

Cross-Database Micro-Expression Recognition: A Benchmark

Yuan Zong

School of Biological Science and
Medical Engineering, Southeast
University
Nanjing, China
xhzyuan@seu.edu.cn

Wenming Zheng

Key Laboratory of Child Development
and Learning Science of Ministry of
Education, Southeast University
Nanjing, China
wenming_zheng@seu.edu.cn

Xiaopeng Hong

Xi'an Jiaotong University
Xi'an, China
hongxiaopeng@ieee.org

Chuangao Tang

School of Biological Science and
Medical Engineering, Southeast
University
Nanjing, China
tcg2016@seu.edu.cn

Zhen Cui

School of Computer Science and
Engineering, Nanjing University of
Science and Technology
Nanjing, China
zhen.cui@njust.edu.cn

Guoying Zhao

Center for Machine Vision and Signal
Analysis, University of Oulu
Oulu, Finland
guoying.zhao@oulu.fi

ABSTRACT

Cross-database micro-expression recognition (CDMER) is one of recently emerging and interesting problems in micro-expression analysis. CDMER is more challenging than the conventional micro-expression recognition (MER), because the training and testing samples in CDMER come from different micro-expression databases, resulting in inconsistency of the feature distributions between the training and testing sets. In this paper, we contribute to this topic from two aspects. First, we establish a CDMER experimental evaluation protocol and provide a standard platform for evaluating their proposed methods. Second, we conduct extensive benchmark experiments by using NINE state-of-the-art domain adaptation (DA) methods and SIX popular spatiotemporal descriptors for investigating the CDMER problem from two different perspectives and deeply analyze and discuss the experimental results. In addition, all the data and codes involving CDMER in this paper are released on our project website: <http://aip.seu.edu.cn/cdmer>.

CCS CONCEPTS

• **Computing methodologies** → *Biometrics; Computer vision representations.*

KEYWORDS

cross-database micro-expression recognition, micro-expression recognition, domain adaptation, transfer learning, spatiotemporal descriptors.

ACM Reference Format:

Yuan Zong, Wenming Zheng, Xiaopeng Hong, Chuangao Tang, Zhen Cui, and Guoying Zhao. 2019. Cross-Database Micro-Expression Recognition: A Benchmark. In *International Conference on Multimedia Retrieval (ICMR '19)*,

1 INTRODUCTION

Micro-expression is one involuntary facial expression whose duration is usually within 0.5 seconds [7]. Different from ordinary facial expressions, micro-expressions happen as a result of conscious suppression or unconscious repression. They can be viewed as a “leakage” often occurring on someone’s face when that person tries to conceal a genuine emotion. For this reason, understanding micro-expressions has great values in lots of practical applications, *e.g.*, lie detection. Unfortunately, without proper training most people cannot recognize micro-expressions in real time. In order to lower this barrier, researchers from computer vision and affective computing community have been focusing on automatic micro-expression recognition (MER) techniques and proposed various methods [10, 15, 22, 25, 31, 32, 36, 38, 40, 41]. For example, lots of handcrafted spatiotemporal descriptors have been designed to describe micro-expressions, *e.g.*, local binary pattern from three orthogonal planes (LBP-TOP) [32, 44], facial dynamics map (FDM) [41], and fuzzy histogram of optical flow orientation (FHOFO) [10]. Meanwhile, inspired by the success of deep learning techniques in vision tasks, more and more deep neural network based feature learning methods [15, 40] have been also proposed in recent years and performed well on existing micro-expression databases.

Although the above MER methods have shown their promising performance, it should be pointed out that they are still far from the requirements of a real-world MER system. One major reason is that existing MER methods are mostly designed and evaluated without the consideration of the complex scenarios encountered in practice. For example, the training and testing micro-expression samples provided for MER system may be recorded by different cameras (*e.g.*, high-speed camera *v.s.* near-infrared camera). In this scenario, the performance of the above MER methods may sharply drop due to the significantly different feature distributions existing between the training and testing micro-expression samples caused by the heterogeneous video qualities. It thus brings us an emerging topic in micro-expression analysis, *i.e.*, **cross-database micro-expression**

recognition (CDMER), in which the training and testing samples come from two different micro-expression databases. CDMER offers a good way to mimic the scenarios the MER system would encounter in reality. Therefore, it is worthy to deeply investigate this challenging topic.

Prior to the research of MER, the cross-database emotion recognition problems have been extensively studied in many other modalities such as speech emotion recognition (SER) [33, 54], facial expression recognition (FER) [42, 47], image emotion recognition [45, 46], and EEG emotion recognition [21, 48]. However, there are several limitations commonly existing in all cross-database emotion recognition research. First, in cross-database emotion recognition, there is a lack of unified standard evaluation protocol. Researchers often choose their preferred experiment materials including emotion databases, emotion features, classifiers, and evaluation metrics to set up their own experimental evaluation protocol. It raises the barriers of entry to this topic because it would cost lots of time to apply protocols for new researchers who are interested in but not familiar with this topic. Second, it can be found that cross-database emotion recognition was purely viewed as a domain adaptation (DA) [30] task in almost all existing works. It should be noted that cross-database emotion recognition including CDMER is not a simple DA problem which usually aims at making use of the machine learning methods to alleviate the feature distribution difference between training and testing sets. How to design excellent micro-expression features is also an important cue to guide us to solve CDMER, which means designing a robust feature used for describing micro-expressions would improve the recognition performance of classifier in the CDMER tasks as well.

Based on the above considerations, we will investigate the CDMER problem to break through the above two limitations. To this end, we make two contributions in this paper. Firstly, we build a CDMER experimental evaluation protocol and design a set of CDMER experiments on two public available micro-expression databases, which can be used to evaluate the CDMER methods from not only the perspective of DA but also the micro-expression features. Secondly, we use NINE representative DA methods and SIX spatiotemporal descriptors used for describing micro-expressions to conduct extensive experiments under the designed protocol and deeply discuss the experimental results. The major motivation of this work is to attract and encourage more researches to join this challenging but interesting topic and provide standardized protocol for them to get started. For this reason, we released all the data and codes involving CDMER in this paper on our project website: <http://aip.seu.edu.cn/cdmer>.

2 BENCHMARK DETAIL

As described previously, most existing cross-database emotion recognition problems including CDMER are often viewed as a DA task and solved by DA methods. In this way, followed by feature extraction, DA technique is used to relieve the feature distribution mismatch between the source (training) and target (testing) micro-expression samples. Then, we are able to learn a classifier based on the labeled source micro-expression database to predict the micro-expression categories of samples from target database. A picture is drawn to illustrate the detailed process of using DA methods to

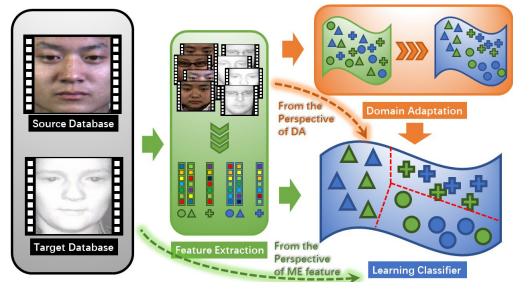


Figure 1: An illustration of how to solve the CDMER problem from the perspectives of DA methods and micro-expression features, respectively.

deal with the CDMER problem, which is shown following the sequence of the orange dash arrow line in Fig. 1. It should be pointed out that besides DA solution, designing robust (database invariant) micro-expression features is also an effective way to solve CDMER problem. By resorting to well-designed robust micro-expression features, CDMER can be actually solved as a typical pattern recognition task which only needs two major steps including feature extraction and classifier learning. We draw its detail in what the green dash arrow directs in Fig. 1. Following the ideas of these two solutions for CDMER, we would like to design a CDMER evaluation protocol which can be used for evaluating both DA methods and micro-expression features, respectively.

