

# An Overview of Facial Micro-Expression Analysis: Data, Methodology and Challenge

Hong-Xia Xie, Ling Lo, Hong-Han Shuai, and Wen-Huang Cheng, *Senior Member, IEEE*

**Abstract**—Facial micro-expressions indicate brief and subtle facial movements that appear during emotional communication. In comparison to macro-expressions, micro-expressions are more challenging to be analyzed due to the short span of time and the fine-grained changes. In recent years, micro-expression recognition (MER) has drawn much attention because it can benefit a wide range of applications, e.g. police interrogation, clinical diagnosis, depression analysis, and business negotiation. In this survey, we offer a fresh overview to discuss new research directions and challenges these days for MER tasks. For example, we review MER approaches from three novel aspects: macro-to-micro adaptation, recognition based on key apex frames, and recognition based on facial action units. Moreover, to mitigate the problem of limited and biased ME data, synthetic data generation is surveyed for the diversity enrichment of micro-expression data. Since micro-expression spotting can boost micro-expression analysis, the *state-of-the-art* spotting works are also introduced in this paper. At last, we discuss the challenges in MER research and provide potential solutions as well as possible directions for further investigation.

**Index Terms**—Facial Micro-expression, Recognition, Spotting, Action Units, Deep Learning, Survey.

## 1 INTRODUCTION

FACIAL micro-expression (ME) is a result of *conscious suppression* (intentional) or *unconscious repression* (unintentional), which can be viewed as a “leakage” of people’s true feelings [1]. MEs are brief involuntary facial expressions that usually appear when people are trying to conceal their true feelings, especially in high-stake situations. Hence, micro-expression recognition (MER) research enables greater awareness and sensitivity to subtle facial behaviors, and is an important subject for human emotion and affective phenomena understanding, which has been explored by various disciplines such as psychology, sociology, neuroscience, computer vision, etc. Such skills are useful for psychotherapists, interviewers, and anyone working in communications.

### 1.1 The Difference between Macro-expression and Micro-expression

Micro-expressions occur when people are trying to conceal or repress their true feelings. On the contrary, macro-expressions are easy to be perceived in daily interactions. The major difference between macro-expressions and micro-expressions lies in their duration. Although there is no strict rule of the threshold to distinguish one from the

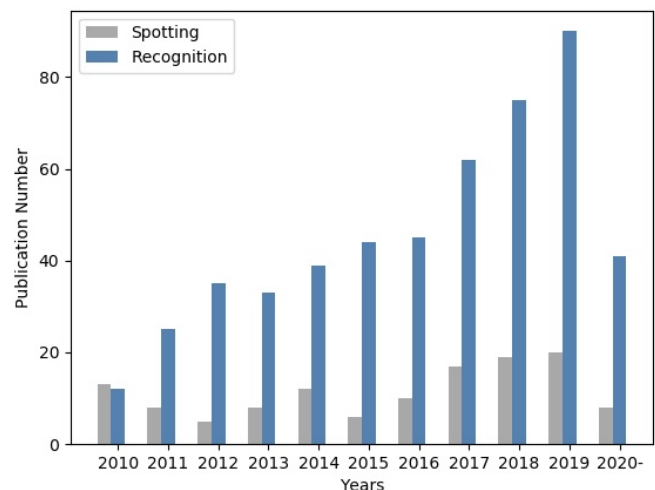


Fig. 1. The publication number of Micro-expression Recognition and Spotting from 2010 to November 2020 (Source: Web of Science).

other, most agreed that macro-expression usually last from 0.5 to 4 seconds while micro-expressions should be no longer than 0.5 seconds [5]. Most researches [6], [7], [8] set 0.5 seconds as the threshold, while 0.2 seconds duration is also regarded as a boundary for differentiating micro- and macro-expressions as validated in [9]. Besides, the neural mechanisms underlying the recognition of micro-expression and macro-expression are different, showing different electroencephalogram (EEG) and event-related potentials (ERPs) characteristics. The brain regions responsible for their differences might be the inferior temporal gyrus and the frontal lobe [9]. General speaking, the macro-expression shows higher intensity and visibility when compared to the micro-expression.

- H.-X. Xie, L. Lo are with the Institute of Electronics, National Chiao Tung University, Hsinchu, 300 Taiwan.  
E-mail: {hongxiaxie.ee08g, lynn97.ee08g}@nctu.edu.tw.
- H.-H. Shuai is with the Department of Electrical and Computer Engineering, National Chiao Tung University, Hsinchu, 300 Taiwan.  
E-mail: hhshuai@nctu.edu.tw.
- W.-H. Cheng is with the Institute of Electronics, National Chiao Tung University, Hsinchu, 300 Taiwan, and the Artificial Intelligence and Data Science Program, National Chung Hsing University, Taichung, 400 Taiwan.  
E-mail: whcheng@nctu.edu.tw.

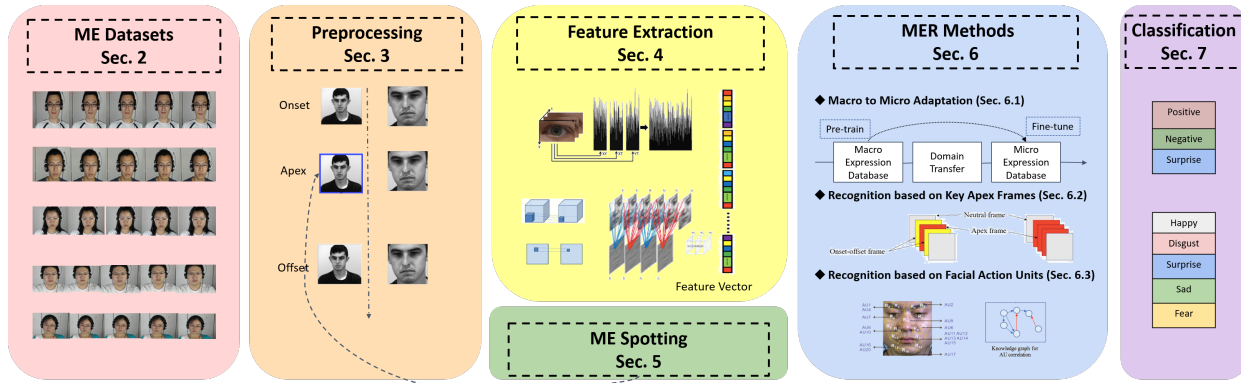


Fig. 2. The organization of this survey is structured according to the general MER pipeline.

TABLE 1  
Academic challenges for micro-expression recognition and spotting.

Challenge	Year	Dataset	Task	Evaluation Metric	Event
MEGC2018 [2]	2018	CASME II, SAMP	Cross-database (HDE, CDE)	UAR, WAR, F1-score	FG <sup>1</sup>
MEGC2019 [3]	2019	CAS(ME) <sup>2</sup> , SAMP	Recognition, Spotting (CDE)	UAR, WAR, F1-score	FG <sup>1</sup>
MEGC2020 [4]	2020	CAS(ME) <sup>2</sup> , SAMP Long Video	Spotting	F1-score	FG <sup>1</sup>
MER2020*	2020	Synthetic Data	Recognition	F1-score	ICIP <sup>2</sup>

\* MER2020: <http://mer2020.tech/>

<sup>1</sup> FG: IEEE International Conference on Automatic Face and Gesture Recognition

<sup>2</sup> ICIP: IEEE International Conference on Image Processing

## 1.2 Facial Action Coding System

Facial Action Coding System (FACS) was first proposed by Ekman and Friesen and updated in 2002 [10], [11] to factorize the composition of micro-expressions. It is the most widely used coding scheme for decomposing facial expressions into individual muscle movements, called Action Units (AUs). With FACS, every possible facial expression can be described as a combination of AUs. There are 32 facial muscle-related actions, and 6 extra unspecific miscellaneous Action Descriptors (ADs) [12].

Yet, the professional training of FACS encoding experts is time-consuming. Professional encoders usually need to receive 100 hours of training, and in practice the encoding process takes 2 hours to encode a 1-minute video on average [13]. Thus, an automatic MER system with high accuracy can be very helpful and valuable.

Since AUs are descriptive for certain facial configurations, specific systems are proposed to explore the relationship between facial muscle movements (AUs) and human emotions, e.g., Emotional Facial Action Coding System (EMFACS-7) [14], Facial Action Coding System Affect Interpretation Dictionary (FACSAID) [15] and the System for Identifying Affect Expressions by Holistic Judgments (Af-fex) [16]. It can be found that various mapping strategies for AUs and emotions are adopted by existing facial expression datasets due to the lack of standard guidelines [17].

## 1.3 Academic Challenges for MER

Academic challenge is a competition created by academic experts for leveraging the power of open innovation to advance the *state-of-the-art* in a particular field. Within the

computer vision domain, several academic challenges related to micro-expression recognition and spotting have been proposed. Some known events in recent years are summarized in Table 1.

## 1.4 Outline of this Paper

Numerous works have been accomplished to mitigate the challenging MER task. The publication counting from 2010 to November 2020 is shown in Fig. 1. Several surveys on MER have been published in recent years [18], [19], [20], [21], [22], [23], [24]. However, most of them have been focused on traditional image-processing methods and a fresh overview to discuss new directions and challenges the MER faces today is necessary. In this survey paper, therefore, we provide a more comprehensive and in-depth study of existing approaches, e.g. including over 100 papers for MER until November 2020 which have not been reviewed in the previous survey papers. Based on the nature of MEs, e.g., the short span of time and the fine-grained changes, we introduce MER approaches from three novel aspects: macro-to-micro adaptation, recognition based on key apex frames, and recognition based on facial action units. Moreover, to mitigate the problem of limited and biased ME data, synthetic data generation algorithms are surveyed for the diversity enrichment of ME data. To the best of our knowledge, this is the first review to summarize the synthetic data generation solutions based on spontaneous ME datasets. Since micro-expression spotting can boost micro-expression analysis, the *state-of-the-art* spotting works are also introduced. At last, we discuss open challenges and provide potential solutions as well as future research directions.

TABLE 2  
A list of public datasets on spontaneous micro-expression.

Dataset	Macro/Micro	Videos	FPS	Resolution	FACS	Emotion	Subjects	AU	Index*
CASME	micro	195	60	640×480,1280×720	Yes	8	19	No	On,Apex,Off
HS		164	100				16		On,Off
SMIC VIS	micro	71	25	640×480	No	3	8	No	N/A
NIR		71	25				8		N/A
CASME II	micro	247	200	640×480	Yes	5	26	Yes	On,Apex,Off
SAMM	micro	159	200	2040×1088	Yes	7	32	Yes	On,Apex,Off
CAS(ME) <sup>2</sup>	macro & micro	300 & 57	30	640×480	No	4	22	N/A	On,Apex,Off
SAMM Long Videos	macro & micro	343 & 159	200	2048×1088	Yes	N/A	30	Yes	On,Apex,Off

\* On, Apex, Off: Onset frame, Apex frame, Offset frame, respectively.

The whole systematic framework of MER is present in Fig. 2. The outline of this paper adheres as follows: Section 2 introduces ME data collection, current public ME datasets, and popular data synthesis methods. Section 3 describes pre-processing techniques commonly used for ME data. Section 4 provides a detailed review of feature extraction methods based on handcrafted features and deep learning-based methods. Discussions of ME spotting are present in Section 5. Section 6 summarizes the commonly-used techniques on MER recently. The loss function widely used for MER is introduced in Section 7. Section 8 presents the overall performance comparison. Open MER challenges and the possible solutions are arranged in Section 9 and some potential future directions for research are suggested in Section 10.

## 2 DATASETS

### 2.1 ME Data Collection

**Emotion Classes.** In discrete emotion theory, there are many different basic emotion definition systems, of which the most popular one adopted in the computer vision community is conducted by Ekman and Friesen [25], where the basic emotions are divided into six categories of anger, disgust, fear, happiness, sadness, and surprise. An extra emotion of contempt is added in later researches.

Current public ME datasets divide micro-expression emotions into different classes according to their collection strategy. The emotion labels should take into account AUs, participants' self-report as well as the stimuli video content, etc. To reduce unpredictability and bias of emotion classes, Spontaneous Micro-expression Corpus (SMIC) [26] group emotion categories from CASME II and SAMM [6] into three main classes (positive, negative, surprise) [27] based on the AUs with FACS coding. Fig. 3 shows sample frames from current ME datasets.

#### 2.1.1 Apparatus Setup

Since the frame rate and resolution could affect the performance of MER [28], existing ME datasets are usually collected in a strictly lab-controlled environment. The observers are kept out of sight to maximise the chances of natural suppression by making participants as comfortable as possible. The participants are asked to watch video clips in front of a screen and avoid any body movement, which can help to exclude habitual movements.

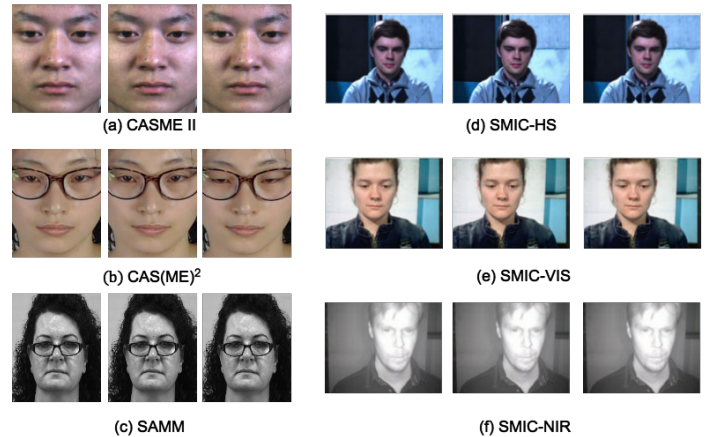


Fig. 3. Sample frames from different ME datasets.

#### 2.1.2 Emotional stimuli

Spontaneous ME datasets were usually elicited by emotional stimuli, e.g., images, movies, music [29]. MEs are more likely to occur under high-arousal stimuli, so video clips with high emotional valence are proved to be effective materials for eliciting MEs [26].

**Emotion Label Rating Quantification.** To explicitly quantify the personal emotional response to the emotional stimulus, a predefined scale that a subject can assign to the perceived emotion response is usually needed. SMIC [26], CASME [7], and CASME II [8] databases assigned an emotion label to stimuli videos based on self-reports completed by participants. Relatively, SAMM adopted Self-Assessment Manikins (SAM) [30], which contains the valence, dominance, and arousal judgments in rating emotional response.

#### 2.1.3 Micro-expression Inducement

**Suppressing.** To evoke MEs, there must be enough pressure to motivate participants to conceal their true feelings. SMIC, CASME, and SAMM asked participants to fully suppress facial movements during the whole experiment so that MEs may occur. For CASME II, half participants were asked to keep neutral faces when watching video clips while other participants were enforced to suppress the facial movements when they realize there is a facial expression.

