

Hierarchical Space-Time Attention for Micro-Expression Recognition

Haihong Hao¹, Shuo Wang^{1*}, Huixia Ben², Yanbin Hao¹, Yansong Wang³ and Weiwei Wang³

¹University of Science and Technology of China

²Hefei University of Technology

³Chery HuiYin Motor Finance Service Co.,Ltd

haohaihong@mail.ustc.edu.cn, {shuowang.edu,huixiabn}@gmail.com, haoyanbin@hotmail.com,
{wangyansong,wangweiwei}@cheryfs.cn

Abstract

Micro-expression recognition (MER) aims to recognize the short and subtle facial movements from the Micro-expression (ME) video clips, which reveal real emotions. Recent MER methods mostly only utilize special frames from ME video clips or extract optical flow from these special frames. However, they neglect the relationship between movements and space-time, while facial cues are hidden within these relationships. To solve this issue, we propose the Hierarchical Space-Time Attention (HSTA). Specifically, we first process ME video frames and special frames or data parallelly by our cascaded Unimodal Space-Time Attention (USTA) to establish connections between subtle facial movements and specific facial areas. Then, we design Crossmodal Space-Time Attention (CSTA) to achieve a higher-quality fusion for crossmodal data. Finally, we hierarchically integrate USTA and CSTA to grasp the deeper facial cues. Our model emphasizes temporal modeling without neglecting the processing of special data, and it fuses the contents in different modalities while maintaining their respective uniqueness. Extensive experiments on the four benchmarks show the effectiveness of our proposed HSTA. Specifically, compared with the latest method on the CASME3 dataset, it achieves about 3% score improvement in seven-category classification. Code is available at https://github.com/OceanSummerDay/HSTA_MER.

1 Introduction

Micro-expression recognition (MER) is a challenging task in affective computing, due to the subtlety and brief duration (typically 1/25 to 1/3 second) of micro-expressions (MEs), making them difficult to capture [Van Quang *et al.*, 2019]. However, MEs are brief, subtle, spontaneous, and involuntary emotional expressions that convey genuine emotions. Although humans can hide their emotions in certain situations, MEs are difficult to hide and can inadvertently reveal true

feelings [Ekman, 2009]. Thus, the MER is critically important in many specific scenes, such as criminal interrogation, clinical diagnosis, and financial risk control.

Recent advances in deep learning have shown promising results in MER [Verma *et al.*, 2019], surpassing traditional hand-crafted methods and emerging as the dominant technique. However, the majority of open-source methods overly depend on special frames or specific data, failing to fully utilize the fundamental temporal characteristics of MEs as special short video sequences. Techniques like CapsuleNet [Van Quang *et al.*, 2019], MMNet [Li *et al.*, 2022a], and OFF-ApexNet [Gan *et al.*, 2019] rely on special frames (Apex and Onset frame), neglecting the most intrinsic video nature of ME data. Similarly, methods such as Dual-ATME [Liong and Wong, 2017], DSSN [Khor *et al.*, 2019], and Bi-WOOF [Liong *et al.*, 2018] rely heavily on special optical flow, which is highly sensitive to environmental factors and is unable to fully capture all the details of non-rigid and complex facial muscle movements. These methods commonly lack the ability to model temporal information effectively, leading to a disconnection between facial movements and specific facial areas. The model only associate facial expressions with static facial features, rather than interpreting ME as a series of continuous facial motions. Actually, MER requires a more comprehensive modeling of space-time information to better capture the correspondence between dynamic facial muscle movements and specific facial regions, which is essential for accurately recognizing micro-expressions [Li *et al.*, 2022a].

For temporal (space-time) modeling, we note that existing video processing methods [Wang *et al.*, 2023a; Rehman *et al.*, 2022; Qian *et al.*, 2021] have achieved great success on conventional action recognition datasets. For both action recognition and MER tasks, it is essential to establish the relationship between movements and time, as well as to identify the connections between specific actions and certain areas or scenes. Inspired by these action recognition work, we incorporate these temporal processing concepts into the MER task. Specifically, we extract both spatial and temporal information by our designed Unimodal Space-Time Attention (USTA) from a small number of uniformly sampled micro-expression frames. This is because the key is to establish connections between these tiny muscle movements and specific small areas of face [Li *et al.*, 2022a]. By incorporating and modeling temporal information, USTA can capture tiny

*Corresponding author

movements over time since all tokens (features) are able to interact with each other [Vaswani *et al.*, 2017] in our USTA. In other words, tokens of the same facial area can interact with their counterparts at different moments and capture subtle facial movements over time. Meanwhile, tokens for different facial areas can also interact with each other and associate minute movements with specific areas of the face.

Based on the studies from past MER methods, we also recognize that the special frames or data are important for recognizing the MEs. Thus, we design the **Crossmodal Space-Time Attention (CSTA)** to effectively fuse the information from different modalities after the unimodal USTA calculations. Specifically, we utilize a symmetrical structure based on cross-attention to capture the inner connection between ME video frames and special data (such as special frames or optical flow). In CSTA, two different modalities of data (temporal video data and special data) are captured and integrated by our cross-attention calculations. Meanwhile, these data also complement each other to assist in the expression of emotions. After passing through our CSTA, the class tokens contain information from the other type of data, while the remaining tokens retain their original data. Thus, the combination of USTA and CSTA achieves effective fusion while maintaining the distinctiveness of different modalities. To adequately grasp the deeper facial cues of motion and time, we extend this combination into a hierarchical structure, named **Hierarchical Space-Time Attention (HSTA)**. In HSTA, the USTA and CSTA are stacked in an orderly manner to capture the MEs effectively. Meanwhile, the adaptable layer design of HSTA enhances its generalization capabilities. Our contributions can be summarized as threefold:

- We design unimodal space-time attention (USTA) to demonstrate the significance of temporal information in MER. Meanwhile, we propose crossmodal space-time attention (CSTA) to complement data of different types, thereby enriching the content within each modality data.
- We extend the USTA and CSTA into the hierarchical structure to fuse the contents in different modalities and grasp the deeper facial cues of motion and temporal data.
- We demonstrate the effectiveness of our proposed HSTA on four MER datasets. Our method outperforms existing methods and achieves state-of-the-art (SOTA) results.

2 Related Work

In this section, we first briefly introduce common solutions for MER tasks and then we list the related attention calculations. Finally, we enumerate the differences between our methods and those of related methods.