2.1 A Standard Evaluation Protocol for CDMER

2.1.1 Data Preparation. Two publicly available spontaneous micro-expression databases are adopted for building the benchmark evaluation experiments, *i.e.*, CASME II [43] and SMIC [20]. CASME II was built by Yan et al. from Institute of Psychology, Chinese Academy of Sciences. It consists of 257 micro-expression samples from 26 subjects. Among these 257 samples, the micro-expression label of one sample is not provided. Each of the other 256 samples is assigned one of seven micro-expression labels including *Happy, Disgust, Repression, Surprise, Sad, Fear, and Others*. Different from CASME II, Li et al. from University of Oulu, Finland considered the image quality diversity of micro-expression samples and hence employed three cameras, *i.e.*, a high-speed (HS) camera, a normal visual (VIS) camera, and a near-infrared (NIR) camera, to collect three subsets to obtain the SMIC (HS, VIS, and NIR) database. The HS subset has 164 samples belonging to 16 subjects and 71 samples of eight subjects from these 16 subjects compose VIS and NIR subsets. All the samples in SMIC are categorized into three types of micro-expressions, *i.e.*, *Positive, Negative, and Surprise*. To make CASME II and SMIC have the same micro-expression labeling, we select the samples of *Happy, Disgust, Surprise, Sad, and Fear* from CASME II and relabel them according to the labeling rule in SMIC, where *Happy* samples are given the *Positive* labels, *Disgust, Sad, and Fear* samples are relabeled with *Negative* micro-expression, and the labels of *Surprise* sample keep unchanged. The sample constitution with respect to consistent categories of the selected CASME II and SMIC databases is shown in Table 1.

Table 1: The Sample Constitution of the selected CASME II and SMIC Databases used in the benchmark.

Micro-Expression Database	Positive	Negative	Surprise	Total
Selected CASME II	32	73	25	130
SMIC (HS)	51	70	43	167
SMIC (VIS)	23	28	20	71
SMIC (NIR)	23	28	20	71

2.1.2 CDMER Tasks. We design two kinds of CDMER tasks based on the selected CASME II and SMIC databases for our CDMER protocol. The first type of tasks denoted by TYPE-I is the one between either two datasets of SMIC (HS, VIS, and NIR). The second type of tasks chooses the selected CASME II and one dataset of SMIC (HS, VIS, and NIR) to serve as source and target micro-expression databases, alternatively, which is denoted by TYPE-II. This leads to totally 12 CDMER experiments, where one CDMER task is denoted by $Exp.i : S \rightarrow T$, where $Exp.i$ is the number of this experiment and S and T are the source and target micro-expression databases, respectively. We summarize all the CDMER experiments in the designed protocol in Table 2.

Table 2: The Detailed Information of Two Types of CDMER Tasks in the Designed Evaluation Protocol.

TYPE	CDMER Task	Source Database	Target Database
TYPE-I	$Exp.1 : H \rightarrow V$	SMIC (HS)	SMIC (VIS)
	$Exp.2 : V \rightarrow H$	SMIC (VIS)	SMIC (HS)
	$Exp.3 : H \rightarrow N$	SMIC (HS)	SMIC (NIR)
	$Exp.4 : N \rightarrow H$	SMIC (NIR)	SMIC (HS)
	$Exp.5 : V \rightarrow N$	SMIC (VIS)	SMIC (NIR)
	$Exp.6 : N \rightarrow V$	SMIC (NIR)	SMIC (VIS)
TYPE-II	$Exp.7 : C \rightarrow H$	Selected CASME II	SMIC (HS)
	$Exp.8 : H \rightarrow C$	SMIC (HS)	Selected CASME II
	$Exp.9 : C \rightarrow V$	Selected CASME II	SMIC (VIS)
	$Exp.10 : V \rightarrow C$	SMIC (VIS)	Selected CASME II
	$Exp.11 : C \rightarrow N$	Selected CASME II	SMIC (NIR)
	$Exp.12 : N \rightarrow C$	SMIC (NIR)	Selected CASME II

2.1.3 Performance Metrics. In the works of [50, 52, 53], weighted average recall (WAR) and unweighted average recall (UAR) are employed to serve as the performance metrics, where WAR is the normal recognition *Accuracy* while UAR is the mean accuracy of each class divided by the number of the classes without the consideration of sample number of each class. The main reason of introducing UAR is due to the class imbalanced problem which widely exists in CASME II and SMIC databases. As shown in Table 1, the number of *Negative* samples in the selected CASME II is 73, which is significantly larger than the numbers of the remaining two types of micro-expression samples (32 for *Positive* and 25 for *Surprise*).

Mean F1-score is another recommended metric, which has been widely used to avoid the bias in performance measurement caused by the class imbalanced problem in MER literatures [16–18, 26, 51]. For this reason, in our benchmark, we adopt the combination of *mean F1-score* and *Accuracy* to serve as the metrics, where *mean F1-score* is the main metric and recognition accuracy as the secondary one. The *mean F1-score* is calculated according to $mean\ F1\text{-score} = \frac{1}{c} \sum_{i=1}^c \frac{2p_i \times r_i}{p_i + r_i}$, where p_i and r_i mean the precision and recall of the i^{th} micro-expression, respectively, and c is

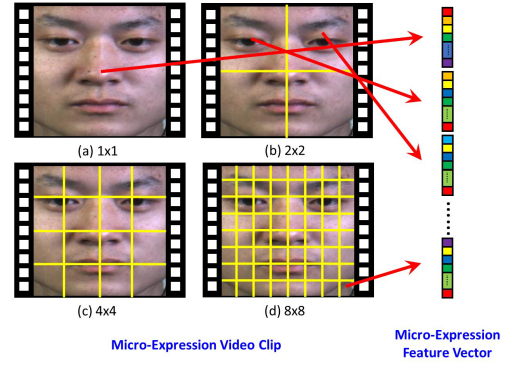


Figure 2: Multi-Scale Grid Based Spatial Division Scheme for Micro-Expression Feature Extraction.

the number of micro-expressions. The *Accuracy* is calculated by $Accuracy = \frac{T}{N} \times 100$, where T and N are the number of correct predictions and the number of target micro-expression samples.

2.1.4 Preprocessing and Feature Extraction. Before feature extraction, preprocessing operations, e.g., face alignment and face cropping, are performed on the micro-expression samples. For convenience, we directly adopted the image sequence data preprocessed by the collectors of CASME II and SMIC for the benchmark evaluation experiments. Then, we employ the temporal interpolation model (TIM) [49] to normalize the frame number of all the micro-expression video clips to 16 and resize each frame image to 112×112 , which allows us to extract spatiotemporal descriptor with specific parameter settings for micro-expression samples. Furthermore, we compute the multi-scale spatiotemporal descriptors using four types of spatial grids (1×1 , 2×2 , 3×3 , and 4×4) shown in Fig. 2 to serve as the micro-expression features. This is expected to extensively cover micro-expression related facial local regions and increase the discriminative power of the extracted spatiotemporal descriptor [51]. As Fig. 2 shows, given a micro-expression sample \mathcal{M} , the spatiotemporal descriptor corresponding to each facial block denoted by \mathbf{x}_i ($i = 1, \dots, K$), where the facial block number $K = 85$, is first extracted one by one and then compose the final micro-expression feature vectors, which can be formulated as $\mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_K]^T$.

As described previously, we attempt to investigate CDMER problem from two different perspectives including domain adaptation (DA) and micro-expression feature extraction. By using effective DA methods and micro-expression features, the large feature distribution difference between the source and target micro-expression databases in CDMER would be relieved. Hence, in our benchmark, we will respectively evaluate the performance of state-of-the-art DA methods and spatiotemporal descriptors used for describing micro-expressions in all the above 12 designed CDMER experiments. Note that in the experiments of using DA methods, we suggest to employ a baseline spatiotemporal descriptor, uniform LBP-TOP [44] with fixed parameters (neighboring radius R and number of the neighboring points P for LBP operator on three orthogonal planes are fixed at 3 and 8, respectively), as shown in Fig. 2 to serve as the micro-expression feature vector \mathbf{x}^v . In the case of using various

features, we extract different \mathbf{x}^v corresponding to different types of spatiotemporal descriptors used for describing micro-expression and then conduct CDMER experiments.

2.1.5 Classifier. To offer a fair comparison, the linear SVM is suggested for serving as the classifier in our benchmark. Without specific description, LibSVM [2] is used in the implementation of SVM for both micro-expression features and DA evaluation experiments.

2.2 Evaluated Methods

2.2.1 DA Methods. For the evaluation of using DA methods to deal with CDMER, the following baseline and state-of-the-art unsupervised DA methods are employed.

Baseline (SVM without DA) [2]: A support vector machine (SVM) without any domain adaptation is served as the baseline method. In the evaluation experiments, we directly learn the linear SVM on the source micro-expression database and then use it to predict the micro-expression labels of samples from the target database.