**Not All Subjects Show Micro-expressions.** According to Ekman's research, when people are telling lies, about half

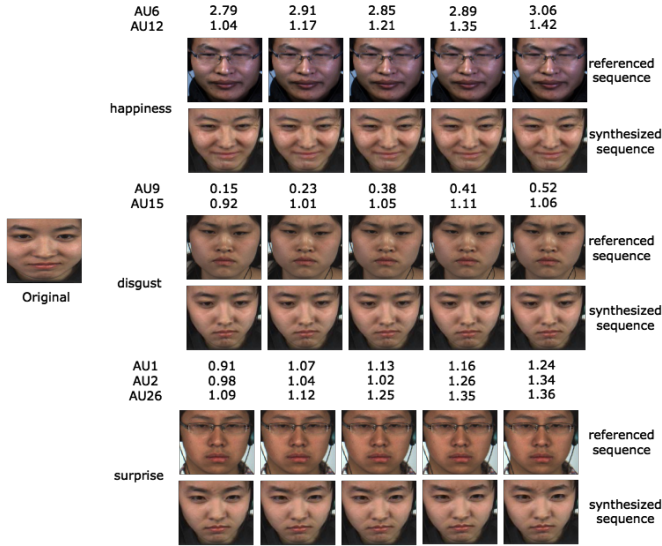


Fig. 4. Synthetic faces generated by AU-ICGAN from CASME II dataset. The number above each image indicates the target AU intensities of chosen AUs [32].

of them might show MEs, while the other half do not [31]. The reason why only some people show MEs is still unclear. In SMIC, four participants did not show any ME at all throughout the course of 35-minutes video watching.

## 2.2 Spontaneous Micro-Expression Dataset

Since posed ME datasets are collected by intentionally controlled, it contradicts the natural occurrence of MEs. In this paper, we only focus on **spontaneous ME datasets**. Table 2 provides an overview of spontaneous ME datasets.

**Chinese Academy of Sciences Micro-Expression (CASME)** [7] dataset consists of 195 video samples from 19 valid subjects with a frame rate of 60 fps. The samples for 8 emotions are highly imbalanced (5 happiness, 6 sadness, 88 disgust, 20 surprise, 3 contempt, 2 fear, 40 repression and 28 tense). Notably, in CASME, each sample is recorded with two different cameras and environmental settings, namely Class A and Class B. Class A was recorded by BenQ M31 camera with natural light, while class B was by GRAS-03K2C camera and two LED lights.

**Chinese Academy of Sciences Micro-Expression II (CASME II)** [8] is an improved version of CASME collected in a well-controlled lab environment. It contains 247 ME sequences from 26 subjects with 7 categories, including happiness, disgust, surprise, fear, sadness, repression and others, which were labeled based on AUs, participants' self-report and the content of stimuli videos. All subjects are Chinese.

**Spontaneous Micro-expression Corpus (SMIC)** [26] is composed of 164 ME sequences from 16 subjects filmed at 100 fps. It is one of the first to include spontaneous MEs through emotional inducement experiments. Most subjects are Asians, and some are from other ethnicity. Ethnicity is more diverse than previous datasets with ten Asians, nine Caucasians and one African participant. For data reliability, ME clips in SMIC are classified into three classes: positive, negative and surprise. While positive (happy) and surprise

only include one target emotion each, negative includes four emotions (sad, anger, fear and disgust).

**Spontaneous Actions and Micro-Movements (SAMM)** [6] has 159 ME sequences from 29 subjects. The seven classes include contempt, disgust, fear, anger, sadness, happiness and surprise. The ethnicities of participants are diverse and the gender split is even.

## 2.3 Macro- and Micro-Expression Datasets

Considering real-world scenarios, macro- and micro-expressions could co-occur, there are two public datasets combining macro and micro facial expressions, i.e., SAMM Long Videos [6] and CAS(ME)<sup>2</sup> [33]. Both of them can be employed for macro-expression and micro-expression spotting from long videos.

**The Chinese Academy of Science Macro- and Micro-expression (CAS(ME)<sup>2</sup>)** dataset was established by the Chinese Academy of Science. The dataset contains 300 macro-expressions and 57 micro-expressions, with four different emotional labels: positive, negative, surprise and others from 22 participants (13 females and 9 males). The expression samples are coded with the onset, apex, and offset frames, with AUs marked and emotions labeled.

**SAMM Long Videos** consists of 147 long videos with 343 macro-expressions and 159 micro-expressions. The frame rate is 200 fps, and 0.5 seconds is set as the threshold for classifying macro- ( $\geq 0.5$  seconds) and micro-expressions ( $< 0.5$  seconds). Emotion labels are not provided.

## 2.4 Synthetic Data Generation

While building ME databases, it is not only challenging to trigger an ME but also very difficult to label one. Labeling ME data requires human labor as well as professional domain knowledge. Even with professional training, it is reported that up to only 47% labeling accuracy can be achieved for a human expert [34]. Thus the size of existing databases is usually very limited. Synthetic datasets are introduced in this situation when annotating the ground-truth is a time-consuming work. In general, synthetic data have been successfully used in many areas in computer vision from learning low-level visual features [35], [36], [37] to high-level tasks [38], [39], [40]. Especially, synthesizing facial data is more challenging than other basic objects as the fidelity of human faces is hard to preserve [41]. Queiroz *et al.* [42] presented a methodology for generation of facial ground truth with synthetic faces. The 3D face model database can control face actions as well as illumination conditions, allowing to generate animation with different facial expressions and eye motions with the ground truth of the facial landmark points provided at each frame. The resulting Virtual Human Faces Database (VHuF) dataset can simulate realistic skin textures extracted from real photos. Abbasnejad *et al.* [43] established a large-scale dataset of facial expression using a 3D face model. It consists of shape and texture models to create different subjects with different expressions. Their experimental results showed that the synthesized dataset enables efficiently deep network training for expression analysis.

Apart from synthetic databases generated by 3D morphable models, an alternative way to produce synthetic data

is using Generative Adversarial Network (GAN). Zhang *et al.* [44] took advantage of GAN and proposed a network exploiting different poses and expressions jointly for facial image synthesis and pose-invariant facial expression recognition. The generated face images with different expressions under arbitrary poses can enlarge and enrich expression dataset and benefit the recognition accuracy. Cai *et al.* [45] synthesized images with a same face in different expressions using a conditional generative model. The resulting dataset consists of sets of images and each image set contains a same identity in different synthetic expressions, which benefits the identity-free expression recognition. Different from large-scale facial expression, Xie *et al.* [32] proposed AU Intensity Controllable Generative Adversarial Networks (AU-ICGAN) for micro-scale expression data synthesis. To enrich the limited data samples, their AU-ICGAN aims to generate face images with specific intensities of action unit to simulate real-world ME data. In addition, the image structure similarity along with image sequence authenticity are taken into consideration so that the generated ME image sequences can be more realistic. Their experimental results show the merits of using an auxiliary synthetic dataset for training deep recognition networks. Fig. 4 shows the synthetic results for enriching ME datasets.

### 3 PRE-PROCESSING

The pre-processing stage in MER consists of all steps required before the extraction of meaningful features can commence. One of the important aim of the pre-processing stage is to detect and align faces into a common reference, so that the features extracted from each face correspond to the same semantic locations. It removes rigid head motion and, to some extent, the anthropomorphic variations among people. After alignment, the subtle micro muscle movements can be further magnified in order to enhance the discriminative characteristic.

#### 3.1 Face Detection

The first step of any face analysis method is to detect the face region. The Viola&Jones (V&J) face detector [53] can give real-time and robust near-frontal face in an image and thus is one of the most popular face detector. Their cascade classifier is based on Haar feature, and its reliability and computational simplicity make it a widely employed approach. Another widely-used face detector is based on Histogram of Gradient (HOG) [54], [55]. HOG-based face detectors share the merit of computational efficiency with V&J face detectors, but they also have same short of being unable to deal with non-frontal face images. With the rapid development of Convolution Neural Network (CNN), more CNN-based face detectors are adopted in popular open source libraries. For example, Max-Margin Object Detection (MMOD) [56] is used in dlib<sup>1</sup>, and a well-trained single-Shot-Multibox detector using ResNet-10 as the backbone is provided in OpenCV<sup>2</sup>. Also, the CNN-based face detectors are reported more accurate and robust under occlusions [57].

1. <http://dlib.net/>  
2. <https://opencv.org/>

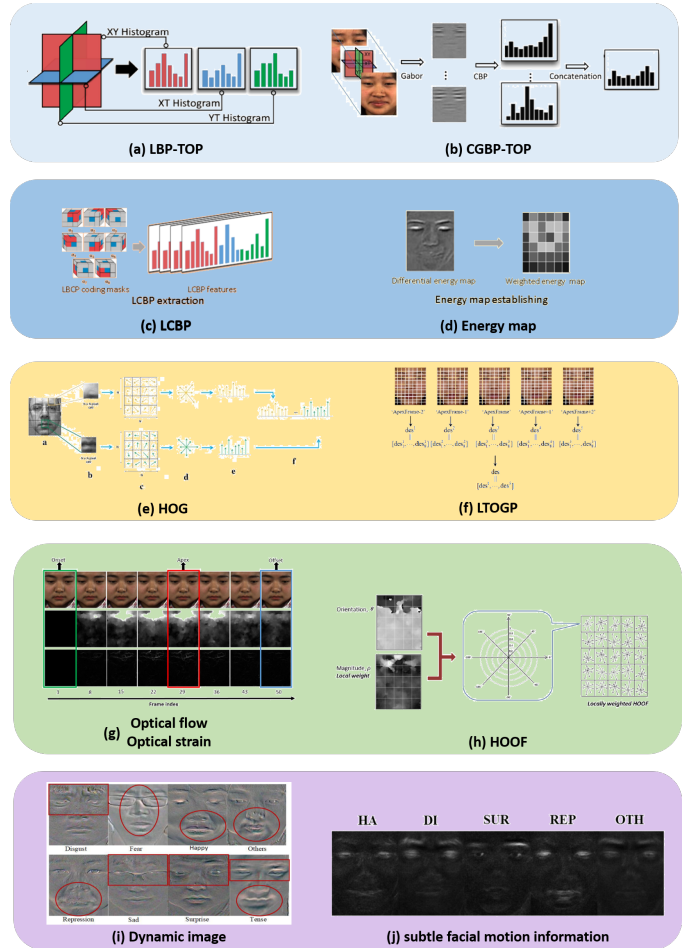


Fig. 5. An illustration of handcrafted features: (a) LBP-TOP [46], (b) CGBP-TOP [47], (c) LCBP [46], (d) energy map [46], (e) HOG [48], (f) LTOGP [49], (g) optical flow (optical strain) [50], (h) HOOF [50], (i) dynamic image [51] and (j) subtle facial motion information [52].

#### 3.2 Facial Landmark Detection

While recognizing face regions in images is the vital beginning, aligning faces in a video is even more crucial to be used in most ME recognition frameworks. To align faces in different frames, the coordinates of localized facial landmarks are first detected. Then, the faces in different frames are aligned into a referenced face according to the location of the landmark keypoints. Proper face alignment can improve the recognition accuracy significantly because the subtle motion of faces over time can be captured. Active Shape Model (ASM) [66] is one of the face models that is frequently used. Discriminative response map fitting (DRMF) [67] is a regression approach with strong ability in general face fitting. With lower computational consumption and real-time capabilities, the DRMF model can handle occlusions, dynamic backgrounds and various illumination conditions.

In recent years, deep-learning methods have been widely used for facial landmarks localization. Cascaded-CNN, which predicts the facial landmarks in a cascaded fashion, has become a *state-of-the-art* method for its high accuracy and speed. Tweaked Convolutional Neural Networks (TCNN) [68] can harness the robustness of CNNs for landmark detection in an appearance-sensitive manner.

TABLE 3  
Performance comparison among micro-expression spotting methods

Methods	Year	Features	Protocol	Datasets	Experimental Results		
					ACC	F1-score	recall
[58]	2020	Specific Pattern	-	SAMM CAS(ME) <sup>2</sup>	0.090 0.029	0.133 0.055	0.258 0.456
[59]	2019	HOOF&RNN	LOSO	SAMM	0.045	0.082	0.465
[60]	2019	LTP-ML	LOSO	SAMM CAS(ME) <sup>2</sup>	0.017 0.009	0.032 0.018	0.296 0.281
[61]	2019	HIGO-TOP & HOG-TOP	LOSO	SMIC-VIS-E CASME II	- -	0.620 0.860	- -
[62]	2019	HOG	LOSO	CASME II	0.823	-	-
[63]	2019	time-constrasted feature	-	CASME II	0.815	-	-
[64]	2018	collaborative feature difference (CFD)	-	SMIC-E CASME II	- -	- -	0.942(AUC) 0.971(AUC)
[65]	2018	reiesz pyramid	LOSO	SMIC-HS CASME II	- -	- -	0.898(AUC) 0.951(AUC)

Besides, for facial landmark alignment or face registration, generic method like Kanade-Lucas-Tomasi (KLT) [69] is also proposed to track point features in different frames among a whole video.

### 3.3 Motion Magnification

Through motion magnification, subtle facial expressions become more recognizable. Local Amplitude-based Eulerian Motion Magnification (EMM) [70] is commonly used to magnify the micro-level movements by exaggerating the differences in the brightness and color of pixels at the same spatial location across two consecutive frames. Compared to local magnification, global-scale Lagrangian Motion Magnification (LMM) [71] was proposed for explicitly tracking all the displacements within a video in both the spatial and temporal domains. In LMM, each pixel can be traced back and forth through the timeline. Also, Principal Component Analysis (PCA) is employed in LMM to learn the statistically dominant displacements for all video frames. Lei *et al.* [72] first used the learning-based method for video motion magnification to extract shape representations from the intermediate layer of neural networks.

## 4 FEATURE EXTRACTION

Feature representation plays an important role for MER. With proper extraction, the raw input data of a ME video clip can be represented in a simple and concise form. The two main challenges for the feature descriptor of ME are 1) It should be able to capture the difference in both spatial and temporal domains and 2) It should be able to capture the micro-level differences. In our survey, we divide the commonly used feature representations into two main categories: handcrafted feature and learning-based feature.

### 4.1 Handcrafted Feature

Handcrafted feature for face data can be further divided into appearance-based feature and geometry-based feature [73]. Appearance-based feature represents the intensity or the

texture information of face region data while geometry-based feature describes the face geometrics such as the location of each facial landmark. In the literature, appearance-based feature has shown its effectiveness on dealing with illumination difference and unaligned images [74], [75], [76], [77]. An illustration of handcrafted features is shown in Fig 5.

**Local Binary Pattern (LBP)-based Feature.** LBP-based feature is commonly used in the literature due to the computational simplicity. The LBP was first proposed in [78] as a texture operator via thresholding the eight neighbors of each pixel and representing the result with a binary code. An extended version of LBP is described in [79] to meet the need of rotation invariance.

However, for an ME video clip, the dynamic information between different frames is crucial. As a time-domain extension of LBP, LBP-TOP was proposed to deal with dynamic feature analysis [80]. Usually, LBP-TOP features are applied with region-based algorithms to improve the robustness of misalignment. Currently, LBP-TOP is reported by most of the existing ME datasets as the baseline evaluation. Due to its computational simplicity, LBP-TOP features are utilized in a variety of different MER frameworks including cross-domain MER and ME spotting [81], [82], [61].