2.1 Micro-expression Recognition

MER task is to classify the facial MEs in a video. Related technologies are mainly divided into two categories. The **first category** relies solely on special data (e.g., special frames, optical flow) and feeds them into a 2D CNN. This approach is the one most commonly adopted in most current open-source work. They are highly sensitive to environmental factors [Zhou *et al.*, 2019] such as changes in lighting, shadows. In the **other category**, temporal MEs are input into

the model and then learned by a time series network or a 3D CNN ([Reddy *et al.*, 2019; Khor *et al.*, 2018]). However, due to information redundancy, it becomes difficult to focus on the most important features. Consequently, the performance of these methods is probably not as effective as that of those using only special data, leading to their relative underestimation. Although there have been attempts to model temporal sequences of optical flow [Li *et al.*, 2019], due to the subtlety of MEs and minimal changes between adjacent frames, extracting optical flow results in significant noise, leading to poor performance. In reality, we can model temporal sequences while minimizing redundancy, without neglecting the importance of special data.

Modeling spatial and temporal information is key to processing sequential data. Currently, the main methods include 3D CNNs, Video Vision Transformers [Arnab *et al.*, 2021] (ViTs), or a combination of both [Wodajo and Atnafu, 2021]. A 3D CNN extends the standard convolution operation from 2D to 3D, allowing it to capture motion information embedded within consecutive frames of a video by analyzing a sequence as a whole, rather than in isolation. ViTs capture global dependencies in video frames, providing a broader understanding of the scene compared to the local focus of 3D CNNs. With their self-attention mechanism, ViTs dynamically concentrate on relevant parts of the input, which is significant for processing inputs from various modalities. However, a direct application of these video models presents drawbacks: they only model temporal sequences of videos, neglecting the importance of special frames. Compared to methods utilizing optical flow, these methods also lack the capture of the big picture.

2.2 Attention Calculation

As a method similar to [Tong *et al.*, 2022], employing self-attention mechanisms on temporal data enables the capture of subtle movements and their temporal relationships. As methods such as [Khor *et al.*, 2019] have shown, simply adding or concatenating feature vectors from different modalities does not significantly improve performance, so we should find a better way to fuse different modalities. Cross-attention has already been applied to non-sequential data [Wei *et al.*, 2020], exploiting crossmodal relationships and achieving tremendous success. Cross-attention between two different modalities offers a higher quality fusion, integrating the different modalities more effectively.

Based on the analysis of related work, the methods most related to ours are the recently proposed cross-attention [Chen *et al.*, 2021] and dual learning in different modalities [Khor *et al.*, 2019]. Our method differs from theirs in two aspects. First, we design a **new attention strategy** for different temporal MEs data and process them parallel, to establish the relationship between motions and time. Second, we introduce a **new cross-attention** to achieve high-quality fusion for sequential data, still keeping their uniqueness to model connections between different modalities.

3 Approach

In this section, we first briefly revisit the preliminaries of the MER tasks and give an overview of our framework.

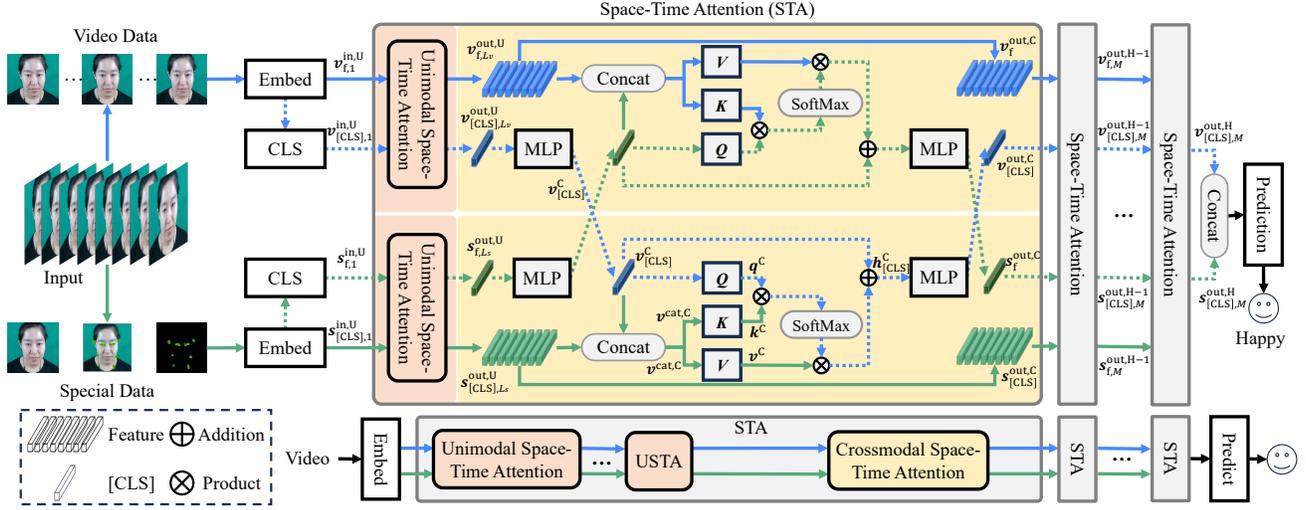


Figure 1: The overview of our Hierarchical Space-Time Attention (HSTA). Our model is a hierarchical structure that captures the underlying emotions in micro-expression videos through multiple attention modules.

Then, we illustrate our Hierarchical Space-Time Attention (HSTA) containing Unimodal Space-Time Attention (USTA) and Crossmodal Space-Time Attention (CSTA). Meanwhile, we introduce the design of hierarchical structures. Finally, we describe the training and inference procedures of our method.

3.1 Preliminaries

The data of the MER task can be divided into two parts: **video frames set** \mathcal{D}_v and **special frames set** \mathcal{D}_s . Specifically, \mathcal{D}_v consists of multiple video clips, and each clip provides the label of its emotional expression. \mathcal{D}_s contain some special frames, such as Apex, Onset, and Offset frames, where the Apex, Onset, and Offset frames represent the moments of maximum expression amplitude, the start, and the end of a micro-expression video. They can be used to assist in the recognition of facial expression content. The goal of the MER task is to use these different formats of data to capture the emotional content.