IW-SVM [11]: Importance-weighted SVM (IW-SVM) was proposed by Hassan et al. to deal with cross-database speech emotion recognition tasks. In this method, a transfer learning method, e.g., unconstrained least-squares importance fitting (uLSIF) [13], is first used to learn a group of importance weights for source samples and then these weights are incorporated into the SVM classifier to eliminate the feature distribution difference between the source and target samples.

TCA [29]: Transfer component analysis (TCA) was proposed by Pan et al., which is to seek some transfer components across domains in a reproducing kernel Hilbert space (RKHS). By using these transfer components, a subspace can be spanned in which the sample distributions from different domains would be close to each other.

GFK [9]: Geodesic flow kernel (GFK) was proposed by Gong et al. and has been a widely-used baseline comparison method in DA research. GFK aims to bridge two domains and narrow their gaps with a well-designed geodesic flow kernel on a Grassmann manifold.

SA [8]: Subspace alignment (SA) is another popular unsupervised DA method and have been served as the baseline comparison method in lots of DA literatures. The target of SA method is to seek a mapping function which can aligns the subspace the source samples lie in with respect to the target ones.

STM [3, 4]: Selective transfer machine (STM) was originally proposed to cope with personalized facial action unit (AU) detection problem. STM makes use of an instance-wise weighted SVM to model the relationship between the training samples and its AU information and meanwhile KMM to eliminate the difference between the AU samples from testing subject and training subjects.

TKL [24]: Long et al. proposed a novel unsupervised DA method called transfer kernel learning (TKL), which aims to learn a domain invariant kernel for eliminating the feature distribution mismatch between the source and target domains.

TSRG [50]: TSRG is short for target sample re-generator and was proposed to deal with CDMER problem. The aim of TSRG is to learn a sample regenerator for the target micro-expression samples and the feature distribution gap between the source and

target micro-expression databases would be narrowed after the regeneration operation.

DRFS-T [53]: TSRG is further extended to a generalized framework called domain regeneration (DR), which inherits the basic idea of TSRG. DRFS-T is one new designed sample regenerator under the DR framework and means domain regeneration in the original feature space with unchanged target samples and it shares the similar idea with TSRG. Their only difference is that DRFS-T keeps the target samples unchanged and regenerates the source samples to have the same or similar feature distributions of target samples, while TSRG is opposite.

DRLS [53]: DRLS is another newly designed sample regenerator based on the DR framework. Different from TSRG and DRFS-T, the subspace used in performing regeneration in DRLS is the label space spanned by the label information provided in the source micro-expression database instead of original feature space.

2.2.2 Micro-Expression Features. For exploring the performance of different existing micro-expression features in coping with CDMER problem, we collect following SIX representative handcrafted spatiotemporal descriptors and ONE deep spatiotemporal descriptor to conduct the benchmark evaluation experiments:

LBP-TOP [44]: LBP-TOP is an spatiotemporal extension of LBP by performing LBP coding on three orthogonal planes. It is originally proposed to deal with dynamic texture recognition tasks and recently has been widely used for describing micro-expressions [22, 32, 38, 51].

LBP-SIP [39]: LBP-SIP is short for LBP with six intersection points. Wang et al. designed it in order to reduce the redundancy in LBP-TOP patterns and provide a more compact and lightweight representation. Compared with LBP-TOP, LBP-SIP uses six intersection points in the intersection lines surrounding the center points for LBP coding and hence its computational complexity is significantly reduced.

LPQ-TOP [28]: Following the manner of LBP-TOP, local phase quantization (LPQ) [27], which quantifies the Fourier transform phase in local neighborhoods, is also extended to the spatiotemporal version called LPQ from three orthogonal planes (LPQ-TOP).

HOG-TOP [19]: Histograms of oriented gradients (HOG) [5] is earliest proposed for human detection and subsequently applied on lots of vision tasks. In the work of [19], Li et al. extends HOG to a 3D version called HOG-TOP, which borrows the basic idea of LBP-TOP.

HIGO-TOP [19]: Histogram of image gradient orientation (HIGO) was proposed in the work of Li et al. [19] by degenerating HOG. Compared with HOG, HIGO simply uses vote rather than weighted vote in counting the responses of the histogram bins. As the name suggests, HIGO-TOP is the 3D extension of HIGO by using the manner of LBP-TOP.

C3D [37]: Recently, the research of spatiotemporal deep feature learning models have made great progress. Three-dimensional convolutional neural network (C3D) is one of excellent representatives and has gained promising performance in video based action recognition tasks. C3D can be actually viewed as a 3D extension of VGG network [34], which replaces 2D convolution and pooling operations with 3D ones.

3 EVALUATION RESULTS

3.1 Implementation Details

We conduct the benchmark CDMER experiments under the designed protocol described in Section 2. For the evaluated methods, we implement them using the original source codes provided by the authors or by ourselves. To offer a fair comparison, for the experiments of different micro-expression features, the parameters of each spatiotemporal descriptors are fixed throughout the experiments, while for the experiments of DA, we follow the strategy used in current mainstream unsupervised DA evaluation experiments that is reporting the best result (in term of *mean F1-score* in our CDMER benchmark experiments) corresponding to the optimal parameter setting for a DA method [1, 6, 23, 24, 50].

3.1.1 Parameter Setting for DA Methods. In this section, we give the parameter searching space for different DA methods in the CDMER evaluation experiments.

SVM [2]: To offer a fair comparison, we use linear kernel for SVM throughout the evaluation experiments. As to the penalized coefficient C , we fix it at $C = 1$. Note that for the experiments of all the DA methods, we use the linear SVM with $C = 1$ to serve as the classifier (if needed).

IW-SVM [11]: In our experiments, we choose uLSIF [13], which has shown its excellent performance in CDMER [50, 53], to learn the importance weights for IW-SVM. Following the suggestion of [50, 53], we search the trade-off parameter λ for uLSIF from a parameter space $[1:1:100] \times t$ ($t = 1, 10, 100, 1000, 10000, 100000$).

TCA [29], **GFK** [9], and **SA** [8]: Among the experiments of these three methods, principal component analysis (PCA) is used to construct the subspace for GFK and SA. For all of them, we search the optimal dimension k (the number of eigenvectors for composing the projection matrix) by trying all possible dimensions, i.e., searching $k \in [1, 2, \dots, k_{max}]$.

STM [3, 4]: STM is originally a binary classification model. In our experiments, we extend it to a multi-class version by using one-against-rest strategy. Similar with SVM used in the benchmark evaluation, its penalized coefficient is set as $C = 1$. As to its second trade-off parameter λ , which is used to balance the KMM regularization term with the SVM objective function, the searching space is set as $[0.001 : 0.001 : 0.009, 0.01 : 0.01 : 0.09, 0.1 : 0.1 : 1, 2 : 1 : 100, 1000, 10000]$.

TKL [24]: According to the work of [24], TKL has one important parameter called the eigenspectrum damping factor ζ . In the evaluation experiments, we determine its optimal value by searching from the parameter space $[0.1 : 0.1 : 5]$.

TSRG [50], **DRFS-T** [53], and **DRLS** [53]: TSRG, DRFS-T, and DRLS have two important trade-off parameters, i.e., λ and μ . Following the works of [50, 53], the optimal values of these two parameters are determined by searching from $[0.001, 0.01, 0.1, 1, 10, 100, 1000]$ for λ and $[0.001 : 0.001 : 0.009, 0.01 : 0.01 : 0.09, 0.1 : 0.1 : 1, 2 : 1 : 10]$ for μ .

3.1.2 Parameter Setting for Micro-Expression Features. We set the parameters for the evaluated spatiotemporal descriptors as follows:

LBP-TOP [44]: LBP-TOP has two important parameters. One is the neighboring radius R and the other is the number of the neighboring points P . In the evaluation experiments, we set the

values of these two parameters as the ones of $R \in \{1, 3\}$ and $P \in \{4, 8\}$, respectively, and hence we report FOUR experimental results of LBP-TOP with different parameter setting, i.e., LBP-TOP_(R1P4), LBP-TOP_(R1P8), LBP-TOP_(R3P4), and LBP-TOP_(R3P8). In addition, the uniform pattern is used for LBP coding. We use Hong et al.'s fast LBP-TOP [12] source code to implement the LBP-TOP feature extraction.

