sequence is done by a vector of features. Motion history images [1] obtained utilizing Hu moments are utilized as global space-time shape descriptors.

2.2 Deep feature based approaches

With the deep learning networks, most of the computer vision tasks have reached remarkable accuracy gains. Human action recognition can be solved more effectively if the entire 3D pose comprehensive information is being exploited in the deep networks [10, 24]. In [25], authors utilize the co-occurrence of the skeletal joints in a sequence learning framework. However, since there is a joint optimization framework involved it is computationally expensive. In [20], authors investigate a two-stream ConvNet architecture and jointly model the spatio-temporal information using separate networks to capture these. They also demonstrate that multi-frame dense optical flow results in improved performance gains in spite of limited training data. In [23], authors utilize trajectory constrained pooling to aggregate the discriminative convolutional features into trajectory-pooled deep-convolutional descriptors. They construct two normalization methods for spatiotemporal and channel information. Trajectory-pooled deep-convolutional descriptor performs sampling and pooling on aggregated deep features thus improving the discriminative ability.

2.3 Pose feature based approaches

Neural Network (CNN) architectures [3, 17] can provide excellent result for 2D human pose estimation with images. However, handling occlusion is quite challenging task. In [18], authors obtain pose proposals by localizing the pose classes using anchor poses. Then, the pose proposals are scored to get the classification score and regressed independently. Since, pose regression techniques result in a single pose estimation thus cannot handle multimodal pose distributions. In order to capture such information, the pose estimates can be discretized into various bins followed by pose classification. In [16], authors construct a classification-regression framework that utilizes classification network to generate a discretized multimodal pose estimate followed by regression network to refine the discretized estimate to a continuous one. The architecture incorporates various loss functions to enable different models. In [26], authors incorporate pose, motion and color channel information into a

Markovian model to incorporate the spatial and temporal information for localization and classification.

3. PROPOSED METHODOLOGY

In this section, the proposed framework and its main components are discussed in detail including the recognition of a dance move from a sequence of frames in a video using LSTM networks and feature extraction for the video frames. First, we extract features using uniform frame skipping in a sequence of frames such that the skipping of frame does not affect the sequence of the action in the video. The feature vector comprises of multiple components. We perform feature extraction using a 3-tier framework combining Inception V3 features, 3D CNN features and novel pose signatures. The second component trains our model for dance classification. It takes a chunk of features as input where a chunk is defined by a collection of features for the selected frames in the video. The feature chunks are then fed to a LSTM network in order to classify various dance forms. The output of the LSTM layer is connected to 2 fully connected layers with Batch Normalization at each stage and is finally connected to a softmax layer of size 6, corresponding to the 6 ICD dance forms.

3.1 Feature Construction: Inception V3 Features, 3D CNN, Pose Signature

The feature vector for training our LSTM model is an embedding of various deep features (both 2D and 3D) and geometric cues in a single vector learning to classify the classical dance type. The 2D spatial cues are learnt using InceptionV3 network weights trained on ImageNet dataset [8]. The 3D temporal cues are described by using ResNext-101 with cardinality 32 on Kinetics dataset [12] which contains 400 video classes. In order to construct the pose signature, geometric information is used to construct features comprising of pose, flow and other structural information. The variation in Indian classical dance forms in terms of hand and body postures, head orientation makes pose estimation an extremely challenging task. To circumvent this problem, we construct a pose-oblivious shape signature which is fed to a sequence learning framework. The pose estimation is done as a two-fold process. The pose estimation is done using Alphapose framework. First, we represent person-pose in a frame of a dance video using a symmetric Spatial Transformer Networks (STN) to extract good person object proposals and CNN-based parallel single person pose estimator (SPPE). Next, the pose basis are converted to pose flows by assigning a similarity score between successive poses followed by non-maximal suppression. The video frames are fed into Alphapose framework computes the skeletal pose comprising of anchor joints at various body parts of a human being. The anchor joints are used to further compute various geometric features, resulting into a pose signature (as explained in Sec. 3.2. The feature generation strategy is described in detail further below.