An overview of our method is depicted in Figure 1. First, we use available embed methods to capture the features and the “[CLS]s” of video and special frames, where “[CLS]” can be regarded as features containing certain classification information. Then we design the unimodal and crossmodal space-time attentions to fuse these different contents. Meanwhile, we expand these attentions into hierarchical structures to obtain a deep emotional expression. Finally, we concatenate the “[CLS]s” from different modalities for prediction.

3.2 Unimodal Space-Time Attention

In our method, video and special frames are calculated simultaneously. Since the Unimodal Space-Time Attention (USTA) calculations of these two parts are similar, we first introduce the operation of USTA and then describe its application in these two parts’ calculations to simplify the description.

As shown in Figure 2(a), our USTA is a **cascaded** structure. To illustrate it in detail, we use the l^{th} layer of calcula-

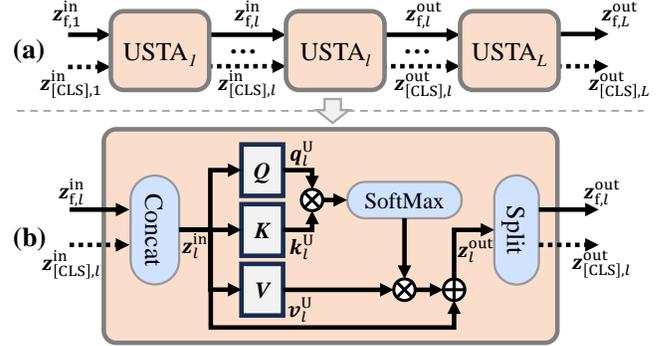


Figure 2: The details of USTA. (a) The cascaded structure of USTA. (b) The calculations of l^{th} layer USTA.

tions as an example to describe each USTA unit. Details are depicted in Figure 2(b), denoting the inputted features and “[CLS]” of the l^{th} layer as $z_{f,l}^{in} \in \mathbb{R}^{N \times d}$ and $z_{[CLS],l}^{in} \in \mathbb{R}^d$, respectively, where N is the number of features and d is the size of the embedding, we first concatenate them into one feature:

$$z_l^{in} = [z_{f,l}^{in} || z_{[CLS],l}^{in}] \in \mathbb{R}^{(N+1) \times d}. \quad (1)$$

Then we design three attention matrices q_l^U , k_l^U , and v_l^U to capture self-similarities from the concatenated feature z_l^{in} :

$$q_l^U = z_l^{in} W_{q,l}^U, k_l^U = z_l^{in} W_{k,l}^U, v_l^U = z_l^{in} W_{v,l}^U, \quad (2)$$

where $W_{q,l}^U$, $W_{k,l}^U$, and $W_{v,l}^U$ are learnable parameters in $\mathbb{R}^{d \times d}$, and q_l^U , k_l^U , and v_l^U are same size to the input z_l^{in} . Third, following the operations in [Vaswani *et al.*, 2017], we employ self-attention to achieve the refinement feature z_l^{out} :

$$z_l^{out} = \text{LayerNorm}(z_l^{in} + \text{SoftMax}(\frac{q_l^U (k_l^U)^T}{\sqrt{d}}) v_l^U), \quad (3)$$

where “LayerNorm” and “SoftMax” are normalization function and activation function. Finally, we split the refined fea-

ture $z_l^{\text{out}} \in \mathbb{R}^{(N+1) \times d}$ into two parts as the outputs of the l^{th} layer USTA:

$$[z_{f,l}^{\text{out}} || z_{[\text{CLS}],l}^{\text{out}}] = \text{Split}(z_l^{\text{out}}), \quad (4)$$

where $z_{f,l}^{\text{out}} \in \mathbb{R}^{N \times d}$ and $z_{[\text{CLS}],l}^{\text{out}} \in \mathbb{R}^d$ are the features and “[CLS]”, respectively. During this operation, the contents of the frames along the time dimension can interact with each other. To conveniently describe our cascaded USTA, we define the calculation of USTA in the l^{th} layer as USTA_l . The different layer calculations can be defined as:

$$[z_{f,l}^{\text{out}} || z_{[\text{CLS}],l}^{\text{out}}] = \text{USTA}_l([z_{f,l}^{\text{in}} || z_{[\text{CLS}],l}^{\text{in}}]). \quad (5)$$

For the cascade structure, we calculate the input and output of different USTA layers as $z_L^{\text{out}} = [z_{f,L}^{\text{out}} || z_{[\text{CLS}],L}^{\text{out}}]$:

$$z_L^{\text{out}} = \begin{cases} [z_{f,1}^{\text{in}} || z_{[\text{CLS}],1}^{\text{in}}], & L = 1, \\ \text{USTA}_L([z_{f,L-1}^{\text{out}} || z_{[\text{CLS}],L-1}^{\text{out}}]), & L > 1, \end{cases} \quad (6)$$

where L is the total layer of the cascaded USTA. Then we illustrate the USTA calculations in different modalities.

For given video frames set \mathcal{D}_v and special frames set \mathcal{D}_s , we use different 3D CNN [Tong *et al.*, 2022] for feature extraction. With setting “[CLS]” token, the embedding of \mathcal{D}_v and \mathcal{D}_s can be represented as $[v_{f,1}^{\text{in,U}} || v_{[\text{CLS}],1}^{\text{in,U}}] \in \mathbb{R}^{(N_v+1) \times d}$ and $[s_{f,1}^{\text{in,U}} || s_{[\text{CLS}],1}^{\text{in,U}}] \in \mathbb{R}^{(N_s+1) \times d}$, where N_v and N_s are the number of features of \mathcal{D}_v and \mathcal{D}_s respectively. Then we define L_v and L_s as the total layer of USTA in these two calculations in Eq. (6). We process them **parallel**. Thus, the USTA outputs from \mathcal{D}_v and \mathcal{D}_s can be calculated as:

$$\begin{aligned} [v_{f,L_v}^{\text{out,U}} || v_{[\text{CLS}],L_v}^{\text{out,U}}] &= \text{USTA}_{L_v}([v_{f,1}^{\text{in,U}} || v_{[\text{CLS}],1}^{\text{in,U}}]), \\ [s_{f,L_s}^{\text{out,U}} || s_{[\text{CLS}],L_s}^{\text{out,U}}] &= \text{USTA}_{L_s}([s_{f,1}^{\text{in,U}} || s_{[\text{CLS}],1}^{\text{in,U}}]). \end{aligned} \quad (7)$$

In this operation, we use the same embedding size d to facilitate our subsequent calculations. The outputs of USTA modules ($[v_{f,L_v}^{\text{out,U}} || v_{[\text{CLS}],L_v}^{\text{out,U}}]$ and $[s_{f,L_s}^{\text{out,U}} || s_{[\text{CLS}],L_s}^{\text{out,U}}]$) are then used as the input for crossmodal space-time attention.

