

DYNAMIC TEXTURE RECOGNITION USING PDV HASHING AND DICTIONARY LEARNING ON MULTI-SCALE VOLUME LOCAL BINARY PATTERN

Ruxin Ding, Jianfeng Ren, Heng Yu, Jiawei Li

School of Computer Science
University of Nottingham Ningbo China
199 Taikang East Road, Ningbo, 315100 China

ABSTRACT

Spatial-temporal local binary pattern (STLBP) has been widely used in dynamic texture recognition. STLBP often encounters the high-dimension problem as its dimension increases exponentially, so that STLBP could only utilize a small neighborhood. To tackle this problem, we propose a method for dynamic texture recognition using PDV hashing and dictionary learning on multi-scale volume local binary pattern (PHD-MVLBP). Instead of forming very high-dimensional LBP-histogram features, it first uses hash functions to map the pixel difference vectors (PDVs) to binary vectors, then forms a dictionary using the derived binary vector, and encodes them using the derived dictionary. In such a way, the PDVs are mapped to feature vectors of the size of the dictionary, instead of LBP histograms of very high dimension. Such an encoding scheme could extract the discriminant information from videos in a much larger neighborhood effectively. The experimental results on two widely-used dynamic textures datasets, DynTex++ and UCLA, show the superior performance of the proposed approach over the state-of-the-art methods.

Index Terms— Dynamic texture recognition, Volume LBP, Hashing, Dictionary learning, Multi-scale LBP

1. INTRODUCTION

Dynamic textures (DTs) refer to sequences of the image that consists of repeated patterns related to time and space. DT has been widely used in various applications such as video retrieval [1], fire detection [2], and micro-expression analysis [3]. Compared to static textures, DT classification poses increased challenges because their appearance, organization and motion information are not static but time-varying [4]. Consequently, a descriptive feature representation is a key to the success of dynamic texture recognition.

Many methods have been developed for dynamic texture recognition, *e.g.*, geometric properties computation [5], local spatio-temporal filtering [6], dictionary learning [7], deep convolutional neural network [8, 9], and various spatial-temporal local binary patterns [10–15]. Among these, STLBP

is most widely used in DT recognition. Zhao et al. developed Volume Local Binary Pattern (VLBP) [10], which combines the motion and appearance features together, instead of analyzing each frame individually. However, the feature dimension of VLBP increases exponentially with the number of neighbors, which prevents VLBP from utilizing the information in a large neighborhood. To reduce the feature dimension, LBP-TOP [10] was developed to extract LBP features in three orthogonal planes. Inspired by LBP-TOP, various improved STLBP methods have been developed, *e.g.*, MBSIF-TOP [11], LPQ-TOP [12], and ASF-TOP [16].

LBP-TOP and its variants partially address the high-dimension problem of VLBP, but the resulting dimension may be still high, *e.g.*, the dimension of LBP-TOP is $3 \times 2^{P^2-1}$ for a VLBP neighborhood of size $P \times P \times P$. In addition, as LBP features are extracted independently from the three orthogonal planes, the correlation information among the three planes is lost, which leads to possible performance degradation.

To address these challenges, we propose a method using PDV hashing and dictionary learning on multi-scale VLBP (PHD-MVLBP). Firstly, the pixel difference vector (PDV) of neighbors w.r.t. the center pixel is mapped to a binary vector using hash functions. Both the hash functions and binary vectors are jointly optimized so that the resulting binary vectors are evenly distributed. Human knowledge on image/video such as the uniformity of LBP codes [17] is embedded in the optimization function. Then, clustering is performed on the binary codes to construct a dictionary, and a histogram feature for each video clip is derived using the learned codebook. Multi-scale histogram features are extracted for neighborhoods of different sizes. In such a way, the learned feature representations aggregate the discriminant information using all neighbors in the LBP neighborhoods, at different scales, from all videos in the dataset.

The proposed PHD-MVLBP is compared with the state-of-the-art methods on DynTex++ and UCLA datasets for dynamic texture recognition. It consistently outperforms all the compared methods.

This work was supported in part by the National Natural Science Foundation of China under Grant 72071116, and in part by the Ningbo Municipal Bureau Science and Technology under Grants 2019B10026.