**LBP-TOP variants.** Aside from the use of original LBP-TOP operator, several variants were proposed to meet different needs for MER [52], [90], [91], [92], [47], [46], [93]. Huang *et al.* [52] combined the idea of integral projection and texture descriptor like LBP-TOP to bone texture characterization and face recognition. Guo *et al.* [91] extended the LBP-TOP operator into two novel binary descriptors: Angular Difference LBP-TOP (ADLBPTOP) and Radial Difference LBP-TOP (RDLBPTOP), respectively. Hu *et al.* [92] combined LBP-TOP and learning based features to form a feature fusion for multi-task learning. Hu *et al.* [47] also took advantage of Gabor filter and proposed Centralized Gabor Binary Pattern from Three Orthogonal Panels (CGBP-TOP). Yu *et al.* [46] proposed a new Local Cubes Binary Patterns (LCBP) especially for ME spotting. Instead of using three different plane combinations to represent the spatio-

TABLE 4  
Performance comparison among *macro-to-micro adaptation* methods

Methods	Year	Macro Datasets	Features	Protocol	Micro Datasets	Experimental Results		
						ACC	F1-score	UAR
[83]	2020	CK+ MMI Oulu-CASIA	ResNet18	LOSO (cross-dataset)	CASME II	0.756	0.701	0.872
					SAMM	0.741	0.736	0.819
					SMIC	0.768	0.744	0.861
[84]	2019	Oulu-CASIA	3D CNN	5-fold cross validation	CASME II	0.976	-	-
					SAMM	0.974	-	-
[85]	2019	BU-3DFE	ResNet	LOSO (cross-dataset)	CASME II	-	0.761	0.755
					SAMM	-	0.448	0.487
					SMIC	-	0.551	0.546
[86]	2018	CK+	LBP-TOP	-	CASME II	0.836	0.856	0.863
[87]	2018	CK+ Oulu-CASIA Jaffe, MUGFE	CNN	LOSO	CASME II	0.659	0.539	0.584
					SAMM	0.485	0.402	0.559
					SMIC	0.494	0.496	-
[88]	2018	CK+, Oulu-CASIA Jaffe, MUGFE	ResNet10	LOSO (cross-dataset)	CASME II	0.757	0.650	-
					SAMM	0.706	0.540	-
[89]	2018	CK+	DCP-TOP & HWP-TOP	-	CASME II	0.607	-	-

temporal feature [94], [95], they utilized a cube mask to encode the information of eight directional angles for both the space and time domains.

**Gradient-Based Feature.** One of the key challenges for the feature descriptor is to describe the subtle changes in ME sequences. Aside from LBP, local patterns based on gradient have been used for the property that high order gradient would represent the detailed structure information of an image. Dalal and Triggs proposed the histogram of gradients (HOG) in 2005 [54], which is one of the most commonly used feature descriptors when it comes to object recognition for the ability to specialize the edges in an image. The HOG descriptor has the property of geometric invariance and optical invariance. With HOG descriptors, the expression contour feature can be well captured.

Histogram of Image Gradient Orientation (HIGO) is a variant of HOG. It ignores the magnitude weighting in the original HOG and thus can suppress the illumination effect. Zhang *et al.* [82] utilized both HOG-TOP and HIGO-TOP as the feature descriptors in their work of MER. Tran *et al.* [61] built a spotting network based on LSTM using HOG-TOP and HIGO-TOP features as the input of the network. To better capture the structural changes in ME videos using gradient information, Niu *et al.* [49] proposed a local pattern of Local Two-Order Gradient Pattern (LTOGP).

**Optical Flow Based Feature.** The feature descriptors based on optical flow infer the relative motion information between different frames in order to capture subtle muscle movements for MER. The idea of optical flow was first introduced by Horn *et al.* [96] to describe the movement of brightness patterns in an image. By utilizing the pixel-wise difference between consecutive frames in a video clip, the motion information of the object in the video can thus be obtained. The basic concept is to find the distance of an identical object in different frames.

Owing to the fact that optical flow could capture temporal patterns between consecutive frames, one of the most employed architecture is to combine optical flow feature

with CNN to further recognize spatial patterns [97], [98], [99], [100], [101], [102].

Verburg *et al.* [59] utilized Histogram of Oriented Optical Flow (HOOF) to encode the subtle changes in the time domain for selected face regions. Li *et al.* [103], [104] revisited the HOOF feature descriptor and proposed an enhanced version to reduce the redundant dimensions in HOOF. Liong *et al.* [50] proposed another feature descriptor based on optical flow, Bi-Weighted Oriented Optical Flow (Bi-WOOF). Bi-WOOF represents a sequence of subtle expressions using only two frames. In contrast to HOOF, both the magnitude and optical strain values are used as weighting schemes to highlight the importance of each optical flow so the noisy optical flows with small intensities are reduced.

However, using histogram of optical flow as feature descriptors may have some flaws. When the histogram is used as a feature vector for a classifier, even a slightly shifted version of a histogram will create a huge difference as most classification algorithms use euclidean distance to measure the difference of two images. Happy *et al.* [105] proposed Fuzzy Histogram of Optical Flow Orientations (FHOFO) to collect the motion directions into angular bins based on the fuzzy membership function. It ignores the subtle motion magnitudes during the MEs and only takes the motion direction into consideration.

**Other Handcrafted Feature.** Apart from the abovementioned features, there are other descriptors to capture the distinctive properties embedded in ME videos.

Li *et al.* [106] defined the local and temporal patterns (LTP) of facial movement. LTP could be extracted from a projection in the PCA space. Instead of using LBP-TOP, Pawar *et al.* [107], [108] introduced a computationally efficient 3D version of Harris corner function, called Spatio Temporal Texture Map (STTM). Lin *et al.* [109] proposed a method using spatio-temporal Gabor filters.

While some approaches divided a face into regions spatially to better capture the low-intensity expression, it is hard to choose an ideal size of division because the division grids

directly affect the discriminative feature attained. Zong *et al.* [110] proposed a hierarchical spatial division scheme for spatio-temporal feature descriptor to address this issue.

To leverage the advantage of multiple feature types, there are a few approaches that take more than one type of feature descriptors. Wang *et al.* [111] first acquired facial LBP-TOP as the base feature, and then calculated the optical-flow as different feature weights. Zhao *et al.* [112], on the other hand, combined LBP-TOP and necessary morphological patches (NMPs).

## 4.2 Learning based Feature

In recent years, deep learning has shown a substantial improvement in computer vision tasks [116], [117], [118], [119], [120], [121], [122], [123]. Convolutional Neural Networks (CNNs, or ConvNet) have been considered as the most widely used ways for visual feature extraction with discriminative ability [116].

Since MEs are sequences of images, 2D CNN and 3D CNN are useful for their representative feature learning. 2D CNN is generally used on image data, which was first introduced in LeNet5 [124]. It is called two dimensional CNN because the kernel slides along two dimensions on the data. Input and output data in 3D CNN are four-dimensional since the kernel moves in three directions. 3D CNN is mostly used on 3D image data, e.g. MRI, CT scans, and videos [125], [126], [127].

**Shallow learning.** Deep learning-based methods require large-scale data to train the model. However, current public ME datasets are mostly limited and imbalanced, and tend to cause over-fitting issues when directly applying ConvNet. Most works usually adopt shallow and lightweight layers for MER. Dual Temporal Scale Convolutional Neural Network (DTSCNN) [128] was the first end-to-end middle-size neural network for MER. The DTSCNN designs two temporal channels, where each channel only has 4 convolutional layers and 4 pooling layers to partially avoid over-fitting. Peng *et al.* [88] adopted a simplified version of ResNet, ResNet10 for representative learning. Takalkar *et al.* [129] employed five convolutional layers, three max-pooling layers, and three fully-connected layers. Zhao *et al.* [101] used four convolutional layers and three pooling layers for capturing discriminative and high-level ME features. The  $1 \times 1$  convolutional layer is added right after the input layer to increase the non-linear expression of input data without increasing the computational load of the model. Micro-attention [87] was built with 10 residual blocks, whose micro-attention unit is an extension of ResNet. The shortcut connection designed for identity mapping can reduce the degradation problem. By learning the spatial attention of feature maps, the network can focus on the facial subtle movements.

**Two-step learning.** There are many works that first extract spatial features among all frames, then use recurrent convolutional layer or LSTM module to explore their temporal correlation. Due to the small amount of training samples, many learning based works [97], [131], [59] include handcrafted features (e.g., optical flow, HOOF) to give a higher signal-to-noise (SNR) ratio in comparison to using raw pixel data.

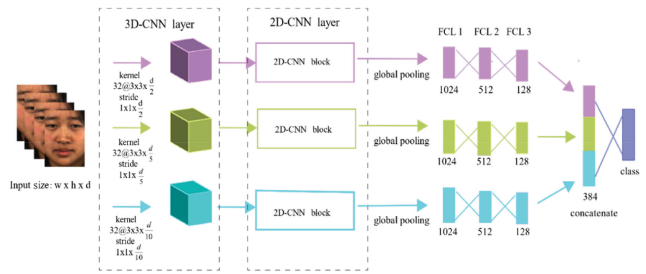


Fig. 6. The combination structure of the 2D CNN and 3D CNN in TSNN-LF for MER [130].

The optical flow map and its related extension algorithms, e.g., HOOF, are commonly used for extracting motion deformations of facial regions and have shown good performance [59], [97], [131], [59]. RCN-F [97] first extracted an optical flow map from the located onset and apex frames followed by a recurrent convolutional network (RCN). The RCN is composed of three parameter-free modules (i.e., wide expansion, shortcut connection, and attention unit) to enhance the representational ability. Enriched Long-term Recurrent Convolutional Network (ELRCN) [131] also used the optical flow as input data for two learning modules: dimension enrichment and temporal dimension enrichment.

**3D CNN.** In 2D CNN, convolution and pooling operations are applied on 2D feature maps which is lacking of motion information. To better preserve temporal information of the input signals effectively, 3D convolution has been proved its ability for capturing more representative features in the spatio-temporal aspect [125]. Considering that MEs occur only in a short period among dynamic facial movements and can be explicitly analyzed by detecting their constituent temporal segments, several MER works [132], [84], [133], [130] used 3D CNN to extract the facial temporal movements.

Li *et al.* [132] applied the 12-layer 3D flow-based CNNs model for MER, which extracts motion flow information arising from subtle facial movements. To better represent the dynamic and appearance features of MEs, optical flow and gray scale frames are combined as input data. The 3D convolution adopts  $3 \times 3 \times 3$  kernels to represent the spatial changes at each local region. Zhi *et al.* [84] employed 3D CNN for self-learning feature extraction. All convolutional kernels are set to the size of  $3 \times 5 \times 5$ , where 3 is the temporal depth and  $5 \times 5$  is the size of the spatial receptive field by taking consideration of efficiency and computational complexity. Reddy *et al.* [133] proposed two 3D CNN models for spontaneous MER, i.e., MicroExpSTCNN and MicroExpFuseNet, by exploiting the spatio-temporal information in the CNN framework.

Nevertheless, applying deep 3D CNN from scratch significantly increases the computational cost and memory demand. Besides, most popular 3D networks, e.g., C3D and P3D ResNet, are trained on Sports 1-M, that is an action database and is very different from MEs. Hence, it still needs to investigate the trained deep neural networks for ME feature extraction, which would boost the performance of the MER task [81]. TSNN [130] designed a three-stream combining 2D and 3D CNN to extract expression sequence feature. The advantages of single 3D kernel sizes and mul-



TABLE 5  
Performance comparison among recognition based on key apex frames methods

Methods	Year	Key Frames	Features	Protocol	Datasets	Experimental Results		
						ACC	F1-score	UAR
[113]	2020	Key Frame (SSIM)	dual-cross patterns (DCP)	-	CASME II	0.687	-	-
[97]	2020	Onset & Apex	RCN	LOSO	CASME II	-	0.856	0.812
					SAMM	-	0.765	0.677
					SMIC	-	0.658	0.660
[114]	2019	OFF & Apex	optical flow & CNN	LOSO	CASME II*	0.883	0.870	-
					SAMM*	0.682	0.542	-
					SMIC*	0.677	0.671	-
[101]	2019	Onset & Apex	optical flow	-	CASME II	0.870	-	-
					SAMM	0.698	-	-
					SMIC	0.702	-	-
[115]	2018	Apex	Bi-WOOF	-	CASME II	0.589	0.610	-
					SAMM	0.622	0.620	-
[85]	2019	Apex	ResNet	LOSO (cross-dataset)	CASME II	-	0.761	0.755
					SAMM	-	0.448	0.487
					SMIC	-	0.551	0.546
[88]	2018	Apex	ResNet10	LOSO (cross-dataset)	CASME II	0.757	0.650	-
					SAMM	0.706	0.540	-

\* on upper-right of the name of datasets means the labels of positive, negative and surprise were used in the experiments.

multiple 3D kernel combination have been made full use in the proposed framework to improve the MER performance. The architecture of TSNN is shown in Fig. 6.

## 5 MICRO-EXPRESSION SPOTTING

Micro-expression spotting refers to locate the segments of micro movements for a given video. There are less publications compared to MER in the literature [134] due to the difficulty of discovering subtle difference in a short duration (see Fig. 1). However, ME spotting is a vital step for automatic ME analysis and has attracted more and more attention recently, since proper spotting can decrease the redundant information for further recognition. Table 3 compares the *state-of-the-art* micro-expression spotting methods. Several earlier studies on automatic facial ME analysis primarily focused on distinguishing facial micro-expressions from macro-expressions [135], [136], [137]. Shreve *et al.* [136], [137] used an optical flow method for automatic ME spotting on their own database. However, their database contains 100 clips of posed MEs, which were obtained by asking participants to mimic some example videos that contain micro-expressions.

Tran *et al.* [59] first introduced deep sequence model for ME spotting. LSTM shows its effectiveness on using both local and global correlations of the extracted features to predict the score of the ME apex frame. Li *et al.* [19] first proposed an ME spotting method which was effective on spontaneous ME datasets. A spotting ME convolutional network [138] was designed for extracting features from video clips, which is the first time that deep learning is used in ME spotting.

The co-occurrence of macro- and micro-expressions are common in real life. An automatic spotting system for micro- and macro-expressions was designed by [58]. Based on the fact that micro- and macro-expressions have different

intensities, a multi-scale filter is used to improve the performance.

**Onset and offset detection.** While it is relatively easier to identify the peaks and valleys of facial movements, the onset and offset frames are much more difficult to determine. Locating the onset and offset frames is crucial for real-life situations where facial movements are continuously changing. CASME II and SAMM provide the onset and offset frame ground-truth, which can be helpful for model training.