3.3 Crossmodal Space-Time Attention

Crossmodal Space-Time Attention (CSTA) is designed to capture the inner connection between the features of one modal data and the “[CLS]” of another. Since the calculations in our CSTA from different modalities are **symmetrical**, we only discuss the calculation between the features of special data $s_{f,L_s}^{\text{out,U}}$ and “[CLS]” of video data $v_{[\text{CLS}],L_v}^{\text{out,U}}$ for brevity description (in the **lower** part of the golden area of Figure 1, and the other part of the operations are similar).

Firstly, we project “[CLS]” of video data by a MLP to match the embedding dimension of the features of special data and then concatenate them for the subsequent operation:

$$v_{[\text{CLS}]}^{\text{C}} = \text{MLP}_1(v_{[\text{CLS}],L_v}^{\text{out,U}}), v^{\text{cat,C}} = [v_{[\text{CLS}]}^{\text{C}} || s_{f,L_s}^{\text{out,U}}]. \quad (8)$$

Similarly to the calculation in the attention of the USTA, we then also design three attention matrices q^{C} , k^{C} , and v^{C} to capture cross-similarities between “[CLS]” of video data $v_{[\text{CLS}]}^{\text{C}}$ and the concatenated feature $v^{\text{cat,C}}$:

$$q^{\text{C}} = v_{[\text{CLS}]}^{\text{C}} W_q^{\text{C}}, k^{\text{C}} = v^{\text{cat,C}} W_k^{\text{C}}, v^{\text{C}} = v^{\text{cat,C}} W_v^{\text{C}}, \quad (9)$$

where the definition of q^{C} , k^{C} , v^{C} , W_q^{C} , W_k^{C} , and W_v^{C} are similar to that of the USTA. Thus, the refinement “[CLS]” can be calculated by the attention and residual operations:

$$h_{[\text{CLS}]}^{\text{C}} = v_{[\text{CLS}]}^{\text{C}} + \text{SoftMax}\left(\frac{q^{\text{C}}(k^{\text{C}})^{\top}}{\sqrt{d}}\right)v^{\text{C}}. \quad (10)$$

Finally, we map $h_{[\text{CLS}]}^{\text{C}}$ by using a new MLP to keep the output dimensions consistent:

$$v_{[\text{CLS}]}^{\text{out,C}} = \text{MLP}_2(h_{[\text{CLS}]}^{\text{C}}). \quad (11)$$

Meanwhile, the output features of special data is consistent with the original inputted features:

$$s_f^{\text{out,C}} = s_{f,L_s}^{\text{out,U}}. \quad (12)$$

Another part has similar operations. Thus, given the features of video data $v_{f,L_v}^{\text{out,U}}$ and “[CLS]” of special data $s_{[\text{CLS}],L_s}^{\text{out,U}}$, the output of CSTA can be calculated as $v_f^{\text{out,C}}$ and $s_{[\text{CLS}]}^{\text{out,C}}$, respectively. Thus, the calculations of the whole CSTA can be summarized as:

$$\begin{aligned} (v_{[\text{CLS}]}^{\text{out,C}}, v_f^{\text{out,C}}, s_{[\text{CLS}]}^{\text{out,C}}, s_f^{\text{out,C}}) &= \\ \text{CSTA}(v_{[\text{CLS}],L_v}^{\text{out,U}}, v_{f,L_v}^{\text{out,U}}, s_{[\text{CLS}],L_s}^{\text{out,U}}, s_{f,L_s}^{\text{out,U}}). \end{aligned} \quad (13)$$

3.4 Hierarchical Learning

We refer to a single-layer structure combining USTA and CSTA as Space-Time Attention (STA). In our **hierarchical** structure, we define each Hierarchical Space-Time Attention (HSTA) module as containing multiple USTA modules and one CSTA module:

$$\text{HSTA}_m = [(\text{USTA}_{L_v}, \text{USTA}_{L_s}); \text{CSTA}] \quad (14)$$

where m^{th} is one calculation layer of HSTA. Thus, the outputs set ($z_M^{\text{H}} = [v_{[\text{CLS}],M}^{\text{out,H}}, v_{f,M}^{\text{out,H}}, s_{[\text{CLS}],M}^{\text{out,H}}, s_{f,M}^{\text{out,H}}]$) of the whole operations of our M layer HSTA are summarized as:

$$z_M^{\text{H}} = \begin{cases} v_{[\text{CLS}]}^{\text{out,C}}, v_f^{\text{out,C}}, s_{[\text{CLS}]}^{\text{out,C}}, s_f^{\text{out,C}}, & M = 1, \\ \text{HSTA}_M(z_{M-1}^{\text{H}}), & M > 1. \end{cases} \quad (15)$$

For training, we select the “[CLS]s” of two different modalities to predict the MEs. Specifically, given the outputs of the last HSTA of video and special data as $v_{[\text{CLS}],M}^{\text{out,H}}$ and $s_{[\text{CLS}],M}^{\text{out,H}}$, respectively, we first concatenate them and use **Mean Squared Error** (MSE) loss to measure distance between the prediction and the label Y :

$$\mathcal{L} = \text{MSE}(P_{\theta}([v_{[\text{CLS}],M}^{\text{out,H}} || s_{[\text{CLS}],M}^{\text{out,H}}]), Y), \quad (16)$$

where P_{θ} is the prediction function with parameter θ .

4 Experiments

In this section, we conduct experiments to evaluate the effectiveness of our proposed method. First, we introduce the experimental settings. Then analyze the effects of different

modules of our method. Finally, we compare other state-of-the-art methods with ours. Our experiments are intended to address the following research questions (RQs):

RQ1: What are the effects of unimodal and crossmodal space-time attentions?

RQ2: How does hierarchical learning influence the micro-expression recognition results?

RQ3: How does the performance comparison between our method and the state-of-the-art methods?