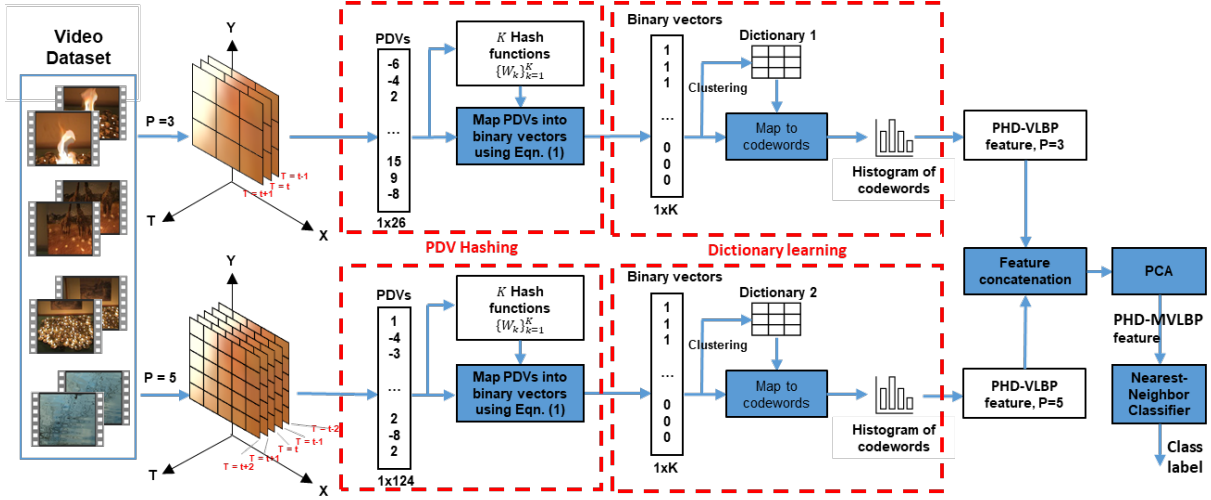


Fig. 1. Overview of the proposed method. Pixel difference vectors (PDVs) are first generated by comparing neighbors with the center pixel. The PDVs are then mapped to binary vectors using hash functions, where the binary vectors and hash functions are jointly optimized. The generated binary vectors are then clustered to form a codebook, and each vector is encoded using the derived dictionary. Finally, the nearest-neighbor classifier with cosine distance is used to classify dynamic textures.

2. PROPOSED METHOD

2.1. Problem Analysis of VLBP and LBP-TOP

Local Binary Pattern is a popular descriptor to represent the local texture of an image [17–21]. To extend it to the spatial-temporal textures, volume local binary pattern (VLBP), was introduced by Zhao and Pietikainen [10]. The captured LBP features are extended from the two-dimensional plane to three-dimensional space, to capture the dynamic texture information. VLBP works well when analyzing video data since it combines appearance and motion information [22]. But VLBP suffers from the problem of high dimensionality. Given a volume neighborhood of $P \times P \times P$, its feature dimension is as high as 2^{P^3-1} . To tackle this problem, LBP-TOP [10] was developed, but its dimension is still as high as $3 \times 2^{P^2-1}$. Several variants such as MBSIF-TOP [11], LPQ-TOP [12], and ASF-TOP [16] may have similar high-dimension issues. As a result, P is often limited to 3 in practice, which limits the power of LBP feature descriptors. Furthermore, as LBP features are extracted in three orthogonal planes independently, the correlation information among these three planes is lost, which may lead to a degradation in classification performance, as evidenced later in experiments.

2.2. Overview of Proposed PHD-MVLBP

The proposed PHD-MVLBP aims to tackle the challenges of VLBP and LBP-TOP, *i.e.*, effectively extracting the discriminant information from a volume of neighborhoods, without exponentially increasing the feature dimension and losing the discriminant information by splitting into three orthogonal

planes. Towards this end, instead of constructing the LBP histogram features of high dimension, the whole pipeline is redesigned, as shown in Fig. 1. Firstly, pixel difference vectors (PDV) are extracted from local volume neighborhood by comparing neighbors with the center pixel. Instead of directly thresholding the PDVs into binary vectors in traditional LBP [10], PDVs are mapped to binary vectors using hash functions, and binary vectors and hash functions are jointly optimized. The optimization functions are carefully designed so that human knowledge on LBPs are well incorporated into the formulation. Secondly, the derived binary vectors are clustered to form the dictionary and each vector is encoded using the derived dictionary. Thirdly, multi-scale features are extracted using the neighborhood of different sizes. Finally, a simple nearest-neighbor classifier with cosine distance is used to evaluate different LBP feature descriptors.

