


---

This is the **accepted version** of the journal article:

Wang, Ziqi; Fu, Wenwen; Zhang, Yue; [et al.]. «MCNet : meta-clustering learning network for micro-expression recognition». Journal of electronic imaging, Vol. 33, issue 2 (March 2024), art. 23014. DOI 10.1117/1.JEI.33.2.023014

---

This version is available at <https://ddd.uab.cat/record/311811>

under the terms of the  <sup>IN</sup> COPYRIGHT license

# MCNet: Meta-Clustering Learning Network for Micro-Expression Recognition

Ziqi Wang<sup>a</sup>, Wenwen Fu<sup>a</sup>, Yue Zhang<sup>a</sup>, Jiarui Li<sup>a</sup>, Wenjuan Gong<sup>a,\*</sup>, Jordi González<sup>b</sup>

<sup>a</sup>China University of Petroleum (East China), Qingdao Institute of Software, College of Computer Science and Technology, Computer Technology, Changjiang West Road, Qingdao, China, 266580

<sup>b</sup>Universitat Autònoma de Barcelona, Computer Vision Center, Edifici O, Bellaterra, Spain, 08193

**Abstract.** Facial micro-expressions are categorized into various types based on different criteria, and typically each major category is further divided into multiple subcategories of expressions. For traditional micro-expression recognition problems, multiple subcategories of the same emotions are indiscriminately learned and verified, leading to potential mis-classification, especially with negative emotions. To address the issue of intra-class variation in micro-expressions, we propose a novel meta-clustering learning network for micro-expression recognition called MCNet. This approach integrates the ideas of meta-learning and clustering, hierarchically clustering subcategories within a micro-expression class to generate multiple class centers for metric-based classification. The proposed method diverges from the common strategy of metric-based meta-learning algorithms, which typically use the mean feature of all samples within the same class as the class center. Furthermore, we incorporate transfer learning into the meta-learning process to jointly alleviate overfitting caused by the scarcity of micro-expression data. We conduct extensive comparative experiments based on the Leave-One-Subject-Out protocol on three widely used micro-expression datasets. The experimental results demonstrate the competitive performance and strong generalization ability of the proposed MCNet approach.

**Keywords:** meta-learning network, hierarchical clustering, few-shot classification, micro-expression recognition.

\*Wenjuan Gong, [wenjuangong@upc.edu.cn](mailto:wenjuangong@upc.edu.cn)

## 1 Introduction

Facial expressions play a crucial role in conveying people’s emotional states and psychological fluctuations, making them integral to interpersonal communication.<sup>1</sup> Therefore, facial expression recognition has consistently remained a prominent area of research in computer vision. Micro-expressions, characterized by their short duration and subtle facial muscle movements,<sup>2</sup> are often imperceptible to the naked eyes. However, they offer a fresh perspective for investigating human emotional states, as they can accurately reflect emotions, even when individuals deliberately conceal their true feelings. Moreover, micro-expression recognition holds immense value across various domains, including lie detection, psychological assessment, negotiation, and interrogation.

31 Despite the potential applications of micro-expression recognition, detecting and interpreting  
32 micro-expressions accurately with the naked eyes remains a formidable challenge. Even experts  
33 find it difficult to manually identify these subtle facial movements. To alleviate this challenge,  
34 psychologists have developed auxiliary tools such as the Micro Expression Training Tool (METT)<sup>3</sup>  
35 and the Facial Action Coding System (FACS)<sup>4</sup> to facilitate the manual identification of micro-  
36 expressions. However, manual micro-expression recognition still faces several obstacles. Firstly,  
37 the ability to discern micro-expressions is limited to extensively trained professionals, which incurs  
38 high costs. Secondly, even with intensive training, humans can only recognize approximately  
39 47% of micro-expressions.<sup>5</sup> Moreover, human analysis of micro-expressions is time-consuming,  
40 expensive, and prone to errors. Consequently, the development of automated systems for micro-  
41 expression analysis has become imperative.

42 In recent years, the availability of several spontaneous micro-expression datasets<sup>6-9</sup> has sparked  
43 increased interest among researchers in the computer vision field to tackle the challenge of micro-  
44 expression recognition, resulting in the proposal of numerous recognition algorithms. Automatic  
45 micro-expression recognition primarily focuses on extracting salient and distinctive features from  
46 micro-expression clips. In one of the early works on micro-expression recognition, Polikovsky et  
47 al.<sup>10</sup> introduced the 3D Histogram of Oriented Gradients (3DHOG) descriptor to extract appearance-  
48 based micro-expression features. To enhance the compactness and robustness of the features un-  
49 der varying lighting conditions, Huang et al.<sup>11</sup> proposed the SpatioTemporal Completed Local  
50 Quantization Patterns (STCLQP) approach. This method extracts information encompassing sign,  
51 magnitude, and orientation differences, and efficiently applies vector quantization and codebook  
52 selection for each component in both spatial and temporal domains. Finally, the features from  
53 the three components are combined into one feature based on the feature subspace. Additionally,

54 dynamic micro-expression features based on optical flow have demonstrated notable performance  
55 in micro-expression recognition tasks.<sup>12-15</sup> Optical flow captures the motion patterns of objects or  
56 scenes in an image sequence by detecting pixel intensity changes between frames. For example,  
57 Liu et al.<sup>14</sup> proposed the Main Direction Mean Optical flow (MDMO) feature, which integrates the  
58 magnitude and direction of the primary optical flow vectors across 36 non-overlapping Regions of  
59 Interest (ROIs) in the face, effectively reducing the feature dimensionality.

60 Different from the aforementioned handcrafted feature extraction methods, deep learning ap-  
61 proaches adopt the end-to-end manner by integrating feature extraction and classification. Patel  
62 et al.<sup>16</sup> were the first to introduce Convolutional Neural Networks (CNNs) to micro-expression  
63 recognition and mitigated overfitting due to the lack of sufficient micro-expression samples by  
64 using transfer learning and feature selection algorithms. Similarly, Kim et al.<sup>17</sup> employed a frame-  
65 work based on CNN and Long Short-Term Memory (LSTM) networks to learn spatial and temporal  
66 features of micro-expressions and incorporated expression states into the objective function to en-  
67 hance category separability. However, both CNN and LSTM models require a substantial number  
68 of training samples, which poses a challenge due to the limited availability of micro-expression  
69 datasets and increases the risk of overfitting. To alleviate the issue of overfitting, many works<sup>18-22</sup>  
70 have adopted transfer learning and data augmentation. For example, Zhi et al.<sup>21</sup> employed transfer  
71 learning by pre-training a 3DCNN on the macro-expression dataset Oulu-CASIA<sup>23</sup> and transferring  
72 the learned knowledge to the micro-expression domain. Sun et al.<sup>22</sup> proposed a novel knowledge  
73 transfer method that extracts and transfers knowledge from facial action units for micro-expression  
74 recognition. Takalkar et al.<sup>18</sup> and Zhang et al.<sup>24</sup> expanded the dataset by combining and horizon-  
75 tally flipping the data, respectively.

76 Despite the advancements made in deep learning-based approaches for micro-expression recog-

77 nition, these works have not achieved remarkable results and even performed worse than some tra-  
78 ditional methods. To address this issue, Peng et al.<sup>25</sup> introduced a Dual Temporal Scale Convolu-  
79 tional Neural Network (DTSCNN) that can handle micro-expression segments with different frame  
80 rates. They incorporated optical flow sequences extracted from the original micro-expression clips  
81 as inputs to the network, leading to a significant improvement in the overall recognition accuracy.  
82 Inspired by this success, subsequent works<sup>26-30</sup> have also adopted the strategy of extracting optical  
83 flow information as input to their CNN models. In addition to handling discriminative represen-  
84 tations of micro-expression sequences from a dynamic feature perspective, deep learning methods  
85 based on attention mechanisms have also demonstrated success<sup>31-34</sup>. For instance, Zhou et al.<sup>32</sup>  
86 proposed the Dual-branch Attention Network (Dual-ATME), comprising Hand-crafted Attention  
87 Region Selection (HARS) and Automated Attention Region Selection (AARS). The HARS man-  
88 ually extracted features from the Region of Interest (ROI) using prior knowledge, while AARS  
89 automatically extracted hidden information from the sequence based on attention mechanisms.  
90 This method effectively learned micro-expression representations by the dual-scale features, and  
91 demonstrated that the influence of facial visual features and the visual relationships between fa-  
92 cial ROIs in distinguishing micro-expressions. Therefore, Thuseethan et al.<sup>35</sup> introduced a novel  
93 end-to-end facial micro-expression detection framework named Deep3DCANN. This method em-  
94 ployed a deep 3D convolutional neural network to learn spatiotemporal features from facial ex-  
95 pression sequences and a deep artificial neural network to capture visual associations between  
96 different facial sub-regions. The combination of output features from both pipelines allowed for  
97 the identification of micro-expression changes on frame-level and showed excellent performance  
98 in multi-class classification tasks.

99 Furthermore, integrating deep features with handcrafted features in a dynamic facial micro-

100 expression recognition architecture has proven effective in enhancing performance. For example,  
101 Wang et al.<sup>36</sup> proposed the Local Cube Binary Pattern Spatial-Temporal Graph Convolutional Net-  
102 work (LCBP-STGCN), where STGCN consists of a Spatial Graph Convolutional Network (SGCN)  
103 for acquiring spatial information and a Temporal Convolutional Network (TCN) for capturing tem-  
104 poral information related to micro-expressions. This method introduced low-level spatio-temporal  
105 features obtained from Local Cube Binary Pattern, and then the spatio-temporal graph convolu-  
106 tional network was used to extract high-level features of micro-expression sequences. The combi-  
107 nation of these structures preserved potential discriminative information in the sequence, improv-  
108 ing the robustness of feature extraction. Although notable progress has been made by these methods  
109 in the field of micro-expression recognition, the limited availability of data samples and the risk  
110 of overfitting still pose challenges during the training process. Therefore, further improvements in  
111 the performance of micro-expression recognition can still be pursued.

112 To solve this problem, the meta-learning approaches have been proposed. Chelsea Finn et al.  
113<sup>37</sup> introduced the idea of Model-Agnostic Meta-Learning (MAML). The goal of MAML is to en-  
114 able the model to quickly adapt to new tasks through a small number of training steps and data.  
115 By performing gradient updates across multiple tasks, MAML finds initialized parameters, from  
116 which the model can achieve good generalization using a small amount of sample data within a few  
117 update steps. MAML, with its model-independent advantages, achieves the effect that can be used  
118 with any model structure, and has become an important tool in the field of meta-learning. Ravi et  
119 al.<sup>38</sup> proposed a novel meta-learning method by training an optimizer to enable it to quickly adapt  
120 to new tasks and achieve good performance in a few-shot learning task. This method is known  
121 as learning optimization (Learning to Optimize) or LSTM-meta-optimizer. The core idea is to let  
122 the model itself learn how to perform optimization, rather than relying on manually designed opti-

123 mization algorithms such as stochastic gradient descent (SGD). By using recurrent neural network  
124 to update the parameters of the model, the optimizer receives the information of the loss func-  
125 tion during learning and outputs the parameter values to be changed. In the meta-learning phase  
126 on a specific task, the LSTM optimizer will decide how to do the next parameter update based  
127 on the past parameter update history and the gradient of the target function. In the new task, the  
128 LSTM optimizer can generate parameter updates when the model needs rapid adaptation, which  
129 significantly improves the performance of few-shot learning.