4.1 Experimental Settings

Dataset

We evaluate our method on four benchmark datasets. Specifically, **CASME3** [Li *et al.*, 2022b] is the largest MER dataset. It includes about 1,000 manually annotated MEs with seven expressions (“Happiness”, “Anger”, “Sad”, “Surprise”, “Fear”, “Disgust” and “Others”). **CASME II** [Yan *et al.*, 2014], **SMIC** [Li *et al.*, 2013], and **SAMM** [Davison *et al.*, 2016] contain 247, 161, and 159 MEs videos, respectively. For the CASME II, SMIC, and SAMM datasets, we follow the experimental settings in the MEGC2019 Challenge [See *et al.*, 2019] to map these datasets into three general categories: “Negative”, “Positive”, and “Surprise”. More details of experimental datasets can be found in our cited work.

Evaluation Metrics

Following the metrics in MER2019 challenge [See *et al.*, 2019], we use three common evaluation metrics **unweighted F1-scores** (UF1), **unweighted average recall** (UAR), and extra **accuracy** (ACC) to measure the effectiveness of our method. These metrics are calculated as:

$$\text{UF1} = \sum_{i=1}^C ((2\text{TP}_i) / (2\text{TP}_i + \text{FP}_i + \text{FN}_i)) / C,$$

$$\text{UAR} = \sum_{i=1}^C (\text{TP}_i / N_i) / C, \quad (17)$$

$$\text{ACC} = (\sum_{i=1}^C \text{TP}_i) / (\sum_{i=1}^C N_i),$$

where C is the total number of the micro-expression categories, TP_i , FP_i , and FN_i is the number of true positive, false positive, and false negative samples of the i^{th} category, respectively, N_i is the sample number of the i^{th} category.

During the testing phase, we employ the **leave-one-subject-out cross-validation** strategy (LOSO) [Li *et al.*, 2022a] to compare with other methods. Specifically, for each iteration, one of the subsets is randomly selected to be used as the test set, and the remaining subsets are used as the training set. However, this strategy is time-consuming. Therefore, we use the more efficient and less resource-intensive **K -fold cross-validation** strategy [Zhao *et al.*, 2023] in our ablation studies to help us find the appropriate parameters.

Implementation Details

Due to the significant variation in sample numbers across different categories in MEs datasets, we employ balanced sampling during training to ensure uniform exposure of the model to approximately the same number of samples from each category. In our experiments, we use a batch size of 32, conduct 150 training epochs, set the base learning rate is $5e-5$, and the weight decay is 0.05. To better train the parameters, we employ a warm-up learning strategy in the first 5 epochs. Specifically, we set the initial learning rate $1e-6$ and gradually increase until it reaches the basic learning rate.

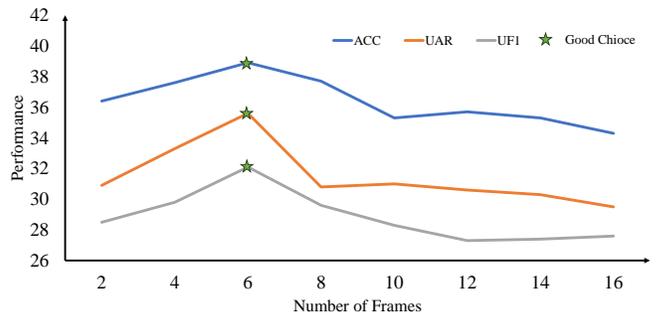


Figure 3: Performances of frames with different numbers in USTA.

#L	L_v			L_s		
	ACC	UAR	UF1	ACC	UAR	UF1
0	15.8	14.3	16.1	15.8	14.3	16.1
1	37.5	32.4	31.1	37.9	32.1	29.4
2	36.5	32.6	29.8	38.8	32.6	30.8
4	38.0	33.1	29.8	39.3	32.9	31.0
6	41.6	33.6	32.0	39.2	33.9	31.0
8	42.1	36.7	33.5	39.5	35.5	32.2
10	40.0	35.5	32.0	40.9	34.1	31.8
12	40.1	35.6	31.8	39.5	34.0	31.4

Table 1: Performance of different L_s and L_v .

4.2 Ablation Studies

In the ablation study, we use CASME3 in seven-category classification to evaluate the effectiveness of the different components of HSTA. To minimize computational costs, all experiments in this study are conducted using the **5-fold cross-validation** strategy. Within each fold, we randomly designate 20% of the MEs as the test set, while the remaining 80% constitute the training set, each MEs video appears in the test set once and only once.

The Effectiveness of USTA (RQ1.a)

We conduct experiments solely using USTA. For efficiency and to expedite subsequent computational processes, we initially employ a single-layer USTA to determine the optimal number of uniformly sampled frames for video frames set \mathcal{D}_v . The results are presented in Figure 3. As illustrated in the figure, varying the frame count significantly impacts the experimental outcomes. Performance initially improves with an increase in frame count but subsequently diminishes. An excessive number of frames not only increases computational overhead but also reduces performance. This can be attributed to the introduction of redundant information by too many frames, which may skew the capture of finer details. Therefore, in our subsequent experiments, we validate our model using **six** frames.

Next, we verify the impact with and without USTA. We use the Apex frame and Onset frame as special frames in our experiments. Without USTA means we set the total cascaded USTA layer count $L_v = 0$ for video frames set \mathcal{D}_v in Eq. (7) and $L_s = 0$ for special frames set \mathcal{D}_s . We also investigate

Setting	CSTA	L_v	L_s	ACC	UAR	UF1
a1	×	0	0	15.8	14.3	16.1
a2	✓	0	0	35.5	30.9	28.2
a3	×	1	0	37.5	32.4	31.1
a4	×	0	1	37.9	32.5	29.4
a5	×	1	1	37.1	32.1	29.7
a6	✓	1	1	40.2	33.9	31.6

Table 2: Performances of CSTA and its coordination.

the relationship of different USTA layer counts by setting the number of layers from 1 to 12. The specific results are shown in Table 1, # L is the number of L_v or L_s . When no USTA layers are included, the performance is extremely poor. However, introducing just a single-layer USTA led to a significant performance improvement (ACC from 16% \rightarrow 38%). The improvement observed after the inclusion of USTA indicates the critical importance of temporal modeling in MER. It also demonstrates the effectiveness of our USTA. As L_v or L_s increase, performance improves continuously. However, when L_v or L_s become too large, it leads to overfitting, resulting in an initial increase followed by a decrease.