LBP-SIP [39]: Only one parameter needs to be set for LBP-SIP, i.e., the neighboring radius R . Similar to LBP-TOP, we set R as 1 and 3 for LBP-SIP, respectively, and report the experimental results of LBP-SIP_(R1) and LBP-SIP_(R3). In the experiments, LBP-SIP is implemented by ourselves.

LPQ-TOP [28]: For LPQ-TOP, we set its parameters following the suggestion of [28]. Specifically, the size of the local window in each dimension is set as the default one ([5, 5, 5]). The parameters $[\rho_s, \rho_t]$ for a correlation model used in LPQ-TOP are set as $[0.1, 0.1]$ and $[0, 0]$ (without correlation), respectively. The results of LPQ-TOP_{decorr=0.1} and LPQ-TOP_{decorr=0} are reported in the evaluation experiments.

HOG-TOP and **HIGO-TOP** [19]: HOG-TOP and HIGO-TOP both have one important parameter, i.e., the number of bins p to be set, which controls the dimensional histogram of image gradient orientations for three orthogonal planes. In the evaluation experiments, we fix p at 4 and 8 for HOG-TOP and HIGO-TOP and they are denoted by HOG-TOP_{p=4}, HOG-TOP_{p=8}, HIGO-TOP_{p=4} and HIGO-TOP_{p=8}, respectively.

C3D [37]: In the evaluation experiments, we choose the publicly released C3D pretrained on Sports-1M [14] and UCF101 [35] as the feature extractor and use the last two fully connected layers to serve as the micro-expression features, which are denoted by C3D-FC1 (Sports-1M) and C3D-FC2 (Sports-1M), respectively.

3.2 Results and Discussions

3.2.1 Results at A Glance. In this section, we report the benchmark results of the evaluated methods including various DA methods and micro-expression features. All the experimental results are depicted in Tables 3, 4, 5, and 6. Before deeply comparing and analyzing the results obtained by different DA methods and micro-expressions, we would like to probe into the CDMER tasks based on the obtained results. We calculate the average results of all the methods in each experiment and all the experiments for each method, which are given in the last line and last column of these tables. From the comparison between them, we are able to reach the following conclusions.

Firstly, we observe that for our designed benchmark, the second type of experiments (TYPE-II) are significantly more difficult than TYPE-I, which can be clearly revealed by the remarkable differences between the average results of each method in these two types of experiments. For example, the baseline method (SVM without DA), GFK (a well-performing representative of DA methods), and LPQ-TOP_{decorr=0.1} (a well-performing representative of micro-expression features) achieve the average *mean F1-score / Accuracy* of 0.6003 / 61.62%, 0.7223 / 72.44%, and 0.6157 / 63.79% in the first type of experiments, which are much higher than their achieved results (0.4112 / 45.55%, 0.5161 / 54.31%, and 0.3236 / 38.51)

Table 3: Experimental results (mean F1-score / Accuracy) of various domain adaptation methods for CDMER, where the source and target databases are two subsets of SMIC (HS, VIS, and NIR). The micro-expression categories (3 classes) are *Negative*, *Positive*, and *Surprise*. The best results in each experiment are highlighted in bold.

DA Method	Exp.1: H → V	Exp.2: V → H	Exp.3: H → N	Exp.4: N → H	Exp.5: V → N	Exp.6: N → V	Average
Baseline	0.8002 / 80.28	0.5421 / 54.27	0.5455 / 53.52	0.4878 / 54.88	0.6186 / 63.38	0.6078 / 63.38	0.6003 / 61.62
IW-SVM [11]	0.8868 / 88.73	0.5852 / 58.54	0.7469 / 74.65	0.5427 / 54.27	0.6620 / 69.01	0.7228 / 73.24	0.6911 / 68.07
TCA [29]	0.8269 / 83.10	0.5477 / 54.88	0.5828 / 59.15	0.5443 / 57.32	0.5810 / 61.97	0.6598 / 67.61	0.6238 / 64.01
GFK [9]	0.8448 / 84.51	0.5957 / 59.15	0.6977 / 70.42	0.6197 / 62.80	0.7619 / 76.06	0.8142 / 81.69	0.7223 / 72.44
SA [8]	0.8037 / 80.28	0.5955 / 59.15	0.7465 / 74.65	0.5644 / 56.10	0.7004 / 71.83	0.7394 / 74.65	0.6917 / 69.44
STM [3, 4]	0.8253 / 83.10	0.5059 / 51.22	0.6628 / 66.20	0.5351 / 56.10	0.6427 / 67.61	0.6922 / 70.42	0.6440 / 65.78
TKL [24]	0.7742 / 77.46	0.5738 / 57.32	0.7051 / 70.42	0.6116 / 62.20	0.7558 / 76.06	0.7579 / 76.06	0.6964 / 69.92
TSRG [50]	0.8869 / 88.73	0.5652 / 56.71	0.6484 / 64.79	0.5770 / 57.93	0.7056 / 70.42	0.8116 / 81.69	0.6991 / 70.05
DRFS-T [53]	0.8643 / 85.92	0.5767 / 57.32	0.7179 / 71.83	0.6163 / 61.59	0.7286 / 73.24	0.7732 / 77.46	0.7128 / 71.23
DRLS [53]	0.8604 / 85.92	0.6120 / 60.98	0.6599 / 66.20	0.5599 / 55.49	0.6620 / 69.01	0.5771 / 61.97	0.6552 / 66.60
Average	0.8373 / 83.80	0.5700 / 56.96	0.6714 / 67.18	0.5658 / 57.25	0.6819 / 69.86	0.7156 / 72.82	-

Table 4: Experimental results (mean F1-score / Accuracy) of various domain adaptation methods for CDMER, where the source and target databases are CASME II or one subset of SMIC (HS, VIS, and NIR). The micro-expression categories (3 classes) are *Negative*, *Positive*, and *Surprise*. The best results in each experiment are highlighted in bold.

DA Method	Exp.7: C → H	Exp.8: H → C	Exp.9: C → V	Exp.10: V → C	Exp.11: C → N	Exp.12: N → C	Average
Baseline	0.3697 / 45.12	0.3245 / 48.46	0.4701 / 50.70	0.5367 / 53.08	0.5295 / 52.11	0.2368 / 23.85	0.4112 / 45.55
IW-SVM [11]	0.3541 / 41.46	0.5829 / 62.31	0.5778 / 59.15	0.5537 / 54.62	0.5117 / 50.70	0.3456 / 36.15	0.4876 / 50.73
TCA [29]	0.4637 / 46.34	0.4870 / 53.08	0.6834 / 69.01	0.5789 / 59.23	0.4992 / 50.70	0.3937 / 42.31	0.5177 / 53.45
GFK [9]	0.4126 / 46.95	0.4776 / 50.77	0.6361 / 66.20	0.6056 / 61.50	0.5180 / 53.52	0.4469 / 46.92	0.5161 / 54.31
SA [8]	0.4302 / 47.56	0.5447 / 62.31	0.5939 / 59.15	0.5243 / 51.54	0.4738 / 47.89	0.3592 / 36.92	0.4877 / 50.90
STM [3, 4]	0.3640 / 43.90	0.6115 / 63.85	0.4051 / 52.11	0.2715 / 30.00	0.3523 / 42.25	0.3850 / 41.54	0.3982 / 45.61
TKL [24]	0.3829 / 44.51	0.4661 / 54.62	0.6042 / 60.56	0.5378 / 53.08	0.5392 / 54.93	0.4248 / 43.85	0.4925 / 51.93
TSRG [50]	0.5042 / 51.83	0.5171 / 60.77	0.5935 / 59.15	0.6208 / 63.08	0.5624 / 56.34	0.4105 / 46.15	0.5348 / 56.22
DRFS-T [53]	0.4524 / 46.95	0.5460 / 60.00	0.6217 / 63.38	0.6762 / 68.46	0.5369 / 56.34	0.4653 / 50.77	0.5498 / 57.65
DRLS [53]	0.4924 / 53.05	0.5267 / 59.23	0.5757 / 57.75	0.5942 / 60.00	0.4885 / 49.83	0.3838 / 42.37	0.5102 / 53.71
Average	0.4226 / 46.77	0.5084 / 57.54	0.5761 / 59.71	0.5500 / 55.46	0.5011 / 51.46	0.3852 / 41.08	-

in the second type of experiments. In fact, the differences in difficulty between the TYPE-I and TYPE-II experiments stand to reason. The three datasets (HS, VIS, and NIR) used in TYPE-I involve the same subjects, stimulus materials, and recording environments and just different cameras, which results in the relatively small dataset difference. However, another dataset used in TYPE-II experiments, CASME II, corresponds to substantially different subjects, stimulus material, and recording environments compared with SMIC (HS, VIS, and NIR).