Inception V3 features: These features widely are used for classification of images. The input size is 299x299x3 and the output layer gives a feature of size 1x1x2048. We utilize the pre-trained ImageNet model used for image classification. This has several advantages over the other networks. It incorporates Batch Normalization to the previous architectures, replaces the 5x5 convolution layers in Inception V2 by two 3x3 convolution layers and has the benefit of an

added BN-Auxiliary layer i.e. the fully connected layer of the auxiliary classifier is also normalized, not just the convolutional layers.

ResNext-101 Kinetics features: It is pre-trained on Kinetic dataset with 400 video classes. We have used 16 consecutive frames for generating a 3D feature. The architecture comprises of 101 convolutional layers with skip connections and cardinality of 32. The output can be represented using a feature vector of size 2048 or 16x128 for 16 frames. Thus, it gives a 128-dimensional feature vector per frame.

3.2 Pose Signature

After the AlphaPose estimation, we construct novel pose signature which consists of the normalized distances and angles between the anchor joints in the pose vector. The pose vector constitutes 16 anchor joints namely, 1: *foot_{right}*, 2: *knee_{right}*, 3: *hip_{right}*, 4: *hip_{left}*, 5: *knee_{left}*, 6: *foot_{left}*, 7: *hip_{center}*, 8: *spine*, 9: *shoulder_{center}*, 10: *head*, 11: *hand_{right}*, 12: *elbow_{right}*, 13: *shoulder_{right}*, 14: *shoulder_{left}*, 15: *elbow_{left}*, 16: *hand_{left}* as shown in Fig. 3. We utilize joint 7 as the reference point for taking the normalized distances to all other anchor joints. The distance metric used is Euclidean distance norm. We also incorporate the angles between the key anchor joints to characterize the dance stances. The angles are considered between these ordered anchor joint pairs: {1, 3}, {2, 7}, {4, 6}, {5, 7}, {11, 13}, {9, 12}, {9, 15}, {14, 16}. We embed the normalized distances between the leg joints ({1, 6}, {2, 5}) and hand joints ({11, 16}, {12, 15}) as well in the pose signature. Next, we calculate the flow vectors between the anchor joints in successive frames to characterize the temporal dependencies. We also compute the flow directionality for the anchor joints across successive frames. Thus, in total the pose signature constitutes a 75D vector. Instead of simply feeding a chain of joints in the sequence learner which generally hinders the network performance we constitute a pose signature which captures the adjacency configuration in the skeletal pattern. Hence, the kinematic relationship amongst the body joints across the frames using pose estimation step helps in better establishing the spatio-temporal dependencies.

3.3 Sequence Learning Framework

The feature vectors obtained from 3DCNN, pose signature and ImageNet features are concatenated to be fed into a sequence learning framework. We utilize Long Short Term Memory networks for dance sequence classification. Usually, Recurrent Neural Networks (RNNs) are difficult to train with various activation functions such as tanh and sigmoid due to the problems of vanishing gradient and error exploding [9]. LSTMs can be used to learn the long-term dependencies in place of RNNs and it solves the aforementioned problems. It consists of one self-connected memory cell c and three multiplicative units, input gate i , forget gate f and output gate o .

Given an input sequence $x = (x_0, \dots, x_T)$, the activations of the memory cell and three gates are given as follows:

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \quad (1)$$

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \quad (2)$$

$$c_t = f_t c_{t-1} + i_t (\tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c)) \quad (3)$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \quad (4)$$