130 To mitigate the problem of limited data samples in micro-expression recognition, a meta-  
131 learning-based multi-model fusion network (Meta-MMFNet) was proposed in,<sup>24</sup> which was the  
132 first application of meta-learning in the field of micro-expression recognition. This method em-  
133 ployed a Prototypical Network (PN) which represented each class by its average feature, and  
134 fused the transfer learning stream through a weighted sum approach. The Meta-MMFNet method  
135 achieves comparable performance to state-of-the-art methods even with a small number of training  
136 samples, thus alleviating the challenge of overfitting to some extent. However, it is worth not-  
137 ing that this metric-based meta-learning micro-expression recognition algorithm treats the mean  
138 feature as the class centroid and adopts the nearest centroid method for classifying unseen sam-  
139 ples. From the perspective of micro-expression categories, certain classes may have more than one  
140 centroid under specific situations. For instance, the category of “negative” in the SMIC dataset  
141 comprises multiple sub-emotions such as “disgust”, “repression” and “sadness”, each with its own  
142 centroid. Moreover, individual differences among participants introduce additional factors unre-  
143 lated to expressions, such as race and gender. Furthermore, due to variances in facial muscle  
144 movement habits, there exist substantial differences between different samples of the same expres-  
145 sion category. Collectively, these factors contribute to a considerable intra-class distance within

146 micro-expressions.

147 To tackle this issue, we present a novel meta-clustering learning network named MCNet, draw-  
148 ing inspiration from hierarchical clustering method.<sup>39</sup> Hierarchical clustering is an algorithm that  
149 constructs a nested clustering tree with hierarchies by assessing the similarity between data points  
150 of various categories, including bottom-up merging and top-down splitting strategies. In this study,  
151 we focus on learning the bottom-up merging concept, where the two most similar data points are  
152 successively merged. We set the stopping criterion for clustering when the data points to be merged  
153 belong to different classes.

154 The proposed MCNet replaces the previous practice of utilizing a single mean class prototype  
155 (a centroid) by incorporating a hierarchical clustering module. This module employs multiple clus-  
156 tering prototypes to represent each category or sub-category. Specifically, when a category en-  
157 compasses multiple distinct sub-categories, each with unique features, the hierarchical clustering  
158 module clusters micro-expression samples based on the feature differences among sub-categories,  
159 resulting in the formation of multiple centroids. This ensures a comprehensive modeling of intra-  
160 category differences, addressing the limitation of relying solely on average features as centroids  
161 that may inaccurately represent the distinctive characteristics of each category. Consequently, this  
162 enhances MCNet’s ability to capture fine-grained features. The hierarchical clustering module fa-  
163 cilitates the identification and separation of similar and dissimilar sub-groups within datasets, and  
164 the integration of multiple prototypes enhances MCNet’s ability to capture the complexity and  
165 nuances of each category. This alleviates the discrepancies in characteristics among samples in  
166 sub-emotion clusters, thereby leading to more precise sentiment analysis. Additionally, MCNet in-  
167 troduces prior knowledge from both micro-expressions and macro-expressions, further enhancing  
168 recognition accuracy on unknown samples.

169 The main contributions of this study are as follows:

- 170 • This study identifies and addresses the mis-classification problem caused by the large intra-  
171 class differences in micro-expression samples. To tackle this, we propose a novel meta-  
172 clustering learning network for micro-expression recognition, introducing a hierarchical  
173 clustering module that utilizes multiple clustering prototypes instead of a single mean pro-  
174 totype, breaking the simple strategy of using centroids as class prototypes in previous meta-  
175 learning methods.
- 176 • This study combines transfer learning, meta-learning and clustering principles, which uti-  
177 lizes the sufficient macro-expression data along with meta-learning to alleviate the the accu-  
178 racy limitations caused by the scarcity of micro-expression data.
- 179 • The proposed method is evaluated on three widely-used micro-expression datasets, and the  
180 experimental results demonstrate highly competitive performance of the MCNet method.

## 181 2 The Meta-clustering Learning Network

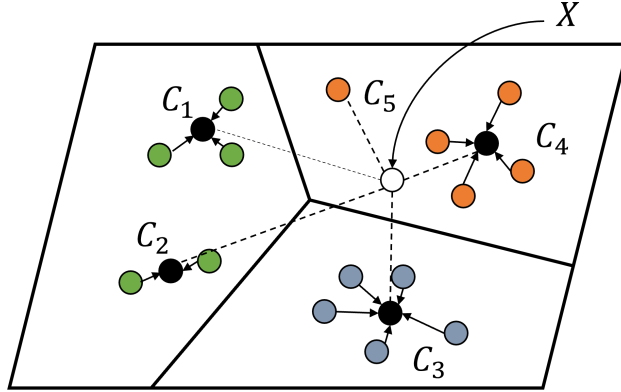
182 Following the meta-learning formulation for micro-expression recognition,<sup>24</sup> we also adopt the  $N$ -  
183 way  $K$ -shot classification task of few-shot learning, where the dataset is typically divided into a  
184 meta-training set and meta-testing set. In meta-learning, the dataset is typically divided into a train-  
185 ing (also called meta-training) set and a testing (also called meta-testing) set. In the training (also  
186 called meta-training) stage,  $N$  unique categories are randomly sampled from the meta-training set,  
187 from which  $K$  and  $M$  unique samples are randomly selected from each category to compose the  
188 support set  $\Phi_S$  and query set  $\Phi_Q$ , respectively. These support sets and query sets constitute dis-  
189 tinct meta-tasks, allowing the model to be trained in each iteration. In the meta-testing stage, the

190 support set and query set are sampled from the testing set using a similar mechanism as the one  
191 used in the meta-training stage. By training on different meta-tasks, the model learns generalizable  
192 knowledge, enabling it to quickly adapt to new tasks.

193 Due to the limited availability of data in micro-expression classification, we make adjustments  
194 to the traditional meta-tasks using non-overlapping categories, according to the distribution of  
195 participants. Specifically, following the Leave-One-Subject-Out (LOSO) evaluation protocol, we  
196 set the query set  $\Phi_Q^{test}$  of the meta-testing set  $\Omega_{test}$  to include all samples from the participant set  
197  $Z_{test}$  for testing. The support set  $\Phi_S$  of the meta-training set  $\Omega_{train}$  and the meta-testing set  $\Omega_{test}$   
198 are composed of samples from the remaining participant set  $Z_{train}$ , where  $Z_{train} \cap Z_{test} = \emptyset$ . This  
199 strategy ensures a comprehensive evaluation by sampling from different participants for training  
200 and testing, leading to robust and reliable results.

201 We use a metric-based prototypical network as the meta-learning model. The standard ap-  
202 proach has some drawbacks when directly applied for micro-expression classification. The prob-  
203 lem lies in that the query samples may be far away from the class prototype when the gap between  
204 the samples within the class is too large, thus affecting the classification accuracy of the model.  
205 It should be noted that a micro-expression class can be further divided into multiple sub-groups,  
206 and there are subtle variations among samples of the same class due to slight muscle movement  
207 differences. Therefore, it is difficult to achieve a favorable classification using only the center of  
208 samples. To address this issue, we incorporate hierarchical clustering, as illustrated in Fig. 1, in  
209 the proposed approach.

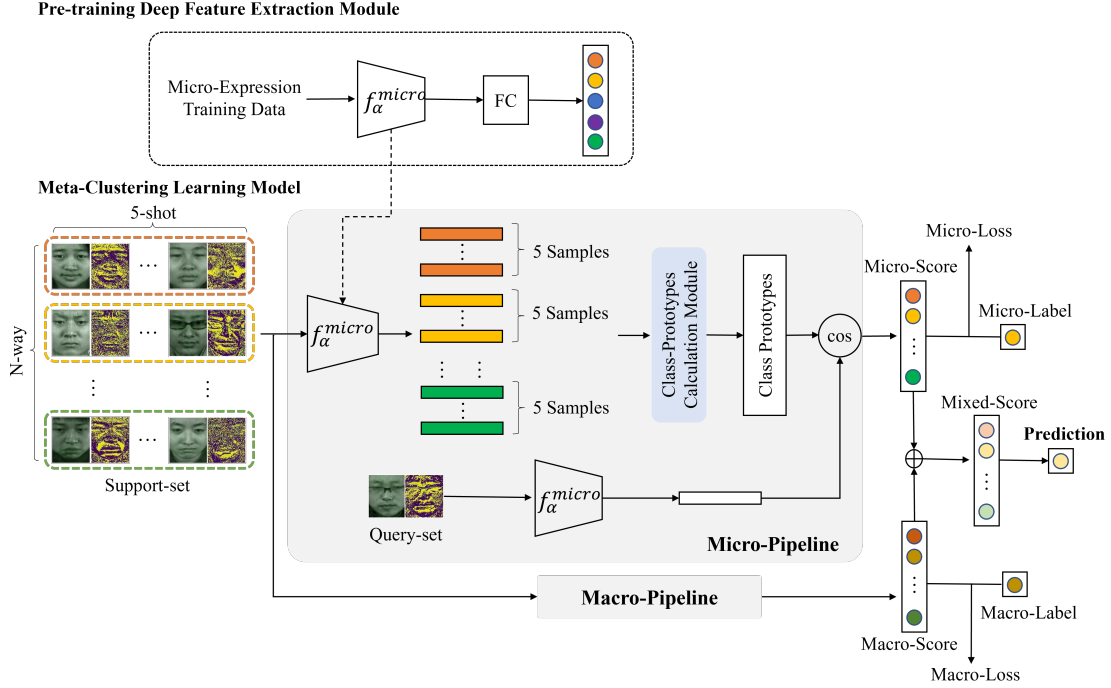
210 The proposed MCNet method aims to address the challenge of large intra-class differences  
211 in micro-expression recognition by integrating the idea of hierarchical clustering with a fusion  
212 network based on meta-learning. Meta-learning network is trained and tested on a per-task basis,



**Fig 1** The process of predicting sample categories using the proposed meta-clustering method. Instead of using a single mean feature prototype for each class, the method determines the class prototypes by clustering the two nearest samples within the same class in the support set, resulting in  $N$  or more class prototypes, where  $N$  represents the number of classes in the support set. To predict the label of a query sample  $X$  during meta-testing, the meta-clustering method compares the distances between the sample and each class prototype and identifies the closest class prototype to the sample. The label of the test sample is obtained by mapping the nearest class prototype back to its original category.