2.3. Mapping PDVs into Binary Vectors Using Hash Functions

Firstly, PDVs are extracted by comparing neighbors with the center pixel within a local neighborhood. Formally, give the neighborhood of size $P \times P \times P$, the center pixel is I_c and the neighboring pixels are $I_1, I_2, \dots, I_{P^3-1}$, respectively, the pixel difference vector $\mathbf{x} = [I_1, I_2, \dots, I_{P^3-1}] - I_c \in \mathcal{R}^{P^3-1}$. Iterating this process for all neighborhoods in all videos could generate the PDVs $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$, where \mathbf{x}_n is the n -th pixel difference vector and N is the number of PDVs in the dataset. PDVs are capable of encoding significant micro-patterns, such as edges and lines.

Secondly, we aim to map the PDVs into binary vectors using hash functions. Inspired by Duan et al. [21], K hash

functions are used to map each \mathbf{x}_n into a binary vector $\mathbf{b}_n = [b_{1n}, \dots, b_{kn}]^T \in \{0, 1\}^{K \times 1}$. The k -th binary code b_{kn} of \mathbf{x}_n can be computed as,

$$b_{kn} = 0.5 \times (\text{sgn}(\mathbf{w}_k^T \mathbf{x}_n) + 1), \quad (1)$$

where $\mathbf{w}_k \in \mathcal{R}^{P^3-1}$ is the projection vector for the k -th function, $\text{sgn}(v)$ equals to 1 if $v \geq 0$ and -1 otherwise.

Based on the concept of uniformity in LBP [17], the number of transitions between code 0 and 1 should be minimized, and hence the adjacent bits should be as equal as possible in the generated binary vector. However, such restraint could compel the learned binary code to all-zeros or all-ones, reducing the discriminant power of binary codes. To address this issue, we limit the sum of bitwise 0 or 1 switches in each binary code. Even if small alternatives occur in the original code \mathbf{X} , the learned binary vector could still be stable with this restriction. The objective function of feature representation is generated as follow,

$$\begin{aligned} \min_{\mathbf{w}_k} J &= J_1 + \lambda_1 J_2 + \lambda_2 J_3 + \lambda_3 J_4 \quad (2) \\ &= \sum_{n=1}^N \left\| \sum_{k=1}^{K-1} \|b_{kn} - b_{(k+1)n}\|^2 - 1 \right\|^2 \\ &+ \lambda_1 \sum_{n=1}^N \sum_{k=1}^K \|(b_{kn} - 0.5) - \mathbf{w}_k^T \mathbf{x}_n\|^2 \\ &+ \lambda_2 \sum_{k=1}^K \left\| \sum_{n=1}^N (b_{kn} - 0.5) \right\|^2 \\ &- \lambda_3 \sum_{n=1}^N \sum_{k=1}^K \|b_{kn} - \mu_k\|^2 \end{aligned}$$

where μ_k is the mean of the k -th bit of all N PDVs, and λ_1 , λ_2 and λ_3 are three parameters to balance the weight of different terms. In our experiments, we set $\lambda_1 = 1000$, $\lambda_2 = 100$ and $\lambda_3 = 1000000$ empirically. The minimization of the first term, J_1 , attempts to make the adjacent bits of the learned binary codes as equal as possible and prevents all zeros or ones from showing up in the codes. By doing so, the learned codes could be more robust to noise. J_2 is designed to minimize the loss of energy during the projection process via reducing the quantization loss. In order to make more information to be present, the goal for J_3 is to evenly distribute each feature bit in the learned binary code. Furthermore, J_4 could maximize the variance of binary codes to improve the independency of projection vectors, so that the generated binary vectors could carry as diversified information as possible. To find the most appropriate $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$, a gradient descent method is applied with the curvilinear search algorithm. More specifically, the objective function in Eqn. (2) is optimized iteratively, by fixing b_{kn} while optimizing \mathbf{w}_k and by fixing \mathbf{w}_k while optimizing b_{kn} . In the beginning, $\{\mathbf{w}_k\}_{k=1}^K$ is initialized as the top K eigenvectors of $\mathbf{X}\mathbf{X}^T$. The algorithm terminates after a fixed number of iterations.