Secondly, it should be pointed out that the class imbalanced problem existing in the source or target database remarkably degrades the performance of either DA method or micro-expression features in dealing with the CDMER tasks. For example, in the cases with SMIC (NIR) as the target database, *i.e.*, Exp.11, Exp.3, and Exp.5, we can observe that the average performance in terms of *mean F1-score* and *Accuracy* of all the DA methods can reach 0.6881 and 70.42% in Exp.5 whose the source database, SMIC (VIS), is relatively class-balanced. These two metrics drop to 0.6782 and 67.86% in Exp.3 and 0.5011 and 51.39% in Exp.11, where the source databases of Exp.3 and Exp.11 are SMIC (HS) and CASME II, respectively and very class-imbalanced. Similarly, the performance of all the methods is also affected by the class-imbalanced target database. As shown in Exp.6, Exp.3, and Exp.12 whose source database is fixed, *i.e.*, SMIC (NIR), it can be seen that with class-imbalanced database as target one, the average *mean F1-score / Accuracy* decrease from the level of 0.7272 / 73.88% (Exp.6: class-balanced) to

0.5668 / 57.26% (Exp.4: class-imbalanced) and 0.3928 / 41.96% (Exp.4: class-imbalanced), respectively.

Thirdly, we can also observe that the heterogeneous problem existing between source and target databases raises the level of difficulty of the CDMER tasks. It is known that the samples in SMIC (NIR) are recorded by a near-infrared camera, whose image-quality is considerably different from the samples recorded by high-speed camera (used in CASME II and SMIC (HS)) and visual camera (used in SMIC (VIS)). Therefore, it is intuitive that the CDMER tasks involving SMIC (NIR) would be more difficult than others. In order to check this point, we first fix the source database as SMIC (HS) and observe two experiments, *i.e.*, Exp.1 and Exp.3, where the target database in Exp.1 is SMIC (VIS) and Exp.3 corresponds to SMIC (NIR). It can be seen that the average results among all the DA methods and micro-expression features (0.8405 / 84.12% and 0.6459 / 67.52%) in Exp.1 (homogeneous) are significantly higher than the results (0.6782 / 67.86% and 0.4420 / 49.65%) in Exp.3 (heterogeneous). We further check the opposite case if the heterogeneous image-quality samples exist in the source database. We observe Exp.2 and Exp.4, whose target databases are the same, *i.e.*, SMIC (HS) and source databases are different, *i.e.*, SMIC (VIS) *v.s.* SMIC (NIR). From the results, we notice the performance difference between the Exp.2 (homogeneous case) and Exp.4 (heterogeneous case). Specifically, the average performance achieved by DA methods and micro-expressions are 0.5764 / 57.60% and 0.4318 / 45.35% in Exp.2 (homogeneous case) *v.s.* 0.5668 / 57.26% and 0.3616 / 39.18% in Exp.4 (heterogeneous case).

Table 5: Experimental results (mean F1-score / Accuracy) of various domain adaptation methods for CDMER, where the source and target databases are two subsets of SMIC (HS, VIS, and NIR). The micro-expression categories (3 classes) are *Negative*, *Positive*, and *Surprise*. The best results in each experiment are highlighted in bold.

Spatiotemporal Descriptors	Exp.1: H \rightarrow V	Exp.2: V \rightarrow H	Exp.3: H \rightarrow N	Exp.4: N \rightarrow H	Exp.5: V \rightarrow N	Exp.6: N \rightarrow V	Average
Baseline	0.8002 / 80.28	0.5421 / 54.27	0.5455 / 53.52	0.4878 / 54.88	0.6186 / 63.38	0.6078 / 63.38	0.6003 / 61.62
LBP-TOP _(R1P4) [44]	0.7185 / 71.83	0.3366 / 40.24	0.4969 / 49.30	0.3457 / 40.24	0.5480 / 57.75	0.5085 / 59.15	0.4924 / 53.32
LBP-TOP _(R1P8) [44]	0.8561 / 85.92	0.5329 / 53.66	0.5164 / 57.75	0.3246 / 35.37	0.5124 / 57.75	0.4481 / 50.70	0.5318 / 56.86
LBP-TOP _(R3P4) [44]	0.4656 / 49.30	0.4122 / 45.12	0.3682 / 40.85	0.3396 / 40.85	0.5069 / 59.15	0.5144 / 60.56	0.4345 / 49.31
LBP-SIP _(R1) [39]	0.6290 / 63.38	0.3447 / 40.85	0.3249 / 33.80	0.3490 / 42.07	0.5477 / 60.56	0.5509 / 60.56	0.4577 / 50.20
LBP-SIP _(R3) [39]	0.8574 / 85.92	0.4886 / 50.00	0.4977 / 54.93	0.4038 / 42.68	0.5444 / 59.15	0.3994 / 46.48	0.5319 / 56.53
LPQ-TOP _(decorr=0.1) [28]	0.9455 / 94.37	0.5523 / 54.88	0.5456 / 61.97	0.4729 / 47.56	0.5416 / 57.75	0.6365 / 66.20	0.6157 / 63.79
LPQ-TOP _(decorr=0) [28]	0.7711 / 77.46	0.4726 / 48.78	0.6771 / 67.61	0.4701 / 48.17	0.7076 / 71.83	0.6963 / 70.42	0.6325 / 64.05
HOG-TOP _(p=4) [19]	0.7068 / 71.83	0.5649 / 57.32	0.6977 / 70.42	0.2830 / 29.27	0.4569 / 49.30	0.3218 / 36.62	0.4554 / 48.47
HOG-TOP _(p=8) [19]	0.7364 / 74.65	0.5526 / 56.10	0.3990 / 46.48	0.2941 / 32.32	0.4137 / 46.48	0.3245 / 38.03	0.4453 / 49.01
HIGO-TOP _(p=4) [19]	0.7933 / 80.28	0.4775 / 50.61	0.4023 / 47.89	0.3445 / 35.98	0.5000 / 53.52	0.3747 / 40.85	0.4821 / 51.52
HIGO-TOP _(p=8) [19]	0.8445 / 84.51	0.5186 / 53.66	0.4793 / 54.93	0.4322 / 43.90	0.5054 / 54.93	0.4056 / 46.48	0.5309 / 56.40
C3D-FC1 (Sports1M) [37]	0.1577 / 30.99	0.2188 / 23.78	0.1667 / 30.99	0.3119 / 34.15	0.3802 / 49.30	0.3032 / 36.62	0.2564 / 34.31
C3D-FC2 (Sports1M) [37]	0.2555 / 36.62	0.2974 / 29.27	0.2804 / 33.80	0.3239 / 36.59	0.4518 / 47.89	0.3620 / 38.03	0.3285 / 37.03
C3D-FC1 (UCF101) [37]	0.3803 / 46.48	0.3134 / 34.76	0.2697 / 47.89	0.3440 / 34.76	0.3916 / 47.89	0.2433 / 29.58	0.3404 / 40.23
C3D-FC2 (UCF101) [37]	0.4162 / 46.48	0.2842 / 32.32	0.3053 / 42.25	0.2531 / 28.05	0.3937 / 47.89	0.2489 / 32.39	0.3169 / 38.23
Average	0.6459 / 67.52	0.4318 / 45.35	0.4420 / 49.65	0.3613 / 39.18	0.5013 / 55.28	0.4341 / 48.50	-

Table 6: Experimental results (mean F1-score / Accuracy) of various domain adaptation methods for CDMER, where the source and target databases are CASME II or one subset subset of SMIC (HS, VIS, and NIR). The micro-expression categories (3 classes) are *Negative*, *Positive*, and *Surprise*. The best results in each experiment are highlighted in bold.