213 where each task comprises a support set and a query set. For each micro-expression sample in  
 214 meta-tasks, the fused motion feature obtained from the feature fusion module is fed into a pre-  
 215 trained deep feature extraction module to acquire corresponding deep feature vector. In the metric-  
 216 based meta-learning module, these feature vectors are utilized to generate multiple class prototypes  
 217 using hierarchical clustering, allowing for the prediction of class membership for each sample in  
 218 the query set.

219 The overall structure of the proposed method is shown in Fig. 2. The proposed method mainly  
 220 consists of three components: the feature processing module, the meta-learning module and the  
 221 meta-clustering module. To effectively capture motion features and reduce redundancy in micro-  
 222 expression clips, we fuse the optical flow features and frame difference features, which are then  
 223 processed by a deep feature extraction network to obtain comprehensive and discriminative deep  
 224 features. To tackle the problem of limited micro-expression data, we introduce transfer learning  
 225 during the pre-training phase and obtain two independent feature extraction networks. Subse-



**Fig 2** The overall architecture of the meta-clustering learning network. This method comprises of two main stages: pre-training the deep feature extraction module and learning the meta-clustering classification model. In the pre-training stage, two independent deep feature extraction networks based on ResNet18 are trained on the micro-expression dataset  $\Omega_{train}$  and the macro-expression dataset CK+. In the meta-clustering stage, the fused dynamic features are fed into the micro-pipeline and macro-pipeline and the deep features of micro-expressions are then extracted in each pipeline. The cosine distance between the query samples and the class prototypes obtained from the class prototype calculation module is computed, yielding the micro-score and macro-score, respectively. Based on these two score vectors, we optimize the micro-pipeline and macro-pipeline separately using cross-entropy loss. during the validation process, we combine the score vectors generated from both pipelines to merge the micro- and macro-predictions to predict the final micro-expression category.

226 quently, the metric-based meta-learning module learns micro-expression classification on a task-  
 227 based level under two pathways: the micro-pipeline and the macro-pipeline. In the testing stage,  
 228 the meta-clustering module generates multiple class prototypes for each task to address the chal-  
 229 lenge of mis-classification caused by intra-class variations. Furthermore, to fully utilize priors  
 230 from both micro-expressions and macro-expressions, we fuse the two predictions in the last step,  
 231 further enhancing the recognition accuracy.

### 232 2.0.1 Feature processing module

233 Due to subtle variations and brief durations of micro-expression motions, micro-expression se-  
234 quences often contain redundant information. Previous studies<sup>15,40</sup> have demonstrated that extract-  
235 ing features from the onset and apex frames of a video are more effective. Therefore, frame dif-  
236 ference features and optical flow features are extracted from the two frames and then concatenated  
237 to obtain the fused feature  $x_i$ . Based on this, we represent the support set and the query set of the  
238 meta-task as  $\Phi_S = \{(x_s^1, y_s^1), \dots, (x_s^i, y_s^i), \dots, (x_s^{NK}, y_s^{NK})\}$  and  $\Phi_Q = \{(x_q^1, y_q^1), \dots, (x_q^i, y_q^i), \dots, (x_q^{NM}, y_q^{NM})\}$ ,  
239 where  $y_i$  represents the ground-truth label of the  $i$ -th sample and  $M$  denotes the number of samples  
240 in each category of the query set. And  $x_s^i$  represents the dynamic feature of the  $i$ -th micro-expres-  
241 sion sequence in the support set, while  $x_q^i$  is the dynamic micro-expression feature of the  $i$ -th  
242 sample in the query set.

243 The fused feature serves as input to the pre-training deep feature extraction module for further  
244 processing. Collecting micro-expressions is a time-consuming and labor-intensive process, and  
245 manual annotations from experts are required, resulting in a severe lack of available data. The  
246 success of previous works<sup>41,42</sup> have shown the feasibility of leveraging the abundance of macro-  
247 expressions to recognize micro-expressions. To address the data scarcity issue, we introduce the  
248 idea of transfer learning into the deep feature extraction module. In this module, we use ResNet-18  
249 as the backbone network and pre-train it on both the micro-expression training set and the macro-  
250 expression dataset, respectively. This process yields two trained deep feature extraction networks  
251 denoted as  $f_\alpha^{micro}$  and  $f_\beta^{macro}$ , where  $\alpha$  and  $\beta$  represent the network parameters pre-trained on the  
252 corresponding datasets. These two networks are trained independently and have no interference  
253 with each other, only being weighted and fused in the final recognition stage to achieve better

254 predictions.

### 255 2.0.2 Meta-learning module

256 During the training of the meta-learning model, each meta-task consists of a support set with  $N$   
257 classes, where each class is sampled with  $K$  samples, and a query set with randomly selected  $M$   
258 samples from each class. For each task, the meta-learning module processes the fused feature  
259 vectors extracted for each sample using the previous module. Specifically, the fused features are  
260 processed by networks  $f^*$ ,  $*$  = {micro, macro} to adapt to micro-expression meta-tasks based  
261 on the priors from micro-expressions and macro-expressions, respectively. The updated feature  
262 matrices of a meta-task for each pipeline are represented as  $\Psi_S = [\psi_s^1, \dots, \psi_s^i, \dots, \psi_s^{NK}]$  and  $\Psi_Q =$   
263  $[\psi_q^1, \dots, \psi_q^i, \dots, \psi_q^{NM}]$ .

264 To optimize the prototype networks using stochastic gradient descent mechanism, we calculate  
265 the average feature vector for each class  $c$  as the representation of the corresponding class prototype  
266  $r_c$  during the meta-training phase of each data stream. Then, we compute the cosine distance  
267 between each query sample feature  $\psi_q^i$  and each class prototype to obtain the predicted probabilities  
268 of the  $i$ -th sample belonging to each class. The prediction is denoted as  $V_i = [v_i^1, \dots, v_i^c, \dots, v_i^C]$ ,  
269 where  $C$  represents the total number of classes in the support set. The predicted probability is  
270 formulated as follows:

$$v_i^c = \text{softmax}(\cos(\psi_q^i, r_c)), \quad (1)$$

271 where  $\cos(\cdot)$  calculates the cosine distance between vectors, and  $\text{softmax}(\cdot)$  represents the ac-  
272 tivation function. Next, we update the meta-learning processes for the micro-pipeline and the

273 macro-pipeline using cross-entropy loss, respectively:

$$L = - \sum_{i=1}^{N \times M} \sum_{c=1}^C y_q^{i,c} \log v_i^c, \quad (2)$$

274 where  $y_q^{i,c}$  represents the ground-truth label indicating whether the  $i$ -th query sample belongs to  
 275 class  $c$ , and  $v_i^c$  represents the predicted probability of the  $i$ -th query sample belonging to class  $c$ .

### 276 2.0.3 Meta-clustering module

277 In the meta-testing stage of a standard prototypical network based solution, the feature vectors of  
 278 the samples in the testing task are extracted from the deep feature extraction network. Then, the  
 279 mean feature of all samples of each category in the support set is computed as the corresponding  
 280 class prototype.<sup>24</sup>

$$r_c = \frac{1}{|S_c|} \sum_{x \in S_c} f^*(x), * = \{micro, macro\}, \quad (3)$$

281 where  $S_c$  is the sample cluster of class  $c$  in the support set  $\Psi_S$ . After obtaining the class prototypes,  
 282 the distances between the feature vectors of the query samples and class prototypes are calculated  
 283 for classification.

284 To address the issue of large intra-class distances and mitigate prediction errors, we propose an  
 285 improved clustering method in the testing phase. Specifically, we adopt a hierarchical clustering  
 286 approach, which aggregates feature vectors from the bottom up based on their similarity. Fig. 1  
 287 illustrates the process of the proposed hierarchical clustering module. The module generates a  
 288 number of prototypes equal to or greater than the number of classes in the support set, which are  
 289 then used for distance metric calculations. Taking the micro-pipeline for example, each sample in  
 290 the support set is initially considered as a cluster, i.e., treated as a prototype. For all the feature

291 vectors in the  $\Psi_S^{test}$ , the distances between each pair of vectors are computed, and the two most  
 292 similar micro-expression data are sequentially identified and merged:

$$x_{new} = \frac{f_{\alpha}^{micro}(x_i) \times w_i + f_{\alpha}^{micro}(x_j) \times w_j}{w_i + w_j} \quad (4)$$

$$w_{new} = w_i + w_j$$

293 where  $f_{\alpha}^{micro}(\cdot)$  is the deep feature extraction network pre-trained for micro-pipeline and  $\alpha$  rep-  
 294 resents the parameter set of this neural network. And, the two feature vectors  $x_i$  and  $x_j$  from the  
 295 sample cluster of category  $c$  have the same ground-truth class label and the closest distance among  
 296 all pairs of vectors in the support set. Each vector is assigned a weight,  $w_i$  and  $w_j$ , respectively.  
 297 And the aggregated new vector,  $x_{new}$ , is obtained by calculating the weighted average of the two  
 298 vectors. The weight of the new vector is denoted as  $w_{new}$ . In this study, initial weights of all  
 299 samples are 1.

300 The above process is iterated, where the closest two clusters are merged until two clusters from  
 301 different classes are encountered, indicating the end of clustering. Each cluster is treated as a  
 302 separate class, and the distances between the generated class prototypes and the query samples are  
 303 metricized. There is a mapping relationship between the newly generated classes and the original  
 304 classes in the support set. The label of the test sample is obtained by mapping the nearest class  
 305 prototype back to its original category.

306 Additionally, knowledge from macro-expressions domain can provide valuable insights for  
 307 micro-expression recognition. To leverage this prior, we utilize a feature extraction network,  
 308  $f_{\beta}^{macro}$ , which has been pre-trained on macro-expressions, and transfer it to the micro-expression  
 309 domain. In the metric-based meta-learning module, we calculate two vectors of cosine similar-

ities through both the macro-pipeline and micro-pipeline, referred to as macro-score and micro-score, respectively. During the testing process, we further integrate the prior knowledge from both micro-expressions and macro-expressions by combining the outputs of these two data streams in a weighted sum, which defines the final metric distances between samples and the class prototypes. The process is formulated as follows:

$$\begin{aligned}
 H_{total}^i &= H_{micro}^i + \gamma H_{macro}^i, \\
 P_i &= \text{softmax}(H_{total}^i),
 \end{aligned}
 \tag{5}$$

where  $\gamma$  is a hyperparameter that balances the weight of micro-expression and macro-expression prior knowledge,  $H_{micro}^i$  and  $H_{macro}^i$  represent the similarity scores obtained from the micro-pipeline and macro-pipeline, respectively, and  $P_i$  represents the probabilities of the  $i$ -th query sample matching each prototype. The predicted category is the prototype with the highest probability. This fusion of micro- and macro-priors enhances the model’s classification performance and helps to improve the accuracy of micro-expression recognition.