2.4. Dictionary Learning for Binary Vectors

After deriving the binary vectors, they are clustered to form the dictionary of D codewords. Then, each vector is mapped to the nearest codeword. The histogram of codewords is used as the feature representation. In such a way, instead of being mapped to high-dimensional LBP-histogram features, PDVs are mapped to binary vectors using hash functions first and then mapped to histogram of codewords using the derived dictionary. Finally, principal component analysis (PCA) is applied to further compress the features and reduce the feature dimension. The derived features are less sensitive to illumination variations and local alterations with stronger discriminative power.

2.5. Multi-scale Feature Extraction

The proposed PHD-VLBP provides an effective scheme to encode the PDVs of a large neighborhood. To make full use of the discriminant information embedded in the dynamic texture videos, we propose a multi-scale scheme. We extract features using neighborhoods of different P . The PDVs of different scales are extracted and mapped to binary vectors of each scale correspondingly. Then, a dictionary is learned for each scale, and features extracted at different scales are concatenated. PCA is applied to the combined features to derive the final multi-scale features.

3. EXPERIMENTAL RESULTS

The proposed method is compared with the state-of-the-art methods on DynTex++ and UCLA datasets for dynamic texture recognition. The evaluation results are given in the following subsections.

3.1. DynTex++ dataset

The DynTex++ dataset [23] consists of a set of DT videos, sampling from videos of the original DynTex dataset [24]. It categorizes the dynamic textures into 36 classes, where each contains 100 videos. The size of the video texture is $50 \times 50 \times 50$. For a fair comparison to existing methods, the experimental settings in [25] are followed, where half of the dataset is randomly selected as the training set and the rest for testing. The experiment is repeated 5 times and the average result is reported. VLBP [10] and LBP-TOP [10] are closely related to the proposed method and chosen as the baseline methods. MBSIF-TOP [25] achieves previous best results and CVLBC [26] was published in a reputed journal recently, and hence they are chosen for comparison as well.

Two scales, $P = 3$ and $P = 5$, are used for our method, corresponding to a neighborhood of size $3 \times 3 \times 3$ and $5 \times 5 \times 5$, respectively. We report the results for these two scales and PHD-MVLBP. The dictionary size and PCA dimension are empirically set to 1500 and 500. The comparison results

Table 1. Comparison results between the proposed method and other approaches on the DynTex++ dataset.

Method	Accuracy
Distance Learning [23]	63.70%
DFA [27]	89.90%
VLBP [10]	87.35%
LBP-TOP [10]	93.20%
CVLBC [26]	91.31%
MBSIF-TOP [25]	97.17%
Proposed PHD-VLBP, $P = 3$	97.51%
Proposed PHD-VLBP, $P = 5$	97.10%
Proposed PHD-MVLBP	97.77%

are summarized in Table 1. Compared to VLBP [10] and LBP-TOP [10], the performance gain of the proposed PHD-MVLBP is 10.42% and 4.57%, respectively, which demonstrates the effectiveness of the proposed method in extracting the discriminant information in videos. Compared to the previous best method, MBSIF-TOP [25], the performance gain is 0.6%. The proposed PHD-MVLBP also outperforms two single-scale PHD-VLBP methods.

3.2. UCLA Dataset

The UCLA dataset [28] is a popular dataset for dynamic texture recognition. It consists of 200 DT sequences in total, including 50 scenes, with 4 sequences for each scene. The example dynamic textures are waterfalls, plants, swaying flowers, fire, boiling water, and fountains. For each sequence, there are 75 frames with 160×110 pixels. In our experiment, the dataset is cropped to contain the most motion among all videos with the size of $75 \times 48 \times 48$. Because of the insufficient data in the UCLA dataset to train a robust projection matrix and dictionary, these two trained on the DynTex++ dataset, a much larger dataset, are used on the UCLA dataset.

3.2.1. 50-Class Breakdown

4-fold cross-validation is used, same as in [11, 23, 25]. As shown in Table 2, many existing methods achieve almost perfect classification results on the UCLA-50 dataset. The proposed methods, PHD-VLBP for $P = 3$ and $P = 5$, and PHD-MVLBP, all achieve the perfect classification accuracy and outperform all the compared methods.

3.2.2. 9-Class Breakdown

In the UCLA database, each scene is often captured several times. Thus, the data in the UCLA dataset can be categorized into nine classes: 8 videos for boiling water, 8 for fire, 12 for flowers, 20 for fountains, 108 for plants, 12 for sea, 4 for smoke, 12 for water, and 16 for waterfall. Similarly, as in [25], half of the data are randomly chosen as the training,

Table 2. Comparisons between the proposed methods and other approaches on the UCLA dataset with 50-class setting.