Spatiotemporal Descriptors	Exp.7: C \rightarrow H	Exp.8: H \rightarrow C	Exp.9: C \rightarrow V	Exp.10: V \rightarrow C	Exp.11: C \rightarrow N	Exp.12: N \rightarrow C	Average
Baseline	0.3697 / 45.12	0.3245 / 48.46	0.4701 / 50.70	0.5367 / 53.08	0.5295 / 52.11	0.2368 / 23.85	0.4112 / 45.55
LBP-TOP _(R1P4) [44]	0.3358 / 44.51	0.3260 / 47.69	0.2111 / 35.21	0.1902 / 26.92	0.3810 / 43.66	0.2492 / 26.92	0.2823 / 37.49
LBP-TOP _(R1P8) [44]	0.3680 / 43.90	0.3339 / 54.62	0.4624 / 49.30	0.5880 / 57.69	0.3000 / 33.80	0.1927 / 23.08	0.3742 / 43.73
LBP-TOP _(R3P4) [44]	0.3117 / 43.90	0.3436 / 44.62	0.2723 / 39.44	0.2356 / 28.46	0.3818 / 49.30	0.2332 / 25.38	0.2964 / 38.52
LBP-SIP _(R1) [39]	0.3580 / 45.12	0.3039 / 44.62	0.2537 / 38.03	0.1991 / 26.92	0.3610 / 46.48	0.2194 / 26.92	0.2825 / 38.02
LBP-SIP _(R3) [39]	0.3772 / 42.68	0.3742 / 56.15	0.5846 / 59.15	0.6065 / 60.00	0.3469 / 35.21	0.2790 / 27.69	0.4279 / 46.81
LPQ-TOP _(decorr=0.1) [28]	0.3060 / 42.07	0.3852 / 48.46	0.2525 / 33.80	0.4866 / 47.69	0.3020 / 35.21	0.2094 / 23.85	0.3236 / 38.51
LPQ-TOP _(decorr=0) [28]	0.2368 / 43.90	0.2890 / 51.54	0.2531 / 38.03	0.3947 / 40.77	0.2369 / 35.21	0.4008 / 41.54	0.3019 / 41.83
HOG-TOP _(p=4) [19]	0.3156 / 34.76	0.3502 / 47.69	0.3266 / 35.21	0.4658 / 46.92	0.3219 / 35.21	0.2163 / 27.46	0.3327 / 37.91
HOG-TOP _(p=8) [19]	0.3992 / 43.90	0.4154 / 52.31	0.4403 / 45.07	0.4678 / 47.69	0.4107 / 40.85	0.1390 / 20.77	0.3787 / 41.77
HIGO-TOP _(p=4) [19]	0.2945 / 39.63	0.3420 / 53.85	0.3236 / 40.85	0.5590 / 55.38	0.2887 / 29.58	0.2668 / 31.54	0.3458 / 41.81
HIGO-TOP _(p=8) [19]	0.2978 / 41.46	0.3609 / 50.00	0.3679 / 43.66	0.5699 / 54.62	0.3395 / 33.80	0.1743 / 22.31	0.3517 / 40.98
C3D-FC1 (Sports1M) [37]	0.1994 / 42.68	0.2394 / 56.15	0.1631 / 32.39	0.1075 / 19.23	0.1631 / 32.39	0.2397 / 56.15	0.1854 / 39.83
C3D-FC2 (Sports1M) [37]	0.1994 / 42.68	0.1317 / 24.62	0.1631 / 32.39	0.1075 / 19.23	0.1631 / 32.39	0.2397 / 56.15	0.1674 / 34.58
C3D-FC1 (UCF101) [37]	0.1581 / 31.10	0.1075 / 19.23	0.1886 / 39.44	0.1075 / 19.23	0.1886 / 39.44	0.2397 / 56.15	0.1650 / 34.10
C3D-FC2 (UCF101) [37]	0.1994 / 42.68	0.1705 / 19.23	0.1631 / 32.39	0.1075 / 19.23	0.1631 / 32.39	0.1075 / 19.23	0.1414 / 27.53
Average	0.2954 / 41.88	0.2959 / 44.95	0.3060 / 40.32	0.3581 / 38.94	0.3049 / 37.94	0.2277 / 31.80	-

3.2.2 Results of CDMER Experiments by Using DA Methods. Tables 3 and 4 show the experimental results of different DA methods corresponding to the TYPE-I and TYPE-II experiments, respectively. From these two tables, it can be seen that in both two types of designed CDMER experiments, nearly all the DA methods can achieve promisingly better results in terms of both mean F1-score and Accuracy than the baseline method (SVM without any DA). More importantly, some well-performing DA methods, e.g., GFK [9], DRFS-T [53], and the proposed RSTR, have significant improvements of at least 0.1000 (average *mean F1-score*) and 10% (average *Accuracy*) compared with the baseline results in either TYPE-I or TYPE-II experiments. Based on the above observations, we are able to reach the conclusion that considering CDMER as an DA problem is no doubt an effective solution for the CDMER problem. It is a good choice to develop excellent DA methods to relieve the feature distribution mismatch between the samples from different micro-expression databases.

3.2.3 Results of CDMER Experiments by Using Micro-Expression Features. The detailed *mean F1-score / Accuracy* achieved by various micro-expression features are depicted in Tables 5 and 6. From these two tables, it is clear to see that LPQ-TOP with *decorr* = 0 and LBP-SIP with *R* = 3 achieve the best average performance in terms of *mean F1-score* and *Accuracy* in TYPE-I and TYPE-II CDMER experiments, respectively. More importantly, we notice that LPQ-TOP_(decorr=0) obtains the average *mean F1-score / Accuracy* of 0.6325 / 64.05% in the first type of experiments, which are even competitive among the results of DA methods. By referring Table 3 and Tables 5, we can find that LPQ-TOP_(decorr=0) outperforms TCA (0.6238 / 64.01%) and is very competitive against STM (0.6440 / 65.78%) and DRLS (0.6552 / 66.60%). In addition, we also observe that several micro-expression features achieve very promising results. For example, LPQ-TOP_(decorr=0.1) achieves the *mean F1-score* of 0.9455 and *Accuracy* of 94.37% in Exp.1, which are significantly better than all the DA methods. In Exp.9 and Exp.10, LBP-SIP_(R=3)

Table 7: Experimental results (mean F1-score / Accuracy) of finetuned C3D features.

Method	Exp.1: H → V	Exp.3: H → N	Exp.8: H → C
Baseline	0.8002 / 80.28	0.5455 / 53.52	0.3245 / 48.46
TCA	0.8269 / 83.10	0.5828 / 59.15	0.4870 / 53.08
C3D-FC1 (Sports1M)	0.1577 / 30.99	0.1667 / 30.99	0.2394 / 56.15
C3D-FC2 (Sports1M)	0.2555 / 36.62	0.2804 / 33.80	0.1317 / 24.62
C3D-FC1 (Finetune)	0.4858 / 56.34	0.2751 / 35.21	0.3276 / 34.62
C3D-FC2 (Finetune)	0.4144 / 47.89	0.3758 / 40.85	0.3390 / 36.92

is a very good competitor against DA methods. Based on the above observations, it is believed that developing database-invariant spatiotemporal descriptors for robustly describing micro-expressions also provide a promising and feasible way to solve the CDMER problem.

However, we have to admit that most of the existing spatiotemporal descriptors including the above well-performing ones are still not satisfactory and cannot completely meet the requirement in CDMER problem. More specifically, several limitations still exist in current micro-expression feature research. Firstly, the performance of most spatiotemporal descriptors is not stable. In other words, one spatiotemporal descriptor performs well in Task A but very poor in Task B. An example is LPQ-TOP. It can be seen that LPQ-TOP with two different parameters both achieve satisfactory results beating the baseline method. But the performance of both two LBP-TOP based micro-expression features decrease sharply and even under the level of baseline method. Secondly, it is clear to see that the performance of most of handcrafted spatiotemporal descriptors is strongly affected by its parameters. As the results of LBP-TOP showed, its performance varies very sensitively in nearly all the experiments of our CDMER benchmark with the changes of its parameter R and P . Therefore, in the future study of spatiotemporal descriptors used for describing micro-expressions, we should also consider to reduce the sensitiveness of its parameters such that it will be simpler and more convenient to use in dealing with the CDMER problem. Finally, it should be pointed out that at present it is not enough for solving CDMER problem to simply use existing spatiotemporal descriptors as micro-expression features because its average performance is still far from the DA methods. From another point of view, it is believed that there is still very large development space in this direction, *i.e.*, developing robust (database-invariant) micro-expression features.