### 3 Experimental results

In this study, we conducted rich experiments on common micro-expressions datasets (including SMIC, CASME, and CASME II datasets) to validate the effectiveness of the metaclustering learning model for micro-expressions classification. The SMIC dataset consists of 164 micro-expression videos from 16 subjects. The database categorizes micro-expressions into three categories, “positive”, “negative”, and “surprise”. The SMIC database includes micro-expression image sequences captured at 25 frames per second using a camera, as well as infrared image sequences. The CASME dataset comprises of 195 video sequences of micro-expressions

329 captured in natural situations from 26 participants. It covers seven main types of expressions, i.e.,  
330 "happiness", "disgust", and others. The classification task in this experiment mainly  
331 focused on four expression types, i.e., "disgust", "surprise", "repression", and  
332 "tense". The CASME II dataset contains 256 micro-expression events collected from 35  
333 Asian participants who were asked to watch specific videos to elicit micro-expression. CASME  
334 II has the advantages of high resolution, high sampling rate, large sample size, and detailed la-  
335 beling. The classification task in this experiment mainly focused on four expression types, i.e.,  
336 "happiness", "disgust", "surprise", and "repression".

337 Following the LOSO protocol, we compared the effectiveness of the MCNet with several other  
338 state-of-the-art methods, utilizing prediction accuracy and F1-score as the evaluation metrics. The  
339 results clearly demonstrated that the proposed MCNet method achieved competitive performance.  
340 In the experimental process, the model was initially trained using meta-learning methods with spe-  
341 cific task settings. Each meta-task is divided into a support set and a query set, with the support  
342 set used for model training and the query set used for evaluating its performance. In the assess-  
343 ment of the model's performance, accuracy was used as the primary performance metric. Accuracy  
344 represents the proportion of correct classifications by the model for a given task.

$$Accuracy = \frac{TP + FN}{TP + TN + FP + FN}, \quad (6)$$

345 where TP stands for True Positives, TN stands for True Negatives, FP stands for False Positives,  
346 and FN stands for False Negatives. Additionally, a confusion matrix was employed to provide a  
347 more detailed breakdown of the classification results. The confusion matrix reveals the model's  
348 classification performance for each category, including true positives, true negatives, false posi-

349 tives, and false negatives. In addition to accuracy and the confusion matrix, we also employ the F1  
350 score as an evaluation metric.

$$\begin{aligned} Precision &= \frac{TP}{TP + FP}, Recall = \frac{TP}{TP + FN}, \\ F1 &= \frac{2 \times (Precision \times Recall)}{Precision + Recall}, \end{aligned} \quad (7)$$

351 where P represents Precision, R represents Recall, and F1 represents the F1 Score. The F1 score  
352 combines precision and recall into a weighted average, providing a comprehensive assessment of  
353 the model’s performance. Precision represents the proportion of true positives among samples pre-  
354 dicted as positive, while recall represents the proportion of true positives among all actual positive  
355 samples. In the experimental results, a detailed report of accuracy, confusion matrix, and F1 scores  
356 was provided for the relevant test tasks, allowing for a comprehensive understanding of the model’s  
357 performance across tasks and evaluating its strengths and weaknesses.

### 358 3.1 Experimental setup

359 For the meta-learning sampling process for each dataset, we set  $K = 5$ , and the value of  $N$  is  
360 determined based on the number of classes in the dataset. For example, for the SMIC dataset,  
361 we sampled 3-way 5-shot tasks, and for the CASME and CASME II datasets, we sampled 4-way  
362 5-shot tasks. In the pre-processing phase of the deep feature extraction network, we initialized  
363 the parameters with ResNet18 pre-trained on the ImageNet dataset, and trained the model for 50  
364 epochs with a batch size of 128, using an NVIDIA GeForce RTX 3090 GPU. During the meta-  
365 learning phase, we kept the learning rate fixed at 0.01 and used a batch size of 8 meta-tasks.  
366 For the SMIC dataset, where the vertex frames were not annotated, we used the middle frame of  
367 each micro-expression sequence as the vertex frame. For fusing the priors of micro- and macro-

368 expressions, the hyperparameter  $\gamma$  was learned independently for each dataset. In this study, we  
369 randomly selected a subject set  $Z_{val}$  from the dataset  $Z_{total}$  to learn the optimal  $\gamma$ , which maximized  
370 the prediction accuracy on the samples from  $Z_{val}$ . The value of  $\gamma$  was chosen between 0.1 and 1  
371 with a step size of 0.1. Notably, once the subject was picked for a dataset, the samples of  $Z_{val}$  were  
372 not involved in the accuracy validation during the testing phase, i.e.,  $Z_{total} - Z_{val} = Z_{train} + Z_{test}$ .

### 373 *3.2 Results and analysis*

374 To validate the effectiveness of the proposed meta-clustering network, we conducted comprehen-  
375 sive comparisons with the state-of-the-art methods using LOSO protocol. The methods of com-  
376 parison follow various evaluation protocols, including LOSO and Leave-One-Video-Out (LOVO).  
377 LOSO and LOVO are two cross-validation strategies that are commonly used to evaluate the perfor-  
378 mance of machine learning models, especially in sentiment analysis, facial expression recognition,  
379 or other human behavior analysis LOSO is a method for cross-validation that is particularly useful  
380 when dealing with datasets involving multiple participants. The LOVO strategy, on the other hand,  
381 involves using the data from one video at a time as a test set and the data from all the remaining  
382 videos as a training set in a dataset containing multiple video samples. Again, this step is repeated  
383 once for each video sample to ensure that each sample has a chance to be used as test data. In this  
384 paper, LOSO is used because LOSO uses all the videos of a participant at a time as the test set, and  
385 the model does not learn the micro-expression distribution of that participant, so LOSO validation  
386 better reflects the model’s adaptability to unknown objects.

387 The comparison results on the SMIC, CASME, and CASME II datasets were listed in Ta-  
388 bles 1, 2, and 3, respectively. To enhance the clarity of the comparison results, we only retained  
389 works that employed the same classification criteria as MCNet, excluding those methods evaluat-

**Table 1** Performance comparisons with the existing methods on SMIC dataset

Method	Protocol	Task	Accuracy(%)	F1-score
Baseline <sup>6</sup>	LOSO	Negative,Positive,Surprise	45.70	0.46
LBP-SIP <sup>43</sup>	LOSO	Negative,Positive,Surprise	42.12	0.42
Selective <sup>16</sup>	LOSO	Negative,Positive,Surprise	53.60	-
STCLQP <sup>11</sup>	LOSO	Negative,Positive,Surprise	64.02	0.63
FHOFO <sup>13</sup>	LOSO	Negative,Positive,Surprise	51.83	0.52
VGGMag <sup>44</sup>	LOSO	Negative,Positive,Surprise	59.75	0.58
Bi-WOOF <sup>15</sup>	LOSO	Negative,Positive,Surprise	62.80	-
3D-FCNN <sup>28</sup>	LOVO	Negative,Positive,Surprise	55.49	-
SSSN <sup>26</sup>	LOSO	Negative,Positive,Surprise	63.41	0.63
DSSN <sup>26</sup>	LOSO	Negative,Positive,Surprise	63.41	0.65
STRCN-A <sup>45</sup>	LOSO	Negative,Positive,Surprise	53.10	0.51
SVM+revised HOOFF <sup>46</sup>	LOSO	Negative,Positive,Surprise	60.67	-
LGCconD <sup>47</sup>	LOSO	Negative,Positive,Surprise	63.41	0.62
FR <sup>48</sup>	LOSO	Negative,Positive,Surprise	57.90	-
Meta-MMFNet <sup>24</sup>	LOSO	Negative,Positive,Surprise	63.13	-
Dual-ATME <sup>32</sup>	LOSO	Negative,Positive,Surprise	64.60	-
SVM+LCBP-STGCN <sup>36</sup>	LOSO	Negative,Positive,Surprise	75.51	0.74
MFVAN <sup>33</sup>	LOSO	Negative,Positive,Surprise	<b>79.87</b>	<b>0.8009</b>
Ours	LOSO	Negative,Positive,Surprise	<b>65.63</b>	<b>0.65</b>

**Table 2** Performance comparisons with the existing methods on CASME dataset

Method	Protocol	Task	Accuracy(%)	F1-score
Baseline <sup>6</sup>	LOSO	Disgust, Surprise, Repression, Tense	40.35	0.26
LBP-SIP <sup>43</sup>	LOSO	Disgust, Surprise, Repression, Tense	36.84	0.33
STCLQP <sup>11</sup>	LOSO	Disgust, Surprise, Repression, Tense	57.31	0.50
FHOFO <sup>13</sup>	LOSO	Disgust, Surprise, Repression, Tense	65.99	0.54
VGGMag <sup>44</sup>	LOSO	Disgust, Surprise, Repression, Tense	60.23	0.58
3D-FCNN <sup>28</sup>	LOVO	Disgust, Surprise, Repression, Tense	54.44	-
LGCconD <sup>47</sup>	LOSO	Disgust, Surprise, Repression, Tense	57.31	0.54
LGCcon <sup>47</sup>	LOSO	Disgust, Surprise, Repression, Tense	60.82	0.60
Meta-MMFNet <sup>24</sup>	LOSO	Disgust, Surprise, Repression, Tense	69.59	-
GPT-4V <sup>49</sup>	-	Disgust, Surprise, Repression, Tense	36.93	-
SVM+LCBP-STGCN <sup>36</sup>	LOSO	Disgust, Surprise, Repression, Tense	<b>81.26</b>	<b>0.77</b>
Ours	LOSO	Disgust, Surprise, Repression, Tense	<b>70.27</b>	<b>0.70</b>

390 ing using different sets of emotions from the proposed method. From these tables, it was able to  
391 be observed that the proposed MCNet achieves competitive recognition performance compared to  
392 baselines and other existing methods. Specifically, the prediction accuracy reached 65.63%, and  
393 the F1-score reached 0.65 on the SMIC dataset, showing 1.03% improvement over the accuracy