$$h_t = o_t \tanh(c_t), \quad (5)$$

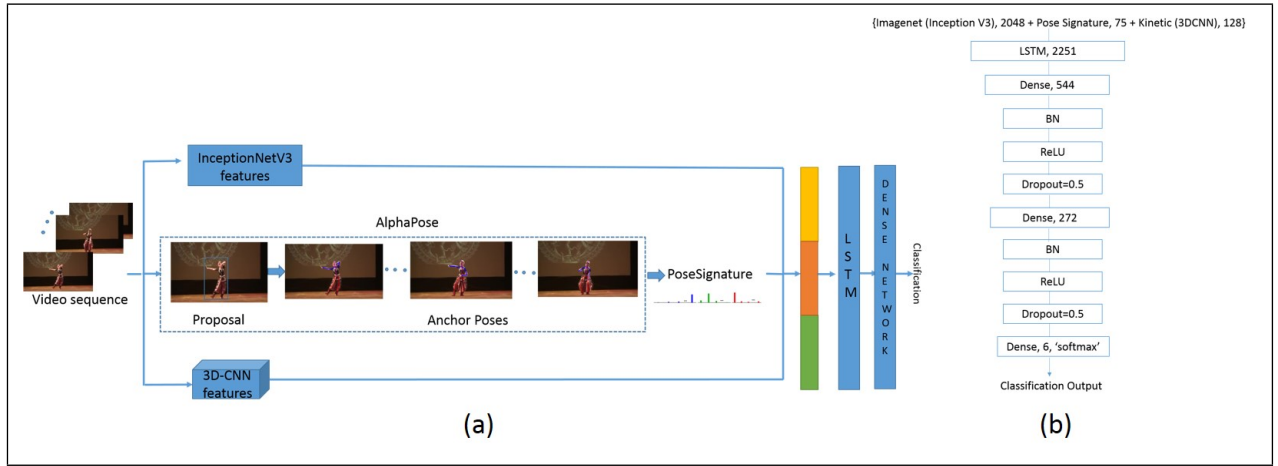


Figure 2: (a) Achitecture of the Nrityantar framework (b) Block Diagram of the sequence learning framework in Nrityantar. The number in the layers indicate the number of features.



Figure 3: (a) Pose of a Kathak dancer (b) Visualization of the anchor joints

where $\sigma(\cdot)$ is the sigmoid function, W are the connection weights between two units and h denotes the output values of a LSTM cell.

We have fed input feature vector into the LSTM network. The concatenated feature is 2251D vector. We have utilized a dense network after an LSTM layer followed by softmax layer for dance sequence classification as shown in Fig. 2(b). The constructed feature vector is able to characterize the spatio-temporal dependencies between the anchor joints. We utilize the cross-entropy loss function given as below.

Training Loss: We use Adam optimizer with Categorical cross entropy for training. The equation for categorical cross entropy is:

$$\frac{-1}{N} \sum_{i=1}^N \sum_{c=1}^C 1_{y_i \in C_c} \log p_{model}[y_i \in C_c] \quad (6)$$

The double sum is over the observations (video features) i , whose number is N , and the categories c , whose number is C . The term $1_{y_i \in C_c}$ is the indicator function of the i^{th} observation belonging to the c^{th} category. The $p_{model}[y_i \in C_c]$ is the probability predicted by the model for the i^{th} observation to belong to the c^{th} category. When there are more than two categories, the neural network outputs a vector of C probabilities, each giving the probability that the network input should be classified as belonging to the respective category.

4. EXPERIMENTAL RESULTS

In this section, we evaluate our model and compare with other five different architectures on a benchmark dataset for Indian classical dance: ICD dataset [19]. We also discuss the overfitting issues and the computational efficiency of the proposed model.

4.1 Evaluation Dataset and Parameter Settings

ICD dataset mainly consists of curated videos from YouTube of primarily 3 oldest dance forms Bharatnatyam, Kathak and Odissi with video class annotation. Each class consists of 30 video clips of maximum resolution 400Å350 and of 25 sec maximum duration. We have done data augmentation with these videos for classes Manipuri, Kuchipudi and Mohiniattam from YouTube. During data processing, we further clipped the video segments into 5-6 seconds chunks of frames at 25 fps to generate a maximum of 150 frames. The train to test ratio for evaluation has been selected as 7 : 3. The resultant dataset posed several challenges including varying illumination changes, shadow effects of dancers on stage, similar dance stances etc. The low accuracy of the skeleton joint coordinates and the partial body parts missing in some sequences makes this dataset very challenging. Tab. 4.1 shows the parameter settings of our proposed model on the evaluated dataset. For training, we utilize Adam optimizer with 0.0004 as learning rate and 0.000001 as decay. The loss function is categorical cross-entropy. Training is done for a maximum of 100 epochs with early stopping if validation error is not improved for consecutive 5 epochs.