Lastly, we discuss the deep features evaluated in the benchmark. From the results of deep features extracted by C3D pretrained on Sports1M and UCF101, we can observe that these four deep features perform poorly in nearly all the CDMER experiments and cannot reach the level of most handcrafted features. In fact, this is foreseeable because the Sports1M and UCF101 served for C3D pretraining are both action databases, whose samples are quite different from the micro-expression ones. In fact, to enable the C3D to gain the micro-expression information, we can finetune the pretrained models based on the source database and then use the finetuned model to serve as feature extractor. To this end, we finetune the C3D pretrained on Sports1M based on SMIC (HS) and then extract the deep features to conduct the experiments of using SMIC (HS) as source database including Exp.1, Exp.3, and Exp.8.

Note that for finetuning, we augment the samples in SMIC (HS) to 3350 and divide the samples into a training set whose sample number is 2550 and a validation set containing 800 samples. Finally, we randomly extract one frame from all the parts to obtain a micro-expression sample with frame length of 16. The experimental results are shown in Table 7, where we also list the results of C3D pretrained on Sports1M, the baseline method (LBP-TOP with $R = 3, P = 8$), and TCA (a representative DA method). From Table 7, it is clear that by resorting to finetuning strategy, the performance of C3D features can be improved significantly compared with the original model. It is also interesting to see that in Exp.8, the deep features extracted by finetuned C3D even outperform the baseline method (LBP-TOP with $R = 3, P = 8$) in term of mean F1-score. However, compared with TCA (a representative DA method) and the baseline method, the performance of deep features is still very poor and not satisfactory. We think there may be two possible reasons. Firstly, it may attribute to the problem of lacking enough good data for finetuning the C3D models. It is clear that the training or finetuning of deep learning models requires large numbers of samples while the sample numbers of the CASME II and SMIC are too small. Although we can use some methods to augment the samples, the augmented micro-expression samples are not satisfactory. Therefore, more micro-expression samples need be collected such that a better-performing C3D can be finetuned to deal with CDMER problem. The second reason may be that it seems not enough for learning the database-invariant deep features to simply finetune the deep model based on the source database. Leveraging the idea of domain adaptation methods, *i.e.*, reducing the difference between source and target domains, to finetune deep models may offer a feasible solution to improve the performance of deep features.

4 CONCLUSION

In this paper, we have investigated and expanded the research on cross-database micro-expression recognition (CDMER) by conducting a standard benchmark evaluation from two different perspectives including domain adaptation (DA) and spatiotemporal features used for describing micro-expressions. First, under a well-designed evaluation protocol, we make use of two widely-used micro-expression databases, *i.e.*, CASME II and SMIC, to set up two types of CDMER experiments. Then, we perform benchmark evaluation experiments by adopting NINE state-of-the-art DA methods and SIX excellent spatiotemporal descriptors under the designed CDMER protocol. Finally, comprehensive discussions for the experimental results are provided. In the future, we will evaluate more well-performing DA methods and micro-expression features for updating the benchmark evaluation results in this paper.

5 ACKNOWLEDGMENTS

This work is supported by the National Key R&D Program of China under Grant 2018YFB1305200, the National Natural Science Foundation of China under Grant 61572009 and Grant 61772419, the Fundamental Research Funds for the Central Universities under Grant 2242018K3DN01 and Grant 2242019K40047, the Tencent AI Lab Rhino-Bird Focused Research Program under Grant JR201922, Academy of Finland, Tekes Fidipro Program, and Infotech Oulu.