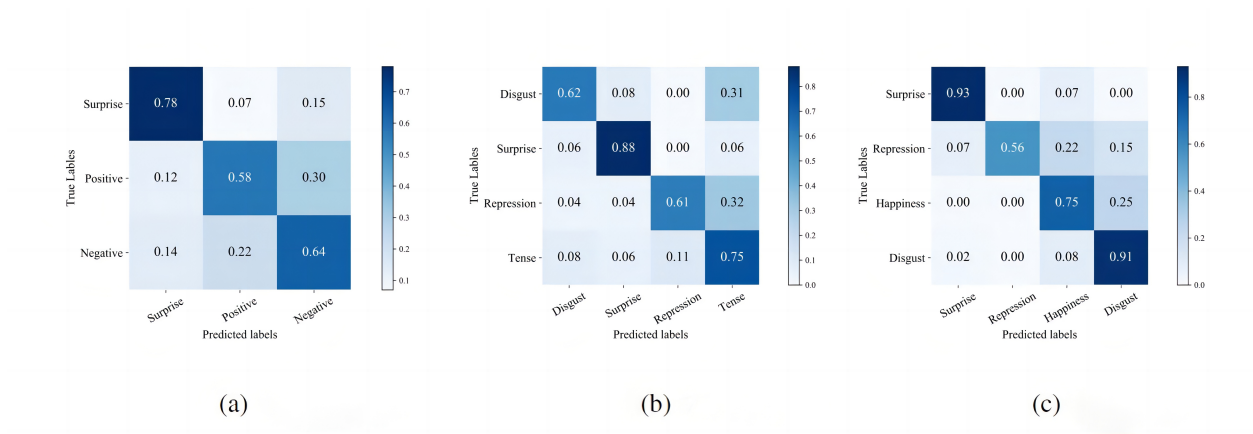
**Table 3** Performance comparisons with the existing methods on CASME II dataset

Method	Protocol	Task	Accuracy(%)	F1-score
DTCM <sup>50</sup>	LOSO	Happiness, Surprise, Disgust, Repression	72.06	-
Meta-MMFNet <sup>24</sup>	LOSO	Happiness, Surprise, Disgust, Repression	80.95	-
Ours	LOSO	Happiness, Surprise, Disgust, Repression	<b>81.63</b>	<b>0.81</b>

394 (64.60%) achieved by the latest Dual-ATME method.<sup>33</sup> For the CASME dataset, the prediction  
395 accuracy was 70.27%, and the F1-score was 0.7, indicating an improvement of 0.68% in accuracy  
396 over the Meta-MMFNet method<sup>24</sup> and an increase of 0.1 in F1-score over the LGCconD method<sup>47</sup>.  
397 Moreover, on the CASME II dataset, the prediction accuracy achieved 81.63%, and the F1-score  
398 achieved 0.81, which represents an accuracy improvement of 0.68% over the performance of the  
399 Meta-MMFNet.

400 -While the MCNet method on SMIC dataset was less effective than the optimal MFVAN  
401 method<sup>33</sup>, it achieved significantly higher recognition scores than MFVAN method on the CASME II  
402 dataset. Additionally, on the CASME dataset, the SVM+LCBP-STGCN method<sup>36</sup> outperformed  
403 the proposed method in accuracy and F1-score by 10.99% and 0.07, respectively. However, the  
404 proposed MCNet method achieved the best recognition results for the four classification tasks of  
405 “happiness”, “surprise”, “disgust”, and “repression” on the CASME II dataset. The evaluation re-  
406 sults across all three datasets demonstrated that the micro-expression recognition method proposed  
407 in this study exhibited good generalization ability and achieved competitive performance on these  
408 datasets.

409 The results also indicated that the proposed MCNet method achieved the most promising per-  
410 formance on the CASME II dataset. We argued that this outcome is attributed to the balanced distri-  
411 bution of the four categories in the CASME II dataset. In contrast, the SMIC and CASME datasets  
412 suffered from severe class imbalance issues, which negatively impact the recognition performance.  
413 Additionally, the absence of annotated apex frames in the SMIC dataset poses a challenge. The



**Fig 3** Confusion matrices of the proposed method on three validation datasets. (a) Confusion matrix on the SMIC dataset. (b) Confusion matrix on the CASME dataset. (c) Confusion matrix on the CASME II dataset.

414 apex frame, which exhibited the greatest variation in the micro-expression sequence, was crucial  
 415 for determining the type and strength of micro-expressions. Although the middle frame between  
 416 the onset and offset frames was used as an approximation, it introduced some errors, leading to a  
 417 relatively lower recognition rate for micro-expressions in the SMIC dataset.

418 To further analyze the performance of the proposed method, we visualized the confusion ma-  
 419 trices for each dataset in Fig. 3. These matrices provided insights into the recognition accuracy  
 420 of each micro-expression category and the probability of mis-classification into other categories.  
 421 In these figures, the horizontal axis represented the predicted class labels, while the vertical axis  
 422 represented the true class labels. From the confusion matrices, it was evident that the “surprise”  
 423 expression consistently achieves the highest recognition rate in all three datasets, reaching 78%,  
 424 88%, and 93% on the SMIC, CASME, and CASME II datasets, respectively. On the other hand,  
 425 mis-classifications often occurred between positive and negative emotions. For example, the “sur-  
 426 prise” and “positive” expressions in the SMIC dataset were easily mis-predicted as “negative”,  
 427 primarily because the majority of samples belong to the negative category.

## 428 4 Conclusion

429 This study proposed a novel micro-expression recognition method based on the meta-clustering  
430 network, aiming to address the challenge of mis-classification caused by large intra-class distances  
431 in micro-expressions. Different from the standard centroid-based strategy used in previous meta-  
432 learning methods, the proposed method effectively took into account the variations within each  
433 class and introduced a hierarchical clustering module that produced multiple prototypes instead  
434 of the mean prototypes. This effectively reduced the impact of intra-class differences in micro-  
435 expression samples. Additionally, The proposed approach not only integrated two dynamic fea-  
436 tures but also extracted deep feature representations of micro-expressions, enabling a more com-  
437 prehensive and robust understanding of the facial expressions. To alleviate the problem of limited  
438 micro-expression datasets, the proposed MCNet method leveraged transfer learning by pre-training  
439 on a macro-expression dataset with ample data and then adapting to the micro-expression target  
440 domain. Moreover, the fused priors of micro-expression and macro-expression in the meta-testing  
441 phase further contributed to improving the classification accuracy of micro-expressions. Com-  
442 prehensive experiments were conducted on multiple datasets, including the SMIC, CASME, and  
443 CASME II datasets, and the comparative results highlighted the effectiveness of the proposed  
444 model in the field of micro-expression recognition.

445 In future works, the proposed method holds potential for enhancement. The current methodol-  
446 ogy involves pre-training deep feature extraction networks,  $f_{\alpha}^{micro}$  and  $f_{\beta}^{macro}$ , on micro-expression  
447 and CK+ datasets, respectively, as an attempt to integrate prior knowledge of micro-expressions  
448 and macro-expressions. This approach hinders the end-to-end training of the proposed MCNet  
449 method. Additionally, the challenge of collecting and annotating micro-expression training data

450 persists, leading to issues of data insufficiency and class imbalance in recognition tasks. Although  
451 the proposed MCNet, employing few-shot learning, mitigates the problem of limited amount of  
452 data to some extent, its accuracy remains to be further improved. Furthermore, the high similarity  
453 among sub-classes in negative micro-expressions contributes to an elevated mis-classification rate.  
454 In the forthcoming research, we will introduce information on facial action unit combinations and  
455 leverage attention mechanisms to discern more nuanced negative micro-expression classes, aiming  
456 for further enhancement in recognition performance.

### 457 **Data Used In This Manuscript**

458 Availability of data used in this study is limited, which were used under licence for this study. The  
459 "SMIC" data that support the findings of the study are available on request at [https://www.oulu.fi/en/university/faculty-and-units/faculty-information-technology-and-electrical-engineering/center-machine-vision-and-signal-](https://www.oulu.fi/en/university/faculty-and-units/faculty-information-technology-and-electrical-engineering/center-machine-vision-and-signal-analysis)  
460 [analysis](https://www.oulu.fi/en/university/faculty-and-units/faculty-information-technology-and-electrical-engineering/center-machine-vision-and-signal-analysis), the "CASME" are available on request at <http://casme.psych.ac.cn/casme/c1>, and the  
461 CASME II are available on request at <http://casme.psych.ac.cn/casme/c2>

### 463 **Acknowledgment**

464 The authors acknowledgment the support by the Natural Science Foundation of Shandong Province  
465 under the project **No. ZR2023MF041**, the Spanish Ministry of Economy and Competitiveness  
466 (MINECO) and the European Regional Development Fund (ERDF) under Grants PID2020-120311RB-  
467 I00 funded by MCIN/AEI/10.13039/501100011033.

### 468 *References*

469 1 R. Arya, J. Singh, and A. Kumar, "A survey of multidisciplinary domains contributing to  
470 affective computing," *Computer Science Review* **40**, 100399 (2021).

- 471 2 W.-J. Yan, Q. Wu, J. Liang, *et al.*, “How fast are the leaked facial expressions: The duration  
472 of micro-expressions,” *Journal of Nonverbal Behavior* **37**, 217–230 (2013).
- 473 3 T. A. Russell, E. Chu, and M. L. Phillips, “A pilot study to investigate the effectiveness of  
474 emotion recognition remediation in schizophrenia using the micro-expression training tool,”  
475 *British journal of clinical psychology* **45**(4), 579–583 (2006).
- 476 4 P. Ekman and E. L. Rosenberg, Eds., *What the face reveals: Basic and applied studies of*  
477 *spontaneous expression using the Facial Action Coding System (FACS)*, New York (2005).
- 478 5 M. Frank, M. Herbasz, K. Sinuk, *et al.*, “I see how you feel: Training laypeople and profes-  
479 sionals to recognize fleeting emotions,” in *The annual meeting of the international communi-*  
480 *cation association. Sheraton New York, New York City*, 1–35 (2009).
- 481 6 X. Li, T. Pfister, X. Huang, *et al.*, “A spontaneous micro-expression database: Inducement,  
482 collection and baseline,” in *2013 10th IEEE International Conference and Workshops on*  
483 *Automatic face and gesture recognition (fg)*, 1–6 (2013).
- 484 7 W.-J. Yan, Q. Wu, Y.-J. Liu, *et al.*, “Casm database: A dataset of spontaneous micro-  
485 expressions collected from neutralized faces,” in *2013 10th IEEE international conference*  
486 *and workshops on automatic face and gesture recognition (FG)*, 1–7 (2013).
- 487 8 W.-J. Yan, X. Li, S.-J. Wang, *et al.*, “Casm ii: An improved spontaneous micro-expression  
488 database and the baseline evaluation,” *PloS one* **9**(1), e86041 (2014).
- 489 9 A. K. Davison, C. Lansley, N. Costen, *et al.*, “Samm: A spontaneous micro-facial movement  
490 dataset,” *IEEE transactions on affective computing* **9**(1), 116–129 (2016).
- 491 10 S. Polikovsky, Y. Kameda, and Y. Ohta, “Facial micro-expressions recognition using high  
492 speed camera and 3d-gradient descriptor,” *IET Conference Proceedings* , 16–16(1) (2009).