The Effectiveness of CSTA (RQ1.b)

We utilize only CSTA structure and skip the USTA to verify our CSTA. Then we use a structure where both L_v and L_s are 1 to compare with the structure USTA coordinates with CSTA. In Table 2, the comparison of settings a1 and a2 reveals that employing only the CSTA structure yields a certain level of performance compared to without CSTA. The single application of CSTA alone does not surpass using USTA alone as observed in settings a3 and a4. Whether it is USTA or CSTA, single use results in unsatisfactory performance. Furthermore, the setting a5, which utilizes both video frames and special frames without a robust fusion mechanism, leads to a decline in performance. In contrast, the combination of CSTA with USTA in setting a6, leads to a significant improvement. This not only demonstrates the effectiveness of CSTA but also indicates that CSTA and USTA work well together. We believe that only after USTA establishes connections between subtle facial movements and specific facial areas, CSTA can achieve higher-quality fusion. Here different modalities also learn aspects of content from each other, enriching their “[CLS]” token’s content.

The Effectiveness of Hierarchical Learning (RQ2)

Here we evaluate the performance of hierarchical HSTA, as detailed in Table 3. Comparisons with different HSTA layer settings b1, b2, and b3 in Table 3 reveal that HSTA methods outperform the single-layer STA unit. The performance of the USTA is influenced by L_v and L_s . Maintaining L_s constant while increasing L_v initially enhances performance then decreases (as shown from b3 to b7). In contrast, keeping L_v constant while increasing L_s leads to a continuous decrease in performance. Our experiments suggest that performance is better when $L_v > L_s$ compared to $L_v < L_s$. $L_v = L_s$ also provides satisfactory results, likely due to the larger volume of information from video frames matched by greater computational capacity when $L_v \geq L_s$. Then we explore the rela-

Setting	HSTA	L_s	L_v	ACC	UAR	UF1
b1	1	1	1	40.2	33.9	31.6
b2	2	■ 1	1	40.8	34.1	32.2
b3	3	■ 1	1	41.3	34.9	31.3
b4	3	● 1	2	40.7	35.1	32.0
b5	3	1	3	41.0	36.1	33.0
b6	3	1	4	40.3	36.0	33.2
b7	3	1	5	40.7	35.8	32.6
b8	3	2	1	39.8	34.3	31.6
b9	3	▲ 2	2	42.4	36.6	34.8
b10	3	3	1	41.2	34.5	32.2
b11	3	3	2	41.0	35.3	30.8
b12	3	3	3	40.8	35.2	32.8

Table 3: Performances of HSTA and its coordination.

M	■ $L_v=1$ $L_s=1$			● $L_v=2$ $L_s=1$			▲ $L_v=2$ $L_s=2$		
	ACC	UAR	UF1	ACC	UAR	UF1	ACC	UAR	UF1
1	40.2	33.9	31.6	41.0	35.7	31.2	40.7	33.4	31.1
2	40.8	34.1	32.2	42.2	36.4	33.5	40.0	34.5	31.2
3	41.3	34.9	31.3	40.7	35.1	32.0	42.4	36.6	34.8
4	40.5	34.8	31.9	42.6	38.0	35.1	42.2	36.9	33.6
5	41.2	36.4	33.6	40.9	36.0	32.5	43.1	37.4	33.6
6	40.7	36.4	32.9	41.9	36.2	32.9	42.1	35.9	33.0
7	41.1	36.0	33.1	42.8	37.8	35.1	41.1	36.1	32.4
8	42.8	38.6	35.1	41.0	36.0	33.8	40.6	36.3	32.6
10	42.6	35.7	33.7	41.8	36.8	33.6	41.5	36.4	33.0
12	42.4	38.5	35.1	41.3	36.1	33.4	42.2	37.0	33.3

Table 4: Performances with different numbers of HSTA.

tionship between the HSTA layers value M and performance using three combinations of L_v and L_s : $\langle \blacksquare L_v = 1, L_s = 1 \rangle$, $\langle \bullet L_v = 2, L_s = 1 \rangle$ and $\langle \blacktriangle L_v = 2, L_s = 2 \rangle$ highlighted in blue in Table 3. Subsequently, we determine the optimal HSTA layers M for these combinations, with results presented in Table 4. An interesting observation is that our method performs best on the CASME3 dataset when the total value of L_v or L_s approximates 8, i.e., when $M \times L \approx 8$ performs best. Here L is the number of L_v or L_s . For instance, configurations such as $\langle M = 8, L = 1 \rangle$ and $\langle M = 4, L = 2 \rangle$ exhibit good performance, without requiring both L_v and L_s to be set at 8. Considering that L_v corresponds to the processing of larger video frames as mentioned above, we prefer configurations where $L_v > L_s$. Balancing computational load and performance, a configuration like $\langle L_v = 2, L_s = 1, M = 4 \rangle$, highlighted in blue in Table 4, achieves an effective balance between accuracy and computational effort. We plan to apply these optimally determined parameter values in our subsequent experiments.

4.3 Comparisons with Other Methods (RQ3)

We compare our approach with traditional hand-crafted methods, mainstream approaches in recent years and some latest comprehensive methods on the classic MER datasets detailed in three-category classification in Table 5. It is observed that our model outperforms any hand-crafted method like LBP-TOP [2014], Bi-WOOF [2018] and those that rely solely

Method	SMIC		SAMM		CASMEII	
	UF1	UAR	UF1	UAR	UF1	UAR
LBP-TOP [2014]	20.0	52.8	39.5	41.0	70.3	74.3
Bi-WOOF [2018]	57.3	58.3	52.1	51.4	78.1	80.3
CapsuleNet [2019]	58.2	58.8	62.1	59.9	70.7	70.2
MMNet [2022a]	44.1	43.8	32.6	34.2	71.9	89.9
Dual-Incep [2019]	57.1	57.1	49.3	49.6	75.4	77.4
Dual-ATME [2023]	64.6	65.8	56.2	53.8	76.5	75.1
STSTNet [2019]	68.0	70.1	65.9	68.1	83.8	86.9
MERSiamC3D [2021]	73.6	76.0	74.8	72.8	89.2	88.7
FRL-DGT [2023]	74.3	74.9	77.2	75.8	91.9	90.3
Micro-BERT [2023]	85.5	83.8	83.9	84.8	90.3	89.1
HSTA(Ours)	84.7	78.0	84.7	83.9	92.5	92.2

Table 5: Comparisons on SMIC, SAMM, and CASME II.