**Apex frame spotting.** Liong *et al.* [139] introduced an automatic apex frame spotting method and this strategy was also adopted in [85], [115], [102]. The LBP feature descriptor was first employed to encode the features of each frame, then a *divide-and-conquer* methodology was exploited to detect the frames with peak facial changes as apex frame.

Most of the aforementioned ME spotting methods were conducted on public ME lab-controlled datasets. However, ME spotting in real-world scenes with different environmental factors is still an open issue. More ME datasets containing long video and real-world scenes are essential for further research.

## 6 MER METHODOLOGY

### 6.1 Macro-to-Micro Adaptation

As mentioned in the previous section, how to solve automatic labeling and recognition problems for MEs is a challenging task under the condition of small number of labeled training ME samples available. Meanwhile, there are large amounts of macro-expression databases [140], [141], each of which consists of vast labeled training samples compared with micro-expression databases. Thus, how to take advantage of the macro-expression databases for MER has become an important direction for research. Table 4 summarizes the *macro-to-micro adaptation* methods.

**Transfer-learning based method** has proved to be efficient in applying deep CNN on small databases [142]. Thus, by using the idea of transfer learning, it is reasonable to take advantage of the quantitative superiority of macro-expression to recognize the micro-expression [88], [86], [87], [84], [85]. In [86], the macro-expression features and micro-expression features are considered as training matrix, where the gallery features can be transformed to the probe features. The two different features were then projected into a joint subspace, where they are associated with each other. The nearest neighbor (NN) classifier was used to classify the probe micro-expression samples at the last step. Zhi *et al.* [87] pretrained the 3D CNN on macro-expression database Oulu-CASIA [143], and then the pre-trained model was transferred to the target micro-expression domain. The experimental results show 3.4% and 1.6% in MER performance higher than the model without transfer learning, respectively.

Xia *et al.* [83] imposed a loss inequality regularization to make the output of MicroNet converge to that of MacroNet. In [89], macro-expression images and micro-expression sequences were encoded by proposed hot wheel patterns (HWP), Dual-cross patterns (DCP-TOP) and HWP-TOP, respectively. The coupled metric learning algorithm was employed to model the shared features between micro-expression samples and macro-information. In [87], the original residual network (without micro-attention units) was first initialized with the ImageNet database. Then, to narrow the gap between object recognition (ImageNet) and facial expression recognition, the network was further pre-trained on several popular macro-expression databases, including CK+ [144], Oulu-CASIA NIR & VIS [143], Jaffe [145], and MUGFE [146]. Finally, the residual network together with micro-attention units is fine-tuned with micro-expression databases, including CASME II, SAMM and SMIC. Similarly, Jia *et al.* [88] proposed a macro-to-micro transformation network. ResNet10 pre-trained on ImageNet dataset was fine-tuned on macro-expression datasets first and then on the micro-expression datasets (CASME II, SAMM).

**Knowledge distillation** strategy is also used for MER. Sun *et al.* [51] proposed a multi-task teacher network containing AUs recognition and facial view classification on FERA2017 dataset [147], then the learned knowledge was distilled to a shallow student network for MER. In SA-AT [85], the larger scale of macro-expression data were served as auxiliary database to train a teacher model, then the teacher model was transferred to train the student model on micro-expression databases with limited samples. To narrow the gap between macro- and micro-expressions, a style aggregated strategy was used to transform micro-expression samples from different macro-expression datasets to generate an aggregated style via CycleGAN [148].

## 6.2 Recognition based on Key Apex Frames

Apex frame, which is the instant indicating the most expressive emotional state in a video sequence, is effective to classify the emotion in that particular frame. The apex frame portrays the highest intensity of facial motion among all frames. The apex occurs when the change in facial muscle

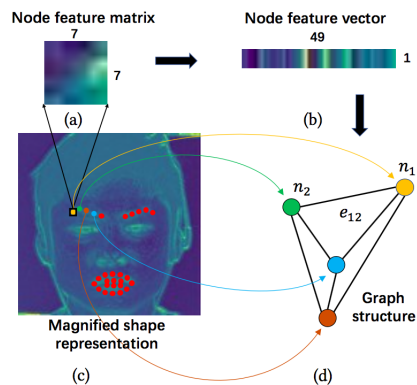


Fig. 7. Sample facial graph built in Graph-TCN [72].

reaches the peak or the highest intensity of the facial motion. Many works [114], [97] extracted features of the apex frame as feature descriptor.

Existing ME datasets, e.g., CASME II and SAMM, provide the index annotation of apex frames, while SMIC-HS database does not release the apex information. Thus, automatic apex frame spotting [139], [85], [115], [102] is necessary to be applied to each video to acquire the location of the apex frame. Table 5 summarizes the *state-of-the-art* key apex frames based MER approaches.

**Using single apex frame.** Considering the subtle motion change of ME frames in video, several works emphasized the use of a single apex frame to minimize the redundancy of the repetitive input frames. Using apex frame to represent the whole ME sequence was expected to reduce the computational complexity for feature learning [88], [102], [115], [85]. The experimental results of [51] showed that the effect of only using apex frame is better than onset-apex-offset sequence and the whole video.

**The optical flows from the onset and apex frames** are commonly used for extracting motion deformations of facial regions and can achieve good performance in subject independent evaluation [154], [98], [101], [155]. Zhou *et al.* [152] took the mid-position frames in ME samples based on the observation that most apex frame existing in the middle segment of a data sequence. OFF-ApexNet [114] adopted onset and apex frames to represent the ME details and extracted optical flow to encode the motion flow features based on two chosen frames. The two streams are combined in a fully-connected layer for MER. The similar strategy has been used in [97]. The obtained onset and apex frames are combined with the corresponding optical flow map as input of the proposed MER model. Zhao *et al.* [101] used the Total Variation-L1 (TV-L1) optical flow algorithm to calculate the motion information between the onset frame and the apex frame. Then, the corresponding optical flow feature map is obtained and fed into subsequent deep network learning.

**Others.** Key frames of each local region are computed by SSIM in [113], and then the dual-cross patterns (DCP) are applied to get the final feature vectors.

## 6.3 Recognition based on Facial Action Units

As physical subtle variations in facial expressions, AUs analysis are crucial because they constitute essential signals

TABLE 6  
Performance comparison on CASME II dataset

Methods	Year	Features	Classifier	Protocol	Experimental Results		
					ACC	F1-score	UAR
[91]	2019	ELBPTOP	SVM	LOSO	0.739	0.690	-
[91]*	2019	ELBPTOP	SVM	LOSO	0.796	0.660	-
[47]	2019	CGBP-TOP	Gaussian kernel function	LOSO	0.658	-	-
[103]	2019	revised HOOF + CNN	MLP	LOSO	0.580	-	-
[149]	2019	LMP	SVM	LOSO	0.702	-	-
[107]	2019	spatiotemporal texture map	SVM	k-fold cross validation	0.800	0.890	-
[150]*	2019	landmark point feature	MLP	LOSO	0.793	-	-
[52]	2018	discriminative STLBP	SVM + Gentle adaboost	LOSO	0.648	-	-
[151]	2018	FMBH	SVM	LOSO	0.643	-	-
[110]	2018	hierarchical ST descriptor	GSL model	LOSO	0.639	0.613	-
[97] *	2020	optical flow + RCN	MLP	LOSO	-	0.856	0.813
[98]*	2019	optical flow+STSTNet	MLP	LOSO	-	0.680	0.701
[99]*	2019	optical flow + CNN	MLP	LOSO	0.883	0.897	-
[100]*	2019	optical flow + CNN	MLP	LOSO	-	0.829	0.820
[152]*	2019	optical flow + CNN	MLP	LOSO	-	0.862	0.856
[111]	2019	LBP-TOP + optical flow	SVM	LOSO	0.696	-	-
[101]	2019	optical flow + CNN	MLP	LOSO	0.870	-	-
[153]	2019	dynamic imaging + CNN	MLP	k-fold cross validation	0.766	-	-
[113]	2019	local ROI+DCP	chi-square distance	-	0.687	-	-
[92]	2018	LGBP-TOP + CNN	MLP	LOSO	0.662	-	-
[109]	2018	frame difference + ST Gabor filter	SVM	LOSO	0.553	-	-
[32]	2020	3D CNN	Softmax	LOSO & LOVO	0.561	0.394	-
[72]	2020	Graph-TCN	Softmax	LOSO	0.740	0.725	-
[83]	2020	ResNet18	Triplet loss	LOSO	0.756	0.701	0.872
[132]	2019	3D flow-based CNN	SVM	LOSO	0.591	-	-
[84]	2019	3D CNN	SVM	5-fold cross validation	0.976	-	-
[88]	2018	ResNet10	-	LOSO	0.757	0.650	-
[87]	2018	CNN	Softmax	LOSO	0.659	0.539	0.584

The \* on upper-right of the methods mean only three labels of positive, negative and surprise were used in the experiments .

to be understood. Emotion states may be conceptually similar in terms of the difference in facial muscle movements (e.g., fear and disgust). There are many works focusing on facial AU recognition [156], [157], including AU occurrence detection [158], [159], [160], [161] and AU intensity estimation [162], [163]. Some surveys compare recent AUs recognition works, e.g., [164], [165].

Since the occurrence of AUs are strongly correlated, AU detection is usually considered as a multi-label learning problem. Several works considered the relationships among AUs and modeled AU interrelations to improve recognition accuracy [166], [167], [168]. However, most works rely on probabilistic graphical models with manually extracted features [169], [170]. Given that the graph has the ability of handling multi-relational data [168], Liu *et al.* [166] proposed the first work that employed GCN to model AU relationship. The cropped AU regions by EAC-Net [161] are fed into GCN as nodes, after that the propagation of the graph is determined by the relationship of AUs.

AU intensities are annotated by appending five-point ordinal scale (A-E for minimal-maximal intensity) [165]. The intensity levels are not uniform intervals. To infer the co-occurrences of AUs, it can be formulated into a heatmap regression-based framework. Fan *et al.* [162] modeled AU

co-occurring patterns among feature channels, where semantic descriptions and spatial distributions of AUs are both encoded.

Further insights into locating AUs could possibly provide even better discrimination between types of MEs. Some public ME datasets, e.g., CASME II and SAMM, contain emotion classes as well as AUs annotation, which makes the relationship construction between ME classes and AUs become achievable. Until now, there is few literature exploring automatic MER based on AUs [12], [171], [167], [172]. Wang *et al.* [12] defined 16 ROIs based on AUs and features extracted for each ROI can further boost the MER performance. Liu *et al.* [171] partitioned the facial region into 36 ROIs from 66 feature points and proposed a ROI-based optical feature, MDMO for MER. However, such predefined rules for modeling relationships among AUs may lead to limited generalization. Li *et al.* [167] used a structured knowledge-graph for AUs and integrated a Gated GNN (GGNN) to generate enhanced AU representation. This is the first work that integrates AU detection with MER. Lo *et al.* [172] proposed MER-GCN, which was the first end-to-end AU oriented MER architecture based on GCN, where GCN layers are able to discover the dependency laying between AU nodes for MER. Graph-TCN [72] utilized the

TABLE 7  
Performance comparison on SMIC dataset

Methods	Year	Features	Classifier	Protocol	Experimental Results		
					ACC	F1-score	UAR
[91]	2019	ELBPTOP	SVM	LOSO	0.691	0.650	0.660
[47]	2019	CGBP-TOP	Gaussian kernel function	LOSO	0.594	-	-
[149]	2019	LMP	SVM	LOSO	0.674	-	-
[52]	2018	discriminative STLBP	SVM + Gentle adaboost	LOSO	0.634	-	-
[151]	2018	FMBH	SVM	LOSO	0.683	-	-
[50]	2018	Bi-WOOF	SVM	LOSO	0.620	-	-
[105]	2018	FHOFO	SVM	LOSO	0.512	0.518	-
[110]	2018	hierarchical ST descriptor	GSL model	LOSO	0.604	0.613	-
[97]	2020	optical flow + RCN	MLP	LOSO	-	0.658	0.660
[111]	2019	LBP-TOP + optical flow	SVM	LOSO	0.717	-	-
[98]	2019	optical flow + STSTNet	MLP	LOSO	-	0.659	0.681
[100]	2019	optical flow + CNN	MLP	LOSO	-	0.746	0.753
[152]	2019	optical flow + CNN	MLP	LOSO	-	0.664	0.673
[99]	2019	optical flow + CNN	MLP	LOSO	0.677	0.671	-
[101]	2019	optical flow + CNN	MLP	LOSO	0.698	-	-
[153]	2019	dynamic imaging + CNN	MLP	k-fold	0.911	-	-
[92]	2018	LGBP-TOP + CNN	MLP	LOSO	0.651	-	-
[109]	2018	frame difference + ST Gabor filter	SVM	LOSO	0.545	-	-
[83]	2020	ResNet18	Triplet loss	LOSO	0.768	0.744	0.861
[132]	2019	3D flow-based CNN	SVM	LOSO	0.555	-	-
[133]	2019	3DCNN	Softmax	-	0.650	-	-
[87]	2018	CNN	Softmax	LOSO	0.494	0.496	-

graph structure for node and edge feature extraction, where the facial graph construction is shown in Fig. 7. Sun *et al.* [51] proposed a knowledge transfer strategy that distills and transfers AUs information for MER.

## 7 LOSS FUNCTION

### 7.1 Ranking Loss

Ranking Loss, such as *Contrastive Loss*, *Margin Loss*, *Hinge Loss* or *Triplet Loss*, is widely used in MER [82], [173], [153], [60], [174], [175], [83], [176]. The objective of ranking loss is to predict relative distances between inputs, which is also called metric learning [177]. Lim and Goh [176] proposed Fuzzy Qualitative Rank Classifier (FQRC) to model the ambiguity in MER task by using a multi-label rank classifier. Xia *et al.* [83] imposed an loss inequality regularization in triplet loss to make the output of MicroNet converge to that of MacroNet. The hinge loss is used for maximum-margin classification, most notably for Support Vector Machines (SVMs). SVM classify training data points by finding a discriminative hyperplane function according to different kernels such as Linear kernel, Polynomial kernel, and Radial Basis Function (RBF) kernel. Each different kernel has its own advantage for different separability and order of the dataset. Li *et al.* [60] split each video into 12 feature ensembles, which represent the ME local movements, and SVM was employed to classify these local ROIs. LEARNet [153] employed RankSVM [178] to compute the frame scores in a video. To reduce computation time, Sequential Minimal Optimization (SMO), one of the computationally fastest methods of evaluating linear SVMs, was used in [174].

Overall, SVM is one of the most well-applied classifiers for MER. However, they show poor performances when the feature dimension is far greater than the training sample number.

### 7.2 Cross-Entropy Loss

Unlike traditional methods, where the feature extraction step and the feature classification step are independent, deep networks can perform MER in an end-to-end way. Face recognition system, e.g., Deepface [179], first adopted cross-entropy based softmax loss for facial feature learning. Specifically, softmax loss is the most commonly used function that minimizes the cross-entropy between the estimated class probabilities and the ground-truth distribution. Most works directly applied softmax loss in the MER network [59], [97], [114]. For example, the softmax function was employed followed by an LSTM network which consists of two LSTM layers, each with 12 dimensions [59].