- 493 11 X. Huang, G. Zhao, X. Hong, *et al.*, “Spontaneous facial micro-expression analysis using  
494 spatiotemporal completed local quantized patterns,” *Neurocomputing* **175**, 564–578 (2016).
- 495 12 C. Zach, T. Pock, and H. Bischof, “A duality based approach for realtime tv-l 1 optical flow,”  
496 in *Pattern Recognition: 29th DAGM Symposium, Heidelberg, Germany, September 12-14,*  
497 *2007. Proceedings 29*, 214–223 (2007).
- 498 13 S. Happy and A. Routray, “Fuzzy histogram of optical flow orientations for micro-expression  
499 recognition,” *IEEE Transactions on Affective Computing* **10**(3), 394–406 (2017).
- 500 14 Y.-J. Liu, J.-K. Zhang, W.-J. Yan, *et al.*, “A main directional mean optical flow feature for  
501 spontaneous micro-expression recognition,” *IEEE Transactions on Affective Computing* **7**(4),  
502 299–310 (2015).
- 503 15 S.-T. Liong, J. See, K. Wong, *et al.*, “Less is more: Micro-expression recognition from video  
504 using apex frame,” *Signal Processing: Image Communication* **62**, 82–92 (2018).
- 505 16 D. Patel, X. Hong, and G. Zhao, “Selective deep features for micro-expression recognition,”  
506 in *2016 23rd international conference on pattern recognition (ICPR)*, 2258–2263 (2016).
- 507 17 D. H. Kim, W. J. Baddar, and Y. M. Ro, “Micro-expression recognition with expression-  
508 state constrained spatio-temporal feature representations,” in *Proceedings of the 24th ACM*  
509 *international conference on Multimedia*, 382–386 (2016).
- 510 18 M. A. Takalkar and M. Xu, “Image based facial micro-expression recognition using deep  
511 learning on small datasets,” in *2017 international conference on digital image computing:*  
512 *techniques and applications (DICTA)*, 1–7 (2017).
- 513 19 M. Peng, Z. Wu, Z. Zhang, *et al.*, “From macro to micro expression recognition: Deep learn-

- 514 ing on small datasets using transfer learning,” in *2018 13th IEEE International Conference*  
515 *on Automatic Face & Gesture Recognition (FG 2018)*, 657–661 (2018).
- 516 20 N. Van Quang, J. Chun, and T. Tokuyama, “Capsulenet for micro-expression recognition,”  
517 in *2019 14th IEEE International Conference on Automatic Face & Gesture Recognition (FG*  
518 *2019)*, 1–7 (2019).
- 519 21 R. Zhi, H. Xu, M. Wan, *et al.*, “Combining 3d convolutional neural networks with transfer  
520 learning by supervised pre-training for facial micro-expression recognition,” *IEICE TRANS-*  
521 *ACTIONS on Information and Systems* **102**(5), 1054–1064 (2019).
- 522 22 B. Sun, S. Cao, D. Li, *et al.*, “Dynamic micro-expression recognition using knowledge distil-  
523 lation,” *IEEE Transactions on Affective Computing* **13**(2), 1037–1043 (2020).
- 524 23 G. Zhao, X. Huang, M. Taini, *et al.*, “Facial expression recognition from near-infrared  
525 videos,” *Image and vision computing* **29**(9), 607–619 (2011).
- 526 24 W. Gong, Y. Zhang, W. Wang, *et al.*, “Meta-mmfnnet: Meta-learning based multi-model fusion  
527 network for micro-expression recognition,” *ACM Transactions on Multimedia Computing,*  
528 *Communications, and Applications (TOMM)* **20**(39), 1–20 (2022).
- 529 25 M. Peng, C. Wang, T. Chen, *et al.*, “Dual temporal scale convolutional neural network for  
530 micro-expression recognition,” *Frontiers in psychology* **8**, 1745 (2017).
- 531 26 H.-Q. Khor, J. See, S.-T. Liong, *et al.*, “Dual-stream shallow networks for facial micro-  
532 expression recognition,” in *2019 IEEE international conference on image processing (ICIP)*,  
533 36–40 (2019).
- 534 27 H.-Q. Khor, J. See, R. C. W. Phan, *et al.*, “Enriched long-term recurrent convolutional net-

- 535 work for facial micro-expression recognition,” in *2018 13th IEEE international conference*  
536 *on automatic face & gesture recognition (FG 2018)*, 667–674 (2018).
- 537 28 J. Li, Y. Wang, J. See, *et al.*, “Micro-expression recognition based on 3d flow convolutional  
538 neural network,” *Pattern Analysis and Applications* **22**, 1331–1339 (2019).
- 539 29 Y. S. Gan, S.-T. Liong, W.-C. Yau, *et al.*, “Off-apexnet on micro-expression recognition sys-  
540 tem,” *Signal Processing: Image Communication* **74**, 129–139 (2019).
- 541 30 S. Zhao, H. Tao, Y. Zhang, *et al.*, “A two-stage 3d cnn based learning method for spontaneous  
542 micro-expression recognition,” *Neurocomputing* **448**, 276–289 (2021).
- 543 31 M. Tang, M. Ling, J. Tang, *et al.*, “A micro-expression recognition algorithm based on feature  
544 enhancement and attention mechanisms,” *Virtual Reality* **27**, 2405–2416 (2023).
- 545 32 H. Zhou, S. Huang, J. Li, *et al.*, “Dual-atme: Dual-branch attention network for micro-  
546 expression recognition,” *Entropy* **25**(3), 460 (2023).
- 547 33 H. Pan, H. Yang, L. Xie, *et al.*, “Multi-scale fusion visual attention network for facial micro-  
548 expression recognition,” *Frontiers in Neuroscience* **17** (2023).
- 549 34 D. Hao, M. Zhu, C. Zhang, *et al.*, “A lightweight attention-based network for micro-  
550 expression recognition,” *Multimedia Tools and Applications* , 1–22 (2023).
- 551 35 S. Thuseethan, S. Rajasegarar, and J. Yearwood, “Deep3dcann: A deep 3dcnn-ann framework  
552 for spontaneous micro-expression recognition,” *Information Sciences* **630**, 341–355 (2023).
- 553 36 Y. Wang, J. Han, and Z. Guo, “Lcbp-stgcn: A local cube binary pattern spatial temporal graph  
554 convolutional network for micro-expression recognition,” *Journal of Intelligent & Fuzzy Sys-*  
555 *tems* **44**(2), 1601–1611 (2023).

- 556 37 C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep  
557 networks,” 1126–1135 (2017).
- 558 38 S. Ravi and H. Larochelle, “Optimization as a model for few-shot learning,” (2016).
- 559 39 F. Murtagh and P. Contreras, “Algorithms for hierarchical clustering: an overview,” *Wiley*  
560 *Interdisciplinary Reviews: Data Mining and Knowledge Discovery* **2**(1), 86–97 (2012).
- 561 40 H. Lu, K. Kpalma, and J. Ronsin, “Motion descriptors for micro-expression recognition,”  
562 *Signal Processing: Image Communication* **67**, 108–117 (2018).
- 563 41 C. Wang, M. Peng, T. Bi, *et al.*, “Micro-attention for micro-expression recognition,” *Neuro-*  
564 *computing* **410**, 354–362 (2020).
- 565 42 X. Jia, X. Ben, H. Yuan, *et al.*, “Macro-to-micro transformation model for micro-expression  
566 recognition,” *Journal of computational science* **25**, 289–297 (2018).
- 567 43 Y. Wang, J. See, R. C.-W. Phan, *et al.*, “Lbp with six intersection points: Reducing redundant  
568 information in lbp-top for micro-expression recognition,” in *Computer Vision–ACCV 2014:*  
569 *12th Asian Conference on Computer Vision, Singapore, Singapore, November 1-5, 2014,*  
570 *Revised Selected Papers, Part I 12*, 525–537 (2015).
- 571 44 Y. Li, X. Huang, and G. Zhao, “Can micro-expression be recognized based on single apex  
572 frame?,” in *2018 25th IEEE International Conference on Image Processing (ICIP)*, 3094–  
573 3098 (2018).
- 574 45 Z. Xia, X. Hong, X. Gao, *et al.*, “Spatiotemporal recurrent convolutional networks for rec-  
575 ognizing spontaneous micro-expressions,” *IEEE Transactions on Multimedia* **22**(3), 626–640  
576 (2020).

- 577 46 Q. Li, J. Yu, T. Kurihara, *et al.*, “Deep convolutional neural network with optical flow for  
578 facial micro-expression recognition,” *Journal of Circuits, Systems and Computers* **29**(01),  
579 2050006 (2020).
- 580 47 Y. Li, X. Huang, and G. Zhao, “Joint local and global information learning with single apex  
581 frame detection for micro-expression recognition,” *IEEE Transactions on Image Processing*  
582 **30**, 249–263 (2020).
- 583 48 L. Zhou, Q. Mao, X. Huang, *et al.*, “Feature refinement: An expression-specific feature learn-  
584 ing and fusion method for micro-expression recognition,” *Pattern Recognition* **122**, 108275  
585 (2022).
- 586 49 Z. Lian, L. Sun, H. Sun, *et al.*, “Gpt-4v with emotion: A zero-shot benchmark for multimodal  
587 emotion understanding,” *arXiv preprint arXiv:2312.04293* (2023).
- 588 50 Z. Lu, Z. Luo, H. Zheng, *et al.*, “A delaunay-based temporal coding model for micro-  
589 expression recognition,” in *Computer Vision-ACCV 2014 Workshops: Singapore, Singapore,*  
590 *November 1-2, 2014, Revised Selected Papers, Part II 12*, 698–711 (2015).

## 591 **Biography**

592 **Ziqi Wang** received the B.S. degree from Shandong University of Science and Technology, China,  
593 in 2021. He is currently pursuing the master’s degree with the Qingdao Institute of software, col-  
594 lege of computer science and technology, China University of Petroleum (East China), China. His  
595 research interests include micro-expression recognition and multimedia information fusion pro-  
596 cessing.

597 **Wenwen Fu** received the B.E. degree from University of Jinan, China, in 2020. She is currently  
598 pursuing the master’s degree with the Qingdao Institute of Software, College of Computer Science

599 and Technology, China University of Petroleum (East China), China. Her research interests in-  
600 clude micro-expression recognition and cross-modality facial expression generation.

601 **Yue Zhang** received her Bachelor of Engineering degree from Jinan University, China in 2020 and  
602 her Master of Engineering degree from China University of Petroleum (East China) with 2023.  
603 Her research interest is micro-expression recognition.

604 **JiaRui Li** received the B.E. degree from Shandong University of Science and Technology, China,  
605 in 2023. He is currently pursuing the master's degree with the Qingdao Institute of Software, Col-  
606 lege of Computer Science and Technology, China University of Petroleum (East China), China.  
607 His research interests include Multimodal fusion (images, text, etc.) and Cross-modal alignment  
608 tasks.