Method	Accuracy
Distance Learning [23]	99.0%
KDT-MD [28]	97.5%
CVLBC [26]	99.5%
MBSIF-TOP [25]	99.5%
Proposed PHD-VLBP, $P = 3$	100.0%
Proposed PHD-VLBP, $P = 5$	100.0%
Proposed PHD-MVLBP	100.0%

Table 3. Comparisons between the proposed method and other approaches on the UCLA dataset with 9-class setting.

Method	Recognition accuracy rate
Distance Learning [23]	95.60%
VLBP [10]	96.30%
LBP-TOP [10]	96.00%
KDT-MD [28]	97.50%
MBSIF-TOP [25]	98.75%
Proposed PHD-VLBP, $P = 3$	98.65%
Proposed PHD-VLBP, $P = 5$	98.50%
Proposed PHD-MVLBP	98.90%

and the rest are used for testing. The experiment is repeated 20 times. To make full use of the information embedded in videos, a video is divided into 5 non-overlapping sub-videos with 15 frames in each sub-video. Majority vote is used to combine the classification results of sub-videos. The proposed PHD-VLBP for $P = 3$ achieves an accuracy of 98.65%, which significantly outperforms most of the compared approaches such as VLBP [10] and LBP-TOP [10], and it is comparable to the multi-scale feature descriptor, MBSIF-TOP [25]. The proposed PHD-MVLBP slightly outperforms the previous almost best-performed method, MBSIF-TOP [25], by 0.15%. All these results demonstrate the superior performance of the proposed feature descriptors over the state-of-the-art methods for DT recognition.

4. CONCLUSION

To tackle the problems of high dimensionality of VLBP and the potential information loss of LBP-TOP, we propose PHD-MVLBP for dynamic texture recognition. Firstly, PDVs are extracted and encoded into binary vectors using hash functions, where binary vectors and hash functions are jointly optimized. Then, a dictionary is derived from the binary vectors and used to encode the vectors. Multi-scale features are extracted for the neighborhood of different sizes. The proposed method could effectively extract spatial-temporal discriminant information from videos. It is evaluated on DynTex++ and UCLA datasets and demonstrates superior performance compared with the state-of-the-art methods.