### 7.3 Others

Since the MER task suffers from not only high inter-class similarity but also high intra-class variation, instead of simply employed SVM or softmax loss, several works have proposed novel loss layers for MER [180], [181], [182], [183].

Inspired by the center loss, which penalizes the distance between deep features and their corresponding class centers, Lalitha *et al.* [180] combined cross-entropy loss and center loss to improve the discriminative energy of the deeply learned features. The focal loss [181] was utilized

TABLE 8  
Performance comparison on SAMM dataset

Methods	Year	Features	Classifier	Protocol	Experimental Results		
					ACC	F1-score	UAR
[91]	2019	ELBPTOP	SVM	LOSO	0.634	0.480	-
[91]*	2019	ELBPTOP	SVM	LOSO	-	0.780	0.720
[150]*	2019	landmark point feature	MLP	LOSO	0.746	-	-
[97]*	2020	optical flow + RCN	MLP	LOSO	-	0.765	0.677
[98]*	2019	optical flow + STSTNet	MLP	LOSO	-	0.838	0.869
[100]*	2019	optical flow + CNN	MLP	LOSO	-	0.775	0.715
[152]*	2019	optical + CNN	MLP	LOSO	-	0.587	0.566
[99]*	2019	optical + CNN	MLP	LOSO	0.682	0.542	-
[101]	2019	optical + CNN	MLP	LOSO	0.702	-	-
[83]	2020	ResNet18	Triplet loss	LOSO	0.741	0.736	0.819
[32]	2020	3D CNN	Softmax	LOSO & LOVO	0.523	0.357	-
[84]	2019	3D CNN	SVM	5-fold cross validation	0.974	-	-
[87]	2018	CNN	Softmax	LOSO	0.485	0.402	0.559

The \* on upper-right of the methods mean only three labels of positive, negative and surprise were used in the experiments .

in MACNN [182] to overcome the sample imbalance challenge. Compared with cross-entropy loss, focal loss is more suitable to solve the problem of small sample classification or imbalance between samples. Besides, Xie *et al.* [183] proposed a *feature loss* metric to use the complementary information of handcrafted and deep features during training.

## 8 EXPERIMENTS

### 8.1 Evaluation

#### 8.1.1 Evaluation Protocols

The commonly used evaluation protocols for MER are  $k$ -fold cross-validation, leave-one-subject-out (LOSO) and leave-one-video-out (LOVO).  $k$ -fold cross-validation repeats random sub-sampling validation. For LOSO validation, the model leaves out all samples of one single subject for the performance evaluation, and all other data are used as training data. The overall performance is then averaged from all folds. Similar to LOSO, LOVO validation protocol requires the model to pick up the frames from one video for validation purpose while all other data are sampled for training. From our review, the LOSO is the most widely used.

There are limitations for the above protocols. With severe class-imbalanced ME data, the effectiveness of  $k$ -fold is easily influenced. The same problem goes with LOVO, which can introduce additional biases on certain subjects that have more representations during the evaluation process. Moreover, the evaluation results may be over-estimated due to the large training data. For LOSO, the intrinsic ME dynamics of each subject may be limited since the intensity and manner of MEs may differ from person to person.

#### 8.1.2 Evaluation Metrics

The popularly used MER evaluation metrics include Accuracy (ACC) and F1-score. ACC shows the average hit rate across all classes and is susceptible to bias data, and thus it can only reflect the partial effectiveness of an MER classifier.

F1-score can remedy the bias issue by computing on the total true positives, false positives and false negatives to reveal the ME classes.

For cross-dataset evaluation, unweighted average recall (UAR) and weighted average recall (WAR) are commonly used [2], [184]. WAR refers to the number of correctly classified samples divided by the total number of samples. UAR is defined as mean accuracy of each class divided by the number of classes without consideration of samples per class, which can reduce the bias caused by the class imbalance.

### 8.2 Cross-dataset MER

Cross-dataset MER is one of the recently emerging while challenging problems in ME analysis. The training sets and testing sets are from different ME datasets, resulting in the inconsistency of feature distributions [185].

The cross-database recognition in Micro-Expression Grand Challenge (MEGC 2018) [2] used the CASME II and SAMM datasets to serve as the source and target ME database. Zong *et al.* [184] built a cross-database micro-expression recognition (CDMER) benchmark. There are two types of CDMER tasks. The TYPE-I is conducted between any two subsets of SMIC, i.e., HS, VIS, and NIR. And the TYPE-II uses the selected CASME II and one subset of SMIC, which is proved more difficult than TYPE-I [184]. The experiments showed the variant of Local Phase Quantization (LPQ), i.e., LPQ-TOP, and LBP with six intersection points (LBP-SIP) performed the best in terms of F1-score and ACC in TYPE-I and TYPE-II. Also, deep features perform rather poorly than most handcrafted features in the CDMER task. The possible reason may be the C3D was pre-trained on Sports1M and UCF101 datasets, whose samples are quite different from ME data. After finetuning the C3D on SMIC (HS), the performance of C3D features can be improved significantly. However, with a simple finetuning strategy based on the source database, it seems not enough for learning the database-invariant deep features.

The heterogeneous problem existing between source and target databases raises the level of difficulty of the CDMER task. For example, SMIC (NIR) is collected by a near-infrared camera, whose image quality is considerably different from images recorded by a high-speed camera (as in CASME II and SMIC (HS)) and visual camera (as in SMIC (VIS)). Li *et al.* [186] proposed the target-adapted least-squares regression (TALSUR) to learn a regression coefficient matrix between source samples and the provided labels. The idea is that such coefficient matrix could also reflect the sample-label relation in the target database domain. Zhang *et al.* [185] proposed a structure of the super wide regression network (SWiRN) for unsupervised cross-database MER. The *state-of-the-art* domain adaptation methods [187], [188] can be further exploited to reduce the differences between source and target domains to improve the performance of CDMER methods.

### 8.3 Overall Comparison

A comparison of MER methodologies with handcrafted features and learning-based features is provided in Table 6, Table 7 and Table 8, respectively. Several variants of optical flow and LBP-TOP have also been proposed in recent years. Also, optical flow related features are most frequently used handcrafted features. The experimental results show that the combination of optical flow features and CNN can achieve fine recognition accuracy. Among all methods using handcrafted feature, LEARNet [153] achieved averagely the best performance, proving the effectiveness of dynamic imaging and the feasibility of using key frames to represent a video sequence.

Although there are many approaches using handcrafted features, we can observe that the trend of MER is changing. In several approaches, low-level handcrafted feature extraction is considered as a pre-stage before high-level feature extraction and fine performances can be obtained. On the other hand, although the end-to-end learning based approaches are reported to have higher performances, they are still restricted by limited labeled data compared to other classification tasks in the literature, encouraging the development of data augmentation or transfer learning works.

While the obtained results from most works are encouraging, there are still some restrictions. Since there is no standard evaluation protocol and class setting, it is hard to reach the conclusions which method performs the best for MER. For example, there are some works only considering three or four emotion labels (i.e., Positive, Negative, Surprise, and Others) [50], [103]. The reduction of emotion classes makes the MER task simpler but also introduces the class bias towards negative categories since there is only one positive category (i.e., Happiness). Another problem is that different methods have different dataset splitting or composite strategies [97], [189]. For example, Lai *et al.* [182] mixed CASME and CASME II and split the training/testing set as 70/30 percentage. While SMIC, CASME II, and SAMM were combined into a composite dataset in [189].

Generally, the problem of real-world automatic micro-expression recognition and spotting still remain challenging, since the existing datasets are collected under the well-controlled environment and the data diversity is very limited.

The recognition performances are insufficient and there remain still many topics unexplored.

## 9 CHALLENGES AND POTENTIAL SOLUTIONS

### 9.1 In-the-wild Scenarios

Current MER researches have limited scenarios since existing publicly annotated datasets are all collected from lab-controlled environments. In the real-world scenario, more complex and natural emotion including various environmental factors, e.g., illumination interference, 3D head rotation, interaction or interrogation scenarios where more persons involved should be taken into account. Lai *et al.* [182] took a step forward to avoid excessive computation and established a real-time MER framework with 60 fps running on Intel Xeon 2.10 GHz CPU, 32 GB memory and Ubuntu 14.04 operating system. However, it was experimented on CASME and CASME II, which both are lab-controlled datasets. The MER in-the-wild still remains an open challenge.

### 9.2 Uncertainty Modeling

To figure out why the *state-of-the-art* methods reached a bottleneck in the recognition rate, several works introduced fuzzy set theory [190] into MER due to its ability to model the uncertainties. For ambiguous movements of different emotions, Fuzzy Qualitative can model the ambiguity by using the fuzzy membership function to represent each feature dimension with respect to each emotion class. Lim and Goh [176] considered ME as a non-mutual exclusive case. Instead of conducting crisp classification, Fuzzy Qualitative Rank Classifier (FQRC) was proposed to model the ambiguity in MER task by using a multi-label rank classifier. Chen *et al.* [191] proposed weighted fuzzy classification for analyzing emotion and achieved promising accuracy. Happy and Routray [105] constructed fuzzy histograms of orientation features based on optical flow.

### 9.3 Machine Bias

Previous studies have indicated that ingroup advantage in macro-expression recognition exists in various kinds of social groups, e.g., cultural groups, racial groups, and religious groups [192], [193], [194], [195], [196], [197], [198].

However, it remains unclear whether the social category of the target influences MER. Xie *et al.* [199] conducted the intergroup bias experiments among Chinese. The results showed that there is an ingroup disadvantage for the Chinese participants: The recognition accuracy of MEs of outgroup members (White targets) was actually higher than that of ingroup members (Asian targets). And the results also showed that such an intergroup bias is unaffected by the duration of MEs and the ingroup disadvantage remains the same even after the participants had received the training of Micro Expression Training Tool (METT). Regardless of the duration of MEs, there is an ingroup disadvantage for Chinese participants and such a bias still exists even after the training of MER training program.

One possible solution is to discover the hidden group-invariant features. Investigating the potential factors that may affect the MER help us to develop efficient MER training programs and efficient automatic recognition tools.

## 10 CONCLUSIONS AND FUTURE RESEARCH DIRECTIONS

The past decade has witnessed the development of many new MER algorithms. This paper provides a comprehensive review of recent advances in MER technology. Future MER works can expand from certain research directions below.

### 10.1 Enriching Limited and Biased Datasets

Current ME databases are usually collected from young subjects, especially undergraduate students. The age range should be extended in the future because MEs decoding is varying among different ages [200]. Meanwhile, the male/female percentage and ethnic groups should also be well considered. Also, to reduce the labeling process, synthetic data generation algorithms can be further exploited in MER as introduced in Section 2.

Due to the subjectivity of human annotators and the ambiguous nature of the expression labels, there might exist annotation inconsistency [201], [202], how to reduce label noises is also needs to consider.

### 10.2 Facial Asymmetrical Phenomenon for MER

*Chirality* is a chemistry term used to describe two objects that appear identical but not symmetrical when folded over onto themselves. Since human faces are naturally asymmetrical [203], *facial chirality* demonstrates the asymmetry as the left and the right side of faces differ in emotional communication when people experience multiple emotions at the same time or when there is an attempt to hide an emotion. The general observation is that emotional expressions are more intense on the left side of faces since the right cerebral hemisphere is dominant for the expression [204]. On the other hand, AUs may be coded as symmetrical or asymmetrical. Some existing datasets have left and right AUs annotation that can reveal such phenomenon. Hypothetically, the facial chirality implies that the left side of faces might already include sufficient features that can distinguish one emotion from another. However, there are few MER works extend their researches based on this hypothesis.

### 10.3 Multimodality for MER

Auxiliary metadata, e.g., words, gestures, voices, can serve as important cues for MER in real world scenarios. Among them, the body gesture is proved to be capable of conveying emotional information. A *Micro Gesture* dataset [205] containing 3,692 manually labeled gesture clips was released in 2019. The dataset collects subtle body movements that are elicited when hidden expressions are triggered in unconstrained situations. Their experiments verified the latent relation between one's micro-gestures and the hidden emotional states. Therefore, MEs can be fused with micro-gesture and other physiological signals for a finer level emotion understanding. Moreover, human physiological features, e.g., ECG, EDA, are very informative features for affective revealing, which are less considered in MER works.

## ACKNOWLEDGMENTS

This work was supported in part by the CTBC Bank under Industry-Academia Cooperation Project and the Ministry of Science and Technology of Taiwan under Grants MOST-108-2218-E-002-055, MOST-109-2223-E-009-002-MY3, MOST-109-2218-E-009-025 and MOST-109-2218-E-002-015.

## REFERENCES

- [1] P. Ekman, *Telling lies: Clues to deceit in the marketplace, politics, and marriage (revised edition)*. WW Norton & Company, 2009.
- [2] W. Merghani, A. Davison, and M. Yap, "Facial micro-expressions grand challenge 2018: evaluating spatio-temporal features for classification of objective classes," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2018, pp. 662–666.
- [3] J. See, M. H. Yap, J. Li, X. Hong, and S.-J. Wang, "Megc 2019—the second facial micro-expressions grand challenge," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2019, pp. 1–5.
- [4] J. Li, S. Wang, M. H. Yap, J. See, X. Hong, and X. Li, "Megc2020—the third facial micro-expression grand challenge," in *IEEE International Conference on Automatic Face and Gesture Recognition*, 2020, pp. 234–237.
- [5] X.-b. Shen, Q. Wu, and X.-l. Fu, "Effects of the duration of expressions on the recognition of microexpressions," *Journal of Zhejiang University Science B*, vol. 13, no. 3, pp. 221–230, 2012.
- [6] C. H. Yap, C. Kendrick, and M. H. Yap, "Samm long videos: A spontaneous facial micro-and macro-expressions dataset," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2020, pp. 194–199.
- [7] W.-J. Yan, Q. Wu, Y.-J. Liu, S.-J. Wang, and X. Fu, "Casme database: a dataset of spontaneous micro-expressions collected from neutralized faces," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2013, pp. 1–7.
- [8] W.-J. Yan, X. Li, S.-J. Wang, G. Zhao, Y.-J. Liu, Y.-H. Chen, and X. Fu, "Casme ii: An improved spontaneous micro-expression database and the baseline evaluation," *PLoS one*, vol. 9, no. 1, p. e86041, 2014.
- [9] X. Shen, Q. Wu, K. Zhao, and X. Fu, "Electrophysiological evidence reveals differences between the recognition of microexpressions and macroexpressions," *Frontiers in psychology*, vol. 7, p. 1346, 2016.
- [10] R. Ekman, *What the face reveals: Basic and applied studies of spontaneous expression using the Facial Action Coding System (FACS)*. Oxford University Press, USA, 1997.
- [11] E. L. Rosenberg and P. Ekman, *What the face reveals: Basic and applied studies of spontaneous expression using the facial action coding system (FACS)*. Oxford University Press, 2020.
- [12] S.-J. Wang, W.-J. Yan, X. Li, G. Zhao, C.-G. Zhou, X. Fu, M. Yang, and J. Tao, "Micro-expression recognition using color spaces," *IEEE Transactions on Image Processing*, vol. 24, no. 12, pp. 6034–6047, 2015.
- [13] M. Pantic, "Machine analysis of facial behaviour: Naturalistic and dynamic behaviour," *Philosophical Transactions of the Royal Society B: Biological Sciences*, vol. 364, no. 1535, pp. 3505–3513, 2009.
- [14] W. V. Friesen and P. Ekman, *EMFACS-7: Emotional Facial Action Coding System, Version 7*. Unpublished manual, 1984.
- [15] P. Ekman, W. V. Friesen, and J. C. Hager, "Facial action coding system: The manual on cd rom," *A Human Face*, Salt Lake City, pp. 77–254, 2002.
- [16] R. Cairns, G. Elder Jr, and E. Costello, "Cambridge studies in social and emotional development," *Developmental science*. Cambridge University Press, 1996.
- [17] J. I. Durán, R. Reisenzein, and J.-M. Fernández-Dols, "Coherence between emotions and facial expressions," *The science of facial expression*, pp. 107–129, 2017.
- [18] W.-J. Yan, S.-J. Wang, Y.-J. Liu, Q. Wu, and X. Fu, "For micro-expression recognition: Database and suggestions," *Neurocomputing*, vol. 136, pp. 82–87, 2014.
- [19] X. Li, X. Hong, A. Moilanen, X. Huang, T. Pfister, G. Zhao, and M. Pietikäinen, "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, 2017.