609 **Wenjuan Gong** received the Ph.D. (cum laude) degree from the Autonomous University of Barcelona,  
610 in 2013. She was a Postdoctoral Research Assistant with Oxford Brookes University, in 2014. She  
611 is currently an associate Professor with the China University of Petroleum (East China). She  
612 participated in the European Project Consolider Ingenio 2010 and the EPSRC Project Tensorial  
613 Modeling of Dynamical Systems for Gait and Activity Recognition. She has led the CCF-Tencent  
614 Funds and the Natural Science Foundation of China of Shandong Province. Her research interests  
615 include computer vision and machine learning.

616 **Jordi González** received the Ph.D. degree in computer engineering from the Universitat Autònoma  
617 de Barcelona (UAB), in 2004. He is currently a Professor in computer science with the Department  
618 of Computer Science, UAB. He is also a Research Fellow with Computer Vision Center, where he  
619 has cofounded three spin-offs, namely Cloud Size Services, Visual Tagging, and Care Respite, and  
620 the Image Sequence Evaluation (ISE Lab) Research Group. His research interest includes machine  
621 learning techniques for the computational interpretation of social images, or visual hermeneutics.

## 622 **List of Figures**

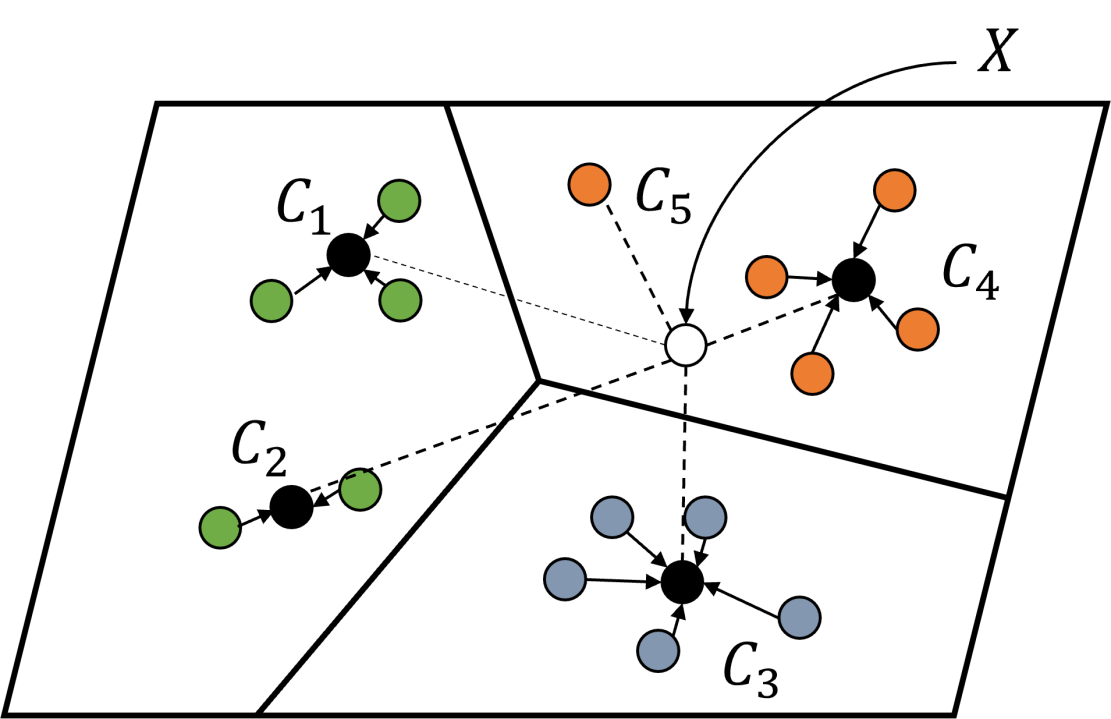
623 1 The process of predicting sample categories using the proposed meta-clustering  
624 method. Instead of using a single mean feature prototype for each class, the method  
625 determines the class prototypes by clustering the two nearest samples within the  
626 same class in the support set, resulting in  $N$  or more class prototypes, where  $N$   
627 represents the number of classes in the support set. To predict the label of a query  
628 sample  $X$  during meta-testing, the meta-clustering method compares the distances  
629 between the sample and each class prototype and identifies the closest class proto-  
630 type to the sample. The label of the test sample is obtained by mapping the nearest  
631 class prototype back to its original category.

632 2 The overall architecture of the meta-clustering learning network. This method  
633 comprises of two main stages: pre-training the deep feature extraction module and  
634 learning the meta-clustering classification model. In the pre-training stage, two  
635 independent deep feature extraction networks based on ResNet18 are trained on  
636 the micro-expression dataset  $\Omega_{train}$  and the macro-expression dataset CK+. In the  
637 meta-clustering stage, the fused dynamic features are fed into the micro-pipeline  
638 and macro-pipeline and the deep features of micro-expressions are then extracted  
639 in each pipeline. The cosine distance between the query samples and the class pro-  
640 totypes obtained from the class prototype calculation module is computed, yielding  
641 the micro-score and macro-score, respectively. Based on these two score vectors,  
642 we optimize the micro-pipeline and macro-pipeline separately using cross-entropy  
643 loss. during the validation process, we combine the score vectors generated from  
644 both pipelines to merge the micro- and macro-predictions to predict the final micro-  
645 expression category.

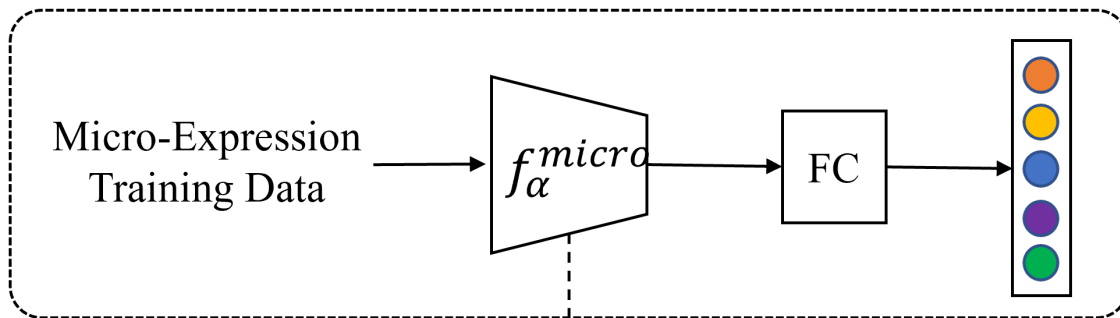
646 3 Confusion matrices of the proposed method on three validation datasets. (a) Con-  
647 fusion matrix on the SMIC dataset. (b) Confusion matrix on the CASME dataset.  
648 (c) Confusion matrix on the CASME II dataset.

## 649 List of Tables

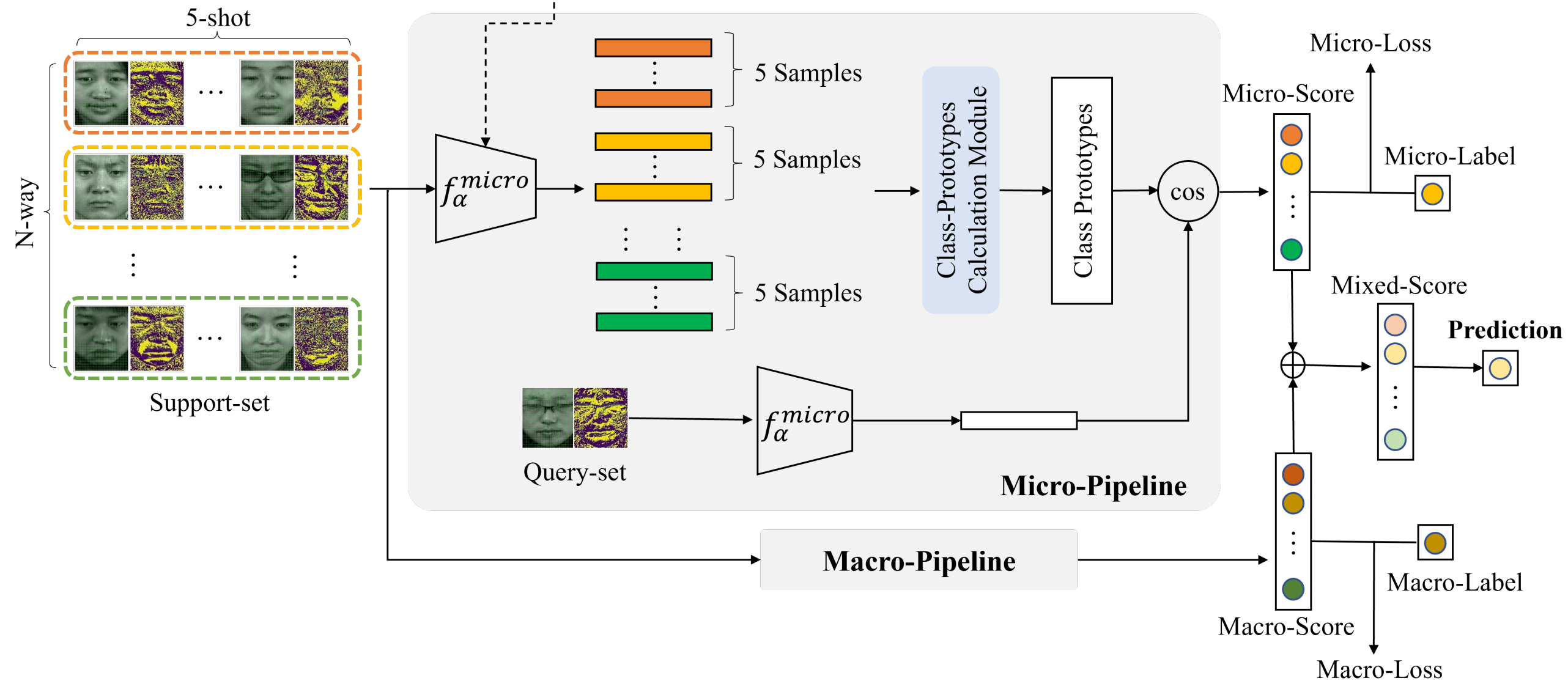
- 650 1 Performance comparisons with the existing methods on SMIC dataset
- 651 2 Performance comparisons with the existing methods on CASME dataset
- 652 3 Performance comparisons with the existing methods on CASME II dataset