Method	Classes	UF1	UAR
FR [2022]	3	34.9	34.1
HTNet [2023b]	3	57.7	54.2
Micro-BERT [2023]	3	56.0	61.3
HSTA(Ours)	3	59.3	61.8
RGB [2022b]	7	17.6	18.0
RGB-D [2022b]	7	17.7	18.3
Micro-BERT [2023]	7	32.6	32.5
HSTA(Ours)	7	34.1	35.8

Table 6: Comparisons with others on CASME3.

on special frames as CapsuleNet [2019], MMNet [2022a] or optical flow Dual-Incep [2019], Dual-ATME [2023]. Compared to other models utilizing temporal information STSTNet [2019], MERSiamC3D [2021], our model also demonstrates higher performance across all datasets. Our approach also outperforms some of the latest comprehensive methods, such as FRL-DGT [2023]. When compared with Micro-BERT [Nguyen *et al.*, 2023], there are instances, such as in smaller-scale datasets like SMIC, our model’s performance is not as high. However, our computing costs are only about 1/96 of Micro-BERT’s. We speculate poor performance here is due to the small data size and the consequential larger impact of randomness. Therefore, we conduct further tests on the larger and more diverse CASME3 dataset, as shown in Table 6. Compared to the benchmark method RGB-D[2022b] used as our baseline for CASME3, which neither considers temporal information nor incorporates the fusion of cross-modal data, our model exhibits 16.4% performance improvement (UF1 17.7% \rightarrow 34.1%). So when the data size is sufficiently large, our model’s performance significantly surpasses other methods, without the need for the extensive and resource-intensive pre-training required by Micro-BERT.

4.4 Additional Exploration

We discover that benchmark datasets contain a wealth of additional data, such as macro-expressions (MaE) and objective classes (OC) which utilize objective facial muscle motion blocks - action units (AUs) - as proposed by [Davison *et al.*, 2018] to categorize MEs rather than relying on annotator

SF	OF	MaE	OC	Classes	ACC	UAR	UF1
✓	×	×	×	7	42.6	38.0	35.1
×	✓	×	×	7	40.1	36.4	35.4
×	×	×	✓	7	52.3	43.2	43.0
×	✓	✓	✓	7	48.6	39.3	41.8
✓	×	✓	✓	7	69.8	51.7	52.7

Table 7: Additional exploration on other data in CASME3 dataset.

or subject’s subjective judgment in CASME3. However, the use and exploration of these data have been limited in previous studies, which we aim to investigate their effectiveness. Recent MER methods tend to favor the utilization of optical flow (OF), but it is not essential. Due to the high flexibility of our model’s input, we replace the input of the special frames (SF) with optical flow to compare their effectiveness. Based on these considerations mentioned above, we conduct the following experiments: firstly, we use SF as the input as our baseline and then replace them with the OF; secondly, we employ a more rational label categorization method called objective classes to divide MEs categories, avoiding subjective judgment; thirdly, we explore the effect of incorporating MaE as additional data for training. The performance of these three approaches is presented in Table 7.

Using optical flow data extracted from special frames leads to a decrease in performance compared to directly using special frames in Table 7. When we utilize more objective classes based on AUs for a seven-category classification, there is a significant performance improvement (UF1 35.1% \rightarrow 43.0%). This further demonstrates the precision of our model in capturing subtle, objective facial movements. Given the limited amount of ME data, the model does not fully realize its potential, particularly when encountering unfamiliar facial types or expressions. However, by utilizing larger data from the CASME3 dataset for auxiliary training, we overcome the limitations (UF1 35.1% \rightarrow 52.7%). Our model excels with both subjective labels and objective classes. In summary, we can attain even higher performance by the often-neglected data in datasets which are also of great value, and our method is a universally applicable one.

5 Conclusion

In this paper, we propose a hierarchical attention strategy for different modalities to tackle the MER problem. Specifically, (1) Unimodal space-time attention (USTA) is used to capture temporal information in MER. (2) Crossmodal space-time attention (CSTA) is designed to fuse the different modalities while maintaining their uniqueness. (3) The hierarchical structure based on USTA and CSTA is proposed to grasp deeper facial cues. The extensive experiments have demonstrated the effectiveness of our proposed method. In addition, we verify the generalizability of our method on different types of additional data contained in a benchmark dataset.

Note that our method relies on a cascaded structure of USTA with CSTA, which may increase the computational complexity. In future work, we will focus on more efficient space-time attention methods to accelerate the recognition process of MEs.