- [20] K. M. Goh, C. H. Ng, L. L. Lim, and U. U. Sheikh, "Micro-expression recognition: an updated review of current trends, challenges and solutions," *The Visual Computer*, vol. 36, no. 3, pp. 445–468, 2018.
- [21] Y.-H. Oh, J. See, A. C. Le Ngo, R. C.-W. Phan, and V. M. Baskaran, "A survey of automatic facial micro-expression analysis: databases, methods, and challenges," *Frontiers in Psychology*, vol. 9, p. 1128, 2018.
- [22] W. Merghani, A. K. Davison, and M. H. Yap, "A review on facial micro-expressions analysis: datasets, features and metrics," *arXiv preprint arXiv:1805.02397*, 2018.
- [23] M. Takalkar, M. Xu, Q. Wu, and Z. Chaczko, "A survey: facial micro-expression recognition," *Multimedia Tools and Applications*, vol. 77, no. 15, pp. 19 301–19 325, 2018.
- [24] L. Zhou, X. Shao, and Q. Mao, "A survey of micro-expression recognition," *Image and Vision Computing*, p. 104043, 2020.
- [25] P. Ekman, "Facial expressions of emotion: an old controversy and new findings," *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, vol. 335, no. 1273, pp. 63–69, 1992.
- [26] X. Li, T. Pfister, X. Huang, G. Zhao, and M. Pietikäinen, "A spontaneous micro-expression database: Inducement, collection and baseline," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2013, pp. 1–6.
- [27] A. K. Davison, W. Merghani, and M. H. Yap, "Objective classes for micro-facial expression recognition," *Journal of Imaging*, vol. 4, no. 10, p. 119, 2018.
- [28] W. Merghani, A. K. Davison, and M. H. Yap, "The implication of spatial temporal changes on facial micro-expression analysis," *Multimedia Tools and Applications*, vol. 78, no. 15, pp. 21 613–21 628, 2019.
- [29] J. A. Coan and J. J. Allen, *Handbook of emotion elicitation and assessment*. Oxford university press, 2007.
- [30] P. Lang, "Behavioral treatment and bio-behavioral assessment: Computer applications," *Technology in Mental Health Care Delivery Systems*, pp. 119–137, 1980.
- [31] P. Ekman, "Lie catching and microexpressions," *The Philosophy of Deception*, vol. 1, no. 2, p. 5, 2009.
- [32] H.-X. Xie, L. Lo, H.-H. Shuai, and W.-H. Cheng, "Au-assisted graph attention convolutional network for micro-expression recognition," in *Proceedings of the 28th ACM International Conference on Multimedia*. ACM, 2020, pp. 2871–2880.
- [33] F. Qu, S.-J. Wang, W.-J. Yan, H. Li, S. Wu, and X. Fu, "Cas (me)<sup>2</sup>: A database for spontaneous macro-expression and micro-expression spotting and recognition," *IEEE Transactions on Affective Computing*, vol. 9, no. 4, pp. 424–436, 2017.
- [34] M. Frank, M. Herbasz, K. Sinuk, A. Keller, and C. Nolan, "I see how you feel: Training laypeople and professionals to recognize fleeting emotions," in *The Annual Meeting of the International Communication Association*. ICA, 2009, pp. 1–35.
- [35] W. T. Freeman, E. C. Pasztor, and O. T. Carmichael, "Learning low-level vision," *International Journal of Computer Vision*, vol. 40, no. 1, pp. 25–47, 2000.
- [36] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black, "A naturalistic open source movie for optical flow evaluation," in *European Conference on Computer Vision*. Springer, 2012, pp. 611–625.
- [37] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. Van Der Smagt, D. Cremers, and T. Brox, "Flownet: Learning optical flow with convolutional networks," in *IEEE International Conference on Computer Vision*. IEEE, 2015, pp. 2758–2766.
- [38] Q. Wang, J. Gao, W. Lin, and Y. Yuan, "Learning from synthetic data for crowd counting in the wild," in *IEEE International Conference on Computer Vision and Pattern Recognition*. IEEE, 2019, pp. 8198–8207.
- [39] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, "Playing for data: Ground truth from computer games," in *European Conference on Computer Vision*. Springer, 2016, pp. 102–118.
- [40] J. Papon and M. Schoeler, "Semantic pose using deep networks trained on synthetic rgb-d," in *IEEE International Conference on Computer Vision*. IEEE, 2015, pp. 774–782.
- [41] G. J. Qi and J. Luo, "Small data challenges in big data era: A survey of recent progress on unsupervised and semi-supervised methods," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2020.
- [42] R. Queiroz, M. Cohen, J. L. Moreira, A. Braun, J. C. J. Júnior, and S. R. Musse, "Generating facial ground truth with synthetic faces," in *IEEE Conference on Graphics, Patterns and Images*. IEEE, 2010, pp. 25–31.
- [43] I. Abbasnejad, S. Sridharan, D. Nguyen, S. Denman, C. Fookes, and S. Lucey, "Using synthetic data to improve facial expression analysis with 3d convolutional networks," in *IEEE International Conference on Computer Vision Workshops*. IEEE, 2017, pp. 1609–1618.
- [44] F. Zhang, T. Zhang, Q. Mao, and C. Xu, "Joint pose and expression modeling for facial expression recognition," in *IEEE International Conference on Computer Vision and Pattern Recognition*. IEEE, 2018, pp. 3359–3368.
- [45] J. Cai, Z. Meng, A. S. Khan, Z. Li, J. O'Reilly, and Y. Tong, "Identity-free facial expression recognition using conditional generative adversarial network," *arXiv preprint arXiv:1903.08051*, 2019.
- [46] M. Yu, Z. Guo, Y. Yu, Y. Wang, and S. Cen, "Spatiotemporal feature descriptor for micro-expression recognition using local cube binary pattern," *IEEE Access*, vol. 7, pp. 159 214–159 225, 2019.
- [47] C. Hu, J. Chen, X. Zuo, H. Zou, X. Deng, and Y. Shu, "Gender-specific multi-task micro-expression recognition using pyramid cgbp-top feature," *Computer Modeling in Engineering and Sciences*, vol. 118, no. 3, pp. 547–559, 2019.
- [48] P. Carcagni, M. Del Coco, M. Leo, and C. Distanto, "Facial expression recognition and histograms of oriented gradients: a comprehensive study," *SpringerPlus*, vol. 4, no. 1, p. 645, 2015.
- [49] M. Niu, Y. Li, J. Tao, and S.-J. Wang, "Micro-expression recognition based on local two-order gradient pattern," in *IEEE Asian Conference on Affective Computing and Intelligent Interaction*. IEEE, 2018, pp. 1–6.
- [50] S.-T. Liang, J. See, K. Wong, and R. C.-W. Phan, "Less is more: Micro-expression recognition from video using apex frame," *Signal Processing: Image Communication*, vol. 62, pp. 82–92, 2018.
- [51] B. Sun, S. Cao, D. Li, J. He, and L. Yu, "Dynamic micro-expression recognition using knowledge distillation," *IEEE Transactions on Affective Computing*, 2020.
- [52] X. Huang, S.-J. Wang, X. Liu, G. Zhao, X. Feng, and M. Pietikäinen, "Discriminative spatiotemporal local binary pattern with revisited integral projection for spontaneous facial micro-expression recognition," *IEEE Transactions on Affective Computing*, vol. 10, no. 1, pp. 32–47, 2017.
- [53] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *IEEE International Conference on Computer Vision and Pattern Recognition*, vol. 1. IEEE, 2001, pp. I–I.
- [54] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *IEEE International Conference on Computer Vision and Pattern Recognition*. IEEE, 2005, pp. 886–893.
- [55] J. Sanchez-Riera, K. Srinivasan, K.-L. Hua, W.-H. Cheng, M. A. Hossain, and M. F. Alhamid, "Robust rgb-d hand tracking using deep learning priors," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 9, pp. 2289–2301, 2017.
- [56] D. E. King, "Max-margin object detection," *arXiv preprint arXiv:1502.00046*, 2015.
- [57] I. Masi, Y. Wu, T. Hassner, and P. Natarajan, "Deep face recognition: A survey," in *2018 31st SIBGRAPI conference on graphics, patterns and images (SIBGRAPI)*. IEEE, 2018, pp. 471–478.
- [58] L.-w. Zhang, J. Li, S. Wang, X. Duan, W. Yan, H. Xie, and S. Huang, "Spatio-temporal fusion for macro-and micro-expression spotting in long video sequences," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2020, pp. 245–252.
- [59] M. Verburg and V. Menkovski, "Micro-expression detection in long videos using optical flow and recurrent neural networks," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2019, pp. 1–6.
- [60] J. Li, C. Soladie, R. Seghier, S.-J. Wang, and M. H. Yap, "Spotting micro-expressions on long videos sequences," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2019, pp. 1–5.
- [61] T.-K. Tran, Q.-N. Vo, X. Hong, and G. Zhao, "Dense prediction for micro-expression spotting based on deep sequence model," *Electronic Imaging*, vol. 2019, no. 8, pp. 401–1, 2019.
- [62] K. X. Beh and K. M. Goh, "Micro-expression spotting using facial landmarks," in *IEEE International Colloquium on Signal Processing and Its Applications*. IEEE, 2019, pp. 192–197.