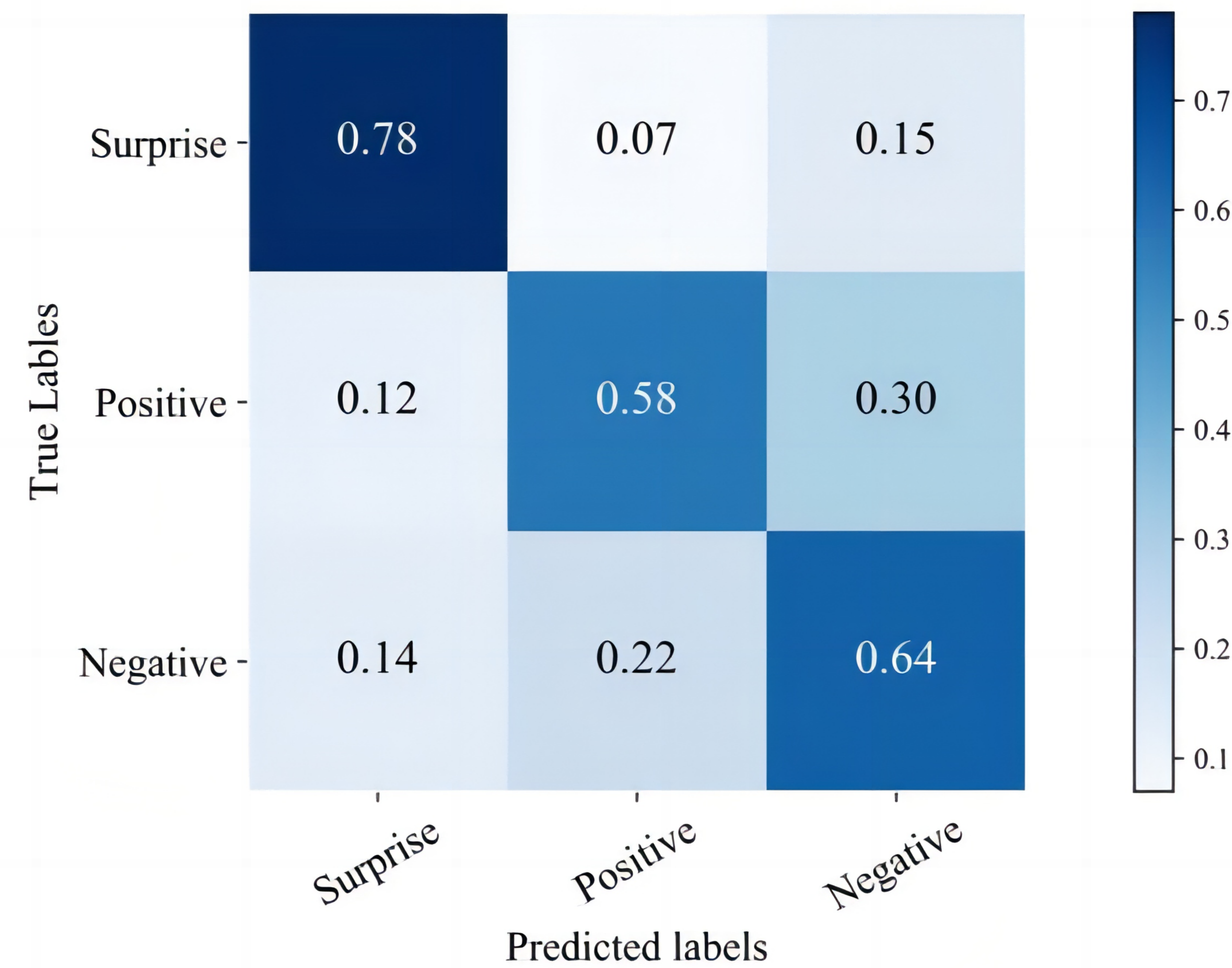


## Pre-training Deep Feature Extraction Module

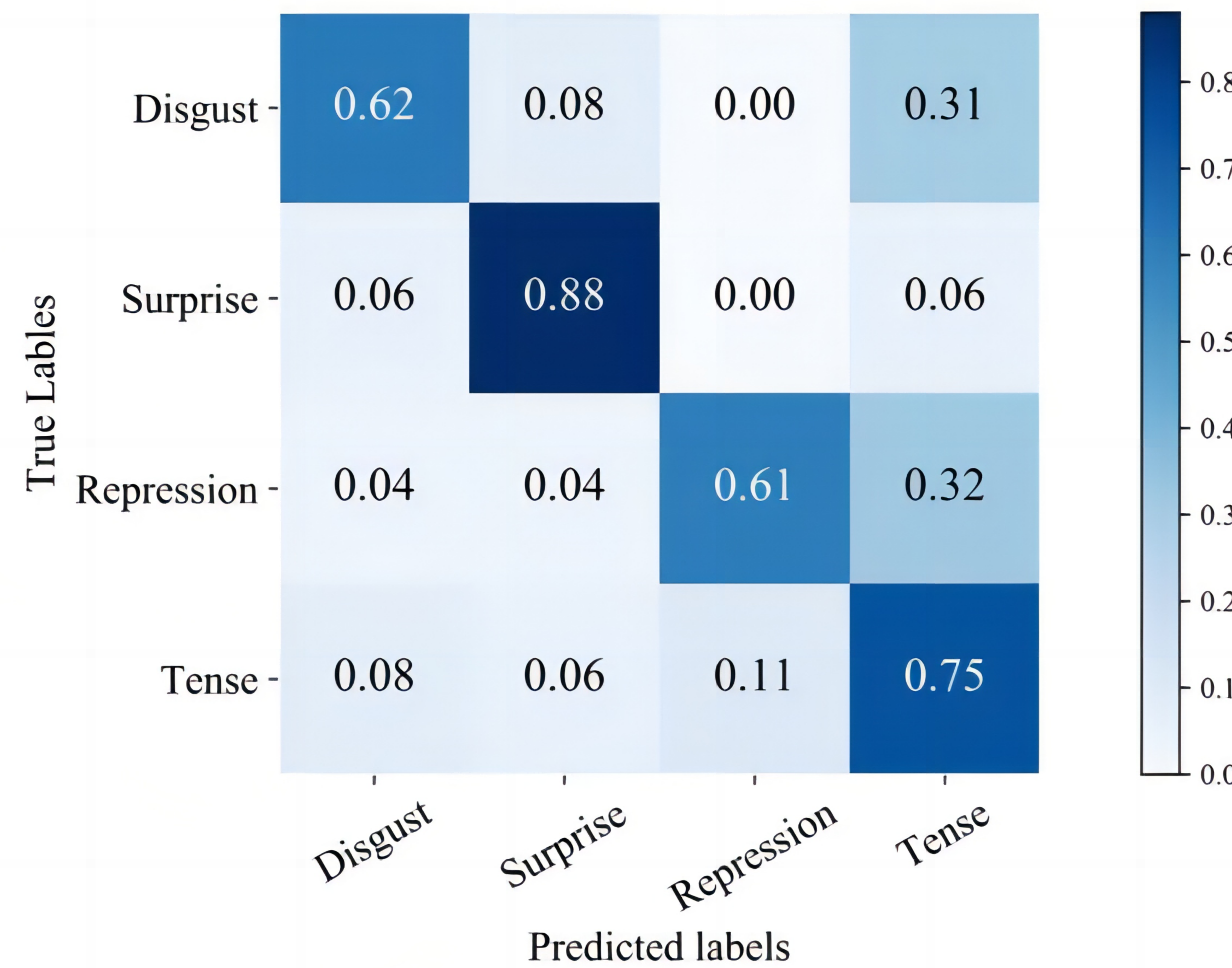


## Meta-Clustering Learning Model

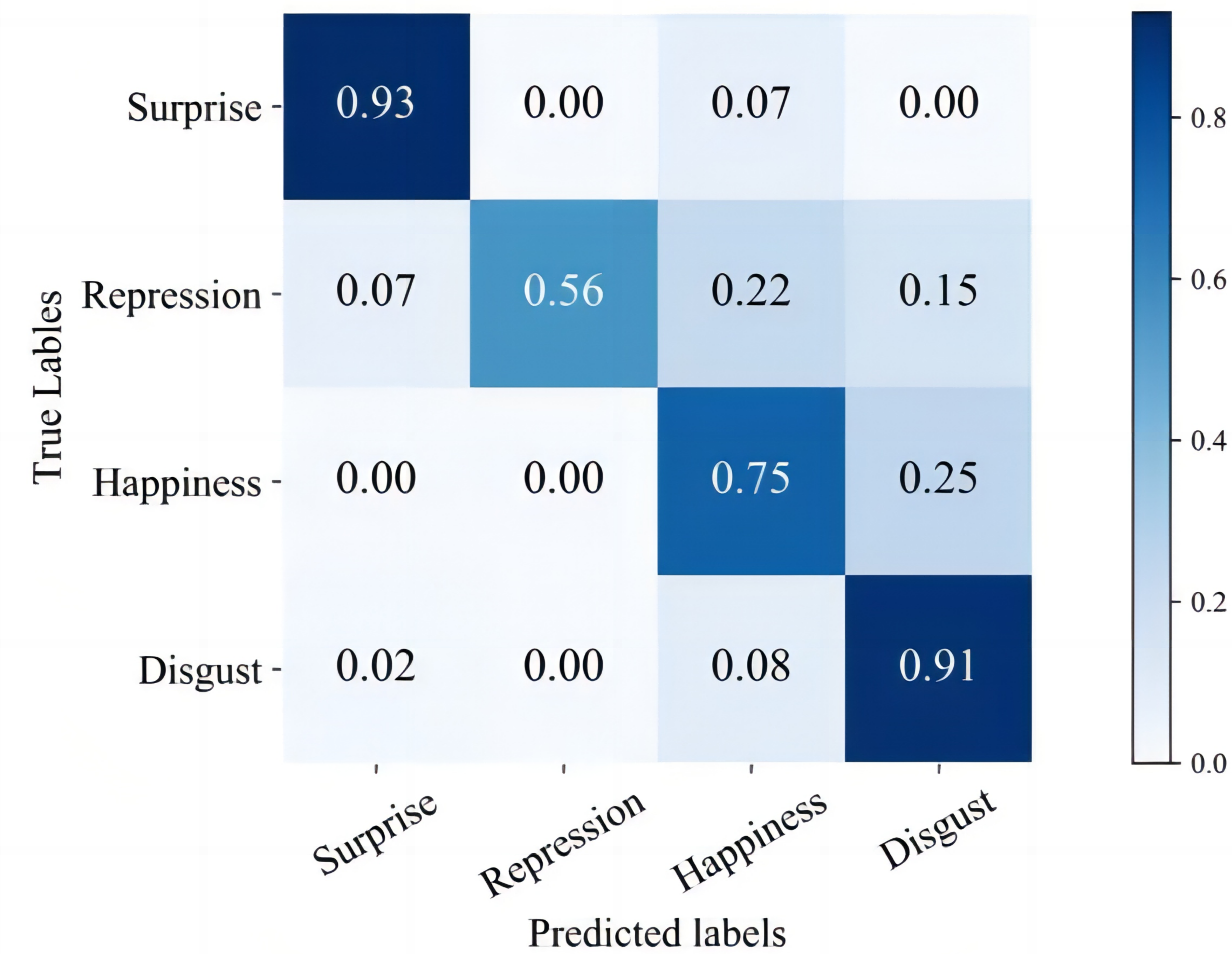




(a)



(b)



(c)