4.2 Evaluation Results and Discussion

In order to verify the effectiveness of the proposed network, we compare with other four deep learning-based architectures: Convolutional LSTM (LRCN), CONV3D, Multilayer Perceptron (MLP) and LSTM. We also provide a baseline with the various features from Inception V3 ImageNet features, Kinetics 3D-CNN features and a combination of these with Pose signature using LSTM framework. The LSTM output is forwarded to 2 fully connected layers in all cases. The size of the layer varied as the feature length varied. The first dense layer is of quarter size to the LSTM output. The second dense layer is 1/2 of the first one. This

Table 1: Parameters used for training the LSTM

TRAINING PARAMETERS	
Optimizer	Adam
Loss	Categorical Cross-entropy
Learning rate	0.0001
Decay	0.000001
Feature length	2251
Output length	6
Sequence Length (One training sample)	48 frames (Uniformly Distributed)
Batch Size	32
Maximum Epoch	100

Table 2: Final results

Dance Class	Precision	Recall	F1-score	Support	Class Accuracy
Bharatnatyam	0.54	0.87	0.66	76	86.84
Kathak	0.65	0.88	0.74	58	87.93
Kuchipudi	0.82	0.71	0.76	126	70.63
Manipuri	0.62	0.25	0.35	53	24.52
Mohiniattam	0.97	0.94	0.95	63	93.63
Odissi	0.83	0.64	0.72	69	63.76
Average	0.75	0.72	0.71	445	72.35

is done for both optimality as well as for a fair comparison between all features. The final layer is a dense softmax layer with size 6. Tab. 4.2 provides the classification results of the proposed architecture. Tab. 4.2 shows the confusion matrix for the 6 ICD forms. For training, we utilize Inception V3 pretrained model on ImageNet and extract image features. We randomly sample 48 frames from every video and pass those frames through ImageNet pretrained model to create a stacked feature set for each video. The initial input size to the model is (batch size, sample size, feature size) i.e. (32,48,2048). We add 2 fully connected layers to the output of LSTM and feed it to softmax layer for evaluation. We achieved best results with LSTM framework as compared to other evaluated deep learning architectures. Second best results are obtained using Multi-layer Perceptron approach with reasonable accuracy (40-50%) but lagging behind LSTM (60-70%).

We have also shown the evaluation on various possible permutations of the feature sets to validate the efficacy of the concatenation on improving discriminative ability without hindering the network’s training performance. Tab. 4 provides an evaluation on various possible combination of the feature sets. In order to further improve the feature representation and model the spatio-temporal dependencies, we incorporate Residual Network with multiple cardinalities (also known as ResNext) using 3D CNN architecture. The ResNext-101 was trained on Kinetics dataset, which comprises of 400 video classes. The ResNext architecture utilizes 16 consecutive frames as input and outputs 3D feature vector with size 2048 (or 16x128). We implemented ResNext for all frames in the videos (with sequence length>48). The output is then saved in a JSON format with segment length 16 per video. We convert this data into a feature vector of 128 dimensionality for the input images and then selected 48 frames again uniformly for training. The training input is of size (32,48,128). The training is computationally less expensive due to reduced feature vector. The results were better than InceptionV3 which utilised only spatial information, but since ResNext captures flow i.e. temporal information,

the results were better even with reduced feature length.

Next, we utilize AlphaPose framework to compute pose of the dancing person. The pose constitutes 16 anchor joints of different body parts. We then utilize these 16 anchor joints to construct a pose signature capturing the temporal flow and flow directionality between successive frames and embed the geometric constraints with the normalized distances between various body anchor joints. The resultant pose signature is of dimensionality $75 - D$ feature vector. LSTM achieved the best results which is comparable to the result using ResNext-101 3D CNN feature used independently. Since both the architectures reflect temporal information, the results differ by close margins.

After training these features individually, we concatenate these features to construct a novel feature set for training via LSTM networks. The InceptionV3 combined with kinetics gave slightly better results ($\sim 1 - 2\%$) than combined with pose signature. The combined feature descriptor (2251 length) gives the best performance.

5. CONCLUSION

We have presented a novel pose signature in a sequential learning framework for Indian Classical Dance (ICD) classification. We incorporated pose, flow and spatio-temporal dependencies in the pose signature to capture the adjacency relationship between anchor joints in the skeletal pattern of the dancer. We performed exhaustive experiments and demonstrated the effectiveness of the proposed methodology on dance classification. We showed that deep descriptors with handcrafted pose signature outperformed on ICD dataset. We also showed that due to high similarities between dance moves and dressing attires it is highly challenging to classify dance sequences. In future works, we plan to incorporate facial gestures into the classification pipeline.