- [63] S. Nag, A. K. Bhunia, A. Konwer, and P. P. Roy, "Facial micro-expression spotting and recognition using time contrasted feature with visual memory," in *IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2019, pp. 2022–2026.
- [64] Y. Han, B. Li, Y.-K. Lai, and Y.-J. Liu, "Cfd: a collaborative feature difference method for spontaneous micro-expression spotting," in *IEEE International Conference on Image Processing*. IEEE, 2018, pp. 1942–1946.
- [65] C. A. Duque, O. Alata, R. Emonet, A.-C. Legrand, and H. Konik, "Micro-expression spotting using the riesz pyramid," in *IEEE Winter Conference on Applications of Computer Vision*. IEEE, 2018, pp. 66–74.
- [66] S. Milborrow and F. Nicolls, "Locating facial features with an extended active shape model," in *European conference on computer vision*. Springer, 2008, pp. 504–513.
- [67] A. Asthana, S. Zafeiriou, S. Cheng, and M. Pantic, "Robust discriminative response map fitting with constrained local models," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2013, pp. 3444–3451.
- [68] Y. Wu, T. Hassner, K. Kim, G. Medioni, and P. Natarajan, "Facial landmark detection with tweaked convolutional neural networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 12, pp. 3067–3074, 2017.
- [69] S. Birchfield, "Derivation of kanade-lucas-tomasi tracking equation," *unpublished notes*, 1997.
- [70] A. C. Le Ngo and R. C.-W. Phan, "Seeing the invisible: Survey of video motion magnification and small motion analysis," *ACM Computing Surveys (CSUR)*, vol. 52, no. 6, pp. 1–20, 2019.
- [71] A. C. Le Ngo, A. Johnston, R. C.-W. Phan, and J. See, "Micro-expression motion magnification: Global lagrangian vs. local eulerian approaches," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2018, pp. 650–656.
- [72] L. Lei, J. Li, T. Chen, and S. Li, "A novel graph-tcn with a graph structured representation for micro-expression recognition," in *Proceedings of the 28th ACM International Conference on Multimedia*. ACM, 2020, pp. 2237–2245.
- [73] Z. Zeng, M. Pantic, G. I. Roisman, and T. S. Huang, "A survey of affect recognition methods: Audio, visual, and spontaneous expressions," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 1, pp. 39–58, 2008.
- [74] Y.-C. Lin, M.-C. Hu, W.-H. Cheng, Y.-H. Hsieh, and H.-M. Chen, "Human action recognition and retrieval using sole depth information," in *Proceedings of the 20th ACM international conference on Multimedia*, 2012, pp. 1053–1056.
- [75] M.-C. Hu, C.-W. Chen, W.-H. Cheng, C.-H. Chang, J.-H. Lai, and J.-L. Wu, "Real-time human movement retrieval and assessment with kinect sensor," *IEEE transactions on cybernetics*, vol. 45, no. 4, pp. 742–753, 2014.
- [76] I.-C. Shen and W.-H. Cheng, "Gestalt rule feature points," *IEEE Transactions on Multimedia*, vol. 17, no. 4, pp. 526–537, 2015.
- [77] A. R. Rivera, J. R. Castillo, and O. O. Chae, "Local directional number pattern for face analysis: Face and expression recognition," *IEEE transactions on image processing*, vol. 22, no. 5, pp. 1740–1752, 2012.
- [78] T. Ojala, M. Pietikainen, and D. Harwood, "Performance evaluation of texture measures with classification based on kullback discrimination of distributions," in *IEEE International Conference on Computer Vision and Pattern Recognition*. IEEE, 1994, pp. 582–585.
- [79] T. Ojala, M. Pietikainen, and T. Maenpaa, "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.
- [80] 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.
- [81] Y. Zong, W. Zheng, X. Huang, J. Shi, Z. Cui, and G. Zhao, "Domain regeneration for cross-database micro-expression recognition," *IEEE Transactions on Image Processing*, vol. 27, no. 5, pp. 2484–2498, 2018.
- [82] Y. Zhang, H. Jiang, X. Li, B. Lu, K. M. Rabie, and A. U. Rehman, "A new framework combining local-region division and feature selection for micro-expressions recognition," *IEEE Access*, vol. 8, pp. 94 499–94 509, 2020.
- [83] B. Xia, W. Wang, S. Wang, and E. Chen, "Learning from macro-expression: a micro-expression recognition framework," in *Proceedings of the 28th ACM International Conference on Multimedia*. ACM, 2020, pp. 2936–2944.
- [84] R. Zhi, H. Xu, M. Wan, and T. Li, "Combining 3d convolutional neural networks with transfer learning by supervised pre-training for facial micro-expression recognition," *IEICE Transactions on Information and Systems*, vol. 102, no. 5, pp. 1054–1064, 2019.
- [85] L. Zhou, Q. Mao, and L. Xue, "Cross-database micro-expression recognition: a style aggregated and attention transfer approach," in *IEEE International Conference on Multimedia and Expo Workshops*. IEEE, 2019, pp. 102–107.
- [86] X. Jia, X. Ben, H. Yuan, K. Kpalma, and W. Meng, "Macro-to-micro transformation model for micro-expression recognition," *Journal of Computational Science*, vol. 25, pp. 289–297, 2018.
- [87] C. Wang, M. Peng, T. Bi, and T. Chen, "Micro-attention for micro-expression recognition," *Neurocomputing*, vol. 410, pp. 354–362, 2020.
- [88] M. Peng, Z. Wu, Z. Zhang, and T. Chen, "From macro to micro expression recognition: Deep learning on small datasets using transfer learning," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2018, pp. 657–661.
- [89] X. Ben, X. Jia, R. Yan, X. Zhang, and W. Meng, "Learning effective binary descriptors for micro-expression recognition transferred by macro-information," *Pattern Recognition Letters*, vol. 107, pp. 50–58, 2018.
- [90] L. Mao, N. Wang, L. Wang, and Y. Chen, "Classroom micro-expression recognition algorithms based on multi-feature fusion," *IEEE Access*, vol. 7, pp. 64 978–64 983, 2019.
- [91] C. Guo, J. Liang, G. Zhan, Z. Liu, M. Pietikainen, and L. Liu, "Extended local binary patterns for efficient and robust spontaneous facial micro-expression recognition," *IEEE Access*, vol. 7, pp. 174 517–174 530, 2019.
- [92] C. Hu, D. Jiang, H. Zou, X. Zuo, and Y. Shu, "Multi-task micro-expression recognition combining deep and handcrafted features," in *IEEE International Conference on Pattern Recognition*. IEEE, 2018, pp. 946–951.
- [93] C. Arango Duque, O. Alata, R. Emonet, H. Konik, and A.-C. Legrand, "Mean oriented riesz features for micro expression classification," *arXiv*, pp. arXiv–2005, 2020.
- [94] B. Wu, T. Mei, W.-H. Cheng, and Y. Zhang, "Unfolding temporal dynamics: Predicting social media popularity using multi-scale temporal decomposition," 2016.
- [95] B. Wu, W.-H. Cheng, Y. Zhang, and T. Mei, "Time matters: Multi-scale temporalization of social media popularity," in *Proceedings of the 24th ACM international conference on Multimedia*, 2016, pp. 1336–1344.
- [96] B. K. Horn and B. G. Schunck, "Determining optical flow," in *Techniques and Applications of Image Understanding*, vol. 281. International Society for Optics and Photonics, 1981, pp. 319–331.
- [97] Z. Xia, W. Peng, H.-Q. Khor, X. Feng, and G. Zhao, "Revealing the invisible with model and data shrinking for composite-database micro-expression recognition," *IEEE Transactions on Image Processing*, vol. 29, pp. 8590–8605, 2020.
- [98] S.-T. Liong, Y. Gan, J. See, H.-Q. Khor, and Y.-C. Huang, "Shallow triple stream three-dimensional cnn (ststnet) for micro-expression recognition," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2019, pp. 1–5.
- [99] Y. Gan, S.-T. Liong, W.-C. Yau, Y.-C. Huang, and L.-K. Tan, "Off-apexnet on micro-expression recognition system," *Signal Processing: Image Communication*, vol. 74, pp. 129–139, 2019.
- [100] Y. Liu, H. Du, L. Zheng, and T. Gedeon, "A neural micro-expression recognizer," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2019, pp. 1–4.
- [101] Y. Zhao and J. Xu, "A convolutional neural network for compound micro-expression recognition," *Sensors*, vol. 19, no. 24, p. 5553, 2019.
- [102] Y. Gan and S.-T. Liong, "Bi-directional vectors from apex in cnn for micro-expression recognition," in *IEEE International Conference on Image, Vision and Computing*. IEEE, 2018, pp. 168–172.
- [103] Q. Li, S. Zhan, L. Xu, and C. Wu, "Facial micro-expression recognition based on the fusion of deep learning and enhanced optical flow," *Multimedia Tools and Applications*, vol. 78, no. 20, pp. 29 307–29 322, 2019.
- [104] Q. Li, J. Yu, T. Kurihara, and S. Zhan, "Micro-expression analysis by fusing deep convolutional neural network and optical flow," in *IEEE International Conference on Control, Decision and Information Technologies*. IEEE, 2018, pp. 265–270.

- [105] S. Happy and A. Routray, "Fuzzy histogram of optical flow orientations for micro-expression recognition," *IEEE Transactions on Affective Computing*, vol. 10, no. 3, pp. 394–406, 2017.
- [106] J. Li, C. Soladie, and R. Segulier, "Ltp-ml: Micro-expression detection by recognition of local temporal pattern of facial movements," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2018, pp. 634–641.
- [107] S. S. Pawar, M. Moh, and T.-S. Moh, "Micro-expression recognition using motion magnification and spatiotemporal texture map," in *International Conference on Ubiquitous Information Management and Communication*. Springer, 2019, pp. 351–369.
- [108] J. Sanchez-Riera, K.-L. Hua, Y.-S. Hsiao, T. Lim, S. C. Hidayati, and W.-H. Cheng, "A comparative study of data fusion for rgb-d based visual recognition," *Pattern Recognition Letters*, vol. 73, pp. 1–6, 2016.
- [109] C. Lin, F. Long, J. Huang, and J. Li, "Micro-expression recognition based on spatiotemporal gabor filters," in *IEEE International Conference on Information Science and Technology*. IEEE, 2018, pp. 487–491.
- [110] Y. Zong, X. Huang, W. Zheng, Z. Cui, and G. Zhao, "Learning from hierarchical spatiotemporal descriptors for micro-expression recognition," *IEEE Transactions on Multimedia*, vol. 20, no. 11, pp. 3160–3172, 2018.
- [111] L. Wang, H. Xiao, S. Luo, J. Zhang, and X. Liu, "A weighted feature extraction method based on temporal accumulation of optical flow for micro-expression recognition," *Signal Processing: Image Communication*, vol. 78, pp. 246–253, 2019.
- [112] Y. Zhao and J. Xu, "An improved micro-expression recognition method based on necessary morphological patches," *Symmetry*, vol. 11, no. 4, p. 497, 2019.
- [113] W. Zhong, X. Yu, L. Shi, and Z. Xie, "Facial micro-expression recognition based on local region of the key frame," in *International Symposium on Multispectral Image Processing and Pattern Recognition*, vol. 11430. International Society for Optics and Photonics, 2020, p. 114301L.
- [114] Y. Gan, S.-T. Liang, W.-C. Yau, Y.-C. Huang, and L.-K. Tan, "Off-apexnet on micro-expression recognition system," *Signal Processing: Image Communication*, vol. 74, pp. 129–139, 2019.
- [115] S.-T. Liang, J. See, K. Wong, and R. C.-W. Phan, "Less is more: Micro-expression recognition from video using apex frame," *Signal Processing: Image Communication*, vol. 62, pp. 82–92, 2018.
- [116] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *IEEE conference on Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 248–255.
- [117] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [118] H.-C. Wang, Y.-C. Lai, W.-H. Cheng, C.-Y. Cheng, and K.-L. Hua, "Background extraction based on joint gaussian conditional random fields," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 11, pp. 3127–3140, 2017.
- [119] K.-L. Hua, C.-H. Hsu, S. C. Hidayati, W.-H. Cheng, and Y.-J. Chen, "Computer-aided classification of lung nodules on computed tomography images via deep learning technique," *Oncotargets and therapy*, vol. 8, 2015.
- [120] F. S. Abousaleh, T. Lim, W.-H. Cheng, N.-H. Yu, M. A. Hossain, and M. F. Alhamid, "A novel comparative deep learning framework for facial age estimation," *EURASIP Journal on Image and Video Processing*, vol. 2016, no. 1, p. 47, 2016.
- [121] C.-W. Hsieh, C.-Y. Chen, C.-L. Chou, H.-H. Shuai, J. Liu, and W.-H. Cheng, "Fashionon: Semantic-guided image-based virtual try-on with detailed human and clothing information," in *Proceedings of the 27th ACM International Conference on Multimedia*. ACM, 2019, pp. 275–283.
- [122] I. Y. Susanto, T.-Y. Pan, C.-W. Chen, M.-C. Hu, and W.-H. Cheng, "Emotion recognition from galvanic skin response signal based on deep hybrid neural networks," in *International Conference on Multimedia Retrieval*. ACM, 2020, pp. 341–345.
- [123] W.-H. Cheng, J. Liu, M. S. Kankanhalli, A. El Saddik, and B. Huët, "Ai+ multimedia make better life?" in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 1455–1456.
- [124] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [125] 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*, vol. 35, no. 1, pp. 221–231, 2012.
- [126] Y.-C. Chen, D. S. Tan, W.-H. Cheng, and K.-L. Hua, "3d object completion via class-conditional generative adversarial network," in *International Conference on Multimedia Modeling*. Springer, 2019, pp. 54–66.
- [127] K.-L. Hua, C.-H. Hsu, S. C. Hidayati, W.-H. Cheng, and Y.-J. Chen, "Computer-aided classification of lung nodules on computed tomography images via deep learning technique," *Oncotargets and therapy*, vol. 8, 2015.
- [128] M. Peng, C. Wang, T. Chen, G. Liu, and X. Fu, "Dual temporal scale convolutional neural network for micro-expression recognition," *Frontiers in Psychology*, vol. 8, p. 1745, 2017.
- [129] M. A. Takalkar and M. Xu, "Image based facial micro-expression recognition using deep learning on small datasets," in *IEEE International Conference on Digital Image Computing: Techniques and Applications*. IEEE, 2017, pp. 1–7.
- [130] C. Wu and F. Guo, "Tsn: Three-stream combining 2d and 3d convolutional neural network for micro-expression recognition," *IEEE Transactions on Electrical and Electronic Engineering*, 2020.
- [131] H.-Q. Khor, J. See, R. C. W. Phan, and W. Lin, "Enriched long-term recurrent convolutional network for facial micro-expression recognition," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2018, pp. 667–674.
- [132] J. Li, Y. Wang, J. See, and W. Liu, "Micro-expression recognition based on 3d flow convolutional neural network," *Pattern Analysis and Applications*, vol. 22, no. 4, pp. 1331–1339, 2019.
- [133] S. P. T. Reddy, S. T. Karri, S. R. Dubey, and S. Mukherjee, "Spontaneous facial micro-expression recognition using 3d spatiotemporal convolutional neural networks," in *IEEE International Joint Conference on Neural Networks*. IEEE, 2019, pp. 1–8.
- [134] P. Gupta, B. Bhowmick, and A. Pal, "Exploring the feasibility of face video based instantaneous heart-rate for micro-expression spotting," in *IEEE International Conference on Computer Vision and Pattern Recognition workshops*. IEEE, 2018, pp. 1316–1323.
- [135] T.-K. Tran, Q.-N. Vo, X. Hong, X. Li, and G. Zhao, "Micro-expression spotting: A new benchmark," *arXiv preprint arXiv:2007.12421*, 2020.
- [136] M. Shreve, S. Godavarthy, V. Manohar, D. Goldgof, and S. Sarkar, "Towards macro-and micro-expression spotting in video using strain patterns," in *IEEE Workshop on Applications of Computer Vision*. IEEE, 2009, pp. 1–6.
- [137] M. Shreve, S. Godavarthy, D. Goldgof, and S. Sarkar, "Macro-and micro-expression spotting in long videos using spatio-temporal strain," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2011, pp. 51–56.
- [138] Z. Zhang, T. Chen, H. Meng, G. Liu, and X. Fu, "Smeconvnet: A convolutional neural network for spotting spontaneous facial micro-expression from long videos," *IEEE Access*, vol. 6, pp. 71 143–71 151, 2018.
- [139] S.-T. Liang, J. See, K. Wong, A. C. Le Ngo, Y.-H. Oh, and R. Phan, "Automatic apex frame spotting in micro-expression database," in *IEEE Asian conference on pattern recognition*. IEEE, 2015, pp. 665–669.
- [140] S. Li and W. Deng, "Deep facial expression recognition: A survey," *IEEE Transactions on Affective Computing*, 2020.
- [141] R. Weber, J. Li, C. Soladie, and R. Segulier, "A survey on databases of facial macro-expression and micro-expression," in *International Joint Conference on Computer Vision, Imaging and Computer Graphics*. Springer, 2018, pp. 298–325.
- [142] C. Tan, F. Sun, T. Kong, W. Zhang, C. Yang, and C. Liu, "A survey on deep transfer learning," in *International conference on artificial neural networks*. Springer, 2018, pp. 270–279.
- [143] G. Zhao, X. Huang, M. Taini, S. Z. Li, and M. Pietikäinen, "Facial expression recognition from near-infrared videos," *Image and Vision Computing*, vol. 29, no. 9, pp. 607–619, 2011.
- [144] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews, "The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression," in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition-Workshops*. IEEE, 2010, pp. 94–101.
- [145] M. Lyons, S. Akamatsu, M. Kamachi, and J. Gyoba, "Coding facial expressions with gabor wavelets," in *IEEE international conference on automatic face and gesture recognition*. IEEE, 1998, pp. 200–205.
- [146] N. Aifanti, C. Papachristou, and A. Delopoulos, "The mug facial expression database," in *11th International Workshop on Image*