REFERENCES

- [1] Maruan Al-Shedivat, Jim Jing-Yan Wang, Majed Alzahrani, Jianhua Z Huang, and Xin Gao. 2014. Supervised transfer sparse coding. In *Proceedings of the Twenty-Eighth AAAI conference on artificial intelligence*. The AAAI Press.
- [2] Chih-Chung Chang and Chih-Jen Lin. 2011. LIBSVM: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)* 2, 3 (2011), 27.
- [3] Wen-Sheng Chu, Fernando De la Torre, and Jeffery F Cohn. 2013. Selective transfer machine for personalized facial action unit detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3515–3522.
- [4] Wen-Sheng Chu, Fernando De la Torre, and Jeffrey F Cohn. 2017. Selective transfer machine for personalized facial expression analysis. *IEEE transactions on pattern analysis and machine intelligence* 39, 3 (2017), 529–545.
- [5] Navneet Dalal and Bill Triggs. 2005. Histograms of oriented gradients for human detection. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, Vol. 1. IEEE, 886–893.
- [6] Zhengming Ding, Ming Shao, and Yun Fu. 2014. Latent Low-Rank Transfer Subspace Learning for Missing Modality Recognition. In *AAAI*. 1192–1198.
- [7] Paul Ekman and Wallace V Friesen. 1969. Nonverbal leakage and clues to deception. *Psychiatry* 32, 1 (1969), 88–106.
- [8] Basura Fernando, Amaury Habrard, Marc Sebban, and Tinne Tuytelaars. 2013. Unsupervised visual domain adaptation using subspace alignment. In *Proceedings of the IEEE international conference on computer vision*. 2960–2967.
- [9] Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. 2012. Geodesic flow kernel for unsupervised domain adaptation. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2066–2073.
- [10] SL Happy and Aurobinda Routray. 2017. Fuzzy Histogram of Optical Flow Orientations for Micro-expression Recognition. *IEEE Transactions on Affective Computing* (2017).
- [11] Asif Hassan, Robert Dampier, and Mahesan Niranjan. 2013. On acoustic emotion recognition: compensating for covariate shift. *IEEE Transactions on Audio, Speech, and Language Processing* 21, 7 (2013), 1458–1468.
- [12] Xiaopeng Hong, Yingyue Xu, and Guoying Zhao. 2016. LBP-TOP: a Tensor Unfolding Revisit. In *Asian Conference on Computer Vision*. Springer, 513–527.
- [13] Takafumi Kanamori, Shohei Hido, and Masashi Sugiyama. 2009. A least-squares approach to direct importance estimation. *The Journal of Machine Learning Research* 10 (2009), 1391–1445.
- [14] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. 2014. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. 1725–1732.
- [15] Dae Hoe Kim, Wissam J Baddar, and Yong Man Ro. 2016. Micro-Expression Recognition with Expression-State Constrained Spatio-Temporal Feature Representations. In *Proceedings of the 2016 ACM on Multimedia Conference*. ACM, 382–386.
- [16] Anh Cat Le Ngo, Yee-Hui Oh, Raphael C-W Phan, and John See. 2016. Eulerian emotion magnification for subtle expression recognition. In *Proceedings of the 41st IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1243–1247.
- [17] Anh Cat Le Ngo, Raphael Chung-Wei Phan, and John See. 2014. Spontaneous Subtle Expression Recognition: Imbalanced Databases and Solutions. In *Proceedings of the 12th Asian Conference on Computer Vision (ACCV)*. Springer, 33–48.
- [18] Anh Cat Le Ngo, John See, and C-W Raphael Phan. 2016. Sparsity in Dynamics of Spontaneous Subtle Emotion: Analysis & Application. *IEEE Transactions on Affective Computing* (2016).
- [19] Xiaobai Li, Xiaopeng Hong, Antti Moilanen, Xiaohua Huang, Tomas Pfister, Guoying Zhao, and Matti Pietikäinen. 2017. Towards Reading Hidden Emotions: A Comparative Study of Spontaneous Micro-expression Spotting and Recognition Methods. *IEEE Transactions on Affective Computing* (2017).
- [20] Xiaobai Li, Tomas Pfister, Xiaohua Huang, Guoying Zhao, and Matti Pietikäinen. 2013. A spontaneous micro-expression database: Inducement, collection and baseline. In *Automatic face and gesture recognition (fg), 2013 10th IEEE international conference and workshops on*. IEEE, 1–6.
- [21] Yang Li, Wenming Zheng, Yuan Zong, Zhen Cui, Tong Zhang, and Xiaoyan Zhou. 2018. A Bi-hemisphere Domain Adversarial Neural Network Model for EEG Emotion Recognition. *IEEE Transactions on Affective Computing* (2018).
- [22] Yong-Jin Liu, Jin-Kai Zhang, Wen-Jing Yan, Su-Jing Wang, Guoying Zhao, and Xiaolan Fu. 2016. A Main Directional Mean Optical Flow Feature for Spontaneous Micro-Expression Recognition. *IEEE Transactions on Affective Computing* 7, 4 (2016), 299–310.
- [23] Mingsheng Long, Jianmin Wang, Guiguang Ding, Jiaguang Sun, and Philip S Yu. 2014. Transfer joint matching for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1410–1417.
- [24] Mingsheng Long, Jianmin Wang, Jiaguang Sun, and S Yu Philip. 2015. Domain invariant transfer kernel learning. *IEEE Transactions on Knowledge and Data Engineering* 27, 6 (2015), 1519–1532.
- [25] Ping Lu, Wenming Zheng, Ziyang Wang, Qiang Li, Yuan Zong, Minghai Xin, and Lenan Wu. 2016. Micro-Expression Recognition by Regression Model and Group Sparse Spatio-Temporal Feature Learning. *IEICE TRANSACTIONS on Information and Systems* 99, 6 (2016), 1694–1697.
- [26] Yee-Hui Oh, Anh Cat Le Ngo, John See, Sze-Teng Liong, Raphael C-W Phan, and Huo-Chong Ling. 2015. Monogenic riesz wavelet representation for micro-expression recognition. In *Proceedings of the 20th IEEE International Conference on Digital Signal Processing (DSP)*. IEEE, 1237–1241.
- [27] Ville Ojansivu and Janne Heikkilä. 2008. Blur insensitive texture classification using local phase quantization. In *International conference on image and signal processing*. Springer, 236–243.
- [28] Juhani Päivärinta, Esa Rahtu, and Janne Heikkilä. 2011. Volume local phase quantization for blur-insensitive dynamic texture classification. In *Scandinavian Conference on Image Analysis*. Springer, 360–369.
- [29] Sinno Jialin Pan, Ivor W Tsang, James T Kwok, and Qiang Yang. 2011. Domain adaptation via transfer component analysis. *IEEE Transactions on Neural Networks* 22, 2 (2011), 199–210.
- [30] Sinno Jialin Pan and Qiang Yang. 2010. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering* 22, 10 (2010), 1345–1359.
- [31] Wei Peng, Xiaopeng Hong, Yingyue Xu, and Guoying Zhao. 2019. A Boost in Revealing Subtle Facial Expressions: A Consolidated Eulerian Framework. In *Proceedings of the 14th IEEE International Conference on Automatic Face and Gesture Recognition (FG)*.
- [32] Tomas Pfister, Xiaobai Li, Guoying Zhao, and Matti Pietikäinen. 2011. Recognising spontaneous facial micro-expressions. In *International Conference on Computer Vision*. IEEE, 1449–1456.
- [33] Bjorn Schuller, Bogdan Vlasenko, Florian Eyben, Martin Wollmer, Andre Stuhlsatz, Andreas Wendemuth, and Gerhard Rigoll. 2010. Cross-corpus acoustic emotion recognition: Variances and strategies. *IEEE Transactions on Affective Computing* 1, 2 (2010), 119–131.
- [34] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [35] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. 2012. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402* (2012).
- [36] Lei Tian, Xiaopeng Hong, Chunxiao Fan, Yue Ming, Matti Pietikäinen, and Guoying Zhao. 2018. Sparse Tikhonov-Regularized Hashing for Multi-Modal Learning. In *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 3793–3797.
- [37] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. 2015. Learning spatiotemporal features with 3d convolutional networks. In *Computer Vision (ICCV), 2015 IEEE International Conference on*. IEEE, 4489–4497.
- [38] Su-Jing Wang, Wen-Jing Yan, Xiaobai Li, Guoying Zhao, Chun-Guang Zhou, Xiaolan Fu, Minghao Yang, and Jianhua Tao. 2015. Micro-expression recognition using color spaces. *IEEE Transactions on Image Processing* 24, 12 (2015), 6034–6047.
- [39] Yandan Wang, John See, Raphael C-W Phan, and Yee-Hui Oh. 2014. Lbp with six intersection points: Reducing redundant information in lbp-top for micro-expression recognition. In *Asian Conference on Computer Vision*. Springer, 525–537.
- [40] Zhaoliang Xia, Xiaopeng Hong, Xingyu Gao, Xiaoyi Feng, and Guoying Zhao. 2019. Spatiotemporal recurrent convolutional networks for recognizing spontaneous micro-expressions. *arXiv preprint arXiv:1901.04656* (2019).
- [41] Feng Xu, Junping Zhang, and James Z Wang. 2017. Microexpression identification and categorization using a facial dynamics map. *IEEE Transactions on Affective Computing* 8, 2 (2017), 254–267.
- [42] Keyu Yan, Wenming Zheng, Zhen Cui, and Yuan Zong. 2016. Cross-Database Facial Expression Recognition via Unsupervised Domain Adaptive Dictionary Learning. In *International Conference on Neural Information Processing*. Springer, 427–434.
- [43] Wen-Jing Yan, Xiaobai Li, Su-Jing Wang, Guoying Zhao, Yong-Jin Liu, Yu-Hsin Chen, and Xiaolan Fu. 2014. CASME II: An improved spontaneous micro-expression database and the baseline evaluation. *PLoS one* 9, 1 (2014), e86041.
- [44] Guoying Zhao and Matti Pietikäinen. 2007. Dynamic texture recognition using local binary patterns with an application to facial expressions. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29, 6 (2007), 915–928.
- [45] Sicheng Zhao, Chuang Lin, Pengfei Xu, Sendong Zhao, Yuchen Guo, Ravi Krishna, Guiguang Ding, and Kurt Keutzer. 2019. CycleEmotionGAN: Emotional Semantic Consistency Preserved CycleGAN for Adapting Image Emotions. In *AAAI Conference on Artificial Intelligence*.
- [46] Sicheng Zhao, Xin Zhao, Guiguang Ding, and Kurt Keutzer. 2018. EmotionGAN: unsupervised domain adaptation for learning discrete probability distributions of image emotions. In *ACM International Conference on Multimedia*. 1319–1327.
- [47] Wenming Zheng, Yuan Zong, Xiaoyan Zhou, and Minghai Xin. 2018. Cross-domain color facial expression recognition using transductive transfer subspace learning. *IEEE transactions on Affective Computing* 9, 1 (2018), 21–37.
- [48] Wei-Long Zheng and Bao-Liang Lu. 2016. Personalizing EEG-based affective models with transfer learning. In *Proceedings of the Twenty-Fifth International*

- Joint Conference on Artificial Intelligence*. AAAI Press, 2732–2738.
- [49] Ziheng Zhou, Xiaopeng Hong, Guoying Zhao, and Matti Pietikäinen. 2014. A compact representation of visual speech data using latent variables. *IEEE transactions on pattern analysis and machine intelligence* 36, 1 (2014), 181–187.
- [50] Yuan Zong, Xiaohua Huang, Wenming Zheng, Zhen Cui, and Guoying Zhao. 2017. Learning a Target Sample Re-Generator for Cross-Database Micro-Expression Recognition. In *Proceedings of the 2017 ACM on Multimedia Conference*. ACM, 872–880.
- [51] Yuan Zong, Xiaohua Huang, Wenming Zheng, Zhen Cui, and Guoying Zhao. 2018. Learning from hierarchical spatiotemporal descriptors for micro-expression recognition. *IEEE Transactions on Multimedia* 20, 11 (2018), 3160–3172.
- [52] Yuan Zong, Wenming Zheng, Zhen Cui, Guoying Zhao, and Bin Hu. 2019. Towards Bridging Micro-Expressions from Different Domains. *IEEE Transactions on Cybernetics* (2019).
- [53] Yuan Zong, Wenming Zheng, Xiaohua Huang, Jingang Shi, Zhen Cui, and Guoying Zhao. 2018. Domain Regeneration for Cross-Database Micro-Expression Recognition. *IEEE Transactions on Image Processing* 27, 5 (2018), 2484–2498.
- [54] Yuan Zong, Wenming Zheng, Tong Zhang, and Xiaohua Huang. 2016. Cross-corpus speech emotion recognition based on domain-adaptive least-squares regression. *IEEE signal processing letters* 23, 5 (2016), 585–589.