5. REFERENCES

- [1] Micha Haas, Joachim Rijsdam, Bart Thomee, and Michael Lew, “Relevance feedback: perceptual learning and retrieval in bio-computing, photos, and video,” 2004, pp. 151–156.
- [2] G. Zhao, M. Barnard, and M. Pietikainen, “Lipreading with local spatiotemporal descriptors,” *IEEE Transactions on Multimedia*, vol. 11, no. 7, pp. 1254–1265, 2009.
- [3] X. Li, X. Hong, A. Moilanen, X. Huang, T. Pfister, G. Zhao, and M. Pietikainen, “Towards reading hidden emotions: A comparative study of spontaneous micro-expression spotting and recognition methods,” *IEEE Transactions on Affective Computing*, vol. 9, no. 4, pp. 563–577, 2018.
- [4] X. Zhao, Y. Lin, L. Liu, J. Heikkilä, and W. Zheng, “Dynamic texture classification using unsupervised 3D filter learning and local binary encoding,” *IEEE Transactions on Multimedia*, vol. 21, no. 7, pp. 1694–1708, 2019.
- [5] Yuhui Quan, Yuping Sun, and Yong Xu, “Spatiotemporal lacunarity spectrum for dynamic texture classification,” *Computer Vision and Image Understanding*, vol. 165, pp. 85–96, 2017.
- [6] Adin Ramirez Rivera and Oksam Chae, “Spatiotemporal directional number transitional graph for dynamic texture recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 10, pp. 2146–2152, 2015.
- [7] Yuhui Quan, Chenglong Bao, and Hui Ji, “Equiangular kernel dictionary learning with applications to dynamic texture analysis,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 308–316.
- [8] Xianbiao Qi, Chun-Guang Li, Guoying Zhao, Xiaopeng Hong, and Matti Pietikainen, “Dynamic texture and scene classification by transferring deep image features,” *Neurocomputing*, vol. 171, pp. 1230–1241, 2016.
- [9] Shervin Rahimzadeh Arashloo, Mehdi Chehel Amirani, and Ardeshir Noroozi, “Dynamic texture representation using a deep multi-scale convolutional network,” *Journal of Visual Communication and Image Representation*, vol. 43, pp. 89–97, 2017.
- [10] G. Zhao and M. Pietikainen, “Dynamic texture recognition using local binary patterns with an application to facial expressions,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 29, no. 6, pp. 915–928, 2007.
- [11] Shervin Arashloo and Josef Kittler, “Dynamic texture recognition using multiscale binarized statistical image features,” *IEEE Transactions on Multimedia*, vol. 16, no. 8, pp. 2099–2109, 2014.
- [12] Esa Rahtu, Janne Heikkilä, Ville Ojansivu, and Timo Ahonen, “Local phase quantization for blur-insensitive image analysis,” *Image and Vision Computing*, vol. 30, no. 8, pp. 501–512, 2012.
- [13] Jianfeng Ren, Xudong Jiang, and Junsong Yuan, “Dynamic texture recognition using enhanced LBP features,” in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2013, pp. 2400–2404.
- [14] Jianfeng Ren, Xudong Jiang, and Junsong Yuan, “A Chi-squared-transformed subspace of LBP histogram for visual recognition,” *IEEE Transactions on Image Processing*, vol. 24, no. 6, pp. 1893–1904, 2015.
- [15] Jianfeng Ren, Xudong Jiang, and Junsong Yuan, “Learning LBP structure by maximizing the conditional mutual information,” *Pattern Recognition*, vol. 48, no. 10, pp. 3180–3190, 2015.
- [16] Sungeun Hong, Jongbin Ryu, and Hyun S. Yang, “Not all frames are equal: aggregating salient features for dynamic texture classification,” *Multidimensional Systems and Signal Processing*, vol. 29, no. 1, pp. 279–298, 2018.
- [17] Jianfeng Ren, Xudong Jiang, and Junsong Yuan, “Noise-resistant local binary pattern with an embedded error-correction mechanism,” *IEEE Transactions on Image Processing*, vol. 22, no. 10, pp. 4049–4060, 2013.
- [18] T. Ojala, M. Pietikainen, and T. Maenpää, “Multiresolution gray-scale and rotation invariant texture classification with local binary patterns,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 7, pp. 971–987, 2002.
- [19] Jianfeng Ren, Xudong Jiang, and Junsong Yuan, “Quantized fuzzy LBP for face recognition,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 1503–1507.
- [20] Jianfeng Ren, Xudong Jiang, Junsong Yuan, and Nadia Magnenat-Thalmann, “Sound-event classification using robust texture features for robot hearing,” *IEEE Transactions on Multimedia*, vol. 19, no. 3, pp. 447–458, 2017.
- [21] Y. Duan, J. Lu, J. Feng, and J. Zhou, “Context-aware local binary feature learning for face recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 5, pp. 1139–1153, 2018.
- [22] di Huang, Caifeng Shan, Mohsen Ardabilian, and Liming Chen, “Local binary patterns and its application to facial image analysis: A survey,” *IEEE Transactions on Systems, Man, and Cybernetics, Part C*, vol. 41, no. 6, pp. 765–781, 2011.
- [23] Bernard Ghanem and Narendra Ahuja, “Maximum margin distance learning for dynamic texture recognition,” in *Computer Vision – ECCV 2010*, Kostas Daniilidis, Petros Maragos, and Nikos Paragios, Eds., Berlin, Heidelberg, 2010, pp. 223–236, Springer Berlin Heidelberg.
- [24] Renaud Péteri, Sándor Fazekas, and Mark J. Huiskes, “Dyntex: A comprehensive database of dynamic textures,” *Pattern Recognition Letters*, vol. 31, no. 12, pp. 1627–1632, 2010.
- [25] S. R. Arashloo and J. Kittler, “Dynamic texture recognition using multiscale binarized statistical image features,” *IEEE Transactions on Multimedia*, vol. 16, no. 8, pp. 2099–2109, 2014.
- [26] X. Zhao, Y. Lin, and J. Heikkilä, “Dynamic texture recognition using volume local binary count patterns with an application to 2D face spoofing detection,” *IEEE Transactions on Multimedia*, vol. 20, no. 3, pp. 552–566, 2018.
- [27] Yong Xu, Yuhui Quan, Haibin Ling, and Hui Ji, “Dynamic texture classification using dynamic fractal analysis,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2011, pp. 1219–1226.
- [28] A. B. Chan and N. Vasconcelos, “Classifying video with kernel dynamic textures,” in *2007 IEEE Conference on Computer Vision and Pattern Recognition*, 2007, pp. 1–6.