6. REFERENCES

- [1] A. F. Bobick and J. W. Davis. The recognition of human movement using temporal templates. *IEEE Transactions on pattern analysis and machine intelligence*, 23(3):257–267, 2001.
- [2] M. Breconzio, S. Gong, and T. Xiang. Recognising action as clouds of space-time interest points. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 1948–1955. IEEE, 2009.
- [3] A. Bulat and G. Tzimiropoulos. Human pose estimation via convolutional part heatmap regression. In *European Conference on Computer Vision*, pages 717–732. Springer, 2016.
- [4] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015.
- [5] G. Chéron, I. Laptev, and C. Schmid. P-cnn: Pose-based cnn features for action recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 3218–3226, 2015.
- [6] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE*

Table 3: Comparison of various features and their combinations for Dance Classification

Method	Bharatnatyam	Kathak	Kuchipudi	Manipuri	Mohiniattam	Odissi	Average Accuracy
InceptionV3	80.48	62.71	28.90	30.64	98.41	53.62	59.1
Pose Signature	88.15	79.31	43.65	41.50	96.82	71.01	67.41
Kinetics	84.21	68.96	56.34	30.18	96.82	57.97	65.61
InceptionV3+ Pose Signature	51.31	67.24	65.87	67.92	93.65	73.91	68.98
InceptionV3+ Kinetics	85.52	68.96	63.49	24.53	96.82	57.97	67.19
InceptionV3+ Pose Signature+ Kinetics	86.84	87.93	70.63	24.52	93.65	63.76	72.35

- Computer Society Conference on*, volume 1, pages 886–893. IEEE, 2005.
- [7] N. Dalal, B. Triggs, and C. Schmid. Human detection using oriented histograms of flow and appearance. In *European conference on computer vision*, pages 428–441. Springer, 2006.
- [8] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. Ieee, 2009.
- [9] S. Hochreiter, Y. Bengio, P. Frasconi, J. Schmidhuber, et al. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies, 2001.
- [10] S. Ji, W. Xu, M. Yang, and K. Yu. 3d convolutional neural networks for human action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35(1):221–231, 2013.
- [11] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014.
- [12] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- [13] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [14] I. Laptev. On space-time interest points. *International journal of computer vision*, 64(2-3):107–123, 2005.
- [15] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [16] S. Mahendran, H. Ali, and R. Vidal. A mixed classification-regression framework for 3d pose estimation from 2d images. *arXiv preprint arXiv:1805.03225*, 2018.
- [17] A. Newell, K. Yang, and J. Deng. Stacked hourglass networks for human pose estimation. In *European Conference on Computer Vision*, pages 483–499. Springer, 2016.
- [18] G. Rogez, P. Weinzaepfel, and C. Schmid. Lcr-net: Localization-classification-regression for human pose. In *CVPR 2017-IEEE Conference on Computer Vision & Pattern Recognition*, 2017.
- [19] S. Samanta, P. Purkait, and B. Chanda. Indian classical dance classification by learning dance pose bases. 2012.
- [20] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *Advances in neural information processing systems*, pages 568–576, 2014.
- [21] H. Wang and C. Schmid. Action recognition with improved trajectories. In *Proceedings of the IEEE international conference on computer vision*, pages 3551–3558, 2013.
- [22] L. Wang, Y. Qiao, and X. Tang. Action recognition and detection by combining motion and appearance features. *THUMOS14 Action Recognition Challenge*, 1(2):2, 2014.
- [23] L. Wang, Y. Qiao, and X. Tang. Action recognition with trajectory-pooled deep-convolutional descriptors. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4305–4314, 2015.
- [24] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *European Conference on Computer Vision*, pages 20–36. Springer, 2016.
- [25] W. Zhu, C. Lan, J. Xing, W. Zeng, Y. Li, L. Shen, X. Xie, et al. Co-occurrence feature learning for skeleton based action recognition using regularized deep lstm networks. In *AAAI*, volume 2, page 6, 2016.
- [26] M. Zolfaghari, G. L. Oliveira, N. Sedaghat, and T. Brox. Chained multi-stream networks exploiting pose, motion, and appearance for action classification and detection. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 2923–2932. IEEE, 2017.