- Analysis for Multimedia Interactive Services WIAMIS 10*. IEEE, 2010, pp. 1–4.
- [147] M. F. Valstar, E. Sánchez-Lozano, J. F. Cohn, L. A. Jeni, J. M. Girard, Z. Zhang, L. Yin, and M. Pantic, “Fera 2017-addressing head pose in the third facial expression recognition and analysis challenge,” in *IEEE international conference on automatic face and gesture recognition*. IEEE, 2017, pp. 839–847.
- [148] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, “Unpaired image-to-image translation using cycle-consistent adversarial networks,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [149] B. Allaert, I. M. Bilasco, and C. Djeraba, “Micro and macro facial expression recognition using advanced local motion patterns,” *IEEE Transactions on Affective Computing*, 2019.
- [150] R. Asmara, P. Choirina, C. Rahmad, A. Setiawan, F. Rahutomo, R. Yusron, and U. Rosiani, “Study of drmf and asm facial landmark point for micro expression recognition using klt tracking point feature,” in *Journal of Physics: Conference Series*, vol. 1402, no. 7. IOP Publishing, 2019, p. 077039.
- [151] H. Lu, K. Kpalma, and J. Ronsin, “Motion descriptors for micro-expression recognition,” *Signal Processing: Image Communication*, vol. 67, pp. 108–117, 2018.
- [152] L. Zhou, Q. Mao, and L. Xue, “Dual-inception network for cross-database micro-expression recognition,” in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2019, pp. 1–5.
- [153] M. Verma, S. K. Vipparthi, G. Singh, and S. Murala, “Learnnet: Dynamic imaging network for micro expression recognition,” *IEEE Transactions on Image Processing*, vol. 29, pp. 1618–1627, 2019.
- [154] Z. Xia, X. Hong, X. Gao, X. Feng, and G. Zhao, “Spatiotemporal recurrent convolutional networks for recognizing spontaneous micro-expressions,” *IEEE Transactions on Multimedia*, vol. 22, no. 3, pp. 626–640, 2019.
- [155] Z. Xia, H. Liang, X. Hong, and X. Feng, “Cross-database micro-expression recognition with deep convolutional networks,” in *ACM International Conference on Biometric Engineering and Applications*. ACM, 2019, pp. 56–60.
- [156] L. Lu, L. Tavabi, and M. Soleymani, “Self-supervised learning for facial action unit recognition through temporal consistency,” in *British Machine Vision Conference*. BMVA, 2020.
- [157] C. Wang, J. Zeng, S. Shan, and X. Chen, “Multi-task learning of emotion recognition and facial action unit detection with adaptively weights sharing network,” in *IEEE International Conference on Image Processing*. IEEE, 2019, pp. 56–60.
- [158] C. Ma, L. Chen, and J. Yong, “Au r-cnn: Encoding expert prior knowledge into r-cnn for action unit detection,” *Neurocomputing*, vol. 355, pp. 35–47, 2019.
- [159] I. O. Ertugrul, J. F. Cohn, L. A. Jeni, Z. Zhang, L. Yin, and Q. Ji, “Cross-domain au detection: Domains, learning approaches, and measures,” in *IEEE International Conference on Automatic Face and Gesture Recognition*, 2019, pp. 1–8.
- [160] W. Li, F. Abtahi, and Z. Zhu, “Action unit detection with region adaptation, multi-labeling learning and optimal temporal fusing,” in *IEEE International Conference on Computer Vision and Pattern Recognition*. IEEE, 2017, pp. 1841–1850.
- [161] W. Li, F. Abtahi, Z. Zhu, and L. Yin, “Eac-net: A region-based deep enhancing and cropping approach for facial action unit detection,” in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2017, pp. 103–110.
- [162] Y. Fan, J. C. Lam, and V. O. K. Li, “Facial action unit intensity estimation via semantic correspondence learning with dynamic graph convolution,” in *AAAI Conference on Artificial Intelligence*. AAAI, 2020, pp. 12701–12708.
- [163] Y. Li, J. Zeng, and S. Shan, “Learning representations for facial actions from unlabeled videos,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [164] R. Zhi, M. Liu, and D. Zhang, “A comprehensive survey on automatic facial action unit analysis,” *The Visual Computer*, vol. 36, no. 5, pp. 1067–1093, 2020.
- [165] B. Martínez, M. F. Valstar, B. Jiang, and M. Pantic, “Automatic analysis of facial actions: A survey,” *IEEE Transactions on Affective Computing*, vol. 10, no. 3, pp. 325–347, 2019.
- [166] Z. Liu, J. Dong, C. Zhang, L. Wang, and J. Dang, “Relation modeling with graph convolutional networks for facial action unit detection,” in *International Conference on Multimedia Modeling*. Springer, 2020, pp. 489–501.
- [167] G. Li, X. Zhu, Y. Zeng, Q. Wang, and L. Lin, “Semantic relationships guided representation learning for facial action unit recognition,” in *AAAI Conference on Artificial Intelligence*, vol. 33. AAAI, 2019, pp. 8594–8601.
- [168] Z.-M. Chen, X.-S. Wei, P. Wang, and Y. Guo, “Multi-label image recognition with graph convolutional networks,” in *IEEE International Conference on Computer Vision and Pattern Recognition*. IEEE, 2019, pp. 5177–5186.
- [169] W.-S. Chu, F. De la Torre, and J. F. Cohn, “Learning spatial and temporal cues for multi-label facial action unit detection,” in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2017, pp. 25–32.
- [170] K. Zhao, W.-S. Chu, F. De la Torre, J. F. Cohn, and H. Zhang, “Joint patch and multi-label learning for facial action unit and holistic expression recognition,” *IEEE Transactions on Image Processing*, vol. 25, no. 8, pp. 3931–3946, 2016.
- [171] Y.-J. Liu, J.-K. Zhang, W.-J. Yan, S.-J. Wang, G. Zhao, and X. Fu, “A main directional mean optical flow feature for spontaneous micro-expression recognition,” *IEEE Transactions on Affective Computing*, vol. 7, no. 4, pp. 299–310, 2015.
- [172] L. Lo, H.-X. Xie, H.-H. Shuai, and W.-H. Cheng, “Mer-gcn: Micro-expression recognition based on relation modeling with graph convolutional networks,” in *IEEE Conference on Multimedia Information Processing and Retrieval*. IEEE, 2020, pp. 79–84.
- [173] Q. Li, J. Yu, T. Kurihara, H. Zhang, and S. Zhan, “Deep convolutional neural network with optical flow for facial micro-expression recognition,” *Journal of Circuits, Systems and Computers*, vol. 29, no. 01, p. 2050006, 2020.
- [174] W. Merghani and M. H. Yap, “Adaptive mask for region-based facial micro-expression recognition,” in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2020, pp. 428–433.
- [175] M. A. Takalkar, M. Xu, and Z. Chaczko, “Manifold feature integration for micro-expression recognition,” *Multimedia Systems*, vol. 26, no. 5, pp. 535–551, 2020.
- [176] C. H. Lim and K. M. Goh, “Fuzzy qualitative approach for micro-expression recognition,” in *Asia-Pacific Signal and Information Processing Association Annual Summit and Conference*. IEEE, 2017, pp. 1669–1674.
- [177] T.-Y. Pan, Y.-Z. Dai, M.-C. Hu, and W.-H. Cheng, “Furniture style compatibility recommendation with cross-class triplet loss,” *Multimedia Tools and Applications*, vol. 78, no. 3, pp. 2645–2665, 2019.
- [178] A. J. Smola and B. Schölkopf, “A tutorial on support vector regression,” *Statistics and computing*, vol. 14, no. 3, pp. 199–222, 2004.
- [179] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf, “Deepface: Closing the gap to human-level performance in face verification,” in *IEEE International Conference on Computer Vision and Pattern Recognition*. IEEE, 2014, pp. 1701–1708.
- [180] S. Lalitha and K. Thyagarajan, “Micro-facial expression recognition based on deep-rooted learning algorithm,” *International Journal of Computational Intelligence Systems*, vol. 12, pp. 903–913, 2019.
- [181] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal loss for dense object detection,” in *IEEE International Conference on Computer Vision*. IEEE, 2017, pp. 2980–2988.
- [182] Z. Lai, R. Chen, J. Jia, and Y. Qian, “Real-time micro-expression recognition based on resnet and atrous convolutions,” *Journal of Ambient Intelligence and Humanized Computing*, pp. 1–12, 2020.
- [183] W. Xie, L. Shen, and J. Duan, “Adaptive weighting of handcrafted feature losses for facial expression recognition,” *IEEE Transactions on Cybernetics*, 2019.
- [184] T. Zhang, Y. Zong, W. Zheng, C. P. Chen, X. Hong, C. Tang, Z. Cui, and G. Zhao, “Cross-database micro-expression recognition: A benchmark,” *IEEE Transactions on Knowledge and Data Engineering*, 2020.
- [185] X. Zhang, T. Xu, W. Sun, and A. Song, “Multiple source domain adaptation in micro-expression recognition,” *Journal of Ambient Intelligence and Humanized Computing*, pp. 1–16, 2020.
- [186] L. Li, X. Zhou, Y. Zong, W. Zheng, X. Chen, J. Shi, and P. Song, “Unsupervised cross-database micro-expression recognition using target-adapted least-squares regression,” *IEICE Transactions on Information and Systems*, vol. 102, no. 7, pp. 1417–1421, 2019.
- [187] G. Wilson and D. J. Cook, “A survey of unsupervised deep domain adaptation,” *ACM Transactions on Intelligent Systems and Technology*, vol. 11, no. 5, pp. 1–46, 2020.

- [188] L. Zhang, "Transfer adaptation learning: A decade survey," *arXiv preprint arXiv:1903.04687*, 2019.
- [189] A. J. R. Kumar, R. Theagarajan, O. Peraza, and B. Bhanu, "Classification of facial micro-expressions using motion magnified emotion avatar images." in *IEEE International Conference on Computer Vision and Pattern Recognition Workshops*. IEEE, 2019, pp. 12–20.
- [190] L. A. Zadeh, "Fuzzy sets," *Information and Control*, vol. 8, no. 3, pp. 338–353, 1965.
- [191] M. Chen, H. T. Ma, J. Li, and H. Wang, "Emotion recognition using fixed length micro-expressions sequence and weighting method," in *IEEE International Conference on Real-time Computing and Robotics*. IEEE, 2016, pp. 427–430.
- [192] H. A. Efenbein and N. Ambady, "Is there an in-group advantage in emotion recognition?" *Psychological Bulletin*, vol. 128, pp. 243–9, 2002.
- [193] S. G. Young and J. P. Wilson, "A minimal ingroup advantage in emotion identification confidence," *Cognition and Emotion*, vol. 32, no. 1, pp. 192–199, 2018.
- [194] E. R. Tuminello and D. Davidson, "What the face and body reveal: In-group emotion effects and stereotyping of emotion in african american and european american children," *Journal of Experimental Child Psychology*, vol. 110, no. 2, pp. 258–274, 2011.
- [195] P. Thibault, P. Bourgeois, and U. Hess, "The effect of group-identification on emotion recognition: The case of cats and basketball players," *Journal of Experimental Social Psychology*, vol. 42, no. 5, pp. 676–683, 2006.
- [196] S. G. Young and K. Hugenberg, "Mere social categorization modulates identification of facial expressions of emotion." *Journal of Personality and Social Psychology*, vol. 99, no. 6, p. 964, 2010.
- [197] S. Huang and S. Han, "Shared beliefs enhance shared feelings: religious/irreligious identifications modulate empathic neural responses," *Social Neuroscience*, vol. 9, no. 6, pp. 639–649, 2014.
- [198] C. Prado, D. Mellor, L. K. Byrne, C. Wilson, X. Xu, and H. Liu, "Facial emotion recognition: a cross-cultural comparison of chinese, chinese living in australia, and anglo-australians," *Motivation and Emotion*, vol. 38, no. 3, pp. 420–428, 2014.
- [199] Y. Xie, C. Zhong, F. Zhang, and Q. Wu, "The ingroup disadvantage in the recognition of micro-expressions," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2019, pp. 1–5.
- [200] M. Fölster, U. Hess, and K. Werheid, "Facial age affects emotional expression decoding," *Frontiers in psychology*, vol. 5, p. 30, 2014.
- [201] S. Chen, J. Wang, Y. Chen, Z. Shi, X. Geng, and Y. Rui, "Label distribution learning on auxiliary label space graphs for facial expression recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020*, pp. 13984–13993.
- [202] K. Wang, X. Peng, J. Yang, S. Lu, and Y. Qiao, "Suppressing uncertainties for large-scale facial expression recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020*, pp. 6897–6906.
- [203] M. K. Mandal, H. S. Asthana, and R. Pandey, "Asymmetry in emotional face: Its role in intensity of expression," *The Journal of Psychology*, vol. 129, no. 2, pp. 235–241, 1995.
- [204] W. G. Dopson, B. E. Beckwith, D. M. Tucker, and P. C. Bullard-Bates, "Asymmetry of facial expression in spontaneous emotion," *Cortex*, vol. 20, no. 2, pp. 243–251, 1984.
- [205] H. Chen, X. Liu, X. Li, H. Shi, and G. Zhao, "Analyze spontaneous gestures for emotional stress state recognition: A micro-gesture dataset and analysis with deep learning," in *IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 2019, pp. 1–8.



**Hong-Xia Xie** received the B.S. degree in Internet of Things from the Zhengzhou University of Aeronautics in 2016 and received the M.S. degree in communication and information systems, Fujian Normal University, China, in 2019. She is now pursuing a Ph.D. degree in Institute of Electronics, National Chiao Tung University, Taiwan. Her research interests include micro-expression recognition, emotion recognition and deep learning.



**Ling Lo** received the B.S. degree from Department of Electronics Engineering, National Chiao Tung University (NCTU), Hsinchu, Taiwan, R.O.C., in 2019, and now she is pursuing a Ph.D degree in Institute of Electronics, NCTU. Her current research interests include deep learning and computer vision. Recently her work focuses specifically on facial and micro-expression recognition.



**Hong-Han Shuai** received the B.S. degree from the Department of Electrical Engineering, National Taiwan University (NTU), Taipei, Taiwan, R.O.C., in 2007, the M.S. degree in computer science from NTU in 2009, and the Ph.D. degree from Graduate Institute of Communication Engineering, NTU, in 2015. He is now an assistant professor in NCTU. His research interests are in the area of multimedia processing, machine learning, social network analysis, and data mining. His works have appeared in top-tier

conferences such as MM, CVPR, AAAI, KDD, WWW, ICDM, CIKM and VLDB, and top-tier journals such as TKDE, TMM and JIOT. Moreover, he has served as the PC member for international conferences including MM, AAAI, IJCAI, WWW, and the invited reviewer for journals including TKDE, TMM, JVCi and JIOT.



**Wen-Huang Cheng** is Professor with the Institute of Electronics, National Chiao Tung University (NCTU), Hsinchu, Taiwan. He is also Jointly Appointed Professor with the Artificial Intelligence and Data Science Program, National Chung Hsing University (NCHU), Taichung, Taiwan. Before joining NCTU, he led the Multimedia Computing Research Group at the Research Center for Information Technology Innovation (CITI), Academia Sinica, Taipei, Taiwan, from 2010 to 2018. His current research interests include multimedia, artificial intelligence, computer vision, and machine learning.

He has actively participated in international events and played important leading roles in prestigious journals and conferences and professional organizations, like Associate Editor for IEEE Transactions on Multimedia, General co-chair for IEEE ICME (2022) and ACM ICMR (2021), Chair-Elect for IEEE MSA technical committee, governing board member for IAPR. He has received numerous research and service awards, including the 2018 MSRA Collaborative Research Award, the 2017 Ta-Yu Wu Memorial Award from Taiwan's Ministry of Science and Technology (the highest national research honor for young Taiwanese researchers under age 42), the 2017 Significant Research Achievements of Academia Sinica, the Top 10% Paper Award from the 2015 IEEE MMSP, and the K. T. Li Young Researcher Award from the ACM Taipei/Taiwan Chapter in 2014. He is IET Fellow and ACM Distinguished Member.