References

- [Arnab *et al.*, 2021] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *ICCV*, pages 6836–6846, 2021.
- [Chen *et al.*, 2021] Chun-Fu Richard Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *ICCV*, pages 357–366, 2021.
- [Davison *et al.*, 2016] Adrian K Davison, Cliff Lansley, Nicholas Costen, Kevin Tan, and Moi Hoon Yap. Samm: A spontaneous micro-facial movement dataset. *TAC*, 9(1):116–129, 2016.
- [Davison *et al.*, 2018] Adrian K Davison, Walied Merghani, and Moi Hoon Yap. Objective classes for micro-facial expression recognition. *Journal of imaging*, 4(10):119, 2018.
- [Ekman, 2009] Paul Ekman. *Telling lies: Clues to deceit in the marketplace, politics, and marriage (revised edition)*. WW Norton & Company, 2009.
- [Gan *et al.*, 2019] Yee Siang Gan, Sze-Teng Liong, Wei-Chuen Yau, Yen-Chang Huang, and Lit-Ken Tan. Off-apexnet on micro-expression recognition system. *Signal Processing: Image Communication*, 74:129–139, 2019.
- [Khor *et al.*, 2018] Huai-Qian Khor, John See, Raphael Chung Wei Phan, and Weiyao Lin. Enriched long-term recurrent convolutional network for facial micro-expression recognition. In *IEEE FG 2018*, pages 667–674, 2018.
- [Khor *et al.*, 2019] Huai-Qian Khor, John See, Sze-Teng Liong, Raphael CW Phan, and Weiyao Lin. Dual-stream shallow networks for facial micro-expression recognition. In *ICIP*, pages 36–40, 2019.
- [Li *et al.*, 2013] Xiaobai Li, Tomas Pfister, Xiaohua Huang, Guoying Zhao, and Matti Pietikäinen. A spontaneous micro-expression database: Inducement, collection and baseline. In *IEEE FG*, pages 1–6, 2013.
- [Li *et al.*, 2019] Jing Li, Yandan Wang, John See, and Wenbin Liu. Micro-expression recognition based on 3d flow convolutional neural network. *Pattern Analysis and Applications*, 22:1331–1339, 2019.
- [Li *et al.*, 2022a] Hanting Li, Mingzhe Sui, Zhaoqing Zhu, and Feng Zhao. Mmnet: Muscle motion-guided network for micro-expression recognition. *IJCAI*, 2022.
- [Li *et al.*, 2022b] Jingting Li, Zizhao Dong, Shaoyuan Lu, Su-Jing Wang, Wen-Jing Yan, Yinhan Ma, Ye Liu, Changbing Huang, and Xiaolan Fu. Cas (me) 3: A third generation facial spontaneous micro-expression database with depth information and high ecological validity. *TPAMI*, 45(3):2782–2800, 2022.
- [Liong and Wong, 2017] Sze-Teng Liong and KokSheik Wong. Micro-expression recognition using apex frame with phase information. In *APSIPA ASC*, pages 534–537, 2017.
- [Liong *et al.*, 2018] Sze-Teng Liong, John See, KokSheik Wong, and Raphael C-W Phan. Less is more: Micro-expression recognition from video using apex frame. *Signal Processing: Image Communication*, 62:82–92, 2018.
- [Liong *et al.*, 2019] Sze-Teng Liong, Yee Siang Gan, John See, Huai-Qian Khor, and Yen-Chang Huang. Shallow triple stream three-dimensional cnn (ststnet) for micro-expression recognition. In *IEEE FG*, pages 1–5, 2019.
- [Nguyen *et al.*, 2023] Xuan-Bac Nguyen, Chi Nhan Duong, Xin Li, Susan Gauch, Han-Seok Seo, and Khoa Luu. Micron-bert: Bert-based facial micro-expression recognition. In *CVPR*, pages 1482–1492, 2023.
- [Qian *et al.*, 2021] Rui Qian, Tianjian Meng, Boqing Gong, Ming-Hsuan Yang, Huisheng Wang, Serge Belongie, and Yin Cui. Spatiotemporal contrastive video representation learning. In *CVPR*, pages 6964–6974, 2021.
- [Reddy *et al.*, 2019] Sai Prasanna Teja Reddy, Surya Teja Karri, Shiv Ram Dubey, and Snehasis Mukherjee. Spontaneous facial micro-expression recognition using 3d spatiotemporal convolutional neural networks. In *IJCNN*, pages 1–8, 2019.
- [Rehman *et al.*, 2022] Yasar Abbas Ur Rehman, Yan Gao, Jijun Shen, Pedro Porto Buarque de Gusmao, and Nicholas Lane. Federated self-supervised learning for video understanding. In *ECCV*, pages 506–522, 2022.
- [See *et al.*, 2019] John See, Moi Hoon Yap, Jingting Li, Xiaopeng Hong, and Su-Jing Wang. Megc 2019—the second facial micro-expressions grand challenge. In *2019 14th IEEE FG*, pages 1–5, 2019.
- [Tong *et al.*, 2022] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *NIPS*, 35:10078–10093, 2022.
- [Van Quang *et al.*, 2019] Nguyen Van Quang, Jinhee Chun, and Takeshi Tokuyama. Capsulenet for micro-expression recognition. In *IEEE FG*, pages 1–7, 2019.
- [Vaswani *et al.*, 2017] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NIPS*, 30, 2017.
- [Verma *et al.*, 2019] Monu Verma, Santosh Kumar Vipparthi, Girdhari Singh, and Subrahmanyam Murala. Learnnet: Dynamic imaging network for micro expression recognition. *IEEE Transactions on Image Processing*, 29:1618–1627, 2019.
- [Wang *et al.*, 2023a] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinhan He, Yi Wang, Yali Wang, and Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking. In *CVPR*, pages 14549–14560, 2023.
- [Wang *et al.*, 2023b] Zhifeng Wang, Kaihao Zhang, Wenhan Luo, and Ramesh Sankaranarayanan. Htnet for micro-expression recognition. *arXiv preprint arXiv:2307.14637*, 2023.

- [Wei *et al.*, 2020] Xi Wei, Tianzhu Zhang, Yan Li, Yongdong Zhang, and Feng Wu. Multi-modality cross attention network for image and sentence matching. In *CVPR*, pages 10941–10950, 2020.
- [Wodajo and Atnafu, 2021] Deressa Wodajo and Solomon Atnafu. Deepfake video detection using convolutional vision transformer. *arXiv preprint arXiv:2102.11126*, 2021.
- [Yan *et al.*, 2014] Wen-Jing Yan, Xiaobai Li, Su-Jing Wang, Guoying Zhao, Yong-Jin Liu, Yu-Hsin Chen, and Xiaolan Fu. Casme ii: An improved spontaneous micro-expression database and the baseline evaluation. *PLoS one*, 9(1):e86041, 2014.
- [Zhai *et al.*, 2023] Zhijun Zhai, Jianhui Zhao, Chengjiang Long, Wenju Xu, Shuangjiang He, and Huijuan Zhao. Feature representation learning with adaptive displacement generation and transformer fusion for micro-expression recognition. In *CVPR*, pages 22086–22095, 2023.
- [Zhao *et al.*, 2021] Sirui Zhao, Hanqing Tao, Yangsong Zhang, Tong Xu, Kun Zhang, Zhongkai Hao, and Enhong Chen. A two-stage 3d cnn based learning method for spontaneous micro-expression recognition. *Neurocomputing*, 448:276–289, 2021.
- [Zhao *et al.*, 2023] Sirui Zhao, Huaying Tang, Xinglong Mao, Shifeng Liu, Yiming Zhang, Hao Wang, Tong Xu, and Enhong Chen. Dfme: A new benchmark for dynamic facial micro-expression recognition. *TAC*, 2023.
- [Zhou *et al.*, 2019] Ling Zhou, Qirong Mao, and Luoyang Xue. Database micro-expression recognition. In *IEEE FG*, pages 1–5, 2019.
- [Zhou *et al.*, 2022] Ling Zhou, Qirong Mao, Xiaohua Huang, Feifei Zhang, and Zhihong Zhang. Feature refinement: An expression-specific feature learning and fusion method for micro-expression recognition. *Pattern Recognition*, 122:108275, 2022.
- [Zhou *et al.*, 2023] Haoliang Zhou, Shucheng Huang, Jingtong Li, and Su-Jing Wang. Dual-atme: Dual-branch attention network for micro-expression recognition. *Entropy*, 25(3):460, 2023.